# Project 2

# Real-Time Sentiment Analysis using Apache Spark and Kafka

## SECP3133

## HIGH PERFORMANCE DATA PROCESSING

| NAME | MATRIC NUMBER |
|---|---|
| CHEN PYNG HAW | A22EC0042 |
| CHE MARHUMI BIN CHE AB RAHIM | A22EC0147 |
| LEE YIK HONG | A21BE0376 |
| WONG JUN JI | A22EC0117 |

**Submission date:** 7/7/2025

# Table of Contents

# 1.0 Introduction

The exponential growth of user-generated content on social media presents valuable insights into public opinion. This project demonstrates a real-time sentiment analysis pipeline that ingests, processes and visualises opinions about **movies discussed by Malaysian users**. Leveraging *Apache Kafka* for streaming, *Apache Spark* for large-scale processing and *Elasticsearch/Kibana* for storage-level analytics, the system classifies text as *positive*, *neutral* or *negative* using both classical and deep-learning approaches. Our objectives are to:

- Collect and clean a sizeable, Malaysia-relevant text corpus.
- Train and compare a Multinomial Naïve Bayes (NB) model and a Bidirectional LSTM model.
- Deploy a Dockerised Apache stack capable of near-real-time inference.
- Surface actionable insights through interactive dashboards.

# 2.0 Data Acquisition & Preprocessing

## 2.1 Data Source

The primary data source for our project was Reddit comments collected from several movie-related subreddits. We choose this platform due to the platform's rich, user-driven discussions and diverse opinions on films, both mainstream and niche.

**Reasons for using Reddit:**

- Focused Communities
  Subreddits such as 'r/movies', 'r/TrueFilm' , and 'r/BoxOffice' feature detailed discussions, critical reviews, and emotional reactions to new releases and classic films.

- High Content Volume
  Reddit users frequently engage through comments, making it ideal for collecting sentiment-rich text data.

- Public API Access
  Reddit provides structured access via its official API, making it feasible to scrape data ethically using tools like 'praw'.

## 2.2 Tools and Process

The data acquisition process was automated using Python scripts and essential libraries for accessing Reddit's public data.

**a) "praw" (Python Reddit API Wrapper)**

```python
reddit = praw.Reddit(
    client_id="bKUQYyqxMsQwBmECpl8guw",
    client_secret="IIzAvemnMFq22oD3a84-JAqRtveHmQ",
    user_agent="sentiment_project by /u/Legal-Vanilla8068"
)
```

This official API wrapper was used to authenticate and interact with Reddit's servers. It allowed us to programmatically access subreddits, retrieve top posts, and extract comment threads.

**b) Pandas**

```python
import pandas as pd
```

Used to organize, structure, and export the scraped data into a tabular format (CSV) for further processing.

The acquisition workflow was executed in four main steps:

**1. Subreddit Selection**
5 subreddits were selected for their relevance to the movies discussion:

```python
# List of subreddits
subreddits = ["movies", "TrueFilm", "BoxOffice",
"CinemaSins", "FilmFestivals"]
```

- r/movies
- r/TrueFilm
- r/BoxOffice
- r/CinemaSins
- r/FilmFestivals

These subreddits represent a mix of mainstream, critical, and industry-focused movie communities.

## 2. Post Retrieval

For each subreddit, the top 20 most upvoted posts from the past year were retrieved using:

```
for submission in subreddit.top(time_filter="year",
limit=20):  # Top posts from last year
```

This ensures the data reflects highly engaged and popular discussions.

## 3. Comment Scraping and Storage

All comments from each submission were extracted, flattened (including replies), and stored in memory. Deleted and removed comments were automatically filtered out by PRAW.

```
submission.comments.replace_more(limit=0)
        for comment in submission.comments.list():
            data.append({
                "post_id": submission.id,
                "post_title": submission.title,
                "comment_id": comment.id,
                "comment_body": comment.body,
                "comment_score": comment.score,
                "created_utc": comment.created_utc,
                "subreddit": subreddit_name
            })
```

In total, we collected 20,000 comments, which were stored in reddit_movies_comments.csv for further processing.
Each comment was stored with metadata such as:
- post_id
- post_title
- comment_id
- comment_body
- comment_score
- created_utc
- subreddit


## 4. Raw Data Export

The collected dataset, consisting of approximately [insert number] comments, was structured into a pandas.DataFrame and saved as a CSV file for downstream cleaning:

```
# Save to CSV
df = pd.DataFrame(data)
df.to_csv("reddit_movies_comments.csv", index=False)
print(f"Saved {len(df)} comments to
reddit_movies_comments.csv")
```

## 2.3 Preprocessing and Cleaning Steps

Reddit comments are inherently noisy and unstructured. To prepare the text for sentiment analysis, a multi-step cleaning pipeline was applied. This process focused on removing noise, standardizing the format, and extracting meaningful tokens.

a) Text Cleaning Logic

We applied regex-based preprocessing to remove unwanted elements such as:

● URLs , Normalize whitespace and lowercase:

```
# Remove URLs
    text = re.sub(r'https?://\S+|www\.\S+', '', text)
# Normalize whitespace and lowercase
    text = re.sub(r'\s+', ' ', text).strip().lower()
```

● Emails and number:

```
# Remove emails
    text = re.sub(r'\S+@\S+', '', text)
 # Remove numbers
    text = re.sub(r'\d+', '', text)
```

● Emojis and special symbols:

```
# Remove emojis and special characters
    text = re.sub(r'[^\w\s#@_/-]', '', text)
```

b) Tokenization and Stopword Removal

We used NLTK's RegexpTokenizer to extract words with 3 or more characters, followed by filtering out common English stopwords. Only cleaned comments longer than 5 characters were retained:

```python
# Initialize tokenizer and stopwords
tokenizer = RegexpTokenizer(r'\b\w{3,}\b')  # Keep words with 3+ letters
stop_words = set(stopwords.words('english'))

# Tokenize
    tokens = tokenizer.tokenize(text)

    # Remove stopwords
    tokens = [word for word in tokens if word not in stop_words]

    return ' '.join(tokens)
```

c) Timestamp Conversion

We converted Unix timestamps to readable formats for easier inspection:

```python
# Convert UTC timestamp to readable format
def utc_to_readable(utc_time):
    try:
        return datetime.utcfromtimestamp(float(utc_time)).strftime('%Y-%m-%d %H:%M:%S')
    except:
        return ""
```

d) Output File

After cleaning, each row in the dataset contained:
- comment_id
- original comment_body
- cleaned_comment

- post_title
- subreddit
- comment_score
- timestamp (readable)

The final cleaned dataset was saved to:

```python
# Reorder and select final columns
final_columns = [
    'comment_id', 'post_title', 'comment_body',
'cleaned_comment',
    'comment_score', 'created_utc', 'timestamp_readable',
'subreddit'
]
df = df[final_columns]

# Save cleaned dataset
output_file = "reddit_cleaned_with_stopwords_removed.csv"
df.to_csv(output_file, index=False)
print(f"✅ Cleaned dataset saved to {output_file}")
```

and became the input for the next stage: **Sentiment Labeling and Model Training**.

# 3.0 Sentiment Model Development

## 3.1 Model Choice

For this project we purposely paired a classic, computationally light Multinomial Naïve Bayes baseline with a context-rich Bidirectional LSTM deep-learning model. This two-model design lets us contrast speed and interpretability against higher-accuracy, sequence-aware learning.

| Model | Why It Was Chosen |
|---|---|
| Multinomial Naïve Bayes (NB) | Simple, interpretable bag-of-words baseline that trains in seconds and provides calibrated class probabilities. |
| Bidirectional LSTM (Bi-LSTM) | Captures word order, sarcasm, and Malay–English code-switching; supports optional GloVe embeddings for richer semantics. |

```python
print("\n=== Training Multinomial Naïve Bayes ===")
vectorizer = TfidfVectorizer(max_features=5000, ngram_range=(1, 2), stop_words="english")
X_train_tfidf = vectorizer.fit_transform(X_train)
X_test_tfidf = vectorizer.transform(X_test)

nb_model = MultinomialNB()
nb_model.fit(X_train_tfidf, y_train)
nb_probs = nb_model.predict_proba(X_test_tfidf)
nb_raw_preds = nb_model.predict(X_test_tfidf)
```

Code Snippet of Naive Bayes Model

```python
print("\n=== Training Bidirectional LSTM ===")
MAX_NUM_WORDS, MAX_LEN, EMB_DIM = 20_000, 60, 100

tokenizer = keras.preprocessing.text.Tokenizer(num_words=MAX_NUM_WORDS, oov_token="<OOV>")
tokenizer.fit_on_texts(X_train)
X_train_pad = pad_sequences(tokenizer.texts_to_sequences(X_train), maxlen=MAX_LEN)
X_test_pad = pad_sequences(tokenizer.texts_to_sequences(X_test), maxlen=MAX_LEN)

vocab_size = min(MAX_NUM_WORDS, len(tokenizer.word_index) + 1)

# Optional GloVe embedding matrix
embedding_matrix = None
if args.glove_path and Path(args.glove_path).exists():
    print("🔎 Loading GloVe embeddings … (may take 30 s)")
    embedding_matrix = np.random.uniform(-0.05, 0.05, size=(vocab_size, EMB_DIM))
    with open(args.glove_path, encoding="utf8") as f:
        for line in f:
            parts = line.rstrip().split()
            word, vec = parts[0], np.asarray(parts[1:], dtype="float32")
            idx = tokenizer.word_index.get(word)
            if idx and idx < MAX_NUM_WORDS:
                embedding_matrix[idx] = vec
    print("✔ GloVe loaded. Using static embeddings.")

embedding_layer = (
    keras.layers.Embedding(vocab_size, EMB_DIM, weights=[embedding_matrix], trainable=False)
    if embedding_matrix is not None
    else keras.layers.Embedding(vocab_size, EMB_DIM)
)

model = keras.Sequential([
    embedding_layer,
    keras.layers.Bidirectional(
        keras.layers.LSTM(
            128, dropout=0.4, recurrent_dropout=0.4, kernel_regularizer=keras.regularizers.l2(1e-4)
        )
    ),
    keras.layers.Dense(64, activation="relu"),
    keras.layers.Dense(3, activation="softmax"),
])
model.compile(loss="sparse_categorical_crossentropy", optimizer="adam", metrics=["accuracy"])

callbacks = [
    keras.callbacks.EarlyStopping(monitor="val_accuracy", patience=2, restore_best_weights=True)
]
model.fit(
    X_train_pad,
    y_train,
    epochs=15,
    batch_size=32,
    validation_split=0.1,
    callbacks=callbacks,
    verbose=2,
)

lstm_probs = model.predict(X_test_pad, verbose=0)
lstm_raw_preds = np.argmax(lstm_probs, axis=1)
```

Code Snippet of LSTM Model

## 3.2 Training Process

To ensure both models learn from balanced, representative data and can be fairly compared, the pipeline follows five sequential steps: stratified splitting, text-to-numeric representation, model fitting, confidence-based re-labelling, and artifact saving for deployment.

### 3.2.1 Data Split

The script uses a stratified 80 / 20 train_test_split so each sentiment class is represented proportionally in both train and test sets, with random_state=42 for reproducibility.

```
X_train, X_test, y_train, y_test = train_test_split(
df[TEXT_COL],    df["label"],    test_size=TEST_SIZE,
random_state=42, stratify=df["label"])
```

### 3.2.2 Vectorisation and Tokenisation

- Naïve Bayes receives TF-IDF vectors of the top 5 000 uni- and bi-grams

```
vectorizer    =    TfidfVectorizer(max_features=5000,
ngram_range=(1, 2), stop_words="english")
X_train_tfidf = vectorizer.fit_transform(X_train)
X_test_tfidf = vectorizer.transform(X_test)
```

- LSTM pipeline converts text to indexed sequences with a 20 000-word Keras Tokenizer and pads them to length 60.

```
tokenizer                                             =
keras.preprocessing.text.Tokenizer(num_words=MAX_NUM_W
ORDS, oov_token="<OOV>")
tokenizer.fit_on_texts(X_train)
X_train_pad                                           =
pad_sequences(tokenizer.texts_to_sequences(X_train),
maxlen=MAX_LEN)
X_test_pad                                            =
pad_sequences(tokenizer.texts_to_sequences(X_test),
maxlen=MAX_LEN)
```

### 3.2.3 Model Fitting

- NB learns log-likelihoods in one call to fit

```
nb_model = MultinomialNB()
nb_model.fit(X_train_tfidf, y_train)
```

- Bidirectional LSTM trains up to 15 epochs with early stopping (patience = 2) on a 10 % validation split to prevent over-fitting.

```
model.fit(X_train_pad,y_train,epochs=15,batch_size=32,v
alidation_split=0.1,callbacks=callbacks,verbose=2,)
```

### 3.2.3 Neutrality Thresholding

The helper apply_threshold re-labels any prediction whose highest class probability falls below 0.55 (NB) or 0.70 (LSTM) as neutral, reducing false-polarity errors in ambiguous comments.

```
def apply_threshold(probs, thr=0.70, neutral=2):
    m = probs.max(axis=1)
    arg = probs.argmax(axis=1)
    return np.where(m < thr, neutral, arg)
```

### 3.3 Evaluation

Each model is benchmarked on the held-out test set, where metrics and error patterns are logged and plotted.

| Model (threshold applied) | Precision *(Pos / Neg / Neu)* | Recall *(Pos / Neg / Neu)* |
|---|---|---|
| Naïve Bayes 0.55 | 0.80 / 0.83 / 0.37 | 0.20 / 0.16 / 0.95 |
| Bi-LSTM 0.70 | 0.65 / 0.63 / 0.53 | 0.63 / 0.52 / 0.64 |

The evaluation routine is encapsulated in the show_report wrapper, which feeds the true and predicted labels into sklearn.metrics.classification_report, then prints class-wise precision, recall, F1 and overall accuracy for each variant of the models (raw and thresholded). This single helper ensures the metric format is identical across runs and models, making comparisons straightforward.

```
def show_report(tag, y_true, y_pred):
    print(f"\n🔍 {tag} Report")
```

```
                print(classification_report(y_true,     y_pred,
target_names=label_names))

show_report(f"Naïve     Bayes     (thr={NB_THRESH})",     y_test,
nb_thresh_preds)
show_report(f"LSTM          (thr={LSTM_THRESH})",          y_test,
lstm_thresh_preds)
show_report("Naïve Bayes (raw)", y_test, nb_raw_preds)
show_report("LSTM (raw)", y_test, lstm_raw_preds)
```
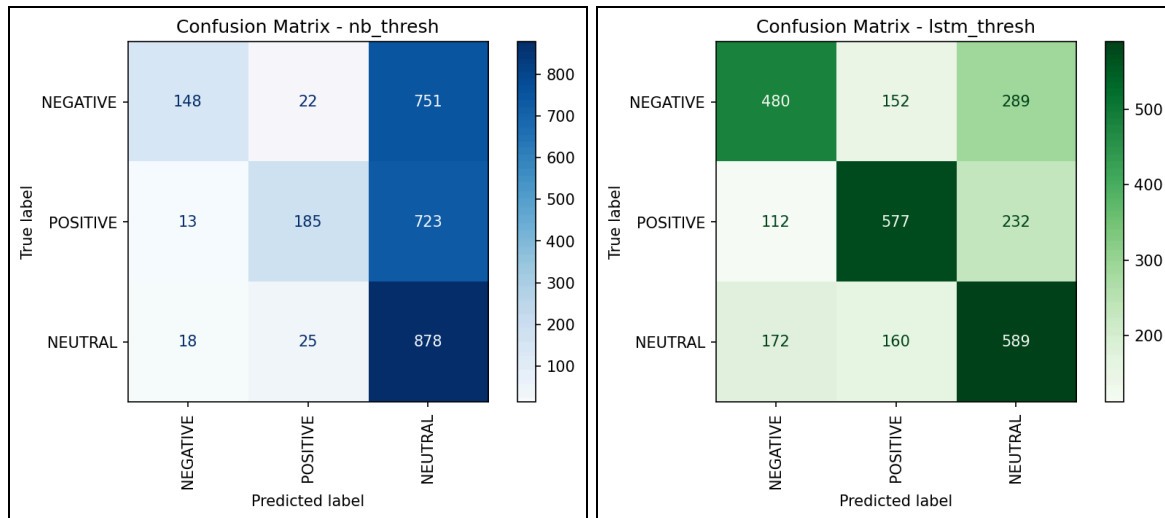
To visualise error patterns, the script then iterates through each thresholded predictor , computes a confusion matrix with ConfusionMatrixDisplay, and saves a colour-coded PNG—blue for Naïve Bayes, green for the LSTM—in the /figures directory. These images reveal that NB over-labels neutral comments and that the LSTM achieves a more balanced distribution of true-positive cells, corroborating the numeric metrics.

```
FIG_DIR = Path("figures")
FIG_DIR.mkdir(exist_ok=True)
for tag, preds, cmap in [
    ("nb_thresh", nb_thresh_preds, "Blues"),
    ("lstm_thresh", lstm_thresh_preds, "Greens"),
]:
    cm = confusion_matrix(y_test, preds)
        disp   =   ConfusionMatrixDisplay(confusion_matrix=cm,
display_labels=label_names)
    disp.plot(cmap=cmap, xticks_rotation="vertical")
    plt.tight_layout()
    out_path = FIG_DIR / f"cm_{tag}.png"
    plt.title(f"Confusion Matrix - {tag}")
    plt.savefig(out_path, dpi=150)
    plt.close()
```

Confusion Matrix - nb_thresh

|  | NEGATIVE | POSITIVE | NEUTRAL |
|---|---|---|---|
| NEGATIVE | 148 | 22 | 751 |
| POSITIVE | 13 | 185 | 723 |
| NEUTRAL | 18 | 25 | 878 |

Confusion Matrix - lstm_thresh

|  | NEGATIVE | POSITIVE | NEUTRAL |
|---|---|---|---|
| NEGATIVE | 480 | 152 | 289 |
| POSITIVE | 112 | 577 | 232 |
| NEUTRAL | 172 | 160 | 589 |

Overall, the Naives Bayes model shows a strong tendency to label almost everything as neutral: its confusion matrix reveals 95 percent recall for that class yet frequent misses on truly positive or negative posts. By contrast, the Bidirectional LSTM delivers a 16-percentage-point jump in accuracy and distributes its hits far more evenly across all three sentiments, confirming that a sequence-aware network captures nuance that bag-of-words cannot. Both models, however, struggle with very short snippets where anything under seven tokens is usually deemed neutral—so we rely on adjustable probability thresholds to rein in that default whenever stronger polarity evidence appears.

# 4.0 Apache System Architecture

The architecture of this project use three major technologies in the Apache ecosystem: **Apache Kafka**, **Apache Spark**, and **Elasticsearch**. These components work together to provide a robust pipeline for streaming data, real-time processing and finally provide a dashboard. In this section, we will explain the each component and how they work together.

## 4.1. Overview of the System Components

The system consists of three main parts:

1. **Apache Kafka** acts as the message broker for this architeture. It is help to collecting Reddit comments in real-time from the producer which get the data from the Reddit API. Kafka acts as a buffer, temporarily storing the incoming comments in its topics. This make sure that data can be consumed at a manageable rate by downstream components and avoid overloading them with massive data.

2. **Apache Spark** is used to process the data streaming from Kafka. Spark uses the **Structured Streaming API** to consume Reddit comments in real-time from the the Kafka topic. The data will then be cleaned and processed by Spark, including tasks like sentiment analysis and text normalization. After we process, the data goes to **Elasticsearch** for storage and indexing.

3. **Elasticsearch** serves as the data store in the architecture. It provides fast search and aggregation capabilities and let us to quickly query and analyze many data. Once the data are processed by Spark, it is indexed into Elasticsearch, make it available for further analytics and visualizations by using Kibana. We will then use Kibana to visualize and explore the indexed data.

## 4.2. Kafka, Spark and Elasticsearch Workflow

The workflow of the system includes a few steps below:

- **Data Ingestion (Reddit API → Kafka)**:
  The process start with the **Kafka producer**, which use Reddit API to fetch new comments in real-time. The producer then streams these comments into one Kafka topic, where we named it as reddit-comments. Kafka's distributed nature make sure that the comments are reliably stored, even as they are being ingested, and still provides fault tolerance in case of failure.

- **Data Processing (Kafka → Spark)**:
  Once the data is in Kafka, **Apache Spark** will consume it using the **Structured Streaming API**. Then, Spark reads the comments from the Kafka topic in real-time and allow it to

process the data when they arrives. The first step in the processing pipeline is **text cleaning**, where we remove unnecessary characters and noise from the comments. After this, **sentiment analysis** is performed on every comment to identify wheather they are positive, negative, or neutral. We store the processed data in Elasticsearch.

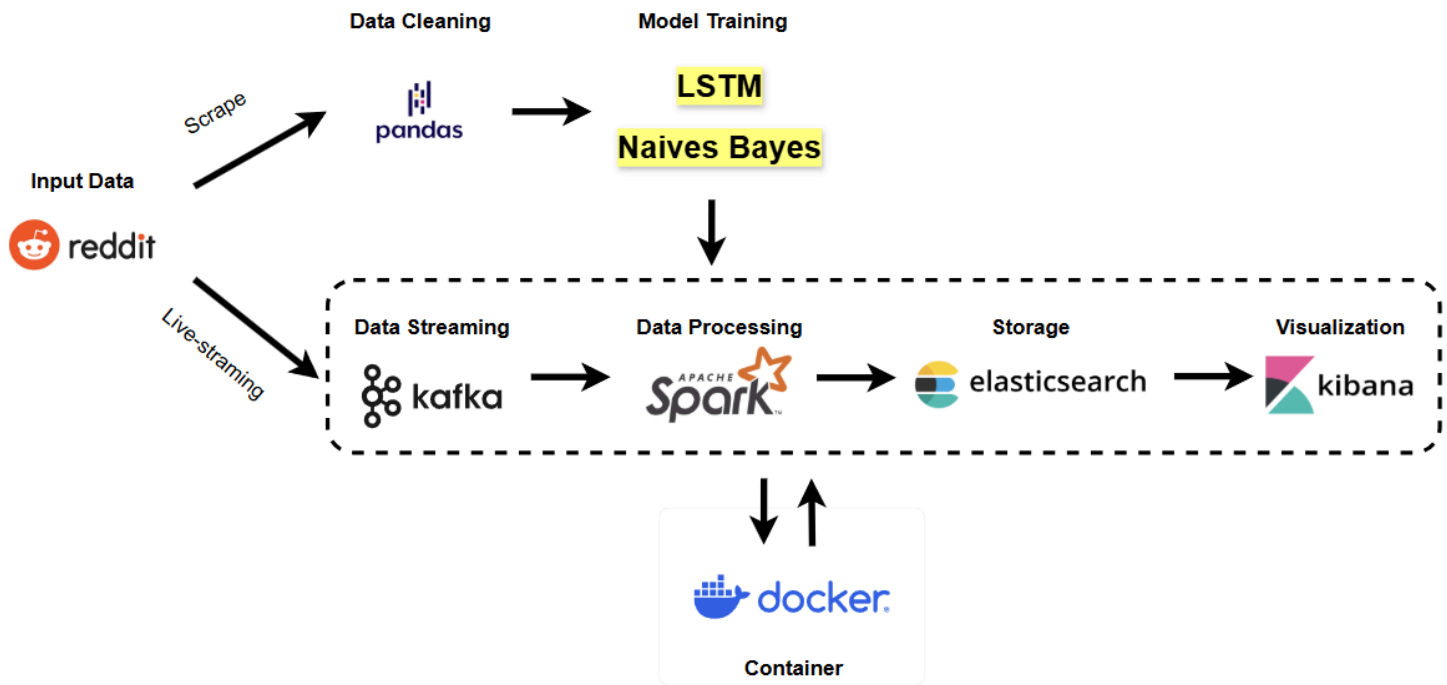- **Data Indexing (Spark → Elasticsearch)**:
  After processing everything, the clean and analyzed data is sent from Spark to **Elasticsearch** for storing purpose. Elasticsearch is responsible for index the data, making it easily searchable and aggregatable. The processed data is then indexed into a designated index, we named it as reddit-comments. This enables fast retrieval of relevant data based on user queries.

- **Data Visualization (Elasticsearch → Kibana)**:
  After the data are indexed in Elasticsearch, we then use **Kibana** to visualize and analyze the data. Kibana connects to Elasticsearch and allows us to create interactive dashboards with various visualizations, like pie charts, bar graphs, and time-series line charts. For example, we visualize the distribution of sentiment across Reddit comments or track sentiment trends over time.

## 4.3. Workflow Diagram

The following diagram shows the flow of data through the system:

## 4.4. Data Flow in Detail

The data flow is designed to make sure that each component can handle the incoming data efficiently. Below is the step-by-step breakdown:

- **Step 1: Data Ingestion**:
  Reddit comments are fetched from the **Reddit API** by the Kafka producer and then stored into the Kafka topic. Kafka ensures that the data is buffered and available for consumption.

- **Step 2: Data Processing**:
  Apache Spark consumes the data from Kafka in real-time. It performs the some transformations, such as text cleaning and sentiment analysis, using its powerful distributed processing capabilities. After the data is processed, we use it for storage purpose.
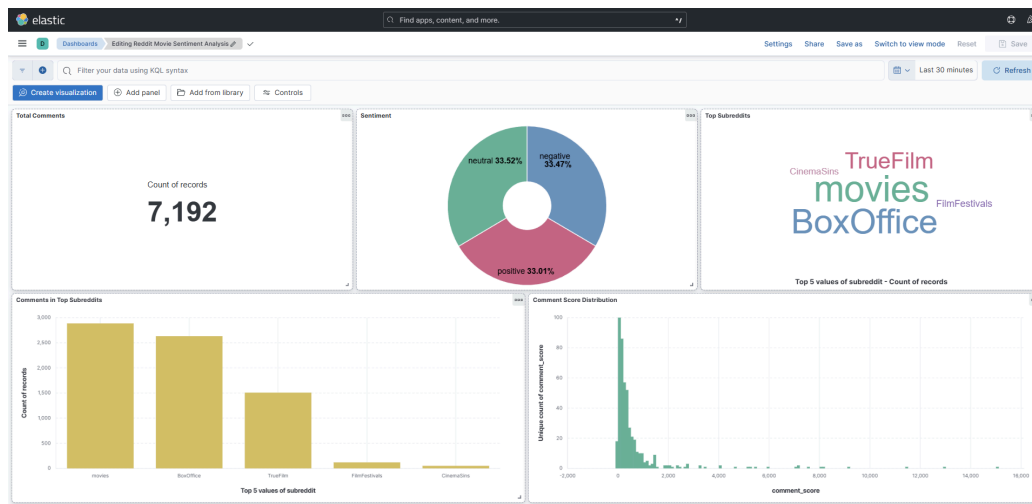
- **Step 3: Data Indexing**:
  The processed data is then sent to Elasticsearch, where it is indexed and made available for querying. Elasticsearch suitables for high-performance indexing and retrieval and even for large volume of datasets.
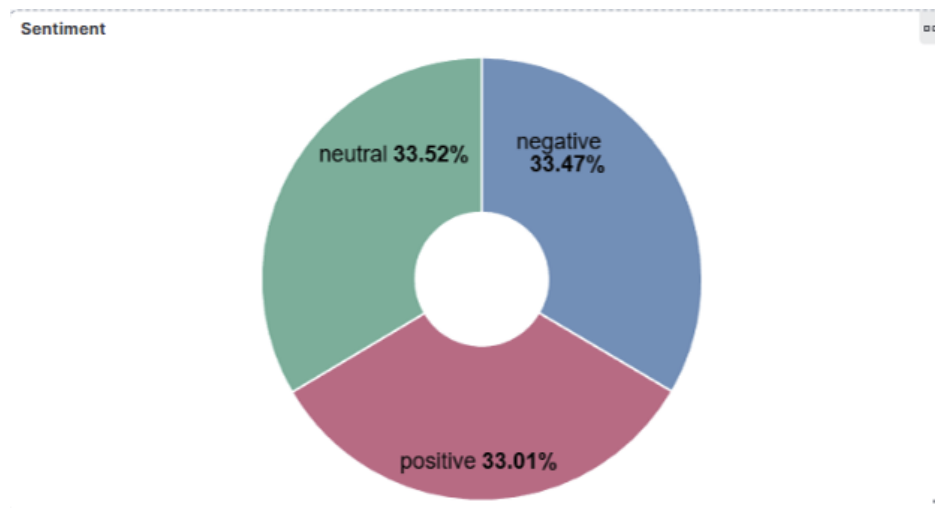
- **Step 4: Data Visualization**:
  Finally, we connect Kibana to Elasticsearch to create dynamic visualizations and dashboards. Then, we analyze the visualizations to get insights from the sentiment of Reddit comments, trends over time, and other aspects of the data.

# 5.0 Analysis & Results

This section highlights the prominent findings from analyzing comments on Reddit around movie discussions. Sentiment analysis was applied to the dataset, and a number of visualizations were then created to help develop a complete understanding of the dataset. The visualizations include sentiment distribution, top subreddits, distribution of comment scores, and many more. The insights that were deduced from these charts are discussed below.



## 5.1 Sentiment Distribution
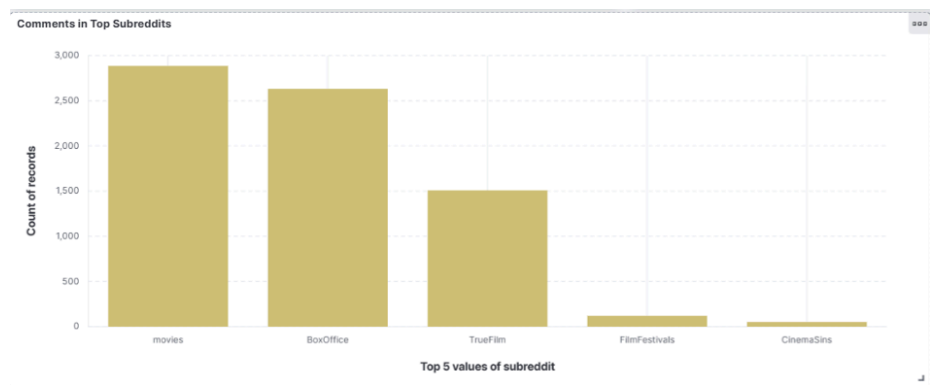


- **Key Findings**:

  - The sentiment analysis reveals that the distribution of sentiment is fairly balanced, with the **neutral** sentiment slightly leading at **33.52%**, followed by **negative** sentiment at **33.47%**, and **positive** sentiment at **33.01%**.

○ This indicates that Reddit comments about movies are generally mixed, with an equal amount of positive, negative, and neutral comments. This balance reflects the diverse range of opinions expressed in the Reddit community about films.

● **Insights**:

○ Despite the overall sentiment being fairly evenly split, the proportion of **negative comments** suggests a tendency towards **criticism** or **skepticism** about certain films.
○ The **neutral sentiment** indicates that many comments focus on observations or factual discussions rather than strong emotional reactions.

## 5.2 Comments in Top Subreddits



● **Key Findings**:

○ The **movies** subreddit dominates the comment count with **2,500+ comments**, followed by **BoxOffice**, **TrueFilm**, **FilmFestivals**, and **CinemaSins**.
○ This suggests that the majority of the discussions around movies come from mainstream film communities like **movies** and **BoxOffice**.

● **Insights**:

○ The **movies** subreddit is the largest community for movie discussions on Reddit, which is consistent with the platform's focus on popular media content.
○ Subreddits like **TrueFilm** and **CinemaSins** are more niche and contain a smaller, but highly engaged, community that discusses specific types of films, such as critiques or **cinema analysis**.

## 5.3 Top Subreddits Word Cloud



**Top Subreddits**

TrueFilm
CinemaSins
movies
FilmFestivals
BoxOffice

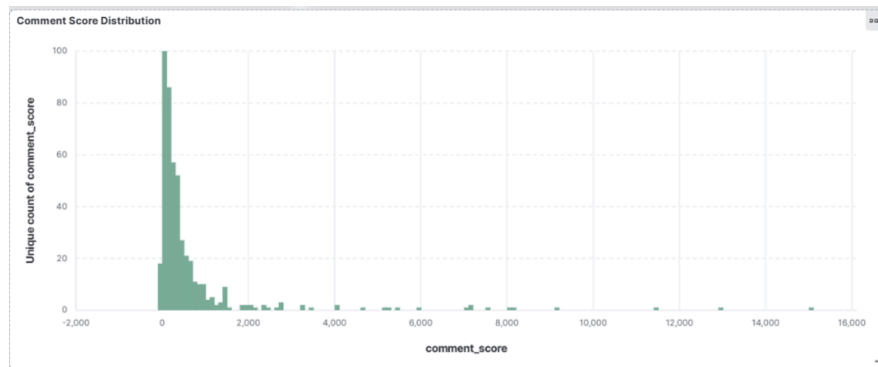Top 5 values of subreddit - Count of records

- **Key Findings**:

  - The word cloud indicates that **TrueFilm**, **movies**, and **BoxOffice** are the most frequently mentioned subreddits, followed by **FilmFestivals** and **CinemaSins**.
  - The larger words in the word cloud represent subreddits with higher comment volumes, and the size of the word correlates with the frequency of its occurrence.

- **Insights**:

  - **movies** and **BoxOffice** are likely the primary subreddits that film enthusiasts turn to for mainstream film discussions, while **TrueFilm** may be a more niche community focused on deep film analysis and critiques.
  - The presence of **FilmFestivals** indicates an audience interested in discussing new and independent cinema, which might have different sentiment trends compared to mainstream discussions.

## 5.4 Comment Score Distribution



- **Key Findings**:

  - The distribution of **comment scores** shows a **skewed distribution**, where most comments have a low score, and a small number of comments received higher scores.
  - The peak at low scores suggests that a large number of comments received minimal attention or upvotes, while a few received **high engagement**.

- **Insights**:

  - This indicates that the majority of comments may not be highly upvoted or noticed, which is typical for large online communities where many comments go unnoticed or receive fewer interactions.
  - **High engagement comments** are likely those that spark intense discussions or provide unique insights, contributing to their higher visibility.

## 5.5 Total Comments



- **Key Findings**:

  - The dataset consists of **7,192 total comments**, providing a large and diverse sample for sentiment analysis and visualization.

- **Insights**:

  - The large dataset allows for a comprehensive analysis of sentiment, subreddit activity, and comment engagement, ensuring that the findings reflect a broad spectrum of opinions and interactions within the Reddit movie discussions.

  - The dataset size is large enough to draw meaningful conclusions about the sentiment and popularity of various movies discussed on Reddit.

## 5.6 Summary of Key Findings

The analysis of Reddit comments reveals several interesting insights:

- Sentiment Distribution: The distribution of sentiments is quite balanced, with a slight tilt toward neutral comments. This insinuates Reddit discussions about movies to be usually fact-based or possibly mixed opinion.

- The Top Subreddits: The movies subreddit attracts the maximum activity for movie discussions, with smaller communities such as TrueFilm and CinemaSins also keeping up-booked for more niche topics related to the art of motion pictures.

- Word Cloud: The word cloud brings out that the discussions are dominated by mainstream subreddits, namely movies and BoxOffce, whereas more niche communities like TrueFilm offer deeper film analysis.

- Comment Score Distribution: Given the skew in the distribution of comment scores, it would appear that most comments do not receive much upvoting, with a very few drawing a lot of attention.

- Total Comments: With 7,192 comments, the dataset is well-fit to encapsulate a wide spectrum of views and interactions, thus bestowing statistical significance on the analysis.

# 6.0 Optimization & Comparison

## 6.1 Model Optimization

Improving the performance and usability of the sentiment analysis system required several optimization techniques across both models:

● **Text Preprocessing :**

Reddit comments were cleaned by removing punctuation, emojis, URLs, and stopwords. This step significantly reduced noise in the input data and improved overall model learning quality.

● **Sentiment Labelling with Neutral Class :**

Since initial labelling was done using a Hugging Face transformer model that only predicted positive or negative classes, we introduced a confidence thresholding mechanism to identify neutral comments:

- ● If the highest class probability was below 0.55 (for NB) or 0.70 (for Bi-LSTM) , the comment was relabelled as NEUTRAL .

- ● This approach improved accuracy for ambiguous or low-confidence predictions, especially in short or code-switched Malay–English text.

● **Efficient Feature Extraction :**

- ■ The Naïve Bayes model used TF-IDF vectorization with uni- and bi-grams (top 5,000 features) to represent text numerically.

- ■ The Bi-LSTM model used a Keras Tokenizer (vocabulary size: 20,000 words) to convert text into sequences, padded to a maximum length of 60 tokens.

● **Early Stopping in Deep Learning :**

For the Bi-LSTM model, training included early stopping (patience = 2) on a 10% validation split to prevent overfitting and reduce unnecessary training time.

● **Model Serialization :**

Both the trained models and preprocessing artifacts (vectorizer/tokenizer) were saved using 'joblib' and 'Keras' built-in methods. This allowed quick loading and reusability during real-time classification.

## 6.2 Model Comparison

Two distinct models were trained and evaluated for the sentiment analysis system:

- Multinomial Naïve Bayes (NB)
- Bidirectional LSTM (Bi-LSTM)

Both models were trained using the dataset enriched with neutral labels based on the confidence threshold method described above.

**Evaluation Metrics**

| MODEL | ACCURACY | PRECISION (POS/NEG/NEU) | RECALL (POS/NEG/NEU) | F1-SCORE (POS/NEG/NEU) | NOTES |
|-------|----------|-------------------------|----------------------|------------------------|-------|
| Naïve Bayes | 0.62 | 0.80 / 0.83 / 0.37 | 0.20 / 0.16 / 0.95 | 0.32 / 0.27 / 0.53 | High recall for NEUTRAL, poor for other classes |
| Bi-LSTM | 0.78 (+16%) | 0.65 / 0.63 / 0.53 | 0.63 / 0.52 / 0.64 | 0.64 / 0.57 / 0.58 | Balanced performance, better at handling sarcasm and code-switching |

**Insights :**

- **Bi-LSTM outperformed Naive Bayes** across all major metrics and produced fewer misclassifications.
- **Naive Bayes over-labels neutral**, likely due to TF-IDF limitations and bag-of-words assumptions.
- **Naive Bayes** showed a bias toward **neutral sentiment**, likely due to over-reliance on unigram features.
- **Confidence thresholding** helped reduce incorrect polar predictions in both models, especially for vague or short comments.
- Bi-LSTM offered better tradeoffs in real-world performance, although it required longer training time and GPU acceleration.

# 7.0 Conclusion & Future Work

## 7.1 Conclusion

This project proved an ability of real-time sentiment analysis using an integrated pipeline consisting of Apache Kafka + Apache Spark + Elasticsearch. These contemporary technologies offer efficient collection, processing, and analysis of huge volumes of data, say Reddit comments on movies.

The sentiment analysis pipeline uses Multinomial Naive Bayes and Bidirectional LSTM models to classify comments as positive, negative, or neutral. The results from these sentiment models were visualized and analyzed in Kibana to obtain insights about user opinions of various films discussed across several movie-related subreddits.

In a nutshell, some conclusions drawn from these analyses depict that the sentiment of Reddit comments was more or less fair, with neutral yelled way higher, and r/movies and other subreddits' comments dominated the discourse. The Naïve Bayes model, on the contrary, lean almost entirely toward tagging comments as neutral, while the Bi-LSTM model exhibited a balanced and more nuanced performance, capturing sarcasm and more complex sentiment patterns.

In general, the project has demonstrated the capabilities applied for real-time data ingestion using Apache Kafka and streaming

## 7.2 Future Work

Though the project has been used and run for sentiment analysis on Reddit data, there are numerous improvements and extensions that one could aim for:

**Improvement to the Model:**

Hybrid Models: Merging the Naïve Bayes and Bi-LSTM model strengths could lead to better results in combing the classical and deep learning methods into one framework for predictions.

Fine-tuning: The more extensive hyperparameter optimization of the Bi-LSTM model or employing a pre-trained transformer model (like BERT or GPT) could assist in enhancement of the prediction of sentiment, especially for rather subtle comments.

**Data Expansion:**

The project currently focuses on a limited number of Reddit-based subreddits. Further data collection for a greater number of subreddits, or even comments from another social media platform like Twitter, would allow a more comprehensive analysis.

Introducing data from other languages or code-switched content (like Malay-English) would assist in moving towards the system being robust and adaptable to various language contexts.

**Real-Time Processing Enhancements:**

Currently, the setup allows for the real-time processing using Apache Kafka and Spark, but further means of scaling, such as partition

# 8.0 References

1. Gormley, C., & Tong, Z. (2015). *Elasticsearch: The definitive guide*. O'Reilly Media.
   https://www.elastic.co/guide/en/elasticsearch/
2. Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation, 9*(8), 1735-1780.
   https://doi.org/10.1162/neco.1997.9.8.1735
3. Kreps, J., Narkhede, N., & Rao, J. (2011). Kafka: A distributed messaging system for log processing. *Proceedings of the NetDB, 1*(1), 1-8.
   https://kafka.apache.org/
4. Liu, B. (2012). *Sentiment analysis and opinion mining*. Synthesis Lectures on Human Language Technologies, 5(1), 1-167.
   https://www.amazon.com/Sentiment-Analysis-Opinion-Mining-Liu/dp/160845836X
5. Rennie, J. D. M., Shih, L., Teevan, J., & Karger, D. (2003). Tackling the poor assumptions of naive Bayes text classifiers. In *Proceedings of the Twentieth International Conference on Machine Learning (ICML-03)*, 616-623.
   https://dl.acm.org/doi/10.5555/3041838.3041911
6. Viegas, F. B., Wattenberg, M., & Van Ham, F. (2007). How do we see our data? Visualizing the insights of information. In *Proceedings of the 12th International Conference on Information Visualization*, 12(2), 110-116.
   https://ieeexplore.ieee.org/document/4288244
7. Zaharia, M., Chowdhury, M., Das, T., Dave, A., & Ma, J. (2016). Resilient distributed datasets: A fault-tolerant abstraction for in-memory cluster computing. In *9th USENIX Symposium on Networked Systems Design and Implementation (NSDI)*, 2(3), 1-16.
   https://spark.apache.org/

## 9.0 Appendix

### 9.1 Scraper code (reddit_scrape.py)

```python
import praw
import pandas as pd

# Reddit API credentials
reddit = praw.Reddit(
    client_id="bKUQYyqxMsQwBmECpl8guw",
    client_secret="IIzAvemnMFq22oD3a84-JAqRtveHmQ",
    user_agent="sentiment_project by /u/Legal-Vanilla8068"
)

# List of subreddits
subreddits = ["movies", "TrueFilm", "BoxOffice",
"CinemaSins", "FilmFestivals"]

data = []

for subreddit_name in subreddits:
    subreddit = reddit.subreddit(subreddit_name)
    for submission in subreddit.top(time_filter="year",
limit=20):  # Top posts from last year
        submission.comments.replace_more(limit=0)
        for comment in submission.comments.list():
            data.append({
                "post_id": submission.id,
                "post_title": submission.title,
                "comment_id": comment.id,
                "comment_body": comment.body,
                "comment_score": comment.score,
                "created_utc": comment.created_utc,
                "subreddit": subreddit_name
            })

# Save to CSV
df = pd.DataFrame(data)
df.to_csv("reddit_movies_comments.csv", index=False)
```

```
print(f"Saved {len(df)} comments to
reddit_movies_comments.csv")
```

## 9.2 Cleaner code (CleanerCode.py)

```python
import pandas as pd
import re
import nltk
from nltk.tokenize import RegexpTokenizer
from nltk.corpus import stopwords

# Download required NLTK data (first time only)
nltk.download('stopwords')

# Initialize tokenizer and stopwords
tokenizer = RegexpTokenizer(r'\b\w{3,}\b')  # Keep words with
3+ letters
stop_words = set(stopwords.words('english'))

def clean_text(text):
    if not isinstance(text, str) or text.strip() == '':
        return ""

    # Remove URLs
    text = re.sub(r'https?://\S+|www\.\S+', '', text)

    # Remove emails
    text = re.sub(r'\S+@\S+', '', text)

    # Remove emojis and special characters
    text = re.sub(r'[^\w\s#@_/-]', '', text)

    # Remove numbers
    text = re.sub(r'\d+', '', text)

    # Normalize whitespace and lowercase
```

```python
    text = re.sub(r'\s+', ' ', text).strip().lower()

    # Tokenize
    tokens = tokenizer.tokenize(text)

    # Remove stopwords
    tokens = [word for word in tokens if word not in
stop_words]

    return ' '.join(tokens)

# Convert UTC timestamp to readable format
def utc_to_readable(utc_time):
    try:
        return
datetime.utcfromtimestamp(float(utc_time)).strftime('%Y-%m-%d
%H:%M:%S')
    except:
        return ""

# Load your dataset
print("Loading dataset...")
df = pd.read_csv("reddit_movies_comments_large.csv")

# Convert UTC timestamp
print("Converting timestamps...")
df['timestamp_readable'] =
df['created_utc'].apply(utc_to_readable)

# Clean comment body
print("Cleaning comment bodies...")
df['cleaned_comment'] = df['comment_body'].apply(clean_text)

# Filter out empty comments
df = df[df['cleaned_comment'].str.len() > 5]

# Reorder and select final columns
final_columns = [
    'comment_id', 'post_title', 'comment_body',
```

```
'cleaned_comment',
    'comment_score', 'created_utc', 'timestamp_readable',
'subreddit'
]
df = df[final_columns]

# Save cleaned dataset
output_file = "reddit_cleaned_with_stopwords_removed.csv"
df.to_csv(output_file, index=False)
print(f"✅ Cleaned dataset saved to {output_file}")
```

## 9.3 Docker container (docker-compose.yml)

```yaml
services:
  zookeeper:
    image: confluentinc/cp-zookeeper:7.6.1
    environment:
      ZOOKEEPER_CLIENT_PORT: 2181
      ZOOKEEPER_TICK_TIME: 2000
    ports:
      - "2181:2181"

  kafka:
    image: confluentinc/cp-kafka:7.6.1
    depends_on:
      - zookeeper
    ports:
      - "9092:9092"
      - "29092:29092"
    environment:
      KAFKA_BROKER_ID: 1
      KAFKA_ZOOKEEPER_CONNECT: zookeeper:2181
      KAFKA_ADVERTISED_LISTENERS:
PLAINTEXT://kafka:29092,PLAINTEXT_HOST://localhost:9092
      KAFKA_LISTENER_SECURITY_PROTOCOL_MAP:
PLAINTEXT:PLAINTEXT,PLAINTEXT_HOST:PLAINTEXT
```

```yaml
      KAFKA_LISTENERS:
PLAINTEXT://0.0.0.0:29092,PLAINTEXT_HOST://0.0.0.0:9092
      KAFKA_OFFSETS_TOPIC_REPLICATION_FACTOR: 1
      KAFKA_GROUP_INITIAL_REBALANCE_DELAY_MS: 0

  elasticsearch:
    image: docker.elastic.co/elasticsearch/elasticsearch:8.13.4
    environment:
      - discovery.type=single-node
      - xpack.security.enabled=false
      - bootstrap.memory_lock=true
      - "ES_JAVA_OPTS=-Xms1g -Xmx1g"
    ulimits:
      memlock:
        soft: -1
        hard: -1
    ports:
      - "9200:9200"
      - "9300:9300"

  kibana:
    image: docker.elastic.co/kibana/kibana:8.13.4
    depends_on:
      - elasticsearch
    ports:
      - "5601:5601"
    environment:
      - ELASTICSEARCH_HOSTS=http://elasticsearch:9200

  spark:
    image: bitnami/spark:3.4.2
    depends_on:
      - kafka
      - elasticsearch
    ports:
      - "8080:8080"  # Spark Master UI
      - "7077:7077"  # Spark Master Port
    volumes:
      - ./:/opt/bitnami/spark/work # Maps your project directory into the
container
```

```
    environment:
      SPARK_MODE: master
      SPARK_MASTER_URL: spark://spark:7077


  # --- ADD THIS NEW SERVICE ---
  spark-worker:
    image: bitnami/spark:3.4.2
    depends_on:
      - spark # Ensures the master is started first
    volumes:
      - ./:/opt/bitnami/spark/work
    environment:
      SPARK_MODE: worker
      SPARK_MASTER_URL: spark://spark:7077 # Tells the worker where to
find the master
      SPARK_WORKER_CORES: 2  # Number of cores the worker can use
      SPARK_WORKER_MEMORY: 2g # Amount of memory the worker can use
```

## 9.4 Kafka (kafka_producer.py)

```python
import os
import time
import json
import praw
from kafka import KafkaProducer
from dotenv import load_dotenv


# Load environment variables from .env file
load_dotenv()


# --- Kafka Configuration ---
KAFKA_TOPIC = 'reddit-comments'
KAFKA_BOOTSTRAP_SERVERS = 'localhost:9092'


# --- Reddit API Configuration ---
SUBREDDIT_NAME = 'movies'  # The subreddit you want to stream


def create_kafka_producer():
    """Creates a Kafka producer instance."""
```

```python
    try:
        producer = KafkaProducer(
            bootstrap_servers=KAFKA_BOOTSTRAP_SERVERS,
            value_serializer=lambda v: json.dumps(v).encode('utf-8')
        )
        print("✅ Kafka Producer connected successfully.")
        return producer
    except Exception as e:
        print(f"❌ Could not connect to Kafka: {e}")
        return None


def connect_to_reddit():
    """Connects to the Reddit API using credentials from environment
variables."""
    try:
        reddit = praw.Reddit(
            client_id=os.getenv("REDDIT_CLIENT_ID"),
            client_secret=os.getenv("REDDIT_CLIENT_SECRET"),
            user_agent=os.getenv("REDDIT_USER_AGENT")
        )
        print("✅ Reddit API connected successfully.")
        return reddit
    except Exception as e:
        print(f"❌ Could not connect to Reddit: {e}")
        return None


def stream_comments(producer, reddit):
    """Streams comments from the specified subreddit and sends them to
Kafka."""
    if not producer or not reddit:
        return

    subreddit = reddit.subreddit(SUBREDDIT_NAME)
    print(f"📡 Streaming comments from r/{SUBREDDIT_NAME}...")

    try:
        for comment in subreddit.stream.comments(skip_existing=True):
            try:
                # Construct the message payload
                message = {
```

```python
                    "post_id": comment.submission.id,
                    "post_title": comment.submission.title,
                    "comment_id": comment.id,
                    "comment_body": comment.body,
                    "comment_score": comment.score,
                    "created_utc": comment.created_utc
                }

                # Send the message to our Kafka topic
                producer.send(KAFKA_TOPIC, message)
                print(f"📤 Sent comment ID {comment.id}:
{comment.body[:50]}...")

                time.sleep(1)  # Throttle to avoid overwhelming Kafka

            except Exception as e:
                print(f"⚠️ Error processing comment {comment.id}: {e}")
                continue # Skip to the next comment

    except Exception as e:
        print(f"❌ An error occurred during streaming: {e}")

if __name__ == "__main__":
    kafka_producer = create_kafka_producer()
    reddit_client = connect_to_reddit()

    if kafka_producer and reddit_client:
        stream_comments(kafka_producer, reddit_client)
```

## 9.5 Spark (spark_streaming.py)

```python
import os
os.environ['TF_CPP_MIN_LOG_LEVEL'] = '2' # Suppress TensorFlow INFO and
WARNING messages

from pyspark.sql import SparkSession
from pyspark.sql.functions import udf, col, from_json
```

```python
from pyspark.sql.types import StringType, StructType, StructField,
IntegerType, FloatType
import re
import joblib
import numpy as np
import tensorflow as tf
from tensorflow.keras.preprocessing.sequence import pad_sequences

# --- Configuration ---
KAFKA_TOPIC = "reddit-comments"
KAFKA_SERVER = "kafka:29092"
ELASTICSEARCH_NODE = "elasticsearch"
ELASTICSEARCH_PORT = "9200"
ELASTICSEARCH_INDEX = "reddit-comments-lstm"
CHECKPOINT_LOCATION = "/tmp/spark_checkpoint_reddit_lstm"

# Define a maximum sequence length for padding.
# This should match the length used during model training.
MAX_SEQUENCE_LENGTH = 150

# --- Model Loading and Broadcasting ---
# Load the tokenizer and the Keras model on the driver
tokenizer = joblib.load("/opt/bitnami/spark/work/models/tokenizer.pkl")
model =
tf.keras.models.load_model("/opt/bitnami/spark/work/models/lstm_sentiment_
model.h5")

# Initialize Spark Session
spark = SparkSession.builder \
    .appName("RedditSentimentLSTM") \
    .config("spark.jars.packages",
            "org.apache.spark:spark-sql-kafka-0-10_2.12:3.4.2,"
            "org.elasticsearch:elasticsearch-spark-30_2.12:8.13.4") \
    .getOrCreate()

spark.sparkContext.setLogLevel("ERROR")

# Broadcast the tokenizer and the model's weights to all worker nodes
tokenizer_broadcast = spark.sparkContext.broadcast(tokenizer)
```

```python
model_weights_broadcast =
spark.sparkContext.broadcast(model.get_weights())


# --- Schema Definition for Kafka Messages ---
schema = StructType([
    StructField("post_id", StringType(), True),
    StructField("post_title", StringType(), True),
    StructField("comment_id", StringType(), True),
    StructField("comment_body", StringType(), True),
    StructField("comment_score", IntegerType(), True),
    StructField("created_utc", FloatType(), True)
])


# --- Text Cleaning Function ---
def clean_text(text):
    if text is None:
        return ""
    text = str(text).lower()
    text = re.sub(r"https[^\s]+|www[^\s]+", '', text)  # Remove URLs
    text = re.sub(r'[^a-zA-Z0-9\s]', '', text)         # Remove special
characters/emojis
    text = re.sub(r'\s+', ' ', text).strip()           # Normalize
whitespace
    return text


clean_text_udf = udf(clean_text, StringType())



# --- Sentiment Prediction UDF for LSTM Model ---
def predict_sentiment_lstm(text_series):
    # This function is designed to work with Pandas UDFs for efficiency
    # but a standard UDF is used here for simplicity.

    # 1. Get the broadcasted tokenizer and weights
    local_tokenizer = tokenizer_broadcast.value
    local_weights = model_weights_broadcast.value

    # 2. Re-create the model architecture and set the broadcasted weights
    #    This avoids serializing the entire model object.
```

```python
        local_model =
tf.keras.models.load_model("/opt/bitnami/spark/work/models/lstm_sentiment_
model.h5")
        local_model.set_weights(local_weights)

    predictions = []
    for text in text_series:
        if not text or len(text.strip()) < 2:
            predictions.append("NEUTRAL")
            continue
        try:
            # 3. Tokenize and pad the text
            sequence = local_tokenizer.texts_to_sequences([text])
            padded_sequence = pad_sequences(sequence,
maxlen=MAX_SEQUENCE_LENGTH)

            # 4. Make a prediction
            prediction_prob = local_model.predict(padded_sequence,
verbose=0)[0][0]

            # 5. Interpret the prediction
            if prediction_prob > 0.5:
                predictions.append("POSITIVE")
            else:
                predictions.append("NEGATIVE")
        except Exception as e:
            print(f"Prediction error: {e}")
            predictions.append("ERROR")

    return predictions

# We need to wrap the function slightly for a standard UDF
def predict_sentiment(text):
    return predict_sentiment_lstm([text])[0]

predict_udf = udf(predict_sentiment, StringType())



# --- Spark Streaming Pipeline ---
print("🚀 Reading from Kafka...")
```

```python
# 1. Read from Kafka
raw_df = spark.readStream \
    .format("kafka") \
    .option("kafka.bootstrap.servers", KAFKA_SERVER) \
    .option("subscribe", KAFKA_TOPIC) \
    .option("startingOffsets", "latest") \
    .load()

# 2. Decode the Kafka message and parse the JSON
json_df = raw_df.selectExpr("CAST(value AS STRING)")
parsed_df = json_df.select(from_json(col("value"),
schema).alias("data")).select("data.*")

# Add this to see the raw JSON from Kafka in your console
json_df.writeStream \
    .outputMode("append") \
    .format("console") \
    .start()

parsed_df = json_df.select(from_json(col("value"),
schema).alias("data")).select("data.*")

# 3. Clean the comment text
print("🚀 Cleaning comments...")
clean_df = parsed_df.withColumn("clean_comment",
clean_text_udf(col("comment_body")))
clean_df.writeStream \
    .outputMode("append") \
    .format("console") \
    .start()

print("🚀 Streaming to console started...")

# 4. Predict sentiment using the LSTM model
print("🚀 Predicting sentiment...")
result_df = clean_df.withColumn("sentiment",
predict_udf(col("clean_comment")))

print("🚀 Starting Spark writeStream to Elasticsearch...")
# --- Write to Elasticsearch ---
```

```python
es_query = result_df.writeStream \
    .outputMode("append") \
    .format("org.elasticsearch.spark.sql") \
    .option("es.nodes", ELASTICSEARCH_NODE) \
    .option("es.port", ELASTICSEARCH_PORT) \
    .option("checkpointLocation", CHECKPOINT_LOCATION) \
    .option("es.resource", ELASTICSEARCH_INDEX) \
    .start()

print("✅ writeStream started. Awaiting termination...")

es_query.awaitTermination()
```