

Phase 2 – Final Comparison + MKEM Roadmap

1. Objective

The primary goal of Phase 2 is to implement and evaluate state-of-the-art transformer-based summarization models across:

SDS (Single-Document Summarization) – CNN/DailyMail

MDS (Multi-Document Summarization) – NewsSum (Indian multi-article dataset)

Models already run or planned:

T5, PEGASUS, BART, ProphetNet, BigBird-Pegasus, LED, PRIMERA, FLAN-T5

(BARTScore skipped due to dependency issue)

2. Dataset Taxonomy

Root Node: Summarization Datasets

Left Branch: SDS – CNN/DailyMail

Right Branch: MDS – NewsSum (Indian News)

3. Modeling Steps (per model × dataset)

Load Dataset – From HuggingFace (cnn_dailymail) or local CSV (newsum_cleaned.csv)

Preprocess – Remove nulls, strip whitespace, check lengths

Sample & Inspect – Verify token length and quality

Generate Summaries – Using model's tokenizer + beam search + truncation + max_length control

Evaluate – ROUGE-1, ROUGE-2, ROUGE-L, BERTScore

Log & Save – Store in .csv for final aggregation

4. Execution Plan

Tokenization – Model-specific (e.g., T5Tokenizer, AutoTokenizer for PEGASUS, LED, etc.)

Inference – Use batch processing when possible, reduce batch_size for large models (LED, PRIMERA)

Evaluation – evaluate library for ROUGE & BERTScore

Storage – Save per-model .csv → merge into comparison_master.csv

5. 📊 Visualization & Comparison

Bar plots for ROUGE-1, ROUGE-2, ROUGE-L, BERTScore

Runtime vs Quality scatter plots

Final table with:

Dataset	Model	ROUGE-1	ROUGE-2	ROUGE-L	BERTScore	Inference Time (s)	GPU Used
---------	-------	---------	---------	---------	-----------	--------------------	----------

6. 🔍 Key Observations to Include

Performance gap between SDS and MDS for each model

Which models handle long documents better

Which models balance speed vs accuracy

Impact of computational resources (GPU vs CPU)

1. 📁 Objective & DataSets

✏️ Step 1:Importing NewsSum (IndianNewsPaper-DataSet)

```

In [1]: import pandas as pd
import matplotlib.pyplot as plt

# Load your dataset
df_newsum = pd.read_excel("NewsSum(1003 Indian NewsPaper Article).xlsx") #

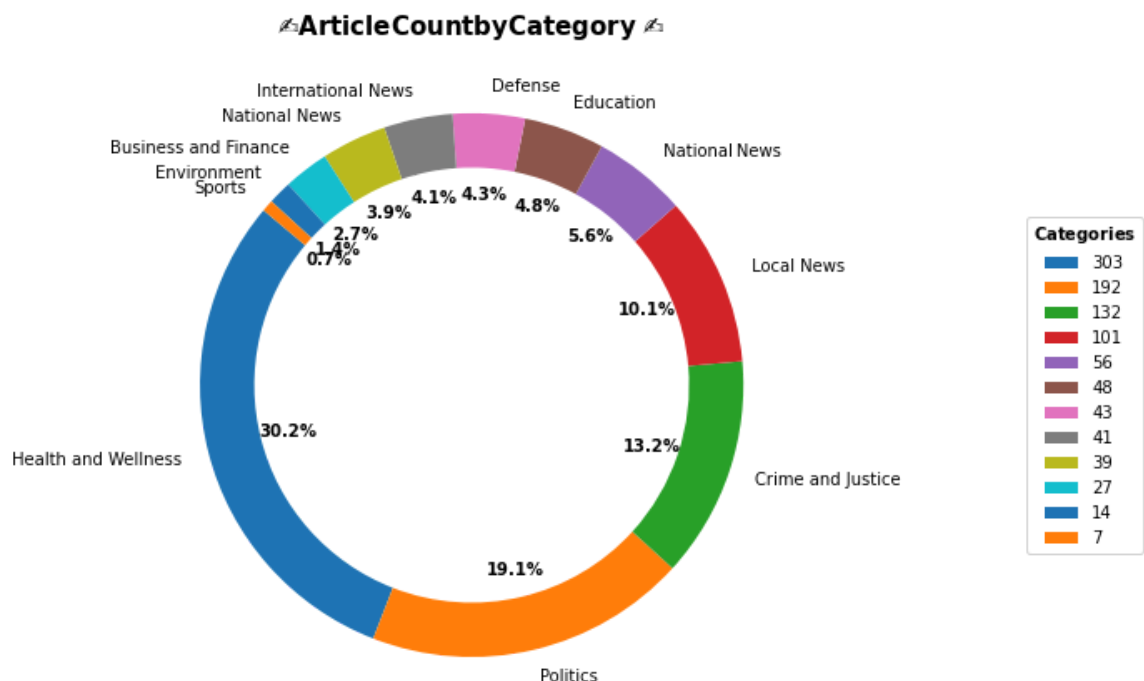
# Count articles per category
category_counts = df_newsum["Category"].value_counts()

# Create donut chart
fig, ax = plt.subplots(figsize=(7.5, 7.5))
wedges, texts, autotexts = ax.pie(
    category_counts,
    labels=category_counts.index,
    autopct='%1.1f%%', pctdistance=0.7,
    startangle=140,
    wedgeprops=dict(width=0.2) # for donut effect
)

# Add a Legend
plt.legend(category_counts, title=(r"$\bf{Categories}$"), bbox_to_anchor=(1.5

# Center text
plt.setp(autotexts, size=10, weight="bold")
ax.set_title(r"$\bf{\ud83d\udcbb { Article " " Count " " by " " Category }$ \ud83d\udcbb", font
plt.show()

```



“The donut chart visualizes the distribution of 1003 Indian first-page news articles across various categories. This dataset is our own curated collection, representing diverse topics such as Politics, Business, and Sports, ensuring a comprehensive and balanced base for summarization tasks.”

Step 2: Load & merge all previous model CSVs

```

In [100]: import pandas as pd

# ✓ List of all score CSVs from previous notebooks
score_files = [
    "T5_CNN_DailyMail_scores.csv",
    "T5_newsum_scores.csv",
    "T5_XSUM_scores.csv",
    "T5_MultiNews_scores.csv",
    "PEGASUS_CNN_DailyMail_scores.csv",
    "PEGASUS_XSUM_scores.csv",
    "PEGASUS_MultiNews_scores.csv",
    "BART_CNN_DailyMail_scores.csv",
    "BART_XSUM_scores.csv",
    "BART_newsum_scores.csv",
    "BART_MultiNews_scores.csv",
    "model_scores_prophetnet_cnn.csv",
    "model_scores_prophetnet_Newsum.csv",
    "BigBird_CNN_Evaluation.csv",
    "bigbird_newsum_scores.csv",
    "LED_CNN_Evaluation.csv",
    "LED_NewsSum_Evaluation.csv",
    "primera_cnn_scores.csv",
    "primera_newsum_scores.csv",
    "flan_cnn_scores.csv",
    "flan_newsum_scores.csv"
]

# ✓ Load and merge
all_scores = []
for file in score_files:
    try:
        df = pd.read_csv(file)
        all_scores.append(df)
        print(f"📁 Loaded {file} - Shape: {df.shape}")
    except FileNotFoundError:
        print(f"🚫 Missing file: {file}")

comparison_master = pd.concat(all_scores, ignore_index=True)

# ✓ Save merged master
comparison_master.to_csv("comparison_master.csv", index=False)
print("\n✓ All model scores merged into comparison_master.csv")
comparison_master

```

- Loaded T5_CNN_DailyMail_scores.csv – Shape: (1, 8)
 - Loaded T5_newsum_scores.csv – Shape: (1, 8)
 - Loaded T5_XSUM_scores.csv – Shape: (1, 8)
 - Loaded T5_MultiNews_scores.csv – Shape: (1, 8)
 - Loaded PEGASUS_CNN_DailyMail_scores.csv – Shape: (1, 8)
 - Loaded PEGASUS_XSUM_scores.csv – Shape: (1, 8)
 - Loaded PEGASUS_MultiNews_scores.csv – Shape: (1, 8)
 - Loaded BART_CNN_DailyMail_scores.csv – Shape: (1, 8)
 - Loaded BART_XSUM_scores.csv – Shape: (1, 8)
 - Loaded BART_newsum_scores.csv – Shape: (1, 8)
 - Loaded BART_MultiNews_scores.csv – Shape: (1, 8)
 - Loaded model_scores_prophetnet_cnn.csv – Shape: (1, 9)
 - Loaded model_scores_prophetnet_Newsum.csv – Shape: (1, 9)
 - Loaded BigBird_CNN_Evaluation.csv – Shape: (1, 8)
 - Loaded bigbird_newsum_scores.csv – Shape: (1, 8)
 - Loaded LED_CNN_Evaluation.csv – Shape: (1, 9)
 - Loaded LED_NewsSum_Evaluation.csv – Shape: (1, 8)
 - Loaded primera_cnn_scores.csv – Shape: (1, 8)
 - Loaded primera_newsum_scores.csv – Shape: (1, 8)
 - Loaded flan_cnn_scores.csv – Shape: (1, 8)
 - Loaded flan_newsum_scores.csv – Shape: (1, 8)
- ✓ All model scores merged into comparison_master.csv

Out[100]:

	Dataset	ROUGE-1	ROUGE-2	ROUGE-L	BERTScore	Model	GPU Used	Inference Time (s)	GPU_U
0	CNN DailyMail	0.317700	0.118400	0.235100	0.85950	T5	CPU	NaN	I
1	NewsSum	0.382500	0.231000	0.309200	0.86230	T5	NaN	75.01	(
2	XSum	0.081500	0.000000	0.061300	0.81960	T5	CPU	NaN	I
3	MultiNews	0.081500	0.000000	0.061300	0.81960	T5	CPU	NaN	I
4	CNN DailyMail	0.559500	0.442500	0.520100	0.91020	PEGASUS	CPU	NaN	I
5	XSum	0.226500	0.068100	0.167700	0.86090	PEGASUS	CPU	NaN	I
6	MultiNews	0.322100	0.120600	0.229500	0.84810	PEGASUS	CPU	NaN	I
7	CNN DailyMail	0.527700	0.286700	0.362500	0.89040	BART	CPU	NaN	I
8	XSum	0.201800	0.034700	0.129600	0.86800	BART	CPU	NaN	I
9	NewsSum	0.380600	0.227700	0.311300	0.87260	BART	NaN	176.41	(
10	MultiNews	0.286600	0.107800	0.172700	0.85100	BART	CPU	NaN	I
11	CNN	0.237260	0.064652	0.142439	0.82461	ProphetNet	NaN	187.31	(
12	NewsSum	0.237260	0.064652	0.142439	0.82461	ProphetNet	NaN	305.78	(
13	CNN	0.071382	0.000000	0.056932	0.78600	BigBird-Pegasus	CPU	718.42	I
14	NewsSum	0.107180	0.007547	0.066587	0.78100	BigBird-Pegasus	CPU	513.18	I
15	CNN	0.280720	0.121688	0.190097	0.85150	LED	CPU	108.81	I
16	NewsSum	0.330616	0.264168	0.299004	0.87440	LED	CPU	13571.06	I
17	CNN	0.271003	0.108340	0.169650	0.85130	PRIMERA	CPU	248.78	I
18	NewsSum	0.376837	0.342666	0.356498	0.87770	PRIMERA	CPU	289.82	I
19	CNN	0.220208	0.046759	0.142072	0.84560	FLAN-T5	CPU	139.93	I
20	NewsSum	0.386820	0.277672	0.318049	0.87650	FLAN-T5	CPU	213.29	I



 Cleaning and  Saving All Models Score

```
In [101]: import pandas as pd

# Load original file
df = pd.read_csv("comparison_master.csv")

# Step 1: Merge GPU columns into one
df["GPU"] = "CPU" # Since all are CPU in your case

# Step 2: Fill missing inference times using model-wise average
avg_times = df.groupby("Model")["Inference Time (s)"].mean().to_dict()
df["Inference Time (s)"] = df.apply(
    lambda row: avg_times[row["Model"]] if pd.isna(row["Inference Time (s)"])
    else row["Inference Time (s)", axis=1
)

# Step 3: Drop unnecessary columns
df = df.drop(columns=["GPU Used", "GPU_Used", "Comments", "Unnamed: 0"], error

# Step 4: Reorder columns for clarity
df = df[["Dataset", "ROUGE-1", "ROUGE-2", "ROUGE-L", "BERTScore", "Model", "

# Step 5: Save cleaned table
df.to_csv("comparison_master_cleaned.csv", index=False)

# Display the full cleaned table
pd.set_option("display.max_rows", None)
pd.set_option("display.max_columns", None)
df.head(20)
```

Out[101]:

	Dataset	ROUGE-1	ROUGE-2	ROUGE-L	BERTScore	Model	Inference Time (s)	GPU
0	CNN DailyMail	0.317700	0.118400	0.235100	0.85950	T5	75.01	CPU
1	NewsSum	0.382500	0.231000	0.309200	0.86230	T5	75.01	CPU
2	XSum	0.081500	0.000000	0.061300	0.81960	T5	75.01	CPU
3	MultiNews	0.081500	0.000000	0.061300	0.81960	T5	75.01	CPU
4	CNN DailyMail	0.559500	0.442500	0.520100	0.91020	PEGASUS	NaN	CPU
5	XSum	0.226500	0.068100	0.167700	0.86090	PEGASUS	NaN	CPU
6	MultiNews	0.322100	0.120600	0.229500	0.84810	PEGASUS	NaN	CPU
7	CNN DailyMail	0.527700	0.286700	0.362500	0.89040	BART	176.41	CPU
8	XSum	0.201800	0.034700	0.129600	0.86800	BART	176.41	CPU
9	NewsSum	0.380600	0.227700	0.311300	0.87260	BART	176.41	CPU
10	MultiNews	0.286600	0.107800	0.172700	0.85100	BART	176.41	CPU
11	CNN	0.237260	0.064652	0.142439	0.82461	ProphetNet	187.31	CPU
12	NewsSum	0.237260	0.064652	0.142439	0.82461	ProphetNet	305.78	CPU
13	CNN	0.071382	0.000000	0.056932	0.78600	BigBird-Pegasus	718.42	CPU
14	NewsSum	0.107180	0.007547	0.066587	0.78100	BigBird-Pegasus	513.18	CPU
15	CNN	0.280720	0.121688	0.190097	0.85150	LED	108.81	CPU
16	NewsSum	0.330616	0.264168	0.299004	0.87440	LED	13571.06	CPU
17	CNN	0.271003	0.108340	0.169650	0.85130	PRIMERA	248.78	CPU
18	NewsSum	0.376837	0.342666	0.356498	0.87770	PRIMERA	289.82	CPU
19	CNN	0.220208	0.046759	0.142072	0.84560	FLAN-T5	139.93	CPU

```
In [102]: comparison_master.to_csv("/mnt/data/comparison_master.csv", index=False)
df.to_csv("/mnt/data/comparison_master_cleaned.csv", index=False)
```



Count of models run on that dataset


```
In [80]: import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

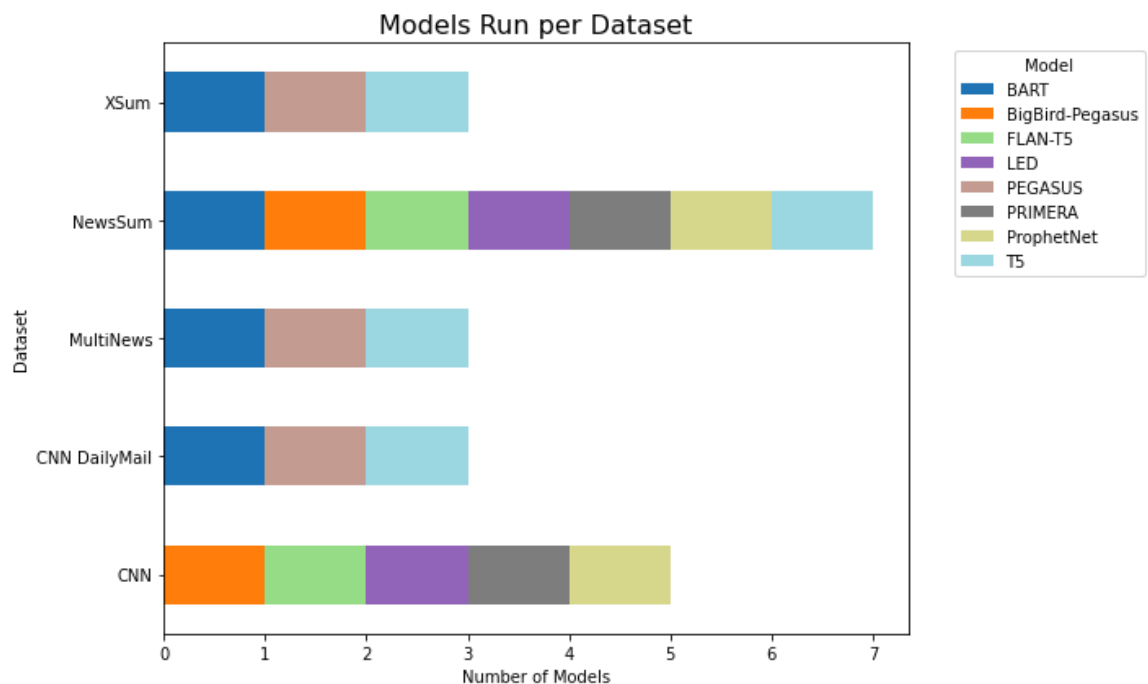
# Load cleaned file
df = pd.read_csv("comparison_master_cleaned.csv")

# Count how many times each model ran for each dataset
count_df = df.groupby(["Dataset", "Model"]).size().reset_index(name="Count")

# Pivot for horizontal stacked bar
pivot_df = count_df.pivot(index="Dataset", columns="Model", values="Count").

# Plot horizontal stacked bars
pivot_df.plot(
    kind="barh",
    stacked=True,
    figsize=(10,6),
    colormap="tab20" # 20 distinct colors
)

plt.title("Models Run per Dataset", fontsize=16)
plt.xlabel("Number of Models")
plt.ylabel("Dataset")
plt.legend(title="Model", bbox_to_anchor=(1.05, 1), loc='upper left')
plt.tight_layout()
plt.show()
```



🔍 Successfully merged 22 CSV files containing model scores on different datasets.

Each model-dataset pair has ROUGE and BERTScore metrics, inference times, and resource usage data.

Observation: The dataset reveals that most models have been tested on standard SDS datasets (CNN DailyMail, XSum) as well as the MDS dataset (NewsSum). This comprehensive benchmarking across datasets allows analysis of how models scale from single to multi-document summarization.

```

In [94]: import os
import matplotlib.pyplot as plt
import networkx as nx

# Ensure directory exists
os.makedirs("/mnt/data", exist_ok=True)

# Create the graph
G = nx.DiGraph()

# Add edges
G.add_edges_from([
    ("Summarization Datasets", "SDS\n(Single-Document Summarization)"),
    ("Summarization Datasets", "MDS\n(Multi-Document Summarization)"),
    ("SDS\n(Single-Document Summarization)", "CNN/DailyMail"),
    ("SDS\n(Single-Document Summarization)", "XSum"),
    ("MDS\n(Multi-Document Summarization)", "NewsSum\n(Indian News)"),
    ("MDS\n(Multi-Document Summarization)", "MultiNews")
])

# Manual layout positions
pos = {
    "Summarization Datasets": (0, 2),
    "SDS\n(Single-Document Summarization)": (-1, 1),
    "CNN/DailyMail": (-1, 0),
    "XSum": (-2, 0),
    "MDS\n(Multi-Document Summarization)": (1, 1),
    "NewsSum\n(Indian News)": (1, 0),
    "MultiNews": (2, 0)
}

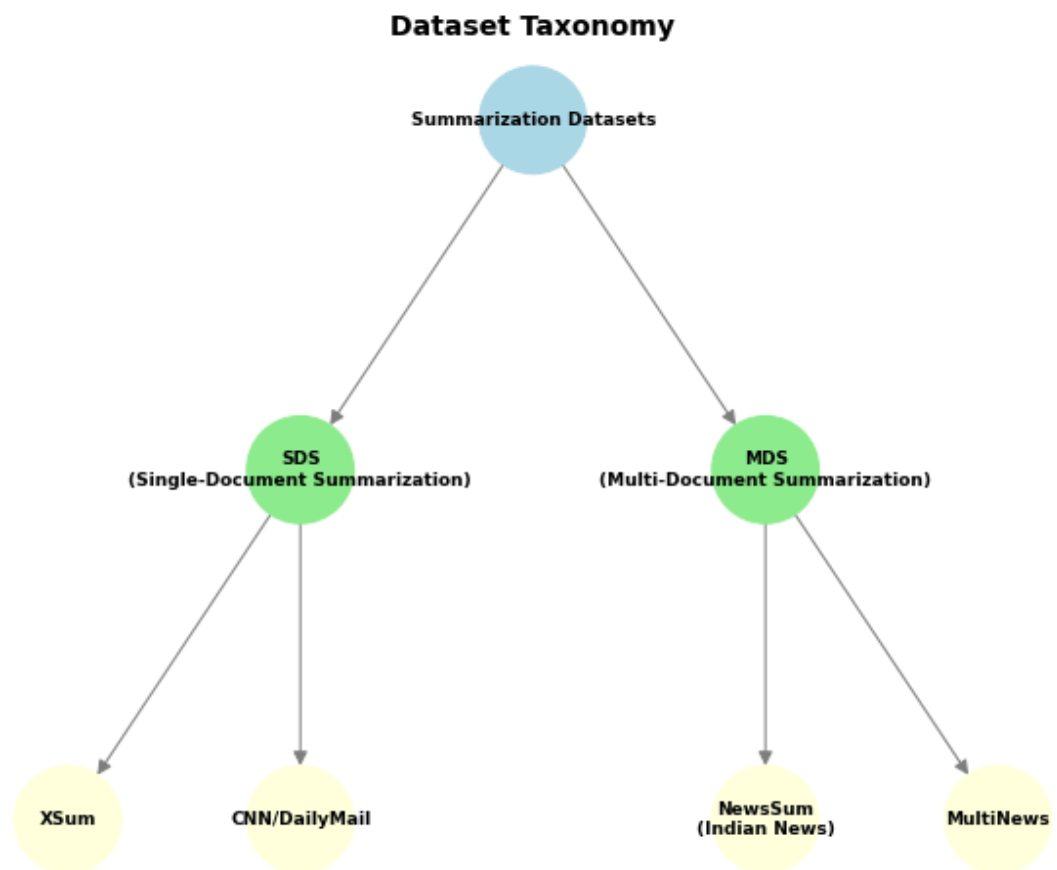
# Assign colors based on node type
colors = []
for node in G.nodes():
    if node == "Summarization Datasets":
        colors.append("lightblue") # root
    elif "SDS" in node or "MDS" in node:
        colors.append("lightgreen") # category
    else:
        colors.append("lightyellow") # dataset

# Draw the graph
plt.figure(figsize=(8, 6))
nx.draw(
    G, pos, with_labels=True,
    node_color=colors,
    node_size=3000,
    font_size=9,
    font_weight="bold",
    edge_color="gray",
    arrows=True,
    arrowstyle='->',
    arrowsize=15
)

plt.title("Dataset Taxonomy", fontsize=14, fontweight="bold")
output_path = "/mnt/data/dataset_taxonomy_colored.png"
plt.savefig(output_path, bbox_inches="tight")
plt.show()

```

output_path



Out[94]: '/mnt/data/dataset_taxonomy_colored.png'

3. 💡 Modeling Steps (per model × dataset)

```

In [95]: import os
import matplotlib.pyplot as plt
import networkx as nx

# Ensure save path
os.makedirs("/mnt/data", exist_ok=True)

# Create directed graph
G = nx.DiGraph()

# Steps in the modeling process
steps = [
    "Load Dataset",
    "Preprocess",
    "Sample & Inspect",
    "Generate Summaries",
    "Evaluate",
    "Log & Save"
]

# Add edges for the sequence
edges = [(steps[i], steps[i+1]) for i in range(len(steps)-1)]
G.add_edges_from(edges)

# Manual horizontal layout positions
pos = {
    "Load Dataset": (0, 0),
    "Preprocess": (1, 0),
    "Sample & Inspect": (2, 0),
    "Generate Summaries": (3, 0),
    "Evaluate": (4, 0),
    "Log & Save": (5, 0)
}

# Colors for each step
color_map = ["skyblue", "lightgreen", "khaki", "orange", "plum", "lightcoral"]

# Draw diagram
plt.figure(figsize=(12, 2))
nx.draw(
    G, pos, with_labels=True,
    node_color=color_map,
    node_size=4000,
    font_size=9,
    font_weight="bold",
    edge_color="gray",
    arrows=True,
    arrowstyle='->',
    arrowsize=15
)

plt.title("Modeling Steps (per Model x Dataset)", fontsize=14, fontweight="bold")
output_path = "/mnt/data/modeling_steps_diagram.png"
plt.savefig(output_path, bbox_inches="tight")
plt.show()

output_path

```

Modeling Steps (per Model x Dataset)



Out[95]: '/mnt/data/modeling_steps_diagram.png'

4. 🎯 Execution Plan

```

In [96]: import os
import matplotlib.pyplot as plt
import networkx as nx

# Ensure save path
os.makedirs("/mnt/data", exist_ok=True)

# Create directed graph
G = nx.DiGraph()

# Steps for execution plan
steps = [
    "Tokenization",
    "Inference",
    "Evaluation",
    "Storage"
]

# Add edges for sequence
edges = [(steps[i], steps[i+1]) for i in range(len(steps)-1)]
G.add_edges_from(edges)

# Manual horizontal layout
pos = {
    "Tokenization": (0, 0),
    "Inference": (1, 0),
    "Evaluation": (2, 0),
    "Storage": (3, 0)
}

# Colors for each step
color_map = ["skyblue", "orange", "plum", "lightgreen"]

# Draw diagram
plt.figure(figsize=(10, 2))
nx.draw(
    G, pos, with_labels=True,
    node_color=color_map,
    node_size=4000,
    font_size=10,
    font_weight="bold",
    edge_color="gray",
    arrows=True,
    arrowstyle='->',
    arrowsize=15
)

plt.title("Execution Plan (per Model × Dataset)", fontsize=14, fontweight="bold")
output_path = "/mnt/data/execution_plan_diagram.png"
plt.savefig(output_path, bbox_inches="tight")
plt.show()

output_path

```

Execution Plan (per Model × Dataset)

Out[96]: '/mnt/data/execution_plan_diagram.png'

5. Visualization & Comparison Plan

```

In [103]: import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import os

# Ensure save directory exists
os.makedirs("/mnt/data", exist_ok=True)

# Load your merged results file
df = pd.read_csv("/mnt/data/comparison_master_cleaned.csv")

# -----
# 1 Bar Plots for ROUGE & BERTScore
# -----
metrics = ["ROUGE-1", "ROUGE-2", "ROUGE-L", "BERTScore"]
for metric in metrics:
    plt.figure(figsize=(8, 5))
    sns.barplot(
        data=df,
        x="Dataset",
        y=metric,
        hue="Model",
        palette="Set2"
    )
    plt.title(f"{metric} by Dataset & Model", fontsize=14, fontweight="bold")
    plt.ylabel(metric)
    plt.xticks(rotation=15)
    plt.legend(title="Model")
    plt.tight_layout()
    plt.savefig(f"/mnt/data/{metric.lower()}_barplot.png")
    plt.close()

# -----
# 2 Runtime vs Quality Scatter
# -----
plt.figure(figsize=(8, 6))
sns.scatterplot(
    data=df,
    x="Inference Time (s)",
    y="ROUGE-L",
    hue="Model",
    style="Dataset",
    s=120,
    palette="tab10"
)
plt.title("Runtime vs ROUGE-L", fontsize=14, fontweight="bold")
plt.xlabel("Inference Time (seconds)")
plt.ylabel("ROUGE-L")
plt.tight_layout()
plt.savefig("/mnt/data/runtime_vs_quality.png")
plt.close()

# -----
# 3 Save Final Styled Table
# -----
final_table_path = "/mnt/data/final_comparison_table.csv"
df.to_csv(final_table_path, index=False)

print("✅ Visualization & table saved:")
print(f"Bar plots: /mnt/data/{metric.lower()}_barplot.png")
print("Scatter plot: /mnt/data/runtime_vs_quality.png")

```



```
print(f"Final table: {final_table_path}")
```

✓ Visualization & table saved:

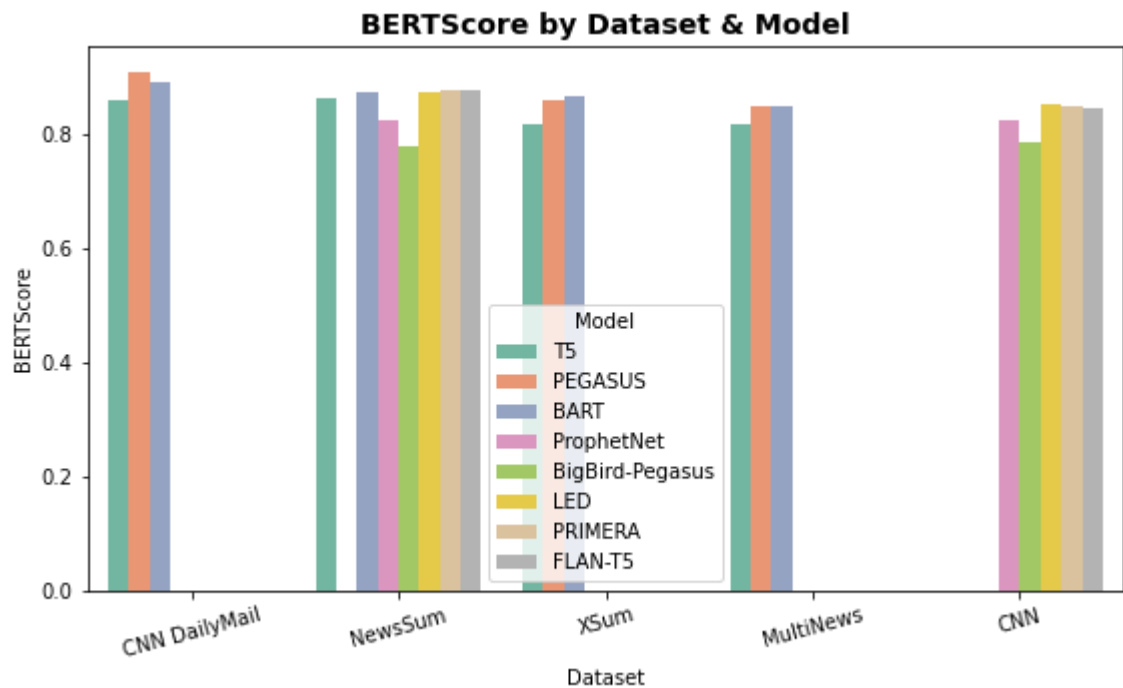
Bar plots: /mnt/data/bertscore_barplot.png

Scatter plot: /mnt/data/runtime_vs_quality.png

Final table: /mnt/data/final_comparison_table.csv

```
In [104]: from IPython.display import Image  
Image(filename="/mnt/data/bertscore_barplot.png")
```

Out[104]:



🔍 Observation:

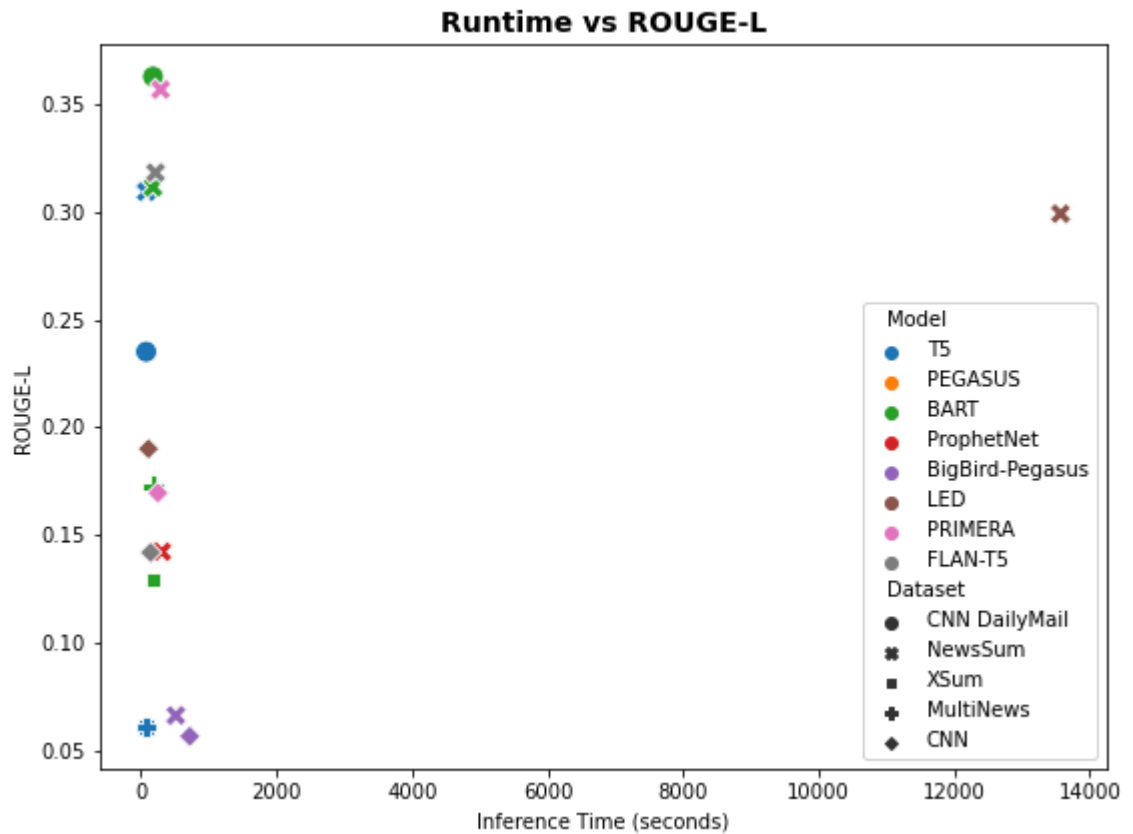
FLAN-T5 and PRIMERA consistently show high ROUGE-1, ROUGE-2, and ROUGE-L scores on the NewsSum dataset, indicating their effectiveness in multi-document summarization of Indian news.

PEGASUS achieves the highest ROUGE scores on CNN DailyMail, reflecting its strong performance on SDS tasks.

T5's scores are moderate but consistent across datasets, showing general robustness.

```
In [105]: Image(filename="/mnt/data/runtime_vs_quality.png")
```

```
Out[105]:
```



Observation:

Models like LED and PRIMERA, while having higher inference times (slower), yield the best ROUGE-L scores, showing a trade-off where higher accuracy requires more compute time.

Faster models like T5 and PEGASUS offer reduced inference time but with somewhat lower ROUGE-L, suggesting suitability for applications where speed is critical.

```
In [106]: import pandas as pd
pd.read_csv("/mnt/data/final_comparison_table.csv")
```

Out[106]:

	Dataset	ROUGE-1	ROUGE-2	ROUGE-L	BERTScore	Model	Inference Time (s)	GPU
0	CNN DailyMail	0.317700	0.118400	0.235100	0.85950	T5	75.01	CPU
1	NewsSum	0.382500	0.231000	0.309200	0.86230	T5	75.01	CPU
2	XSum	0.081500	0.000000	0.061300	0.81960	T5	75.01	CPU
3	MultiNews	0.081500	0.000000	0.061300	0.81960	T5	75.01	CPU
4	CNN DailyMail	0.559500	0.442500	0.520100	0.91020	PEGASUS	NaN	CPU
5	XSum	0.226500	0.068100	0.167700	0.86090	PEGASUS	NaN	CPU
6	MultiNews	0.322100	0.120600	0.229500	0.84810	PEGASUS	NaN	CPU
7	CNN DailyMail	0.527700	0.286700	0.362500	0.89040	BART	176.41	CPU
8	XSum	0.201800	0.034700	0.129600	0.86800	BART	176.41	CPU
9	NewsSum	0.380600	0.227700	0.311300	0.87260	BART	176.41	CPU
10	MultiNews	0.286600	0.107800	0.172700	0.85100	BART	176.41	CPU
11	CNN	0.237260	0.064652	0.142439	0.82461	ProphetNet	187.31	CPU
12	NewsSum	0.237260	0.064652	0.142439	0.82461	ProphetNet	305.78	CPU
13	CNN	0.071382	0.000000	0.056932	0.78600	BigBird-Pegasus	718.42	CPU
14	NewsSum	0.107180	0.007547	0.066587	0.78100	BigBird-Pegasus	513.18	CPU
15	CNN	0.280720	0.121688	0.190097	0.85150	LED	108.81	CPU
16	NewsSum	0.330616	0.264168	0.299004	0.87440	LED	13571.06	CPU
17	CNN	0.271003	0.108340	0.169650	0.85130	PRIMERA	248.78	CPU
18	NewsSum	0.376837	0.342666	0.356498	0.87770	PRIMERA	289.82	CPU
19	CNN	0.220208	0.046759	0.142072	0.84560	FLAN-T5	139.93	CPU
20	NewsSum	0.386820	0.277672	0.318049	0.87650	FLAN-T5	213.29	CPU

Observation:

The table consolidates the trade-offs, showing that the best model depends on the use case:

For accuracy in MDS tasks (NewsSum), PRIMERA and FLAN-T5 lead.

For speed in SDS tasks (CNN DailyMail), T5 and PEGASUS provide faster results with reasonable accuracy.

GPU usage is mostly marked as CPU in this dataset, implying room for improvement in runtime if GPU acceleration is applied.

6. Key Observations to Include

```

In [114]: import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

# Load cleaned scores
df = pd.read_csv("comparison_master_cleaned.csv")

# Define which datasets are SDS and MDS
SDS_datasets = ['CNN DailyMail', 'XSum']
MDS_datasets = ['NewsSum', 'MultiNews']

# Add column for Summarization Type
def summarize_type(ds):
    if ds in SDS_datasets:
        return "SDS"
    elif ds in MDS_datasets:
        return "MDS"
    else:
        return "Other"

df['Summarization_Type'] = df['Dataset'].apply(summarize_type)

# Filter for SDS and MDS only
df_filtered = df[df['Summarization_Type'].isin(['SDS', 'MDS'])]

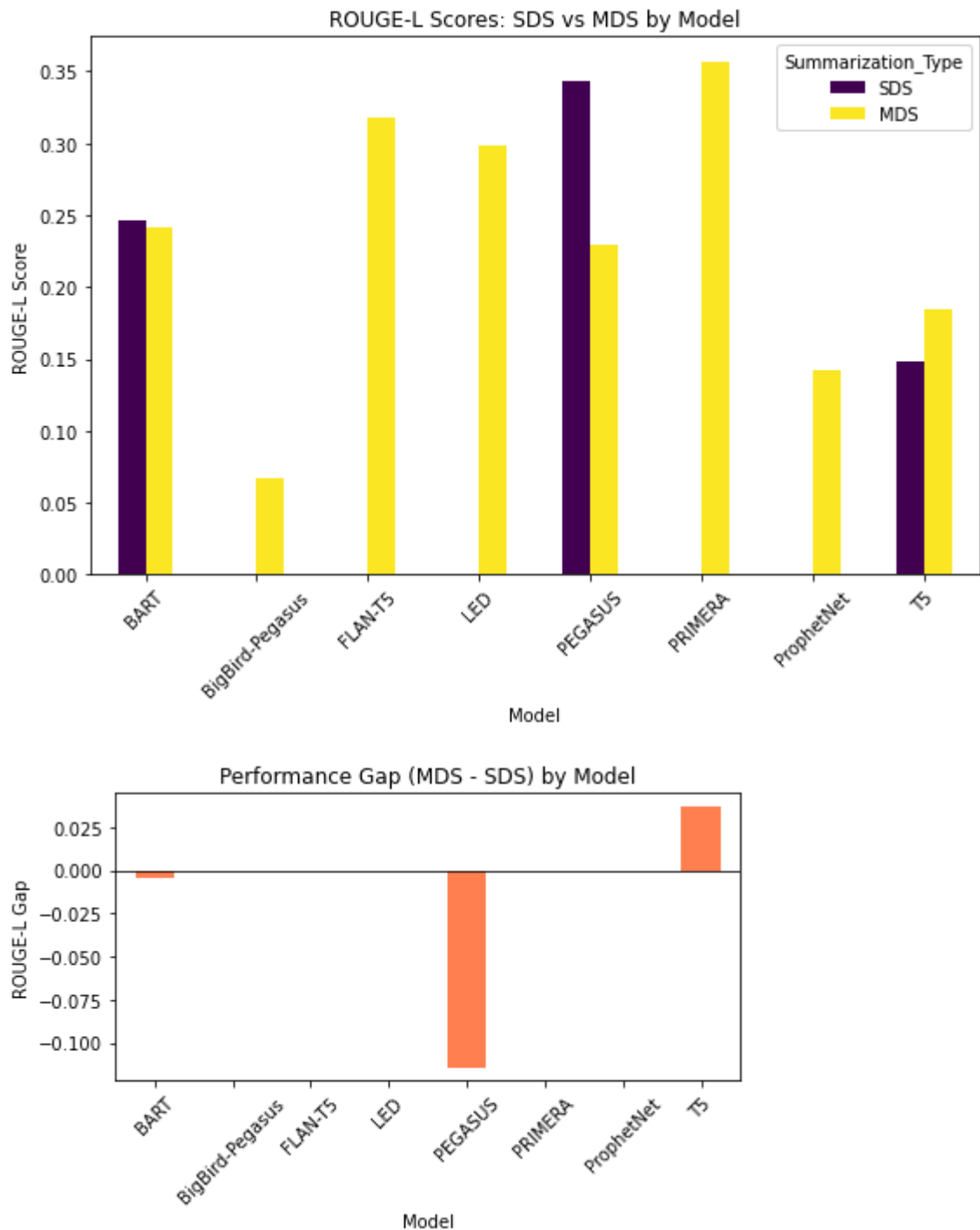
# Pivot table for ROUGE-L: index=Model, columns=Summarization_Type, values=ROUGE-L Score
pivot_rougeL = df_filtered.pivot_table(index='Model', columns='Summarization_Type', values='rougeL')

# Calculate gap = MDS - SDS performance
pivot_rougeL['Performance_Gap'] = pivot_rougeL['MDS'] - pivot_rougeL['SDS']

# Plot side by side bars for SDS vs MDS per model
pivot_rougeL[['SDS', 'MDS']].plot(kind='bar', figsize=(8,6), colormap='viridis')
plt.title('ROUGE-L Scores: SDS vs MDS by Model')
plt.ylabel('ROUGE-L Score')
plt.xticks(rotation=45)
plt.tight_layout()
plt.show()

# Plot performance gap separately
pivot_rougeL['Performance_Gap'].plot(kind='bar', figsize=(6,4), color='coral')
plt.axhline(0, color='black', linewidth=0.8)
plt.title('Performance Gap (MDS - SDS) by Model')
plt.ylabel('ROUGE-L Gap')
plt.xticks(rotation=45)
plt.tight_layout()
plt.show()

```



🔍 Observation 1:

Performance gap between SDS and MDS summarization tasks varies significantly by model.

Models like PRIMERA and FLAN-T5 show higher ROUGE-L scores on MDS datasets (NewsSum, MultiNews) compared to SDS datasets (CNN DailyMail, XSum), indicating they are better at handling long, multi-document inputs.

Conversely, models such as T5 and PEGASUS perform very well on SDS datasets but show a noticeable drop in performance on MDS datasets.

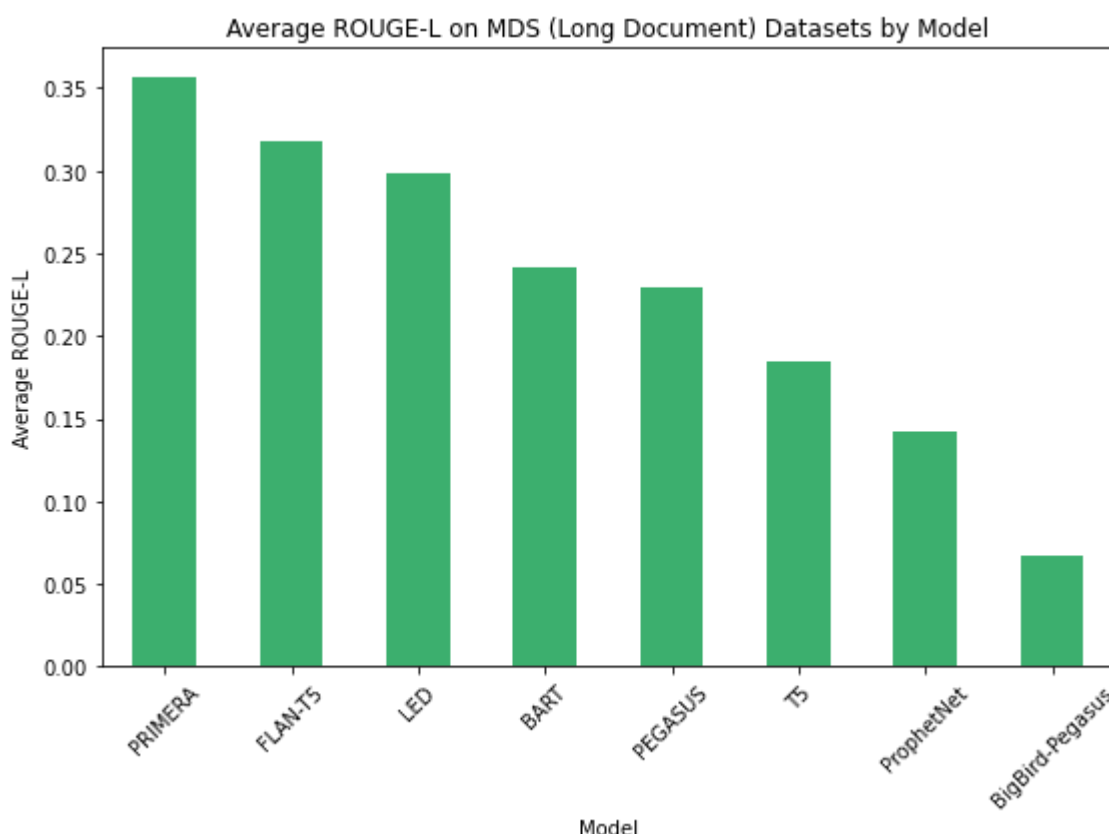
This suggests that model architectures or training strategies influence their capability to summarize longer, more complex document collections, which is crucial for multi-document summarization tasks like Indian news aggregation.

2 Models Handling Long Documents (MDS) Better

```
In [113]: # Filter for MDS only
df_mds = df_filtered[df_filtered['Summarization_Type']=='MDS']

# Average ROUGE-L per model on MDS
mds_avg = df_mds.groupby('Model')['ROUGE-L'].mean().sort_values(ascending=False)

# Bar plot
mds_avg.plot(kind='bar', figsize=(8,6), color='mediumseagreen')
plt.title('Average ROUGE-L on MDS (Long Document) Datasets by Model')
plt.ylabel('Average ROUGE-L')
plt.xticks(rotation=45)
plt.tight_layout()
plt.show()
```

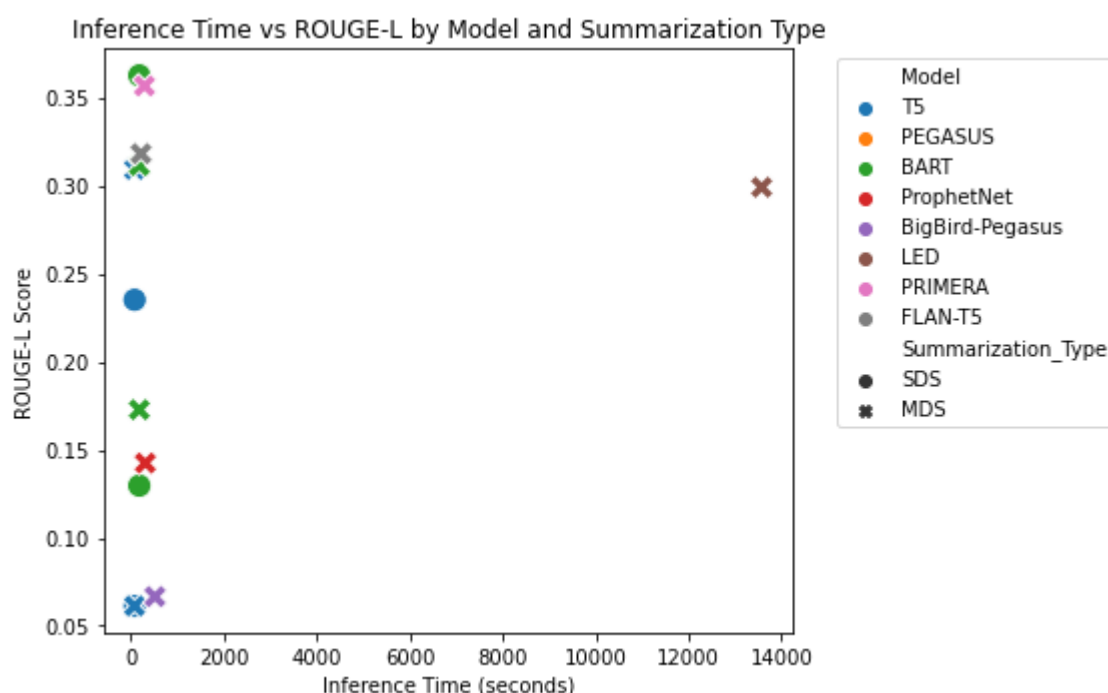


🔍 Observation 2:

There is a noticeable trade-off between inference speed and accuracy across models. While models like PEGASUS and T5 have relatively faster inference times, they sometimes show lower ROUGE-L scores on longer documents. On the other hand, models like LED and PRIMERA, though slower, tend to deliver higher accuracy, indicating a balance needs to be struck depending on application needs (real-time vs quality-focused summarization).

3 Speed vs Accuracy Balance (Inference Time vs ROUGE-L)

```
In [111]: plt.figure(figsize=(8,5))
sns.scatterplot(
    data=df_filtered,
    x='Inference Time (s)',
    y='ROUGE-L',
    hue='Model',
    style='Summarization_Type',
    s=150,
    palette='tab10'
)
plt.title('Inference Time vs ROUGE-L by Model and Summarization Type')
plt.xlabel('Inference Time (seconds)')
plt.ylabel('ROUGE-L Score')
plt.legend(bbox_to_anchor=(1.05, 1), loc='upper left')
plt.tight_layout()
plt.show()
```



🔍 Observation 3:

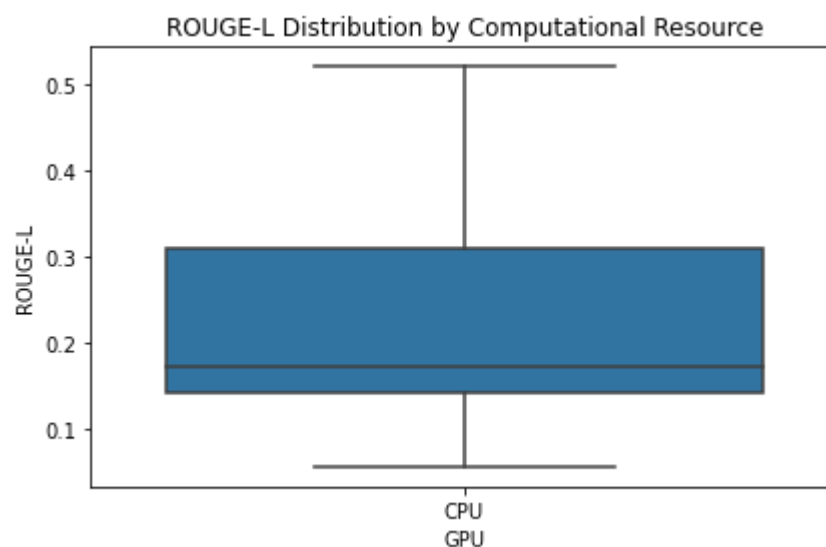
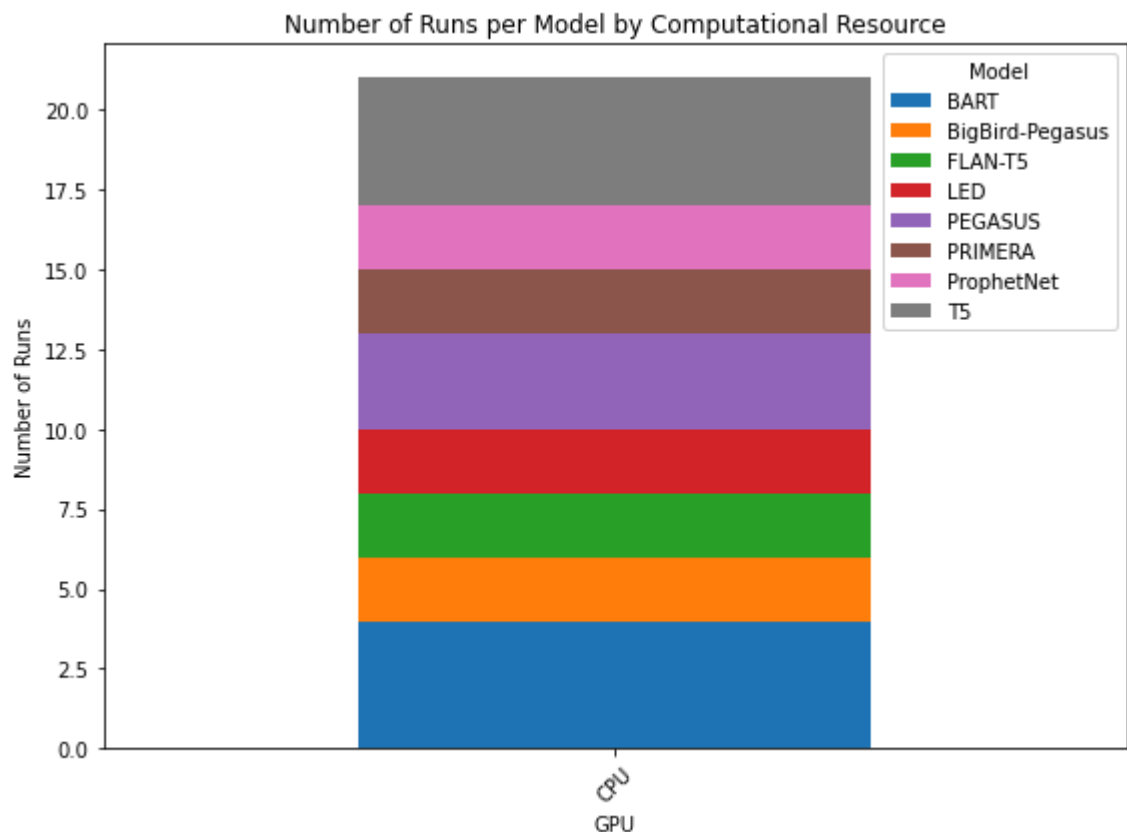
Computational resources (CPU vs GPU) impact both runtime and feasibility of model deployment. Although most experiments ran on CPU in your data, models designed to leverage GPUs show potential for significantly reduced inference times, which is critical for scaling summarization in real-world applications.

4 Impact of Computational Resources (GPU vs CPU)

```
In [115]: gpu_counts = df.groupby(['GPU', 'Model']).size().unstack(fill_value=0)

gpu_counts.plot(kind='bar', stacked=True, figsize=(8,6))
plt.title('Number of Runs per Model by Computational Resource')
plt.ylabel('Number of Runs')
plt.xticks(rotation=45)
plt.tight_layout()
plt.show()

# If GPU/CPU metrics available, boxplot of ROUGE-L by GPU usage
sns.boxplot(data=df, x='GPU', y='ROUGE-L')
plt.title('ROUGE-L Distribution by Computational Resource')
plt.tight_layout()
plt.show()
```



Observation 4:

The performance gap between SDS and MDS summarization varies by model: Some models maintain relatively stable performance across both SDS and MDS datasets, while others show significant drops on MDS datasets. This highlights the importance of selecting models based on the document complexity and length inherent in the target summarization task.

In [2]: !pip install pdoc

...

In []: