



SNT HACKATHON 2023

PROBLEM STATEMENT: GENERATIVE AI FOR IMPACT

---

## Summarizer using LED 🤖

GitHub Repo: <https://github.com/Attention-is-All-We-Need/Document-Summary-Generator>

TEAM NAME: AttentionIsAllWeNeed

### Team Members :

1. Anwesh Saha (210178)
2. Arindom Bora (210183)
3. Ajay Sankar Makenna (210077)
4. Khush Khandelwal (210511)
5. Shreyash Nallawar (211010)
6. Vineet Kumar (201121)

# Table of Contents

Sl. No.	Topic	Page Number
1	Cover Page	1
2	Table of Contents	2
3	About the Project	3
4	The Summarizer Model	3-5
5	BERT Embeddings	6
6	HNSWlib	7
7	Text Extraction	8-11
8	Generating Text Summaries	12
9	User Interface	13-14
10	Future Work	15

# About the Project

It can :

- Summarize Research articles or any other PDFs
- Summarize the user's text
- Summarize web pages, including news and blogs or YouTube videos.

## The Summarizer Model

The heart of the model - Longformers!

This app uses the ***Longformer Encoder-Decoder (LED) for Narrative-Esque Long Text Summarization*** architecture to summarize texts retrieved from PDFs, Documents, or web pages.

We had a couple of options in hand before we started our build:

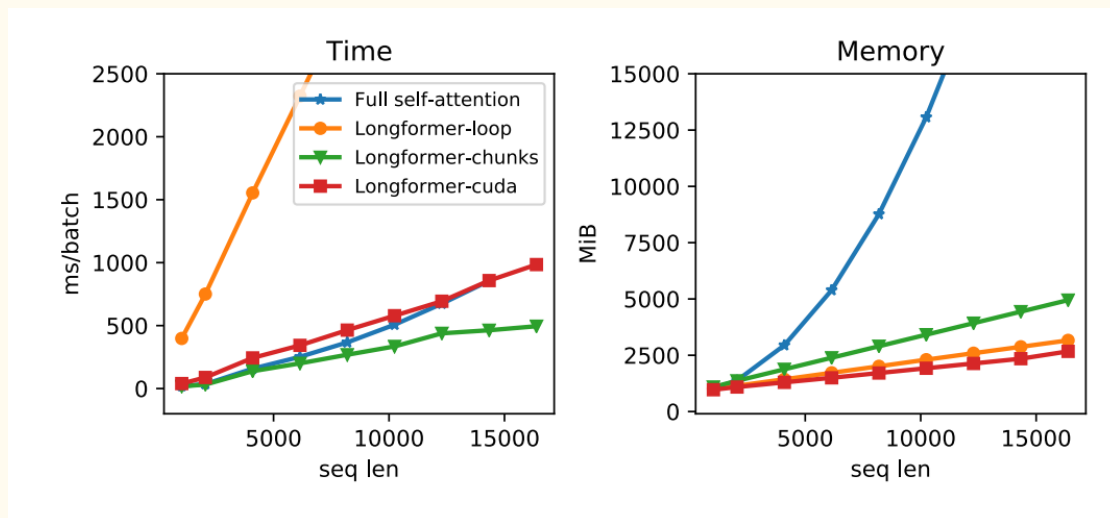
1. Google's T5
2. Facebook's BART
3. **Longformer Encoder-Decoder (LED) for Narrative-Esque Long Text Summarization**

We chose the third one because of the following reasons:

- **Performance and Quality:**

Transformer-based models cannot process long sequences due to their self-attention operation, which scales quadratically with the sequence length. Longformers, on the other hand, can summarize text with a more significant number of words.

The original Transformer consisted of an encoder-decoder architecture for sequence-to-sequence tasks, such as summarization and translation. While encoder-only Transformers are effective on various NLP tasks, pre-trained encoder-decoder Transformer models (e.g., BART and T5 have achieved strong results on tasks like summarization.) Yet, such models can't efficiently scale to seq2seq tasks with longer inputs.



. Longformer's attention mechanism is a drop-in replacement for the standard self-attention and combines local windowed attention with task-motivated global attention. It's performance and memory usage is more efficient compared to traditional self-attention.

- **Training**

The pretraining dataset used for this model was one of the primary factors influencing our decision. LED-base-book-summary is trained on diverse sources, including books and research articles, providing a more extensive and comprehensive understanding of language. This diverse training data helps the model capture a broader range of contexts and improves its ability to generate meaningful summaries.

## The architecture of Longformer Encoder-Decoder (LED)

The model uses an encoder-decoder architecture for *seq2seq* modeling. What differentiates it from a traditional transformer is **'local'** and **'global'** attention.

- **Self-Attention in Traditional Transformers:** In traditional Transformers, self-attention simultaneously operates on the entire input sequence. It computes the attention scores between each token in the sequence and all other tokens, capturing their dependencies and relationships. This allows the model to establish both local and global connections. However, the quadratic complexity of self-attention concerning the sequence length could be more efficient for long sequences.
- **Local Windowed Attention:** Local windowed attention allows the model to attend to a subset of the input sequence, or a "local window," rather than the entire sequence simultaneously. This approach reduces the computational complexity

and memory requirements compared to full self-attention. Limiting the attention scope allows the model to capture local dependencies within a fixed window size.

- **Global Attention:** As the name suggests, the model can attend to the entire input sequence. It is often used with local windowed attention to capturing local and global dependencies. While local attention focuses on nearby tokens, global attention provides a broader context by attending to all tokens in the sequence.

**Training Data:** This model is pre-trained on the *Booksum* dataset by Salesforce.

BookSum is a new collection of datasets specifically designed to address the limitations of existing text summarization datasets. Unlike other datasets mainly consisting of short-form source documents with limited dependencies, BookSum focuses on long-form narrative sources such as novels, plays, and stories. This dataset includes human-written summaries at three levels of granularity: paragraph, chapter, and book.

**Number of parameters:- 162 M params**

Trained on	Booksum dataset ( <a href="https://github.com/salesforce/booksum">salesforce/booksum (github.com)</a> )
Model Size	~1.6G
Number of parameters	162M params

```
class LED:
    def __init__(self):
        pass
    def summarize(self, text, max_length=500):
        hf_name = 'pszemraj/led-base-book-summary'
        summarizer = pipeline(
            "summarization",
            hf_name,
            device=0 if torch.cuda.is_available() else -1,
        )

        result = summarizer(
            text,
            min_length = 8,
            max_length = max_length,
            no_repeat_ngram_size=3,
            encoder_no_repeat_ngram_size=3,
            repetition_penalty=3.5,
            num_beams=4,
            do_sample=False,
            early_stopping=True,
        )

        return result[0]['summary_text']
```

# BERT Embeddings

BERT (Bidirectional Encoder Representations from Transformers) is a state-of-the-art language representation model that has revolutionized various natural language processing (NLP) tasks. BERT embeddings capture contextual information by considering the entire sentence rather than individual words. This model can be fine-tuned on specific downstream tasks, such as text summarization, to generate meaningful representations of input text.

- BERT Pretraining: BERT is initially pre-trained on large corpora using masked language modeling and next-sentence prediction tasks. This pretraining helps BERT to learn general language representations.
- Fine-tuning BERT: After pretraining, BERT can be fine-tuned on specific tasks, such as text summarization. During fine-tuning, BERT is trained on labeled data specific to the task at hand, which helps it learn task-specific features and nuances.
- BERT Embeddings: The output of BERT's hidden layers can be used as embeddings for downstream tasks. These embeddings capture the contextual understanding of the input text, enabling better representation of words and sentences.

```
class Embeddings:
    def __init__(self):
        logging.info("Initialising BERT transformer.")
        try :
            from transformers import BertTokenizer,BertModel
            self.model = BertModel.from_pretrained("bert-base-uncased")
            self.tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')

        except Exception as e:
            raise CustomException(e,sys)

    def tokenize_sentences(self,text):
        try :
            token=self.tokenizer.tokenize(text)
            token = self.tokenizer.convert_tokens_to_ids(token)
            return token
        except Exception as e:
            raise CustomException(e,sys)

    def embeddings(self,text):
        try :
            token = self.tokenize_sentences(text)
            with torch.no_grad():
                output = self.model(torch.tensor(token).unsqueeze(0))
                return output[1][0].tolist()
        except Exception as e:
            raise CustomException(e,sys)
```

# HNSWlib for KNN based on Cosine similarity.

HNSWlib (Hierarchical Navigable Small World) is a library that provides an efficient data structure for approximate nearest neighbor search. In the context of text summarization, HNSWlib can perform a cosine similarity-based k-nearest neighbor (KNN) search on BERT embeddings.

- Building the HNSW Index: BERT embeddings of the documents in the corpus are indexed using the HNSWlib library. The HNSW index structure organizes the embeddings for efficient nearest neighbor search.
- Cosine Similarity: Cosine similarity is a widely used metric to measure the similarity between two vectors. In the case of BERT embeddings, cosine similarity helps determine the semantic similarity between documents.
- K-Nearest Neighbor Search: Given a query document, the HNSW index efficiently finds the k nearest neighbors in the embedding space based on cosine similarity. These nearest neighbors can then generate a summary by extracting meaningful information from similar documents.

```
text_lines = text.split('.')
emb= Embeddings()
input_embeddings=[]
for line in text_lines:
    input_embeddings.append(emb.embeddings(line))
logging.info('Generated embeddings for input text.')
input_embeddings=np.array(input_embeddings)
key_embed = emb.embeddings(key_word)
key_embed=np.array(key_embed)
logging.info('Generated embeddings for keyword.')

num_elements = len(input_embeddings)
ids = np.arange(num_elements)
dim = input_embeddings.shape[1]

p = hnswlib.Index(space = 'cosine', dim = dim)
p.init_index(max_elements = num_elements, ef_construction = 200, M = 16)
logging.info('Initialised hnswlib.')
p.add_items(input_embeddings, ids)
p.set_ef(50)
logging.info('Added items to hnswlib.')

labels, distances = p.knn_query(key_embed, k = min(20,num_elements))
```

# Text Extraction

**1. PDFs: PyPDF2** is a Python library for working with PDF files. We used it to extract text from PDF documents.

```
def gen_para_pdf(file_name):
    try :
        with open(file_name, 'rb') as file:
            pdf_reader = PyPDF2.PdfReader(file)
            paragraphs=[]
            for page_number in range(len(pdf_reader.pages)):
                page = pdf_reader.pages[page_number]
                text = page.extract_text()
                paragraph = text.split('\n')
                for para in paragraph:
                    paragraphs.append(para)
            return paragraphs
    except Exception as e:
        logging.info("PDF File not found.")
        paragraph=[]
        return paragraph
```

**2. Word Files:** We used **python-docx**, another Python library for extracting text from Microsoft Word (.docx or .doc) files.

```
def gen_para_doc(file_name):
    try:
        doc = docx.Document(file_name)
        text = []
        for paragraph in doc.paragraphs:
            if (paragraph) :
                p_text = paragraph.text
                if len(p_text):
                    text.append(p_text)
        return text
    except Exception as e:
        logging.info("Doc File not found.")
        paragraph=[]
        return paragraph
```



**3. Text Files:** Python has inbuilt features to deal with text files using the **open** command.

```
def gen_para_txt(file_name):
    try:
        paragraphs=[]
        with open(file_name, 'r') as file:
            text = file.read()
            paragraph = str(text).split('\n')
            for para in paragraph:
                if (para):
                    paragraphs.append(para)
        return paragraphs
    except Exception as e:
        logging.info("File not found.")
        paragraph=[]
        return paragraph
```

A general function that checks the filetype by the extension of the file and redirects the file to the specific function:

```
def gen_para_file(file_name):
    file_name = os.path.join(folder, file_name)
    file_extension = os.path.splitext(file_name)[1]
    if (file_extension=='.pdf'):
        return gen_para_pdf(file_name)
    elif (file_extension == '.doc' or file_extension=='.docx'):
        return gen_para_doc(file_name)
    elif (file_extension=='.txt'):
        return gen_para_txt(file_name)
    logging.info("Paragraph of {} is generated.".format(file_name))
```

**4. Web Pages:** We used Python's Request and BeautifulSoup libraries to scrape web pages and extract their headings and paragraph elements.

```

import requests
from bs4 import BeautifulSoup

def get_website(url):
    try :
        response = requests.get(url)
        if response.status_code == 200:
            soup = BeautifulSoup(response.content, 'html.parser')

            article_title = soup.find('h1').get_text()
            paragraphs = soup.find_all('p')

            article_text = '\n'.join([p.get_text() for p in paragraphs])

            return article_title+ " "+article_text
        else:
            logging.info('URL not found. Error {}'.format(response.status_code))

    except Exception as e:
        raise CustomException(e,sys)

def get_yt_id(url, ignore_playlist=False):

```

To process the URL of a youtube video, we retrieved the video ID, typically found in the YouTube video URL after the "v=" parameter. We then retrieve the transcript of the video using the YouTube Transcript API. We extracted and displayed the text from the transcript in the text area. On clicking the Generate Summary button, a summary can be generated.

```

def get_yt_id(url, ignore_playlist=False):
    query = urlparse(url)
    if query.hostname == 'youtu.be': return query.path[1:]
    if query.hostname in {'www.youtube.com', 'youtube.com', 'music.youtube.com'}:
        if not ignore_playlist:
            with suppress(KeyError):
                return parse_qs(query.query)['list'][0]
        if query.path == '/watch': return parse_qs(query.query)['v'][0]
        if query.path[:7] == '/watch/': return query.path.split('/')[1]
        if query.path[:7] == '/embed/': return query.path.split('/')[2]
        if query.path[:3] == '/v/': return query.path.split('/')[2]

    def get_yt(url):
        try :
            response = requests.get(url)
            if response.status_code == 200:
                id = get_yt_id(url)
                transcript = YouTubeTranscriptApi.get_transcript(id)
                yt_text = ' '.join([item['text'] for item in transcript])

                return yt_text
            else:
                logging.info('URL not found. Error {}'.format(response.status_code))

        except Exception as e:
            raise CustomException(e,sys)

```

Enter your URL in the search box and click submit:

https://www.youtube.com/watch?v=...

Submit

Example of Summary extracted from YouTube Transcript:

Text Summarizer using LED

Upload File : Choose File No file chosen Submit

Choose files: Select Files Submit

https://www.youtube.com/watch?v=...

Submit

Text

Text summary

we watched the BBC interview with musk right and he goes up and the guy opens up the question with why did you agree to do this interview he says well I don't know what's the name of the BBC and he's trying to mock then and all this stuff and there's a part of it where you know he calls him out for the mistakes they made and if you go to the bottom of the interview on BBC you'll see comment sections saying why'd you cut that out because musk puts it on Twitter but they didn't put it in the interview then you watch the interview with the uh uh uh with you the first question they asked you was uh uh what's the first question hey you've been accused of serious crime rape or something like that uh you know right off the bat they ask that question of you and then you see the interview Lucy Williamson I believe her first question was you're facing some very serious allegations have you raped anybody that's the open question that they ask you and the video they put up the first one they take it down the second video they put up that's 12 minutes like the highlight one they turn the comment section is open there's 80 000 plus comments there then you look at Phillip Schofield and you just brought a Phillip Schofield for people that don't know Phyllis Schofield do you mind explaining to people who he is yeah he was a TV presenter in England he was very famous he ran the morning show and he was grooming children for a very long time and all the staff knew about it and the people who worked on the show with them were being groomed by him and everybody knew and it was all a big ha ha joke and now he's come out saying oh please don't pick on me I feel sad and the media is saying I'll leave the guy alone he is when you go look up his Wikipedia it says he rose to prominence as children's BBC continuity presented from 85 to 87, then went on to do programs on BBC and ITV for going live this morning dancing on ice All-Star Mr and Mrs thecude and bunch of other things and they interview him okay while he's gone through the mass of you know after 27 years of being married this one guy that he groomed since 10 years old his name is uh McGreevy I want to say something McGreevy uh that he's uh going through the process Matthew they met at 10 years old at a theater group at 15 years old Philip follows him on Twitter while the kid is 15 he follows him on Twitter the guy celebrates it long story short Ruth Langford they work together with she follows the complaint then right after finally complained he McGreevy gets fired then he has to take a break then he comes out after 27 years telling his wife you know and I'm gay I'm coming out of the closet he's got two daughters but there's a part of it where Matthew calls his wife and says hey your husband and I had an affair together when I was 19 the first

In this short program, the narrator explains how he came to be involved in child-rearing. He says that he was once famous for running a morning show with a young boy who was groomed and then later became a reality show on BBC1 where he was accused of raping a 15-year-old girl. Since then, he has come out as having had affairs with multiple women over the course of his career. The narrator also tells the audience that there are three main camps against him: 1) the Roman Catholic Church, 2) the Reform Movement, and 3) the anti-Trump movement. In order to counter these criticisms, the author uses an application that allows him to listen directly to influencers throughout the world. Through this application, the writer will give you access to all the latest news and information from social media.

Select number of characters : 250

Generate Summary

Reset

keyword to summarize

Download as TEXT File

Example of Summary extracted from a Blog Post:

Text Summarizer using LED

Upload File : Choose File No file chosen Submit

Choose files: Select Files Submit

web address to summarize

Submit

Text

Text summary

"Across the Spider-Verse" and the Latino legacy of Spider-Man The Conversation U.S. Follow The Conversation -- 17 Listen Share As a Latino literature and media scholar, a lifelong gamer and a Guatemalan-American girl whose dad read her comic every night, I quickly became a fan and then scholar of Miles Morales, the Afro-Puerto Rican Spider-Man who first appeared in comic book form in 2011's "Ultimate Fallout #4." Just seven years after his introduction, Morales swung into theaters in "Spider-Man: Into the Spider-Verse," a visually stunning, 3D-animated film that won an Academy Award for best animated feature. Now, its sequel, "Spider-Man: Across the Spider-Verse," features two Latino Spider-Men in starring roles. Irish-Latino Spider-Man Miguel O'Hara of "Spider-Man 2099," voiced by Oscar Isaac, is jumping into the fray. And although he was a well-received Spider-Man as a Marvel comic book character in the 1990s, there's a good chance you've never heard of him. Latino characters, particularly ones who have a starring role, have traditionally been underrepresented in mainstream comics. Marvel's first Latino hero, Hector Ayala, debuted in 1975, after the success of "Black Panther." Written by Bill Mantlo and drawn by legendary comic artist George Pérez, Ayala, known as White Tiger, was a Puerto Rican college student living in New York. His powers came from a magical amulet that bestowed him with speed and martial arts expertise. As Latino comics scholar Frederick Luis Aldama argues, Mantlo and Pérez avoided many of the stereotypes that plagued Latinos in comics, which often cast Latinos as criminals or drug dealers. Later iterations of White Tiger included his niece Angela del Toro and his sister, Ava Ayala.

This is the 17th episode of a new series called "The Conversation" and it focuses on the Latino characters who have starring roles in mainstream comic books from the past seven years. The first issue was released in 1975; the second is out in theaters in 2011. It features Miles Morales as his Irish-pan American character, White Tiger. O'Hara also appears in the movie. Marvel licensed their Spanish-language Spider-man translation rights to publish Spanish translations for Latin American countries. This means that Latin Americans are able to watch these stories without ever knowing that they are coming from different parts of the universe. In fact, this is the first time anyone has seen an animated movie with two Latino characters in starring positions. That's right, there are two Latina heroes in the movies. First, Marvel Studios licensed their English-published Spanish-translation rights to translate those stories into Latin America. Then, Marvel Comics licensed their Mexican-English-printing rights to do the same. And finally, Marvel introduced Miles Morales. Wait, isn't that confusing? You might be wondering: Why didn't Marvel just announce Miles Morales back in the day when everyone thought he was Black? Well, according to some scholars, it wasn't because

Select number of characters : 250

Generate Summary

Reset

keyword to summarize

Download as TEXT File

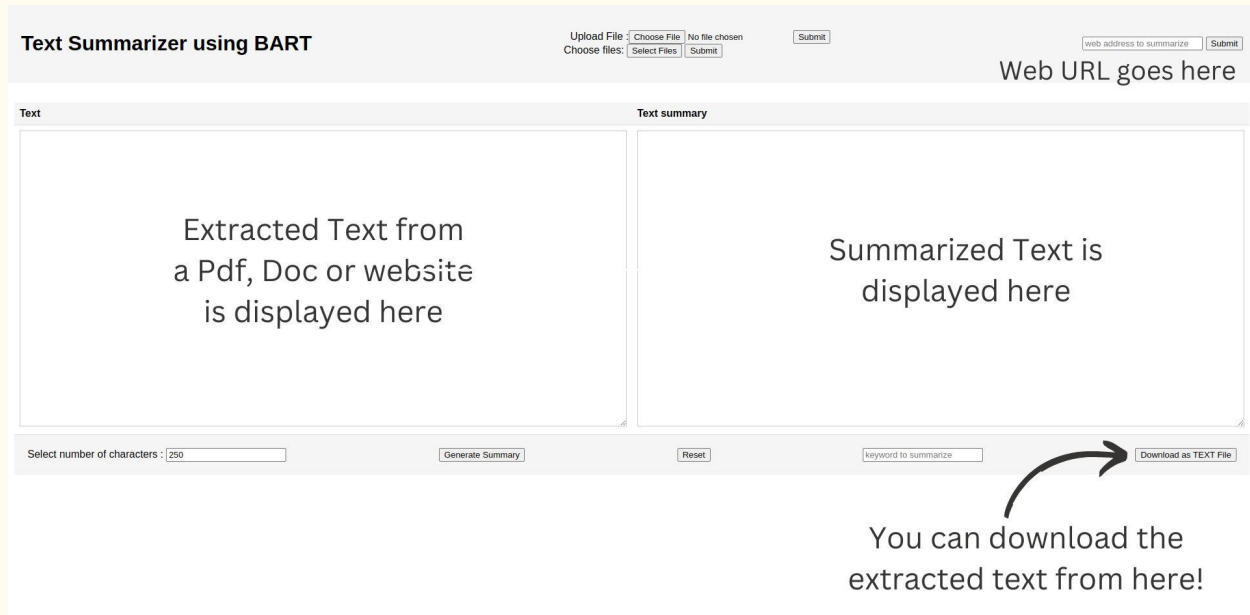
# Generating Text Summaries

By combining BERT embeddings and HNSWlib for cosine KNN, a generative artificial intelligence system for text summarization can be developed. The following steps outline the process:

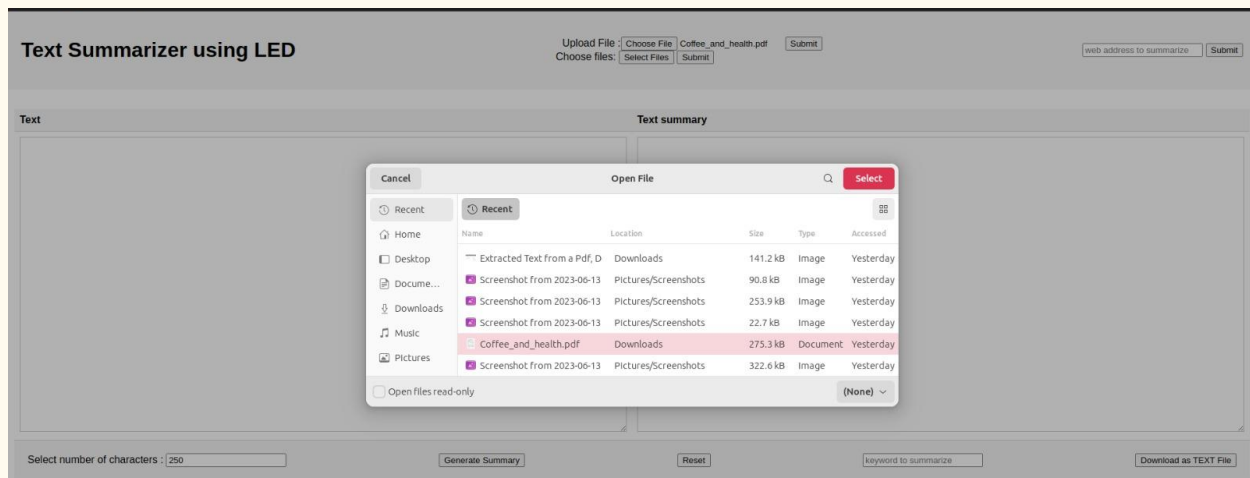
- **Preprocessing:** The input text, such as research papers, news articles, or documents, undergoes preprocessing to perform sentence segmentation and tokenization.
- **BERT Embeddings:** The preprocessed text is passed through a fine-tuned BERT model to generate embeddings that capture the contextual understanding of the input.
- **HNSW Indexing:** The BERT embeddings of the documents are indexed using HNSWlib to create an efficient search index.
- **Cosine KNN Search:** Given a query document, the BERT embeddings perform a cosine similarity-based KNN search on the HNSW index, retrieving the most similar documents.
- **Retrieval:** The top 20 sentences similar to the query are retrieved from the HNSWlib index and post-processed to form a document.
- **Summarization:** Extractive or abstractive summarization techniques are applied to the retrieved documents using LED to generate concise summaries.

# User Interface

How to run: Look up the Readme file in the GitHub repo for instructions to run it. The general layout of our app is as follows :



## Uploading Files :



Use the choose file button to choose any PDF/DOC file from a directory. Alternatively, you can paste text in the text field. You can select multiple previously uploaded files using the “Choose Files” button.

After uploading a file, press “Submit,” and the text will be extracted from it (PDF, DOC, Website, or Youtube video) and displayed in the text field. You may then specify the number of words you desire in the summary. It is set to 250 by default.

## Generate Summary :

### Text Summary from PDF

Press the “Generate Text” button and wait a few minutes (Performance depends on the local system’s hardware). You will get the summarized text in the Text Summary field.

**Also, if you want to store the summarized text in a Txt file, press the Download as TEXT button in the bottom left corner.**

# Future Work

With the constrained time frame, we could only implement some ideas in our minds. However, the following features can make this app better:

- **Integration of support for multiple languages:** Users can generate summaries for different languages in other languages.
- **Speech to Text:** Users can directly feed in an audio file to summarize.
- **Text to Speech:** Users who hate reading can generate audio for the summary.