**This project has been carried out under the supervision of Professor Srinath Srinivasa**

❖ Aparajita Kumari(MT2023037)
❖ Nandini Yadav(MT2023096)

# Nested Summarization with Heading hierarchy

# Contents

# Overview

## Context

- Nowadays people are consumers of short type of content.

- Properly structured headings improve accessibility by providing clear navigation cues.
.

## Problem Statement

- Develop algorithms to analyze the content structure and identify hierarchical relationships between different sections and subsections.

- Develop algorithms to generate concise summaries that capture the main points and essential information by maintaining sense and relevance.

# Objectives

**Collecting Dataset**

Data collection and preparation of the dataset to contain the relevant features appropriately.

**Preprocessing**

Remove any irrelevant metadata or columns that are not needed for analysis.

Remove special characters, punctuation, and non-alphanumeric characters.

**Text Summarization**

Generate a summary by paraphrasing and synthesizing the main points of the text. This typically involves deep learning models such as sequence-to-sequence models with attention mechanisms e.g. BART

**Title Generation**

Review the summarized content and identify the main themes, topics, or ideas that are prominently featured. These could be the central concepts or takeaways from the summary.

# Intuition

- **Text summarization** will condense key information into succinct summaries, saving time and effort for users. Summaries capture essential points while eliminating redundant details, enhancing readability and comprehension.

- **By combining heading hierarchy with summarization**, users can quickly locate relevant sections using headings and access concise summaries to grasp the main ideas, reducing cognitive load and improving efficiency.

# Implementation

- **Lexical Processing**
  - Removing unwanted urls,tags and stopwords using nltk library.
  - We implemented canonicalisation with the help of WordNet Lemmatizer.

- **Syntactic Processing**
  - It performs part-of-speech tagging to identify the grammatical categories of words in the text.
  - Additionally, it uses dependency parsing to analyze the syntactic structure of sentences.

- **Semantic Processing**
  - An important aspect of this processing is fine-tuning the pre-trained BART (Bidirectional and Auto-Regressive Transformers) model for text summarization tasks.

  - Latent Dirichlet Allocation (LDA) is implemented for topic modeling.

# Lexical Processing

- Removing unwanted urls,tags and stopwords using nltk library.
- We implemented **canonicalisation** with the help of WordNet Lemmatizer.
- To streamline the preprocessing process, all the aforementioned techniques are integrated into a cohesive pipeline. This pipeline serves as a standardized workflow, ensuring that each document undergoes the same preprocessing steps consistently.
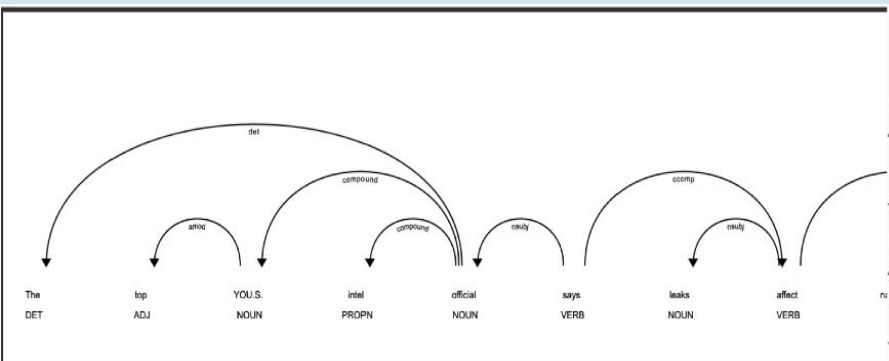


```
∨  Combining into a pipeline

[ ]  def process(stories):
        stories = remove_url(stories)
        stories = remove_html(stories)
        stories = remove_bracket(stories)
        stories = remove_digit(stories)
        stories = remove_underscore(stories)

        processed_list = []
        for story in stories:
            processed = expand_contractions(story)
            processed = tokenize(processed)
            processed = lower_case(processed)
            processed = remove_punctuation(processed)
            processed = remove_stopwords(processed)
            processed_list.append(processed)

        return lemmatizer(processed_list)
```

# Syntactic Processing

- The code utilizes natural language processing (NLP) tools such as NLTK, Stanza, and spaCy for further analysis.
- It performs **part-of-speech tagging** to identify the grammatical categories of words in the text.
- Additionally, it uses **dependency parsing** to analyze the syntactic structure of sentences.

- Finally, it visualizes the **dependency parse tree** using the displaCy library from spaCy

- The code also aims to implement **coreference resolution** to identify and link coreferent mentions in the text, enhancing semantic understanding.



```
tagged = pos_tag(word_tokenize(text))
tagged
```

```
[('Six', 'CD'),
 ('people', 'NNS'),
 ('have', 'VBP'),
 ('been', 'VBN'),
 ('killed', 'VBN'),
 ('in', 'IN'),
 ('three', 'CD'),
 ('incidents', 'NNS'),
 ('over', 'IN'),
 ('the', 'DT'),
 ('past', 'JJ'),
 ('two', 'CD'),
 ('months', 'NNS'),
 ('in', 'IN'),
 ('Honduras', 'NNP'),
 ('.', '.'),
 ('The', 'DT'),
 ('incidents', 'NNS'),
 ('happened', 'VBD'),
 ('during', 'IN'),
 ('the', 'DT'),
 ('course', 'NN'),
```

# Semantic Processing

- An important aspect of this mandate is **fine-tuning the pre-trained BART** (Bidirectional and Auto-Regressive Transformers) model for text summarization tasks.

- **Latent Dirichlet Allocation (LDA)** is implemented for topic modeling.

- LDA plays a crucial role in semantic processing by uncovering latent topics within text data, facilitating semantic analysis.



### Summarization

We use the pre-trained BART model and fine-tuned it for our dataset (CNN-Dailymail).

#### Fine-tuning BART

BART uses a standard Transformer-based neural machine translation architecture. It uses a standard seq2seq/machine translation architecture with a bidirectional encoder (like BERT) and a left-to-right decoder (like GPT). Here, BART is trained over the cnn-dailymail dataset.

```
[ ]
    !pip install bert-score -q
    !pip install blurr

Requirement already satisfied: blurr in /usr/local/lib/python3.10/dist-packages (0.4.1)
Requirement already satisfied: pyyaml in /usr/local/lib/python3.10/dist-packages (from blurr) (6.0.1)
Requirement already satisfied: python-dateutil in /usr/local/lib/python3.10/dist-packages (from blurr) (2.8.2)
Requirement already satisfied: docopt in /usr/local/lib/python3.10/dist-packages (from blurr) (0.6.2)
Requirement already satisfied: boto3 in /usr/local/lib/python3.10/dist-packages (from blurr) (1.34.91)
Requirement already satisfied: smart-open in /usr/local/lib/python3.10/dist-packages (from blurr) (6.4.0)
Requirement already satisfied: botocore<1.35.0,>=1.34.91 in /usr/local/lib/python3.10/dist-packages (from boto3->blurr)
Requirement already satisfied: jmespath<2.0.0,>=0.7.1 in /usr/local/lib/python3.10/dist-packages (from boto3->blurr) (1.0
Requirement already satisfied: s3transfer<0.11.0,>=0.10.0 in /usr/local/lib/python3.10/dist-packages (from boto3->blurr)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.10/dist-packages (from python-dateutil->blurr) (1.16.0)
Requirement already satisfied: urllib3!=2.2.0,<3,>=1.25.4 in /usr/local/lib/python3.10/dist-packages (from botocore<1.35.
```

```
[ ] !pip install ohmeow-blurr
```

```python
[ ] def get_topics(sentence_tokens, num_topics=1, num_words=2):
        dictionary = corpora.Dictionary(sentence_tokens)
        doc_term_matrix = [dictionary.doc2bow(doc) for doc in sentence_tokens]
        Lda = gensim.models.ldamodel.LdaModel
        ldamodel = Lda(doc_term_matrix, num_topics, id2word=dictionary, passes=30)
        return ldamodel.print_topics(num_topics, num_words)
```

# Deliverables

- ❏ Text Summarization
- ❏ Title Generation

# Text Summarization

- In our project, we fine-tune the pre-trained BART model on a specific dataset to adapt it to the summarization task.

- Fine-tuning involves updating the model's parameters using the dataset to optimize its performance for the summarization task.

- This process allows the model to learn task-specific patterns and features from the dataset, thereby improving its ability to generate accurate and coherent summaries.

- BART is a powerful transformer-based model known for its effectiveness in various NLP tasks, including text generation and summarization.

# Text Summarization

- The code snippet uses the BART model and a tokenizer (AutoTokenizer)to generate summaries for a list of input texts. It iterates through the input, processes them, and stores the generated summaries in the 'SystemSummary' list using beam search and specified constraints on length.

- The code generates summaries using BART model for provided text inputs, stores these in a DataFrame, 'Summaries', containing old and BART-generated summaries, then displays the initial rows of the DataFrame.

# Title Generation

- **Preprocessing:**
  - Before Title Generation ,we perform pre-processing on the summarized data in order to get rid of unwanted words such as stopwords , special characters etc.

# Title Generation

- **Sentence Clustering:**
  - It facilitates text summarization by identifying key themes or topics within the document.
  - It enables topic modeling by grouping sentences related to the same latent topic.
  - It improves document organization and comprehension by structuring sentences into coherent groups.

- **TF-IDF Vectorization:**
  - We employ TF-IDF vectorization to represent each sentence numerically based on the importance of its terms relative to the entire document collection.
  - Following TF-IDF vectorization, each sentence is represented as a numerical vector reflecting the importance of its terms within the document collection.

## Sentence Clustering

Since this is an unsupervised task, we use hierarchical agglomerative clustering. Tf-Idf is used to obtain vectors corresponding to the sentences. The optimal number of clusters is found using hierarchy.linkage() which gives the clusters at every iteration of the hierarchical clustering.

### ⌄ Tf-idf

```
[ ]  from sklearn.feature_extraction.text import TfidfVectorizer
```

```
[ ]  tfidf = TfidfVectorizer(analyzer='word', min_df=0)
```

```
[ ]  tfidf_wm = tfidf.fit_transform(preprocessed_sentences)
     tfidf_tokens = tfidf.get_feature_names_out()
     X = pd.DataFrame(data=tfidf_wm.toarray(), columns=tfidf_tokens)
```

# Title Generation

- **Hierarchical Clustering:**
  - With the TF-IDF representations in hand, we then apply hierarchical clustering to group similar sentences together.

  - This clustering approach allows us to iteratively merge clusters based on their similarity until all sentences belong to a single cohesive cluster.

  - **Agglomerative Clustering** is a hierarchical clustering algorithm used for grouping data points into clusters.

  - The algorithm proceeds by continuously joining clusters based on a linkage criterion, such as "ward" linkage, which minimizes the variance of merged clusters.

# Title Generation

- **Topic Modeling:**
  - Here, we are extracting words that frequently occur in the summary and generate heading out of it.
  - **LDA** represents documents as distributions over topics and topics as distributions over words.
    - It defines two probability distributions:
      - The distribution of words in documents given the topics.
      - The distribution of topics in documents.
  - **LDA** iterates over each document in the corpus and:
    - Assigns each word in the document to a topic randomly.
    - Adjusts these assignments based on the observed words in the document and the topics assigned to them.
  - This algorithm assigns a probability score to each word for each topic, indicating the likelihood of that word belonging to the topic.



∨ Title Generation/ Topic Modeling

```
[ ] import nltk
    from nltk.corpus import stopwords
    import gensim
    from gensim import corpora
```

```
[ ]
    nltk.download('stopwords')

    [nltk_data] Downloading package stopwords to /root/nltk_data...
    [nltk_data]   Package stopwords is already up-to-date!
    True
```

```
[ ]
    stopwords = set(stopwords.words('english'))
    def remove_stopwords(tokens):
        return [word for word in tokens if word not in stopwords and word]
```

```
[ ]
    def clean_for_topic_modeling(docs):
        #
        return [remove_stopwords(doc) for doc in docs]
```

```
def topic_for_para(doc, num_topics=1, num_words=2):
    sentences = split_into_sentences(doc)
    sentence_vectors = sentences_to_vectors(sentences)
    sentence_tokens = clean_for_topic_modeling(sentence_vectors)
    return get_topics(sentence_tokens, num_topics, num_words)
```

```
[ ]
    topic_for_para(document, num_topics=1, num_words=3)

    [(0, '0.031*"health" + 0.024*"foods" + 0.021*"junk"')]
```

# Challenges Faced

- We encountered difficulties when attempting to summarize our test dataset using the BART+Seq2Seq tokenizer. The issue arose when trying to generate multiple summaries within the dataframe. Although we could not identify the problem in the code, the Rouge score indicated satisfactory performance.

- In our project, we faced memory constraints, which limited our ability to use a large dataset. As a result, we had to work with a smaller dataset to generate title hierarchies.

- While this approach allowed us to work within our memory limitations, it's important to acknowledge that the results may be influenced by the reduced scope of the dataset.

# Conclusion

- As we near the conclusion of this project, we have successfully condensed lengthy articles into concise paragraphs through summarization techniques.

- Additionally, we have established a structured hierarchy of headings by analyzing clusters of sentences that share similar content.

- This hierarchical organization is achieved by identifying and extracting frequently occurring words, which serve as the basis for categorizing and structuring the summarized content.

```python
for paragraph in paragraphs:
    paragraph['summary'] = summarize_doc(paragraph['text'])
    topics = topic_for_para(paragraph['summary'], 1, 3)
    print(get_first_topic(topics))
    paragraph['h2_heading'] = get_first_topic(topics)
```
```
['health', 'person', 'one']
['high', 'health', 'time']
```

```python
[ ] output_file = open("./output.md", "w")

[ ] def write_h1(f, h1_heading):
        f.write("# ")
        for h in h1_heading:
            f.write(h + " | ")
        f.write("\n")

[ ] def write_sections(f, paragraphs):
        for para in paragraphs:
            f.write("## ")
            for h in para['h2_heading']:
                f.write(h + " | ")
            f.write('\n')
            f.write(para['summary'])
            f.write('\n')

[ ] write_h1(output_file, h1_heading)
    write_sections(output_file, paragraphs)

[ ] output_file.close()

[ ] output_file = open("./output.md", "r")

[ ] print(output_file.read())

    # health | foods | junk |
    ## health | person | one |
    Health has a lot of components that carry equal importance. If even one of them is missing, a person cannot be completely
    ## high | health | time |
    Health was earlier said to be the ability of the body functioning well. However, as time evolved, the definition of heal
```

# References

❏ https://github.com/Jay-Suthar/TEXT-SUMMA RIZATION-USING-BART-T5-PROPHETNET-PEGASUS/blob/master/group-22-ir-project.ip ynb

❏ Topic modeling using LDA

# Collab Notebook

- Our collab notebook:.
- https://colab.research.google.com/drive/1GmdsD8Oe
  V9EnwhOiFCqKDjr9Loz1cxyt?usp=sharing

**Thank You**