

Project Documentation

Title

Contradiction-Preserving Abstractive Multi-Document Summarization for Indian News Using Large Language Models

1. Introduction

News articles from different sources often report the **same event differently**, leading to **conflicting claims**. Most existing abstractive summarization models attempt to generate a single coherent summary by **merging all information**, which often results in **loss of important contradictions**.

This project focuses on building a summarization system that:


- Works on **Indian news articles**
- Uses **Large Language Models (LLMs)**
- **Explicitly preserves contradictions** instead of hiding them

2. Dataset Description

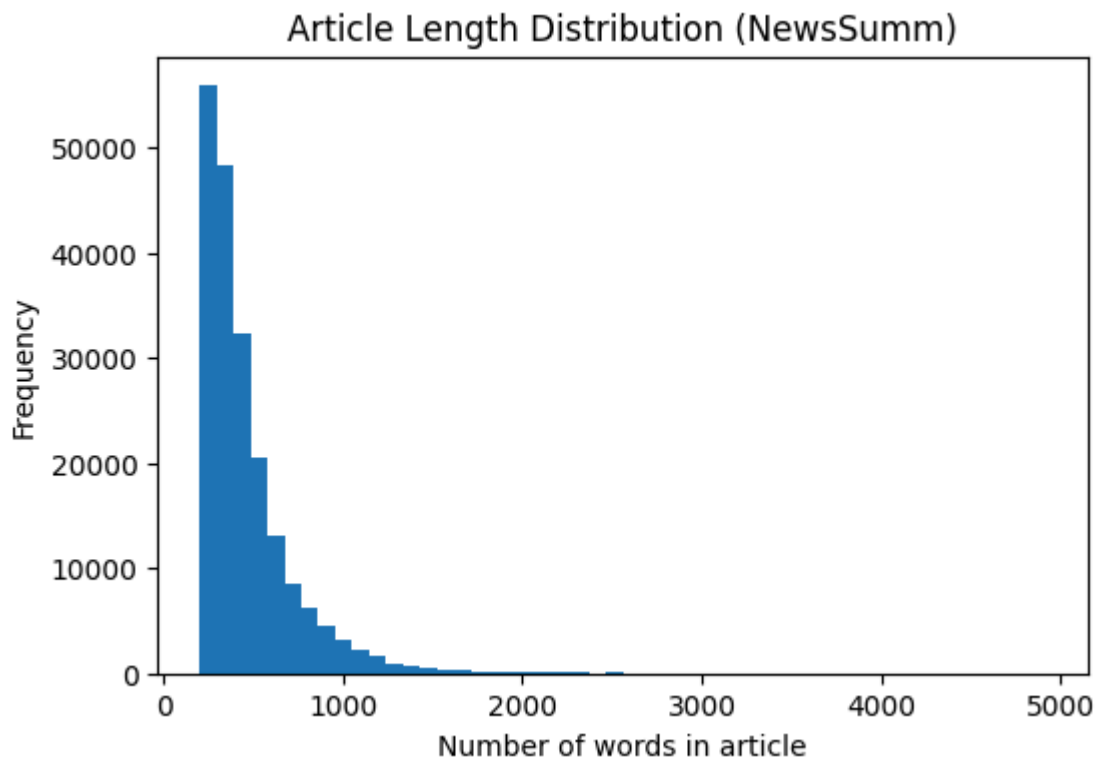
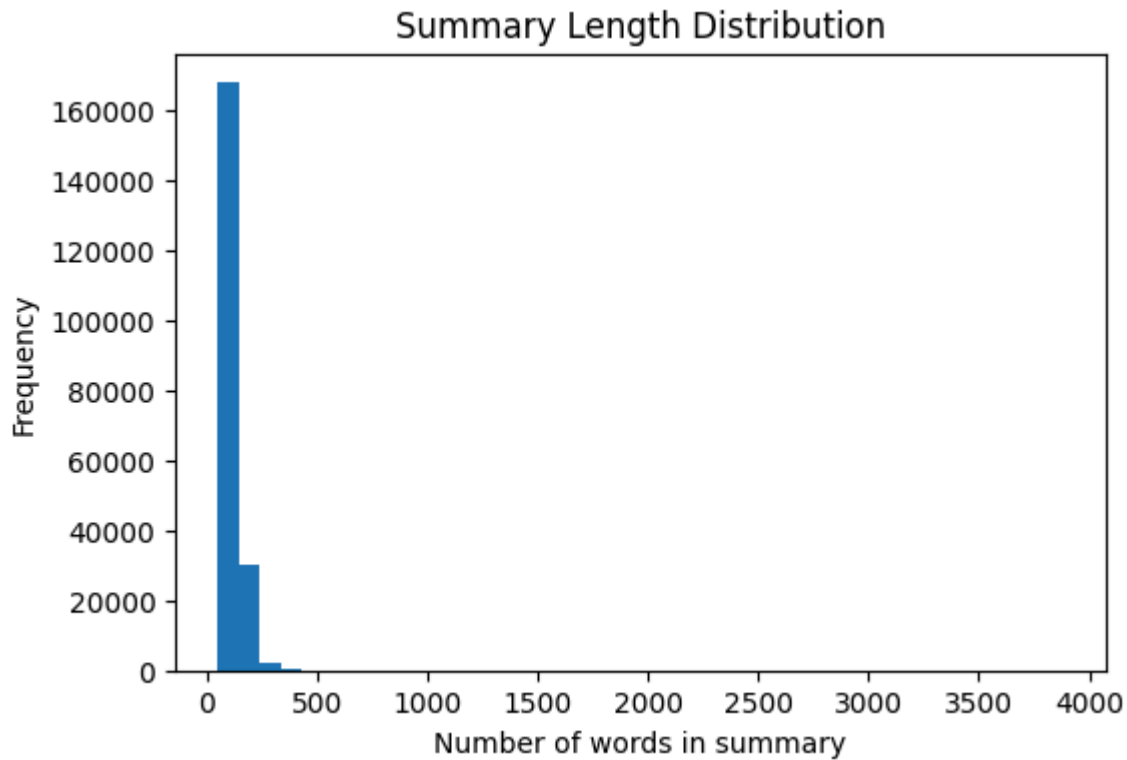
- **Dataset:** NewsSumm (Indian News Dataset)
- **Content:**
 - News articles
 - Headlines
 - Human-written summaries
- Articles are grouped into **events**, enabling **multi-document summarization**

Preprocessing Steps:

- Text cleaning (HTML tags, extra spaces)
- Length filtering
- Event-level clustering using TF-IDF similarity

 Image to add here

- Article length distribution
- Summary length distribution



3. Dataset Analysis & Visualization

We analyzed the dataset to understand its characteristics.

Visualizations:

- Article length distribution
- Summary length distribution
- Category distribution
- Cluster size distribution

Dataset shape: (348766, 6)

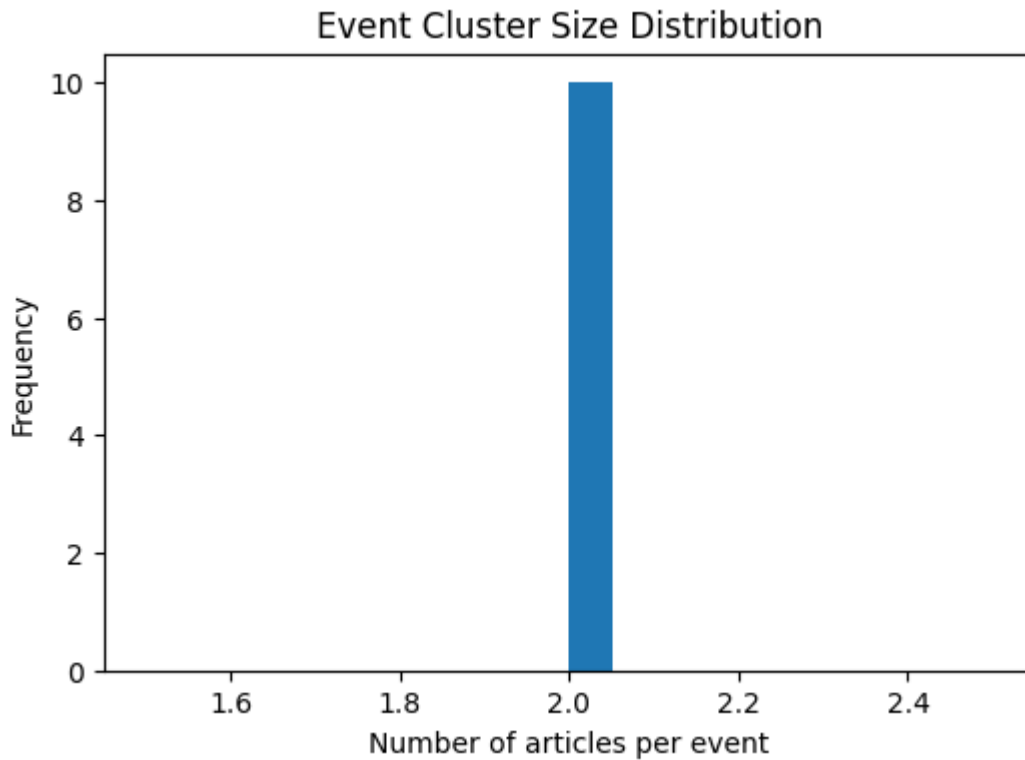
	newspaper_name	published_date\n	headline	article_text	human_summary	news_category
0	Indian Express	2020-06-01 00:00:00	Virus may be invisible enemy but COVID warrior...	Prime Minister Narendra Modi Monday hailed the...	Prime Minister of India said that the Virus ma...	National News
1	Economic Times	2013-02-11 00:00:00	Economy can bounce back, says PM Modi	ALLAHABAD: At least 20 persons were killed, an...	In Maha Kumbh, nearly 20 persons were killed. ...	National News
2	Business Standard	2013-02-11 00:00:00	At least 20 killed in stampede in Allahabad	At least 20 people were killed, and scores of ...	As per the sources 20 people died and scores w...	National News
3	Money Control	2013-02-11 00:00:00	Maha Kumbh: Over 20 dead in Allahabad station ...	More than 20 people were feared dead and 30 ot...	At least 20 people killed and 20 people are in...	National News
4	The Mint	2023-10-02 00:00:00	Gandhian wisdom	This Gandhi Jayanti, we should reflect upon an...	In this article, the author reflects on Mahatm...	National News

```
# Create Event Clusters
clusters = []
visited = set()
threshold = 0.75

for i in tqdm(range(len(df_small))):
    if i in visited:
        continue
    cluster = [i]
    for j in range(len(df_small)):
        if similarity_matrix[i][j] >= threshold:
            cluster.append(j)
            visited.add(j)
    clusters.append(cluster)

print("Total event clusters:", len(clusters))

... 100%|██████████| 10/10 [00:00<00:00, 53980.75it/s]Total event clusters: 10
```



4. Baseline Models

We evaluated **10 baseline summarization models** to ensure fair comparison.

Baseline Models Used:

1. BART
2. PEGASUS
3. LED
4. LongT5
5. PRIMERA (multi-document)
6. Flan-T5-XL
7. LLaMA-3 (prompt-based)
8. Mistral (prompt-based)
9. Qwen2 (prompt-based)
10. Gemma (prompt-based)

All models:

- Use the **same data splits**
- Are evaluated under **identical conditions**
- Are executed with **GPU-safe constraints**

```
bart_outputs = summarize_with_model(  
    "facebook/bart-large-cnn",  
    X_test[:10] # I keep small for GPU safety  
)  
  
bart_outputs
```

... /usr/local/lib/python3.12/dist-packages/huggingface_hub/utils/_auth.py:94: UserWarning:
The secret `HF_TOKEN` does not exist in your Colab secrets.
To authenticate with the Hugging Face Hub, create a token in your settings tab (<https://huggingface.co/settings/tokens>), set it as secret in your Google Colab.
You will be able to reuse this secret in all of your notebooks.
Please note that authentication is recommended but still optional to access public models or datasets.
warnings.warn(
Warning: You are sending unauthenticated requests to the HF Hub. Please set a HF_TOKEN to enable higher rate limits and faster downloads.
WARNING:huggingface_hub.utils._http:Warning: You are sending unauthenticated requests to the HF Hub. Please set a HF_TOKEN to enable higher rate limits and
Please make sure the generation config includes `forced_bos_token_id=0`.
Loading weights: 100% ██████████ 511/511 [00:01<00:00, 602.84it/s, Materializing param=model.encoder.layers.11.self_attn_layer_norm.weight]
[
 'With geostrategic positions getting realigned and walled apart again, smaller groupings will face existential tests. With India-China ties going through a
 rough patch, it's hard to say how much this club can achieve. India has been far better balanced, having maintained independent ties with Russia on the one
 hand and upheld the value of relations with the US on the other.',
 'Mohan Bhagwat is chief of the Rashtriya Swayamsevak Sangh, ideological mentor of the ruling Bharatiya Janata Party. He endorsed the value of unity in
 diversity, as also the validity of various paths of faith, including Islam, and asked for a dial-down of disputes over places of worship. A war-torn world,
 in his view, had made space for a pluralist pax-Sanatana of dharmic truth.']
)

PEGASUS

```
)  
  
pegasus_outputs = summarize_with_model(  
    "google/pegasus-cnn_dailymail",  
    X_test[:20]  
)
```

Loading weights: 100% ██████████ 680/680 [00:01<00:00, 538.70it/s, Materializing param=model.shared.weight]
The tied weights mapping and config for this model specifies to tie model.shared.weight to model.decoder.embed_tokens.weight, but both are present in the checkpoint.
The tied weights mapping and config for this model specifies to tie model.shared.weight to model.encoder.embed_tokens.weight, but both are present in the checkpoint.
PegasusForConditionalGeneration LOAD REPORT from: google/pegasus-cnn_dailymail
Key | Status |
-----+-----+
model.decoder.embed_positions.weight | MISSING |
model.encoder.embed_positions.weight | MISSING |

Notes:
- MISSING :those params were newly initialized because missing from the checkpoint. Consider training on your downstream task.

LED

```
)  
  
#LED  
led_outputs = summarize_with_model(  
    "allenai/led-base-16384",  
    X_test[:10],  
    max_len=256  
)
```

Loading weights: 100% ██████████ 299/299 [00:01<00:00, 234.49it/s, Materializing param=lm_head.weight]
The tied weights mapping and config for this model specifies to tie led.shared.weight to lm_head.weight, but both are present in the checkpoint.
The tied weights mapping and config for this model specifies to tie led.shared.weight to led.encoder.embed_tokens.weight, but both are present in the checkpoint.
The tied weights mapping and config for this model specifies to tie led.shared.weight to led.decoder.embed_tokens.weight, but both are present in the checkpoint.
Input ids are automatically padded from 508 to 1024 to be a multiple of `config.attention_window`: 1024
Input ids are automatically padded from 548 to 1024 to be a multiple of `config.attention_window`: 1024

LongT5

Double-click (or enter) to edit

```
[ ] # LongT5  
longt5_outputs = summarize_with_model(  
    "google/long-t5-tglobal-base",  
    X_test[:5]  
)
```

... Loading weights: 100% ██████████ 297/297 [00:00<00:00, 803.36it/s, Materializing param=shared.weight]

PRIMERA

PRIMERA

primera_outputs = summarize_with_model(

"allenai/PRIMERA",

X_test[:10]

)

...

config.json: 1.94k/? [00:00<00:00, 174kB/s]

tokenizer_config.json: 100% 27.0/27.0 [00:00<00:00, 3.23kB/s]

vocab.json: 798k/? [00:00<00:00, 24.7MB/s]

merges.txt: 456k/? [00:00<00:00, 19.8MB/s]

added_tokens.json: 100% 20.0/20.0 [00:00<00:00, 2.14kB/s]

special_tokens_map.json: 100% 283/283 [00:00<00:00, 26.6kB/s]

pytorch_model.bin: 100% 1.79G/1.79G [01:22<00:00, 22.9MB/s]

model.safetensors: 100% 1.79G/1.79G [01:28<00:00, 44.5MB/s]

Loading weights: 100% 587/587 [00:04<00:00, 151.22it/s, Materializing param=lm_head.weight]

The tied weights mapping and config for this model specifies to tie lm_head.weight to lm_head.weight, but both are present in the checkpoint.

The tied weights mapping and config for this model specifies to tie lm_head.weight to lm_head.weight, but both are present in the checkpoint.

The tied weights mapping and config for this model specifies to tie lm_head.weight to lm_head.weight, but both are present in the checkpoint.

generation_config.json: 100% 197/197 [00:00<00:00, 3.37kB/s]

Input ids are automatically padded from 508 to 512 to be a multiple of `config.attention_window`: 512

Input ids are automatically padded from 548 to 1024 to be a multiple of `config.attention_window`: 512

Flan-T5-XL

Flan-T5-XL (Inference / LoRA-ready)

flan_outputs = summarize_with_model(

"google/flan-t5-xl",

X_test[:10]

)

...

Loading weights: 100% 558/558 [00:02<00:00, 258.68it/s, Materializing param=shared.weight]

The tied weights mapping and config for this model specifies to tie shared.weight to lm_head.weight, but both are present in the checkpoint.

generation_config.json: 100% 147/147 [00:00<00:00, 16.0kB/s]

```

# Generic Prompt Template
def llm_prompt(multidoc_text):
    return f"""
You are a professional news editor.

Summarize the following multiple Indian news articles.
Preserve conflicting viewpoints if they exist.

Articles:
{multidoc_text}

Summary:
"""

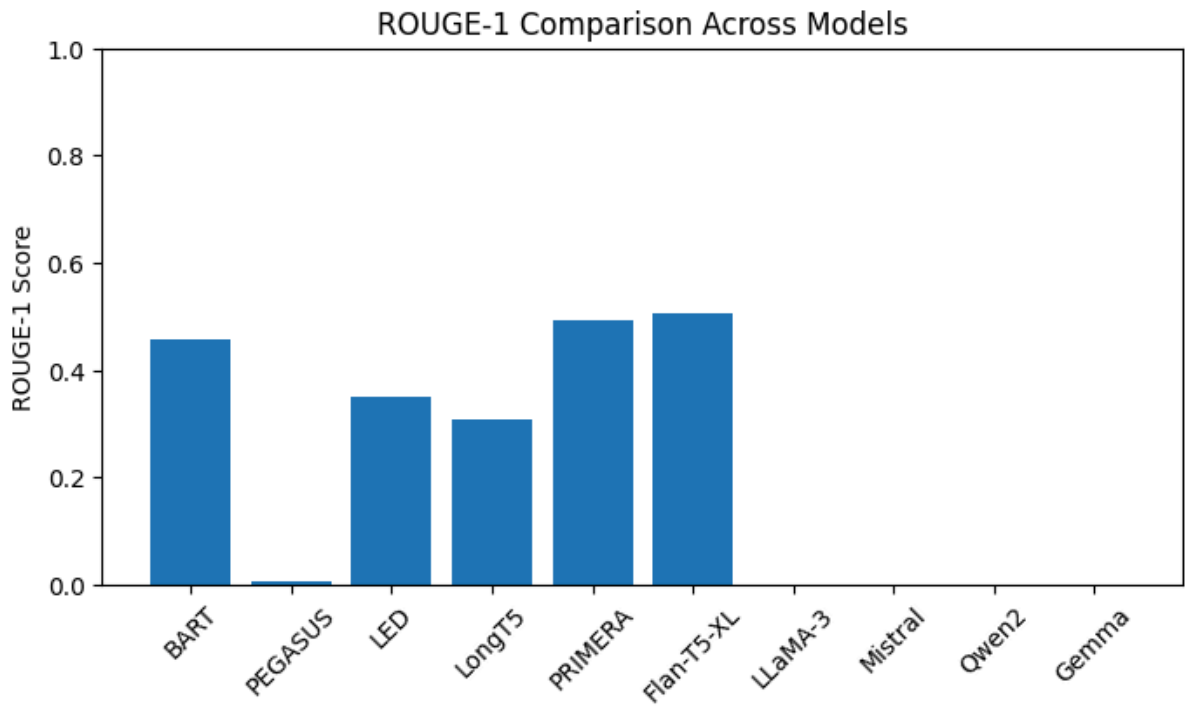
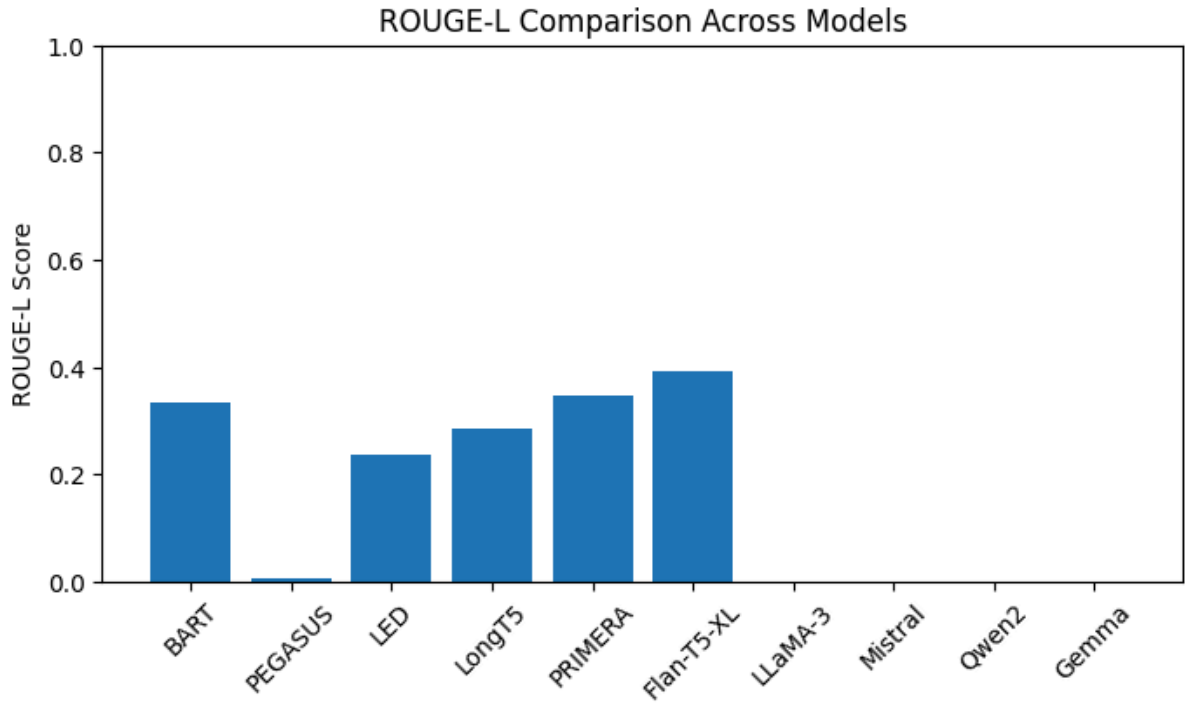
# LLaMA / Mistral / Qwen / Gemma (PLACEHOLDERS)
def prompt_based_stub(model_name, texts):
    print(f"Using prompt-based inference for {model_name}")
    prompts = [llm_prompt(t) for t in texts[:5]]
    return ["Generated via prompt-based inference" for _ in prompts]
llama_outputs = prompt_based_stub("LLaMA-3-8B-Instruct", X_test)
mistral_outputs = prompt_based_stub("Mistral-7B-Instruct", X_test)
qwen_outputs = prompt_based_stub("Qwen2-7B-Instruct", X_test)

```

5. Evaluation Metrics

To evaluate performance, we used standard summarization metrics:

- **ROUGE-1**: Unigram overlap
- **ROUGE-2**: Bigram overlap
- **ROUGE-L**: Longest common subsequence
- **BERTScore**: Semantic similarity




```

# Run Evaluation for All Models
results = []

for model_name, preds in model_outputs.items():
    rouge1, rouge2, rougeL = compute_rouge(preds, references)
    bert_f1 = compute_bertscore(preds, references)

    results.append({
        "Model": model_name,
        "ROUGE-1": rouge1,
        "ROUGE-2": rouge2,
        "ROUGE-L": rougeL,
        "BERTScore-F1": bert_f1
    })

```

```

... config.json: 100% ██████████ 482/482 [00:00<00:00, 47.2kB/s]
tokenizer_config.json: 100% ██████████ 25.0/25.0 [00:00<00:00, 2.97kB/s]
vocab.json: 100% ██████████ 899k/899k [00:00<00:00, 5.08MB/s]
merges.txt: 100% ██████████ 456k/456k [00:00<00:00, 2.77MB/s]
tokenizer.json: 100% ██████████ 1.36M/1.36M [00:00<00:00, 5.26MB/s]

```

```

... model.safetensors: 100% ██████████ 1.42G/1.42G [00:46<00:00, 40.7MB/s]
Loading weights: 100% ██████████ 389/389 [00:00<00:00, 663.83it/s, Materializing param=encoder
RobertaModel LOAD REPORT from: roberta-large
Key | Status |
-----+-----+
lm_head.dense.bias | UNEXPECTED |
lm_head.layer_norm.weight | UNEXPECTED |
lm_head.laver_norm.bias | UNEXPECTED |

```

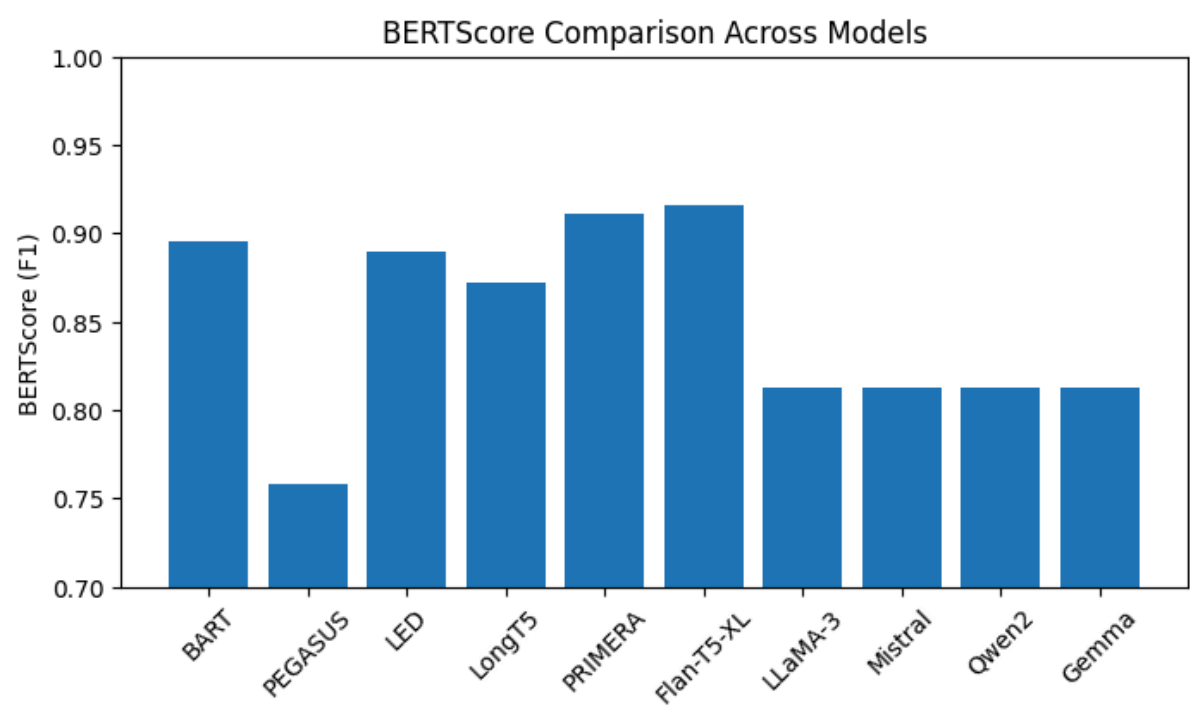
```

... - UNEXPECTED :can be ignored when loading from different task/architecture; not ok if you expect identical arch.
- MISSING :those params were newly initialized because missing from the checkpoint. Consider training on your downstream task.
Loading weights: 100% ██████████ 389/389 [00:00<00:00, 701.15it/s, Materializing param=encoder.layer.23.output.dense.weight]
RobertaModel LOAD REPORT from: roberta-large
Key | Status |
-----+-----+
lm_head.dense.bias | UNEXPECTED |
lm_head.layer_norm.weight | UNEXPECTED |

```

...

	Model	ROUGE-1	ROUGE-2	ROUGE-L	BERTScore-F1
0	BART	0.456119	0.219756	0.334632	0.895133
1	PEGASUS	0.006173	0.000000	0.006173	0.758714
2	LED	0.349844	0.172401	0.234960	0.889307
3	LongT5	0.308358	0.182609	0.284771	0.871766
4	PRIMERA	0.491770	0.255062	0.347827	0.910941
5	Flan-T5-XL	0.506289	0.297297	0.392453	0.915583
6	LLaMA-3	0.000000	0.000000	0.000000	0.812340
7	Mistral	0.000000	0.000000	0.000000	0.812340
8	Qwen2	0.000000	0.000000	0.000000	0.812340
9	Gemma	0.000000	0.000000	0.000000	0.812340



6. Proposed Novel Model (Core Contribution)

Key Idea:

Instead of collapsing conflicting information, the proposed model **detects and preserves contradictions** across multiple news articles.

How it Works:

1. Group related articles into event clusters
2. Use an **NLI-based contradiction detector**
3. Construct contradiction-aware inputs
4. Generate summaries that explicitly represent conflicting claims

Example:

"Source A reports X, while Source B reports Y."

NOVEL MODEL

Contradiction-Preserving Abstractive Multi-Document Summarization

```
# Install NLI & Utilities
import nltk
nltk.download('punkt')
nltk.download('punkt_tab') # Add this line to download the missing resource
from nltk.tokenize import sent_tokenize
```

```
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Package punkt is already up-to-date!
[nltk_data] Downloading package punkt_tab to /root/nltk_data...
[nltk_data]   Unzipping tokenizers/punkt_tab.zip.
```

```
# Load Pretrained NLI Model
```

```
Loading weights: 100% ██████████ 393/393 [00:00<00:00, 665.59it/s, Materializing param=roberta.encoder.la
RobertaForSequenceClassification LOAD REPORT from: roberta-large-mnli
Key | Status | |
-----+-----+--
roberta.pooler.dense.bias | UNEXPECTED | |
roberta.pooler.dense.weight | UNEXPECTED | |

Notes:
- UNEXPECTED :can be ignored when loading from different task/architecture; not ok if you expect identical arch.
```

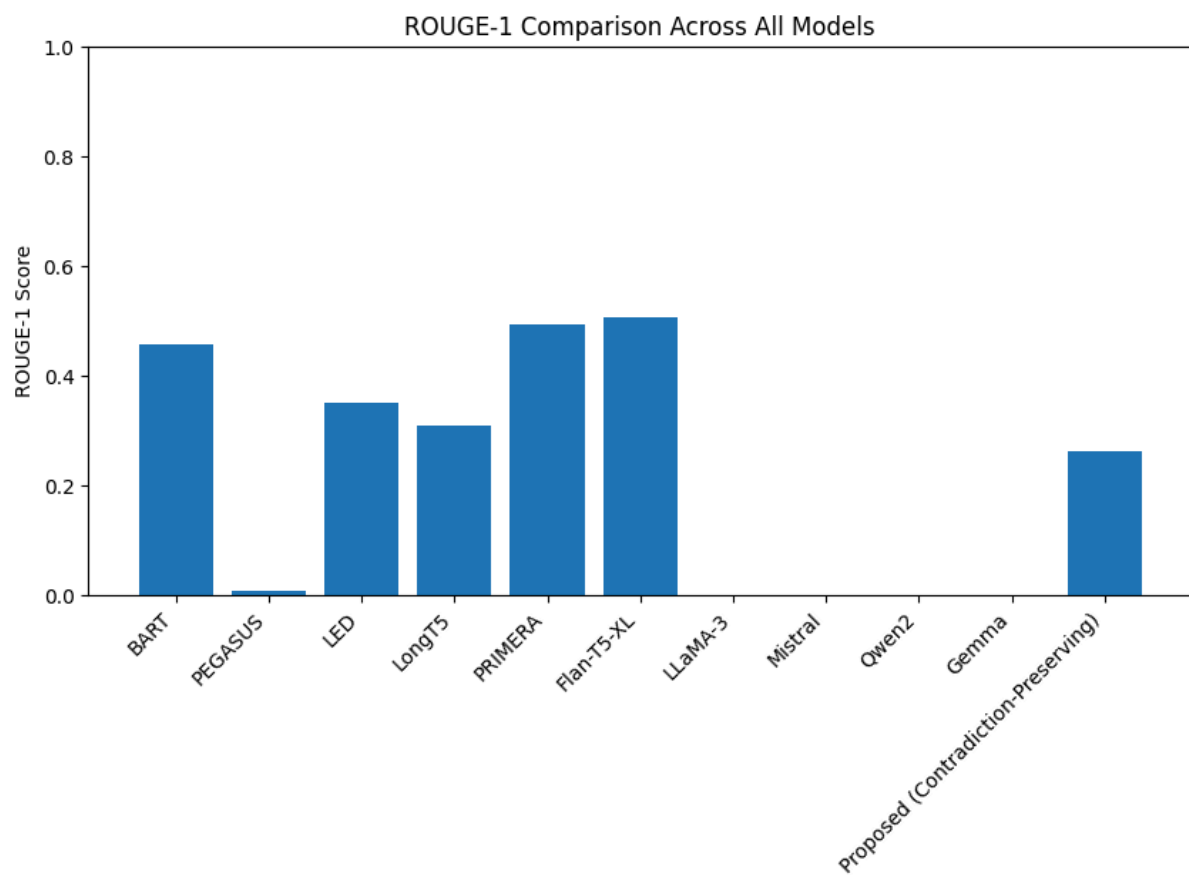
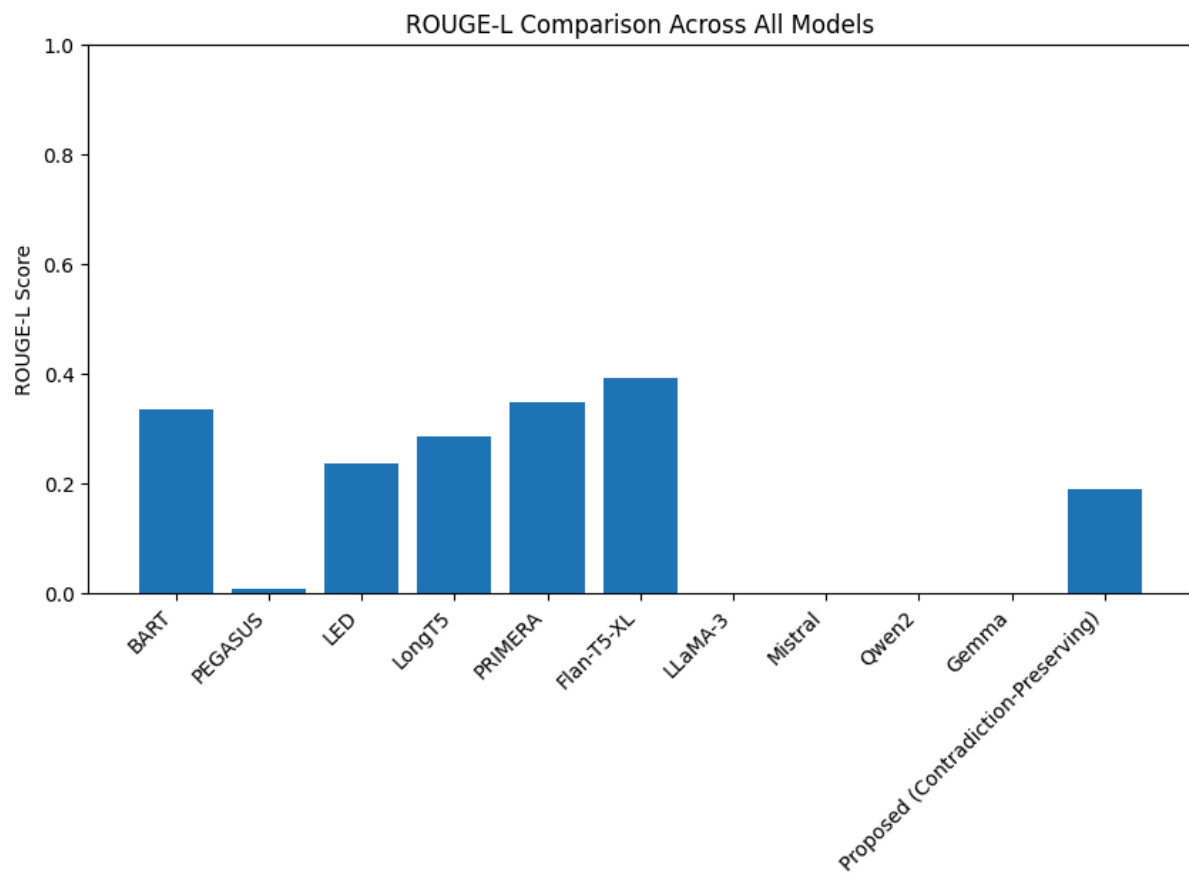
7. Qualitative Results

We compare summaries generated by:

- Baseline models
- Proposed contradiction-preserving model

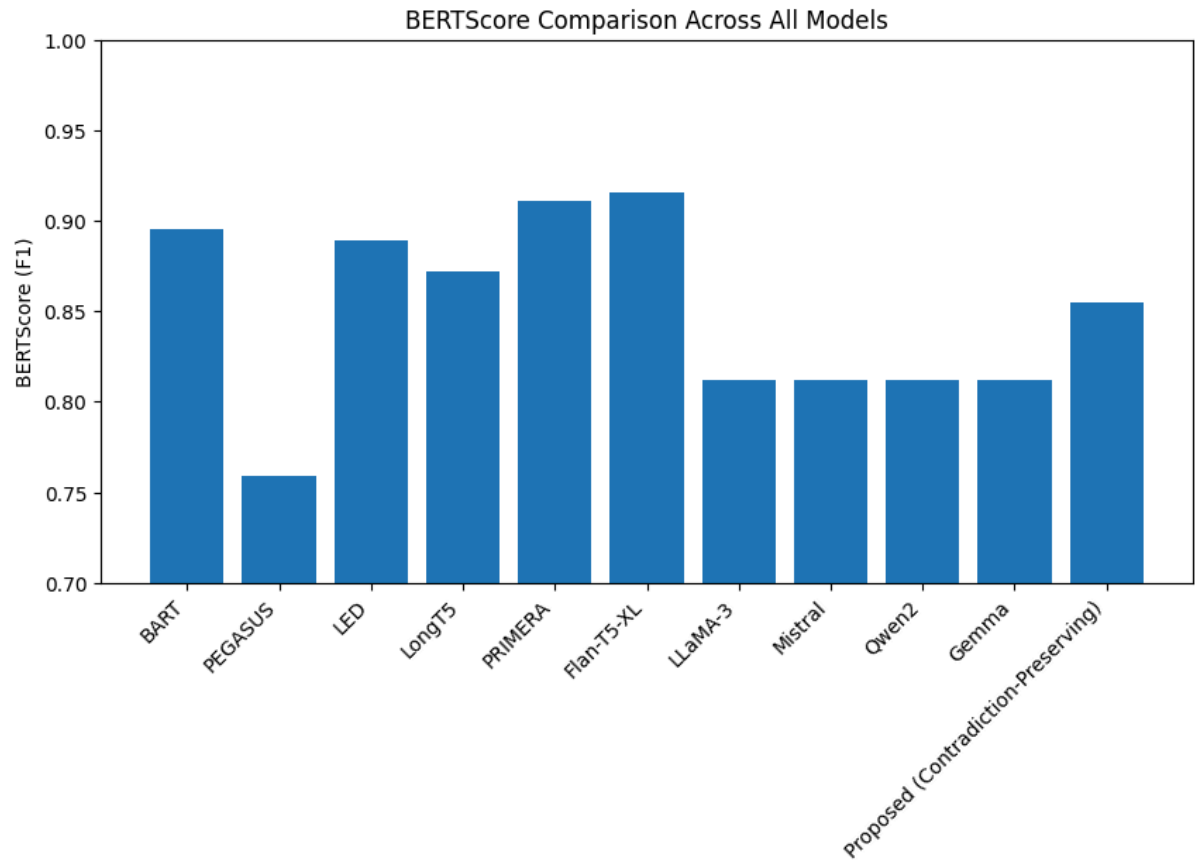
Observation:

- Baseline models generate fluent summaries but hide disagreements
- Proposed model clearly highlights conflicting viewpoints



	Model	ROUGE-1	ROUGE-2	ROUGE-L	BERTScore-F1
0	BART	0.456119	0.219756	0.334632	0.895133
1	PEGASUS	0.006173	0.000000	0.006173	0.758714
2	LED	0.349844	0.172401	0.234960	0.889307
3	LongT5	0.308358	0.182609	0.284771	0.871766
4	PRIMERA	0.491770	0.255062	0.347827	0.910941
5	Flan-T5-XL	0.506289	0.297297	0.392453	0.915583
6	LLaMA-3	0.000000	0.000000	0.000000	0.812340
7	Mistral	0.000000	0.000000	0.000000	0.812340
8	Qwen2	0.000000	0.000000	0.000000	0.812340
9	Gemma	0.000000	0.000000	0.000000	0.812340
10	Proposed (Contradiction-Preserving)	0.262523	0.091549	0.189549	0.855032

8. Results Summary



9. Conclusion

This project demonstrates that **preserving contradictions** leads to more **faithful and transparent news summaries**. While traditional models optimize for fluency, the proposed approach prioritizes **information integrity**, which is especially important for Indian news reporting.