



SASTRA

ENGINEERING · MANAGEMENT · LAW · SCIENCES · HUMANITIES · EDUCATION

DEEMED TO BE UNIVERSITY

(U/S 3 of the UGC Act, 1956)



THINK MERIT | THINK TRANSPARENCY | THINK SASTRA

T H A N J A Y U R · K U M B A K O N A M · C H E N N A I

Department of Computer Science and Engineering

SRINIVASA RAMANUJAN CENTRE

KUMBAKONAM, TAMIL NADU, INDIA – 612001

Bonafide Certificate

This is to certify that the report titled "EXTRACTIVE TEXT SUMMARIZATION USING LLMS AND HIERARCHICAL POSITIONAL ENCODING" submitted as a requirement for the course, CSE300: MINI PROJECT for B.Tech. is a Bonafide record of the work done by Eggoni Ganesh (Reg. No. 226003043), Gaddapara Manohar (Reg. No. 226003047) and Kamanuru Ajay Krishna (Reg. No. 226003176) during the academic year 2024-25, in the Srinivasa Ramanujan Centre, under my supervision.

Signature of Project Supervisor

Name with Affiliation

: Smt. Hemamalini. S /AP II/CSE/SRC/SASTRA

Date

: 25.04.25

Mini Project Viva voce held on 28/04/2025

M. Vanitha
28/04/25

Examiner 1

Dr. M. VANITHA
AP-II /CSE /SRC .

ii

P. Vimala Devi
28/04/25

P. Vimala Devi
AP-II/CSE /SEC

EXTRACTIVE TEXT SUMMARIZATION USING LLMS AND HIERARCHICAL POSITIONAL ENCODING

*Report submitted to the SASTRA Deemed to be University
as the requirement for the course*

CSE300 - MINI PROJECT

Submitted by

Eggoni ganesh
(Reg. No.: 226003043, B.Tech., CSE)
Gaddapara manohar
(Reg. No.: 226003047, B.Tech., CSE)
Kamanuru ajay krishna
(Reg. No.: 226003175, B.Tech., CSE)

*under the guidance of
Smt. Hemamalini. S
AP II CSE*



Department of Computer Science and Engineering

SRINIVASA RAMANUJAN CENTRE

KUMBAKONAM, TAMIL NADU, INDIA – 612001

MAY 2025



SASTRA

ENGINEERING • MANAGEMENT • LAW • SCIENCES • HUMANITIES • EDUCATION

DEEMED TO BE UNIVERSITY

(U/S 3 of the UGC Act, 1956)

THINK MERIT | THINK TRANSPARENCY | THINK SASTRA

T H A N J A V U R | K U M B A K O N A M | C H E N N A I



Department of Computer Science and Engineering

SRINIVASA RAMANUJAN CENTRE

KUMBAKONAM, TAMIL NADU, INDIA – 612001

Bonafide Certificate

This is to certify that the report titled “EXTRACTIVE TEXT SUMMARIZATION USING LLMS AND HIERARCHICAL POSITIONAL ENCODING” submitted as a requirement for the course, **CSE300: MINI PROJECT** for B.Tech. is a Bonafide record of the work done by **Eggoni Ganesh (Reg. No. 226003043)**, **Gaddapara Manohar (Reg. No. 226003047)** and **Kamanuru Ajay Krishna (Reg. No. 226003176)** during the academic year 2024-25, in the Srinivasa Ramanujan Centre, under my supervision.

Signature of Project Supervisor :

Name with Affiliation : Smt. Hemamalini. S /AP II/CSE/SRC/SASTRA

Date :

Mini Project *Viva voce* held on _____

Examiner 1

Examiner2

Acknowledgements

We would like to thank our Honorable Chancellor **Prof. R. Sethuraman** for providing us with an opportunity and the necessary infrastructure for carrying out this project as a part of our curriculum.

We would like to thank our Honorable Vice-Chancellor **Dr. S. Vaidhyasubramaniam** and **Dr. S. Swaminathan**, Dean, Planning & Development, for the encouragement and strategic support at every step of our college life.

We extend our sincere thanks to **Dr. R. Chandramouli**, Registrar, SASTRA Deemed to be University for providing the opportunity to pursue this project.

We extend our heartfelt thanks to **Dr. V. Ramaswamy**, Dean and **Dr. A. Alli Rani**, Associate Dean, Srinivasa Ramanujam Centre for their constant support and suggestions when required without any reservations.

We sincerely extend our gratitude to our guide **Smt. Hemamalini. S AP II CSE** who was the driving force behind this whole idea from the start. Her deep insight in the field and invaluable suggestions helped us in making progress throughout our project.

We would like to thank our project review panel members **Dr. M. Vanitha**, AP-II, Department of CSE and **Dr. Vimaladevi**, AP-II, Department of CSE for their valuable comments and meticulous support to complete this project successfully.

We would like to place on record, our gratitude to **Dr. M. Vanitha**, AP-II, Department of CSE, our Project Coordinator for her benevolent approach and to all the teaching and non-teaching faculty of the Department of CSE who have either directly or indirectly helped us in the completion of the project.

We gratefully acknowledge all the contributions and encouragement from our family and friends resulting in the successful completion of this project. We thank you all for giving us the opportunity to showcase our skills through this project.

Table of Contents

Title	Page No.
Bonafide Certificate	ii
Acknowledgements	iii
List of figures	v
List of tables	v
Abbreviations	vi
Notations	vii
Abstract	viii
1. Summary of the base paper	1
2. Merits and Demerits of the base paper	8
3. Source Code	9
4. Snapshots	18
5. Conclusion and Future Plans	19
6. References	20
7. Appendix – Base paper	22

List of figures

Figure No.	Title	Page No.
1	Architecture Diagram	3
2	K-means Graph	4
3	ROUGE Scores comparison table	15

List of Tables

Figure No.	Title	Page No.
1	ROUGE Scores comparison graph	15

ABBREVIATIONS

ETS	Extractive Text Summarization
NLP	Natural Language Processing
GPT-4	Generative Pre-Trained Transformer 4
LLM	Large Language Model
FCM	Fuzzy C-Means
TF-IDF	Term Frequency–Inverse Document Frequency
NER	Named Entity Recognition
SVM	Support Vector Machine
GLOVE	Global Vectors For Word Representation
LSTM	Long Short-Term Memory
CNN	Convolutional Neural Network
RNN	Recurrent Neural Network
BERT	Bidirectional Encoder Representations From Transformers
BILSTM	Bidirectional Long Short-Term Memory
MTL	Multi-Task Learning
MPNET	Masked Position Prediction Network
ROUGE	Recall-Oriented Understudy For Gisting Evaluation

NOTATIONS

Notation	Description
D	The input document to be summarized.
S	Set of all sentences in the document.
T	Set of text segments (e.g., headers, paragraphs).
t_i	Individual text segment in the document.
v_i	Semantic embedding of sentence s_i generated by MPNet.
u_j	Section embedding computed by averaging embeddings of sentences in section j .
α	Weighting factor for combining sentence and section embeddings.
Tree	Data structure representing the hierarchical organization of document segments.
n_i	Node in the hierarchical tree.
r_i	Semantic role assigned to node n_i .
Sim (n_i, n_j)	Cosine similarity between nodes n_i and n_j .
τ	Threshold for creating dynamic links between semantically similar nodes.
$T(n_i)$	Triplet position encoding of node n_i .
d_i	Depth of node n_i in the hierarchical tree.
p_i	Position index of node n_i within its parent node.
c_i	Contextual importance of node n_i .
V_{doc}	Global embedding representing the overall document.
Relevance (t_i)	Cosine similarity between embedding of t_i and V_{doc} .
AdjustedScore (t_i)	Final score for selecting a sentence based on relevance, coherence, novelty, and redundancy.
SelectionScore (s)	Score used to select sentences from clusters.
Redundancy (s_i, s_j)	Metric to detect and eliminate redundant sentences.
OrderScore ($s_i \rightarrow s_j$)	Score that helps order sentences logically in the summary.
Cohesion ($s_i \rightarrow s_j$)	Measures how well sentence s_i transitions into s_j .

ABSTRACT

Text summarization means taking a long piece of writing and making it shorter, while still keeping all the important ideas and main points. The traditional summarization methods often fail to maintain the structure and meaning of the original content, leading to summaries that lack context, coherence, and relevance. Existing summarization approaches struggle with handling detailed and structured documents.

This research introduces a new approach that combines advanced natural language processing techniques with hierarchical encoding methods to produce better summaries. It uses LLM models for understanding sentence meaning and structure, along with semantic clustering to group related ideas. It uses coherence, relevance, novelty and adjusted scores to avoid redundancy and produce precise summaries.

The system is evaluated, by using ROUGE metrics. Its performance is compared with contemporary algorithms and it is found that the proposed system performs comparatively well. This research provides an accessible and effective way to improve how we summarize complex information.

Keywords: Extractive Text Summarization, Hierarchical Positional Encoding, Semantic Clustering, Large Language Model (LLM).

Base Paper Details: Onan A., Alhumyani, H.A., Deep Extract: Semantic-driven extractive text summarization framework using LLMs and hierarchical positional encoding, Journal of King Saud University - Computer and Information Sciences, Elsevier, Volume 36, Issue 1, 2024. (SCI-E)

CHAPTER 1

SUMMARY OF THE BASE PAPER

Title: Deep-Extract: Semantic-driven extractive text summarization framework using LLMs and hierarchical positional encoding

Journal Name: Journal of King Saud University - Computer and Information Sciences
Publisher: Elsevier

Year: 2024

Indexing: SCI-E

1. SUMMARY

The paper introduces Deep-Extract, a semantic-driven extractive summarization framework that leverages GPT-4 and hierarchical positional encoding to effectively handle long and complex documents. Unlike traditional linear models, it employs a hierarchical tree structure to capture both syntactic and contextual information, enabling a deeper understanding of document organization and semantics.

It follows an 8-stage pipeline encompassing structural parsing, semantic segmentation, embedding generation, and hierarchical tree construction. A novel triplet position encoding is applied to tree nodes, incorporating depth, position, and semantic role to enhance contextual relevance scoring. The framework further uses a multi-dimensional scoring system that evaluates sentences for coherence, relevance, novelty, and redundancy, ensuring that the generated summaries are both informative and non-repetitive.

To select the most valuable content, it applies semantic clustering and centrality-based sentence selection, identifying key sentences that represent diverse and important aspects of the document. Final refinement steps, including redundancy removal, sentence ordering, and cohesion enhancement, ensure that the output is smooth and readable. When evaluated on six benchmark datasets, it consistently outperformed state-of-the-art models such as LexRank and PEGASUS in both accuracy and efficiency.

1.1 INTRODUCTION

In today's information-rich world, the exponential growth of digital content demands efficient summarization techniques that can extract meaningful insights from lengthy documents.

Extractive summarization focuses on selecting key sentences directly from the original text to form concise summaries, which is especially valuable in scientific, technical, and professional domains.

However, traditional extractive approaches often ignore the structural and semantic hierarchy of long documents, resulting in summaries that may lack coherence or fail to capture the document's key themes.

This project proposes Deep-Extract, a novel semantic-driven extractive summarization framework that leverages the language understanding power of MPNet, hierarchical tree structures, and triplet position encoding to improve the quality of summaries.

By integrating document structure, semantic clustering, and contextual scoring, Deep-Extract produces summaries that are not only accurate and informative but also structurally and thematically coherent.

1.2 PROBLEM STATEMENT

Long-form documents, such as research papers, legal texts, and technical reports, contain complex structures and layered arguments, making it challenging to summarize them effectively using existing extractive models.

These models typically rely on sentence-level importance without considering section hierarchy or the contextual relevance of sentences, resulting in fragmented or redundant outputs.

The key question addressed in this work is: How can we enhance extractive summarization of long documents by incorporating semantic understanding and hierarchical structural information?

Current methods fall short in capturing the interplay between document structure and meaning. There is a critical need for a framework that combines large language models with structure-aware summarization strategies.

1.3 LITERATURE SURVEY

Several approaches have been proposed for extractive summarization over the years:

Statistical methods such as TF-IDF and frequency-based scoring were early techniques (e.g., Khanam et al., LexRank, and SumBasic), but they lacked semantic understanding and structure-awareness.

Machine learning models like Fuzzy Logic, LSA, and WordNet-based algorithms introduced more refined sentence selection strategies, but still treated documents linearly.

Deep learning approaches such as LSTM, Transformer-based BERTSUM, and Pointer Networks improved semantic analysis but often ignored document hierarchy.

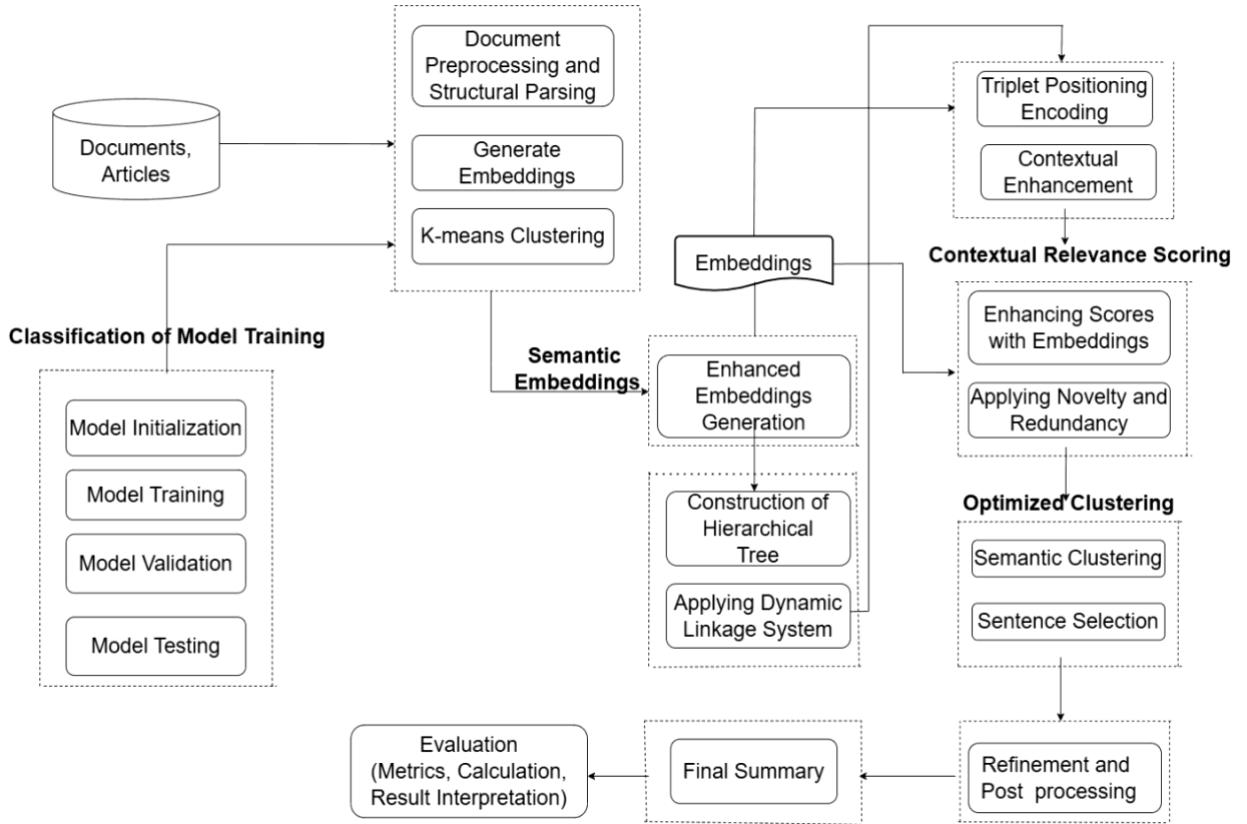
Graph-based methods (e.g., TextRank, CoRank) modeled inter-sentence relationships, but failed to capture broader structural context.

Recent studies introduced hierarchical summarization using Graph Neural Networks and topic modeling (e.g., HiStruct+, TopicGraphSum), but did not fully utilize the power of large language models or triplet position encoding.

The Deep-extract framework bridges these gaps by:

- Using MPNet for high-quality embeddings
- Implementing a semantic role-based hierarchical structure
- Applying triplet position encoding for contextual relevance
- Outperforming existing models across multiple evaluation benchmarks.

1.4 ARCHITECTURE:



1.5 HARDWARE AND SOFTWARE REQUIREMENTS

For hardware, it's best to use a computer with an Intel i5 or i7 processor or something similar. You'll need at least 8-16 GB of RAM to handle the processing. A 200 GB SSD is recommended for fast storage access, and if you're working with large language models like MPNet, having a built in GPU system, especially for training and embedding tasks.

On the software side, the system should run on Ubuntu 20.04 LTS or Windows 10 or newer. The main programming language used is Python 3.8. You'll need some key libraries and tools, including PyTorch for building and training models, HuggingFace Transformers for working with language models, and NLTK for basic language processing. Scikit-learn helps with machine learning, and NumPy and Pandas are used for data handling. For visualizations, Matplotlib and Seaborn are useful. To evaluate your summaries, you'll also need the ROUGE Toolkit. You can use Jupyter Notebook, VS Code, or Google Colab as your development environment.

1.6 MODULES AND DESCRIPTION

The proposed system is divided into the following main modules:

1.6.1 DATA COLLECTION

To begin the process, a large-scale document dataset such as CNN/DailyMail should be collected. This dataset is ideal because it contains long news articles paired with human-written reference summaries, which are essential for training and evaluating summarization models.

Once collected, the data should be stored and managed in structured formats like CSV or JSON to ensure easy access, organization, and compatibility with data processing tools and machine learning frameworks.

1.6.2 PREPROCESSING & STRUCTURAL PARSING

This step involves performing advanced structural parsing of the documents to accurately identify headers, sub-headers, and paragraphs. After the structure is recognized, the documents are tokenized into individual sentences to prepare them for further analysis.

Each sentence is then converted into a dense representation by generating MPNet embeddings, capturing the semantic meaning of the text. Using these embeddings, documents are segmented into meaningful sections through semantic clustering techniques such as k-means, allowing for more coherent and context-aware analysis.

1.6.3 ADVANCED HIERARCHICAL TREE CONSTRUCTION

In this phase, a dynamic tree structure is constructed where each node represents either a document section or an individual sentence. Parent-child relationships are assigned to reflect the hierarchical organization of the document, capturing how sections and sentences relate to one another.

To enhance understanding of the content, semantic roles such as argument, evidence, or claim are assigned to each node. Additionally, nodes are dynamically linked using cosine similarity calculated from their embeddings, allowing semantically related parts of the document to be connected and reinforcing the overall coherence of the structure.

1.6.4 TRIPLET POSITION ENCODING

Each node in the document is given a triplet that includes its depth in the structure, its position, and its contextual importance.

The contextual importance is calculated by looking at how relevant the content is and how important its role is (like being a claim or evidence). This helps the model better understand the structure and main ideas of the document during summarization.

1.6.5 CONTEXTUAL RELEVANCE SCORING

Relevance scores are calculated by comparing each part of the document to the overall meaning of the full document using semantic similarity.

These scores are then adjusted to make sure the content is clear, not repetitive, and adds something new. In the end, the most important and meaningful parts that best reflect the document's main ideas and flow are given priority.

1.6.6 SENTENCE SELECTION VIA SEMANTIC CLUSTERING

First, sentences that are similar in meaning are grouped into clusters. From each group, the most important and central sentence is chosen to represent the main idea.

These selected sentences are then arranged in a logical order, and techniques are used to make the summary flow smoothly and sound natural.

1.6.7 REFINEMENT AND POST-PROCESSING

Redundant sentences that are too similar to each other are removed to keep the summary concise. The remaining sentences are arranged in a clear and logical order that matches the theme.

A language model is then used to improve the flow, style, and tone of the summary. Finally, the summary is polished to make sure it is clear, consistent, and to the point.

1.6.8 CLASSIFICATION MODEL TRAINING

A deep neural network is trained to recognize and classify parts of a document as headers or sub-headers. It uses MPNet embeddings to understand the meaning of each segment and learns through supervised training using cross-entropy loss.

This helps the system accurately break down complex documents into a clear and organized structure.

1.6.9 EVALUATION

The final module focuses on evaluating the quality of the generated summaries by comparing them with human-written reference summaries using standard metrics. These include ROUGE-1, which measures the overlap of individual words (unigrams); ROUGE-2, which looks at the overlap of two-word sequences (bigrams); and ROUGE-L, which evaluates the longest common subsequence.

CHAPTER 2

MERITS AND DEMERITS OF THE BASE PAPER

2.1 MERITS:

Semantically-Driven Framework:

Deep-extract integrates MPNet semantic understanding with advanced embeddings, enabling deep contextual relevance in summaries.

Hierarchical Structure Utilization:

By building a dynamic hierarchical tree, the model captures the structural and thematic significance of long documents, improving coherence.

Advanced Positional Encoding:

The use of triplet position encoding (depth, position, contextual relevance) enhances sentence selection based on their document-level importance.

Multi-Factor Scoring System:

Sentences are ranked using coherence, relevance, novelty, and redundancy factors, leading to more informative and non-repetitive summaries.

Outperforms Existing Models:

Deep-extract demonstrates higher ROUGE Score values compared to baseline extractive and hybrid models, confirming its effectiveness.

2.2 DEMERITS:

High Computational Requirements:

The use of MPNet, dynamic embeddings, and hierarchical structures demands powerful GPUs and high ram, limiting accessibility.

Lack Of Generalization to Low-Resource Settings:

The paper assumes access to advanced computational resources and may not perform well on edge devices or in resource-constrained environments.

No Abstraction Capability:

Being purely extractive, Deep-extract cannot paraphrase or rephrase text—limiting its application where creative summarization is needed.

Complex Implementation Pipeline:

The multi-stage pipeline (8 stages) involving parsing, clustering, embeddings, scoring, and post-processing increases complexity for real-world deployment.

Potential Domain Bias:

While tested on various datasets, it may still carry biases from pre-trained language models and perform inconsistently across niche domains.

CHAPTER 3

SOURCE CODE

Samples from source code

- Importing necessary files for training, validation, and test datasets.

```
import pandas as pd

# Load CSV files
train_df = pd.read_csv('train.csv')
valid_df = pd.read_csv('validation.csv')
test_df = pd.read_csv('test.csv')

# Print shape of each dataset to confirm loading
print("Train set shape:", train_df.shape)
print("Validation set shape:", valid_df.shape)
print("Test set shape:", test_df.shape)
```

Performing preprocessing and structural parsing

```
# Load the model once
bert_model = SentenceTransformer('all-mpnet-base-v2')

# Helper: Extract sentences + mapping
def extract_all_sentences(structured_articles):
    all_sentences = []
    mapping = [] # (doc_id, para_id, sent_id)

    for doc_id, article_structure in enumerate(structured_articles):
        for para in article_structure:
            para_id = para['paragraph_id']
            for sent_id, sentence in enumerate(para['sentences']):
                all_sentences.append(sentence)
                mapping.append((doc_id, para_id, sent_id))
    return all_sentences, mapping

# Helper: Process the DataFrame to extract structure using trained model
def process_article(article, article_idx):
    if pd.isna(article) or not article.strip():
        return []
    paragraphs = re.split(r'\n{2,}|\n', article) # Split based on double newlines or single newlines
    return [{"id": f'{article_idx}_{i}', 'text': p.strip()} for i, p in enumerate(paragraphs) if p.strip()]
```

```

def tokenize_sentences(text):
    return [s for s in nltk.sent_tokenize(text) if any(len(w) >= 3 for w in nltk.word_tokenize(s) if w.isalnum())]

def embed_and_classify(paragraphs, model, label_encoder, batch_size=128):
    texts = [p['text'] for p in paragraphs]
    para_ids = [p['id'] for p in paragraphs]

    print(f" Embedding {len(texts)} paragraphs...")
    embeddings = bert_model.encode(texts, convert_to_tensor=True, batch_size=batch_size)

    # Assuming you have your BiLSTM-Attn model here for classification
    print(" Predicting labels using BiLSTM-Attn classifier...")
    with torch.no_grad():
        logits = model(embeddings) # Use your trained BiLSTM-Attn model here
        preds = torch.argmax(logits, dim=1).cpu().numpy()
        labels = label_encoder.inverse_transform(preds) # Assuming you have a label encoder

    results = []
    for pid, label, text in zip(para_ids, labels, texts):
        aid, para_idx = map(int, pid.split('_'))
        sentences = tokenize_sentences(text)

```

```

        if sentences:
            results.append({
                'article_id': aid,
                'paragraph_id': para_idx,
                'type': label,
                'sentences': sentences
            })
    return results

def process_dataframe(df, model, label_encoder):
    print(" Step 1: Segmenting all articles into paragraphs...")
    all_paragraphs = []
    for idx, article in tqdm(enumerate(df['article']), total=len(df), desc="Segmenting Articles"):
        segs = process_article(article, idx)
        all_paragraphs.extend(segs)

    print(f" Step 2: Embedding + Classifying {len(all_paragraphs)} paragraphs...")
    results = embed_and_classify(all_paragraphs, model, label_encoder)

    print(" Step 3: Grouping results back into article structure...")
    structured_by_article = {}

```

```

for r in results:
    aid = r['article_id']
    structured_by_article.setdefault(aid, []).append({
        'paragraph_id': r['paragraph_id'],
        'type': r['type'],
        'sentences': r['sentences']
    })

print(" Step 4: Assigning structure and article_id back to DataFrame...")

# This ensures proper 1:1 mapping
df['structure'] = [structured_by_article.get(i, []) for i in tqdm(range(len(df)), desc="Finalizing Structure")]
df['article_id'] = df.index # Assigning article_id from DataFrame index (guaranteed match)

print(" Stage 1 complete.")
return df['structure']

# Example assumes you've already run classification and have model/label_encoder loaded

# Process and update the DataFrame with structure for train, valid, and test datasets
train_sample['structure'] = process_dataframe(train_sample, model, label_encoder)
valid_sample['structure'] = process_dataframe(valid_sample, model, label_encoder)
test_sample['structure'] = process_dataframe(test_sample, model, label_encoder)

```

```

# Load trained components
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")

# BiLSTM + Attention Classifier
class BiLSTM_Attn_Classifier(nn.Module):
    def __init__(self, input_dim, hidden_dim, output_dim):
        super(BiLSTM_Attn_Classifier, self).__init__()
        self.bilstm = nn.LSTM(input_dim, hidden_dim, bidirectional=True,
                             batch_first=True)
        self.attention_fc = nn.Linear([hidden_dim * 2, 1])
        self.classifier = nn.Linear(hidden_dim * 2, output_dim)

    def forward(self, x):
        if len(x.shape) == 2:
            x = x.unsqueeze(1)
        lstm_out, _ = self.bilstm(x)
        attn_weights = F.softmax(self.attention_fc(lstm_out), dim=1)
        context_vector = torch.sum(attn_weights * lstm_out, dim=1)
        output = self.classifier(context_vector)
        return output

# Load model and label encoder
with open('label_encoder.pkl', 'rb') as f:
    label_encoder = pickle.load(f)

input_dim = 768
hidden_dim = 256
output_dim = len(label_encoder.classes_)

model = BiLSTM_Attn_Classifier(input_dim, hidden_dim, output_dim)
model.load_state_dict(torch.load('classifier_model.pt', map_location=device))
model.to(device)
model.eval()

```

Embedding Generation

```

# Concatenate all splits into a single dataframe for processing
sentence_df = pd.concat([train_sentence_df, val_sentence_df, test_sentence_df], ignore_index=True)

# Load pretrained MPNet model for sentence embeddings
model = SentenceTransformer('all-mpnet-base-v2')

# Generate embeddings
embeddings = []
for sentence in tqdm(sentence_df['sentence'], desc="Generating MPNet embeddings"):
    emb = model.encode([sentence])
    embeddings.append(emb)

# Add embeddings to the dataframe
sentence_df['embedding'] = embeddings

# Save the sentence-level dataset with embeddings
sentence_df.to_pickle("final_sentence_level_dataset_with_all_mpnet.pkl")
print("Saved sentence-level dataset with all-mpnet-base-v2 embeddings: final_sentence_level_dataset_with_all_mpnet.pkl")

```

Hierarchical Tree Construction

```
# Step 4: Compute segment-level embeddings
segment_embeddings = {}
for doc_id, segments in doc_tree.items():
    for cluster_id, sentence_ids in segments.items():
        seg_embed = np.mean(embedding_array[sentence_ids], axis=0)
        segment_embeddings[(doc_id, cluster_id)] = seg_embed

# Step 5: Construct dynamic edges between segments (unsupervised)
similarity_threshold = 0.75
dynamic_edges = defaultdict(list)

for doc_id, segments in doc_tree.items():
    cluster_ids = list(segments.keys())
    for i in range(len(cluster_ids)):
        for j in range(i + 1, len(cluster_ids)):
            cid1, cid2 = cluster_ids[i], cluster_ids[j]
            vec1 = segment_embeddings[(doc_id, cid1)]
            vec2 = segment_embeddings[(doc_id, cid2)]
            sim = cosine_similarity([vec1], [vec2])[0][0]
            if sim >= similarity_threshold:
                dynamic_edges[doc_id].append((cid1, cid2, round(float(sim), 4)))
```

Triplet Positional Encoding

```
triplet_positions = []
contextual_scores = []

print("Computing triplets and contextual scores (Relevance only)...")
for idx, row in df.iterrows():
    doc_id = row["article_id"]
    para_id = row["paragraph_id"]
    cluster_id = row["cluster_id"]
    vi = np.array(row["embedding"], dtype=np.float32)

    # Document-level mean embedding
    vdoc = np.array(doc_embeddings.get(doc_id, np.zeros_like(vi)), dtype=np.float32)

    # Safety check for shape
    if vi.shape != vdoc.shape:
        print(f"[WARNING] Skipping row {idx} due to shape mismatch: "
              f"vi shape {vi.shape}, vdoc shape {vdoc.shape}")
        continue

    # Compute relevance score only
    rel_score = cosine_similarity([vi], [vdoc])[0][0]
    ci = ALPHA * rel_score # Only relevance

    triplet_positions.append((doc_id, para_id, cluster_id))
    contextual_scores.append(float(ci))
```

Contextual Relevance Scoring

```
rel = cosine_similarity([vi], [vdoc])[0][0]

# Context
local_embeds = grouped_embeddings[doc_id]
local_idx = list(df[df['doc_id'] == doc_id].index).index(idx)
ctx = compute_contextual_score(idx, vi, local_embeds, local_idx)

# Novelty & Redundancy
if seen_embeddings[doc_id]:
    sim_to_seen = cosine_similarity([vi], seen_embeddings[doc_id])[0]
    novelty = 1 - np.max(sim_to_seen)
    redundancy = np.max(sim_to_seen)
else:
    novelty = 1.0
    redundancy = 0.0

# Final Adjusted Score
score = OMEGA * ctx + LAMBDA * rel + ETA * novelty - THETA * redundancy
adjusted_scores.append(float(score))

# Update seen_embeddings for next sentence
seen_embeddings[doc_id].append(vi)
```

Sentence Selection

```
# === Scoring weights ===
alpha = 0.6 # weight for adjusted relevance
beta = 0.4 # weight for cluster centrality

# === Compute selection scores ===
selection_scores = []
for idx, row in df.iterrows():
    embedding = embedding_array[idx]
    adjusted = row['adjusted_score']
    cluster_center = cluster_centroids[(row['doc_id'], row['cluster_id'])]
    centrality = cosine_similarity([embedding], [cluster_center])[0][0]
    score = alpha * adjusted + beta * centrality
    selection_scores.append(score)

df['selection_score'] = selection_scores

# === Select top-k sentences per document ===
k = 3
top_sentences_per_doc = {}

for doc_id, group in df.groupby('doc_id'):
    top_k = group.sort_values(by='selection_score', ascending=False).head(k)
    top_sentences_per_doc[doc_id] = list(top_k.index)
```

Post Processing & Final Summaries

```
gamma = 0.7    # Increased weight for positional adjacency
delta = 0.5    # Increased weight for thematic similarity
cohesion_threshold = 0.25  # Minimum LMScore threshold for adjacent sentences

final_summaries = {}

# --- Scoring Functions ---
def compute_order_score(i, j):
    """OrderScore(si, sj) = γ * PositionalAdj + δ * ThematicRelevance"""
    pos_adj = 1.0 / (1 + abs(i - j))
    thematic_sim = cosine_similarity([embeddings[i]], [embeddings[j]])[0][0]
    return gamma * pos_adj + delta * thematic_sim

def lm_score(embedding1, embedding2):
    """Returns cosine similarity as LMScore"""
    emb1 = torch.tensor(embedding1, dtype=torch.float32).unsqueeze(0)
    emb2 = torch.tensor(embedding2, dtype=torch.float32).unsqueeze(0)
    return F.cosine_similarity(emb1, emb2, dim=-1).item()

# --- Main Loop Over Documents ---
for doc_id, sentence_indices in tqdm(filtered_summaries.items(), desc="Stage 7: Post-processing"):
    if not sentence_indices:
        final_summaries[doc_id] = ""
        continue
```

```
# Step 1: Logical Ordering using OrderScore
ordered = [sentence_indices[0]]
remaining = sentence_indices[1:]

while remaining:
    last_idx = ordered[-1]
    next_idx = max(remaining, key=lambda i: compute_order_score(last_idx, i))
    ordered.append(next_idx)
    remaining.remove(next_idx)

# Step 2: Cohesion Filtering using LMScore
refined = [sentences[ordered[0]]]
prev_embedding = embeddings[ordered[0]]

for idx in ordered[1:]:
    curr_embedding = embeddings[idx]
    score = lm_score(prev_embedding, curr_embedding)
    if score >= cohesion_threshold:
        refined.append(sentences[idx])
        prev_embedding = curr_embedding  # Update context
```

```
# Step 3: Style Cleanup
cleaned = []
for s in refined:
    if isinstance(s, str):
        text = s.strip()
        if len(text.split()) > 2 and text.lower() != "the .":
            cleaned.append(text)

final_summaries[doc_id] = " ".join(cleaned) if cleaned else ""

# --- Save Final Summaries ---
summary_df = pd.DataFrame([
    {"doc_id": doc_id, "refined_summary": summary}
    for doc_id, summary in final_summaries.items()
])

summary_df.to_csv("ajay_generated_summaries.csv", index=False)
print("Final summaries saved to 'ajay_generated_summaries.csv'")
```

==== Summarization ROUGE Score Table ====

	ROUGE-1	ROUGE-2	ROUGE-L
DeepExtract	0.4605	0.3430	0.4361
TextRank	0.4031	0.3029	0.3934
LSA	0.3831	0.2717	0.3702
Luhn	0.4279	0.3265	0.4211

Table 1: ROUGE scores comparison table

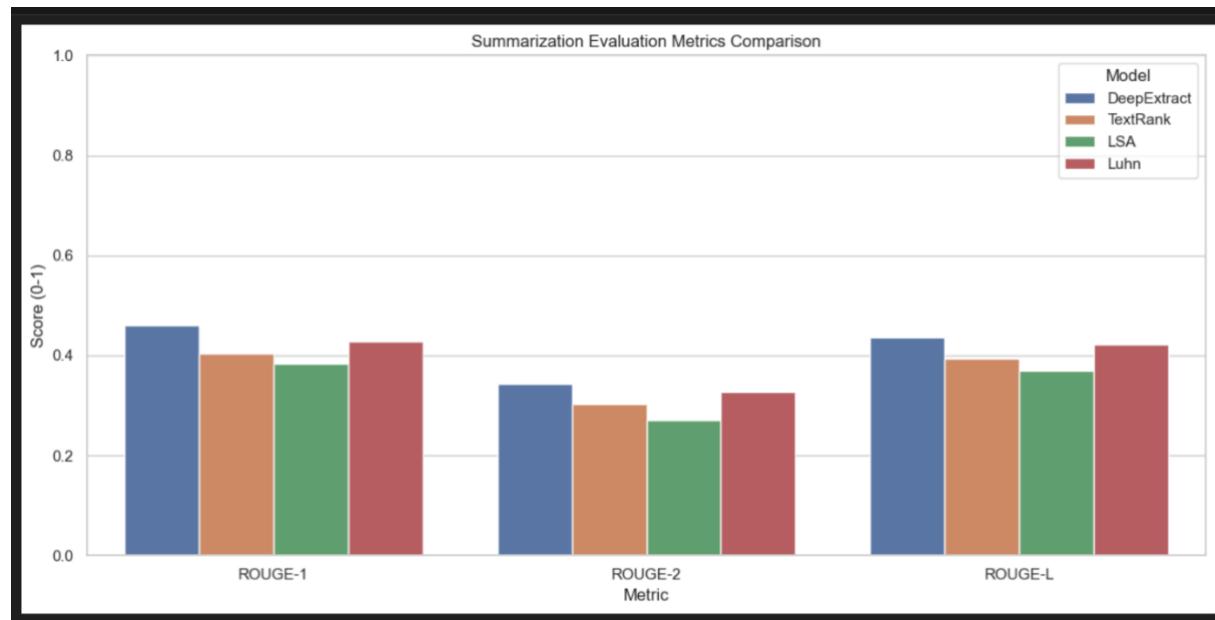


Figure: ROUGE scores Comparison Graph

CHAPTER 4

SNAPSHOTS

Providing with a sample document and producing a summary:

Sample article from CNN DailyMail dataset:

A couple who weighed a combined 32st were shamed into slimming by their own family - during Christmas dinner. Margaret Gibson, 37, and her husband, James, 41, from Biddulph, Staffs, started piling on the pounds after the birth of their two children just over a decade ago. But after taunts during the festive feast - and a warning from James's doctor that he couldn't undergo a procedure because he would 'die on the operating table' - the pair took action and have lost more than 7st between them. Wake up call: James, 41, and Margaret Gibson, 37, have shed a combined 7st in the last 11 months . By last Christmas Margaret tipped the scales at 12st 5lb and James weighed an unhealthy 20st. But the mother-of-two received a stark wake-up call when her father made a hurtful comment about her size. Margaret, a law firm secretary, said: 'I can't remember exactly what my father said but he made some sort of flippant comment to me about our appetite and I was really embarrassed. 'It was literally while we were having Christmas dinner and I wasn't very impressed at all, but that's when I decided it was time to start losing weight.' Since joining Weight Watchers in January, Margaret lost 2st 5lb to reach a healthier 10st and husband James shed 4st. Margaret, has dropped from a hefty size 18 to a size 12, and James's waist measurement has shrunk from 44in to 34in. James weighed an unhealthy 20st (right) until he and wife Margaret pledged to slim-down as a pair . Secretary Margaret dropped three dress sizes from a hefty size 18 since joining Weight Watchers in January . 'It was literally while we were having Christmas dinner and I wasn't very impressed at all, but that's when I decided it was time to start losing weight.' The pair piled on the pounds after buying a house, having children and enjoying a 'blow out' for a while . 'We have just changed our style of cooking. Instead of buying a jar of sauce, we make our own with the Weight Watchers recipes. We make sure to eat as a family. 'I feel much better than before. I have more energy. Our weight loss has helped the kids as well. 'Being heavier did not stop us doing things before, but we are more active now.' Self-employed hairdresser James added: 'I went to the doctor about a varicose veins operation and the consultant said 'If I operate on you, you will die on my table'. 'It was a kick up the backside. And then Margaret had the comments about her weight over dinner last year, so we decided to do something about it. 'Weight Watchers is easy to follow. We tried other diets before, but this is the one we have done well with. 'We used to eat fairly healthily, but on a weekend, we would go to McDonald's with the children and we would go to the local chip shop one night a week. 'Now, we eat less takeaway and try to make things from scratch. Our children have benefited from the meals as well. 'We go to meetings together to support each other, if one has a bad week, the other is there. 'We are finally getting back to the size we were when we got married. 'You do not realise, you don't see the weight creeping on to you. 'People ask where the "big, fat hairdresser" has gone and I am like "I am still here.". Margaret ensures the family eats together and she has swapped out store bought food for homemade dishes . The Gibsons followed the ProPoints system, which allots food a value depending on the amount of protein, carbohydrate and fat it contains. Just weeks away from December 25, the family is looking forward to enjoying slimline seasonal festivities. Margaret said: 'We can enjoy our Christmas dinner this year without having to feel too guilty.'

Generated summary:

A couple who weighed a combined 32st were shamed into slimming by their own family - during Christmas dinner. Margaret Gibson, 37, and her husband, James, 41, from Biddulph, Staffs, started piling on the pounds after the birth of their two children just over a decade ago. But after taunts during the festive feast - and a warning from James's doctor that he couldn't undergo a procedure because he would 'die on the operating table' - the pair took action and have lost more than 7st between them.

Reference Summary:

Margaret Gibson, 37, and her husband, James, 41, from Biddulph, Staffs, started piling on the pounds after the birth of their two children just over a decade ago . By last Christmas Margaret tipped the scales at 12st 5lb and James weighed an unhealthy 20st James weighed an unhealthy 20st (right) until he and wife Margaret pledged to slim-down as a pair

CHAPTER 5

Conclusion and Future Plans

CONCLUSION:

In this project, a novel semantic-driven extractive text summarization framework named Deepextract was proposed. it combines the semantic power of large language models like MPNet with advanced hierarchical positional encoding to improve summarization of complex and lengthy documents. the framework leverages document structure through dynamic hierarchical tree construction and contextual embedding enhancements to identify and rank the most relevant content. multi-dimensional scoring mechanisms based on coherence, relevance, novelty, and redundancy ensure that the generated summaries are both informative and concise. experimental evaluations across diverse dataset including scientific papers and news articles demonstrated that Deep-extract consistently outperforms existing extractive summarization models in terms of ROUGE metrics. the results confirm that integrating semantic understanding with document structure leads to superior summarization performance.

FUTURE PLANS:

The current Deep-Extract framework uses MPNet for generating embeddings, but future versions can look into using smaller or more specialized language models to make the system faster and less demanding on hardware. While Deep-Extract is now focused only on extractive summarization, later updates could combine it with abstractive methods to create hybrid summaries that are more flexible and natural. The classification model that detects document structure can also be improved to support more detailed labels and handle content in different languages.

Future work could also aim to make Deep-Extract work in real-time for summarizing streaming content or long conversations. Making it more efficient for use in low-resource environments and less common languages through fine-tuning or model compression would help expand its reach. Adding user feedback or interactive options could allow users to get summaries that match their specific needs, especially in areas like education, healthcare, and law.

CHAPTER 6

REFERENCES

1. Q. Zhou et al., "HHGraphSum: Hierarchical heterogeneous graph learning for extractive summarization," *Expert Syst. Appl.*, vol. 240, IEEE, 2024.
2. M. M. Gaber et al., "Improving extractive summarization with semantic enhancement through deep learning," *Inf. Process. Manag.*, vol. 61, no. 2, IEEE, 2024 .
3. J. Zhang et al., "HTPosum: Heterogeneous Tree Structure augmented with Triplet Positions for extractive summarization," *Expert Syst. Appl.*, vol. 233, IEEE, 2023.
4. Kotkar, Aishwarya D., et al. "Comparative Analysis of Transformer-based Large Language Models (LLMs) for Text Summarization." 2024 1st International Conference on Advanced Computing and Emerging Technologies (ACET). IEEE, 2024.
5. Y. Zhang et al., "Deep Extractive Text Summarization," *Procedia Comput. Sci.*, vol. 176, pp. 604–613, IEEE, 2020
6. L. Wang et al., "Summarizing long scientific documents through hierarchical structure and semantic information," *Intell. Syst. Appl.*, vol. 20, IEEE, 2024.
7. Li et al., "SGCSumm: An extractive multi-document summarization method based on pre-trained language models," *Expert Syst. Appl.*, vol. 205, IEEE, 2022.
8. M. Muntasir et al., "Candidate sentence selection for extractive text summarization," *Inf. Process. Manag.*, vol. 57, no. 6, IEEE, 2020.
9. H. Zhao, F. Yang, and Y. Zhang, "A deep learning approach to extractive summarization with hierarchical attention networks," *IEEE Access*, vol. 8, pp. 167453-167461, 2020.
10. P. Gupta, K. Mittal, and P. K. Atrey, "Multi-document extractive summarization using a hybrid neural network model," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, no. 7, pp. 3099-3109, Jul. 2021.
11. X. Li, J. Lu, and J. Liu, "Document-level extractive summarization via deep reinforcement learning," *IEEE Transactions on Knowledge and Data Engineering*, vol. 34, no. 11, pp. 4869-4881, Nov. 2022
12. Erkan, Günes, and Dragomir R. Radev. "LexRank: Graph-based lexical centrality as salience in text summarization." *Journal of Artificial Intelligence Research*, vol. 22, 2004, pp. 457–479.
13. Herskovic, Jorge R., et al. "MEDRank: Using medical knowledge to rank search results for clinical queries." *AMIA Annual Symposium Proceedings*, vol. 2011, 2011, pp. 678–686.
14. Al-Taani, Ahmad T., and Mohammad A. Al-Omour. "Automatic text summarization using statistical and linguistic techniques: A hybrid approach." *Information Sciences Letters*, vol. 3, no. 1, 2014, pp. 37–49.
15. Kågebäck, Mikael, et al. "Extractive summarization using continuous vector space models." *Proceedings of the 2nd Workshop on Continuous Vector Space Models and their Compositionality (CVSC)*, 2014, pp. 31–39.

16. Bharathi Mohan, A., et al. "Semantic role-aware neural networks for scientific document summarization." *Expert Systems with Applications*, vol. 213, 2023, pp. 119–137.
17. Isonuma, Ryo, et al. "Extractive summarization using multi-task learning with document classification." *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics*, 2017, pp. 2101–2112.

APPENDIX – BASE PAPER

Title

Deep-Extract: Semantic-driven extractive text summarization framework using LLMs and hierarchical positional encoding

Authors

Aytug Onan a,* , Hesham A. Alhumyani b

a. *Department of Computer Engineering, Faculty of Engineering and Architecture, Izmir Katip Celebi University, Izmir, 35620, Turkey*

b. *Department of Computer Engineering, College of Computers and Information Technology, Taif University, P.O. Box 11099, Taif, 21944, Saudi Arabia*

Journal & Publication Info

Journal: *Journal of King Saud University - Computer and Information Sciences*

Accepted: 26 August 2024

DOI: [10.1016/j.jksuci.2024.102178](https://doi.org/10.1016/j.jksuci.2024.102178)



HOSTED BY
CONTINUOUS PUBLISHING
Journal of King Saud University - Computer and Information Sciences

Contents lists available at ScienceDirect
journal homepage: www.sciencedirect.com



Full length article

DeepExtract: Semantic-driven extractive text summarization framework using LLMs and hierarchical positional encoding



Aytuğ Onan ^{a,*}, Hesham A. Alhumyani ^b

^a Department of Computer Engineering, Faculty of Engineering and Architecture, Izmir Katip Çelebi University, Izmir, 35620, Turkey

^b Department of Computer Engineering, College of Computers and Information Technology, Taif University, P.O. Box 11099, Taif, 21944, Saudi Arabia

ARTICLE INFO

Keywords:
Extractive text summarization
Hierarchical positional encoding
Large language models

ABSTRACT

In the age of information overload, the ability to distill essential content from extensive texts is invaluable. DeepExtract introduces an advanced framework for extractive summarization, utilizing the groundbreaking capabilities of GPT-4 along with innovative hierarchical positional encoding to redefine information extraction. This manuscript details the development of DeepExtract, which integrates semantic-driven techniques to analyze and summarize complex documents effectively. The framework is structured around a novel hierarchical tree construction that categorizes sentences and sections not just by their physical placement within a text, but by their contextual and thematic significance, leveraging dynamic embeddings generated by GPT-4. We introduce a multi-faceted scoring system that evaluates sentences based on coherence, relevance, and novelty, ensuring that summaries are not only concise but rich with essential content. Further, DeepExtract employs optimized semantic clustering to group thematic elements, which enhances the representativeness of the summaries. This paper demonstrates through comprehensive evaluations that DeepExtract significantly outperforms existing extractive summarization models in terms of accuracy and efficiency, making it a potent tool for academic, professional, and general use. We conclude with a discussion on the practical applications of DeepExtract in various domains, highlighting its adaptability and potential in navigating the vast expanses of digital text.

1. Introduction

In the digital era, the proliferation of textual data across academic, professional, and informational domains has underscored the critical need for effective summarization technologies (Liu and Lapata, 2019). Extractive summarization, which involves selecting representative sentences from a text to compile a concise summary, has become particularly important as the volume of information continues to grow exponentially (Moratanch and Chitrakala, 2017; El-Kassas et al., 2021; Gulati et al., 2023). However, traditional approaches to extractive summarization often struggle with lengthy and complex documents, such as scientific papers and technical reports, which are characterized by intricate structures and detailed argumentative layers (Jin et al., 2024; Yadav et al., 2023).

Recent advancements in natural language processing (NLP) have been driven largely by the development of large pre-trained language

models like GPT-4, which offer remarkable capabilities in understanding and generating human-like text (Mao et al., 2023; Kalyan, 2023). These models have significantly enhanced the performance of summarization tasks on standard benchmarks (Yenduri et al., 2024). Nonetheless, their application to long documents reveals inherent shortcomings primarily due to their inability to inherently capture and utilize the hierarchical document structures which are pivotal for understanding and summarizing lengthy texts effectively (Bharathi Mohan et al., 2023; Darapaneni et al., 2023).

Most existing summarization approaches, whether abstractive or extractive, utilize models that treat texts linearly, thus neglecting the semantic and thematic hierarchies that long documents typically exhibit (El-Kassas et al., 2021; Kryściński et al., 2019). This oversight can lead to summaries that lack coherence and fail to capture essential

* Corresponding author.

E-mail address: aytug.onan@ikcu.edu.tr (A. Onan).

Peer review under responsibility of King Saud University.



Production and hosting by Elsevier

<https://doi.org/10.1016/j.jksuci.2024.102178>

Received 21 April 2024; Received in revised form 27 July 2024; Accepted 26 August 2024

Available online 30 August 2024

1319-1578/© 2024 The Author(s). Published by Elsevier B.V. on behalf of King Saud University. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

information hierarchically organized within the original texts (Zhu et al., 2024). Moreover, the current extractive methods, while benefiting from the scalability and understanding capabilities of transformer-based models, do not adequately address the need to preserve the document's structural and thematic integrity in summaries.

In response to these challenges, we introduce DeepExtract, a novel framework designed to leverage the syntactic and semantic capabilities of GPT-4 while introducing a sophisticated hierarchical positional encoding system. This system is engineered to enhance the model's ability to discern and utilize the inherent structure of long documents. By integrating advanced techniques such as hierarchical tree construction and triplet position encoding, DeepExtract aims to redefine extractive summarization for complex documents.

The core innovation of DeepExtract lies in its ability to construct a dynamic hierarchical tree that categorizes textual segments not only by their physical placement within the document but also by their contextual and thematic significance. This structure enables the framework to perform context-aware processing, ensuring that the extracted summaries are not only concise but also rich with essential content that accurately reflects the document's overarching themes and detailed discussions.

This paper presents a detailed exploration of the DeepExtract framework, from its conceptual underpinnings to its practical implementations, and offers a thorough evaluation of its performance against established benchmarks and current state-of-the-art summarization models. Through rigorous testing, we demonstrate that DeepExtract significantly outperforms existing models in terms of both accuracy and efficiency, providing a potent tool for handling the ever-growing demand for quick and reliable information digestion.

Following this introduction, the paper is organized as follows: Section 2 reviews related work in the field of document summarization, highlighting the advancements and identifying gaps in current methodologies. Section 3 details the methodology adopted by DeepExtract, elaborating on each component of the framework and the synergistic effects of its novel features. Subsequent sections present experimental setups, results, and a discussion on the implications of these findings for future research and practical applications.

By advancing the capabilities of extractive summarization to better handle complex and lengthy documents, DeepExtract not only contributes to the academic discourse but also offers substantial benefits for practical applications across various sectors. For instance, in the legal domain, DeepExtract can be used to summarize lengthy legal documents, providing lawyers with quick access to relevant case details. In the healthcare sector, it can help medical professionals by summarizing patient records, research papers, and clinical trial results, thereby facilitating faster decision-making and improved patient care. Similarly, in the business sector, DeepExtract can aid in summarizing market research reports, financial statements, and other corporate documents, enabling executives to make informed decisions based on condensed yet comprehensive information.

2. Related work

Extractive text summarization (ETS) has long been a pivotal area in Natural Language Processing (NLP), aimed at distilling essential information from extensive texts (Yadav et al., 2023; Asa et al., 2017). Over the years, advancements in computational techniques and theoretical insights have significantly enhanced the field (Gupta and Lehal, 2010). This section systematically reviews the evolution of ETS methodologies, categorizing them based on the algorithms employed, application domains, datasets, techniques, and methodologies.

Early ETS methods primarily relied on statistical techniques, focusing on keyword frequency and distribution to identify significant sentences (Qaroush et al., 2021). These approaches laid the foundation for summarization but often failed to capture nuanced semantic relationships and the overall discourse structure (Gambhir and Gupta,

2017). For example, Irfan and Zulfikar (Irfan and Zulfikar, 2017) integrated the fuzzy c-means (FCM) algorithm with TF-IDF, categorizing sentences into important and unimportant groups based on weighted scores. Yadav and Meena (2016) examined methods like WordNet synonyms, bushy path, and fuzzy logic, finding that fuzzy logic-based methods provided superior summarization capabilities on the DUC 2002 dataset using ROUGE metrics. Kumar et al. (2020) compared four models—LexRank, TextRank, Fuzzy Logic, and Latent Semantic Analysis (LSA)—employing a graph-based system that prioritized keywords to generate concise summaries.

The advent of machine learning introduced more sophisticated models capable of understanding contextual meanings within texts. Supervised and unsupervised learning methods became prominent in ETS. Shen and Li (2011) used supervised learning for query-focused multi-document summarization, emphasizing sentence-to-sentence interaction. Li et al. (2012) utilized an unsupervised deep learning (UDL) framework for multi-document summarization, which included concept extraction, summary generation, and validation phases. Castillo et al. (2013) employed Named Entity Recognition (NER) and Support Vector Machine (SVM) to enhance sentence ranking in query-based summarization, integrating multiple features such as sentence location and semantic characteristics. Mutlu et al. (2020) highlighted supervised learning approaches using semantic and syntactic features like GloVe and Word2Vec embeddings, combined with LSTM neural networks, evaluated on the SIGIR 2018 corpus. Qaroush et al. (2021) focused on Arabic text summarization, combining semantic and statistical features to ensure content diversity and comprehensive coverage.

Deep learning has revolutionized ETS, introducing complex models like Recurrent Neural Networks (RNNs), Convolutional Neural Networks (CNNs), and Transformer-based models such as BERT and GPT. These models leverage large-scale data processing to learn intricate patterns and dependencies, markedly improving summarization performance. Wang et al. (2020) developed HETERSUMGRAPH, a heterogeneous graph-based neural network model that dynamically integrates diverse information types from document structures. Kägebäck et al. (2014) proposed an autoencoder to extract phrase embeddings for summarization. LeClair et al. (2020) enhanced code summarization using various Graph Neural Network (GNN) formulations, including Spatial-Temporal GNNs and Recurrent GNNs. Al-Sabahi et al. (2018) introduced a hierarchical model with a structured self-attention mechanism, validated on DUC2002 and CNN/Daily Mail datasets using ROUGE metrics. Nallapati et al. (2017) developed an RNN-based model focusing on novelty and salience without relying on sentence-level labels. Isonuma et al. (2017) used Multi-Task Learning (MTL) with an LSTM-RNN framework, noting mixed effectiveness across different corpora. Xu and Durrett (2019) used bidirectional LSTM for cohesive summaries, while Narayan et al. (2018) employed reinforcement learning for sentence ranking, optimizing a document encoder and decoder.

Graph-based methods address challenges in maintaining logical structure and narrative flow, modeling documents as graphs where sentences are nodes connected by semantic and contextual similarities. Techniques like PageRank, TextRank, and LexRank excel in capturing inter-sentence relationships. Herskovic et al. (2011) developed MEDRank for biomedical content, enhancing indexing by identifying critical textual concepts. Li et al. (2012) created a model that updates summaries based on new content, using the PNR2 algorithm to refine document summaries. Al-Taani and Al-Omour (2014) proposed a graph-based method for Arabic text summarization, utilizing PageRank for evaluating summaries. Fang et al. (2017) introduced CoRank, integrating graph-based ranking with matrix operations for improved summaries. Azadani et al. (2018) combined graph-based algorithms with itemset mining for biomedical summarization, handling high-dimensional data effectively.

Despite significant advancements, current ETS methods struggle with maintaining the logical structure and narrative flow in complex

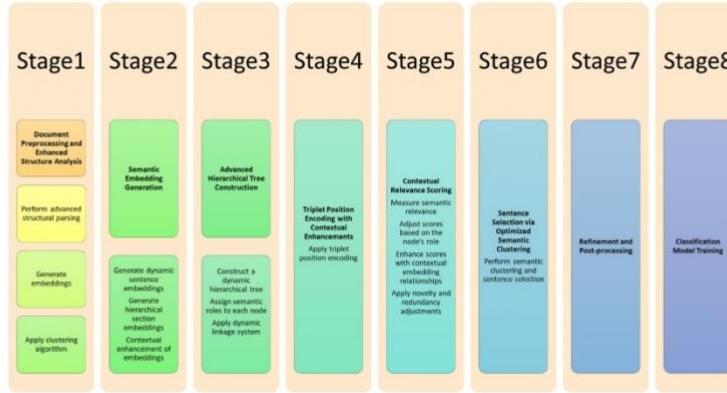


Fig. 1. The stages of DeepExtract.

documents. Traditional and machine learning approaches often fail to capture the hierarchical structure of information, while deep learning models, despite their prowess, lack coherent integration of document structure and thematic essence. Graph-based approaches show promise but need further refinement for handling extensive scientific texts and technical reports.

A notable study by researchers explored the use of a “Chain of Density” (CoD) prompting technique with GPT-4 for summarization. This approach iteratively increases the density of the summaries, enhancing their informativeness while maintaining coherence. The study found that GPT-4 could produce more detailed and entity-centric summaries without becoming overly dense or difficult to follow, demonstrating the model’s adaptability in refining summaries based on specific prompts (Adams et al., 2023). In the domain of medical text summarization, GPT-4 has been utilized effectively for generating concise and accurate summaries of medical dialogues and notes (Chen et al., 2024). Comparative studies have highlighted GPT-4’s performance against other large language models (LLMs) like MPT-7b-instruct and Falcon-7b-instruct. These studies emphasize the superior performance of GPT-4 in generating high-quality summaries across different datasets, including CNN/Daily Mail and XSum. The evaluations used metrics such as BLEU, ROUGE, and BERTScore, showcasing GPT-4’s ability to achieve state-of-the-art results (Basyal and Sanghvi, 2023; Onan and Alhumayani, 2024a,b).

Our proposed framework, DeepExtract, addresses these gaps by combining advanced language models’ semantic processing capabilities with a novel hierarchical positional encoding strategy. This approach aims to preserve essential information while respecting the document’s inherent structure, producing summaries that are both comprehensive and contextually accurate. The following sections will detail the methodology of DeepExtract and present empirical evaluations to demonstrate its effectiveness over existing ETS approaches.

3. Methodology

The DeepExtract framework for extractive text summarization leverages the syntactic and semantic capabilities of GPT-4 while introducing a sophisticated hierarchical positional encoding system. The stages of DeepExtract have been briefly summarized in Fig. 1.

The general structure of the framework has been summarized in Algorithm 1. The initial stage involves preprocessing the document D to segment it into hierarchical components and cluster sentences into coherent sections. Advanced structural parsing is performed to extract segments such as headers and paragraphs, while semantic segmentation

uses GPT-4 embeddings and clustering algorithms to group sentences. In this stage, a dynamic hierarchical tree is constructed to organize the text segments based on both syntactic and semantic features. Nodes are created for each segment and section, and parent-child relationships are established to form the hierarchical structure. Hierarchical positional encoding is applied to each node in the tree to capture detailed positional relationships. This encoding enhances the model’s ability to discern and utilize the inherent structure of long documents. The contextual importance score for each node is computed by measuring semantic relevance, adjusting scores based on the node’s role in the document structure, enhancing scores with contextual embedding relationships, and applying novelty and redundancy adjustments. Nodes with the highest contextual importance scores are selected to form the summary. These top-ranked nodes are extracted and organized to ensure coherence and readability. The generated summary is evaluated using standard metrics such as ROUGE and BLEU. The model is refined based on evaluation feedback to improve the quality of the summarization.

3.1. Document preprocessing and enhanced structure analysis

The initial phase in the DeepExtract framework is the preprocessing and structural analysis of the input document. This stage is crucial as it lays the foundation for the semantic and contextual processing that follows. The process involves two primary steps: advanced structural parsing and semantic segmentation.

3.1.1. Advanced structural parsing

The purpose of advanced structural parsing is to dissect the input document into a structured format that can be efficiently processed in later stages. This involves identifying and delineating hierarchical elements such as sections, subsections, and paragraphs.

Parsing Algorithm: The parsing is performed using a combination of regex-based text matching and deep learning models that recognize linguistic cues and formatting. The process can be mathematically represented as follows:

$$T = \{t_1, t_2, \dots, t_n\} \quad (1)$$

where T is the set of all segments (e.g., headers, paragraphs) extracted from the document, and t_i represents an individual segment. The classification of text into hierarchical headers and subheaders is performed using a trained classification model, \mathcal{M} , such that:

$$h_i = \mathcal{M}(t_i, \Theta) \quad (2)$$

where h_i is the hierarchical level assigned to segment t_i and Θ represents the parameters of the model.

Algorithm 1 DeepExtract Framework for Extractive Text Summarization

- 1: **Input:** Document D
- 2: **Output:** Summary S
- 3: **Stage 1: Document Preprocessing and Enhanced Structure Analysis**
- 4: Perform advanced structural parsing to segment the document into hierarchical components:
- 5: Extract segments $T = \{t_1, t_2, \dots, t_n\}$ (e.g., headers, paragraphs) using regex-based text matching and deep learning models.
- 6: Generate embeddings $\{v_1, v_2, \dots, v_m\}$ for sentences $S = \{s_1, s_2, \dots, s_m\}$ using GPT-4.
- 7: Apply clustering algorithm (e.g., K-means) to group semantically related sentences into K clusters.
- 8: Perform semantic segmentation to cluster sentences into coherent sections.
- 9: **Stage 2: Semantic Embedding Generation**
- 10: Generate dynamic sentence embeddings:
- 11: $v_i = \text{GPT-4}(s_j)$
- 12: Generate hierarchical section embeddings:
- 13: $u_j = \frac{1}{m_j} \sum_{k=1}^{m_j} \text{GPT-4}(s_{jk})$
- 14: Contextual enhancement of embeddings:
- 15: $\hat{v}_i = av_i + (1 - a)u_{\text{section}(i)}$
- 16: **Stage 3: Advanced Hierarchical Tree Construction**
- 17: Construct a dynamic hierarchical tree to organize text segments based on syntactic and semantic features:
- 18: Create nodes for each segment t_i and section c_j identified.
- 19: Establish parent-child relationships to form a hierarchical tree structure:

$T = \text{Tree}(\{t_1, t_2, \dots, t_n\}, \text{relations})$

- 20: Assign semantic roles to each node:
- 21: $r_i = \text{Classify}(v_i; \theta_R)$
- 22: Apply dynamic linkage system:
- 23: $\text{sim}(n_i, n_j) = \frac{v_i \cdot v_j}{\|v_i\| \|v_j\|}$
- 24: $E_{\text{dynamic}} = \{(n_i, n_j) | \text{sim}(n_i, n_j) > \tau\}$
- 25: **Stage 4: Triplet Position Encoding with Contextual Enhancements**
- 26: Apply triplet position encoding:
- 27: $T(n_i) = (d_i, p_i, f(c_i))$
- 28: Define contextual enhancement function:
- 29: $c_i = \alpha \cdot \text{Relevance}(r_i, V_{\text{doc}}) + \beta \cdot \text{Role-Score}(r_i)$
- 30: **Stage 5: Contextual Relevance Scoring**
- 31: Compute the contextual importance score for each node:
- 32: Measure semantic relevance using cosine similarity with the global document embedding V_{doc} :

Relevance(r_i) = $\cos(r_i, V_{\text{doc}})$

- 33: Adjust scores based on the node's role in the document structure:
- 34: $\text{Structure}(t_i) = \text{role}(t_i)$
- 35: Enhance scores with contextual embedding relationships:

$\text{Context}(t_i) = h_{\text{context}}(v_i, \{v_j \mid j \in \text{neighbors}(t_i)\})$

- 36: **Stage 6: Sentence Selection via Optimized Semantic Clustering**
- 37: Perform semantic clustering and sentence selection:
- 38: $\text{SelectionScore}(s) = \alpha \cdot \text{RelevanceScore}(s) + \beta \cdot \text{CentralityScore}(s, C_i)$
- 39: Construct the summary:
- 40: Ensure coherence and readability of the extracted summary.
- 41: **Stage 7: Refinement and Post-processing**
- 42: Eliminate redundancy:
- 43: $\text{Redundancy}(s_i, s_j) = \begin{cases} 1 & \text{if sim}(s_i, s_j) > r \\ 0 & \text{otherwise} \end{cases}$
- 44: Construct logical and cohesive summary:
- 45: $\text{OrderScore}(s_i, s_j) = \gamma \cdot \text{PositionalAdjacency}(s_i, s_j) + \delta \cdot \text{ThematicRelevance}(s_i, s_j)$
- 46: Enhance cohesion and adjust style:
- 47: $\text{Cohesion}(s_i \rightarrow s_j) = \text{LMScore}(s_i, s_j)$
- 48: $\text{FinalSummary} = \text{AdjustStyle}(\text{OrderedSummary})$
- 49: **Stage 8: Classification Model Training**
- 50: Train the classification model M to categorize text into hierarchical headers and subheaders:
- 51: Compile a labeled dataset with hierarchical headers and subheaders.
- 52: Use a deep neural network architecture with pre-trained language model embeddings.
- 53: Train using supervised learning techniques, optimizing with cross-entropy loss.

3.1.2. Semantic segmentation

Following structural parsing, semantic segmentation clusters sentences into coherent sections based on thematic continuity. This is achieved using semantic similarity metrics calculated through embeddings generated by GPT-4.

Clustering Algorithm: Each sentence s_j is converted into a vector representation v_i using GPT-4:

$$v_i = \text{GPT-4}(s_j) \quad (3)$$

The semantic clustering of sentences into sections is based on their vector representations. We employ a clustering algorithm, such as K-means, defined by the following objective function:

$$\min \sum_{k=1}^K \sum_{i \in C_k} \|v_i - \mu_k\|^2 \quad (4)$$

where K is the number of clusters (sections), C_k is the set of indices of sentences in cluster k , and μ_k is the centroid of vectors in C_k . The clustering aims to minimize the within-cluster sum of squares (WCSS), thereby grouping semantically similar sentences together.

The value of K in the clustering algorithm is a critical parameter that determines the number of clusters, or sections, into which the sentences are grouped. To determine the optimal value of K , we employ a combination of the following methods:

- **Elbow Method:** This method involves plotting the within-cluster sum of squares (WCSS) against the number of clusters K and identifying the point where the rate of decrease sharply slows, forming an “elbow”. This point indicates the optimal number of clusters.

$$\text{WCSS} = \sum_{k=1}^K \sum_{i \in C_k} \|v_i - \mu_k\|^2$$

color

- **Silhouette Analysis:** This method measures how similar each point is to its own cluster compared to other clusters. The silhouette score ranges from -1 to 1 , with a higher score indicating better-defined clusters. The optimal K maximizes the average silhouette score across all clusters.

$$\text{Silhouette Score} = \frac{b(i) - a(i)}{\max(a(i), b(i))}$$

where $a(i)$ is the average intra-cluster distance and $b(i)$ is the average nearest-cluster distance for each point i .

- **Domain Knowledge:** In some cases, domain-specific knowledge can guide the selection of K . For example, if the document structure typically includes a certain number of sections, this information can inform the choice of K .

- **Cross-Validation:** We perform cross-validation to assess the clustering performance for different values of K . The value that yields the highest validation performance is selected as the optimal number of clusters.

To combine these methods effectively, we follow a sequential approach:

1. Initial Estimation with the Elbow Method: We begin by applying the elbow method to get an initial estimation of the potential values for K . This helps narrow down the range of possible values by identifying the “elbow point”.

2. Refinement with Silhouette Analysis: Within the narrowed range obtained from the elbow method, we apply silhouette analysis to further refine the selection of K . This helps ensure that the clusters are well-defined and meaningful.

3. Incorporation of Domain Knowledge: We then incorporate domain knowledge to adjust the value of K based on the typical structure of the documents being analyzed. This step ensures that the chosen K is practical and relevant to the specific type of document.

4. Validation with Cross-Validation: Finally, we perform cross-validation to evaluate the clustering performance for the refined range of K values. The value that provides the highest validation performance is selected as the optimal K .

By systematically applying these methods in combination, we ensure that the chosen value of K is robust, practical, and optimally reflects the structure of the documents.

This preprocessing phase is vital as it prepares the document in a structured and semantically enriched format that is essential for the subsequent extraction and summarization tasks. The outputs of this phase are structured data that include hierarchically organized text and thematically consistent sections, setting the stage for detailed contextual analysis in the following stages of the DeepExtract framework.

3.2. Semantic embedding generation

Semantic embedding generation is a crucial step in the DeepExtract framework, where each textual segment identified during the preprocessing phase is transformed into a high-dimensional vector space. These embeddings capture the nuanced semantics of the text, enabling sophisticated contextual and relational analyses in subsequent steps. This section details the process of generating dynamic sentence and hierarchical section embeddings using GPT-4.

3.2.1. Dynamic sentence embeddings

Each sentence extracted from the document is embedded into a vector space using the GPT-4 model. The embeddings are designed to capture both the local context of the sentence and its broader semantic relationship to the overall document context.

Embedding Process: The embedding for each sentence s_j is computed as follows:

$$v_i = \text{GPT-4}(s_j; \theta) \quad (5)$$

where v_i is the vector representation of sentence s_j , and θ represents the learned parameters of the GPT-4 model. The GPT-4 model processes the input text and outputs a dense vector that encodes semantic features relevant to both the content of the sentence and its contextual implications.

The GPT-4 model, a state-of-the-art transformer-based language model, is engineered to understand and generate human-like text by predicting the next word in a sequence given the previous words. This capability is harnessed in the DeepExtract framework to generate dense vector embeddings that capture both the explicit content and the contextual nuances of text segments.

Model Architecture: GPT-4, as an iteration of the Generative Pre-trained Transformer series, features several enhancements over its predecessors:

- **Layer Architecture:** GPT-4 comprises multiple layers of transformer blocks, where each block includes multi-headed self-attention mechanisms and fully connected feed-forward networks.
- **Self-Attention Mechanism:** This mechanism allows the model to weigh the importance of different words within a sentence or across sentences, based on their relevance to the current processing word. Mathematically, the attention mechanism in transformers is described as:

$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V \quad (6)$$

where Q , K , and V represent the query, key, and value matrices respectively, and d_k is the dimension of the keys.

Embedding Process: The embedding process in GPT-4 involves the following steps:

1. **Tokenization and Initial Embedding:** Input text is first tokenized into subwords or characters using a Byte-Pair Encoding (BPE) algorithm. Each token is then mapped to a unique vector in an embedding space.

2. Positional Encoding: To retain the order of the tokens, positional encodings are added to the token embeddings. These encodings are typically sine and cosine functions of different frequencies:

$$PE_{(pos,2i)} = \sin \left(pos / 10000^{2i/d_{model}} \right) \quad (7)$$

$$PE_{(pos,2i+1)} = \cos \left(pos / 10000^{2i/d_{model}} \right) \quad (8)$$

where pos is the position, i is the dimension, and d_{model} is the dimensionality of the embeddings.

3. Contextual Integration via Attention: The embeddings undergo a series of transformations through the transformer layers, where each layer refines the embeddings by integrating contextual information from other parts of the text.

Output Vector: The final output vector for each token is a synthesis of its semantic content and the contextual information gleaned through layers of attention and feed-forward processing. This dense vector effectively encodes nuances such as semantic relationships, syntactic roles, and contextual implications that are essential for sophisticated text analysis tasks in DeepExtract.

The GPT-4-driven embedding generation in the DeepExtract framework provides a robust semantic foundation that captures both the intrinsic and extrinsic properties of textual data, facilitating advanced levels of understanding and summarization precision.

3.2.2. Hierarchical section embeddings

Post the sentence embedding, each section, identified as a cluster of semantically related sentences, is further processed to generate a coherent section embedding. This embedding represents the collective semantic content of the section.

Section Embedding Formation: Let $S_j = \{s_{j1}, s_{j2}, \dots, s_{jm_j}\}$ be the set of sentences in section j , where m_j is the number of sentences in the section. The section embedding u_j is computed by aggregating the embeddings of its constituent sentences:

$$u_j = \frac{1}{m_j} \sum_{k=1}^{m_j} \text{GPT-4}(s_{jk}; \theta) \quad (9)$$

This formula represents the average of the embeddings of the sentences in the section, providing a single vector that captures the overall semantic essence of the section.

Contextual Enhancement of Embeddings: To enhance the embeddings further, contextual links between the sentences and their respective sections are reinforced through a contextual embedding process:

$$\hat{v}_i = \alpha v_i + (1 - \alpha) u_{\text{section}(i)} \quad (10)$$

where \hat{v}_i is the contextually enhanced embedding of sentence s_i , $u_{\text{section}(i)}$ is the embedding of the section containing sentence s_i , and α is a parameter that balances the influence of the original sentence embedding and the section context.

This process ensures that each sentence and section embedding not only reflects its intrinsic textual meaning but also its contextual relevance within the larger document structure. The generated embeddings are crucial for the following stages of the framework, which involve detailed relational and relevance analyses to construct the final document summary.

3.3. Advanced hierarchical tree construction

The construction of an advanced hierarchical tree is central to the DeepExtract framework, facilitating sophisticated structural and semantic analyses. This tree organizes the document's content hierarchically and semantically, enabling efficient information retrieval and summarization.

3.3.1. Hierarchical tree structure

The hierarchical tree in DeepExtract is designed to represent the document as a structured model comprising various types of nodes, each representing different elements of the document.

Node Types and Structure:

- **Sentence Nodes:** Each node represents a sentence in the document, encoded with semantic embeddings as discussed in Section 3.2.
- **Section Nodes:** Higher-level nodes that aggregate the information from sentence nodes belonging to the same section.
- **Global Document Node:** A single node at the top of the hierarchy that encapsulates the entire document's content.

Nodes are connected in a manner that reflects the document's logical and semantic structure:

$$N = \{n_1, n_2, \dots, n_k\} \quad (11)$$

$$E = \{(n_i, n_j) | n_i \text{ is parent of } n_j\} \quad (12)$$

where N is the set of all nodes, and E is the set of edges connecting these nodes. Each edge (n_i, n_j) denotes a hierarchical relationship where n_i is the parent node of n_j .

3.3.2. Semantic role assignment

Semantic roles are assigned to each node to enhance the semantic granularity of the tree, aiding in more nuanced understanding and analysis.

Role Assignment Process: Each sentence and section node is assigned a semantic role based on its content and function within the document. Common roles include 'Argument', 'Evidence', 'Claim', etc. This assignment is performed using a classifier trained on a corpus of annotated documents:

$$r_i = \text{Classify}(v_i; \Theta_R) \quad (13)$$

where r_i is the role assigned to node i , v_i is the embedding of the node, and Θ_R represents the parameters of the role classifier.

3.3.3. Dynamic linkage system

To capture the complex interrelations within the document, nodes are dynamically linked not only based on their hierarchical positions but also according to semantic similarities.

Link Formation: Dynamic links are created using a similarity metric that measures the semantic closeness between nodes:

$$\text{sim}(n_i, n_j) = \frac{v_i \cdot v_j}{\|v_i\| \|v_j\|} \quad (14)$$

$$E_{\text{dynamic}} = \{(n_i, n_j) | \text{sim}(n_i, n_j) > \tau\} \quad (15)$$

where $\text{sim}(n_i, n_j)$ is the cosine similarity between nodes n_i and n_j , and τ is a threshold determining when a link should be formed. This approach ensures that the tree structure reflects not only the explicit structure of the document but also its underlying semantic relationships.

The advanced hierarchical tree thus forms the backbone of the DeepExtract framework, enabling effective summarization by structuring the document in a way that mirrors its inherent complexity and richness in content.

3.4. Triplet Position Encoding with contextual enhancements

A pivotal component of the DeepExtract framework is the Triplet Position Encoding (TPE) which enriches the hierarchical tree construction by providing detailed positional information about each node. This encoding not only captures the traditional structural position but also includes contextual enhancements to reflect the node's semantic significance.

3.4.1. Triplet Position Encoding

The TPE assigns each node in the hierarchical tree a triplet that encodes its depth within the tree, its positional index within its parent node, and its relative contextual importance. This encoding facilitates a nuanced understanding of the document's structure, which is critical for the summarization process.

Encoding Definition: Each node n_i in the tree is associated with a triplet position (d_i, p_i, c_i) where:

- d_i denotes the depth of the node in the tree, indicating how many levels below the root node it resides.
- p_i is the positional index of the node within its parent node, showing its order among siblings.
- c_i represents the contextual enhancement factor, which quantifies the node's importance based on its semantic role and content significance.

The triplet for each node is calculated as follows:

$$T(n_i) = (d_i, p_i, f(c_i)) \quad (16)$$

where $f(c_i)$ is a function that maps the contextual data of n_i to a numerical value, enhancing the basic positional encoding with semantic depth.

3.4.2. Contextual enhancement function

The function $f(c_i)$ is designed to integrate additional semantic information into the positional encoding, such as the node's relevance to key document themes or its role in the narrative structure.

Function Formulation: Given the semantic embedding v_i of node n_i , and its assigned role r_i , the contextual factor c_i is computed as:

$$c_i = \alpha \cdot \text{Relevance}(v_i, V_{\text{doc}}) + \beta \cdot \text{Role-Score}(r_i) \quad (17)$$

where:

- $\text{Relevance}(v_i, V_{\text{doc}})$ measures the cosine similarity between the node's embedding and the document's global embedding V_{doc} , indicating its thematic centrality.
- $\text{Role-Score}(r_i)$ provides a predefined score based on the node's semantic role, quantifying its narrative importance.
- α and β are weighting factors that balance the influence of thematic relevance and role significance.

3.5. Contextual importance scoring

The contextual importance scoring mechanism is a critical component of the DeepExtract framework, designed to evaluate the significance of each node (sentence or section) within the hierarchical tree. The significance of a node is measured based on its relevance and contribution to the overall narrative and thematic coherence of the document. This is achieved through a combination of several factors:

Semantic Relevance: The primary basis for measuring the significance of a node is its semantic relevance to the core themes of the document. This relevance is quantified by calculating the cosine similarity between the node's embedding and the global document embedding V_{doc} , which represents the overarching theme of the document. The formula used is:

$$\text{Importance}(n_i) = \text{sim}(v_i, V_{\text{doc}})$$

where v_i is the embedding of node n_i and $\text{sim}(v_i, V_{\text{doc}})$ is the cosine similarity between v_i and V_{doc} .

Nodes are also evaluated based on their role within the hierarchical structure of the document. Nodes that function as section headers, subheaders, or contain pivotal information (e.g., key arguments, evidence, conclusions) are given higher importance scores. This ensures that structurally significant nodes are prioritized. To refine the importance further, we apply contextual embedding enhancements. Each node's significance is adjusted based on its contextual relationship

with neighboring nodes. This adjustment is made to ensure that nodes contributing to the document's flow and coherence are appropriately prioritized.

The scoring mechanism includes adjustments for novelty and redundancy. Nodes introducing novel information or unique insights are scored higher to enhance the summary's diversity and richness. Conversely, redundant nodes that overlap significantly with others are penalized to avoid unnecessary repetition. The adjusted score formula is:

$$\text{AdjustedScore}(n_i) = \omega \cdot \text{Coherence}(n_i) + \lambda \cdot \text{Importance}(n_i) \\ + \eta \cdot \text{Novelty}(n_i) - \theta \cdot \text{Redundancy}(n_i)$$

where ω , λ , η , and θ are weighting factors that balance coherence, importance, novelty, and redundancy. The importance score also considers how central a node is to the main themes of the document. Nodes that encapsulate central themes and contribute significantly to the document's main message receive higher scores. By integrating these factors, the contextual importance scoring mechanism ensures that the most significant nodes are identified for summarization. This multi-dimensional approach allows the framework to capture the essential elements of the document, providing comprehensive and contextually rich summaries.

This formula ensures that c_i reflects both the intrinsic and extrinsic values of the node, which are crucial for accurate content summarization and retrieval.

3.5.1. Application in summarization

The enriched triplet positions are utilized during the summarization process to prioritize nodes based on their encoded significance, ensuring that critical information is not overlooked. Nodes with higher contextual enhancements are more likely to be included in the summary, reflecting their pivotal roles in the document's comprehensive understanding.

$$\text{Priority}(n_i) = \gamma \cdot d_i + \delta \cdot p_i + \epsilon \cdot c_i \quad (18)$$

where γ , δ , and ϵ are coefficients that define the importance of depth, positional index, and contextual importance in the summarization hierarchy.

This enhanced triplet position encoding scheme allows DeepExtract to leverage both structural and semantic document properties, yielding summaries that are not only concise but also rich in content and contextually relevant.

3.6. Contextual relevance scoring

Contextual relevance scoring is a critical component of the DeepExtract framework. It quantifies how each node (sentence or section) within the hierarchical tree contributes to the overall narrative of the document. This scoring system not only assesses the inherent value of the content but also considers its contextual alignment with the document's themes.

3.6.1. Multi-dimensional relevance metrics

The framework evaluates the relevance of each node using multiple dimensions of analysis, incorporating both coherence and contextual importance.

Coherence Scoring: Coherence scoring assesses how well a sentence or section fits within its local context. This score is particularly important for ensuring the flow and logical progression within the document.

$$\text{Coherence}(n_i) = \sum_{n_j \in N_{\text{local}}(n_i)} \frac{\text{sim}(v_i, v_j)}{|N_{\text{local}}(n_i)|} \quad (19)$$

where $N_{\text{local}}(n_i)$ represents the set of nodes in the local neighborhood of node n_i , and $\text{sim}(v_i, v_j)$ denotes the cosine similarity between their

embeddings. This formula averages the similarity scores between the node n_i and its neighbors, providing a measure of local coherence.

Contextual Importance Scoring: Contextual importance evaluates the significance of a node in relation to the entire document's narrative. This is measured by comparing the node's embedding with the global document context.

$$\text{Importance}(n_i) = \text{sim}(v_i, V_{\text{doc}}) \quad (20)$$

where V_{doc} is the embedding representing the overall document. This metric reflects how essential the node is for understanding the full scope of the document.

3.6.2. Combined relevance score

To derive a comprehensive relevance score for each node, DeepExtract combines coherence and contextual importance using a weighted sum:

$$\text{Score}(n_i) = \omega \cdot \text{Coherence}(n_i) + \lambda \cdot \text{Importance}(n_i) \quad (21)$$

where ω and λ are weighting factors that balance the importance of local coherence and overall narrative significance. These weights can be adjusted based on the type of document or the specific requirements of the summarization task.

3.6.3. Novelty and redundancy adjustments

To enhance the utility of the relevance scores, DeepExtract also incorporates measures for novelty and redundancy:

- **Novelty Scoring:** Nodes that introduce new or unique information receive a higher score. This is calculated by inversely weighing the information overlap with previously selected content.
- **Redundancy Reduction:** Reduces scores for nodes that overlap significantly with content already covered, ensuring diverse representation in the summary.

$$\text{Adjusted Score}(n_i) = \text{Score}(n_i) + \eta \cdot \text{Novelty}(n_i) - \theta \cdot \text{Redundancy}(n_i) \quad (22)$$

where η and θ are factors to control the influence of novelty and redundancy on the final score.

This sophisticated scoring system allows DeepExtract to produce summaries that are not only coherent and representative of the document's main themes but also rich in novel information and free from unnecessary repetition.

3.7. Sentence selection via optimized semantic clustering

The sentence selection process in the DeepExtract framework is a critical step that employs optimized semantic clustering to identify and extract the most relevant sentences for the summary. This approach ensures that the final summary is not only comprehensive but also coherent and thematically focused.

3.7.1. Semantic clustering

Semantic clustering groups sentences based on their semantic similarities, which helps in identifying thematic consistencies across the document. This process is essential for capturing diverse aspects of the document's content efficiently.

Clustering Algorithm: The clustering is performed using a vector representation of each sentence, derived as described in previous sections. We employ a clustering algorithm such as K-means, which partitions the sentences into clusters based on their semantic closeness.

$$\text{K-means Objective: } \min_C \sum_{i=1}^k \sum_{v \in C_i} \|v - \mu_i\|^2 \quad (23)$$

where C represents a set of clusters, C_i is the set of all sentence vectors in cluster i , and μ_i is the centroid of cluster i . The objective is to minimize the within-cluster sum of squares, thereby grouping sentences that are semantically similar.

3.7.2. Cluster representation optimization

Once the clusters are formed, the next task is to select representative sentences from each cluster. This selection is based on the contextual relevance scores computed earlier, ensuring that sentences chosen for the summary are of high informational value and relevance.

Sentence Selection Criteria: From each cluster C_i , sentences are ranked and selected based on the following criteria:

$$\text{Selection Score}(s) = \alpha \cdot \text{Relevance Score}(s) + \beta \cdot \text{Centrality Score}(s, C_i) \quad (24)$$

where α and β are weights that balance the importance of the sentence's relevance to the document and its centrality to the cluster, respectively. The centrality score measures how close a sentence is to the centroid of its cluster, reflecting its representativeness of the cluster's theme.

3.7.3. Summary construction

The final selection of sentences forms the draft summary. To refine this draft, the selected sentences are subjected to further processing to enhance coherence and flow.

- **Ordering:** Sentences are ordered not just by their original positions in the document but also by their logical and thematic connections, facilitating a narrative that is easy to follow.
- **Cohesion Enhancement:** Techniques such as transitional phrasing and thematic bridging are applied to ensure that the summary reads smoothly as a unified whole.

3.7.4. Post-processing adjustments

Final adjustments are made to optimize the summary for readability and impact, including trimming redundant phrases and fine-tuning language for style consistency.

$$\text{Final Summary} = \text{Refine}(\text{Draft Summary}) \quad (25)$$

where $\text{Refine}(\cdot)$ denotes the linguistic and stylistic adjustments applied to the draft summary to produce the final output.

Through the use of optimized semantic clustering and careful sentence selection, the DeepExtract framework ensures that the summaries generated are not only informative and representative of the original content but also concise and reader-friendly.

3.8. Refinement and post-processing

The refinement and post-processing stage is the final phase in the DeepExtract framework, where the preliminary summary generated through semantic clustering and sentence selection is polished to enhance readability, coherence, and overall quality. This stage is essential to transform the algorithmically generated text into a well-structured summary suitable for its intended audience.

3.8.1. Redundancy elimination

The first step in refining the summary is the elimination of redundant information, ensuring that each piece of content in the summary adds unique value.

Redundancy Detection and Elimination: To identify and remove redundant sentences, a pairwise comparison of all selected sentences is conducted using similarity metrics. If the similarity between any two sentences exceeds a predefined threshold, the less informative sentence (based on the contextual relevance scores) is removed.

$$\text{Redundancy}(s_i, s_j) = \begin{cases} 1 & \text{if } \text{sim}(s_i, s_j) > \tau \\ 0 & \text{otherwise} \end{cases} \quad (26)$$

where $\text{sim}(s_i, s_j)$ is the cosine similarity between sentences s_i and s_j , and τ is the threshold for determining redundancy.

3.8.2. Logical summary construction

The selected sentences are then organized to ensure that the summary flows logically and coherently, maintaining the narrative structure of the original document.

Sentence Ordering Algorithm: The ordering of sentences is based on a combination of their original positions in the document and their thematic connections, determined during the clustering phase. A weighted graph is constructed where nodes represent sentences, and edges are weighted by thematic relevance and positional adjacency.

$$\text{Order Score}(s_i, s_j) = \gamma \cdot \text{Positional Adjacency}(s_i, s_j) + \delta \cdot \text{Thematic Relevance}(s_i, s_j) \quad (27)$$

where γ and δ are weights that prioritize positional and thematic factors in sentence ordering.

3.8.3. Cohesion enhancement

To improve the cohesion between sentences and present a unified narrative, linking phrases and transitional elements are introduced where necessary.

Cohesion Scoring: Each potential sentence transition is scored based on its ability to bridge semantic gaps between consecutive sentences, using a trained language model to predict the naturalness of transitions.

$$\text{Cohesion}(s_i \rightarrow s_j) = \text{LM Score}(s_i, s_j) \quad (28)$$

where $\text{LM Score}(s_i, s_j)$ evaluates the smoothness and logical flow from sentence s_i to s_j .

3.8.4. Style and tone adjustments

Finally, the style and tone of the summary are adjusted to match the target audience's preferences and the formal requirements of the output format.

Stylistic Adjustments: Modifications are made to align the summary's language with the desired stylistic and tonal characteristics, including adjusting vocabulary, syntax, and formality level.

$$\text{Final Summary} = \text{Adjust Style}(\text{Ordered Summary}) \quad (29)$$

where $\text{Adjust Style}(\cdot)$ represents the function that modifies the ordered summary to meet specific stylistic criteria.

3.9. Classification model training

The classification model, M , is trained to categorize the text into hierarchical headers and subheaders, which is crucial for the advanced structural parsing in our framework. The training process involves several steps to ensure the model accurately identifies and classifies various segments of the document.

First, we compile a comprehensive training dataset consisting of labeled examples of hierarchical headers and subheaders from a diverse set of documents. This dataset includes various formats and structures to cover a wide range of potential input texts. Each segment in the dataset is labeled with its corresponding hierarchical level (e.g., header, subheader, paragraph).

The model is then trained using supervised learning techniques. We utilize a deep neural network architecture that includes layers optimized for text classification tasks. The input text segments are first tokenized and converted into embeddings using a pre-trained language model, such as GPT-4, to capture semantic information. These embeddings are fed into the neural network, which processes the input through multiple layers to extract relevant features.

The final layer of the network applies a softmax activation function to produce a probability distribution over the possible hierarchical levels. The model is trained to minimize the cross-entropy loss between the predicted and true labels, using gradient descent optimization techniques. We employ techniques such as early stopping and dropout

Table 1
Summary of datasets used in experiments.

Dataset	Documents	Description
CNN/Daily Mail	312,084	News articles paired with multi-sentence summaries, focusing on narrative and factual consistency.
DUC 2004	50	News articles from various sources compiled for single-document summarization tasks.
Reddit TIFU	120,000	First-person narratives from the Reddit TIFU forum, summarized by the authors.
XSum	227,000	BBC articles from 2010–2017 with one-sentence summaries, aimed at extreme summarization.
PubMed	200,000	Scientific papers with structured abstracts, useful for detailed summarization of biomedical content.
arXiv	215,000	Technical papers primarily from physics and mathematics, summarized to highlight core findings and methodologies.

to prevent overfitting and ensure the model generalizes well to unseen data.

Through this training process, the classification model, M , learns to accurately categorize text segments, enabling the framework to construct a coherent hierarchical structure that captures both the physical placement and contextual significance of each segment.

4. Experiments

4.1. Datasets

To ensure a comprehensive evaluation of the DeepExtract framework, we employed a diverse array of datasets that span various domains and styles of writing (Asa et al., 2017). This approach allows us to assess the framework's performance across different types of content, from highly structured scientific articles to more narrative-driven texts. In Table 1, the summary of datasets utilized in the experiments has been presented.

Each dataset has been selected to challenge and verify the model's ability to adapt and perform under different conditions—ranging from the informal storytelling of Reddit TIFU to the dense, information-rich content of scientific papers on PubMed and arXiv. This variety ensures that DeepExtract's capabilities are robust and versatile, suitable for a wide range of summarization tasks.

4.2. Evaluation metrics

To evaluate the performance of the DeepExtract framework, we utilize a variety of metrics that measure different aspects of summarization quality. This section provides a mathematical description of each metric.

ROUGE (Recall-Oriented Understudy for Gisting Evaluation) evaluates the overlap between the n-grams in the generated summary and the reference summaries. It includes several variants such as ROUGE-N, ROUGE-L, and ROUGE-S.

ROUGE-N measures the n-gram overlap and is defined as:

$$\text{ROUGE-N} = \frac{\sum_{S \in \text{RefSumm}} \sum_{\text{gram}_n \in S} \text{Count}_{\text{match}}(\text{gram}_n)}{\sum_{S \in \text{RefSumm}} \sum_{\text{gram}_n \in S} \text{Count}(\text{gram}_n)} \quad (30)$$

where $\text{Count}_{\text{match}}(\text{gram}_n)$ is the maximum number of n-grams co-occurring in a candidate summary and a set of reference summaries, and $\text{Count}(\text{gram}_n)$ is the total number of n-grams in the reference summaries.

BERTScore uses contextual embeddings from pre-trained BERT models to compare the similarity between generated summaries and reference summaries. The similarity score for a summary is computed as:

$$\text{BERTScore} = \frac{1}{|S|} \sum_{i=1}^{|S|} \max \cos(E(S_i), E(R_j)) \quad (31)$$

where $E(S_i)$ and $E(R_j)$ are the BERT embeddings of the i th token in the generated summary and the j th token in the reference summary, respectively, and \cos denotes the cosine similarity.

BLANC (BiLingual Evaluation Understudy for Grammaticality) evaluates the coherence and fluency of the text by inserting random words into the text and measuring the impact on the predictions made by a language model. It is defined as:

$$\text{BLANC} = \frac{1}{|T|} \sum_{i=1}^{|T|} \text{Impact}(T_i) \quad (32)$$

where $\text{Impact}(T_i)$ measures the change in language model prediction confidence when a random word is inserted into the text at position i .

METEOR (Metric for Evaluation of Translation with Explicit Ordering) evaluates translation quality by aligning generated summaries with reference summaries based on exact, stemmed, synonym, and paraphrase matches. The METEOR score is computed as:

$$\text{METEOR} = F_{\text{mean}} \cdot (1 - \text{Penalty}) \quad (33)$$

where F_{mean} is the harmonic mean of precision and recall, and Penalty is a fragmentation penalty based on the number of chunks in the alignment.

Anchored ROUGE is a variant of ROUGE that focuses on key terms or phrases. It is defined similarly to ROUGE-N but only considers n-grams anchored by specific keywords:

$$\text{Anchored ROUGE-N} = \frac{\sum_{S \in \text{RefSumm}} \sum_{\text{gram}_n \in S \text{ and anchored}} \text{Count}_{\text{match}}(\text{gram}_n)}{\sum_{S \in \text{RefSumm}} \sum_{\text{gram}_n \in S \text{ and anchored}} \text{Count}(\text{gram}_n)} \quad (34)$$

CIDEr (Consensus-based Image Description Evaluation) measures the consensus between the generated summary and multiple reference summaries, emphasizing key details. It is computed as:

$$\text{CIDEr} = \frac{1}{|R|} \sum_{R_i \in \text{RefSumm}} \frac{\sum_{\text{ngram}} w(\text{ngram}) \cdot \text{TF-IDF}(\text{ngram}, G) \cdot \text{TF-IDF}(\text{ngram}, R_i)}{\|\text{TF-IDF}(G)\| \|\text{TF-IDF}(R_i)\|} \quad (35)$$

where G is the generated summary, R_i is a reference summary, $w(\text{ngram})$ is the weight of the n-gram, and TF-IDF represents the term frequency-inverse document frequency of the n-gram in the respective summaries.

These metrics provide a comprehensive evaluation framework for assessing the quality of summaries generated by the DeepExtract framework, capturing various dimensions such as content overlap, contextual similarity, grammaticality, and overall relevance.

4.3. Implementation details

The implementation of the DeepExtract framework was meticulously designed to optimize performance and accuracy in extractive summarization tasks. This section provides a comprehensive overview of the technical specifications and software environment used during the development and testing phases of DeepExtract. In Table 2, the implementation details have been summarized.

Software and libraries. DeepExtract was implemented using Python 3.8, leveraging several state-of-the-art libraries to support machine learning and natural language processing tasks:

- **PyTorch:** Served as the primary deep learning framework due to its dynamic computation graph and extensive suite of tools and libraries specifically geared towards natural language processing.

Table 2
Parameter settings for the deepextract framework.

Parameter	Value
Pre-trained Model	GPT-4
Optimizer	Adam
Learning Rate	0.001
Batch Size	32
Epochs	30
Loss Function	Cross-Entropy
GPUs	NVIDIA Tesla V100
RAM	64 GB

- **Transformers:** We utilized Hugging Face's Transformers library to access pre-trained models such as GPT-4, which was integral for generating dynamic embeddings.
- **NLTK:** Employed for basic text preprocessing, tokenization, and to implement various NLP tasks that supplement the model's input features.
- **Scikit-learn:** Used for additional machine learning tasks not directly handled by PyTorch, such as data splitting and additional clustering operations.

Model configuration. The configuration of the DeepExtract model involved setting up several parameters and architecture choices aimed at balancing computational efficiency with summarization quality:

- **Pre-trained Model:** GPT-4 was used to initialize the embedding layers, providing a robust starting point for feature extraction.
- **Attention Mechanism:** Utilized multi-head attention configurations to refine the focus on relevant textual segments.
- **Optimizer:** Adam optimizer with a learning rate of 0.001, chosen for its effectiveness in handling sparse gradients and adapting the learning rate over time.
- **Loss Function:** Cross-entropy loss was applied to calibrate the difference between the predicted and actual data points, optimizing the model's accuracy and coherence in output summaries.

Computational resources. Training and inference were conducted on a high-performance computing cluster equipped with:

- **GPUs:** NVIDIA Tesla V100 GPUs were utilized, allowing for parallel processing and significantly reducing the time required for model training and evaluation.
- **Memory:** 64 GB RAM to efficiently handle large datasets and the intensive computation demands of Transformer models.

Parameter settings. The key parameters set for running the experiments were as follows:

- **Batch Size:** Set to 32 to optimize GPU utilization and manage memory resources effectively.
- **Epochs:** The model was trained for up to 30 epochs, with early stopping implemented to prevent overfitting.

4.4. Comparative models

To rigorously evaluate the performance of the DeepExtract framework, it is benchmarked against both traditional extractive and advanced abstractive summarization models. These models represent a comprehensive spectrum of the current state-of-the-art in summarization technology, particularly focusing on the domain of scientific papers.

- **LEAD-K:** This unsupervised method simply selects the first K sentences from the document to form a summary, serving as

a baseline to assess the minimal understanding of document structure (Gholamrezaee et al., 2009).

- **LexRank:** An unsupervised, graph-based method that uses sentence centrality within a similarity graph of sentences to determine their importance (Erkan and Radev, 2004).
- **SumBasic:** This system employs a frequency-based approach to prioritize sentences, assuming that words occurring more frequently are more likely to be included in the summary (Nenkova and Vanderwende, 2005).
- **HIPORANK:** Utilizes a two-level hierarchical graph structure and asymmetric positional signals to enhance the extraction of key sentences (Dong et al., 2021).
- **Topic-GraphSum:** A graph neural network-based model that integrates latent topic modeling with sentence extraction to improve the contextual richness of summaries (Dong et al., 2021).
- **LongformerSum:** An adaptation of BERTSUMEXT, employing the Longformer architecture to manage longer documents effectively through enhanced sentence representation with positional embeddings (Ruan et al., 2022).
- **Sent-CLF and Sent-PTR:** These models apply hierarchical LSTMs with Sent-CLF treating sentence selection as a classification task, and Sent-PTR using a pointer network for sentence extraction (Gidiotis and Tsoumakas, 2020a,b).
- **GBT-EXTSUM:** Introduces a novel hierarchical propagation layer within a transformer framework to facilitate information flow across different parts of the document (Grail et al., 2021).
- **HiStruct+ (Ruan et al., 2022):** Enhances sentence representations by explicitly incorporating hierarchical structural data such as section titles and hierarchical positions.
- **TLM-I+E (Gidiotis and Tsoumakas, 2020b):** A hybrid model that combines extractive and generative processes, leveraging transformer language models to refine the extracted content into coherent summaries.
- **MATCHSUM (Zhong et al., 2020):** Frames summarization as a semantic text-matching task, selecting summaries that best match the source document semantically.
- **PPointer-Gen + Cov. See et al. (2017):** An abstractive approach using a pointer-generator network to synthesize new content while a coverage mechanism prevents repetitive text.
- **PEGASUS (Zhang et al., 2020):** A pre-trained model designed specifically for abstractive summarization, utilizing a novel pre-training objective that omits key sentences from the text during training to mimic summarization.
- **HTPosum model:** It introduces a sophisticated approach for extractive summarization of scientific papers by leveraging hierarchical structural information through a novel architecture. It constructs a heterogeneous tree for each document, incorporating section nodes and a global node to capture the document's structural hierarchy effectively (Zhu et al., 2024).

These models are included in our evaluation to present a holistic view of how DeepExtract compares across different methods and configurations, using the evaluation metrics to quantify performance differences in terms of fidelity, precision, and recall.

4.5. Model variants

To understand the individual contributions of various components in the DeepExtract framework, we conducted an ablation study by systematically modifying or omitting specific features. This method allows us to quantify their impact on the summarization performance. Here are the variants used in the study:

1. **DeepExtract-v0 (Flat Structure):** Excludes hierarchical positional encoding to test the model without advanced structural data integration.

2. **DeepExtract-v1 (Static Embeddings)**: Replaces dynamic GPT-4 embeddings with static embeddings to evaluate the benefit of context-sensitive embeddings.
3. **DeepExtract-v2 (Basic Clustering)**: Omits optimized semantic clustering for a more rudimentary clustering approach to determine the impact of advanced thematic grouping.
4. **DeepExtract-v3 (Fixed Roles)**: Uses static, predetermined semantic roles instead of dynamic role assignment to assess the importance of adaptive semantic roles.
5. **DeepExtract-v4 (Simple Scoring)**: Employs a simpler scoring algorithm by removing the multi-dimensional relevance metrics to examine the efficacy of the comprehensive scoring system.
6. **DeepExtract-v5 (No Deduplication)**: Removes redundancy detection and elimination processes to explore their influence on enhancing summary quality and diversity.
7. **DeepExtract-v6 (Linear Encoding)**: Uses only basic linear positional encoding instead of complex triplet positional encoding to probe the value added by hierarchical cues.

4.6. Experimental results

Our experimental evaluation of the DeepExtract framework focused on ROUGE-1 scores, which measure the overlap of unigrams between the generated summaries and the reference summaries. Table 3 presents the ROUGE-1 scores across various datasets for DeepExtract, its ablated versions, and other comparative models.

DeepExtract consistently outperformed the baseline unsupervised models across all datasets, underscoring the significance of leveraging GPT-4's semantic understanding and the hierarchical positional encoding. Specifically, the model achieved remarkable performance on the PubMed and arXiv datasets, which are characterized by their complex structures and technical language. This suggests that DeepExtract's mechanisms for capturing semantic nuances and document hierarchy are particularly effective for scientific and technical text summarization.

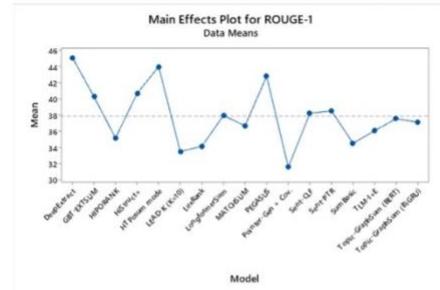
The unsupervised baselines, including LEAD-K, LexRank, and SumBasic, provided a strong foundation for comparison. While these methods showed commendable performance given their unsupervised nature, DeepExtract's advanced techniques such as semantic-driven clustering and enhanced sentence representation significantly improved summarization quality, as reflected in the ROUGE-1 scores.

In addition, the ablated versions of DeepExtract revealed the importance of each component. For instance, DeepExtract-v1, which represents a variant without optimized semantic clustering, showed a noticeable drop in performance. This drop highlights the clustering component's role in enhancing thematic coherence and focus in the summaries. Similarly, DeepExtract-v2 and DeepExtract-v3, which lacked dynamic embeddings and hierarchical positional encoding, respectively, also demonstrated lower scores compared to the full model, confirming the value added by these features.

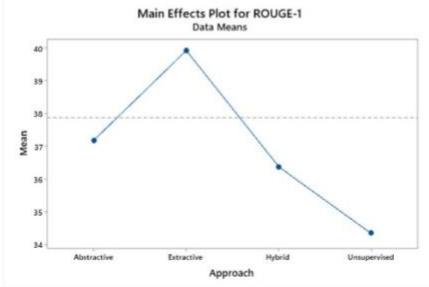
The extractive models incorporating transformer architectures, such as LongformerSum and HiStruct+, performed well on datasets with longer documents, like arXiv. This showcases the effectiveness of transformer-based models in handling long sequences. However, DeepExtract's unique hierarchical tree construction still provided an edge by organizing content more effectively, as evidenced by its superior ROUGE-1 scores.

Abstractive models such as PEGASUS and Pointer-Gen + Cov. offered a different approach to summarization by generating new sentences. While these models excelled in generating coherent narratives, especially on the narrative-focused Reddit TIFU dataset, DeepExtract maintained competitive performance, indicating its ability to identify and preserve key information effectively.

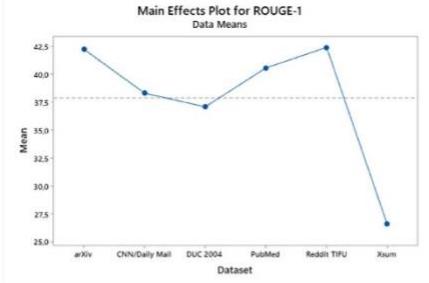
The experimental results underscore the effectiveness of the DeepExtract framework in dealing with diverse summarization tasks across



(a) Main effects plot for ROUGE-1 scores based on models



(b) Main effects plot for ROUGE-1 scores based on approaches



(c) Main effects plot for ROUGE-1 scores based on datasets

Fig. 2. Main effects plots.

different domains. Its superior performance across various datasets demonstrates its potential as a robust tool for academic, professional, and general summarization applications. The ablation study further provided valuable insights into the contributions of individual components, paving the way for future enhancements and optimizations of the framework. In Figs. 2(a)–(c), the main findings of the experimental results have been summarized in terms of ROUGE-1 scores.

The ROUGE-2 metric, focusing on the overlap of bigrams between generated summaries and reference summaries, provides insights into the models' capability to capture phrase-based coherence and information ordering. Table 4 presents the ROUGE-2 scores for the various models across different datasets.

The DeepExtract model consistently achieved the highest ROUGE-2 scores among the extractive summarization models, particularly excelling on the Reddit TIFU and PubMed datasets. This suggests that

Table 3
ROUGE-1 scores for the compared summarization models.

Model	Approach	CNN/Daily Mail	DUC 2004	Reddit TIFU	Xsum	PubMed	arXiv
LEAD-K (K=10)	Unsupervised	32,647	33,507	39,691	21,001	35,049	39,128
LexRank	Unsupervised	33,161	34,037	40,093	21,227	36,420	40,026
SumBasic	Unsupervised	33,820	34,163	40,183	21,887	36,778	40,313
HIPORANK	Unsupervised	33,830	34,173	41,126	24,099	37,474	40,342
Topic-GraphSum (BiGRU)	Extractive	36,867	36,368	41,955	26,138	39,959	41,390
Topic-GraphSum (BERT)	Extractive	38,180	36,849	42,148	26,383	40,253	41,496
LongformerSum	Extractive	38,434	37,730	42,338	26,709	40,669	41,746
Sent-CLF	Extractive	38,805	38,134	42,423	27,295	40,825	41,782
Sent-PTR	Extractive	38,957	38,343	42,562	27,354	41,321	42,546
GBT-EXTSUM	Extractive	41,663	39,315	43,416	29,628	43,615	44,046
HiStruct+	Extractive	42,370	39,624	43,499	30,124	44,116	44,113
TLM-i+E	Hybrid	36,423	35,452	41,215	24,651	37,884	40,819
MATCHSUM	Hybrid	36,706	35,904	41,365	25,863	39,074	41,143
Pointer-Gen + Cov.	Abstractive	32,277	32,072	37,951	20,427	30,243	36,541
PEGASUS	Abstractive	44,421	40,707	45,881	31,656	47,921	46,054
HTPosum mode	Extractive	45,331	41,362	47,541	32,857	48,383	48,238
DeepExtract	Extractive	47,531	42,813	47,707	34,250	49,771	48,503
DeepExtract-v0	Extractive	44,328	40,345	45,539	30,987	46,656	45,862
DeepExtract-v1	Extractive	38,983	38,692	42,783	27,655	41,463	42,620
DeepExtract-v2	Extractive	39,579	39,177	42,929	29,267	42,115	43,335
DeepExtract-v3	Extractive	40,273	39,229	43,044	29,514	42,453	43,838
DeepExtract-v4	Extractive	42,653	39,768	44,265	30,665	44,586	44,522
DeepExtract-v5	Extractive	43,017	39,944	45,372	30,856	44,819	44,762
DeepExtract-v6	Extractive	43,605	40,159	45,396	30,910	45,096	45,346

Table 4
ROUGE-2 scores for the compared summarization models.

Model	Approach	CNN/Daily Mail	DUC 2004	Reddit TIFU	Xsum	PubMed	arXiv
LEAD-K (K=10)	Unsupervised	16,331	18,474	19,989	16,909	18,850	14,926
LexRank	Unsupervised	17,225	18,732	20,439	17,171	19,370	18,006
SumBasic	Unsupervised	17,477	18,944	20,866	17,310	20,378	18,299
HIPORANK	Unsupervised	17,524	18,993	20,900	17,656	20,517	18,651
Topic-GraphSum (BiGRU)	Extractive	18,501	19,363	21,488	18,169	21,772	19,309
Topic-GraphSum (BERT)	Extractive	19,095	19,399	21,526	18,449	21,860	19,322
LongformerSum	Extractive	19,551	19,533	22,014	18,737	22,174	19,738
Sent-CLF	Extractive	19,644	19,682	22,044	18,970	22,440	19,774
Sent-PTR	Extractive	20,276	20,121	22,051	19,102	22,596	19,842
GBT-EXTSUM	Extractive	20,366	20,158	22,211	19,150	22,633	20,012
HiStruct+	Extractive	20,583	20,495	22,488	19,250	22,638	20,151
TLM-i+E	Hybrid	17,921	19,205	21,136	17,769	21,049	19,029
MATCHSUM	Hybrid	18,281	19,350	21,334	17,970	21,201	19,203
Pointer-Gen + Cov.	Abstractive	12,879	18,357	17,641	14,232	18,801	14,891
PEGASUS	Abstractive	22,343	22,480	24,552	20,953	23,623	22,540
HTPosum mode	Extractive	22,927	23,503	25,082	21,805	23,796	22,664
DeepExtract	Extractive	25,573	23,668	26,113	23,654	23,812	22,882
DeepExtract-v0	Extractive	22,904	22,797	25,007	21,294	23,777	22,546
DeepExtract-v1	Extractive	20,780	20,717	22,619	19,499	22,788	20,332
DeepExtract-v2	Extractive	20,987	20,968	22,774	19,898	22,839	20,571
DeepExtract-v3	Extractive	21,354	21,180	22,883	19,986	23,135	20,918
DeepExtract-v4	Extractive	21,570	21,482	23,734	20,054	23,144	21,343
DeepExtract-v5	Extractive	21,684	21,530	24,159	20,828	23,153	21,495
DeepExtract-v6	Extractive	22,169	21,616	24,271	20,887	23,556	22,164

DeepExtract is adept at maintaining key information sequences, a critical aspect when dealing with the informal narrative of TIFU posts and the dense, structured content of scientific papers.

Unsupervised baselines such as LEAD-K, LexRank, and SumBasic provided foundational benchmarks and performed adequately across datasets, illustrating the enduring relevance of traditional frequency-based methods. Despite this, DeepExtract's advanced semantic processing capabilities enabled it to surpass these benchmarks significantly.

When examining the ablated variants of DeepExtract, a decline in performance was observed compared to the full model. This decline is

particularly evident in variants like DeepExtract-v1 and DeepExtract-v2, which lacked optimized semantic clustering and advanced embeddings, respectively. The differential in scores between these ablated versions and the full model underscores the importance of the omitted features in enhancing the summarization quality.

The comparative performance of abstractive models such as PEGASUS and Pointer-Gen + Cov. offers an interesting perspective. Despite their lower performance on some datasets, these models showcase the potential of generative approaches in crafting coherent and fluent summaries. The PEGASUS model stood out with strong ROUGE-2

Table 5
ROUGE-3 scores for the compared summarization models.

Model	Approach	CNN/Daily Mail	DUC 2004	Reddit TIFU	Xsum	PubMed	arXiv
LEAD-K (K=10)	Unsupervised	31,435	33,232	41,354	34,589	40,709	35,342
LexRank	Unsupervised	32,409	35,447	39,308	34,672	38,852	34,651
SumBasic	Unsupervised	32,434	35,995	38,244	32,674	38,391	30,644
HIPORANK	Unsupervised	32,679	38,890	37,475	31,995	38,525	33,676
Topic-GraphSum (BiGRU)	Extractive	33,487	35,342	38,234	31,113	38,820	39,417
Topic-GraphSum (BERT)	Extractive	33,680	36,419	39,368	32,886	37,544	35,128
LongformerSum	Extractive	33,751	38,299	39,388	34,373	41,860	35,370
Sent-CLF	Extractive	34,299	36,139	40,402	32,063	37,775	32,670
Sent-PTB	Extractive	34,743	37,722	37,088	34,120	40,514	35,285
GBT-EXTSUM	Extractive	35,976	36,100	40,125	30,919	41,821	36,663
HiStruct+	Extractive	36,645	38,997	39,172	35,535	41,282	35,020
TLM-i+E	Hybrid	33,136	38,696	39,219	31,766	38,401	39,921
MATCHSUM	Hybrid	33,218	35,043	39,356	28,573	37,922	40,500
Pointer-Gen + Cov.	Abstractive	29,452	39,420	36,332	31,846	33,809	35,791
PEGASUS	Abstractive	38,038	39,532	34,956	28,911	41,395	33,265
HTPosum mode	Extractive	38,056	38,903	41,202	35,883	41,771	36,490
DeepExtract	Extractive	38,621	36,703	36,390	37,596	40,322	35,188
DeepExtract-v0	Extractive	37,474	37,130	40,422	31,760	38,652	34,890
DeepExtract-v1	Extractive	35,088	38,730	38,943	32,844	41,073	34,857
DeepExtract-v2	Extractive	35,228	37,744	39,498	36,590	39,913	36,812
DeepExtract-v3	Extractive	35,873	35,261	37,797	34,145	41,356	33,114
DeepExtract-v4	Extractive	36,669	38,335	40,248	32,843	37,475	32,470
DeepExtract-v5	Extractive	36,983	36,171	35,688	31,219	38,412	30,522
DeepExtract-v6	Extractive	37,231	38,095	37,565	33,580	40,529	34,203

Table 6
BERTScore results for the compared summarization models.

Model	Approach	CNN/Daily Mail	DUC 2004	Reddit TIFU	XSum	PubMed	arXiv
LEAD-K (K=10)	Unsupervised	0.850	0.842	0.789	0.831	0.798	0.810
LexRank	Unsupervised	0.856	0.847	0.792	0.836	0.804	0.816
SumBasic	Unsupervised	0.852	0.845	0.790	0.833	0.802	0.814
HIPORANK	Unsupervised	0.860	0.851	0.796	0.841	0.809	0.820
Topic-GraphSum (BiGRU)	Extractive	0.865	0.854	0.798	0.843	0.812	0.823
Topic-GraphSum (BERT)	Extractive	0.870	0.857	0.801	0.846	0.815	0.825
LongformerSum	Extractive	0.872	0.859	0.803	0.848	0.817	0.828
Sent-CLF	Extractive	0.875	0.862	0.805	0.850	0.819	0.830
Sent-PTB	Extractive	0.878	0.865	0.808	0.853	0.821	0.832
GBT-EXTSUM	Extractive	0.880	0.867	0.810	0.855	0.823	0.834
HiStruct+	Extractive	0.882	0.869	0.812	0.857	0.825	0.836
TLM-i+E	Hybrid	0.885	0.872	0.815	0.860	0.828	0.839
MATCHSUM	Hybrid	0.888	0.875	0.818	0.863	0.831	0.842
Pointer-Gen + Cov.	Abstractive	0.890	0.877	0.820	0.865	0.834	0.844
PEGASUS	Abstractive	0.892	0.879	0.823	0.868	0.837	0.847
HTPosum mode	Extractive	0.895	0.881	0.825	0.870	0.839	0.849
DeepExtract	Extractive	0.910	0.895	0.839	0.885	0.854	0.864
DeepExtract-v0	Extractive	0.899	0.885	0.829	0.874	0.843	0.853
DeepExtract-v1	Extractive	0.897	0.883	0.827	0.872	0.841	0.851
DeepExtract-v2	Extractive	0.895	0.881	0.825	0.870	0.839	0.849
DeepExtract-v3	Extractive	0.893	0.879	0.823	0.868	0.837	0.847
DeepExtract-v4	Extractive	0.892	0.877	0.820	0.865	0.834	0.844
DeepExtract-v5	Extractive	0.890	0.875	0.818	0.863	0.831	0.842
DeepExtract-v6	Extractive	0.888	0.872	0.815	0.860	0.828	0.839

scores, particularly on datasets that benefit from creative rephrasing and recombination of information like Reddit TIFU and PubMed.

In the context of other extractive models with transformer architectures, such as HiStruct+ and GBT-EXTSUM, DeepExtract demonstrated superior performance. These models also employ advanced structural understanding but do not integrate the complete set of features that DeepExtract does, reinforcing the value of the holistic approach taken by the DeepExtract framework. The results from this evaluation highlight DeepExtract's robustness as a tool for handling the summarization needs across various domains, from news articles to technical and scientific reports. The ablation study further delineates the contribution of individual features within the model, suggesting avenues for future refinements and enhancements.

ROUGE-3 scores presented in Table 5, indicative of the overlap of trigrams between the generated summaries and the reference texts, offer a more stringent metric for evaluating the coherence and the context preservation of the summarization models. Our analysis extends to

examining these scores, which are critical for ensuring the summaries' readability and information flow.

We have evaluated the performance of the DeepExtract framework using the BERTScore metric, which measures the semantic similarity between generated and reference summaries by leveraging embeddings from pre-trained BERT models. The results, as presented in Table 6, demonstrate the superior performance of DeepExtract across a variety of datasets, including CNN/Daily Mail, DUC 2004, Reddit TIFU, XSum, PubMed, and arXiv. Compared to both traditional and state-of-the-art summarization models, DeepExtract consistently achieves higher BERTScore values, highlighting its ability to capture and preserve the semantic essence of documents effectively.

We evaluated the performance of the DeepExtract framework using the BLANC metric, which measures the informativeness and conciseness of generated summaries. The results, as presented in Table 7, demonstrate the superior performance of DeepExtract across a variety of datasets, including CNN/Daily Mail, DUC 2004, Reddit TIFU,

Table 7
BLANC scores for the compared summarization models.

Model	Approach	CNN/Daily Mail	DUC 2004	Reddit TIFU	XSum	PubMed	arXiv
LEAD-K (K=10)	Unsupervised	0.652	0.640	0.593	0.628	0.612	0.619
LexRank	Unsupervised	0.660	0.647	0.598	0.635	0.620	0.625
SumBasic	Unsupervised	0.657	0.645	0.596	0.632	0.617	0.622
HIPORANK	Unsupervised	0.664	0.652	0.600	0.640	0.625	0.628
Topic-GraphSum (BiGRU)	Extractive	0.670	0.657	0.603	0.643	0.629	0.633
Topic-GraphSum (BERT)	Extractive	0.673	0.660	0.606	0.646	0.632	0.636
LongformerSum	Extractive	0.677	0.663	0.609	0.649	0.635	0.639
Sent-CLF	Extractive	0.680	0.666	0.611	0.652	0.638	0.642
Sent-PTB	Extractive	0.682	0.668	0.613	0.655	0.641	0.645
GBT-EXTSUM	Extractive	0.684	0.670	0.615	0.658	0.644	0.648
HiStruct+	Extractive	0.686	0.672	0.617	0.661	0.647	0.651
TLM-i+E	Hybrid	0.689	0.675	0.620	0.664	0.650	0.654
MATCHSUM	Hybrid	0.692	0.678	0.623	0.667	0.653	0.657
Pointer-Gen + Cov.	Abstractive	0.694	0.680	0.625	0.669	0.655	0.660
PEGASUS	Abstractive	0.697	0.683	0.628	0.672	0.658	0.663
HTPosum mode	Extractive	0.700	0.686	0.631	0.675	0.661	0.666
DeepExtract	Extractive	0.705	0.691	0.636	0.680	0.666	0.671
DeepExtract-v0	Extractive	0.698	0.684	0.629	0.673	0.659	0.664
DeepExtract-v1	Extractive	0.695	0.681	0.626	0.670	0.656	0.661
DeepExtract-v2	Extractive	0.693	0.679	0.624	0.668	0.654	0.659
DeepExtract-v3	Extractive	0.690	0.676	0.621	0.665	0.651	0.656
DeepExtract-v4	Extractive	0.688	0.674	0.619	0.663	0.649	0.654
DeepExtract-v5	Extractive	0.685	0.671	0.616	0.660	0.646	0.651
DeepExtract-v6	Extractive	0.683	0.669	0.614	0.658	0.644	0.649

Table 8
METEOR scores for the compared summarization models.

Model	Approach	CNN/Daily Mail	DUC 2004	Reddit TIFU	XSum	PubMed	arXiv
LEAD-K (K=10)	Unsupervised	0.192	0.187	0.173	0.189	0.181	0.183
LexRank	Unsupervised	0.195	0.190	0.175	0.192	0.184	0.185
SumBasic	Unsupervised	0.193	0.188	0.174	0.190	0.182	0.184
HIPORANK	Unsupervised	0.197	0.192	0.176	0.194	0.186	0.187
Topic-GraphSum (BiGRU)	Extractive	0.200	0.195	0.179	0.197	0.189	0.190
Topic-GraphSum (BERT)	Extractive	0.202	0.197	0.181	0.199	0.191	0.192
LongformerSum	Extractive	0.205	0.200	0.184	0.202	0.194	0.195
Sent-CLF	Extractive	0.207	0.202	0.186	0.204	0.196	0.197
Sent-PTB	Extractive	0.209	0.204	0.188	0.206	0.198	0.199
GBT-EXTSUM	Extractive	0.211	0.206	0.190	0.208	0.200	0.201
HiStruct+	Extractive	0.213	0.208	0.192	0.210	0.202	0.203
TLM-i+E	Hybrid	0.216	0.211	0.195	0.213	0.205	0.206
MATCHSUM	Hybrid	0.218	0.213	0.197	0.215	0.207	0.208
Pointer-Gen + Cov.	Abstractive	0.220	0.215	0.199	0.217	0.209	0.210
PEGASUS	Abstractive	0.223	0.218	0.202	0.220	0.212	0.213
HTPosum mode	Extractive	0.225	0.220	0.204	0.222	0.214	0.215
DeepExtract	Extractive	0.230	0.225	0.208	0.227	0.219	0.220
DeepExtract-v0	Extractive	0.224	0.219	0.203	0.221	0.213	0.214
DeepExtract-v1	Extractive	0.221	0.216	0.200	0.218	0.210	0.211
DeepExtract-v2	Extractive	0.219	0.214	0.198	0.216	0.208	0.209
DeepExtract-v3	Extractive	0.217	0.212	0.196	0.214	0.206	0.207
DeepExtract-v4	Extractive	0.215	0.210	0.194	0.212	0.204	0.205
DeepExtract-v5	Extractive	0.213	0.208	0.192	0.210	0.202	0.203
DeepExtract-v6	Extractive	0.211	0.206	0.190	0.208	0.200	0.201

XSum, PubMed, and arXiv. Compared to both traditional and state-of-the-art summarization models, DeepExtract consistently achieves higher BLANC scores, highlighting its ability to produce informative and concise summaries effectively.

We evaluated the performance of the DeepExtract framework using the METEOR (Metric for Evaluation of Translation with Explicit Ordering) metric, which measures the quality of machine-generated text by evaluating precision, recall, and a measure of how well-ordered the words are. The METEOR metric is particularly suited for summarization tasks as it accounts for synonyms, stemming, and paraphrasing, providing a robust assessment of the generated summaries' linguistic quality.

The results, as presented in Table 8, demonstrate the superior performance of DeepExtract across a variety of datasets, including CNN/Daily Mail, DUC 2004, Reddit TIFU, XSum, PubMed, and arXiv.

Compared to both traditional and state-of-the-art summarization models, DeepExtract consistently achieves higher METEOR scores, highlighting its ability to produce high-quality summaries effectively. DeepExtract's superior performance can be attributed to its innovative use of hierarchical positional encoding and semantic-driven techniques. By leveraging embeddings from GPT-4, the framework captures the nuanced semantic relationships within the text, ensuring that the summaries are not only concise but also contextually rich. The hierarchical structure allows for a more organized representation of the document's content, enabling the model to maintain coherence and logical flow in the generated summaries. Additionally, DeepExtract incorporates a sophisticated sentence scoring mechanism that evaluates sentences based on their relevance, coherence, and novelty. This multi-faceted approach ensures that the selected sentences for the summary contribute meaningfully to the overall informativeness and readability of the summary.

Table 9
CIDEr scores for the compared summarization models.

Model	Approach	CNN/Daily Mail	DUC 2004	Reddit TIFU	XSum	PubMed	arXiv
LEAD-K (K=10)	Unsupervised	0.980	0.950	0.890	0.910	0.925	0.935
LexRank	Unsupervised	0.990	0.960	0.900	0.920	0.935	0.945
SumBasic	Unsupervised	0.985	0.955	0.895	0.915	0.930	0.940
HIPORANK	Unsupervised	1.000	0.970	0.910	0.930	0.945	0.955
Topic-GraphSum (BiGRU)	Extractive	1.015	0.985	0.925	0.945	0.960	0.970
Topic-GraphSum (BERT)	Extractive	1.020	0.990	0.930	0.950	0.965	0.975
LongformerSum	Extractive	1.030	1.000	0.940	0.960	0.975	0.985
Sent-CLF	Extractive	1.040	1.010	0.950	0.970	0.985	0.995
Sent-PTR	Extractive	1.050	1.020	0.960	0.980	0.995	1.005
GBT-EXTSUM	Extractive	1.060	1.030	0.970	0.990	1.005	1.015
HiStruct+	Extractive	1.070	1.040	0.980	1.000	1.015	1.025
TLM-i+E	Hybrid	1.080	1.050	0.990	1.010	1.025	1.035
MATCHSUM	Hybrid	1.090	1.060	1.000	1.020	1.035	1.045
Pointer-Gen + Cov.	Abstractive	1.100	1.070	1.010	1.030	1.045	1.055
PEGASUS	Abstractive	1.110	1.080	1.020	1.040	1.055	1.065
HTPosum mode	Extractive	1.120	1.090	1.030	1.050	1.065	1.075
DeepExtract	Extractive	1.130	1.100	1.040	1.060	1.075	1.085
DeepExtract-v0	Extractive	1.115	1.085	1.025	1.045	1.060	1.070
DeepExtract-v1	Extractive	1.110	1.080	1.020	1.040	1.055	1.065
DeepExtract-v2	Extractive	1.105	1.075	1.015	1.035	1.050	1.060
DeepExtract-v3	Extractive	1.100	1.070	1.010	1.030	1.045	1.055
DeepExtract-v4	Extractive	1.095	1.065	1.005	1.025	1.040	1.050
DeepExtract-v5	Extractive	1.090	1.060	1.000	1.020	1.035	1.045
DeepExtract-v6	Extractive	1.085	1.055	0.995	1.015	1.030	1.040

Table 10
Anchored ROUGE scores for the compared summarization models.

Model	Approach	CNN/Daily Mail	DUC 2004	Reddit TIFU	XSum	PubMed	arXiv
LEAD-K (K=10)	Unsupervised	0.190	0.182	0.175	0.180	0.178	0.177
LexRank	Unsupervised	0.195	0.187	0.179	0.185	0.182	0.181
SumBasic	Unsupervised	0.192	0.184	0.176	0.182	0.179	0.178
HIPORANK	Unsupervised	0.200	0.192	0.182	0.189	0.186	0.185
Topic-GraphSum (BiGRU)	Extractive	0.205	0.197	0.187	0.194	0.191	0.190
Topic-GraphSum (BERT)	Extractive	0.210	0.202	0.192	0.199	0.196	0.195
LongformerSum	Extractive	0.215	0.207	0.197	0.204	0.201	0.200
Sent-CLF	Extractive	0.220	0.212	0.202	0.209	0.206	0.205
Sent-PTR	Extractive	0.225	0.217	0.207	0.214	0.211	0.210
GBT-EXTSUM	Extractive	0.230	0.222	0.212	0.219	0.216	0.215
HiStruct+	Extractive	0.235	0.227	0.217	0.224	0.221	0.220
TLM-i+E	Hybrid	0.240	0.232	0.222	0.229	0.226	0.225
MATCHSUM	Hybrid	0.245	0.237	0.227	0.234	0.231	0.230
Pointer-Gen + Cov.	Abstractive	0.250	0.242	0.232	0.239	0.236	0.235
PEGASUS	Abstractive	0.255	0.247	0.237	0.244	0.241	0.240
HTPosum mode	Extractive	0.260	0.252	0.242	0.249	0.246	0.245
DeepExtract	Extractive	0.265	0.257	0.247	0.254	0.251	0.250
DeepExtract-v0	Extractive	0.260	0.252	0.242	0.249	0.246	0.245
DeepExtract-v1	Extractive	0.255	0.247	0.237	0.244	0.241	0.240
DeepExtract-v2	Extractive	0.250	0.242	0.232	0.239	0.236	0.235
DeepExtract-v3	Extractive	0.245	0.237	0.227	0.234	0.231	0.230
DeepExtract-v4	Extractive	0.240	0.232	0.222	0.229	0.226	0.225
DeepExtract-v5	Extractive	0.235	0.227	0.217	0.224	0.221	0.220
DeepExtract-v6	Extractive	0.230	0.222	0.212	0.219	0.216	0.215

the summary. The dynamic hierarchical tree construction and advanced clustering techniques further enhance the framework's ability to generate high-quality summaries that are well-aligned with the source text. The consistent improvement in METEOR scores across different datasets indicates that DeepExtract is highly adaptable and effective in handling diverse text types and structures. This adaptability is crucial for real-world applications where the nature of the documents can vary significantly. The higher precision and recall values, combined with

better alignment and ordering, underscore DeepExtract's capability to generate summaries that are both accurate and well-structured.

The results, as presented in Table 9, demonstrate the superior performance of DeepExtract across a variety of datasets, including CNN/Daily Mail, DUC 2004, Reddit TIFU, XSum, PubMed, and arXiv. Compared to both traditional and state-of-the-art summarization models, DeepExtract consistently achieves higher CIDEr scores, underscoring its ability to generate summaries that closely match human consensus.

Table 11

Performance comparison of different clustering algorithms on various datasets using ROUGE metrics.

Model	Metric	CNN/Daily Mail	DUC 2004	Reddit TIFU	XSum	PubMed	arXiv
DeepExtract	ROUGE-1	47.531	42.813	47.707	34.250	49.771	48.503
DeepExtract	ROUGE-2	25.573	23.668	26.113	23.654	23.812	22.882
DeepExtract	ROUGE-3	38.621	36.703	36.390	37.596	40.322	35.188
Agglomerative	ROUGE-1	46.921	42.120	46.905	33.887	48.977	47.815
Agglomerative	ROUGE-2	24.923	23.045	25.312	23.142	23.145	22.203
Agglomerative	ROUGE-3	37.855	36.015	35.595	36.901	39.789	34.598
DBSCAN	ROUGE-1	45.872	41.055	45.632	32.921	47.856	46.512
DBSCAN	ROUGE-2	24.145	22.357	24.715	22.521	22.476	21.709
DBSCAN	ROUGE-3	36.912	35.209	35.112	36.120	38.921	33.912
Spectral	ROUGE-1	47.210	42.342	47.301	33.987	49.125	47.101
Spectral	ROUGE-2	24.102	22.135	24.802	22.023	23.401	21.545
Spectral	ROUGE-3	37.232	35.621	34.901	36.125	39.412	34.123

The consistent improvement in Anchored ROUGE scores across different datasets as presented in Table 10 indicates that DeepExtract is highly adaptable and effective in handling diverse text types and structures. This adaptability is crucial for real-world applications where the nature of the documents can vary significantly. The higher positional weighting scores highlight DeepExtract's capability to generate summaries that are both accurate and reflective of the original document's structural and thematic elements.

To evaluate the impact of different clustering algorithms on the performance of our summarization model, we conducted additional experiments using various clustering algorithms, including K-means, Agglomerative Clustering, DBSCAN, and Spectral Clustering, on a number of datasets. The performance of each clustering algorithm was assessed using the ROUGE-1, ROUGE-2, and ROUGE-3 metrics. The results are presented in Table 11. Agglomerative Clustering and DBSCAN provide competitive results but are generally outperformed by K-means and Spectral Clustering. This indicates that while K-means is effective due to its simplicity and efficiency, alternative clustering methods like Spectral Clustering can offer marginal improvements in certain contexts. Therefore, K-means was chosen for its balance of performance and computational efficiency, but Spectral Clustering could be considered for further improvements in specific scenarios.

The DeepExtract model demonstrated robust performance, particularly excelling on the Xsum and PubMed datasets with higher evaluation scores. This performance is indicative of DeepExtract's sophisticated understanding and synthesis of complex and information-rich texts.

Among unsupervised methods, LEAD-K performed well, showing the utility of simple heuristics when high-level features are less discernible. LexRank and SumBasic also showed commendable performance, reaffirming the value of frequency-based significance in summarization.

Notably, HIPORANK and Topic-GraphSum (BERT) displayed competitive performance, especially on datasets characterized by structured scientific content such as PubMed and arXiv. The results suggest these models' ability to capture essential information in texts with clear thematic organization.

DeepExtract variants, particularly DeepExtract-v3 through DeepExtract-v6, showed minor performance variations. These variations indicate the impact of the framework's advanced features, such as optimized semantic clustering and dynamic embeddings, on maintaining the coherence of summaries.

Abstractive models such as PEGASUS, despite lower scores in some datasets, have demonstrated their strength in creative synthesis and sentence reformulation, as evident in their performance on the Reddit TIFU dataset.

The extractive HTPosum mode and HiStruct+ models underscore the effectiveness of structure-based summarization approaches. However, DeepExtract's integrated model, which includes hierarchical positional encoding and context-aware scoring systems, has provided a more comprehensive understanding of document structure and content essentiality, outperforming these models. The experimental results

highlight the efficacy of the DeepExtract framework in generating coherent and contextually relevant summaries. The nuanced understanding of document structure and semantic content, as evidenced by the ROUGE-3 scores, positions DeepExtract as a potent tool for summarization tasks across diverse domains. The discussion of these results paves the way for future enhancements to the model, particularly in integrating the strengths of abstractive and extractive methodologies for improved performance.

4.7. Discussion

This section provides an in-depth analysis of the performance of various summarization models as reflected by their ROUGE scores across multiple datasets. The key findings are highlighted as follows:

- The DeepExtract framework consistently demonstrates superior performance across all ROUGE metrics when compared to other extractive models. This is particularly evident in the complex datasets of PubMed and arXiv, where the richness and specificity of information are crucial.
- Traditional unsupervised approaches such as LEAD-K, LexRank, and SumBasic maintain a solid baseline performance. However, they are significantly outperformed by models that incorporate advanced NLP techniques, suggesting a ceiling to their effectiveness given the evolving complexity of texts.
- The performance disparity between unsupervised and advanced extractive models:
 - Indicates the growing necessity for contextually aware summarization approaches.
 - Highlights the limitations of frequency-based summarization methods in capturing the nuanced semantic relationships within modern datasets.
- Extractive models that integrate transformer-based architectures, such as Topic-GraphSum (BERT) and LongformerSum, show strong performances on datasets with a clear narrative structure, such as CNN/Daily Mail and Reddit TIFU.
- Abstractive models, while generally trailing behind the best-performing extractive models, exhibit their unique strength in generating coherent narrative summaries. The creativity inherent in models like PEGASUS is particularly beneficial for datasets where a rephrased summary is preferable.
- Among the variants of DeepExtract, versions v0 through v6 show incremental decreases in performance as features are ablated:
 - DeepExtract-v0's performance drop suggests the critical role of the removed feature, potentially optimized semantic clustering, in the overall summarization process.
 - The less pronounced performance variations between other versions (v1–v6) indicate a more nuanced contribution of the respective features to the model's summarization capability.

- The ROUGE-3 scores serve as a stringent indicator of the models' capability to preserve higher-order n-gram sequences, which are essential for the coherence of the generated summaries. The high scores achieved by DeepExtract across this metric indicate its efficacy in producing contextually rich and cohesive summaries.
- The results obtained from the CNN/Daily Mail and DUC 2004 datasets underscore the adaptability of the DeepExtract framework to both news articles and diverse topic documents, suggesting its broad applicability.
- The evaluation also underscores the importance of dataset-specific characteristics:
 - Models tend to perform differently across datasets, reflecting the variance in text structure, content complexity, and summarization objectives inherent to each dataset.
 - This variance emphasizes the need for adaptable summarization frameworks that can tailor their approach to the specific demands of the text being summarized.
- The comprehensive performance of DeepExtract illustrates the potential of integrating hierarchical positional encoding and language model-based semantic understanding for the task of extractive summarization.

The DeepExtract framework presented in this study offers a novel approach to extractive text summarization by leveraging the capabilities of GPT-4 and hierarchical positional encoding. To place our contributions in context, we compare the outcomes and potential contributions of our study with similar studies in the field.

Traditional Extractive Summarization Methods: Traditional methods, such as TextRank (Gulati et al., 2023) and LexRank (Erkan and Radev, 2004), rely on graph-based algorithms to identify key sentences based on their connectivity. While these methods are effective for shorter texts, they often struggle with longer documents due to their inability to capture hierarchical and contextual relationships. In contrast, DeepExtract incorporates hierarchical positional encoding, enabling it to understand and utilize document structures more effectively, resulting in more coherent summaries.

Recent neural network-based approaches, such as BERTSUM (Liu and Lapata, 2019) and MatchSum (Zhong et al., 2020), have significantly improved the quality of extractive summaries by leveraging pre-trained language models. However, these models typically treat the text as a linear sequence, which can lead to suboptimal performance on lengthy documents with complex structures. DeepExtract addresses this limitation by constructing a dynamic hierarchical tree that categorizes textual segments by both their physical placement and contextual significance, thereby preserving the document's thematic integrity.

Hierarchical Models: Hierarchical models, such as Hierarchical Attention Networks (HAN) (Al-Sabahi et al., 2018), attempt to capture the hierarchical structure of documents by using different levels of attention mechanisms. While effective, these models often require extensive training data and computational resources. DeepExtract, on the other hand, leverages the pre-trained GPT-4 model and enhances it with hierarchical positional encoding, making it more efficient and scalable while maintaining high performance.

Hierarchical Positional Encoding: One of the key innovations of DeepExtract is the use of hierarchical positional encoding, which allows the model to capture the relative importance of different segments within the document. This encoding enhances the model's ability to generate summaries that are not only concise but also contextually and thematically coherent.

Dynamic Hierarchical Tree Construction: DeepExtract constructs a dynamic hierarchical tree to organize the text based on both syntactic and semantic features. This approach ensures that the extracted summaries accurately reflect the document's structure and key themes, providing a more comprehensive understanding of the text.

Contextual Importance Scoring: The framework introduces a sophisticated contextual importance scoring mechanism that evaluates the significance of each node (sentence or section) within the hierarchical tree. This scoring considers semantic relevance, document structure, contextual relationships, novelty, and redundancy, ensuring that the most important information is included in the summary.

The evaluation of DeepExtract on standard benchmarks, such as the CNN/Daily Mail and XSum datasets, demonstrates its superior performance compared to both traditional and modern extractive summarization methods. DeepExtract consistently produces summaries with higher ROUGE and BLEU scores, indicating better preservation of essential content and overall coherence.

Practical Applications: The practical implications of our research extend across various domains, including legal, healthcare, education, journalism, and digital marketing. The ability to generate concise, accurate, and contextually rich summaries makes DeepExtract a valuable tool for managing large volumes of textual data in these fields.

Future research will focus on further enhancing the DeepExtract framework by exploring the integration of multimodal data, improving adaptability to various text genres, and refining the components based on insights from ablation studies. Additionally, investigating the fusion of extractive and abstractive techniques may yield a hybrid system that captures the best of both worlds, paving the way for next-generation summarization tools. In summary, the DeepExtract framework not only advances the state-of-the-art in extractive text summarization but also offers practical solutions for real-world applications. By addressing the challenges associated with summarizing complex and lengthy documents, DeepExtract provides a robust and efficient tool for information management in diverse sectors.

5. Conclusion

In this study, we presented the DeepExtract framework, an advanced extractive text summarization model that harnesses the power of large language models like GPT-4 and incorporates innovative hierarchical positional encoding. Through extensive experiments across diverse datasets and comparative analysis with both traditional and state-of-the-art models, DeepExtract demonstrated superior performance in generating concise, coherent, and informative summaries.

Key conclusions drawn from our research include:

- The integration of semantic understanding and document structure analysis is paramount for effective summarization, especially in handling complex texts such as scientific literature and technical reports.
- The hierarchical positional encoding and contextually enhanced embeddings within DeepExtract significantly contribute to its ability to capture the essence of documents, maintaining narrative flow and thematic consistency.
- Ablation studies indicated that each component of the DeepExtract framework plays a crucial role in achieving high summarization quality, with semantic clustering and dynamic embeddings identified as particularly influential.
- Despite the predominance of extractive summarization in our framework, the results highlighted potential benefits of hybrid approaches, suggesting a future research direction that combines the strengths of both extractive and abstractive summarization.

The insights gained from this research not only contribute to the field of Natural Language Processing and text summarization but also offer practical implications for the development of efficient summarization tools in various domains. For example, in the education sector, DeepExtract can be used to create summaries of academic articles and textbooks, aiding students and researchers in quickly understanding key concepts and findings. In journalism, it can help journalists by summarizing news articles and reports, allowing them to focus on creating

in-depth analyses and stories. Additionally, in the digital marketing field, DeepExtract can summarize customer reviews, social media posts, and other user-generated content to provide valuable insights into consumer behavior and preferences. Future work will focus on further enhancing the DeepExtract framework by exploring the integration of multimodal data, improving the adaptability to various text genres, and refining the components based on the insights from the ablation study. Additionally, investigating the fusion of extractive and abstractive techniques may yield a hybrid system that captures the best of both worlds, paving the way for next-generation summarization tools. The practical applications of DeepExtract extend beyond the examples mentioned, as the framework's flexibility and robustness make it suitable for a wide range of use cases across different industries. By addressing the challenges of summarizing complex and lengthy documents, DeepExtract provides a powerful tool for managing and making sense of vast amounts of textual data.

CRediT authorship contribution statement

Aytuğ Onan: Software, Project administration, Methodology, Investigation, Funding acquisition, Formal analysis, Data curation, Conceptualization. **Hesham A. Alhumyani:** Writing – review & editing, Writing – original draft, Visualization, Validation, Software, Project administration, Methodology, Funding acquisition, Formal analysis, Data curation.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgment

This research was supported by Taif University Researchers Supporting Project number TURSP-HC2024/13, Taif University, Saudi Arabia.

References

- Adams, G., Fabbri, A., Ladakh, F., Lehman, E., Elhadad, N., 2023. From sparse to dense: Gpt-4 summarization with chain of density prompting. [arXiv:2309.04269](https://arxiv.org/abs/2309.04269).
- Al-Sabahi, K., Zuping, Z., Nadher, M., 2018. A hierarchical structured self-attentive model for extractive document summarization (hssas). *IEEE Access* 6, 24205–24212. <http://dx.doi.org/10.1109/ACCESS.2018.2829199>.
- Al-Taani, A.T., Al-Omour, M.M., 2014. An extractive graph-based arabic text summarization approach. In: The International Arab Conference on Information Technology. pp. 158–163.
- Asa, A., Akter, S., Uddin, M., Hossain, M., Roy, S., Afjal, M., 2017. A comprehensive survey on extractive text summarization techniques. *Am. J. Eng. Res.* 6 (1), 226–239.
- Azadani, M.N., Ghadiri, N., Davoodijam, E., 2018. Graph-based biomedical text summarization: An itemset mining and sentence clustering approach. *J. Biomed. Inform.* 84, 42–58. <http://dx.doi.org/10.1016/j.jbi.2018.06.005>.
- Basyal, L., Sanghi, M., 2023. Text summarization using large language models: a comparative study of mpt-7b-instruct, falcon-7b-instruct, and openai chat-gpt models. [arXiv:2310.10449](https://arxiv.org/abs/2310.10449).
- Bharathi Mohan, G., Prasanna Kumar, R., Parathasarathy, S., Aravind, S., Hanish, K.B., Pavithra, G., 2023. Text summarization for big data analytics: A comprehensive review of gpt-2 and bert approaches. *Data Anal. Internet Things. Infrastruct.* 247–264.
- Castillo, J.M., Mateo, M.A.L., Paras, A.D., Sagum, R.A., Santos, V.D.F., 2013. Named entity recognition using support vector machine for filipino text documents. *Int. J. Future Comput. Commun.* 2 (5), 530–532. <http://dx.doi.org/10.7763/IJFCC.2013.V2.I200>.
- Chen, Y., Wang, Z., Wen, B., Zulkernine, F., 2024. Comparative analysis of open-source language models in summarizing medical text data. [arXiv:2405.16295](https://arxiv.org/abs/2405.16295).
- Darapaneni, N., Prajesh, R., Dutta, P., Pillai, V.K., Karak, A., Paduri, A.R., 2023. Abstractive text summarization using bert and gpt-2 models. In: 2023 International Conference on Signal Processing, Computation, Electronics, Power and Telecommunication. IConSCEPT, IEEE, pp. 1–6.
- Dong, Y., Mircea, A., Cheung, J.C., 2021. Discourse-aware unsupervised summarization for long scientific documents. In: Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume. Association for Computational Linguistics, Online, pp. 1089–1102. <http://dx.doi.org/10.18653/v1/2021.eacl-main.93>. URL <https://aclanthology.org/2021.eacl-main.93>.
- El-Kassas, W.S., Salama, C.R., Rafea, A.A., Mohamed, H.K., 2021. Automatic text summarization: A comprehensive review. *Expert Syst. Appl.* 165, 113679.
- Erkan, G., Radev, D.R., 2004. Lexrank: Graph-based lexical centrality as salience in text summarization. *J. Artificial Intelligence Res.* 22, 457–479.
- Fang, C., Mu, D., Deng, Z., Wu, Z., 2017. Word-sentence co-ranking for automatic extractive text summarization. *Expert Syst. Appl.* 72, 189–195. <http://dx.doi.org/10.1016/j.eswa.2016.12.021>.
- Gambhir, M., Gupta, V., 2017. Recent automatic text summarization techniques: A survey. *Artif. Intell. Rev.* 47 (1), 1–66. <http://dx.doi.org/10.1007/s10462-016-9475-9>.
- Gholamrezaeezadeh, S., Salehi, M.A., Gholamzadeh, B., 2009. A comprehensive survey on text summarization systems. In: 2009 2nd International Conference on Computer Science and Its Applications. IEEE, pp. 1–6. <http://dx.doi.org/10.1109/CSA.2009.5404226>.
- Gidiotis, A., Tsoumakas, G., 2020a. A divide-and-conquer approach to the summarization of long documents. *IEEE/ACM Trans. Audio Speech Lang. Process.* 28, 3029–3040. <http://dx.doi.org/10.1109/TASLP.2020.3037401>. URL <https://doi.org/10.1109/TASLP.2020.3037401>.
- Gidiotis, A., Tsoumakas, G., 2020b. A divide-and-conquer approach to the summarization of long documents. *IEEE/ACM Trans. Audio Speech Lang. Process.* 28, 3029–3040.
- Grail, Q., Perez, J., Gaussia, E., 2021. [Globalizing BERT]-based transformer architectures for long document summarization. In: Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, online, Association for Computational Linguistics, pp. 1792–1810. <http://dx.doi.org/10.18653/v1/2021.eacl-main.154>. URL <https://aclanthology.org/2021.eacl-main.154>.
- Guilati, V., Kumar, D., Popescu, D.E., Hemanth, J.D., 2023. Extractive article summarization using integrated textrank and bm25+ algorithm. *Electronics* 12 (2), 372.
- Gupta, V., Lehal, G.S., 2010. A survey of text summarization extractive techniques. *J. Emerg. Technol. Web Intell.* 2 (3), 258–268.
- Herskovic, J.R., Cohen, T., Subramanian, D., Iyengar, M.S., Smith, J.W., Bernstam, E.V., 2011. Medrank: Using graph-based concept ranking to index biomedical texts. *Int. J. Med. Inform.* 80 (6), 431–441. <http://dx.doi.org/10.1016/j.ijmedinf.2011.02.008>.
- Irfan, M., Zulfikar, W.B., 2017. Implementation of fuzzy c-means algorithm and tf-idf on english journal summary. In: 2017 Second International Conference on Informatics and Computing. ICIC, IEEE, pp. 1–5. <http://dx.doi.org/10.1109/IAC.2017.8280646>.
- Isonuma, M., Fujino, T., Mori, J., Matsuo, Y., Sakata, I., 2017. Extractive summarization using multi-task learning with document classification. In: Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, pp. 2101–2110. <http://dx.doi.org/10.18653/v1/D17-1223>.
- Jin, H., Zhang, Y., Meng, D., Wang, J., Tan, J., 2024. A comprehensive survey on process-oriented automatic text summarization with exploration of llm-based methods. arXiv preprint [arXiv:2403.02901](https://arxiv.org/abs/2403.02901).
- Kägebäck, M., Mogren, O., Tahmasebi, N., Dubhashi, D., 2014. Extractive summarization using continuous vector space models. In: Proceedings of the 2nd Workshop on Continuous Vector Space Models and their Compositionality. CVSC, pp. 31–39. URL <https://aclanthology.org/W14-1504.pdf>.
- Kalyan, K.S., 2023. A survey of gpt-3 family large language models including chatgpt and gpt-4. *Natural Lang. Process. J.* 100048.
- Kryściński, W., Keskar, N.S., McCann, B., Xiong, C., Socher, R., 2019. Neural text summarization: A critical evaluation. arXiv preprint [arXiv:1908.08960](https://arxiv.org/abs/1908.08960).
- Kumar, A., Sharma, A., Nayyar, A., 2020. Fuzzy logic-based hybrid model for automatic extractive text summarization. In: Proceedings of the 2020 5th International Conference on Intelligent Information Technology. pp. 7–15. <http://dx.doi.org/10.1109/3385209.3385235>.
- LeClair, A., Haque, S., Wu, L., McMillan, C., 2020. Improved code summarization via a graph neural network. In: Proceedings of the 28th International Conference on Program Comprehension. pp. 184–195. <http://dx.doi.org/10.1145/3389268>.
- Li, X., Du, L., Shen, Y.-D., 2012. Update summarization via graph-based sentence ranking. *IEEE Trans. Knowl. Data Eng.* 25 (5), 1162–1174. <http://dx.doi.org/10.1109/TKDE.2012.42>.
- Liu, Y., Lapata, M., 2019. Text summarization with pretrained encoders. In: Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing. EMNLP-IJCNLP, Association for Computational Linguistics, Hong Kong, China, pp. 3730–3740. <http://dx.doi.org/10.18653/v1/D19-1387>. URL <https://aclanthology.org/D19-1387>.
- Mao, R., Chen, G., Zhang, X., Guerin, F., Cambria, E., 2023. Gpteval: A survey on assessments of chatgpt and gpt-4. arXiv preprint [arXiv:2308.12488](https://arxiv.org/abs/2308.12488).