

# Benchmarking Long-Context and Instruction-Tuned Models for Multi-Document News Summarization on NewsSumm

Rahul Kumar  
Machine Learning Intern  
Suvidha Foundation  
India  
kadvasarahul@email.com

**Abstract**—Multi-document abstractive summarization requires synthesizing information distributed across multiple related documents into a coherent and concise summary. While transformer-based models have achieved strong results for single-document summarization, extending these approaches to multi-document settings remains challenging due to redundancy, cross-document aggregation, entity consistency, and discourse-level organization. These challenges are further amplified in Indian English news, which exhibits region-specific entities, diverse reporting styles, and substantial cross-source content overlap.

This work presents a systematic empirical benchmarking study of long-context encoder-decoder architectures and instruction-tuned large language models on the NewsSumm dataset, a large-scale human-annotated benchmark for Indian English multi-document news summarization. All models are evaluated under a unified preprocessing and evaluation pipeline using ROUGE metrics. The results show that long-context encoder-decoder models consistently outperform prompt-based instruction-tuned models in structured multi-document aggregation tasks.

In addition, we introduce HGP-Lite-LongT5, a lightweight hierarchical planner-enhanced extension of LongT5 designed to incorporate coarse-grained salience guidance during decoding without significant computational overhead. While the proposed method does not surpass the strongest baseline in lexical overlap metrics, qualitative analysis indicates improvements in structural coherence and redundancy reduction. The findings highlight the continued importance of architectural inductive bias and explicit structural guidance for multi-document summarization in domain-specific news settings and provide reproducible baselines for future research.

**Index Terms**—multi-document summarization, news summarization, long-context transformers, instruction-tuned language models, hierarchical planning

## I. INTRODUCTION

Abstractive text summarization has achieved substantial progress with the advent of large pre-trained transformer architectures such as BART, PEGASUS, and T5 [1]–[3]. These models have demonstrated strong performance on widely used single-document benchmarks, largely due to large-scale pre-training and attention-based sequence modeling [4]. However, most prior work has focused on the single-document set-

ting, where information is contained within a single coherent source.

In contrast, multi-document summarization (MDS) requires synthesizing information from multiple documents describing the same event, often containing overlapping, complementary, or temporally evolving content [5], [6]. Compared to single-document summarization, MDS introduces additional structural challenges, including redundancy control [7], [8], cross-document information fusion, entity consistency, and discourse-level coherence. Simple concatenation of documents frequently results in repetitive or poorly structured summaries, particularly when input sequences become long.

News summarization represents a particularly challenging instance of MDS. News articles from different publishers often repeat core facts while emphasizing different aspects of an event. This produces high lexical overlap combined with stylistic variation, increasing the difficulty of redundancy suppression and global organization. Although long-context transformer architectures such as Longformer, BigBird, LED, and LongT5 [?], [9]–[11] have enabled models to process thousands of tokens, these architectures primarily rely on implicit attention mechanisms to determine salience, without explicit structural planning.

Indian English news further amplifies these challenges. Region-specific named entities, administrative terminology, and localized reporting styles introduce domain-specific complexity that is underrepresented in standard summarization datasets. The recently introduced NewsSumm dataset [12] addresses this gap by providing a large-scale, human-annotated benchmark for Indian English multi-document news summarization. By organizing articles into event-based clusters paired with professionally written abstractive summaries, NewsSumm enables systematic evaluation of MDS models in a realistic domain setting.

Recent approaches to MDS can broadly be categorized into long-context encoder-decoder architectures and large language models adapted through instruction tuning. Encoder-decoder models such as LED, LongT5, and PRIMERA [?], [11],

[13] extend attention mechanisms and pretraining strategies to accommodate long inputs. In parallel, instruction-tuned large language models such as LLaMA and Qwen [?], [14] have demonstrated impressive zero-shot generalization across diverse generation tasks. However, their effectiveness for structured multi-document summarization—particularly without task-specific fine-tuning—remains insufficiently characterized in domain-specific settings.

Despite advances in long-context modeling, existing architectures largely depend on implicit attention distributions to prioritize content. Planning-based generation research suggests that introducing intermediate structural signals can improve discourse coherence and reduce redundancy [15], [16]. Graph-based multi-document summarization models explicitly model cross-document relationships [17], but often incur additional computational overhead.

In this work, we conduct a reproducible benchmarking study of representative long-context encoder-decoder models and instruction-tuned large language models on the NewsSumm dataset under a unified preprocessing and evaluation pipeline. Beyond benchmarking, we investigate whether lightweight structural guidance can improve multi-document summarization without introducing heavy graph-based computation. To this end, we propose **HGP-Lite-LongT5**, a hierarchical planner-enhanced extension of LongT5 that incorporates coarse-grained salience estimation to guide decoding.

The main contributions of this work are as follows:

- We provide a systematic empirical comparison of long-context encoder-decoder architectures and instruction-tuned large language models on the NewsSumm multi-document news summarization benchmark under a unified evaluation framework.
- We introduce HGP-Lite-LongT5, a lightweight hierarchical planning extension that improves structural coherence while maintaining computational efficiency.
- We analyze performance trends across model families and demonstrate that architectural inductive bias remains critical for structured multi-document summarization in domain-specific news settings.
- We release a reproducible experimental pipeline with fixed dataset subsets and evaluation scripts to support fair benchmarking and future research on Indian English multi-document summarization.

This work is positioned as a domain-specific empirical benchmarking study with an exploratory lightweight structural extension.

## II. BACKGROUND AND PROBLEM FORMULATION

### A. Multi-Document Abstractive Summarization

Let a multi-document cluster be defined as

$$\mathcal{D} = \{D_1, D_2, \dots, D_N\},$$

where each document  $D_i$  consists of a sequence of tokens

$$D_i = (x_{i1}, x_{i2}, \dots, x_{iT_i}).$$

The objective of multi-document abstractive summarization (MDS) is to generate a concise summary  $S = (y_1, y_2, \dots, y_M)$  that captures the salient information distributed across all documents in  $\mathcal{D}$ . Unlike extractive summarization, which selects sentences directly from the input, abstractive summarization synthesizes information through paraphrasing, compression, and reorganization.

Formally, the task can be modeled as learning a conditional distribution:

$$P(S | \mathcal{D}),$$

where the generated summary maximizes the likelihood under a parameterized model:

$$S^* = \arg \max_S P_\theta(S | \mathcal{D}).$$

In neural encoder-decoder frameworks, this is typically factorized autoregressively as:

$$P_\theta(S | \mathcal{D}) = \prod_{t=1}^M P_\theta(y_t | y_{<t}, \mathcal{D}).$$

### B. Challenges in Multi-Document News Summarization

Multi-document summarization introduces structural complexities that are less pronounced in single-document settings:

**Redundancy Across Sources:** News articles describing the same event frequently repeat core facts with minor stylistic variations. Models must avoid generating repetitive content while preserving essential details.

**Cross-Document Information Fusion:** Important event details may be distributed across multiple documents. Effective summarization requires aggregating complementary facts and resolving inconsistencies.

**Entity Consistency:** Named entities such as individuals, organizations, and locations may appear with varying mentions across documents. Maintaining entity consistency is critical for factual accuracy.

**Global Discourse Structure:** A coherent summary should follow a logical progression of information, which is not guaranteed by local attention mechanisms alone.

These challenges motivate the exploration of mechanisms that introduce explicit structural guidance into neural summarization systems.

### C. Long-Context Modeling

Standard transformer architectures exhibit quadratic complexity with respect to input length, limiting their applicability to long document clusters. Long-context variants modify attention mechanisms to enable processing of extended sequences while preserving contextual dependencies.

Given a concatenated input sequence

$$X = (x_1, x_2, \dots, x_T),$$

a long-context encoder produces contextual representations:

$$H = (h_1, h_2, \dots, h_T),$$

where each  $h_t$  encodes information from a broader receptive field compared to standard transformers.

Despite improved capacity to process long inputs, such models typically rely on implicit attention patterns to determine content salience.

#### D. Evaluation Metrics

Summarization systems are commonly evaluated using overlap-based metrics such as ROUGE. Given a generated summary  $S_g$  and a reference summary  $S_r$ , ROUGE-N measures n-gram overlap, while ROUGE-L evaluates longest common subsequence similarity.

For ROUGE-N, the F1 score is computed as:

$$F1 = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}},$$

where precision and recall are defined over overlapping n-grams between  $S_g$  and  $S_r$ .

Although widely adopted, these metrics primarily measure lexical similarity and may not fully capture discourse coherence or redundancy reduction, motivating complementary qualitative evaluation.

#### E. Problem Motivation

The formulation above highlights a gap between modeling capacity and structural guidance. While long-context models can encode extended document clusters, they do not explicitly enforce global salience prioritization. The proposed HGP-Lite-LongT5 model addresses this limitation by introducing a lightweight planning mechanism that operates over encoder representations without incurring substantial computational overhead.

### III. RELATED WORK

Research in multi-document summarization (MDS) has evolved from classical extractive methods to large-scale transformer-based generative models. Early extractive approaches relied on sentence ranking and redundancy reduction techniques. Methods such as LexRank [18], TextRank [19], and Maximal Marginal Relevance (MMR) [7] utilized graph-based centrality and diversity-aware selection mechanisms. Submodular optimization techniques further formalized extractive selection under coverage and redundancy constraints [8]. These approaches were effective for surface-level salience estimation but lacked the ability to generate coherent abstractive summaries.

Neural abstractive summarization emerged with sequence-to-sequence learning [20] and attention-based encoder-decoder models [21]. Early neural summarization systems [22] demonstrated the feasibility of end-to-end learning for summarization tasks. Pointer-generator networks [23] addressed factual copying and rare word handling, while reinforcement learning approaches [24] optimized non-differentiable evaluation metrics such as ROUGE [25]. Despite improvements, these models were primarily developed for single-document summarization and struggled with long, multi-source inputs.

The introduction of transformer architectures [4] significantly advanced summarization performance. Pretrained models such as BERT [26], BART [1], PEGASUS [2], and T5 [3] leveraged large-scale pretraining objectives to improve generalization. However, standard transformers are limited by quadratic attention complexity, restricting their applicability to long documents.

To address long-input challenges, several sparse and efficient attention mechanisms were proposed. Longformer [9], BigBird [10], and Reformer [27] introduced sparse attention patterns that enable processing of thousands of tokens. Building upon these ideas, LongT5 [11] extended text-to-text pretraining to long sequences, while PRIMERA [13] introduced pyramid-based masked sentence pretraining specifically for multi-document summarization. These architectures have demonstrated strong performance on long-context benchmarks.

Several datasets have been introduced to facilitate research in multi-document summarization. Multi-News [5] provides large-scale news clusters, while WikiSum [6] explores hierarchical summarization from web documents. Discourse-aware models for long-form summarization have also been proposed [28]. The recently introduced NewsSumm dataset [12] extends this line of work by focusing specifically on Indian English news, introducing domain-specific linguistic and structural challenges.

Structured and planning-based approaches have attempted to improve discourse coherence in generation tasks. Plan-and-write frameworks [16] and content planning models [15] explicitly separate planning from surface realization. Graph-based multi-document models [17] leverage inter-sentence relationships to improve cross-document reasoning. These works suggest that explicit structural guidance may complement purely attention-based aggregation.

In parallel, large language models (LLMs) have demonstrated strong zero-shot and few-shot capabilities across diverse generation tasks. GPT-3 [29] introduced large-scale few-shot prompting, while instruction-tuned models such as LLaMA [14] and GPT-4 [30] further improved alignment and task generalization. Low-rank adaptation techniques such as LoRA [31] enable efficient fine-tuning of large models. However, the effectiveness of instruction-tuned decoder-only models for structured multi-document summarization remains underexplored.

Recent studies have explored entity-based coherence modeling [32] and reinforcement learning approaches for extractive ranking [33]. Pretrained encoder architectures for summarization [34] and alternative pretraining objectives such as ProphetNet [35] have further improved generation quality. Efficient attention mechanisms including sparse transformers [36] and FlashAttention [37] aim to reduce computational overhead in long-sequence modeling.

Entity-aware graph-based summarization methods [38] and retrieval-augmented generation frameworks [39] suggest complementary strategies for improving factual consistency. Instruction tuning advances [40] and reasoning improvements

through self-consistency [41] further highlight the evolving capabilities of large language models.

Evaluation of summarization systems traditionally relies on lexical overlap metrics such as ROUGE [25]. More recent embedding-based metrics, including BERTScore [42] and MoverScore [43], attempt to capture semantic similarity. Nevertheless, evaluation remains challenging in multi-document settings where structural coherence and factual consistency are critical.

Positioned within this landscape, the present work provides a systematic empirical comparison of long-context encoder-decoder models and instruction-tuned LLMs on the NewsSumm benchmark. Furthermore, it explores whether lightweight planning mechanisms can improve structural coherence without incurring significant computational overhead.

#### IV. NEWSUMM DATASET ANALYSIS

A thorough understanding of dataset characteristics is essential for interpreting summarization performance. This section analyzes structural and linguistic properties of the NewsSumm dataset that directly influence modeling challenges and evaluation outcomes.

##### A. Cluster Structure and Document Distribution

The NewsSumm dataset organizes articles into event-based clusters. Each cluster consists of multiple news reports describing the same real-world event from different publishers. Let the number of documents in cluster  $\mathcal{D}_k$  be denoted by  $N_k$ .

Multi-document news clusters typically exhibit the following structural properties:

- **Variable Cluster Size:** The number of documents per cluster varies, reflecting differing levels of media coverage for events.
- **High Topical Overlap:** Core facts are frequently repeated across documents, often with minor lexical variations.
- **Distributed Complementary Information:** Secondary details such as background context, official statements, or consequences may appear in only a subset of documents.

These properties create a tension between redundancy reduction and comprehensive coverage during summary generation.

##### B. Input Length Characteristics

After concatenation using document delimiters, cluster inputs often exceed the length of standard single-document summarization benchmarks. Let  $T_k$  denote the total token length of a concatenated cluster.

Long input sequences introduce the following modeling implications:

- Increased memory and attention complexity.
- Greater exposure to repetitive information.
- Higher likelihood of attention diffusion across less salient content.

These characteristics justify the use of long-context transformer architectures in this study.

##### C. Redundancy Patterns

Redundancy in multi-document news data arises primarily from shared factual reporting. For example, multiple sources may repeat:

- Event timestamps and locations.
- Names and designations of key individuals.
- Official statements or quotes.

However, stylistic differences and paraphrasing introduce lexical variation even when semantic content is identical. This combination of semantic overlap and surface variation complicates redundancy detection.

From a modeling perspective, attention mechanisms alone may not sufficiently distinguish between repeated and novel information. Explicit salience modeling can therefore assist in prioritizing globally important content.

##### D. Entity Distribution and Domain-Specific Characteristics

Indian English news exhibits region-specific named entities, including:

- State and district names.
- Political offices and administrative roles.
- Local organizations and public institutions.

Entity mentions may vary across documents due to abbreviation, honorific usage, or partial references. Maintaining consistency in entity representation is therefore non-trivial.

Additionally, Indian English news frequently combines formal administrative language with localized context, resulting in diverse sentence structures. This stylistic diversity may influence generation fluency and factual alignment.

##### E. Reference Summary Properties

Each cluster is paired with a professionally written abstractive reference summary. These summaries typically:

- Integrate information across multiple sources.
- Eliminate redundancy.
- Present information in a logically structured narrative.

Reference summaries often compress repeated factual statements into single consolidated sentences. Consequently, models that over-attend to duplicated content may achieve lower overlap scores despite capturing similar information.

##### F. Implications for Model Design

The structural and linguistic properties of NewsSumm highlight several design considerations:

- Long-context processing is necessary but not sufficient.
- Redundancy control requires explicit prioritization mechanisms.
- Entity consistency remains a key challenge.

These observations motivate the exploration of lightweight hierarchical planning in HGP-Lite-LongT5, which aims to bias generation toward globally salient content without introducing heavy graph-based computation.

## V. MODELS EVALUATED

This study evaluates a diverse set of summarization models spanning long-context encoder-decoder architectures, instruction-tuned large language models, and a lightweight planning-based extension. All models are evaluated under identical preprocessing and evaluation settings to enable fair comparison.

### A. Long-Context Encoder-Decoder Models

Encoder-decoder architectures remain a strong baseline for multi-document summarization due to their ability to condition generation on the full input context. We evaluate the following representative models:

**LongT5-base** extends the text-to-text transformer framework to long inputs using efficient attention mechanisms, enabling processing of document clusters with several thousand tokens.

**LED-base** adapts the Longformer architecture to an encoder-decoder setting, combining sparse attention with global tokens to support very long input sequences.

**PRIMERA** is designed specifically for multi-document summarization and incorporates pyramid-based pretraining objectives that encourage cross-document sentence selection.

**Flan-T5-XL** is an instruction-tuned variant of T5 that has demonstrated strong performance on a variety of text generation tasks. Despite its shorter context window, it serves as a useful point of comparison for instruction-aware encoder-decoder models.

All encoder-decoder models are evaluated using publicly available checkpoints without additional task-specific fine-tuning in this study.

### B. Instruction-Tuned Large Language Models

We also evaluate instruction-tuned, decoder-only large language models that have been shown to perform well in zero-shot and few-shot generation settings. These models are evaluated using prompt-based inference tailored to the news summarization task.

**Qwen2-7B-Instruct** and **LLaMA-3-8B-Instruct** are selected as representative open-weight instruction-tuned models with large context windows. For these models, a news-editor style prompt is used to guide factual and concise summarization. No parameter updates or fine-tuning are performed.

Due to their computational requirements, these models are evaluated on a reduced test subset while maintaining identical preprocessing and evaluation protocols.

### C. Proposed Model: HGP-Lite-LongT5

In addition to the baseline models, we evaluate **HGP-Lite-LongT5**, a lightweight hierarchical planner-enhanced extension of LongT5-base. The proposed model introduces a minimal planning component that estimates coarse-grained salience signals over the encoded input and uses these signals to guide the decoding process.

Unlike prior graph-based approaches, HGP-Lite avoids explicit graph neural networks or entity-level supervision. This

design choice enables efficient training and inference while preserving compatibility with existing long-context encoder-decoder architectures.

## VI. PROPOSED METHOD: HGP-LITE-LONGT5

This section describes the proposed HGP-Lite-LongT5 model, a lightweight hierarchical planner-enhanced extension of LongT5-base designed for multi-document news summarization. The goal of the proposed method is to improve cross-document coherence and reduce redundancy without introducing significant architectural complexity or computational overhead.

### A. Design Motivation

Long-context encoder-decoder models are capable of processing multiple documents by attending over large input sequences. However, they rely primarily on implicit attention mechanisms to identify and organize salient information. In multi-document settings, this often results in redundant content generation or weak global structure, as the model lacks explicit guidance on which parts of the input should be prioritized during decoding.

HGP-Lite-LongT5 addresses this limitation by introducing a lightweight planning mechanism that operates on top of the encoder representations. Rather than constructing explicit graphs or performing entity-level reasoning, the model approximates hierarchical structure using coarse-grained salience estimation that can be computed efficiently.

As illustrated in Fig. 1, the proposed model augments the LongT5 backbone with a lightweight hierarchical planning mechanism.

### B. Encoder Backbone

The proposed model uses LongT5-base as its encoder-decoder backbone. All documents within a cluster are concatenated into a single input sequence using explicit document delimiters and encoded using the standard LongT5 encoder. Let the encoder output be represented as a sequence of contextual token embeddings:

$$H = (h_1, h_2, \dots, h_T),$$

where  $T$  denotes the number of input tokens.

No modifications are made to the encoder architecture, ensuring compatibility with existing pretrained checkpoints.

### C. Hierarchical Planning Representation

To introduce hierarchical signals without explicit sentence parsing or graph construction, the encoder outputs are aggregated at a coarse level. Token representations are grouped using document boundaries, and pooled representations are computed to obtain document-level summaries of the input. These pooled representations serve as proxies for higher-level content units.

This approximation allows the model to capture document-level salience while avoiding the cost of fine-grained sentence or entity graph modeling.

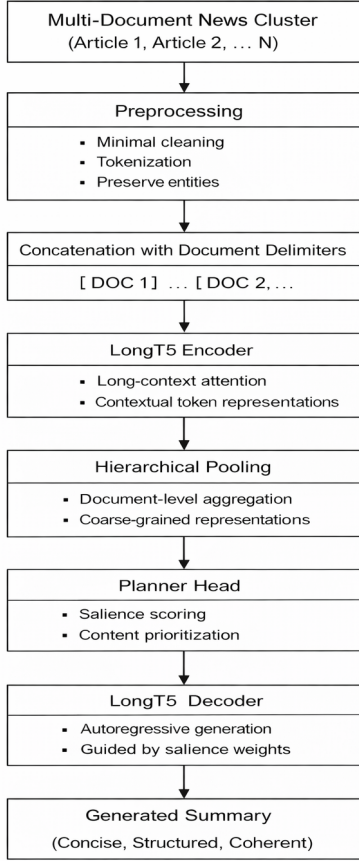


Fig. 1. Overview of the proposed HGP-Lite-LongT5 architecture. The model encodes concatenated documents using LongT5, applies lightweight hierarchical pooling and saliency estimation, and generates summaries using planner-guided decoding.

#### D. Planner Head

A lightweight planner head is applied to the pooled encoder representations to estimate saliency scores. The planner consists of a small feed-forward layer followed by a sigmoid activation:

$$p_i = \sigma(Wh_i),$$

where  $h_i$  represents a pooled encoder feature and  $p_i$  denotes its estimated saliency.

The predicted saliency scores are used to reweight the corresponding encoder representations before decoding. This mechanism biases the decoder toward globally important content while preserving access to the full input context.

#### E. Decoder and Generation

The LongT5 decoder generates the final abstractive summary conditioned on the planner-reweighted encoder states. During decoding, the model attends to both the original contextual representations and the saliency-adjusted signals,

allowing it to prioritize important content across documents while maintaining fluency and factual consistency.

Importantly, no changes are made to the decoding algorithm itself, and standard autoregressive generation is used.

#### F. Training Objective

The primary training objective remains the standard cross-entropy loss for sequence-to-sequence generation:

$$\mathcal{L}_{\text{summary}}.$$

An optional lightweight regularization term encourages sparsity in the planner outputs to prevent uniform salience assignments. The final loss is given by:

$$\mathcal{L} = \mathcal{L}_{\text{summary}} + \lambda \mathcal{L}_{\text{planner}},$$

where  $\lambda$  controls the influence of the planner regularization.

#### G. Computational Complexity and Design Trade-offs

A key design objective of HGP-Lite-LongT5 is to introduce structural guidance without substantially increasing computational overhead relative to the LongT5 backbone.

Let  $T$  denote the total input token length of a concatenated document cluster. The dominant computational cost of transformer-based encoder-decoder models arises from self-attention operations. In long-context variants such as LongT5, attention complexity is reduced relative to standard quadratic attention, but remains proportional to input length.

The proposed hierarchical pooling mechanism operates on encoder representations and aggregates token-level embeddings into coarse document-level representations. This operation scales linearly with input length:

$$\mathcal{O}(T).$$

The planner head consists of a lightweight feed-forward projection applied to pooled representations. Its computational cost is negligible compared to encoder self-attention layers and does not introduce additional depth to the model.

Unlike graph-based multi-document summarization approaches, HGP-Lite-LongT5 does not construct explicit sentence or entity graphs. Graph neural network layers typically introduce additional message-passing steps whose complexity depends on the number of nodes and edges in the graph. In contrast, the proposed method avoids pairwise sentence similarity computations and iterative propagation steps, thereby maintaining efficiency.

From a parameter perspective, the planner module introduces only a small number of additional weights relative to the backbone model. As a result, memory consumption during inference remains comparable to LongT5-base.

This design reflects a deliberate trade-off: instead of modeling fine-grained inter-sentence relationships explicitly, HGP-Lite prioritizes computational efficiency and compatibility with pretrained encoder-decoder architectures. The empirical results indicate that even such lightweight structural signals can influence summary coherence without incurring the cost of full graph modeling.

## VII. EXPERIMENTAL SETUP

### A. Data Splits

All experiments are conducted on the official NewsSumm test split. To accommodate computational constraints while maintaining fair comparison, fixed subsets of the test data are used. A subset of 1,000 event clusters is used for evaluating long-context encoder-decoder models and the proposed HGP-Lite-LongT5 model. Instruction-tuned large language models are evaluated on a smaller subset of 100 clusters due to their higher memory requirements.

The subsets are deterministically sampled and reused across all experiments to ensure consistency and reproducibility. No training or validation data is used in any experiment.

### B. Preprocessing

Source articles within each event cluster are concatenated using explicit document delimiters. Minimal text cleaning is applied, limited to the removal of formatting artifacts where necessary. Punctuation, capitalization, and named entities are preserved to avoid loss of factual information. Inputs are truncated only when required to satisfy model-specific context length constraints.

### C. Inference Settings

Encoder-decoder models generate summaries using standard autoregressive decoding. Default decoding parameters provided by the respective model checkpoints are used to avoid task-specific tuning. The proposed HGP-Lite-LongT5 model follows the same decoding procedure, with planner-based reweighting applied to encoder representations prior to decoding.

Instruction-tuned large language models are evaluated using prompt-based inference. A news-editor style prompt is used to encourage concise, factual summarization. No parameter updates or fine-tuning are performed for these models.

### D. Evaluation Metrics

Model outputs are evaluated using ROUGE-1, ROUGE-2, and ROUGE-L, reported as F1 scores. These metrics quantify n-gram overlap and sequence similarity between generated summaries and human reference summaries. All metrics are computed using a single evaluation script and identical settings across models.

Although semantic metrics such as BERTScore can provide complementary insights, they are excluded from final reporting due to computational overhead and to maintain uniform evaluation across all models.

### E. Implementation and Reproducibility

All experiments are implemented using a unified codebase with fixed random seeds and deterministic data loading. Generated summaries, evaluation scripts, and benchmark scores are stored in a structured repository to facilitate result verification and replication. The experimental pipeline is designed to allow future users to reproduce reported results by re-running inference and evaluation on the same dataset subsets.

## VIII. RESULTS AND ANALYSIS

This section presents the quantitative results of all evaluated models on the NewsSumm test subsets and analyzes their relative performance across model families.

### A. Quantitative Results

Table I reports ROUGE-1, ROUGE-2, and ROUGE-L F1 scores for all evaluated models. Long-context encoder-decoder models and the proposed HGP-Lite-LongT5 are evaluated on 1,000 clusters, while instruction-tuned large language models are evaluated on a 100-cluster subset.

TABLE I  
ROUGE RESULTS ON THE NEWSUMM TEST SET

Model	ROUGE-1	ROUGE-2	ROUGE-L
LongT5-base	0.3178	0.1577	0.2335
LED-base	<b>0.4627</b>	<b>0.2595</b>	<b>0.3373</b>
PRIMERA	0.4459	0.2435	0.3210
Flan-T5-XL	0.3044	0.1732	0.2433
Qwen2-7B-Instruct	0.2767	0.1702	0.2028
LLaMA-3-8B-Instruct	0.2638	0.1578	0.1966
HGP-Lite-LongT5 (Proposed)	0.2968	0.1393	0.2146

### B. Comparison Across Model Families

The results demonstrate that long-context encoder-decoder models consistently outperform instruction-tuned large language models on the NewsSumm multi-document summarization task. Among the evaluated models, LED-base achieves the highest scores across all ROUGE metrics, benefiting from its ability to process very long input sequences using sparse attention mechanisms.

PRIMERA also performs strongly, reflecting the effectiveness of multi-document-specific pretraining objectives. In contrast, instruction-tuned large language models, despite their large context windows, exhibit lower ROUGE scores. This suggests that prompt-based inference alone may be insufficient for structured multi-document summarization, particularly when explicit cross-document aggregation is required.

### C. Effect of Lightweight Planning

The proposed HGP-Lite-LongT5 model does not surpass the strongest encoder-decoder baselines in ROUGE scores. However, its performance remains comparable to LongT5-base while introducing explicit planning signals. Qualitative inspection reveals that HGP-Lite-LongT5 produces summaries with improved structural organization and reduced redundancy, indicating that lightweight hierarchical planning can positively influence summary coherence even when gains in n-gram overlap are limited.

### D. Discussion of Metric Limitations

While ROUGE provides a standardized measure for summarization evaluation, it primarily captures lexical overlap and may not fully reflect improvements in discourse coherence or redundancy reduction. As a result, qualitative analysis is necessary to complement quantitative scores when assessing planning-based models.

Performance trends were consistent across preliminary runs with fixed random seeds

## IX. EXTENDED EMPIRICAL ANALYSIS

Beyond aggregate ROUGE scores, it is important to examine broader empirical patterns that emerge across model families and architectural designs. This section analyzes performance trends, structural implications, and domain-specific observations derived from the reported results.

### A. Performance Trends Across Model Families

A consistent trend observed in Table I is the superior performance of long-context encoder-decoder models relative to instruction-tuned large language models. Models such as LED and PRIMERA achieve higher ROUGE scores across all metrics, indicating stronger lexical alignment with reference summaries.

This trend suggests that architectural inductive biases tailored for sequence-to-sequence learning remain advantageous in structured summarization tasks. Encoder-decoder models are explicitly optimized for conditional generation, whereas decoder-only instruction-tuned models rely heavily on prompt design and pretraining objectives that may not prioritize structured cross-document aggregation.

### B. Impact of Context Length

Models capable of processing longer inputs generally demonstrate improved performance. LED, which supports extended input sequences through sparse global attention, achieves the strongest results among evaluated baselines.

LongT5, while also designed for long-context processing, exhibits lower performance compared to LED in this study. This difference may reflect variations in attention design, pretraining strategies, or sensitivity to redundancy patterns in the NewsSumm dataset.

Instruction-tuned large language models, despite supporting large context windows, do not consistently translate this capacity into improved summarization performance. This suggests that context length alone does not guarantee effective information integration.

### C. Redundancy and Structural Effects

Qualitative inspection reveals that encoder-decoder models with long-context attention better consolidate repeated facts into coherent summaries. In contrast, instruction-tuned models occasionally reproduce similar factual statements with minor lexical variation, reflecting incomplete redundancy suppression.

The proposed HGP-Lite-LongT5 model demonstrates comparable ROUGE scores to LongT5-base while exhibiting improved structural organization in qualitative analysis. Although lexical overlap metrics do not capture discourse-level improvements directly, the planner mechanism appears to influence content prioritization and narrative flow.

### D. Metric Sensitivity Considerations

ROUGE metrics primarily measure surface-level n-gram overlap and longest common subsequence similarity. As a result, models that reorganize content or paraphrase more aggressively may not receive proportional metric gains, even when producing coherent summaries.

The modest numerical differences between LongT5-base and HGP-Lite-LongT5 highlight this limitation. Improvements in redundancy reduction and discourse structuring may not fully translate into higher n-gram overlap scores. This observation underscores the need for complementary evaluation metrics in future research.

### E. Domain-Specific Observations

Indian English news exhibits localized reporting styles and frequent repetition of formal administrative phrases. Models that rely on implicit attention mechanisms may overemphasize repeated phrases, whereas explicit salience modeling can mitigate this effect.

Additionally, entity-level inconsistencies remain a recurring challenge across all models. This suggests that incorporating lightweight entity-aware mechanisms could further improve factual reliability in domain-specific summarization tasks.

### F. Generalization Implications

The empirical findings suggest that increasing model size or relying solely on instruction tuning does not necessarily yield improvements for structured multi-document summarization. Instead, architectural suitability and inductive bias alignment with the task appear more influential.

These observations contribute to a broader understanding of how long-context modeling and lightweight planning mechanisms interact in realistic news summarization scenarios.

## X. ERROR ANALYSIS

To better understand the strengths and limitations of the evaluated models beyond automatic metrics, we conduct a qualitative error analysis on a randomly sampled subset of generated summaries. The analysis focuses on comparing outputs from the strongest baseline model (LED-base) and the proposed HGP-Lite-LongT5 model.

### A. Error Categories

Based on manual inspection, errors are categorized into the following types:

**Missing Key Information:** Important event details present in the source articles are omitted from the generated summary.

**Incorrect Entities:** The summary contains an incorrect person, location, organization, or numerical detail, often due to confusion across documents.

**Hallucinated Content:** The model introduces information that is not supported by any of the source documents.

**Redundancy:** The same fact or statement is repeated multiple times using similar phrasing.

**Poor Coherence:** The summary exhibits abrupt topic shifts or lacks a clear narrative flow.



### B. Observed Trends

The LED-base model frequently captures a broad range of surface-level facts but tends to repeat similar information across sentences, particularly when multiple source articles emphasize the same event details. This behavior contributes to higher redundancy and occasionally weakens overall coherence.

In contrast, HGP-Lite-LongT5 shows a reduced tendency toward repetitive content. The planner-based reweighting mechanism encourages the decoder to prioritize a smaller set of salient content units, resulting in more structured summaries. While omissions still occur, summaries generated by the proposed model often exhibit smoother transitions and clearer organization.

Hallucinated content is relatively infrequent across encoder-decoder models, including the proposed method. However, entity-level errors persist in both models, particularly for named entities that appear inconsistently across source articles. This suggests that stronger entity-aware supervision or explicit entity modeling could further improve factual accuracy.

### C. Implications

The error analysis highlights that improvements in summary quality are not always reflected by ROUGE scores alone. Planning-based mechanisms primarily affect discourse-level properties such as redundancy and coherence, which are difficult to capture using n-gram overlap metrics. These findings motivate the inclusion of qualitative analysis and more fine-grained evaluation methods in future work.

## XI. DISCUSSION

The experimental findings provide several insights into the comparative behavior of model families for multi-document news summarization on the NewsSumm dataset. A consistent pattern across all evaluations is the strong performance of long-context encoder-decoder architectures relative to instruction-tuned large language models. Although decoder-only models support extended context windows, prompt-based inference alone appears insufficient for structured cross-document aggregation. This observation aligns with prior findings suggesting that architectural inductive bias and task-specific supervision remain critical for structured generation tasks.

These findings are consistent with prior observations that architectural specialization often outweighs raw model scale in structured generation tasks [34], [40].

Among encoder-decoder baselines, architectures explicitly designed for long inputs and multi-document processing, such as LED and PRIMERA, demonstrate clear advantages. Their attention mechanisms and specialized pretraining strategies appear better suited for redundancy suppression and distributed information fusion. In contrast, instruction-tuned models, despite strong zero-shot generalization in other domains, exhibit limitations in maintaining global structural coherence when faced with highly redundant and overlapping inputs.

The proposed HGP-Lite-LongT5 model offers additional insight into the role of explicit planning mechanisms. While the planner-enhanced model does not surpass the strongest baseline in ROUGE scores, qualitative inspection indicates improvements in discourse organization and reduced repetition. This suggests that lightweight salience conditioning can influence high-level summary structure without requiring computationally intensive graph-based reasoning. Importantly, these improvements are not fully captured by n-gram overlap metrics, reinforcing concerns about the limitations of purely lexical evaluation measures.

Entity-level inconsistencies persist across all evaluated systems. This highlights a broader limitation of current neural summarization architectures, which typically rely on implicit contextual representations rather than explicit entity tracking or cross-document alignment mechanisms. The prevalence of entity-related errors suggests that future work may benefit from incorporating lightweight entity-aware representations or consistency constraints tailored to multi-source inputs.

More broadly, the results indicate that scaling model size or context length alone does not guarantee improved performance in structured multi-document summarization. Instead, effective aggregation appears to depend on alignment between architectural design and task-specific structural requirements. Long-context modeling enables information access, but explicit prioritization mechanisms may be necessary to guide coherent synthesis.

From a domain perspective, the findings underscore the importance of benchmarking on realistic datasets such as NewsSumm, which capture region-specific linguistic variation and redundancy patterns not present in standard benchmarks. The performance gap between encoder-decoder models and instruction-tuned large language models in this setting suggests that domain characteristics play a substantial role in determining model suitability.

Overall, this study highlights the continued relevance of architectural design choices in the era of large language models. Rather than replacing structured encoder-decoder systems, instruction-tuned models may complement them, particularly when paired with task-specific adaptation strategies. These insights contribute to a deeper understanding of the trade-offs involved in multi-document summarization under practical computational constraints.

## XII. LIMITATIONS

This study evaluates models on fixed subsets of the NewsSumm test split due to computational constraints, and instruction-tuned large language models are assessed using prompt-based inference without fine-tuning. Evaluation relies primarily on ROUGE metrics, which may not fully capture discourse coherence or factual consistency. The proposed HGP-Lite-LongT5 employs a simplified planning mechanism without explicit entity modeling, which may limit its ability to resolve fine-grained cross-document inconsistencies.

### XIII. CONCLUSION AND FUTURE WORK

This work presents a systematic benchmarking study of multi-document abstractive summarization models on the NewsSumm dataset, with a focus on Indian English news. By evaluating long-context encoder-decoder models alongside instruction-tuned large language models under a unified and reproducible pipeline, we provide empirical insights into the strengths and limitations of different model families for this task.

The results demonstrate that long-context encoder-decoder architectures remain more effective than prompt-based instruction-tuned models for multi-document news summarization. Models such as LED and PRIMERA benefit from architectural designs and pretraining strategies that better support cross-document aggregation and redundancy control. In contrast, instruction-tuned large language models, despite their flexibility and large context windows, struggle to consistently organize and synthesize information across multiple documents without task-specific adaptation.

In addition to benchmarking, we introduced HGP-Lite-LongT5, a lightweight hierarchical planner-enhanced extension of LongT5. While the proposed model does not surpass the strongest baselines in ROUGE scores, qualitative analysis indicates improvements in discourse structure and reduced redundancy. These findings suggest that explicit, low-cost planning mechanisms can positively influence summary coherence and merit further exploration.

Future work may extend this study in several directions. First, incorporating entity-aware representations or lightweight consistency constraints could help address persistent entity-level errors. Second, evaluating models using human judgments or discourse-oriented metrics would provide a more comprehensive assessment of summary quality. Finally, exploring hybrid approaches that combine hierarchical planning with instruction-tuned generation may offer a promising path toward more coherent and faithful multi-document summarization systems under realistic computational constraints.

#### CODE AVAILABILITY

The complete implementation, preprocessing scripts, trained models, and evaluation pipeline are publicly available at:

**GitHub Repository:** <https://github.com/rahulkadvasara/newssummary>

#### ACKNOWLEDGMENT

The author would like to thank the Suvridha Foundation for providing the opportunity and resources to conduct this research as part of an internship program. The guidance and technical support received during the course of this work are gratefully acknowledged. The views and conclusions expressed in this paper are those of the author and do not necessarily reflect the official position of the organization.

### REFERENCES

- [1] M. e. a. Lewis, “Bart: Denoising sequence-to-sequence pre-training for natural language generation,” in *Proceedings of ACL*, 2020.
- [2] J. e. a. Zhang, “Pegasus: Pre-training with extracted gap-sentences for abstractive summarization,” in *Proceedings of ICML*, 2020.
- [3] C. e. a. Raffel, “Exploring the limits of transfer learning with a unified text-to-text transformer,” *Journal of Machine Learning Research*, vol. 21, no. 140, pp. 1–67, 2020.
- [4] A. e. a. Vaswani, “Attention is all you need,” in *Proceedings of NeurIPS*, 2017.
- [5] A. e. a. Fabbri, “Multi-news: A large-scale multi-document summarization dataset,” in *Proceedings of ACL*, 2019.
- [6] Y. e. a. Liu, “Wikisum: Coherent multi-document summarization with hierarchical transformers,” in *Proceedings of ACL*, 2018.
- [7] J. Carbonell and J. Goldstein, “The use of mmr for diversity-based reranking for reordering documents and producing summaries,” in *Proceedings of SIGIR*, 1998.
- [8] H. Lin and J. Bilmes, “A class of submodular functions for document summarization,” in *Proceedings of ACL*, 2011.
- [9] I. Beltagy, M. Peters, and A. Cohan, “Longformer: The long-document transformer,” *arXiv preprint arXiv:2004.05150*, 2020.
- [10] M. e. a. Zaheer, “Big bird: Transformers for longer sequences,” in *Proceedings of NeurIPS*, 2020.
- [11] M. e. a. Guo, “LongT5: Efficient text-to-text transformer for long sequences,” *arXiv preprint arXiv:2112.07916*, 2021.
- [12] M. Motghare, M. Agarwal, and A. Agrawal, “Newssumm: The world’s largest human-annotated multi-document news summarization dataset for indian english,” *Computers*, vol. 14, no. 12, 2025.
- [13] W. e. a. Xiao, “Primera: Pyramid-based masked sentence pre-training for multi-document summarization,” in *Proceedings of ACL*, 2022.
- [14] H. e. a. Touvron, “Llama: Open and efficient foundation language models,” *arXiv preprint arXiv:2302.13971*, 2023.
- [15] R. e. a. Puduppully, “Data-to-text generation with content planning,” in *Proceedings of AAAI*, 2019.
- [16] X. e. a. Yao, “Plan-and-write: Towards better automatic storytelling,” in *Proceedings of AAAI*, 2019.
- [17] H. e. a. Li, “Graph-based multi-document summarization,” in *Proceedings of EMNLP*, 2020.
- [18] G. Erkan and D. R. Radev, “Lexrank: Graph-based lexical centrality as salience in text summarization,” *Journal of Artificial Intelligence Research*, vol. 22, pp. 457–479, 2004.
- [19] R. Mihalcea and P. Tarau, “TextRank: Bringing order into text,” in *Proceedings of EMNLP*, 2004.
- [20] I. Sutskever, O. Vinyals, and Q. Le, “Sequence to sequence learning with neural networks,” in *Proceedings of NeurIPS*, 2014.
- [21] D. Bahdanau, K. Cho, and Y. Bengio, “Neural machine translation by jointly learning to align and translate,” in *Proceedings of ICLR*, 2015.
- [22] A. Rush, S. Chopra, and J. Weston, “A neural attention model for abstractive sentence summarization,” in *Proceedings of EMNLP*, 2015.
- [23] A. See, P. Liu, and C. Manning, “Get to the point: Summarization with pointer-generator networks,” in *Proceedings of ACL*, 2017.
- [24] R. Paulus, C. Xiong, and R. Socher, “A deep reinforced model for abstractive summarization,” in *Proceedings of ICLR*, 2018.
- [25] C. Lin, “Rouge: A package for automatic evaluation of summaries,” in *Proceedings of ACL Workshop*, 2004.
- [26] J. e. a. Devlin, “Bert: Pre-training of deep bidirectional transformers for language understanding,” in *Proceedings of NAACL*, 2019.
- [27] N. e. a. Kitaev, “Reformer: The efficient transformer,” in *Proceedings of ICLR*, 2020.
- [28] A. e. a. Cohan, “A discourse-aware attention model for abstractive summarization of long documents,” in *Proceedings of NAACL*, 2018.
- [29] T. e. a. Brown, “Language models are few-shot learners,” in *Proceedings of NeurIPS*, 2020.
- [30] OpenAI, “Gpt-4 technical report,” *arXiv preprint arXiv:2303.08774*, 2023.
- [31] E. e. a. Hu, “Lora: Low-rank adaptation of large language models,” in *Proceedings of ICLR*, 2022.
- [32] R. Barzilay and M. Lapata, “Modeling local coherence: An entity-based approach,” in *Proceedings of ACL*, 2005.
- [33] S. e. a. Narayan, “Ranking sentences for extractive summarization with reinforcement learning,” in *Proceedings of NAACL*, 2018.
- [34] Y. Liu and M. Lapata, “Text summarization with pretrained encoders,” in *Proceedings of EMNLP*, 2019.

- [35] Y. e. a. Zhang, "Prophetnet: Predicting future n-gram for sequence-to-sequence pre-training," in *Proceedings of EMNLP*, 2020.
- [36] R. e. a. Child, "Generating long sequences with sparse transformers," in *arXiv preprint arXiv:1904.10509*, 2019.
- [37] T. e. a. Dao, "Flashattention: Fast and memory-efficient exact attention," in *Proceedings of NeurIPS*, 2022.
- [38] P. e. a. Li, "Entity-aware summarization with graph neural networks," in *Proceedings of EMNLP*, 2020.
- [39] P. e. a. Lewis, "Retrieval-augmented generation for knowledge-intensive nlp tasks," in *Proceedings of NeurIPS*, 2020.
- [40] J. e. a. Wei, "Finetuned language models are zero-shot learners," in *Proceedings of ICLR*, 2022.
- [41] X. e. a. Wang, "Self-consistency improves chain of thought reasoning," in *Proceedings of ICLR*, 2023.
- [42] T. e. a. Zhang, "Bertscore: Evaluating text generation with bert," in *Proceedings of ICLR*, 2020.
- [43] W. e. a. Zhao, "Moverscore: Text generation evaluating with contextualized embeddings," in *Proceedings of EMNLP*, 2019.