

Text Summarization



Mentor: Mr. Narendra Kumar

Group: 4



Introduction

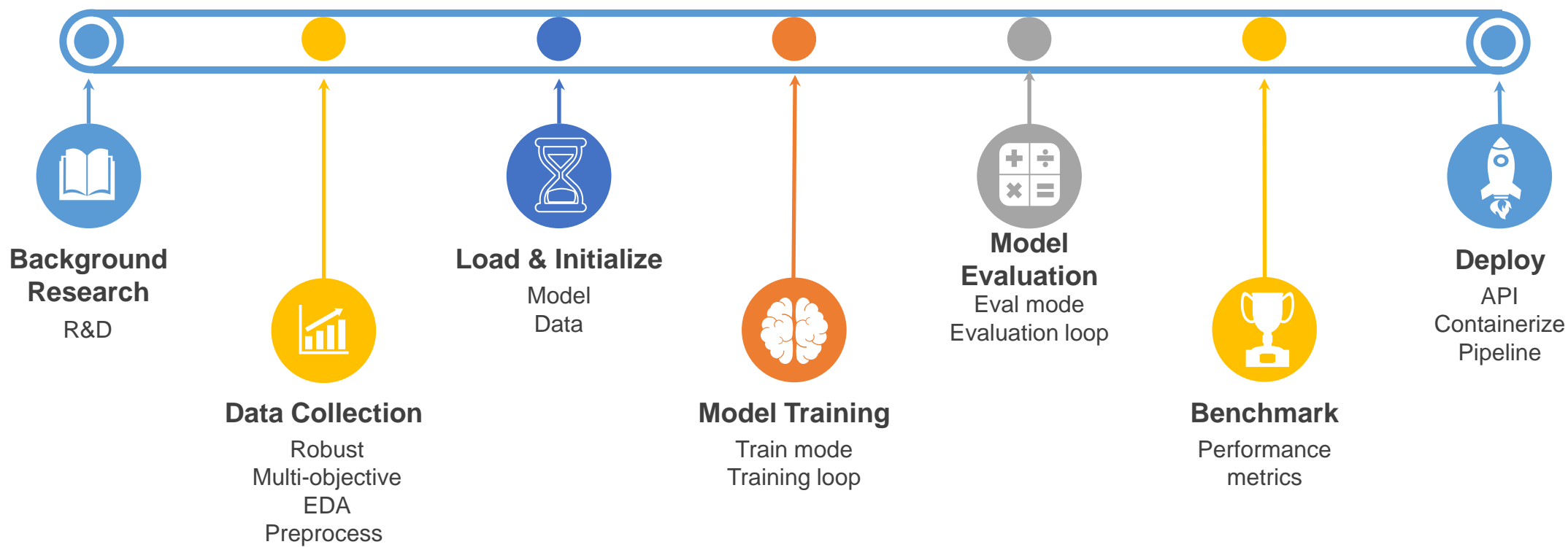
Problem Statement & Planning

Introduction

Problem Statement

- Developing an automated text summarization system that can accurately and efficiently condense large bodies of text into concise summaries is essential for enhancing business operations.
- This project aims to deploy NLP techniques to create a robust text summarization tool capable of handling various types of documents across different domains.
- The system should deliver high-quality summaries that retain the core information and contextual meaning of the original text.

INTENDED PLAN



Background Research

Literature Review & Findings

Background Research

Literature Review

| S. No | Use-Case | Paper Title | Year | Method | Dataset | Results | Limitations |
|-------|----------------------------|---|------|--|----------------------|--|--|
| 1 | General text summarization | Text Summarization Using Deep Learning Techniques: A Review | 2023 | Deep Learning (Seq2Seq, Attention, Transformers) | CNN/Daily Mail, XSum | Improved performance in capturing semantic relationships, better coherence | Computationally expensive, requires large datasets |

[1] Saiyyad, M.M.; Patil, N.N. "Text Summarization Using Deep Learning Techniques: A Review". Eng. Proc. 2023, 59, 194.

Background Research

Literature Review

| S. No | Use-Case | Paper Title | Year | Method | Dataset | Results | Limitations |
|-------|--|---------------------------|------|-------------|--|---|---|
| 2. | Implementation of the Transformer architecture | Attention is all you need | 2023 | Transformer | WMT 2014 English-German, WMT 2014 English-French | Introduced the Transformer architecture, significantly improving the performance of text summarization tasks. | Requires large datasets and computational resources for training. |

[2] A. Vaswani, L. Jones, N. Shazeer, N. Parmar, A. N. Gomez, J. Uszkoreit, Ł. Kaiser, and I. Polosukhin, "Attention Is All You Need," arXiv:1706.03762v7 [cs.CL], Aug. 2, 2023.

Background Research

Literature Review

| S. No | Use-Case | Paper Title | Year | Method | Dataset | Results | Limitations |
|-------|------------------------------|--|------|--|----------------|---|--|
| 3. | Multi-document summarization | Surveying the Landscape of Text Summarization with Deep Learning | 2023 | Deep learning methods. Various techniques like RBMs and fuzzy logic employed for summarization. | CNN/Daily Mail | Incorporating transfer learning enhances summary quality and reduces data demand. | Complex models, high computational resources |

[3] G. Wang and W. Wu, "Surveying the Landscape of Text Summarization with Deep Learning: A Comprehensive Review," arXiv:2310.09411v1 [cs.CL], Oct. 13, 2023.

Background Research

Literature Review

| S. No | U4se-Case | Paper Title | Year | Method | Dataset | Results | Limitations |
|-------|---------------------------|--|------|-----------------------|---------------------------------------|---|--------------------|
| 4. | Abstractive summarization | Pegasus: Pre-training with gap-sentences for abstractive summarization | 2020 | Transformer (Pegasus) | XSum, CNN/Daily Mail, and Reddit TIFU | Significant improvements in abstractive summarization quality | Resource-intensive |

[4] J. Zhang, Y. Zhao, M. Saleh, and P. J. Liu, "PEGASUS: Pre-training with Extracted Gap-sentences for Abstractive Summarization," arXiv:1912.08777v3 [cs.CL], Jul. 10, 2020.

Background Research

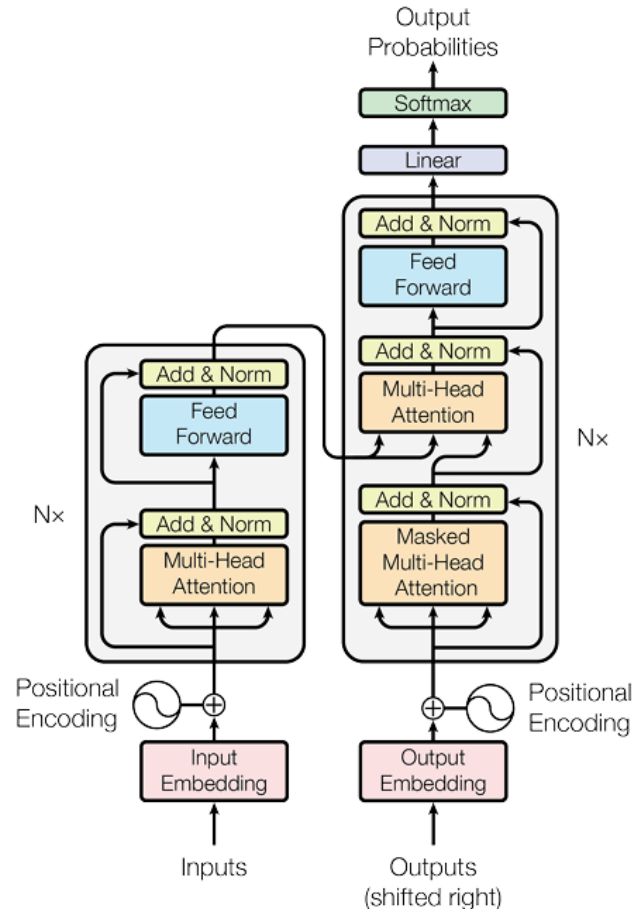
Literature Review

| S. No | Use-Case | Paper Title | Year | Method | Dataset | Results | Limitations |
|-------|--------------------------|---|------|--|--------------------------------------|---|---------------------------------------|
| 5. | Extractive summarization | Text Summarization with Pretrained Encoders | 2019 | Intersentence Transformer layers for summarization | CNN/Daily Mail, NYT, Xsum, DailyMail | BERT-based models outperformed other approaches in abstractive summarization. | High computational resources required |

[5] Y. Liu and M. Lapata, "Text Summarization with Pretrained Encoders," arXiv:1908.08345v2 [cs.CL], Sep. 5, 2019.

Research

Selected Architecture



[2] Fig. :Transformer architecture:

- **Implementation methods:**

- *From Scratch*

- Build Model
 - NN
- Initialize normalized W&B
- Train model with extensive data
- Hence,
 - Computationally Intensive
 - Sub-Optimal usage of resources
 - Out-of-scope

- *Using Pre-trained model*

- Load Model & its parameters
- Re-Train with specific dataset
- Evaluate
- Hence,
 - Innovation can be done at intended tasks
 - Optimal utilization of resources

Proposal

Architecture, Findings & Team Details

Proposal Workflow

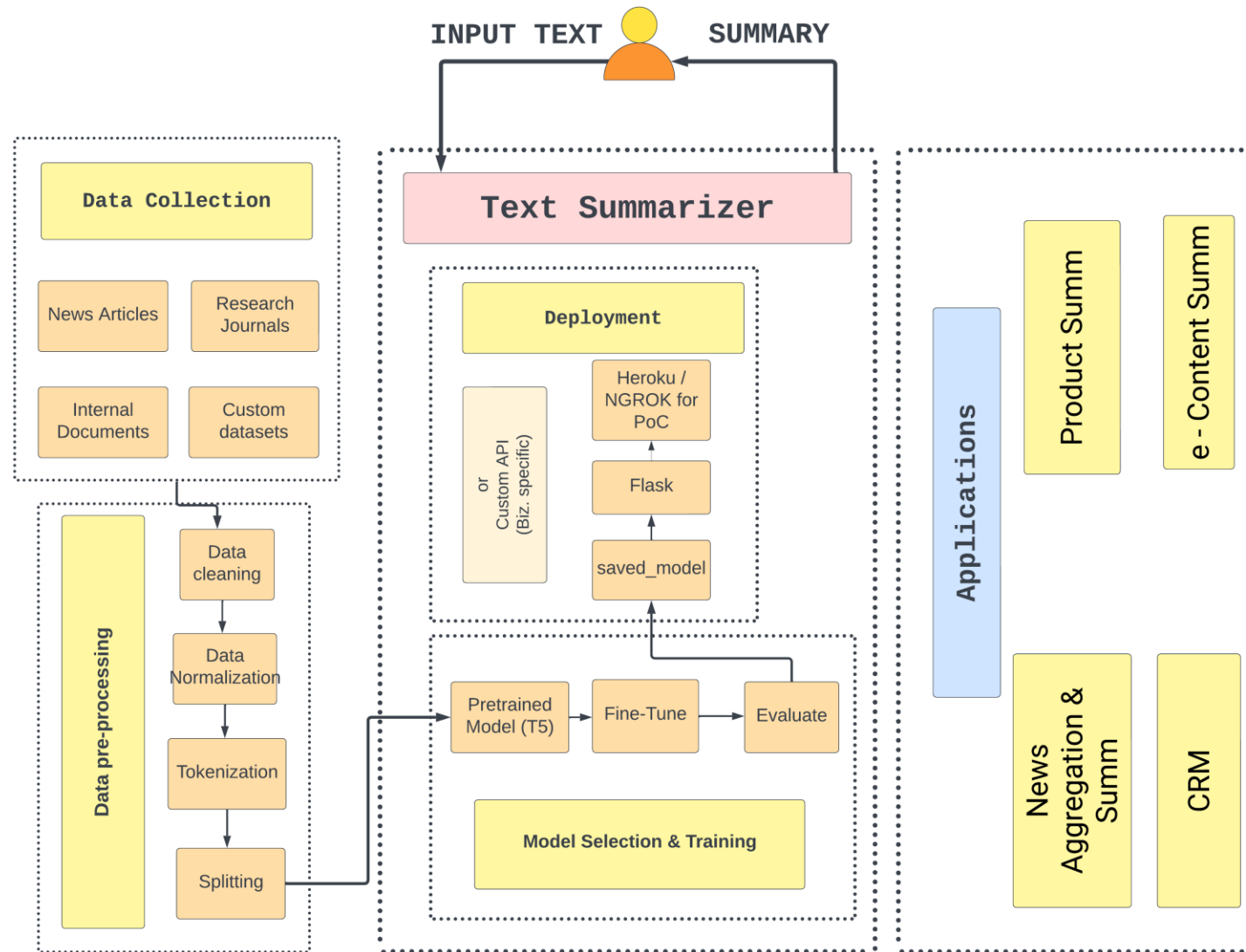


Fig. : Proposed Workflow

Proposal

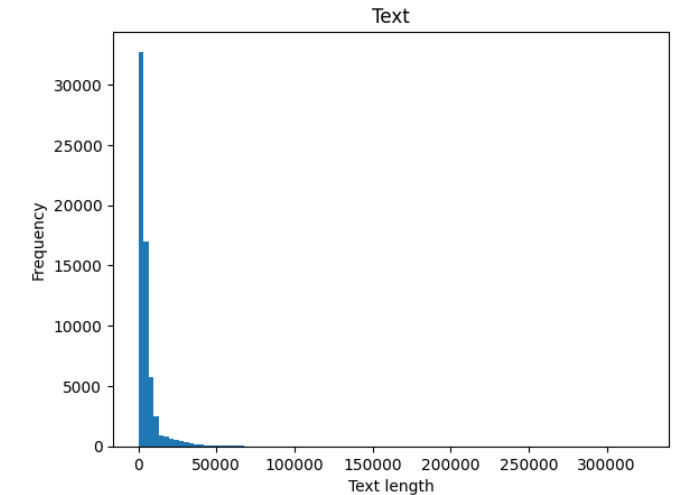
Dataset

- Merged selective dataset from
 - CNN, Daily Mail : News,
 - BillSum: Legal,
 - ArXiv : Scientific
 - Dialoguesum : Conversations.
- Completed - data preprocessing
 - Removed
 - NULL records, punctation, stop-words
 - Lowercasing, lemmatization.

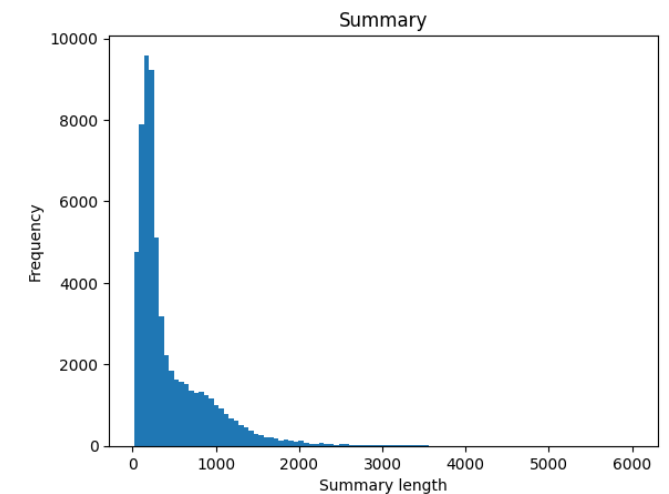
| | text | summary |
|-------|---|---|
| 0 | section 1 liability business entity providing ... | shield business entity civil liability relatin... |
| 1 | section 1 short title act may cited human righ... | human right information act requires certain f... |
| 2 | section 1 short title act may cited jackie rob... | jackie robinson commemorative coin act directs... |
| 3 | section 1 nonrecognition gain rollover small b... | amends internal revenue code provide temporari... |
| 4 | section 1 short title act may cited native ame... | native american energy act sec 3 amends energy... |
| ... | ... | ... |
| 62702 | person1 excuse mr green manchester arent perso... | tan ling pick mr green easily recognized white... |
| 62703 | person1 mister ewing said show conference cent... | person1 person2 plan take underground together... |
| 62704 | person1 help today person2 would like rent car... | person2 rent small car 5 day help person1 |
| 62705 | person1 look bit unhappy today whats person2 w... | person2s mom lost job person2 hope mom wont fe... |
| 62706 | person1 mom im flying visit uncle lee family n... | person1 asks person2s idea packing bag visitin... |

62707 rows x 2 columns

```
count    62707.000000
mean      5211.270975
std       7794.860686
min        83.000000
25%      1275.000000
50%      3176.000000
75%      5684.500000
max     323742.000000
Name: text, dtype: float64
```



```
count    62707.000000
mean      448.081937
std       459.087443
min        16.000000
25%       154.000000
50%       255.000000
75%       618.000000
max     6014.000000
Name: summary, dtype: float64
```



* In characters.

https://drive.google.com/drive/folders/1yH89iZmARdc-R7QY6pwfE8tbOJI_n9K8?usp=sharing
[Infosys_Text-Summarization/src/data_preprocessing.ipynb](#) at main · MohanKrishnaGR/Infosys_Text-Summarization (github.com)

Proposal

Model Training



Fig. : Fine-Tuning Overview

- Proposed implementation – Two – 2 Methods
 - Method 1 – Native PyTorch Method
 - Method 2 – Trainer Class Method

Proposal

Model Training (Method 1)



Hugging Face

PyTorch

- Load pre-trained transformer
 - Facebook's Bart Large
- OOP implementation of Dataset
 - Feature, Target
 - Tokenize
 - Padding, Truncate
 - Convert to Tensor
 - Pass to: DataLoader – with batch size
- Training Loop
 - Adam optimizer
 - Forward pass & compute loss
 - Backward pass
 - Update params – compute gradient
 - Update LR
 - Zero the gradients
 - Update total loss
- Only minimal train loss of 1.3280.
 - But, produced inconsistent results.
 - Cannot be pushed into production.
- Raises the need for optimized training and eval loop for Transformer.

```
import pandas as pd
from torch.utils.data import Dataset, DataLoader
from transformers import BartTokenizer

class SummarizationDataset(Dataset):
    def __init__(self, file_path, tokenizer, max_length=512):
        self.dataset = pd.read_csv(file_path)
        self.tokenizer = tokenizer
        self.max_length = max_length

    def __len__(self):
        return len(self.dataset)

    def __getitem__(self, idx):
        text = self.dataset.iloc[idx, 0]
        summary = self.dataset.iloc[idx, 1]

        inputs = self.tokenizer.encode_plus(
            text,
            max_length=self.max_length,
            padding='max_length',
            truncation=True,
            return_tensors='pt'
        )

        targets = self.tokenizer.encode_plus(
            summary,
            max_length=self.max_length,
            padding='max_length',
            truncation=True,
            return_tensors='pt'
        )

        return [inputs, targets]
```

Fig. : Screenshot

Proposal

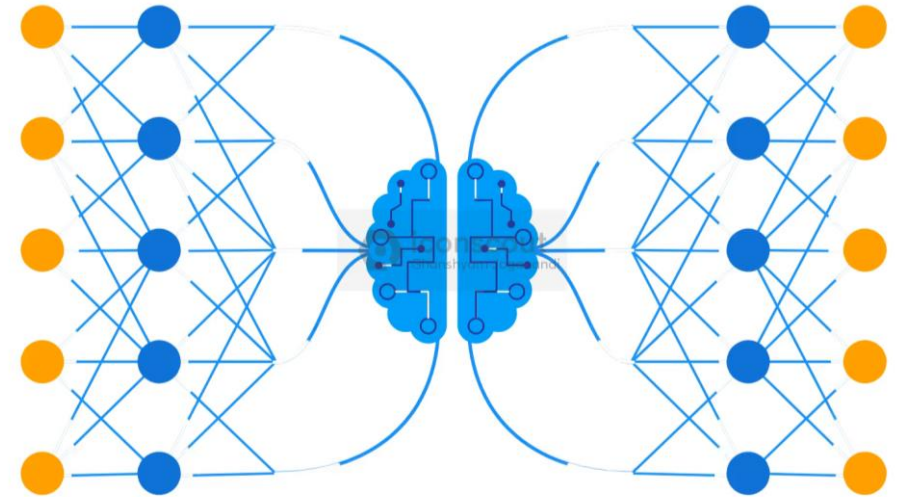
Model Training (Method 2)

- Trainer Method
- Implemented in src/bart.ipynb.
- A function was implemented for the dataset, to convert text data into model inputs and targets.
- Trainer class from transformer package was utilized for training and evaluation. Tainer is a simple but feature-complete training and eval loop for PyTorch, optimized for transformers.
- The model was trained with whole dataset for 10 epochs for 26:24:22 (HH:MM:SS) in 125420 steps.
- Training Loss = 17.4700
- Considered the performance metrics of the models trained by the forementioned methods. After the due analysis, the model trained using 'Method 2' was selected.



Hugging Face

PyTorch



[Infosys_Text-Summarization/src/bart.ipynb at main · MohanKrishnaGR/Infosys_Text-Summarization \(github.com\)](#)

https://drive.google.com/drive/folders/1tNdLI67UTc5es6VB_dml8b5gkRUWzupl?usp=drive_link

Proposal

Model Validation

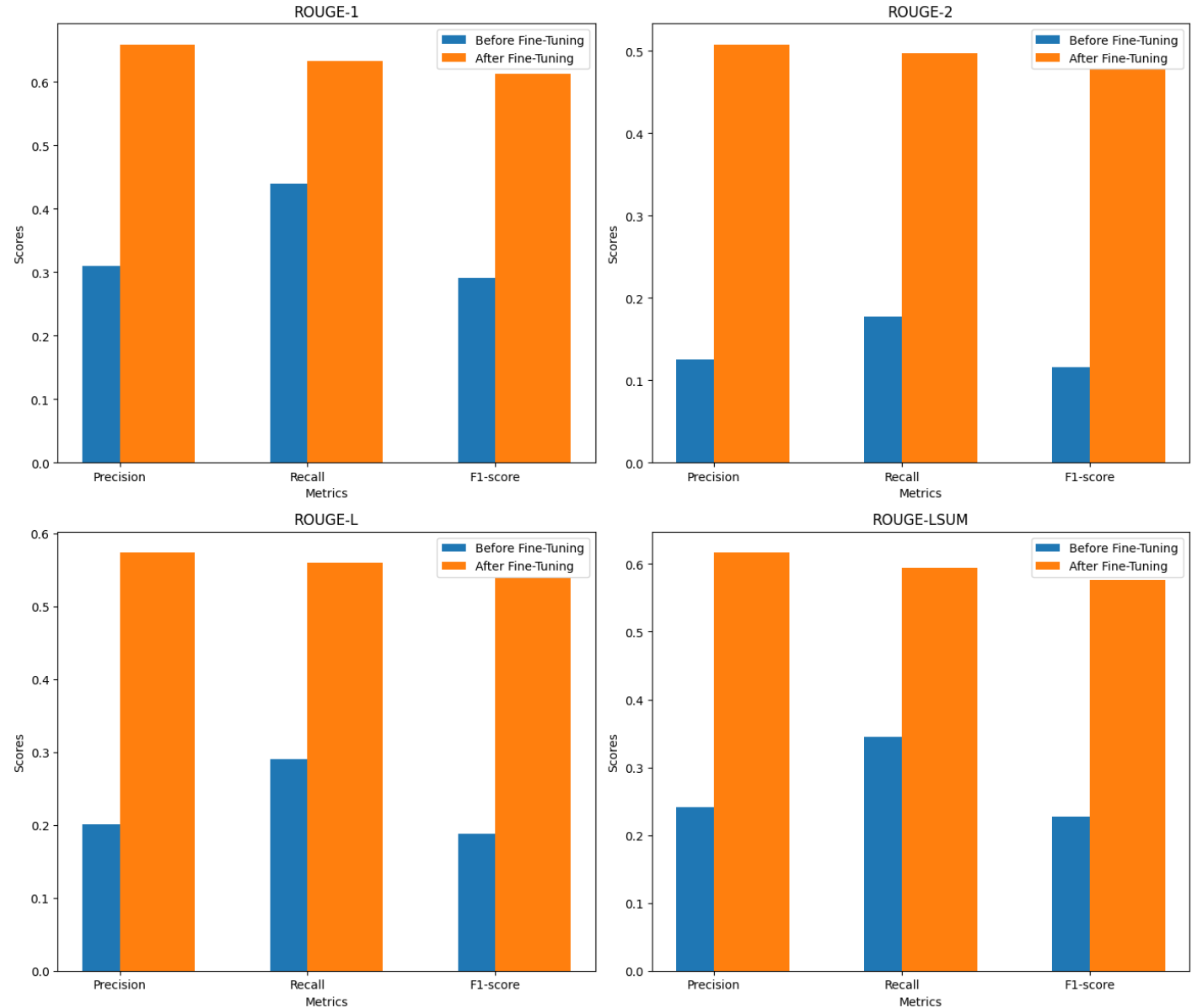
- Performance metrics – **ROUGE** (Recall-Oriented Understudy for Gisting Evaluation)
 - Overlap between generated summary and reference summary.
 - Best suited : evaluating 'Text Summarization' tasks.
 - Other options : **BLEU**.
- Aimed to: implement custom evaluation function.
 - Calc. : ROUGE based on model's inference.
- ROUGE-N: Measures the overlap of n-grams (contiguous sequences of n items) between the candidate summary and the reference summaries.
 - **ROUGE-1:**
 - Overlap of unigrams (single words).
 - **ROUGE-2:**
 - Overlap of bigrams (two-word sequences).
 - **ROUGE-L:**
 - Measures the longest common subsequence (LCS) between the candidate and reference summaries.
 - **ROUGE-LSUM**
 - (LCS Summary) - variant of the ROUGE-L metric, specifically designed to evaluate the quality of summaries.



Proposal

Comparative Analysis

- Analysis of the transformer's performance metrics before and after Fine-Tuning.
- The transformer model shows significant improvements across all ROUGE metrics after fine-tuning.
- The most substantial gains observed in ROUGE-2 scores. This indicates that the fine-tuning process has notably enhanced the model's ability to generate more accurate and relevant summaries.
- The model is now more proficient at generating summaries that are precise, comprehensive, and contextually accurate.
- Will act as a powerful tool for a variety of Business applications that require efficient and effective text summarization.



Proposal Testing

Infosys Springboard - Text Summarizer

A simple and efficient text summarizer. Enter your text in the box below and get a concise summary.

Input text

At least 49 migrant workers, including around 40 Indian citizens, have died in a deadly fire that devastated a building in Kuwait's southern district of Al-Mangaf. The fire that broke out in the apartment building located in Kuwait's Al Ahmadi Governorate early on Wednesday also left more than a dozen injured, who were admitted to nearby hospitals, reported the Kuwait News Agency (KUNA). Prime Minister Narendra Modi and External Affairs Minister S. Jaishankar expressed shock over the incident, and Congress leader Rahul Gandhi expressed 'serious concern' about the condition of Indians in the Gulf region.

"My thoughts are with all those who have lost their near and dear ones. I pray that the injured recover at the earliest. The Indian Embassy in Kuwait is closely monitoring the situation and working with the authorities there to assist the affected," said Mr. Modi in a message. Mr. Modi held a review meeting on Wednesday evening about the condition of the affected Indians in Kuwait. He deputed Minister of State for External Affairs Kirti Vardhan Singh to oversee the help being rendered to the injured and to bring back the remains of the Indians who perished in the incident.

Indian ambassador to Kuwait Adarsh Swaika visited the Mubarak Al-Kabeer Hospital where 11 injured workers were admitted. "Ten of them are expected to be released today and one in hospital is reportedly stable," the Indian embassy said in a statement. The Government of Kuwait has not made any statement officially so far, but Interior Minister Sheikh Fahad Al-Yousuf Al-Sabah has ordered the police to arrest the owner of the building located in Al-Mangaf.

The incident has highlighted the poor living conditions of blue-collar Indian workers in the region.

Clear Submit

Summarized Text

Fire broke out in apartment building located in Kuwait's Al Ahmadi Governorate early on Wednesday. At least 49 migrant workers, including around 40 Indian citizens, have died in the blaze. The Indian embassy in Kuwait is monitoring the situation and working with the authorities there to assist the affected people.

Flag

Powered by Infosys Springboard Intern. Let's connect, [LinkedIn](#).

- Simple interface for the Deep Learning model, developed using Gradio.
- Gradio is an open-source Python package that allows us to quickly build a demo - web-application for the trained models.
- Enables us to test and even deploy the trained model.