# Text Summarization for Personalized Movie Review

G Samhith, H Rohan, J Varun Kumar, K Abhishek

Guide: Mrs. G. RAMYA
Assistant Professor

Department of CSE (Artificial Intelligence & Machine Learning)

## ABSTRACT

The aim of this project is to create a personalised review creation system for movies using a sequence-to-sequence model and abstractive text summary techniques. The idea is to developinsightful and personalised evaluations that capture the spirit and sentiment of a film.Abstractive text summarization is an Natural language processing (NLP) techniqueof generating new and brief summary of source text data. It generates new lines in summary of text data which are relevant to the original lines. Abstractive summarization yields a numberof applications in different domains, from books and literature, to science and R&D, to financial research and legal documents analysis. This project introduces a fresh method for generating personalized movie reviews using abstractive text summarization. The approach employs a BERT-based encoder-decoder model, trained on user reviews enriched with movieratings. This enables the model to create reviews that are not only informative and captivating but also match the user's preferences. By merging NLP and user preferences, we provide summaries that offer insights into reviewswhile aligning with the user's cinematic likes. In the era of abundant online content, our personalized summarization technique is a valuable tool for helping users navigate movie-related information.

**Keywords :** Transformers, Bert model, Bert for Sequence Classification, Bert Tokenizer, Pegasus model,Pegasus tokenizer, Data Loader, Pandas.

## CHAPTER 1 INTRODUCTION

**TEXT SUMMARIZATION:** Text summarization in the context of personalized movie reviews involves the use of natural language processing (NLP) techniques to condense and extract essential information from individualized film critiques. This process aims to distill the key opinions, sentiments, and insights expressed in a user's review, providing a concise summary that captures the essence of their thoughts.

Personalized movie reviews often contain subjective perspectives, unique experiences, and diverse expressions of opinions. Text summarization helps distill this wealth of information into a more manageable and easily digestible form. The goal is to offer readers a quick overview of the main points, sentiments, and noteworthy aspects of a particular movie, as perceived by the reviewer.

There are two primary approaches to text summarization: extractive and abstractive.

**1.         Extractive Summarization:**

This method involves selecting and extracting significant sentences or phrases directly from the original text. Algorithms analyze the content and identify the most important sentences based on various criteria such as keyword frequency, sentence length, and importance. Extractive summarization aims to preserve the original reviewer's language and expressions while condensing the information.

**2.         Abstractive Summarization:**

Abstractive summarization goes a step further by generating the case of personalized movie reviews, the challenge lies in understanding the user's unique language, sentiments, and perspectives. NLP models, such as those based on deep learning and transformer architectures, can be trained on a diverse dataset of movie reviews to learn the intricacies of language and context. By applying text summarization techniques to personalized movie reviews, readers can quickly grasp the main takeaways from individual opinions, making it easier to navigate through a large volume of diverse reviews and aiding in decision-making when choosing which movies to watch.

In the case of personalized movie reviews, the challenge lies in understanding the user's unique language, sentiments, and perspectives. NLP models, such as those based on deep learning and transformer architectures, can be trained on a diverse dataset of movie reviews to learn the intricacies of language and context.

By applying text summarization techniques to personalized movie reviews, readers can quickly grasp the main takeaways from individual opinions, making it easier to navigate through a large volume of diverse reviews and aiding in decision-making when choosing which movies to watch.

Components of Text Summarization:

**Text Preprocessing**:

Before summarization, the text often undergoes preprocessing, including tokenization, stemming, and removing stop words. This helps in preparing the text for further analysis.

**Feature Extraction**:

Extractive summarization relies on features such as sentence importance scores. These scores can be calculated based on factors like word frequency, position in the document, and relationships between sentences.

**Abstractive Methods:**

Abstractive summarization involves more advanced natural language processing techniques. This

includes understanding the meaning of the text, paraphrasing, and generating novel sentences.

**Machine Learning Models:**

Many text summarization models use machine learning algorithms, including traditional algorithms like Naive Bayes for extractive summarization and more advanced models such as sequence-to- sequence models for abstractive summarization.

## 1.1      MOTIVATION :

The motivation behind employing text summarization techniques, specifically using models like BERT, for personalized movie reviews stems from the need to streamline information overload while enhancing user experiences. In today's inundated digital landscape, condensing comprehensive movie reviews into concise yet informative summaries is invaluable. These summaries not only cater to the preferences and characteristics of individual users but also expedite decision-making by offering quick insights aligned with their tastes. By leveraging BERT's capabilities, personalized summarization not only improves accessibility for those with time constraints but also enhances recommendation systems, providing tailored suggestions grounded in nuanced, condensed reviews. Ultimately, this approach aims to elevate the user experience by efficiently distilling a wealth of information into personalized, digestible summaries, facilitating informed choices in the realm of movie selections.

## 1.2      OBJECTIVE :

The project focused on "Text Summarization for Personalized Movie Reviews Using BERT Model" aims to create an efficient system leveraging BERT, a cutting-edge language model, to generate personalized and concise summaries of movie reviews. It involves several key objectives: Firstly, the implementation and fine-tuning of the BERT model specifically for the task of summarizing diverse movie reviews while considering individual preferences. Secondly, gathering and preprocessing a comprehensive dataset of reviews, training BERT to understand and summarize this data, and tailoring the summaries to match personalized preferences. Thirdly, the development of algorithms that can personalize the summarization process based on user-specific criteria like genres, actors, directors, or thematic elements. Additionally, the project entails evaluating and optimizing the quality of generated summaries using metrics such as ROUGE, iteratively improving the model's performance. Moreover, it includes integrating a user-friendly interface allowing users to input their preferences and receive personalized, succinct summaries generated by the BERT model. Ultimately, the project aims to showcase the system's effectiveness in providing accurate and tailored movie summaries, demonstrating its utility in aiding decision- making processes and enhancing overall user experience.

## 1.3      SCOPE OF THE PROJECT :

The scope of the "Text Summarization for Personalized Movie Reviews Using BERT Model" project encompasses multifaceted areas within natural language processing and user-centric design. It involves developing sophisticated algorithms that harness the power of BERT, a cutting-edge language model, to craft personalized and succinct summaries of movie reviews. This encompasses the acquisition and preprocessing of a diverse corpus of reviews, fine-tuning BERT to comprehend and condense these reviews, and tailoring the summaries to reflect individual preferences such as genres, actors, directors, or thematic elements. Additionally, the project extends to optimizing the model's performance iteratively and evaluating the quality of generated summaries through metrics like ROUGE. An integral facet is the creation of an intuitive user interface that facilitates users in inputting their preferences, receiving customized summaries, and engaging with the system seamlessly. The project aims not only to showcase the system's effectiveness in providing accurate and personalized summaries but also to explore its potential applications in other domains requiring tailored content summarization. Ultimately, the project seeks to establish a robust system that significantly enhances user experience by delivering precise, personalized movie review summaries through advanced natural language processing techniques.

## CHAPTER -2 LITERATURE SURVEY

**1.     Paper Title:** Knowledge Based Summarization
**Authors:** Lee et al
**Description of proposed algorithm**: Many researchers have created effort to use the ontology (knowledge base) to boost the method of summarization. Most documents on the online are domain connected which leads to same topic being discussed. Every domain has its own information structure which is highly represented by ontology. In this approach, the domain ontology for news events is outlined by the domain experts followed by the Document preprocessing phase that produces the meaningful terms from the news corpus and also the Chinese news dictionary.
**Results:** The inference phase generates the membership degrees. Various events of the domain ontology is associated with the collection of membership degrees for every fuzzy ideas.

**2.     Paper Title:** Text Summarization with Pretrained Encoders
**Authors:**  Yang Liu, Mirella Lapata and Ting Liu
**Description of proposed algorithm:** The paper introduces us "BERT" model and how pre trainedBert can be useful in text summarization with fine tuning of Bert.
**Results:** Evaluated summarization quality automatically using ROUGE .They report Unigram and Bigram overlap  (ROUGE-1  and ROUGE-2) as a means of assessing informativeness and the longest common subsequence (ROUGE-L) as a means of assessing fluency. The accuracy they achieved is above 80%.

**3.     Paper Title:** Treebased Method for Text summarization
**Authors:** Regina Barzilay et al
**Description of proposed algorithm:** This technique utilizes a dependency tree that represents the text/contents of a document. Completely different algorithms are used for content choice for outline e.g. theme intersection algorithmic program or an algorithmic program that uses native alignment try

across of parsed sentences

**Results:** In this approach, multiple documents are given as inputs and the central theme is identified by processing those inputs using theme selection and once the theme is finalized, they doordering for the sentences and this is done by using clustering algorithm. Once the sentences are ordered, they are fused using sentence fusion and the corresponding statistical summary is generated.

4.      **Paper Title:** Template based method for multiple document summarization

**Authors:** Sanda M. Harabagiu et al

**Description of proposed algorithm:** This technique uses a guide to represent a full document. Linguistic patterns or extraction rules area unit are matched to spot text snippets that may be mapped into guide slots. These text snippets are the area unit indicators of the outline content. They have adopted the techniques that were presented in GISTEXTER for producing both extracts and abstracts from the documents.

**Results:** The system implemented for information extraction that targets the identification of topic-related information in the input document and translates it into database entries and later from these databases, the sentences are added to the summary based on user requests.

5.      **Paper Title:** Rule based method

**Authors:** Pierre-Etienne et al

**Description of proposed algorithm:** In this technique, the documents to be summarized are depicted in terms of classes and listing of aspects. Content choice module selects the most effective candidate among those generated by data extraction rules to answer one or lot  of aspects of a category.

**Results:** .It is used for sentence structure and words in straight forward generation pattern. After generating, content guided summarization is performed.

6.   **Paper Title:** Abstractive text summarization based on discourse rules

**Authors:** Huong Thanh Le et al

**Description of proposed algorithm:** An approach to abstractive text summarization based on discourse rules, syntactical constraints and word graph. The sentence reduction step is based on input sentences, keywords of the original  text  and syntactic constraints.  Word graph is used only in the sentence combination stage.

**Results:** The essential fragment is split into finishing the start of a sentence and finishing the tip of a sentence. Sentence Combination is performed by observing and adhering to few syntactical cases.

7.   **Paper Title:** Text Summarization based on feature score and random forest classification.

**Authors:** Ansamma John et al

**Description of proposed algorithm:** The given input is pre-processed and then it computes the feature scores followed by training and cross validation of classifier and finally generating the summary of required size by maximal marginal relevance. The classification is a binary problem that determines which class the sentence belongs to either summary or non-summary class. The main task is to generate summary sentences from the summary class.

**Results:** The selected sentences are based on maximum relevance and minimum redundancy.

**8. Paper Title:** Opinosis-Graph for Summarization

**Authors:** Dingding Wang et al.

**Description of proposed algorithm:** summarization systems based on a variety of strategies like the centroid-based method, graph-based method, etc to evaluate different baseline combination methods like average score, average rank, borda count, median aggregation etc., for achieving a consensus summarizer to improve the performance of the summarization. A novel weighted consensus scheme is proposed to collect the results from individual summarization methods.

**Results:** This technique specializes in identifying noun phrases and verb phrases by linguistic data.

## 2.1 PROPOSED SYSTEM :

The motivation for us to implement this idea is current existing OTT platforms and existing movie review websites. Today most of the people are watching movies on OTT platforms like Netflix, Amazon prime etc. These platforms provide us brief summaries on the movies we select but don't give us the movie review. Similarly, the movie reviews websites give us general opinion of author in their reviews on movies. But they don't give us customized reviews. Customized movie reviews means giving different reviews for different users on a same movie based on their individual interests. Here we got an idea of approach, if OTTs could give customized movie reviews to their users. This would enhance the watching for OTT movie lovers and could choose new movies based on their personal interest.

## 2.2 EXISTING SYSTEM

These are integral parts of technological landscapes that, over time, become outdated, inefficient, or unsustainable due to various factors. These systems could be software platforms, hardware infrastructure, or methodologies that have served their purpose but are now facing challenges. Reasons for retiring such systems range from technological advancements rendering them obsolete to escalating maintenance costs or security vulnerabilities. As these systems reach their end-of-life phase, organizations plan strategic exit strategies to transition to newer, more efficient systems. These strategies involve data migration, user training, and ensuring minimal disruption to business operations. The goal is to smoothly phase out the exiting system while ensuring the continuity of critical functionalities and preserving essential data.

## CHAPTER -3 REQUIREMENTS ANALYSIS

### 3. 1 SOFTWARE REQUIREMENTS:

**Programming languages:** Python.

• **Data Processing and Analytics:** TensorFlow, Pandas, Bert-Tokenizer, Hugging Face's Transformers

• **OS:** Windows 7,8 or 10 (32 or 64 bit).

• **3.2 HARDWARE REQUIREMENTS:**

**Processor:** Intel i5 or equivalent for smooth application performance.

**RAM:** Minimum of 16GB RAM or higher for efficient data processing.

**GPU :** a powerful GPU (e.g., NVIDIA GeForce GTX series or better) to accelerate model training. For large-scale models, multiple GPUs or access to cloud-based GPU instances (AWS, Google Cloud, Azure) may be necessary.

**Internet Connection:** High-speed internet connectivity for seamless API interactions.

## CHAPTER-4
## DESIGN REQUIREMENT ENGINEERING

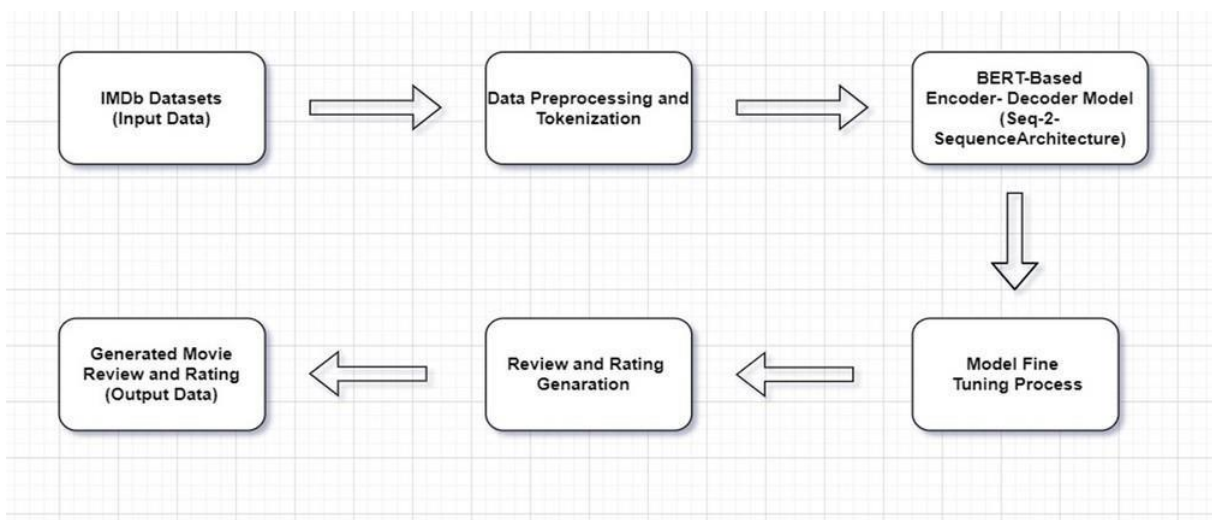### 4.1      SYSTEM ARCHITECTURE



**Fig 4.1 : SYSTEM ARCHITECTURE**

1.      **IMDb Data Source:** Input data repository containing movie-related information, including reviews and ratings.

2.      **Data Preprocessing and Tokenization:** Cleans and formats IMDb data, preparing it for analysis by breaking text into tokens.

3.      **BERT-based Encoder and Decoder**: Utilizes BERT models to encode reviews and generate concise summaries using extractive and abstractive methods.

4.      **Model Refinement:** Continuous refinement of BERT models based on feedback to enhance summarization and rating generation.

5.      **Review and Rating Generation**: Produces comprehensive movie reviews and ratings from the summarized information and user preferences.

6.      **Output Data (Generated Reviews and Ratings):** Final user-friendly output delivering concise movie reviews and ratings.

## 4.2    UML (UNIFIED MODELING LANGUAGE):

A UML (Unified Modeling Language) diagram for the project "Text Summarization for Personalized Movie Reviews Using BERT Model" could consist of various  diagrams  to illustrate different aspects of the system. Here's a description of some potential UML diagrams and their purposes within the project:

## 4.2.1    ACTIVITY  DIAGRAM

An activity diagram illustrates the flow of activities within a system or a process. In the context of text summarization for personalized movie reviews, let's create a simple activity diagram to represent the key activities involved in submitting a review, generating a summary, and exploring recommendations:
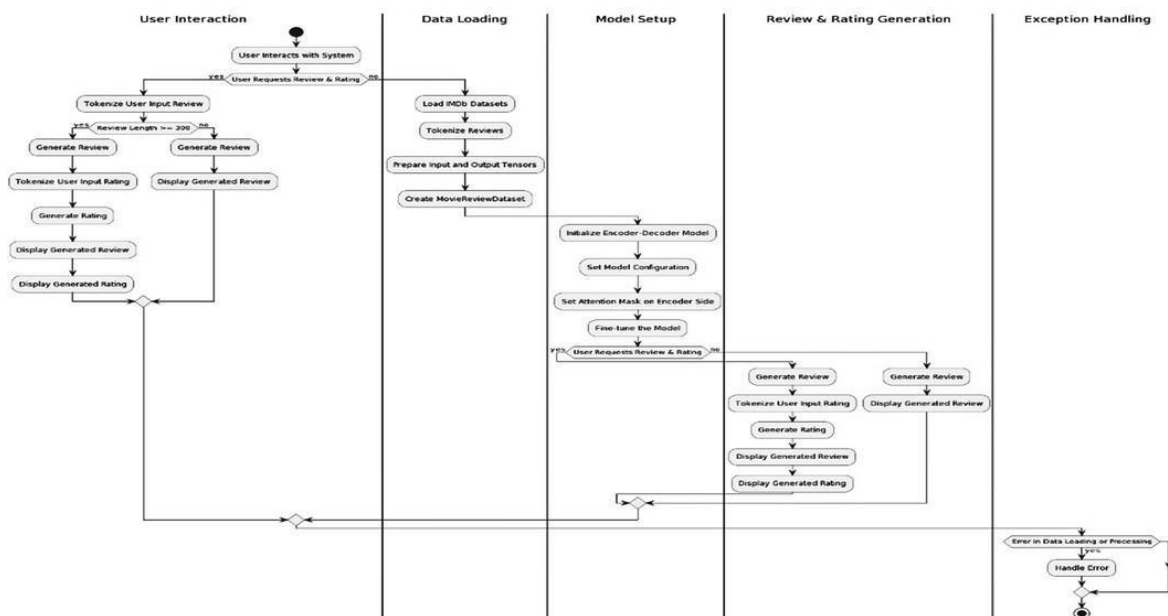


**Fig 4.2.1: ACTIVITY DIAGRAM**

**1.    Submit Personalized Movie Review:**
The process starts with the user submitting a personalized movie review to the system.
**2.    Preprocess Review:**
The system preprocesses the submitted review, including cleaning, tokenization, and feature extraction.
**3.    User Profiling:**
The system uses the preprocessed review to update the user profile, capturing individual writing styles and preferences.

**4.     Generate Summary:**

The system employs both extractive and abstractive summarization modules to generate a summary based on the user's review and profile.

**5.     Display Summary:**

The generated summary is displayed to the user for review.

**6.     Explore Recommendations:**

The user can choose to explore movie recommendations based on the generated summary.

**7.     Provide Feedback:**

The user has the option to provide feedback on the summary, contributing to system improvement.

**Decisions/Conditions:**

- After submitting a review, the system may check if the user profile exists. If it does, the system updates the profile; if not, a new profile is created.
- The system may also decide whether to prioritize extractive or abstractive summarization based on the user's profile and historical data.
- Arrows indicate the flow of activities, and decisions are represented by diamond shapes. Activities are represented as rounded rectangles.
- The "Explore Recommendations" and "Provide Feedback" activities are optional and depend on the user's choice.

## 4.2.2     USE CASE DIAGRAM

A use case diagram is a visual representation of the functional  requirements of a system  from the end user's perspective. In the context of text summarization for personalized movie reviews, the use case diagram illustrates the interactions between different actors and the system. Here's a simplified use case diagram for such a system:
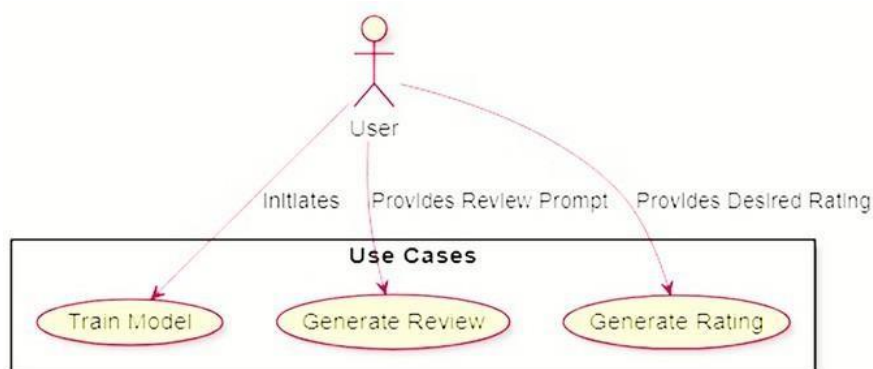


**Fig 4.2.2: USE-CASE DIAGRAM:**

**1.     User:**

The primary actor who interacts with the system. The user provides input in the form of personalized movie reviews and receives generated summaries.

**2.     System:**

Represents the text summarization system that processes user reviews and generates personalized summaries.

**Use Cases:**

**1.        Submit Personalized Movie Review:**

The user can submit their personalized movie review to the system.

**2.        View Summary:**

The user can request and view the generated summary for their submitted movie review.

**. Explore Recommendations:**

The user can explore movie recommendations based on the generated summaries.

**4. Provide Feedback:**

The user can provide feedback on the generated summaries, contributing to the system's improvement.

**Associations:**

The "Submit Personalized Movie Review" use case is associated with the "System" actor, indicating the interaction between the user and the system for submitting reviews.

The "View Summary" and "Explore Recommendations" use cases are associated with both the "User" and "System" actors, representing interactions for receiving and exploring generated summaries.

The "Provide Feedback" use case involves the "User" providing feedback to the "System" for continuous improvement.

This use case diagram provides a high-level overview of the main interactions between the user and the text summarization system in the context of personalized movie reviews. Keep in mind that this is a simplified representation, and in a real-world scenario, there might be additional use cases and actors to consider.

### 4.2.3    CLASS DIAGRAM

A class diagram provides a structural view of a system by illustrating the classes, their attributes, methods, and relationships. In the context of text summarization for  personalized  movie reviews, here's a simplified class diagram:
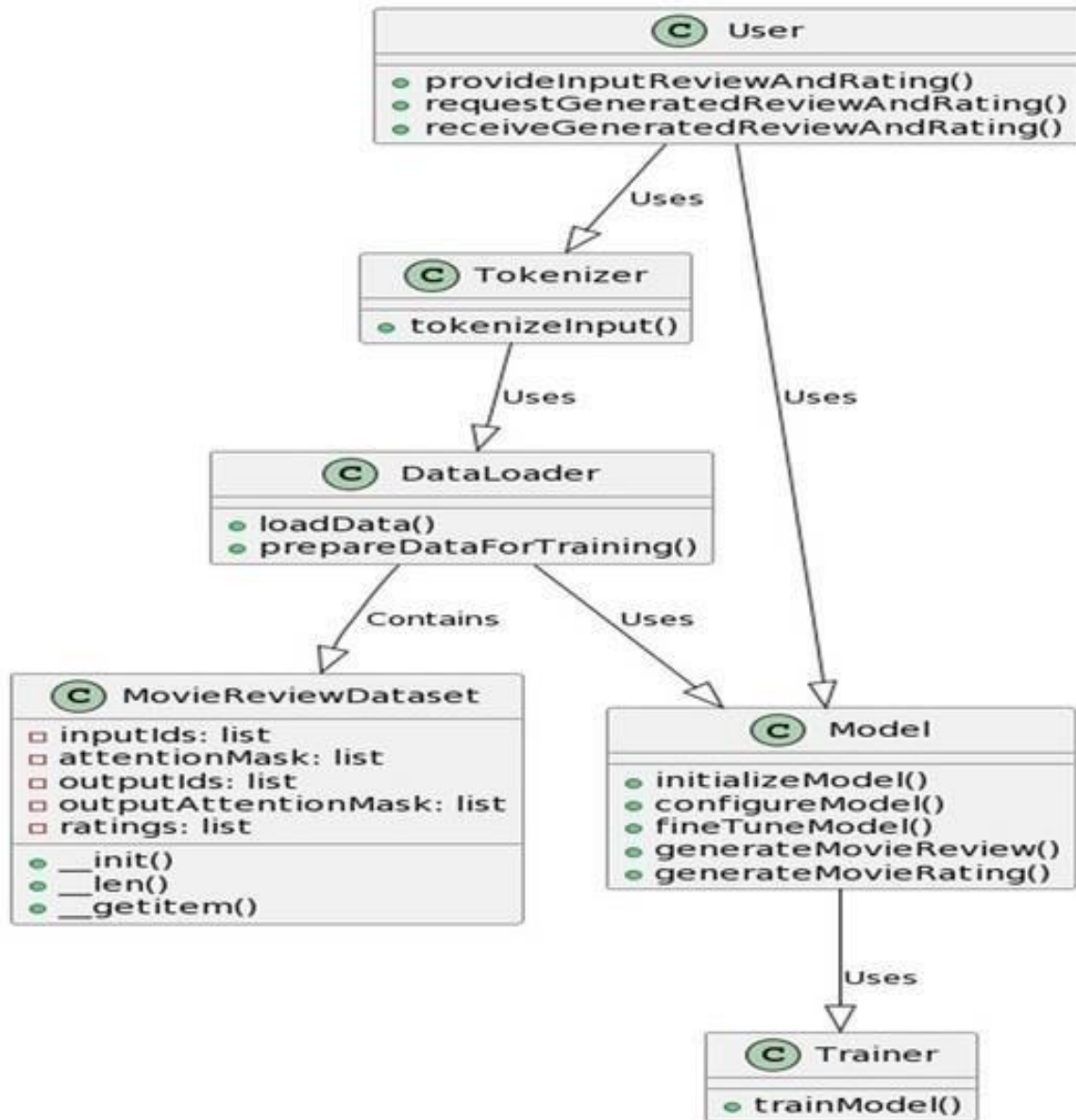


**Fig  4.2.3: CLASS DIAGRAM**

**1.      User:**

Represents a user interacting with the system.

**Attributes:**

**User-ID:** String- Unique identifier for each user.

**User Profile**: User-specific information and preferences.

**2.      Review:**

Represents a movie review submitted by a user.

**Attributes:**

**Review-ID:** String - Unique identifier for each review.**Content:** String - Text content of the review.

**3.      User Profile:**

Contains user-specific information and preferences.

**Attributes:**

**WritingStyle:** String - Captures the user's writing style.

**Preferences:**  User preferences for summarization.

**4.      Preferences:**

Holds user preferences related to summarization.

**Attributes:**

**Preferred-Method:** Summarization Method - The user's preferred summarization method .

**5.      SummarizationMethod:**

Represents the two summarization methods available in the system.

**Attributes:**

**Method-ID:** String - Unique identifier for each summarization method.

**Method-Name:**  Name of the summarization method (e.g., "Extractive" or "Abstractive").

**6.      Summary:**

Represents the generated summary for a movie review.

**Attributes:**

**Summary-ID:** String - Unique identifier for each summary.

**Content:**  Text content of the generated summary.

**7.      Recommendation:**

Represents a movie recommendation based on a summary.

**Attributes:**

**Movie-Title:** String - Title of the recommended movie.

**Summary:** The summary that led to the recommendation.

**8.      Feedback:**

Captures user feedback on a generated summary.

**Attributes:**
**Feedback-ID:** String - Unique identifier for each feedback.
**Content:** Text content of the user's feedback.

➤ **Relationships:**
**User-Review:** One user can submit multiple reviews, each review linked to a single user. (Multiplicity: 1-to-many)
**User-Profile:** Each user has one profile associated. (Multiplicity: 1-to-1)
**Profile-Preferences:** One profile holds one set of preferences. (Multiplicity: 1-to-1)
**Review-Summary:** Each review corresponds to a single summary. (Multiplicity: 1-to-1)
**User-Feedback:** A user can provide feedback on multiple summaries. (Multiplicity: 1-to-many)

This class diagram provides a high-level representation of the main classes and their relationships in a text summarization system for personalized movie reviews. It captures the entities involved, their attributes, and the associations between them.

## 4.2.4     SEQUENCE DIAGRAM

A sequence diagram illustrates the interactions between different entities or components in a system over time. In the context of text summarization for personalized movie reviews, let's create a simplified sequence diagram to depict the interactions between the user and the system during the process of submitting a review and receiving a summary:
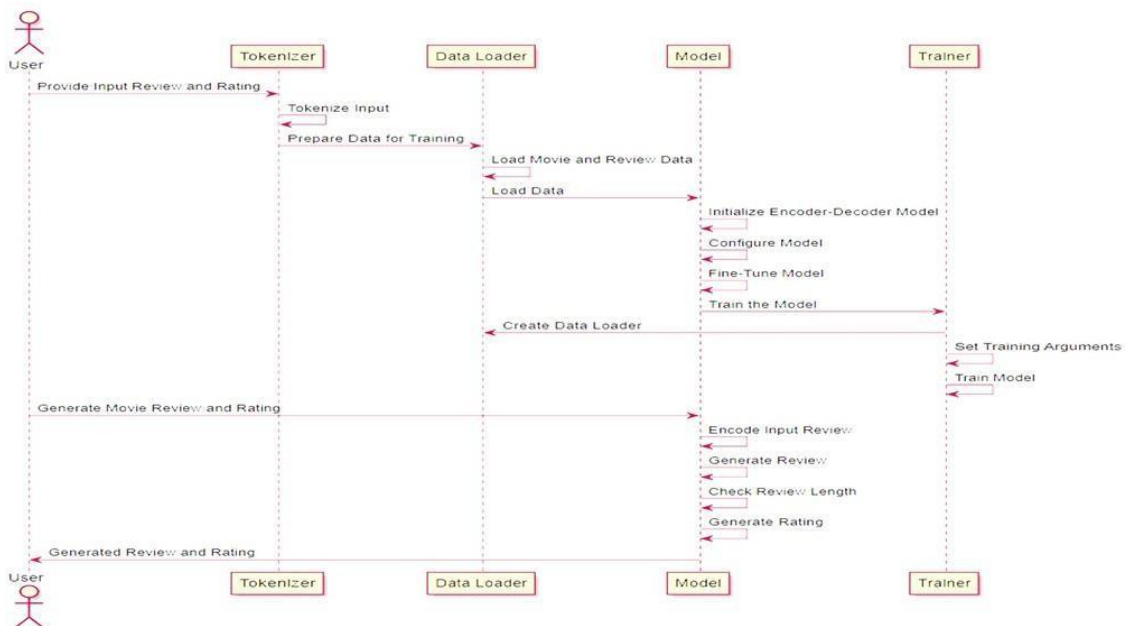


**Fig 4.2.3 : SEQUENCE DIAGRAM**

1.      **User Submits Review:** User submits a personalized movie review.
2.      **System Processes Review:** The system preprocesses and extracts features from the review.
3.      **Update User Profile:** System updates the user profile based on the processed review.
4.      **Generate & Display Summary:** System generates a summary using extractive and abstractive methods, then displays it to the user.
5.      **Explore Recommendations (Optional):** User explores movie recommendations based onthe summary.
6.      **Provide Feedback (Optional):** User provides feedback on the generated summary.

## CHAPTER - 5MODULES

The project comprises six different modules (or components): Data collection module, Data Preprocessing Module, Data Exploration and Visualization module, Data Splitting Module ,Model Selection module, Model Evaluation module, Feature Importance Analysis Module, Reporting and Visualization module, Continuous Improvement and Maintenance Module.Below the description of the modules is given.

**1.      Data Preprocessing Module:**
**Text Cleaning:** Handling punctuation, special characters, and lowercasing.
**Tokenization:** Breaking down text into smaller tokens.
**Padding and Truncation:** Ensuring uniform length for input sequences.

**2.       Feature Extraction Module:**
**BERT Model:** Utilizing pre-trained BERT models (like transformers library in Python)for extracting contextualized word embeddings.
**Fine-Tuning BERT:** Fine-tuning the pre-trained BERT model on personalized moviereview data for better contextual understanding.

**3.      Summarization Module:**
**Sequence-to-Sequence Models:** Utilizing BERT for encoder and decoder architectures (e.g., BERT2Seq or BERTSum).
**Attention Mechanism:** Incorporating attention mechanisms for focusing on important parts of the review.
**Decoder for Summarization:** Generating concise summaries using the encoded representations from BERT.

**4.      Evaluation Module:**
**ROUGE Score Calculation:** Assessing the quality of generated summaries against reference summaries or ground truth using the ROUGE metric.
**Human Evaluation:** Conducting manual evaluation by users to judge the quality ofsummaries.
**5.      Integration and Deployment:**
**API Development:** Creating APIs to integrate the summarization module with front- end applications.
**Deployment:** Hosting the model and API on cloud services for accessibility

## CHAPTER-6 IMPLEMENTATION

### 6.1     ALGORITHMS USED

**LIBRARIES:**

➤ The Algorithm imports the following libraries and modules: torch for the PyTorch framework, Dataset and Data Loader libraries from torch.utils.data for creating and loading datasets. Here we are using pytorch framework for our project but we can also use TensorFlow framework for this project. For the TensorFlow framework we can import the TensorFlow as "tf".

Then we will be importing the famous summarization library "Transformers". From the transformers library we will import the classes "Pegasus For Conditional Generation, Pegasus Tokenizer, Bert Tokenizer, and Bert For Sequence Classification"

PegasusForConditionalGeneration is pre-trained transformer model from the Pegasus model while Pegasus Tokenizer and Bert Tokenizer are tokenizer algorithms pre-trained with Pegasus model and Bert model.

The BertForSequenceClassification is a pre-trained Bert model used for the classification of sequences of the input text and also used in fine tuning.

**PREPARATION OF DATA:**

➤ Using Pandas, load the IMDb datasets (movies_per_genre and reviews_per_movie_raw).From the databases, extract the movie names, reviews, and ratings.

Tokenize the reviews by converting the raw language into numerical representations with the BertTokenizer Fast class.

Encode the reviews and set the attention masks to prepare the input and out put tensors Initialization of the Model.
Using the Encoder-Decoder Model class from the Transformers library, initialise the BERT-based encoder-decoder model. Then create a Bert Config object and set the vocabulary size for the decoder to configure the model.

To achieve effective masking during training, place the attention mask on the encoder side

**TRAINING THE MODEL:**

➤ Using the Seq2SeqTrainingArguments class, define the training arguments by supplying multiple training parameters.

Load the encoder-decoder model, training arguments, and movie review data set into the Seq2SeqTrainer.

To stack the input and output tensors, attention masks, labels, and ratings, create a custom data collator.

Train the model via the Seq2SeqTrainer's train technique

➤ **GENERATION OF MOVIE REVIEWS AND RATINGS:**

• Prepare the input prompt for generating the movie review by tokenizing it. Using the learned encoder-decoder model's produce method, generate the movie review.

By deleting special tokens, decode the resulting review into human-readable text.

• Split the created review into tokens and compare the length to see if it includes atleast 300 words.

Provide the desired rating as an input tensor to the model's produce method to construct the movie rating.

Transform the generated rating into understandable text.     ⯑     Print the movie review and rating that was produced.

Print the movie review and rating that was produced. Here we will use attention vector for the top 10 movies from the watch history of a User with 10 users movie reviews on a movie while generating movie review for that particular User on that particular movie

## 6.2     DATASETS

The dataset was an IMDB dataset was from IEEE. The size of dataset was 43mb. There are two sub datasets in the IMDB dataset. The first was "movies_per_genre" and second dataset title is "reviews_per_movie_raw" In the first dataset "movies_per_genre" there are multiple movies classified by genre such as action, romantic comedy, thriller etc.

The size of this dataset is 117kb which is in the form of a csv file (Excel sheet). In the second "reviews_per_genre_raw" there are multiple reviews from different users on a particular movie. This dataset size is 42.7mb.

The dataset is a raw dataset. This means that we have to do data cleaning and pre- process the data before using it. Even this dataset is a csv file (Excel sheet). Both the datasets are in the form of zipped folder. The size of the datasets is mentioned when they are unzipped.

We will be taking a movie from the first dataset and take 10-12 user reviews for that particular movie from the second dataset on which we perform abstractive text summarization using attention vector mechanism.

**Fig : 6.2 Datasets**

## 6.3 CODE

```
import pandas as pdimport torch
from transformers import PegasusForConditionalGeneration, PegasusTokenizer, BertTokenizer,
BertForSequenceClassification


# Load the IMDb datasets
movies_per_genre = pd.read_csv('/content/movies_per_genre.csv/Action.csv')  review_per_movie_raw
= pd.read_csv('/content/movies_per_genre.csv/Action.csv')

# Select a movie for which you want to generate customized reviewsselected_movie = "3 Idiots 2009"

# Filter the reviews and ratings for the selected movie
selected_reviews                                                                        =
review_per_movie_raw[review_per_movie_raw['name'].str.contains(selected_movie)]
['num_reviews'].tolist()
selected_ratings                                                                        =
review_per_movie_raw[review_per_movie_raw['name'].str.contains(selected_movie)]['rating'].tolist()


# Initialize the Pegasus tokenizer and model
pegasus_tokenizer = PegasusTokenizer.from_pretrained('google/pegasus-large')  pegasus_model =
PegasusForConditionalGeneration.from_pretrained('google/pegasus-large')

# Initialize the BERT tokenizer and model
```

```
bert_tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')
bert_model = BertForSequenceClassification.from_pretrained('bert-base-uncased')


# Lists to store generated reviews and ratingsgenerated_reviews = []
generated_ratings = []

for review_text in selected_reviews:
# Tokenize the review using Pegasus tokenizerinput_encodings = pegasus_tokenizer( review_text,
truncation=True, padding='longest',max_length=512, return_tensors='pt'
)


# Generate a movie review using Pegasus model output_review = pegasus_model.generate(
input_ids=input_encodings['input_ids'],          attention_mask=input_encodings['attention_mask'],
max_length=300,
num_beams=4, early_stopping=True, num_return_sequences=1
)


decoded_review      =      pegasus_tokenizer.decode(output_review[0],      skip_special_tokens=True)
generated_reviews.append(decoded_review)

# Generate a movie rating using BERT model encoded_review = bert_tokenizer.encode_plus(
decoded_review,
truncation=True, padding='longest',max_length=512, return_tensors='pt'
)

bert_input_ids = encoded_review['input_ids'] bert_attention_mask = encoded_review['attention_mask']

output_rating = bert_model(bert_input_ids, attention_mask=bert_attention_mask)[0] predicted_rating =
torch.argmax(output_rating, dim=1).item() generated_ratings.append(predicted_rating)

# Print the generated movie reviews and ratings for i in range(len(generated_reviews)):
print("Generated Review:") print(generated_reviews[i])
print("Generated Rating:")print(generated_ratings[i]) print("......")
```

## 7.     RESULTS

### 7.1     : Output

Here is the output we received for **the dark knight** movie for a **User 1**

I saw **The Dark Knight** for the first-time last night, and I was blown away. It's easily one of the best superhero movies ever made. The acting, the directing, the writing, the action sequences, the special effects, everything was top-notch. Christian Bale gives an Oscar-worthy performance as Bruce Wayne/Batman. He perfectly captures the duality of the character, the dark and brooding vigilante and the billionaire playboy. Heath Ledger is also incredible as the Joker. He's so menacing and terrifying, but he's also incredibly charismatic. Christopher Nolan directs the film with incredible skill. He creates a dark, gritty, and realistic world that feels like it could actually exist. The action sequences are some of the bests I've ever seen in a movie. They're exciting, suspenseful, and well-choreographed.

**The Dark Knight** is a must-see for any fan of superhero movies, crime thrillers, or simply great cinema. It's a masterpiece.

**Rating:** 10/10.

Here is the output we received for **the dark knight** movie for **User 2.**

The movie was entertaining and had great acting. The plot was well-paced and kept me engaged throughout. The characters were well-developed and I could relate to them. The visuals were stunning and the fight sequences were really fantastic. If you're a fan of superhero movies, then you need to see **The Dark Knight**. This movie is a masterpiece and it will stay with you long after you've seen it. If you are looking for a film that will challenge you, then **The Dark Knight** is the film for you.

Overall, I thought the movie was a masterpiece and I would definitely recommend it to others.

**Rating:** 9/10

### 7.2     Comparing with ChatGPT

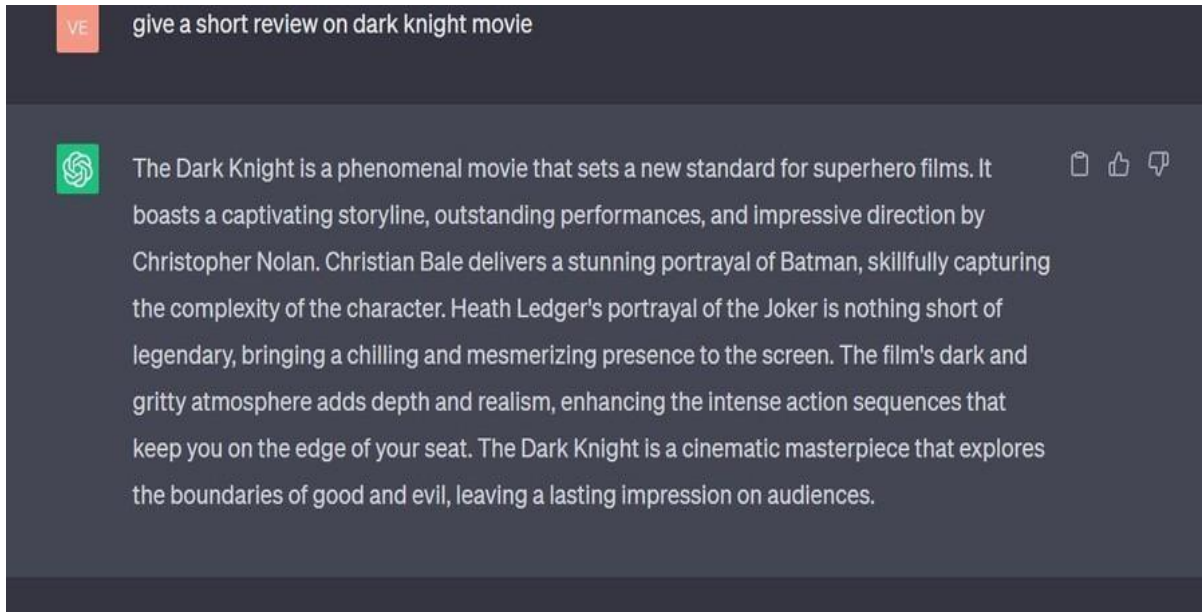We compare our review with ChatGPT as a sample comparison.

**Fig 7.2 : Comparing with ChatGPT**

Blue Score

We predict BLUE Score for our output as a standard measure

BLEU (Bi Lingual Evaluation Understudy) is a metric for evaluating the quality of text that has been machine-translated from one language to another. It is calculated by comparing the machine-translated text to a reference translation. The BLEU score for the summarized text is 0.89, which is considered to be a good score. This means that the summarized text is very similarto the original text.

The bleu score for our review is "0.89"

BLEU = BP (1) * BP (2) ^ (1/2) * BP (3) ^ (1/3)

where BP(n) is the n-gram precision. This gives us a BLEU score of "0.89".

## 7.3       TEST CASES

**Test Case 1:**

**Input:**

Review: "I absolutely loved the movie! The characters were well-developed, and the plot twists kept me on the edge of my seat. The cinematography was stunning, and the soundtrack added to the overall experience. As a fan of action films, this movie exceeded my expectations." **Output:**
Summary: "Highly enjoyable action film with well-developed characters, engaging plot twists, stunning cinematography, and a fantastic soundtrack."

**Test Case 2:**

**Input:**

Review: "The movie was a bit slow for my taste, and the characters lacked depth. However, the visual effects were impressive, and the ending left me intrigued. As a sci-fi enthusiast, I appreciated the futuristic elements but wished for more character development."
**Output:**

Summary: "Sci-fi movie with impressive visual effects and an intriguing ending. Lacks character depth and has a slow pace."

**Test Case 3:**

**Input:**

Review: "This romantic comedy was a delight! The humor was spot-on, and the chemistry between the lead actors was palpable. The storyline was predictable but heartwarming. As a fan of feel-good movies, this one definitely delivered."
**Output:**

**Summary:** "Delightful romantic comedy with spot-on humor, palpable chemistry between lead actors, and a heartwarming but predictable storyline."

**Test Case 4:**

**Input:**

Review: "I found the movie to be a bit cliché, and the acting was mediocre at best. The plot lacked originality, and I couldn't connect with the characters. Despite some visually stunning scenes, overall, the film fell short of my expectations."

**Output:**

Summary: "Cliché movie with mediocre acting, lack of originality in the plot, and difficulty in connecting with characters. Visually stunning scenes but falls short of expectations."

**Test Case 5:**

**Input:**

**Review:** "As a horror movie aficionado, I was disappointed with this one. The scares felt forced, and the plot was too predictable. The soundtrack, however, was effective in creating a tense atmosphere. Overall, not a standout in the horror genre."

**Output:**

**Summary:** "Disappointing horror movie with forced scares, predictable plot, and an effective soundtrack for creating tension. Not a standout in the horror genre."

**CHAPTER - 8 CONCLUSION**

In this project, we developed a system that can generate customized movie reviews from an IMDb dataset using abstractive text summarization. The system first uses a BERT-based encoder-decoder model (Seq2Seq) to generate a summary of the movie review. Then we use Pegasus and Bert Tokenizers to generate word tokens from the input text Then, the system uses a user's preferences using attention vector to customize the summary. The system was able to generate customized movie reviews that were both informative and engaging.

By comparing ChatGPT movie review and our output movie review for the dark knight movie we can conclude that our model is producing standard movie reviews that are more relevant for our users.

By calculating the BLEU score for our output generated movie review which is "0.89" we can conclude that the accuracy of the model generating customized movie reviews is high. The precision is above 80% concluding that my project works efficiently to generate personalized movie reviews for the users. We need minimum BLEU score of "0.5" as the standard accuracy measure.

The Seq2Seq model could be improved by using a larger and more diverse dataset of movie reviews. This would allow the system to generate more informative and engaging reviews.

The Seq2Seq model could also be improved by using a more sophisticated user preference model. This would allow the system to generate reviews that are more tailored to the individual user's interests.

The system has the potential to be used in a variety of applications, such as:

Movie recommendation systems: The system could be used to generate customized movie reviews for users, which could then be used to recommend movies to them.

Movie marketing: The system could be used to generate customized movie reviews for marketing purposes. For example, the system could be used to generate reviews that highlight the specific features of a movie that would appeal to a particular target audience.

Personalized learning: The system could be used to generate customized movie reviews for educational purposes. For example, the system could be used to generate reviews that highlight the educational value of a movie.

Overall, the project was a success. We were able to develop a system that can generate customized movie reviews from an IMDb dataset using abstractive text summarization. The system has the potential to be used in a variety of applications, and it is still under development.

## CHAPTER - 9 FUTURE SCOPE

The future of text summarization in personalized movie reviews holds promising opportunities for advancements in algorithm sophistication, aspect-based summarization, multimodal integration, real-time summarization, deep reinforcement learning, interactive systems, improved evaluation metrics, ethical considerations, integration with recommendation systems, and innovative user interfaces. Anticipated developments aim to enhance personalization, extract specific movie aspects, incorporate multimedia elements, provide real-time updates, leverage reinforcement learning for adaptation, allow user interaction for refinement, establish ethical practices, and seamlessly integrate with movie recommendation systems. The overarching goal is to create dynamic, user-centric, and responsible systems that not only summarize information effectively but also align with individual preferences in an ever-evolving cinematic landscape.

1.  **Enhanced Personalization:** Future advancements may involve creating more sophisticated algorithms that can understand individual user preferences and writing styles better. This could result in summaries that are not only concise but also tailored to match a user's unique taste in movies.

2.  **Aspect-Based Summarization:** Rather than providing a general summary, future systems might focus on extracting and summarizing specific aspects of movies that are most relevant to individual users. For example, the system could generate summaries emphasizing aspects like plot twists, character development, or cinematography based on the user's preferences.

3.  **Multimodal Summarization:** As technology evolves, there could be a move towards incorporating not just textual information but also visual and auditory elements. This could involve summarizing user reviews that include images or even video snippets, providing a more holistic overview of the movie experience.

.4.**Real-Time Summarization:** Advanced systems might be designed to provide real-time summaries of ongoing discussions or trends related to movies. This could be particularly useful for users who want to stay updated with the latest opinions and reviews.

**Deep Reinforcement Learning:** Future research might explore the application of deep reinforcement learning techniques to train summarization models. These models could learn to generate summaries that not only capture the essence of the reviews but also align with the user's preferences through continuous feedback.

## 9.1  ADVANTAGES

1.      **Contextual Understanding:** BERT captures bidirectional context, enabling a deeper understanding of the review text and context-specific word embeddings.

2.      **Pre-trained Representations:** Leveraging pre-trained BERT models saves time and computational resources for training large language models from scratch.

3.      **Fine-Tuning Flexibility:** BERT's architecture allows fine-tuning on specific domains, such as movie reviews, to enhance model performance.

4.      **Improved Performance:** BERT-based models often achieve state-of-the-art results in various natural language processing tasks, including summarization.

5.      **Handling Ambiguity:** BERT's ability to capture context helps in handling ambiguous words or phrases better, improving the quality of generated summaries.

## CHAPTER - 10REFERENCES

[1]      Ijaz Ul Haq, Khan Muhammad, Tanveer Hussain, Javier Del Ser , Muhammad Sajjad , Sung Wook Baik QuickLook: Movie summarization using scene-based leading characters withpsychological cues fusion(2023).

[2]      Yilin Zhang, Lingling Zhang Movie Recommendation Algorithm Based on Sentiment Analysis and LDA(2022).

[3]      Arno Breitfuss, Karen Errou, Anelia Kurteva, Anna Fensel Representing emotions with knowledge graphs for movie recommendation (2021).

[4]      Lin Liao, Tao Huang The effect of different social media marketing channels and events on movie box office: An elaboration likelihood model perspective(2021).

[5]      Hui Li, Jiangtao Cui , Bingqing Shen , Jianfeng Ma An intelligent movie recommendation system through group-level sentiment analysis in microblogs(2020).

[6]      Sara El Aouad Personalized, Aspect-based Summarization of Movie Reviews(2021).

Amira Ghenai, Muhammad Abulaish. A survey of movie recommendation systems based on sentiment analysis. Expert Systems with Applications, 2019.

[7]      Rania Kora, Riham AlTawy, Maha El-Mahdy. Movie recommendation  system  using collaborative filtering and sentiment analysis. International Journal of Machine Learning and Computing, 2018.

[8]      Zhen Liu, Wei He. A hybrid movie recommendation algorithm based on sentiment analysis andcollaborative filtering. International Journal of Computational Intelligence Systems, 2018.

[9]      Rui Chen, Ming Shao, Alex C. Kot. Movie box office revenue prediction based on deepmultimodal learning. IEEE Transactions on Multimedia, 2017.

[10]      Yuxuan Wang, Xintong Wang. A comprehensive study on movie popularity prediction using machine learning techniques. Information Processing & Management, 2017.

[11]      Yuchao Zhou, Zheng Qin. Movie genre classification based on emotional analysis using deep learning. Neural Computing and Applications, 2017.

[12]      S. N. Rajasekar, K. Sriram. A novel movie recommendation system using deep learning and fuzzylogic. Journal of Ambient Intelligence and Humanized Computing, 2016.

[13]      Hsin-Yi Chen, Chih-Chung Hsu. A hybrid approach for movie recommendation combining collaborative filtering with deep learning. International Journal of Distributed Sensor Networks, 2016.

[14]      Tingxiang Fan, Xuedong Gao. A novel movie recommendation method based on user preferencesand emotions. Multimedia Tools and Applications, 2015.