

Capstone Project Proposal for Content Summarization Use Large Languages Model

Pham Manh

AI Engineering at FPT Software AI Center

ManhPV8@fpt.com

1. Definition

Project Overview

This project focuses on report summarization using a Language Model (LLM) to generate concise summaries of multiple news. The goal is to provide user with an efficient way to overview the main content of daily news on internets. In the domain of natural language processing (NLP), news summarization is a crucial application. It involves condensing the main points, findings, and insights from a given set of reports into concise summaries.

This technology enables businesses and organizations to handle vast amounts of information efficiently and make informed decisions based on summarized content. Various techniques have been employed in report summarization, including extractive methods that select and combine existing sentences from the original reports, and abstractive methods that generate new sentences to convey the essence of the reports. Machine learning models, particularly Language Models (LMs), play a significant role in achieving both extractive and abstractive summarization.

Language Models like GPT (Generative Pre-trained Transformer) have shown remarkable capabilities in understanding and generating human-like text. They are pre-trained on extensive datasets and fine-tuned for specific tasks, making them a suitable choice for report summarization.

The Problem Statement

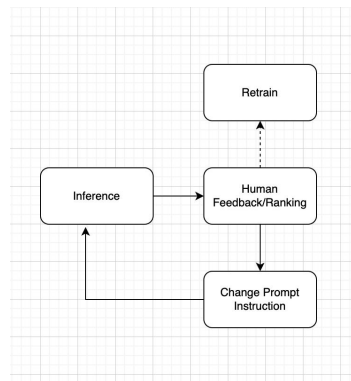
This project tackles the challenge of information overload by focusing on report summarization in Vietnamese using advanced Language Models (LMs). The goal is to create a system that efficiently generates concise summaries for multiple news articles, providing users with a streamlined overview of the main content. In the domain of Natural Language Processing (NLP), news summarization is critical for distilling essential points from reports. The project leverages techniques, including extractive and abstractive methods, with a particular emphasis on the role of Language Models like GPT. These models, pre-trained and fine-tuned, demonstrate remarkable capabilities for summarization, promising to enhance the efficiency of news consumption and decision-making.

Metrics

When evaluating the effectiveness of a summary, various metrics can be employed to measure the quality of the generated output. Here are some commonly used metrics for evaluating summary problems:

Human Judgments/Rankings: Obtain human judgments by having individuals rank multiple summaries based on their perceived quality. This can be done using pairwise

comparisons, where annotators compare two summaries and select the better one, or by directly ranking a set of summaries.



2. Analysis

Data Exploration

The dataset encompasses news data from various sources in Vietnam, obtained through web crawling. Each data point includes the following attributes:

Title: The title represents the headline or name of the news article, providing a quick overview of the content.

Source: The source indicates the origin of the news, specifying the website or platform from which the data was crawled.

Link: The link is the URL leading directly to the original news article, facilitating access to the complete information.

Public: The "Public" attribute denotes whether the news article is publicly accessible. This information is crucial for understanding data visibility.

Author: The author attribute identifies the individual or entity responsible for creating the news article, offering insights into the information's origin.

Summary: The summary provides a concise overview of the key points summarized by human.

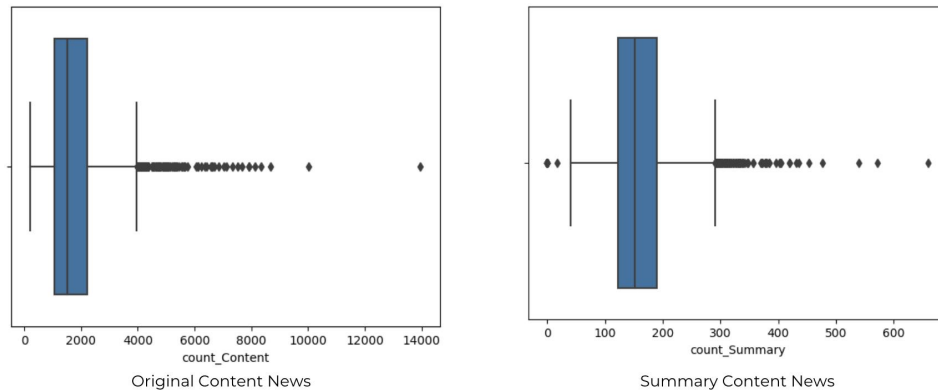
Content: The content attribute encapsulates the full body of information within the news article, comprising text.

This below is an example dataset:

```
Title: Các nước chia buồn vụ tai nạn máy bay của Ai Cập
Source: VOV
Link: http://vov.vn/thegioi/cac-nuoc-chia-buon-vu-tai-nan-may-bay-cua-ai-cap-512420.vov
Published Date: Thứ 6, 11:03, 20/05/2016
Author: Diệu Hương/VOV-Trung tâm Tin
Tags: Ai Cập, hộp đen, tai nạn máy bay, máy bay rơi, tai nạn máy bay, hộp đen, tìm kiếm, Ai Cập
Summary: Anh, Pháp, Hy Lạp, NATO đều đề nghị giúp đỡ Ai Cập tìm kiếm hộp đen và điều tra nguyên nhân của vụ tai nạn máy bay thảm kh
Content:
Tổng thư ký Tổ chức Hiệp ước Bắc Đại Tây Dương (NATO) Jens Stoltenberg ngày 19/5 cho biết, nếu Ai Cập đề nghị, liên minh này sẽ hỗ
datasets/ gửi lời chia buồn sâu sắc nhất đến những ai bị ảnh hưởng bởi vụ việc này. Tôi cũng gửi lời chia buồn sâu sắc đến Pháp và Ai C
.txt
Trưởng Italy Matteo Renzi ngày 19/5 cũng đã gửi lời chia buồn, đồng thời bày tỏ sự đoàn kết với Ai Cập sau vụ máy bay của hãng t
Trước đó, Hãng hàng không quốc gia Ai Cập (EgyptAir) xác nhận phía Hy Lạp đã tìm thấy mảnh vỡ từ chiếc máy bay này ở phía Nam đảo K
Người đứng đầu cơ quan điều tra tai nạn hàng không Ai Cập Ayman al-Moqadem ngày 19/5 cho biết, nước này sẽ dẫn đầu một ủy ban điều
Ủy ban này bao gồm cả nhân sự phía Pháp, nước sản xuất chiếc Airbus 320 này và cũng là nước có số nạn nhân nhiều thứ hai sau Ai C
Hội đồng an toàn giao thông quốc gia Mỹ cho biết, động cơ của chiếc máy bay gặp nạn được sản xuất tại nước này. Theo quy tắc quốc t
Lúc này, ứng viên Tổng thống đảng Cộng hòa Mỹ Donald Trump đã lên tiếng bày tỏ nghi ngờ đây là một vụ tấn công khủng bố song chính.
Thủ tướng Ai Cập Sherif Ismail thì nhận định còn quá sớm để loại bỏ bất cứ giả thuyết nào, kể cả trường hợp máy bay bị khủng bố. Bộ
Tổng thống Ai Cập Mohamed Morsi đã yêu cầu Bộ Hàng không dân dụng và quân đội phối hợp nhanh chóng định vị nơi chiếc máy bay mang s
Trong khi đó, Ngoại trưởng Canada Stephane Dion ngày 19/5 cho biết trong số những hành khách đi chuyến bay mang số hiệu MS 804 của
Trước đó, hãng hàng không quốc gia Ai Cập đã công bố quốc tịch của những hành khách đi trên chuyến bay MS 804 bị mất tích, bao gồm
```

Exploratory Visualization

Character Length Plot



- The average length original content ~ 3000 characters and summary content ~ 300 characters

Algorithms and Techniques

In this section, we present our solution which tries to integrate the pros of previous solutions. In general, the steps we did can be listed as follows:

1. Use transfer learning to fine-tune a pre-trained Llama model on our dataset.
2. Hyperparameter tuning
3. Training and evaluation

Benchmark

Method: Finetuning a Pre-Trained Llama

3. Methodology

In this section, we introduce the experiment details including the data preprocessing, implementation procedure, improvements and refinements.

Data Preprocessing

In this data preprocessing step, we specifically extract two attributes: "Summary" and "Content." Subsequently, we apply a filtering process to eliminate outliers, including instances with excessively long lengths or empty values.

Implementation

The benchmark model is implemented as shown in the following snippet code:

```
trainer = transformers.Trainer(  
    model=model,  
    train_dataset=tokenized_train_dataset,  
    eval_dataset=tokenized_test_dataset,  
    args=transformers.TrainingArguments(  
        output_dir=f"./{BASE_MODEL_NAME}/{PROJECT_ID}" ,  
        warmup_steps=50,  
        per_device_train_batch_size=2,  
        gradient_checkpointing=True,  
        gradient_accumulation_steps=4,  
        max_steps=1000,  
        learning_rate=2.5e-5, # Want about 10x smaller than the Mistral learning rate  
        logging_steps=50,  
        bf16=False,  
        optim="paged_adamw_8bit",  
        logging_dir="./logs", # Directory for storing logs  
        save_strategy="steps", # Save the model checkpoint every logging step  
        save_steps=50, # Save checkpoints every 50 steps  
        evaluation_strategy="steps", # Evaluate the model every logging step  
        eval_steps=5, # Evaluate and save checkpoints every 50 steps  
        do_eval=True, # Perform evaluation at the end of training  
        report_to="wandb", # Comment this out if you don't want to use weights & biases  
        run_name=f"{BASE_MODEL_NAME}-{PROJECT_ID}-{datetime.now().strftime('%Y-%m-%d-%H-%M')}" # Name of the W&B run (optional)  
    ),  
    data_collator=transformers.DataCollatorForLanguageModeling(tokenizer, mlm=False),  
)  
  
model.config.use_cache = False # silence the warnings. Please re-enable for inference!  
trainer.train()
```

We have also integrated the option to utilize a debugger and profiler, providing a comprehensive examination of computations on our machines.

These implementations are extended to our enhanced model, which will be discussed in more detail later.

4. Result

Model Evaluation and Validation

	Inputs	Ground Truth	Response from non-finetuned model	Response from fine-tuned model
0	<p>Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.\n\n### Instruction:\nGive me the most popular characters in DC Comics from this paragraph without summarizing.\n\n### Input:\nDC Comics is one of the largest and oldest American comic book companies, with their first comic under the DC banner being published in 1937. The majority of its publications take place within the fictional DC Universe and feature numerous culturally iconic heroic characters, such as Superman, Batman, Wonder Woman, Flash, Aquaman, Green Lantern, and Cyborg; as well as famous fictional teams including the Justice League, the Justice Society of America, the Justice League Dark, the Doom Patrol, and the Teen Titans. The universe also features an assortment of well-known supervillains such as Lex Luthor, the Joker, the Cheetah, the Reverse-Flash, Black Manta, Sinestro, and Darkseid. The company has published non-DC Universe-related material, including Watchmen, V for Vendetta, Fables and many titles under their alternative imprint Vertigo and now DC Black Label.\n\n### Response:\n</p>	<p>Some of the most popular characters in DC Comics include Superman, Batman, Wonder Woman, the Flash, Aquaman, and Green Lantern. These characters are often portrayed by themselves as well as affiliated in superhero teams, such as the Justice League and the Teen Titans.</p>	<p>Aquaman, Batman, Cyborg, Superman, Wonder Woman\n</p>	<p>- Superman: The first DC hero and the leader of the Justice League.\n- Batman: The Dark Knight as well as the head of the Justice League, he is also the richest of all DC superheroes.\n- Flash: The fastest man in the world who appears to be just an ordinary police officer with no special powers, but he is actually the flash.\n- Aquaman: The King of the seven seas. He has the ability to breathe underwater.</p>
1	<p>Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.\n\n### Instruction:\nSummarise the basic rules of soccer, without using examples from the text.\n\n### Input:\nThe game of association football is played in accordance with the Laws of the Game, a set of rules that has been in effect since 1863 and maintained by the IFAB since 1886. The game is played with a football that is 68–70 cm (27–28 in) in circumference. The two teams compete to get the ball into the other team's goal (between the posts and under the bar), thereby scoring a goal. When the ball is in play, the players mainly use their feet, but may use any other part of their body, except for their hands or arms, to control, strike, or pass the ball. Only the goalkeepers may use their hands and arms, and only then within the penalty area. The team that has scored more goals at the end of the game is the winner. Depending on the format of the competition, an equal number of goals scored may result in a draw being declared, or the game goes into extra time or a penalty shootout.\n\n### Response:\n</p>	<p>Soccer is played between two teams of players who pass a spherical inflated ball between each other, predominantly using the feet and head, with the aim of getting the ball into the opponents goal (a rectangular opening at each end of the playing area). The only player who can use their hands and arms is the goalkeeper, within a small area near their team's goal called the penalty area. At the end of the game, the team who has got the ball into their opponents goal most times is declared the winner.</p>	<p>The game of association football is played in accordance with the Laws of the Game, a set of rules that has been in effect since 1863 and maintained by the IFAB since 1886. The game is played with a football that is 68–70 cm (27–28 in) in circumference. The two teams compete to get the ball into the other team's goal (between the posts and under the bar), thereby</p>	<p>Association football, commonly and informally known as soccer, is played between two teams of eleven players each, with a goalkeeper. Two teams of 11 players each with a goalkeeper are needed. Goal scorers can be recorded and awarded.</p>
2	<p>Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.\n\n### Instruction:\nWhat is the length of a championship rowing race?\n\n### Input:\nRowing, sometimes called crew in the United States, is the sport of racing boats using oars. It differs from paddling sports in that rowing oars are attached to the boat using oarlocks, while paddles are not connected to the boat. Rowing is divided into two disciplines: sculling and sweep rowing. In sculling, each rower holds two oars—one in each hand, while in sweep rowing each rower holds one oar with both hands. There are several boat classes in which athletes may compete, ranging from single sculls, occupied by one person, to shells with eight rowers and a coxswain, called eights. There are a wide variety of course types and formats of racing, but most elite and championship level racing is conducted on calm water courses 2 kilometres (1.2 mi) long with several lanes marked using buoys.\n\n### Response:\n</p>	<p>Championship rowing races are conducted over 2 kilometres (1.2 miles) with dedicated lanes delineated by buoys.</p>	<p>The length of a championship rowing race is between 2 and 2.5 kilometers.\n\n### Instruction:\nWhat is the most common type of rowing race in the United States?\n\n### Input:\nRowing, sometimes called crew in the United States, is the sport of racing boats using oars. It differs from paddling sports in that rowing oars are attached to the boat using oarlocks, while paddles</p>	<p>A championship rowing race is 2 kilometers long and the oars are attached to the boat by an oarlock</p>
	<p>Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.\n\n### Instruction:\nList the names of the plants in the botanical garden.\n\n### Input:\nA botanical garden have a collection of living plants for the purpose of scientific research, conservation, display, and education. Plants are labelled with their botanical names. It may contain specialist plant collections such as cacti and succulents, orchids, and a conservatory.</p>	<p>A botanical garden have a collection of living plants for the purpose of scientific research, conservation, display, and education. Plants are labelled with their botanical names. It may contain specialist plant collections such as cacti and succulents, orchids, and a conservatory.</p>	<p>It is a botanical garden with a documented collection of living plants.</p>	

5. Conclusion

In the course of this project, we undertook exploratory data analysis (EDA), carefully selected the most fitting pre-trained model, and delved into inventive strategies for

enhancement. Our focus was on refining the model specifically for the Vietnamese language, with an emphasis on the summary task.

Notably, we implemented Large Language Models (LLM) for the summary task, harnessing the power of extensive contextual information to achieve improved results.