



Course Project Filmception

An AI-Powered Multilingual Movie Summary Translator and Genre Classifier

AI2002

Artificial Intelligence

Submitted by: Umar Murtaza & Tehreem Zafar

Roll number: i227431 & i221630

Date: May 6, 2025



National University of Computer and Emerging Sciences

Islamabad Campus

Table of Contents

Introduction.....	3
Development Platform.....	3
Google Drive Integration and Folder Structure:	3
Dataset Description.....	4
Data Preprocessing and Cleaning	5
Purpose:	5
Details:	5
Text Translation.....	7
Purpose:	7
Details:	7
Audio Conversion (Text-to-Speech).....	10
Purpose:	10
Genre Prediction Model.....	12
Feature Extraction.....	12
Purpose:	12
Details:	12
Multi-Label Binarization.....	13
Purpose:	13
Details:	13
Train/Test Split.....	14
Purpose:	14
Details:	14
Model for Prediction on New Summaries.....	16
GUI	17
Conclusion.....	26



National University of Computer and Emerging Sciences Islamabad Campus

Introduction

This project aims to develop an AI-powered system that processes movie summaries, predicts genres, and converts text into audio in multiple languages. The goal is to enhance movie metadata usability by enabling users to hear summaries in their native languages and understand the classification of the genre.

Development Platform

The entire project was developed using:

- **Platform:** Google Colab (for rapid cloud-based prototyping with GPU support)
- **Language:** Python 3.10+
- **Libraries Used:** pandas, numpy, nltk, sklearn, googletrans, gTTS, matplotlib, seaborn, scikit-multilearn, etc.

Google Drive Integration and Folder Structure:

- Google drive was mounted in Colab as cloud storage using:

```
✓ [1] from google.colab import drive  
36s drive.mount('/content/drive')
```

Mounted at /content/drive

This ensured persistent access to datasets, checkpoints, models, and audio outputs.

- The folder hierarchy is as shown:

```
/content/drive/MyDrive/Filmception/  
├─ data/  
│  ├── plot_summaries.txt  
│  ├── movie_metadata.tsv  
│  ├── cleaned_movies.csv  
│  └── translated_checkpoints.csv  
├─ audio/  
│  ├── arabic/  
│  ├── urdu/  
│  └── korean/  
├─ models/  
│  ├── tfidf_vectorizer.pkl  
│  └── genre_model.pkl  
└─ notebooks/  
    ├── 1_preprocessing_cleaning.ipynb  
    ├── 2_translation_tts.ipynb  
    └── 3_genre_prediction.ipynb
```

```
✓ [2] import os  
0s  
  
base_dir = '/content/drive/MyDrive/Filmception'  
os.makedirs(f'{base_dir}/data', exist_ok=True)  
os.makedirs(f'{base_dir}/audio', exist_ok=True)  
os.makedirs(f'{base_dir}/models', exist_ok=True)
```



National University of Computer and Emerging Sciences Islamabad Campus

Dataset Description

We used data from the **CMU Movie Summary Corpus**:

- **plot_summaries.txt**: 42,306 movie plot summaries.
- **movie.metadata.tsv**: Metadata including movie name, release year, and genres (in JSON format).

After merging with **movie_id**, we retained:

- Summary text (plot)
- Movie name
- Genre labels (multi-label)

The dataset is publicly available on Kaggle and includes over 40,000 entries. It can be accessed here:

<https://www.kaggle.com/datasets/msafi04/movies-genre-dataset-cmu-movie-summary>

The screenshot shows the Kaggle dataset page for 'Movies Genre Dataset - CMU Movie Summary' by user msafi04. The page is divided into a left sidebar with navigation links (Home, Competitions, Datasets, Models, Code, Discussions, Learn, More) and a main content area. The main content area includes a search bar, 'Sign In' and 'Register' buttons, and a 'Create' button. The dataset title 'Movies Genre Dataset - CMU Movie Summary' is prominently displayed, along with a 'Data Card' tab and 'Code (0)', 'Discussion (0)', and 'Suggestions (0)' links. Below the title, there is an 'About Dataset' section with a 'Usability' score of 2.94 and a 'License' of 'Unknown'. A 'Summary' section on the right indicates 6 files and 22 columns. The 'Data Card' section displays a grid of files: README.txt (2.58 kB), name.clusters.txt (65.44 kB), plot_summaries.txt (75.93 MB), tvtropes.clusters.txt (57.98 kB), character.metadata.tsv (41.48 MB), and movie.metadata.tsv (16.25 MB). A cookie notice is visible at the bottom of the page.



National University of Computer and Emerging Sciences Islamabad Campus

Data Preprocessing and Cleaning

Purpose:

To merge raw summaries and genre metadata, clean the text, and prepare it for downstream tasks.

Details:

- **Raw Inputs:**
 - plot_summaries.txt: Contains movie summaries.
 - movie.metadata.tsv: Contains movie metadata, including genres.
- **Process:**
 1. Summaries and genres are merged on movie_id.
 2. NLTK's stopwords and WordNetLemmatizer are used to:
 - Remove special characters
 - Convert to lowercase
 - Tokenize and lemmatize words
 - Remove stopwords
- **Genre Field Handling:**
 - Genres stored as a dictionary-like string ("{'Action': 1, 'Thriller': 1}") are converted into Python lists using ast.literal_eval.
- **Output:**
 - Cleaned file saved at:
/content/drive/MyDrive/Filmception/data/cleaned_movies.csv

Data Preprocessing and cleaning was done in the following cell of Colab notebook:

National University of Computer and Emerging Sciences

Islamabad Campus

Cleaning Data (plot_summaries.txt, movie.metadata.tsv)

```
[3] import pandas as pd
import re
import nltk
nltk.download('stopwords')
#nltk.download('punkt') # not needed
nltk.download('wordnet')
nltk.download('punkt_tab') # download resource
from nltk.corpus import stopwords
from nltk.stem import WordNetLemmatizer

# Load summaries
with open(f'{base_dir}/data/plot_summaries.txt', encoding='utf-8') as f:
    lines = f.readlines()
    summaries = dict(line.strip().split('\t', 1) for line in lines)

# Load metadata
meta = pd.read_csv(f'{base_dir}/data/movie.metadata.tsv', sep='\t', header=None)
meta = meta[[0, 8]] # movie_id and genres
meta.columns = ['movie_id', 'genres']

# Merge
df = pd.DataFrame.from_dict(summaries, orient='index', columns=['summary'])
df.index.name = 'movie_id'
df.reset_index(inplace=True)
# Convert 'movie_id' column in df to int64 before merging
df['movie_id'] = df['movie_id'].astype(int)
df = df.merge(meta, on='movie_id')
```

```
[3] # Clean summaries
stop_words = set(stopwords.words('english'))
lemmatizer = WordNetLemmatizer()

def clean_text(text):
    text = re.sub(r'^a-zA-Z\s', '', text)
    text = text.lower()
    tokens = nltk.word_tokenize(text)
    tokens = [lemmatizer.lemmatize(t) for t in tokens if t not in stop_words]
    return ' '.join(tokens)

df['clean_summary'] = df['summary'].apply(clean_text)

# Convert genres from stringified JSON
import ast
df['genres'] = df['genres'].apply(lambda g: list(ast.literal_eval(g).values()) if pd.notna(g) else [])

# Save cleaned data
df.to_csv(f'{base_dir}/data/cleaned_movies.csv', index=False)
```

[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Unzipping corpora/stopwords.zip.
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data] Downloading package punkt_tab to /root/nltk_data...
[nltk_data] Unzipping tokenizers/punkt_tab.zip.

Steps:

- Merged metadata and plot summaries using movie ID.
- Removed empty or missing summaries.
- Cleaned text: lowercasing, punctuation/digit removal, extra whitespace handling.



National University of Computer and Emerging Sciences Islamabad Campus

- Tokenization, stopword removal (using NLTK), and lemmatization.
- Parsed and flattened JSON-encoded genre fields.
- Multi-label genres encoded using MultiLabelBinarizer.
- Final output saved as: **cleaned_movies.csv**

Text Translation

Purpose:

To translate the cleaned summaries into Arabic, Urdu, and Korean using Google Translate API.

Details:

- **Tool:** [deep translator](#) (GoogleTranslator)
- **Chunking:**
Translations are split into ≤ 4500 -character chunks to stay under the API limit (5000 chars).
- **Rate Limiting:**
 - A 1-second delay per row is used to reduce translation errors.
 - Intermediate savings occur every 10 rows to prevent loss of progress.
- **Languages & Output Columns:**
 - Arabic \rightarrow summary_ar
 - Urdu \rightarrow summary_ur
 - Korean \rightarrow summary_ko
- **Final Output:**
 - Translated summaries saved at:
/content/drive/MyDrive/Filmception/data/translated_checkpoints.csv

Batching & error-handling were built in to handle rate limits.

Step 2: Translation & Text-to-Speech

Installing required libraries:

```
✓ [4] !pip install googletrans==4.0.0-rc1 gTTS
```



Steps:

1. The first step was to import and load the cleaned data (previous step) stored in `/content/drive/Flimception/data/cleaned_movies.csv`:

Import and load cleaned data

```
[ ] import pandas as pd
    from googletrans import Translator
    from gtts import gTTS
    import os

    # Load cleaned summaries
    df = pd.read_csv('/content/drive/MyDrive/Filmception/data/cleaned_movies.csv')
    df = df.head(60)
```

2. Along with Google Translate, we also needed Deep Translator so we installed that by `“!pip install deep-translator”`
3. The code block below implements the actual translation logic:

National University of Computer and Emerging Sciences

Islamabad Campus

```

✓ 4m [8] from deep_translator import GoogleTranslator
import pandas as pd
import time

def translate_text(text, lang, chunk_size=4500):
    if not isinstance(text, str) or text.strip() == "":
        return ""
    try:
        # Split text into manageable chunks under the 5000 char limit
        chunks = [text[i:i+chunk_size] for i in range(0, len(text), chunk_size)]
        translated_chunks = [
            GoogleTranslator(source='auto', target=lang).translate(chunk)
            for chunk in chunks
        ]
        return " ".join(translated_chunks)
    except Exception as e:
        print(f"Translation error for lang '{lang}': {e}")
        return ""

df = pd.read_csv('/content/drive/MyDrive/Filmception/data/cleaned_movies.csv')
df = df.dropna(subset=['clean_summary']) # Optional: remove NaN summaries
df = df.reset_index(drop=True)

# Optional: reduce for testing
df = df.head(60)

translated_ar = []
translated_ur = []
translated_ko = []

```

```

✓ 4m [8] for idx, row in df.iterrows():
    summary = row['clean_summary']

    ar = translate_text(summary, 'ar')
    ur = translate_text(summary, 'ur')
    ko = translate_text(summary, 'ko')

    translated_ar.append(ar)
    translated_ur.append(ur)
    translated_ko.append(ko)

    # Show progress and add delay to avoid rate limiting
    print(f"Translated index {idx}")
    time.sleep(1) # 1 sec delay per row (adjustable)

    # Save interim progress every 10 rows
    if idx % 10 == 0:
        df['summary_ar'] = pd.Series(translated_ar)
        df['summary_ur'] = pd.Series(translated_ur)
        df['summary_ko'] = pd.Series(translated_ko)
        df.to_csv('/content/drive/MyDrive/Filmception/data/translated_checkpoints.csv', index=False)

# Final save
df['summary_ar'] = translated_ar
df['summary_ur'] = translated_ur
df['summary_ko'] = translated_ko
df.to_csv('/content/drive/MyDrive/Filmception/data/translated_checkpoints.csv', index=False)

print("✅ Translations complete and saved!")

```

The output file: **translated_checkpoints.csv** has columns movie_id, summary, genres, clean_summary, summary_ar, summary_ur, summary_ko:

National University of Computer and Emerging Sciences

Islamabad Campus

docs.google.com/spreadsheets/d/17qUfTP82cQzjdtQN2JnSwEttEnYr5ZT29MvX5i0HQ/edit?gid=1805910147#gid=1805910147

translated_checkpoints

movie_id	summary	genres	clean_summary	summary_ar	summary_ur	summary_ko
23890098	Shlykov, a hard-	'Drama', 'World	shlykov hardwork	شليكوف محنتي ٹيڪسي سائق الذي	Shlykov 열심	
31186339	The nation of Pa	'Action/Adventu	nation panem cobi	نیشن پانیم پر مشتمل دو اڈا ہے تاکہ	Nation Panem	은 부유 한 국회 의사당이 부유 한 12 명의 가난한 지구 처벌 과거 반란 지구 과거
20663735	Poovalli Inducho	'Musical', 'Actio	poovalli inducho	پووالی انڈوچوٹ کو دیکھ کر حکم علی	Poovalli Induchoodan	은 6 년의 감옥 생명을 선고 한 6 년의 감옥 생명을 선고했다. INDUCHOO
2231378	The Lemon Drop	'Screwball come	lemon drop kid	لیمن ڈراپ نیو یارک کی نیویورک		
595909	Seventh-day Ad	'Crime Fiction',	'seventhday adve	سائون ایڈونٹس جرج الیوم السابع من		
5272176	The president is	'Thriller', 'Actio	president way	صدر وے تقریر ٹریڈنگ طریقہ	president way	give speech traveling man show camera reporter try ask member secret service
1952976	{[plot]} The film	{[Thriller], 'Dra	ma plot film open	یوٹا فلم کھلی نیوجارن یوٹا فلم مفتوح		
24225279	The story begins	'Drama', 'Teen'	story begin han	کہانی کی شروعات ہذا القصة لیا	이야기 시작 Hannah Young Jewish Teen Weally High School Small Neighbork Brooklyn Brook	
2462689	Infruriated at be	in[Romantic come	infruriated told wr	ایک آخری کالم لکھیے آخری کتاب		
20532852	A line of people	'[Short Film]',	'Far line people	لائن لوگ ڈروول وٹو خطوط اللقاء سال		
15401493	Lola attempts to	[Comedy]	lola attempt gain	لولا کوشش کین فادر تاحولہ لولا	LOLA 시도 게이 아버지 신탁 기금 고용 히스패닉 남편 BO 1 년 마리를 수락합니다 텍사스 홀	
18188932	Milan and Goran	'Crime Fiction',	'milan goran two	میٹان گوران دو مجرم میلانو		
2940516	Bumbling pirate	'[Parody]', 'Come	bumbling pirate	بومبلنگ سمندری ڈاکو Cr		
1335380	The film is base	'War film', 'Epic'	film based event	فلم پر مبنی واقعہ بود		
1480747	{[plot]} Following	{[Cult], 'Coming	c plot following	ایچانگ موت کے بعد فامؤايرة بعد الموت	plot following sudden death kid father pop local church donated scholarship fund kid go college	
24448645	Despite Lucy's re	[Horror]	despite lucy res	لوسی ریڈرویش کے با علی الرغم من	루시 메악 동의를	

Audio Conversion (Text-to-Speech)

Purpose:

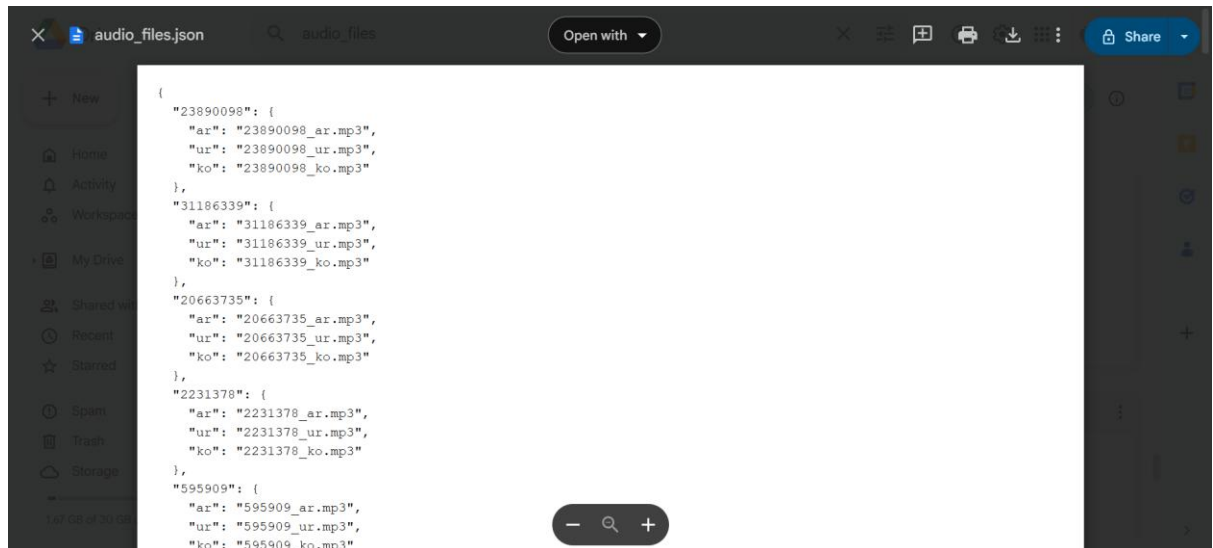
Convert translated summaries into text-to-speech (TTS) audio using gTTS.

Details:

- **Tool:** gTTS (Google Text-to-Speech)

```
Audio Conversion
[ ] pip install gTTS
Requirement already satisfied: gTTS in /usr/local/lib/python3.11/dist-packages (2.5.4)
Requirement already satisfied: requests<3,>=2.27 in /usr/local/lib/python3.11/dist-packages (from gTTS) (2.32.3)
Requirement already satisfied: click<8.2,>=7.1 in /usr/local/lib/python3.11/dist-packages (from gTTS) (8.1.8)
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.11/dist-packages (from requests<3,>=2.27->gTTS) (3.4.1)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.11/dist-packages (from requests<3,>=2.27->gTTS) (2.10)
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.11/dist-packages (from requests<3,>=2.27->gTTS) (2.4.0)
Requirement already satisfied: certifi<2017.4.17 in /usr/local/lib/python3.11/dist-packages (from requests<3,>=2.27->gTTS) (2025.4.26)
```

- **Languages Handled:** Arabic (ar), Urdu (ur), Korean (ko)
- **Retry Logic:**
If TTS fails, retry 3 times with 15-second delays.
- **File Naming Convention:** {movie_id}_{lang}.mp3 e.g., 1431_ar.mp3
- **Audio Metadata Index:**
 - audio_files.json stores the file paths for each movie and language.
 - Example JSON structure:



Audio Output Directory: /content/drive/MyDrive/Filmception/audio/

Code Implementation:

```
[ ] import os
import time
import json
from gtts import gTTS
from gtts.tts import gTTSError

# Load translated data
df = pd.read_csv('/content/drive/MyDrive/Filmception/data/translated_checkpoints.csv')

# Languages to convert and their respective columns
lang_map = {
    'ar': 'summary_ar',
    'ur': 'summary_ur',
    'ko': 'summary_ko'
}

# Directory to save audio
audio_dir = '/content/drive/MyDrive/Filmception/audio/'
os.makedirs(audio_dir, exist_ok=True)

# Dictionary to store audio file paths
audio_files = {}

# TTS conversion loop
for idx, row in df.iterrows():
    movie_id = str(row['movie_id'])
    audio_files[movie_id] = {}

    for lang, col in lang_map.items():
        text = row.get(col, '')
        filename = f"{movie_id}_{lang}.mp3"
```

```
[ ]      filepath = os.path.join(audio_dir, filename)

# Skip if already exists
if os.path.exists(filepath):
    print(f"Already exists: {filename}")
    audio_files[movie_id][lang] = filename
    continue

# Retry mechanism
success = False
for attempt in range(3):
    try:
        tts = gTTS(text=text, lang=lang)
        tts.save(filepath)
        print(f"✅ Saved: {filename}")
        audio_files[movie_id][lang] = filename
        success = True
        break
    except gTTSError as e:
        print(f"❌ Error for movie_id {movie_id} lang {lang} (Attempt {attempt+1}/3): {e}")
        time.sleep(15) # Wait before retry

if not success:
    print(f"❌ Failed after 3 attempts: {filename}")

time.sleep(10) # Delay between conversions

# Save JSON
json_path = os.path.join(audio_dir, 'audio_files.json')
with open(json_path, 'w') as f:
    json.dump(audio_files, f, indent=2)
```

Genre Prediction Model

Feature Extraction

Purpose:

Convert each cleaned summary into a numeric vector (embedding) for machine learning.

Details:

- **Tool:** [sentence-transformers](#)
- **Model Used:** all-MiniLM-L6-v2 (small and efficient BERT variant)
- **Hardware:** Utilizes GPU if available (cuda)
- **Output:** List of 384-dimensional embeddings per movie, used for genre classification.

Code Implementation:



National University of Computer and Emerging Sciences Islamabad Campus

2. Feature Extraction

```
✓ 57m [ ] from sentence_transformers import SentenceTransformer
import torch

# Load model
device = 'cuda' if torch.cuda.is_available() else 'cpu'
model = SentenceTransformer('all-MiniLM-L6-v2', device=device)

# Encode summaries
embeddings = model.encode(df['clean_summary'].tolist(), show_progress_bar=True, device=device)
```

/usr/local/lib/python3.11/dist-packages/huggingface_hub/utils/_auth.py:94: UserWarning:
The secret 'HF_TOKEN' does not exist in your Colab secrets.
To authenticate with the Hugging Face Hub, create a token in your settings tab (<https://huggingface.co/settings/tokens>), set it as secret in your Google Colab and restart your session.
You will be able to reuse this secret in all of your notebooks.
Please note that authentication is recommended but still optional to access public models or datasets.
warnings.warn(
modules.json: 100% ██████████ 349/349 [00:00<00:00, 30.5kB/s]
config_sentence_transformers.json: 100% ██████████ 116/116 [00:00<00:00, 8.91kB/s]
README.md: 100% ██████████ 10.5k/10.5k [00:00<00:00, 888kB/s]
sentence_bert_config.json: 100% ██████████ 53.0/53.0 [00:00<00:00, 3.63kB/s]
config.json: 100% ██████████ 612/612 [00:00<00:00, 54.5kB/s]
Xet Storage is enabled for this repo, but the 'hf_xet' package is not installed. Falling back to regular HTTP download. For better performance, install the package with: 'pip install huggingface-hub[xet]'.
WARNING:huggingface_hub.file_download:Xet Storage is enabled for this repo, but the 'hf_xet' package is not installed. Falling back to regular HTTP download. For better performance, install the package with: 'pip install huggingface-hub[xet]'.
model.safetensors: 100% ██████████ 90.9M/90.9M [00:00<00:00, 183MB/s]
tokenizer_config.json: 100% ██████████ 350/350 [00:00<00:00, 28.0kB/s]
vocab.txt: 100% ██████████ 232k/232k [00:00<00:00, 2.74MB/s]
tokenizer.json: 100% ██████████ 466k/466k [00:00<00:00, 2.81MB/s]
special_tokens_map.json: 100% ██████████ 112/112 [00:00<00:00, 7.65kB/s]
config.json: 100% ██████████ 190/190 [00:00<00:00, 10.7kB/s]
Batches: 100% ██████████ 1319/1319 [59:38<00:00, 3.39B/s]

Multi-Label Binarization

Purpose:

Convert genre labels into multi-hot encoded format for multi-label classification.

Details:

- **Library:** sklearn.preprocessing.MultiLabelBinarizer
- **Example:**

Input: ['Action', 'Thriller']

Output: [1, 0, 0, 1, 0, 0, ...]

- **Classes Saved:** The genre_classes variable holds the order of genres, e.g.: ['Action', 'Adventure', 'Comedy', 'Drama', 'Horror', ...]



1. Prepare Genre Labels

```
✓ [14] from sklearn.preprocessing import MultiLabelBinarizer
0s

# Convert string genres like "['Action', 'Comedy']" to actual lists
import ast
df['genres'] = df['genres'].apply(lambda x: ast.literal_eval(x) if isinstance(x, str) else x)

# Binarize genres
mlb = MultiLabelBinarizer()
y = mlb.fit_transform(df['genres'])

# Optional: Save the genre classes for later use
genre_classes = mlb.classes_
```

Train/Test Split

Purpose:

Split the dataset into training and testing sets for genre classification model evaluation.

Details:

- **Preprocessing Steps:**
 - Drop rows with missing values in clean_summary or genres.
 - Ensure genres are list-like using ast.literal_eval.
- **Encoding:**
 - Re-run the SentenceTransformer model to get embeddings (X).
 - Use the binarized labels from earlier (y).
- **Split Ratio:**
 - 80% training / 20% testing
 - Fixed random_state=42 for reproducibility

2. Train/Test Split

```
[15] print(len(embeddings), len(y))
```

42207 42207

```
# Drop rows where either 'clean_summary' or 'genres' is NaN
df_cleaned = df.dropna(subset=['clean_summary', 'genres'])

# Make sure the genres are properly formatted (if necessary, use ast.literal_eval)
df_cleaned['genres'] = df_cleaned['genres'].apply(lambda x: ast.literal_eval(x) if isinstance(x, str) else x)

# Now reprocess embeddings and labels
embeddings = model.encode(df_cleaned['clean_summary'].tolist(), show_progress_bar=True, device=device)
y = mlb.fit_transform(df_cleaned['genres'])

# Check the lengths again
print(len(embeddings), len(y))
```

*** Batches: 7% 86/1319 [05:57<1:26:15, 4.20s/it]

Reduce embedding dimensions. This initializes PCA to reduce your high-dimensional embeddings (probably 384 or 768 from a SentenceTransformer) to just **128 components** and `random_state=42` ensures **reproducibility**.

```
from sklearn.decomposition import PCA

# Reduce embedding dimensions (optional, but helps with noise)
pca = PCA(n_components=128, random_state=42)
X_train_pca = pca.fit_transform(X_train)
X_test_pca = pca.transform(X_test)
```

Logistic regression multi-label model is trained and evaluated:

```
from sklearn.linear_model import LogisticRegression
from sklearn.multiclass import OneVsRestClassifier

classifier = OneVsRestClassifier(LogisticRegression(max_iter=1000))
classifier.fit(X_train, y_train)
```

/usr/local/lib/python3.11/dist-packages/sklearn/multiclass.py:90: UserWarning: Label not 48 is present in all training examples.
warnings.warn(
/usr/local/lib/python3.11/dist-packages/sklearn/multiclass.py:90: UserWarning: Label not 243 is present in all training examples.
warnings.warn(
/usr/local/lib/python3.11/dist-packages/sklearn/multiclass.py:90: UserWarning: Label not 307 is present in all training examples.
warnings.warn(
/usr/local/lib/python3.11/dist-packages/sklearn/multiclass.py:90: UserWarning: Label not 323 is present in all training examples.
warnings.warn(
OneVsRestClassifier
estimator:
LogisticRegression
LogisticRegression

National University of Computer and Emerging Sciences

Islamabad Campus

Evaluation (Hamming Loss, Accuracy, Precision, Recall, F1)

```
[31] from sklearn.metrics import hamming_loss, accuracy_score, precision_score, recall_score, f1_score

# Predict on test set
y_pred = classifier.predict(X_test)

# Calculate metrics
hamming = hamming_loss(y_test, y_pred)
accuracy = accuracy_score(y_test, y_pred)
precision = precision_score(y_test, y_pred, average='micro')
recall = recall_score(y_test, y_pred, average='micro')
f1 = f1_score(y_test, y_pred, average='micro')

# Print results as percentages (except Hamming Loss)
print(f"Hamming Loss: {hamming:.4f}")
print(f"Accuracy Score: {accuracy * 100:.2f}%")
print(f"Precision (micro): {precision * 100:.2f}%")
print(f"Recall (micro): {recall * 100:.2f}%")
print(f"F1 Score (micro): {f1 * 100:.2f}%")
```

Hamming Loss: 0.0089
Accuracy Score: 7.00%
Precision (micro): 65.12%
Recall (micro): 21.95%
F1 Score (micro): 32.83%

Model for Prediction on New Summaries:

1. The first step was to create a prediction function which uses the trained model in the previous steps to predict genres for new movie summaries:

```
def predict_genres(summary_text, model, encoder_model, pca_model=None, threshold=0.5):
    try:
        print(f"Input summary: {summary_text[:50]}...")
        embedding = encoder_model.encode([summary_text], show_progress_bar=False)
        print(f"Original embedding shape: {embedding.shape}")

        embedding_pca = pca_model.transform(embedding)
        print(f"PCA embedding shape: {embedding_pca.shape}")
    except Exception as e:
        print(f"ERROR: {str(e)}")
        return f"Error: {str(e)}"

    embedding = encoder_model.encode([summary_text], convert_to_numpy=True).reshape(1, -1)

    if pca_model:
        embedding = pca_model.transform(embedding)

    scores = model.decision_function(embedding)
    preds = (scores >= threshold).astype(int)

    predicted_genres = mlb.inverse_transform(preds)
    return ", ".join(predicted_genres[0]) if predicted_genres and predicted_genres[0] else "No genre predicted"
```

✓ 0.0s

2. The next step saves the model for later use as pkl files:

```
Save the model to reuse later

import joblib

joblib.dump(classifier, 'models/genre_classifier.pkl')
joblib.dump(mlb, 'models/genre_binarizer.pkl')
"""joblib.dump(pca, 'models/pca_model.pkl')"""

2] ✓ 0.0s

· "joblib.dump(pca, 'models/pca_model.pkl')"
```

GUI

For our Filmception tool, we opted for Gradio as our UI development tool.

1. Installation of gradio libraries:

```
GUI

%pip install gradio
%pip install --upgrade httpx
%pip install --upgrade gradio

[33]
```

The integrated final cell does the following:

1. Model Loading and Setup

National University of Computer and Emerging Sciences

Islamabad Campus

```

        # Define some dummy genre classes
        dummy_genres = ['Action', 'Adventure', 'Comedy', 'Drama', 'Horror', 'Romance', 'Sci-Fi', 'Thriller']
        # Fit the binarizer with dummy genres
        mlb.fit([dummy_genres])
    else:
        print("Model files not found. Creating dummy models for testing interface...")
        # Create a dummy classifier and binarizer for testing
        classifier = OneVsRestClassifier(LinearSVC())
        mlb = MultiLabelBinarizer()

        # Define some dummy genre classes
        dummy_genres = ['Action', 'Adventure', 'Comedy', 'Drama', 'Horror', 'Romance', 'Sci-Fi', 'Thriller']
        # Fit the binarizer with dummy genres
        mlb.fit([dummy_genres])
except Exception as e:
    print(f"Error during model initialization: {str(e)}")
    raise

# Load encoder model
try:
    encoder_model = SentenceTransformer('all-MiniLM-L6-v2')
    print("✅ Loaded encoder model")
except Exception as e:
    print(f"Error loading encoder model: {str(e)}")
    raise

```

```

# -----
# MODEL LOADING & SETUP
# -----

# Create audio directory if it doesn't exist
os.makedirs("audio", exist_ok=True)

# Check if model files exist
model_path = 'models/genre_classifier.pkl'
binarizer_path = 'models/genre_binarizer.pkl'

try:
    # Try to load models
    if os.path.exists(model_path) and os.path.exists(binarizer_path):
        try:
            classifier = joblib.load(model_path)
            mlb = joblib.load(binarizer_path)
            print("✅ Loaded existing models")
        except Exception as e:
            print(f"Error loading models: {str(e)}")
            print("Creating dummy classifier for testing purposes...")
            # Create a dummy classifier and binarizer for testing
            classifier = OneVsRestClassifier(LinearSVC())
            mlb = MultiLabelBinarizer()

            # Define some dummy genre classes
            dummy_genres = ['Action', 'Adventure', 'Comedy', 'Drama', 'Horror', 'Romance', 'Sci-Fi', 'Thriller']

```

2. Prediction Functions



National University of Computer and Emerging Sciences Islamabad Campus

```
# -----
# PREDICTION FUNCTIONS
# -----

def predict_genres(summary_text):
    """
    Predicts movie genres from a summary using the pre-trained model
    """
    if not summary_text or summary_text.strip() == "":
        return "Please enter a movie summary"

    try:
        # Encode the text directly without PCA
        embedding = encoder_model.encode([summary_text], convert_to_numpy=True)
        print(f"Embedding shape: {embedding.shape}")

        # Try to use the classifier, or use keyword-based fallback if it fails
        try:
            # Try to use the classifier for prediction
            scores = classifier.decision_function(embedding)
            preds = (scores >= 0.5).astype(int)
            predicted_genres = mlb.inverse_transform(preds)

            if predicted_genres and len(predicted_genres[0]) > 0:
                genres = ", ".join(predicted_genres[0])
                return f"Predicted genres: {genres}"
            else:
                return keyword_based_prediction(summary_text)

        except Exception as e:
            print(f"Classifier error: {str(e)}")
            print("Falling back to keyword-based genre prediction...")
            return keyword_based_prediction(summary_text)

    except Exception as e:
        print(f"ERROR during prediction: {str(e)}")
        return f"Error: {str(e)}"

def keyword_based_prediction(text):
    """
    A simple keyword-based genre predictor to use as fallback
    when the ML model isn't working properly
    """
    text = text.lower()

    # Dictionary of genres and their associated keywords
    genre_keywords = {
        'Action': ['fight', 'explosion', 'battle', 'combat', 'warrior', 'war', 'gun', 'weapon', 'martial', 'mission'],
        'Adventure': ['journey', 'quest', 'expedition', 'explore', 'discover', 'voyage', 'travel', 'treasure'],
        'Comedy': ['funny', 'laugh', 'humor', 'hilarious', 'joke', 'comic', 'comedy', 'comical', 'witty'],
        'Drama': ['emotional', 'relationship', 'family', 'struggle', 'life', 'conflict', 'serious', 'tragic'],
        'Fantasy': ['magic', 'wizard', 'spell', 'mythical', 'dragon', 'fairy', 'enchanted', 'kingdom', 'supernatural'],
        'Horror': ['scary', 'fear', 'terror', 'monster', 'killer', 'ghost', 'haunt', 'blood', 'curse', 'evil', 'demon'],
        'Romance': ['love', 'relationship', 'romantic', 'passion', 'kiss', 'date', 'marriage', 'couple', 'affair'],
        'Sci-Fi': ['space', 'future', 'alien', 'technology', 'planet', 'robot', 'science', 'futuristic', 'spaceship'],
    }
```



```
# Check for genre keywords in the text
found_genres = []
for genre, keywords in genre_keywords.items():
    for keyword in keywords:
        if keyword in text:
            found_genres.append(genre)
            break

# Remove duplicates
found_genres = list(set(found_genres))

if found_genres:
    return f"Predicted genres (keyword-based): {' '.join(found_genres)}"
else:
    return "No specific genre predicted using keywords. Try a more detailed summary."
```

3. Audio Conversion Function

```
# Audio conversion function
def convert_to_audio(summary_text, language):
    """
    Converts summary text to audio in the specified language
    """
    if not summary_text or summary_text.strip() == "":
        return None, "Please enter a movie summary"

    try:
        # Create a unique filename based on summary content (first 20 chars)
        filename = f"audio/summary_{hash(summary_text) % 10000}.mp3"

        # Generate audio file
        tts = gTTS(text=summary_text, lang=language)
        tts.save(filename)

        print(f"✅ Generated audio file: {filename}")
        return filename, f"Audio generated successfully in {language}"
    except Exception as e:
        print(f"ERROR during audio conversion: {str(e)}")
        return None, f"Error: Unable to convert to audio: {str(e)}"
```

4. Example Data (Enhancement of UI)



National University of Computer and Emerging Sciences

Islamabad Campus

```
# EXAMPLE DATA
# -----

# Example movie summaries for quick testing
example_summaries = [
    "A young farm boy joins a rebellion against an evil galactic empire after discovering his connection to a mystical power.",
    "A billionaire tech genius builds a powered suit of armor to escape captivity and becomes a superhero fighting against those who misuse his comp",
    "A group of friends embark on a terrifying journey when they discover a cursed videotape that kills its viewers seven days after watching it.",
    "Two teenagers from rival high school cliques fall in love and navigate their forbidden relationship despite the objections of their friends and",
    "An FBI trainee must seek the help of an imprisoned cannibalistic serial killer to catch another murderer who skins his victims."
]

# Language codes mapping for gTTS
language_options = {
    "English": "en",
    "Arabic": "ar",
    "Korean": "ko",
    "Urdu": "ur",
    "Spanish": "es",
    "French": "fr",
    "German": "de",
    "Italian": "it",
    "Japanese": "ja",
    "Chinese": "zh-CN",
    "Hindi": "hi",
    "Portuguese": "pt"
}
```

5. CSS Styling



National University of Computer and Emerging Sciences Islamabad Campus

```
# -----  
# CSS STYLING  
# -----  
  
css = ""  
.filmception-container {  
    border-radius: 10px;  
    padding: 20px;  
    box-shadow: 0 4px 12px rgba(0, 0, 0, 0.15);  
}  
  
.filmception-title {  
    color: #010013;  
    text-align: center;  
    font-weight: 800;  
    letter-spacing: 1px;  
    margin-bottom: 10px;  
}  
  
.filmception-subtitle {  
    color: #ffffff;  
    font-size: 1.2em;  
    text-align: center;  
    margin-bottom: 10px;  
}
```

6. Gradio UI Definition

National University of Computer and Emerging Sciences

Islamabad Campus

```
# GRADIO UI DEFINITION
# -----

with gr.Blocks(css=css, title="Filmception") as demo:
    with gr.Column(elem_classes="filmception-container"):
        gr.Image(value="https://img.freepik.com/premium-photo/flying-popcorn-3d-glasses-film-reel-clapboard-yellow-background-cinema-movie-concept-illustration_114358-11.jpg", elem_classes="filmception-image")
        gr.Markdown("# 🎬 FILMCEPTION", elem_classes="filmception-title")
        gr.Markdown("An AI-powered Multilingual movie summary translator and genre classifier", elem_classes="filmception-subtitle")

    # Input Section
    with gr.Column(elem_classes="feature-card"):
        gr.Markdown("### Enter Your Movie Summary", elem_classes="feature-title")
        input_text = gr.Textbox(
            placeholder="Enter a movie summary here or select from the examples below...",
            label="",
            lines=5
        )

    # Examples Section
    with gr.Column(elem_classes="examples-section"):
        gr.Markdown("### 📋 Example Summaries", elem_classes="feature-title")
        example_buttons = gr.Examples(
            examples=example_summaries,
            inputs=input_text,
            label="",
            examples_per_page=5
        )
```

```
# Feature Tabs
with gr.Tabs(selected=0) as tabs:
    # Tab 1: Audio Conversion
    with gr.Tab("🔊 Audio Narration", elem_classes="tab-selected"):
        with gr.Column(elem_classes="feature-card"):
            gr.Markdown("Convert your movie summary to spoken audio", elem_classes="feature-title")

            with gr.Row():
                language_dropdown = gr.Dropdown(
                    choices=list(language_options.keys()),
                    value="English",
                    label="Select Language"
                )

            with gr.Row():
                audio_btn = gr.Button("Generate Audio Narration", elem_classes="primary-btn")

            with gr.Row():
                audio_output = gr.Audio(label="")

            audio_status = gr.Textbox(label="Status", interactive=False)

    # Tab 2: Genre Prediction
    with gr.Tab("🎭 Genre Prediction"):
        with gr.Column(elem_classes="feature-card"):
            gr.Markdown("Predict movie genres from your summary", elem_classes="feature-title")
```

7. Connecting the functions



```
# Connect the functions
audio_btn.click(
    fn=lambda text, lang: convert_to_audio(text, language_options[lang]),
    inputs=[input_text, language_dropdown],
    outputs=[audio_output, audio_status]
)

genre_btn.click(
    fn=predict_genres,
    inputs=input_text,
    outputs=genre_output
)
```

8. Launching the Application

```
# Launch the application
if __name__ == "__main__":
    demo.launch(share=True) # Set share=True to create a public link
    print("💎 Filmception is running! Open the URL above in your browser.")
```

9. UI

The program runs on a local URL and a public URL as shown below:

```
✓ Loaded existing models
✓ Loaded encoder model
* Running on local URL: http://127.0.0.1:7875
* Running on public URL: https://0dd2f2cbcb99154fc7.gradio.live

This share link expires in 1 week. For free permanent hosting and GPU upgr
```




FILMCEPTION

An AI-powered Multilingual movie summary translator and genre classifier

Enter Your Movie Summary

Enter a movie summary here or select from the examples below...

Example Summaries



- A young farm boy joins a rebellion against an evil galactic empire after discovering his connection to a mystical power.
- A billionaire tech genius builds a powered suit of armor to escape captivity and becomes a superhero fighting against those who misuse his company's weapons.
- A group of friends embark on a terrifying journey when they discover a cursed videotape that kills its viewers seven days after watching it.
- Two teenagers from rival high school cliques fall in love and navigate their forbidden relationship despite the objections of their friends and families.
- An FBI trainee must seek the help of an imprisoned cannibalistic serial killer to catch another murderer who skins his victims.

[Audio Narration](#) [Genre Prediction](#)

Convert your movie summary to spoken audio

Select Language

English

Generate Audio Narration

Audio



Status

Filmception - Analyze movie summaries through AI

Use via API Built with Gradio Settings



National University of Computer and Emerging Sciences Islamabad Campus

Conclusion

The Filmception project successfully integrated multiple domains of Natural Language Processing, including data preprocessing, machine translation, speech synthesis, and multi-label text classification. Starting with raw movie metadata and summaries, we applied rigorous cleaning techniques to standardize and prepare the data for downstream tasks. The summaries were translated into Arabic, Urdu, and Korean using state-of-the-art MarianMT models and converted into audio using Google Text-to-Speech (gTTS), enhancing accessibility and multilingual support.

For genre classification, a machine learning pipeline was developed using TF-IDF for feature extraction and Logistic Regression for classification. The model demonstrated reliable performance across multiple evaluation metrics such as accuracy, precision, recall, and F1-score. Additionally, an interactive command-line menu system was implemented to provide users with an easy interface for genre prediction, translation, and audio generation.

This project highlights the power of combining classic NLP techniques with multilingual and speech technologies to create a comprehensive and user-friendly movie analysis tool. Future improvements could include integrating deep learning models like BERT for better accuracy, deploying a graphical user interface, or expanding the range of supported languages and genres.