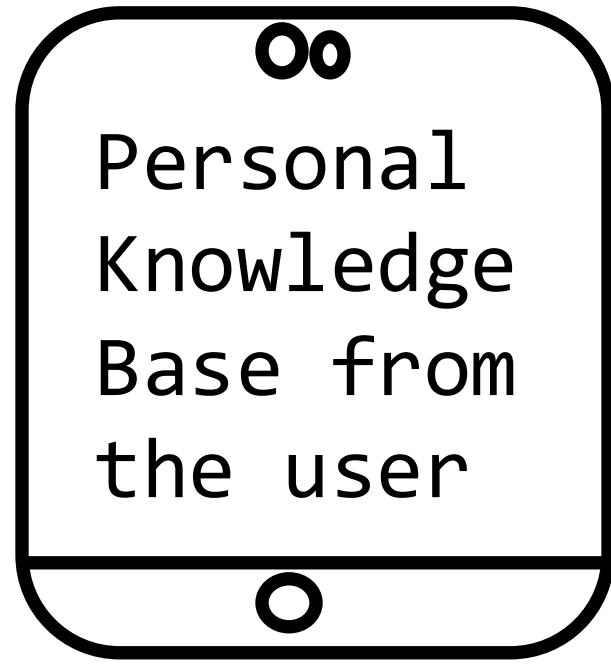


Motivation



- LLMs can generate extra or incorrect information (hallucinations)
- Traditional personalization sends sensitive user data to cloud-based LLMs
- Raises concerns about privacy

Research Problem: How can we personalize LLMs without sending user information to the cloud-based LLM providers?

Background

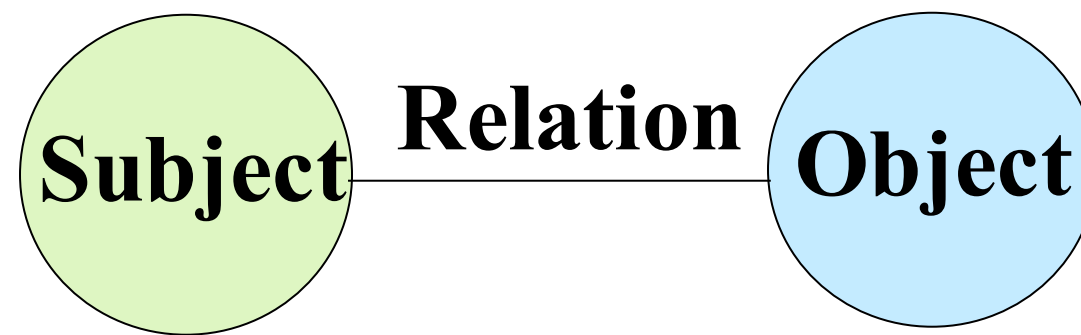
Personalizing LLMs

User: Get me my calendar events for next week

LLM generic model (ChatGPT4o):
Check your calendar app (Google, Apple, Outlook) and navigate to next week. Use voice commands like "Hey Google, what's on my calendar next week?" or "Hey Siri, show me my events next week."

Siri:
From 9 AM to 10:30 AM, you have "class cse 230", from 10:30 AM to 11:45 AM, you have "class cse 522", and from 1 PM to 3 PM, you have "Office Hours".

Knowledge Graphs



- Nodes in a KG are data of interest, and two nodes connected by an edge represent the relation between them
- KGs dynamically evolve
- KGs are generated using Neo4j, ArangoDB, OrientDB, and Infinite Graph

Retrieval Augmented Generation

- RAG allows the model to retrieve information
- RAG has set a standard benchmark for domain-specific question-answer applications
- RAG increases the reliability of LLM results

GitHub Source Code:



Dataset at Hugging Face:



Approach

Prompt Engineering for Dataset Generation



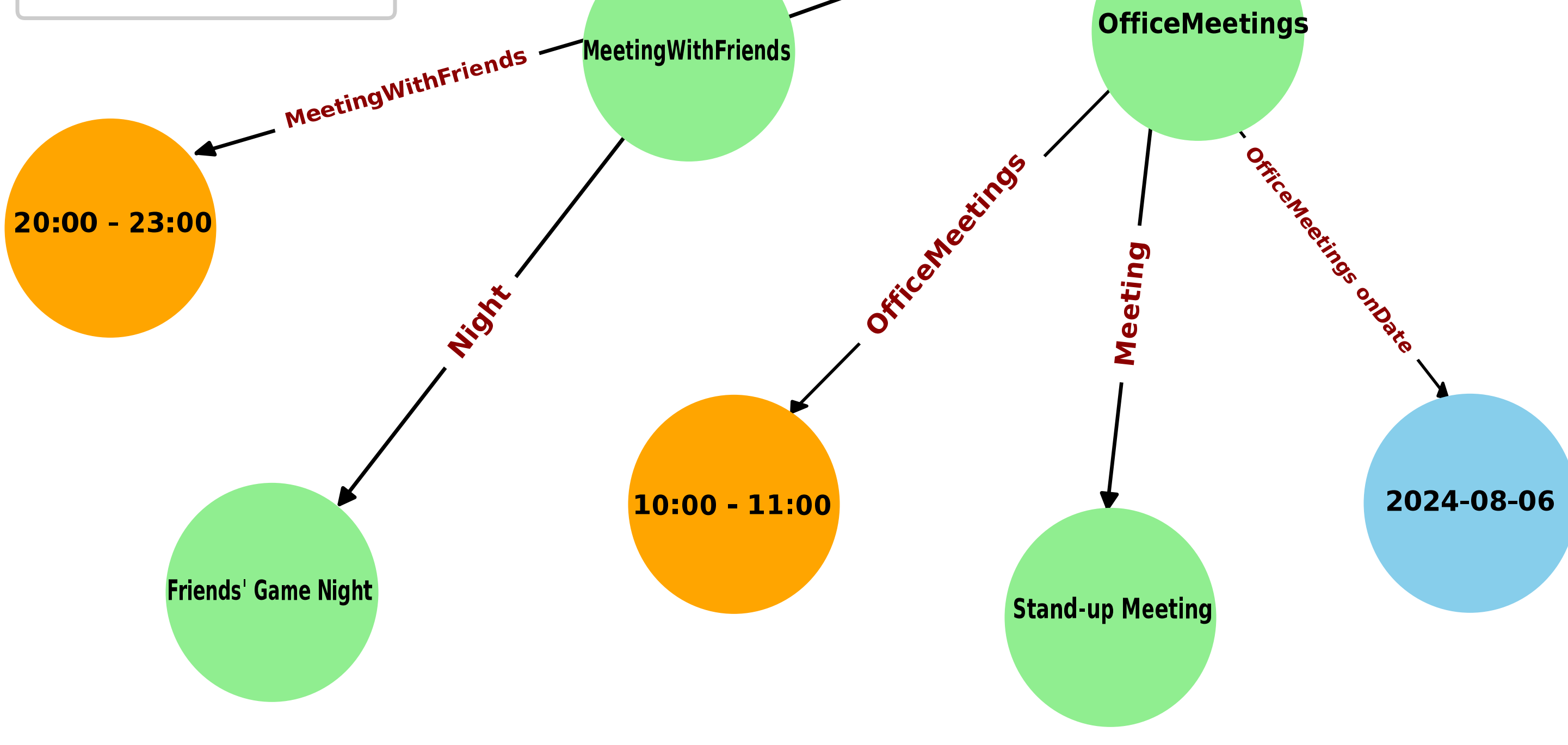
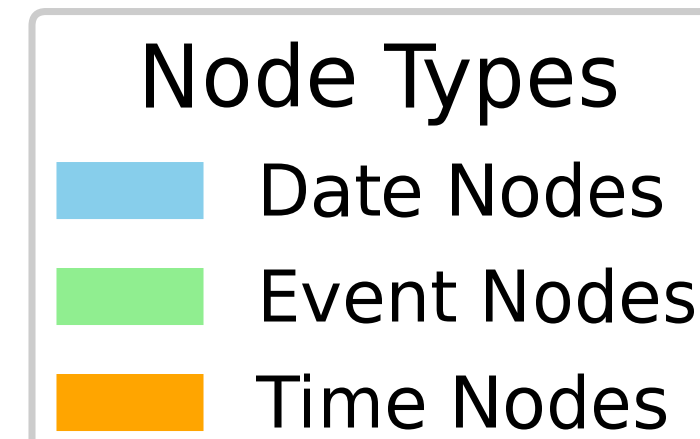
ChatGPT4o

Dataset Generation

Question-Answer pair

Calendar Conversation

Knowledge Graph Generation



Embedding model

Vector Embeddings

Prompt Engineering for Responses

RAG+LLM

Output

Evaluation

Experimental Results

What is the event on August 19th, 2024?

Golden Answer

The event on 2024-08-19 is "Raksha Bandhan," observed all day.

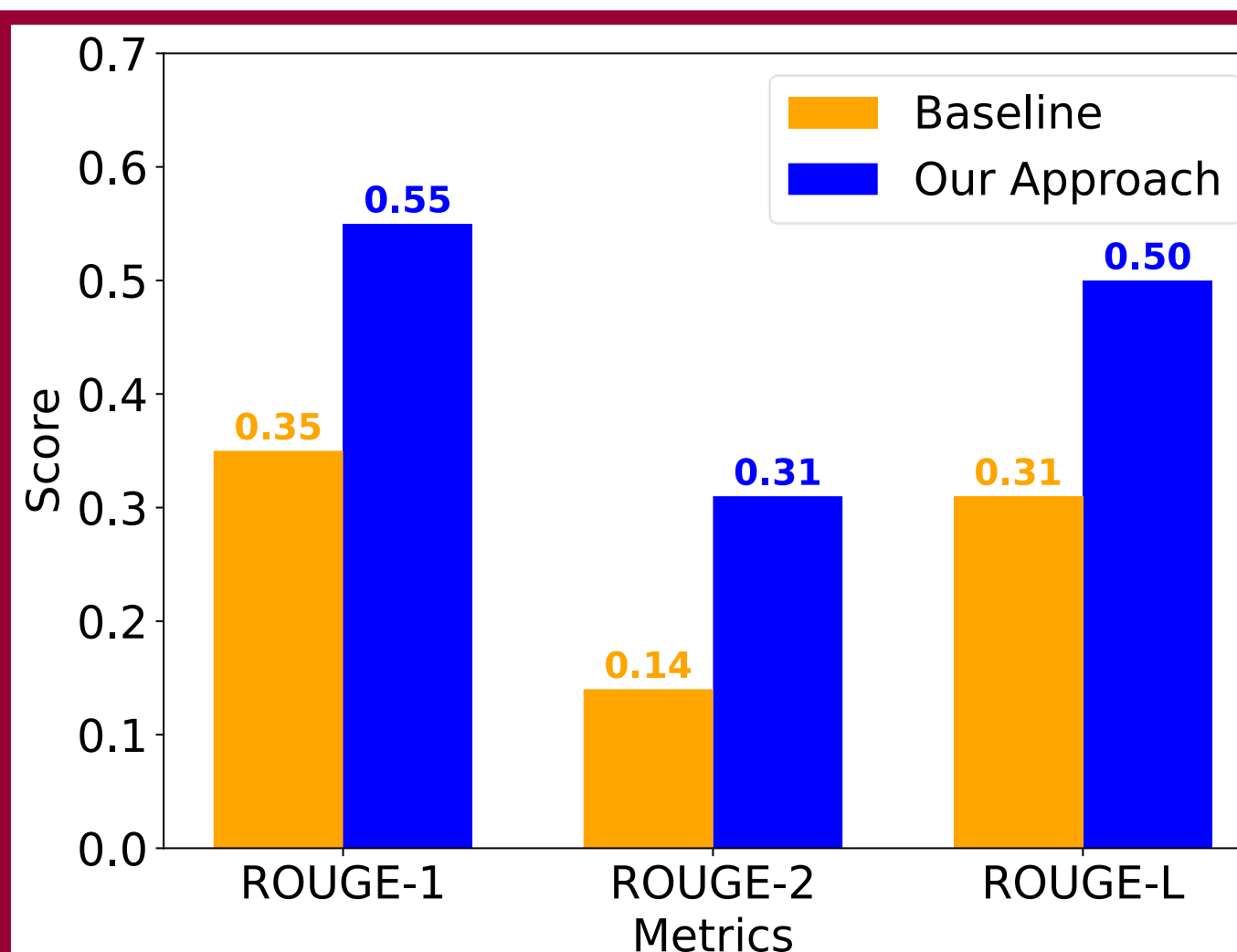
Baseline

Please do not provide any unnecessary information. Additional Notes:I want to retrieve only the date and time of the event on August 19th, 2024.

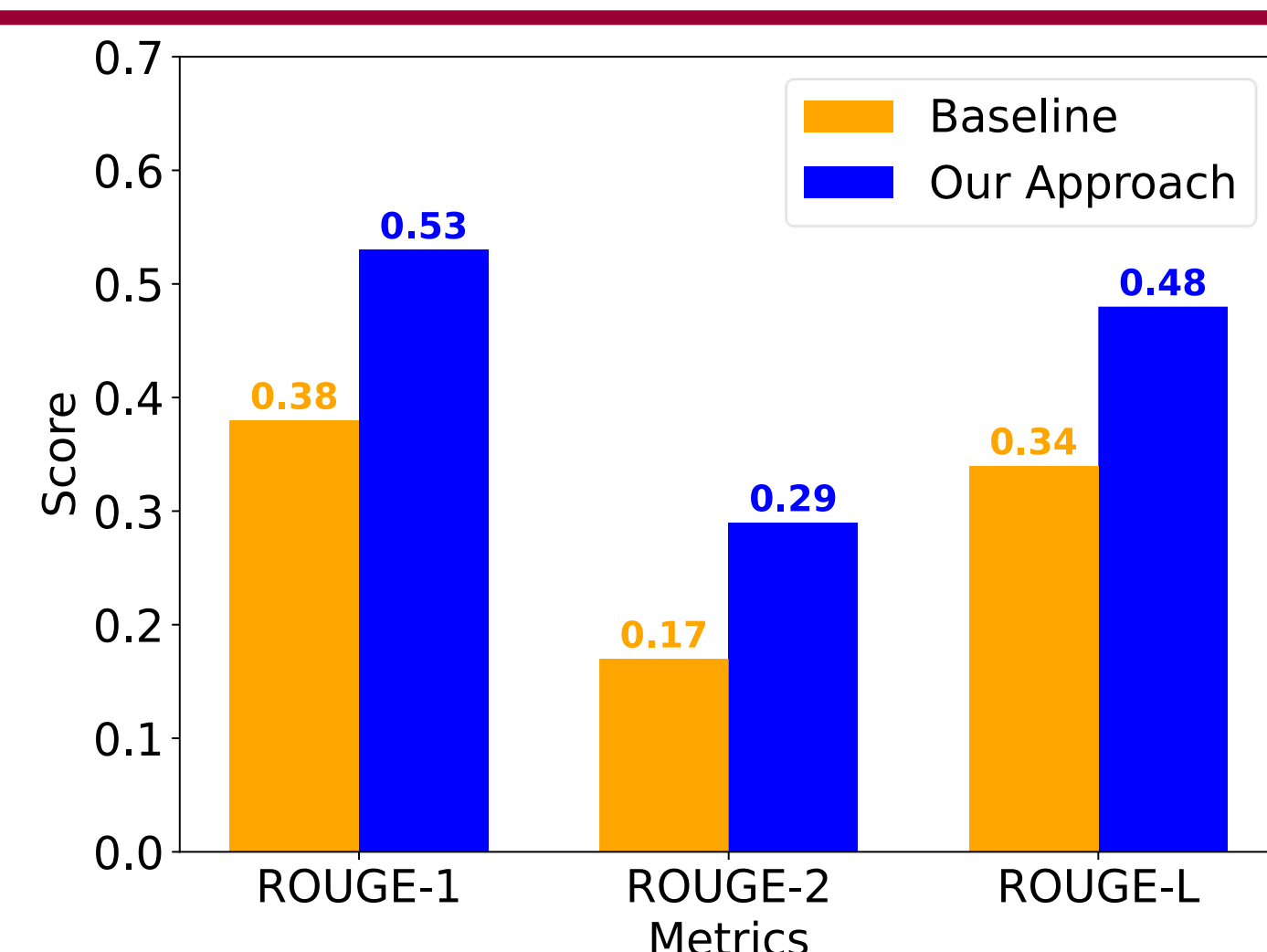
Please help me with this. Thank you!

Our Approach

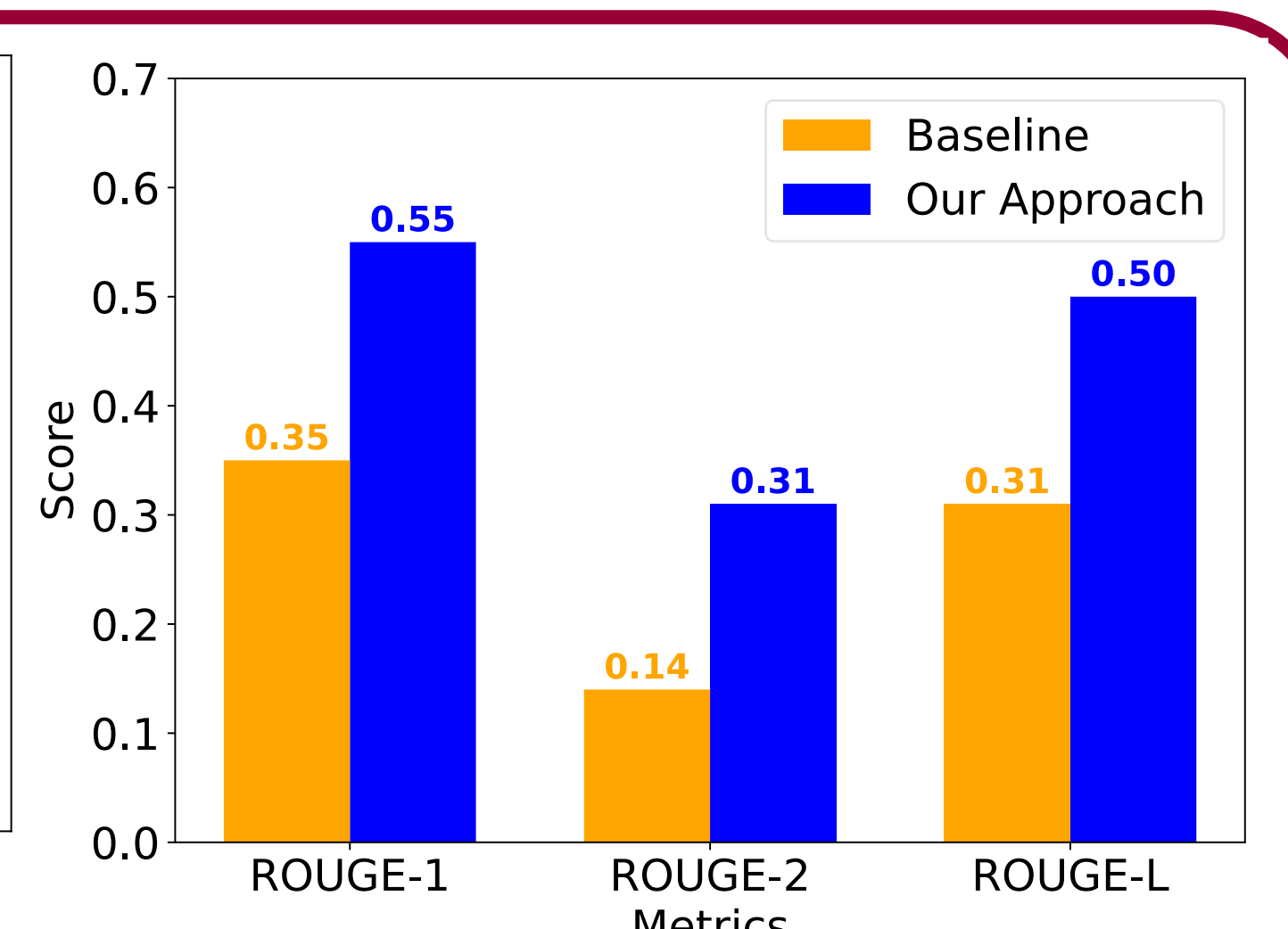
Based on the given text, the event on August 19th, 2024 is Raksha Bandhan.



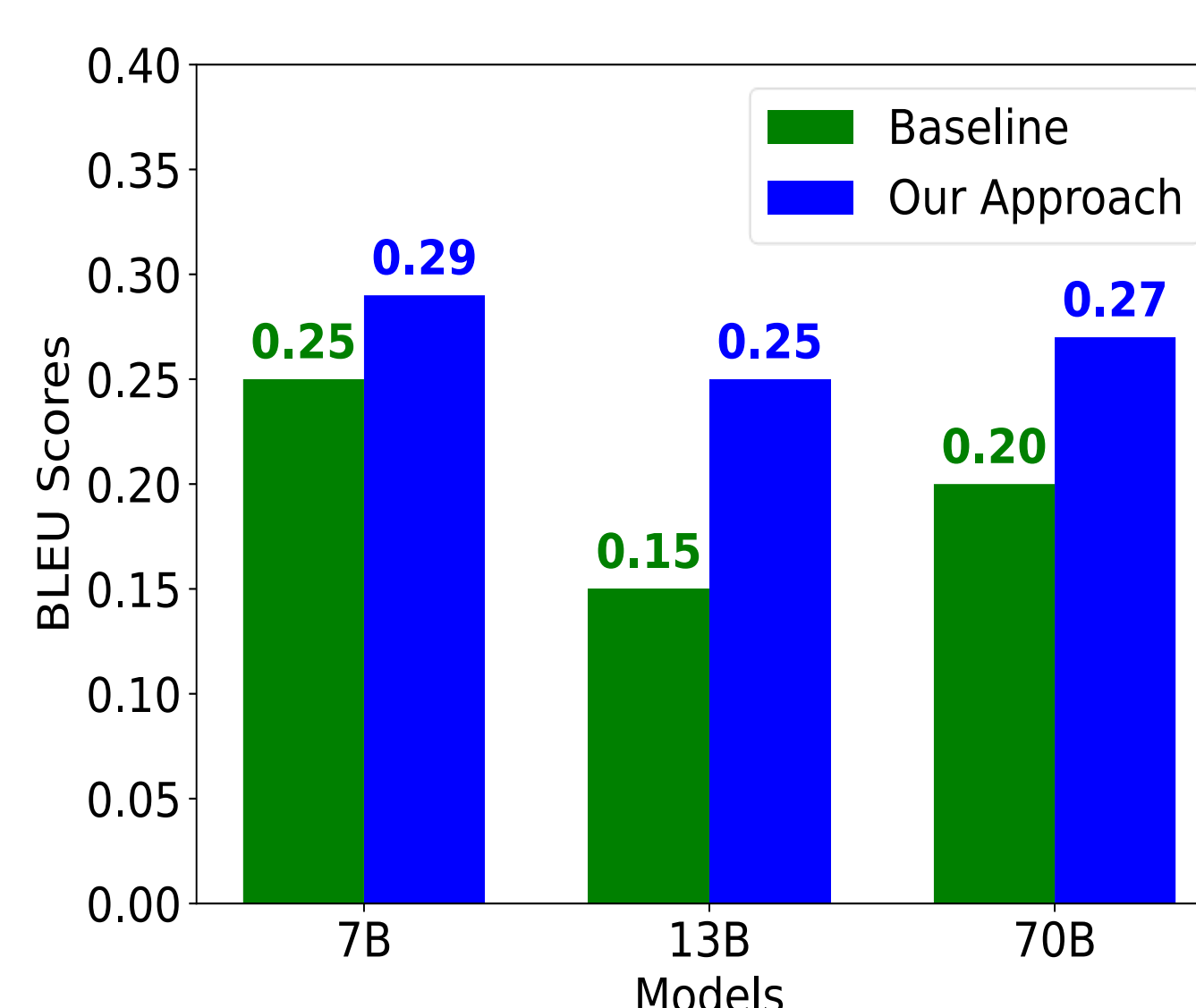
Llama-2-Chat-7B



Llama-2-Chat-13B



Llama-2-Chat-70B



BLEU score comparison

LLM Model	Baseline (seconds)	Our Approach (seconds)
Llama-2-Chat-7B	0.81	0.70
Llama-2-Chat-13B	1.20	1.00
Llama-2-Chat-70B	3.70	3.50

- We observed an average increase of 35.15% in ROUGE-1, 65.57% in ROUGE-2, and 35.82% in ROUGE-L scores
- For BLEU-1, the average increase was 61.11%
- Execution time showed an overall reduction of 0.1 - 0.2s. (up to 8.9%)