# ANALYING SOCIAL MEDIA SHARING PATTERNS IDENTIFY FACTORS DRINVING VIRAL VIDEO PROPAGATION

## A MINI PROJECT REPORT

*Submitted by*

## VENKATA SAI.V(221801060)

## PRAVEEN.B(221801503)

*in partial fulfillment for the award of the degree of*

*Of*

## BACHELOR OF TECHNOLOGY
## IN
## ARTIFICIAL INTELLIGENCE AND DATA SCIENCE



# RAJALAKSHMI ENGINEERING COLLEGE

# ANNA UNIVERSITY, CHENNAI

## NOVEMBER 2024

# BONAFIDE CERTIFICATE

Certified that this Report titled **"Analyzing Social Media Sharing Patterns to Identify Factors Driving Viral Video Propagation"** is the Bonafide work of **VENKATA SAI V (2116221801060), PRAVEEN B (2116221801503)** who carried out the work under my supervision.

SIGNATURE

SIGNATURE

**Dr. J. M. GNANASEKAR, M.E., Ph.D.,**
Professor and Head,
Department of AI&DS
Rajalakshmi Engineering College,
Chennai-602 105.

**Mrs.Y.Nirmala Anandhi,M.E.,**
Assistant Professor(SS),
Deparment of AI&DS,
Rajalakshmi Engineering College,
Chennai-602 105.

Submitted for the project viva-voce examination held on…………………………..

INTERNAL EXAMINER

EXTERNAL EXAMINER

# ACKNOWLEDGEMENT

Initially we thank the Almighty for being with us through every walk of our life and showering his blessings through the endeavor to put forth this report. Our sincere thanks to our Chairman **Mr. S. MEGANATHAN, B.E, F.I.E.,** and our respected Chairperson **Dr. (Mrs.) THANGAMMEGANATHAN, Ph.D.,** and our Vice Chairman **Mr. ABHAY SHANKAR MEGANATHAN, B.E.,M.S**., for providing us with the requisite infrastructure andsincere endeavoring in educating us in their premier institution.

Our sincere thanks to **Dr. S.N. MURUGESAN, M.E., Ph.D.,** our beloved Principal for his kind support and facilities provided to complete our work in time. We express our sincere thanks to **Dr. J. M. GNANASEKAR., M.E., Ph.D.,** Professor and Head of the Department of Artificial Intelligence and Data Science for his guidance and encouragement throughout the project work. We are glad to express our sincere thanks to our supervisor **Mrs.Y.NIRAMALA ANANDHI,M.E., Assistant Professor(SS)**,Department of Artificial Intelligence & Data Science, Rajalakshmi Engineering College and coordinator, **Dr. P. INDIRA PRIYA, M.E., Ph.D.,** Professor, Department of ArtificialIntelligence and Data Science, Rajalakshmi Engineering College for their valuable guidance throughout the course of the project.

Finally we express our thanks for all teaching, non-teaching, faculty and our parents for helping us with the necessary guidance during the time of our project.

# ABSTRACT

In the digital media, viral videos play a pivotal role in shaping audience engagement and extending content reach. This project focuses on analyzing social media interactions and video metrics to identify factors that contribute to the propagation of viral videos. By leveraging machine learning techniques, including logistic regression and neural networks, the system processes data from YouTube videos, extracting features such as views, likes, dislikes, and comments to predict virality. It includes an interactive web interface that allows users to input YouTube video URLs, analyze engagement metrics, and receive insights into key factors influencing virality. Additionally, personalized recommendations are provided to enhance video performance. The proposed system incorporates semantic and visual analysis of video content, community detection, and advanced engagement metrics to gain a holistic understanding of virality drivers. Results demonstrate the efficacy of the models in predicting video virality with high accuracy, offering actionable insights for content creators and marketers to optimize digital strategies. This study presents a scalable and impactful approach to understanding and maximizing content propagation in the digital landscape.

# TABLE OF CONTENTS

# LIST OF FIGURES

# CHAPTER 1

# INTRODUCTION

## 1.1 GENERAL

The rapid growth of digital media and social networking platforms has revolutionized how information is shared and consumed. Among the diverse forms of online content, videos have emerged as a dominant medium for engagement and communication. Viral videos, in particular, have the power to reach millions of viewers in a short span of time, influencing trends, opinions, and even purchasing decisions.

Understanding what makes a video go viral is crucial for digital marketers, content creators, and businesses aiming to maximize audience engagement and enhance their online presence. Virality is often driven by a complex interplay of factors, including user interactions, content characteristics, and the structure of social networks. By analyzing these factors, it becomes possible to derive actionable insights that can guide content creation and dissemination strategies.

This project delves into the dynamics of viral video propagation by leveraging advanced machine learning techniques and data analysis. It aims to identify the key drivers of video virality by studying metrics such as views, likes, comments, and network behaviors. The findings of this study are intended to optimize content strategies and provide a deeper understanding of how digital content spreads in an interconnected world.

## 1.2 NEED FOR THE STUDY

The digital landscape is saturated with content, making it increasingly challenging for creators and businesses to stand out and capture audience attention. In this highly competitive environment, understanding the mechanisms behind viral video propagation has become essential. Viral videos not only generate significant engagement but also amplify brand visibility, shape public opinion, and drive consumer actions.

Traditional approaches to content strategy often rely on intuition or generalized guidelines, which fail to address the dynamic and multifaceted nature of virality. These methods do not account for factors like user behavior, network structure, and the emotional resonance of content, which significantly influence how and why videos spread.

This study addresses the critical gap by providing a data-driven framework to analyze and predict virality. By identifying the key metrics and factors that contribute to video virality, the study aims to empower content creators, marketers, and businesses with actionable insights to optimize their digital strategies. Additionally, the integration of machine learning models enhances the precision of predictions, ensuring that the results are both scalable and adaptable to evolving trends.

The insights gained from this study not only contribute to academic research but also offer practical applications for enhancing content reach, targeting key influencers, and maximizing user engagement in the fast-paced digital ecosystem.

## 1.3 OBJECTIVES OF THE STUDY

The primary objective of this study is to analyze and predict the factors that drive viral video propagation on social media platforms. By leveraging advanced data analytics and machine learning models, the study seeks to provide actionable insights into video virality. The specific objectives are:

1. To Analyze Viral Video Metrics:

    Investigate key metrics such as views, likes, dislikes, comments, and share rates to identify patterns and trends that contribute to the virality of videos.

2. To Examine User Interactions and Network Structures:

    Understand the role of user engagement (e.g., comments, shares, likes) and the architecture of social networks in amplifying content reach.

3. To Develop Predictive Models:

    Employ machine learning algorithms, including logistic regression and neural networks, to predict the virality of videos based on content and engagement metrics.

4. To Identify Key Drivers of Virality:

    Assess qualitative and quantitative factors, such as emotional appeal, topic relevance, and community dynamics, that significantly influence a video's spread.

5. To Provide Actionable Insights and Recommendations:

    Design an interactive system that offers users tailored recommendations for optimizing video content to maximize reach and engagement.

## 1.4 OVERVIEW OF THE PROJECT

The system is designed to leverage advanced machine learning techniques, including logistic regression and neural networks, to process and analyze video data from platforms such as YouTube. Key features include:

1. Video Data Extraction and Preprocessing:
   o Metrics such as views, likes, dislikes, and comments are extracted and normalized for analysis.
   o Feature engineering techniques are employed to derive meaningful insights from raw data.

2. Virality Prediction and Evaluation:
   o Machine learning models are trained and evaluated to predict a video's potential to go viral.
   o Metrics such as accuracy, precision, recall, and F1-score are used to assess model performance.

3. Interactive User Interface:
   o A web-based system allows users to input YouTube video URLs and receive real-time predictions and insights.
   o Recommendations for optimizing video content are generated based on the analysis.

4. Key Insights and Recommendations:
   o The system identifies metrics and content elements that significantly influence virality.
   o Users are provided with actionable suggestions to improve video performance and maximize engagement.

**Technical Components:**

Data Extraction and Preprocessing:

         The system extracts video data via the YouTube API, cleans it, creates features, and scales them for model training.

Machine Learning Models:

         Logistic regression and neural networks predict virality, evaluated with accuracy, precision, recall, and F1-score.

User Interface:

         A Flask-based web interface allows video URL input and displays predictions with visual insights.

Virality Analysis and Insights Generation:

         NLP analyzes comments, computer vision evaluates visual content, and engagement metrics assess virality.

Scalability and Future Enhancements:

         Cloud integration supports larger datasets, and models are periodically updated for improved accuracy.

# CHAPTER 2

# LITERATURE REVIEW

The study of viral video propagation has been enriched by various approaches leveraging social media analytics, network analysis, and machine learning models. Social media platforms such as Hootsuite and BuzzSumo provide essential metrics like likes, shares, and comments, which aid in tracking video performance. However, as Chen et al. (2020) noted, these tools often fail to capture qualitative factors like emotional resonance and contextual nuances. Network analysis further enhances the understanding of content dissemination, with Thompson and Clark (2021) emphasizing the role of influencers and community dynamics in amplifying reach. Despite their insights, interpreting complex network data remains a challenge. Machine learning has emerged as a cornerstone for predicting video virality, with Liu and Kumar (2020) demonstrating the effectiveness of random forest algorithms and Williams and Brown (2021) highlighting the success of logistic regression and neural networks in integrating diverse engagement metrics for more accurate predictions. Nevertheless, as Aphale (2020) pointed out, many existing systems struggle with scalability and adapting to rapidly evolving trends. These findings underscore the need for a holistic framework that integrates advanced analytics and predictive modeling to address the multifaceted nature of viral video propagation.

| S. No | Author Name | Paper Title | Description | Journal | Volume/ Year |
|---|---|---|---|---|---|
| 1 | Waafa Hasan Alwan | Identifyin gInfluenti al Users on Instagram Through Visual Content Analysis | This paper analyzes visual content on Instagram to identify influential users, focusing on how image aesthetics impact engagement metrics like likes and comments. | IEEE Access | 2020 |
| 2 | Li, Meng et al. | Dissecting Virality: A Comprehensiv e Study of Factors Influencing Video Popularity on YouTube | Identifies key factors contributing to video virality, such as content quality, audience engagement, and social network effects. | ACM Transactions on Multimedia | 2021 |
| 3 | Zhang, Yuxiao et al. | The Role of Emotional Arousal in Video Virality: A Large-Scale Empirical Study | Examines the impact of emotional arousal on video sharing and engagement. | Journal of Management Information Systems | 2022 |

# CHAPTER 3

# SYSTEM OVERVIEW

## 3.1 EXISTING SYSTEM

The existing systems for analyzing viral video propagation typically rely on traditional social media analytics platforms and tools, which focus on tracking engagement metrics like views, likes, shares, and comments. These systems offer valuable insights into user interaction but often lack the ability to fully capture the complex factors that contribute to a video's virality. Platforms like Sprout Social and Hootsuite provide basic engagement analytics but do not incorporate deeper analyses of content quality, emotional triggers, or network dynamics. Content tracking systems such as BuzzSumo provide insights into the spread of content across various platforms, allowing users to track which videos are performing well. However, these tools primarily focus on popularity and fail to analyze the qualitative aspects of content that influence virality.

Another major limitation of existing systems is their reliance on influencer-based models, which can overlook the broader, grassroots sharing patterns that are often responsible for viral success. Influencer analytics tools, such as Aspire, identify key influencers who amplify content reach but do not address the role of smaller, everyday users in driving viral trends. Additionally, network analysis tools like Gephi help visualize social network structures, which can reveal how videos spread, but they often require specialized knowledge to interpret effectively. Moreover, the existing systems are typically rigid and do not adapt well to the dynamic and rapidly changing nature of social media trends, making them less scalable and adaptable. These limitations point to the need for a more integrated and comprehensive system that combines quantitative metrics with qualitative content analysis to better predict and understand viral video dynamics.

## 3.2 PROPSED SYSTEM

The proposed system aims to overcome the limitations of existing tools by incorporating a multi-faceted approach that integrates social media analytics, machine learning models, semantic and visual analysis, and network detection to predict viral video propagation. The system will extract video metrics, such as views, likes, shares, and comments, from platforms like YouTube using an API. These metrics will be combined with advanced engagement measures, such as sentiment analysis and user interactions, to build a more complete picture of virality.

In addition to traditional engagement metrics, the proposed system will utilize machine learning algorithms, including logistic regression and neural networks, to analyze and predict the potential virality of videos. By training these models on large datasets of historical video metrics, the system can identify patterns and trends that drive viral success. The system will also incorporate semantic analysis through natural language processing (NLP) techniques, which will allow for the examination of video comments, descriptions, and tags to understand the sentiment and emotional appeal of the content. Furthermore, visual analysis using computer vision algorithms will evaluate visual elements of the video, such as colors, composition, and objects, which contribute to its attractiveness and shareability.

An interactive user interface will allow content creators, marketers, and other users to input   YouTube video URLs, retrieve relevant metrics, and receive predictions about the video's likelihood of going viral. The system will also provide personalized recommendations for enhancing video performance based on its analysis of engagement metrics, content features, and network behavior. This holistic approach ensures a more comprehensive understanding of the factors influencing virality, which can help users optimize their content strategies and maximize reach.

## 3.3 FEASIBILITY STUDY

The feasibility study evaluates the practicality of implementing the proposed based on its technical, operational, and economic aspects.

**TechnicalFeasibility:**

The system will be built using widely available technologies that are both reliable and efficient. The backend will be powered by Python, using popular libraries such as Scikit-learn for machine learning, TensorFlow/Keras for deep learning models, and NLTK for natural language processing. For video data extraction, the YouTube API will be used to collect metrics such as views, likes, and comments. The system's architecture will allow easy integration with cloud services for scalable data storage and processing, ensuring that the system can handle large volumes of video data. The frontend of the system will be developed using Flask, a lightweight Python web framework, which will provide a simple yet powerful interface for user interaction. This combination of technologies ensures that the system will be both scalable and maintainable.

**OperationalFeasibility:**

The system's design is user-friendly and can be operated by both technical and non-technical users. The interactive web interface will allow users to input video URLs, track engagement metrics, and receive predictions with minimal effort. Additionally, the integration of machine learning models ensures that the system can provide accurate and actionable insights in real-time. From an operational perspective, the system's modular design allows for easy updates and enhancements, ensuring its long-term usability and relevance in the fast-changing landscape of social media.

**EconomicFeasibility:**

The proposed system offers a cost-effective solution by utilizing open-source tools and libraries, which reduces the overall cost of development. The use of cloud infrastructure will also enable the system to scale as needed without significant upfront investments. By automating the process of viral video analysis and providing actionable insights, the system has the potential to reduce operational costs for content creators and digital marketers.

# CHAPTER 4

## SYSTEM REQUIREMENTS

### 4.1 SOFTWARE REQUIREMENT

**1. Operating System:**

Windows 10/11

macOS 10.15 or later

Linux distributions with kernel version 4.15 or later

**2. Programming Languages:**

Python 3.8 or later: Python serves as the primary programming language for implementing the backend, machine learning models, and video analysis functionalities. Its extensive library support and ease of use make it ideal for this project.

**3. Frameworks and Libraries**

**Backend Framework:**

**Machine Learning Libraries:**

- Scikit-learn: For implementing logistic regression models and evaluating performance metrics such as precision and recall.
- TensorFlow/Keras: For building and training neural network models to predict video virality.

**Natural Language Processing Libraries:**

- NLTK (Natural Language Toolkit): Used for sentiment analysis and processing text data extracted from video comments.

**Computer Vision Libraries:**

- OpenCV: For visual analysis of video frames to detect features like color schemes and object presence.

**Data Processing Libraries:**

- Pandas: For efficient data manipulation and analysis.
- NumPy: For numerical computations and matrix operations.

**Visualization Libraries:**

- Matplotlib/Seaborn: For creating insightful visualizations of data trends

# CHAPTER 5

## SYSTEM DESIGN

### 5.1 SYSTEM ARCHITECTURE



Fig 1 System Architecture

The system architecture for predicting viral video propagation consists of a modular framework that integrates key components: Data Collection, Data Handling, Data Preparation, Machine Learning, Model Evaluation, and Insights Visualization. Data is initially collected from platforms like YouTube using APIs, capturing video metrics such as views, likes, and comments. This raw data undergoes preprocessing in the Data Handling module, including cleaning, feature engineering, and structuring to ensure quality and relevance. In the Data Preparation stage, the data is scaled and split into training, testing, and validation sets to prepare it for machine learning models.

The Machine Learning module applies Logistic Regression and Neural Network models to analyze and predict video virality, leveraging historical video data to uncover trends and patterns. The Model Evaluation module validates these predictions using metrics like accuracy and F1 scores to ensure reliability. Finally, the system visualizes key insights, offering actionable recommendations for enhancing content performance. This architecture is designed to handle complex video analytics efficiently while delivering accurate and meaningful predictions.

**5.2 MODULE DESCRIPTION**

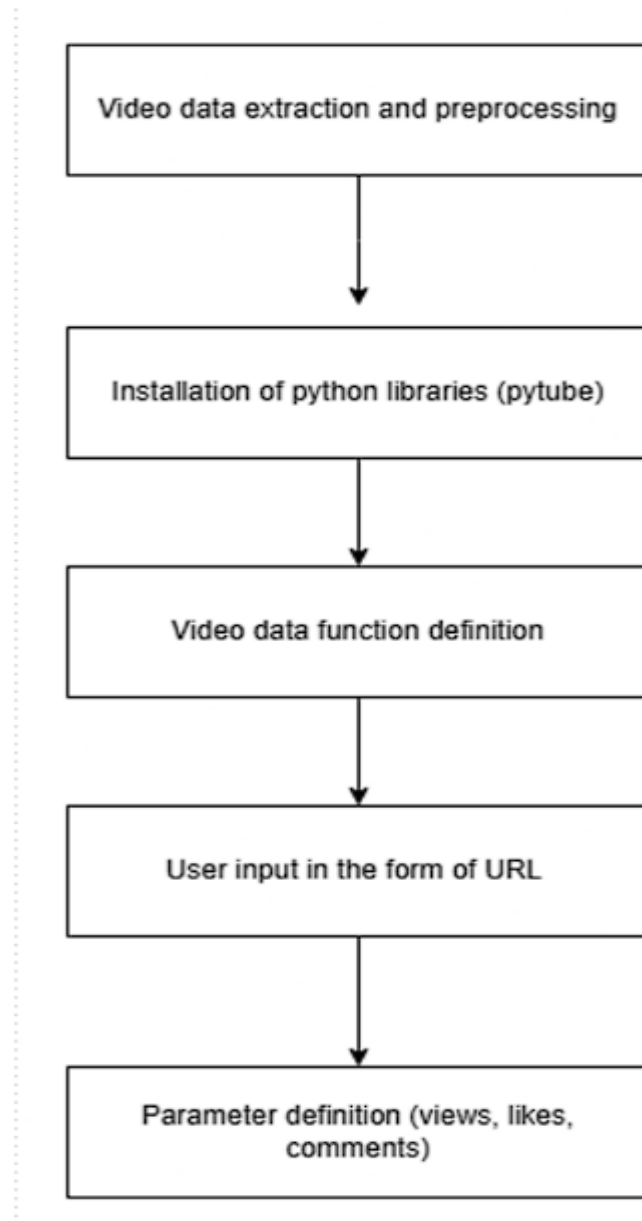**5.2.1 DATA COLLECTION AND PREPROCESSING MODULE:**



Fig5.2.1 Data collection and Pre-processing

**Data Collection:**

This module is crucial for setting the stage in identifying the factors driving viral video propagation. It focuses on extracting, cleaning, and transforming raw video data into a format suitable for analysis. The process involves several essential steps:

1. **Data Extraction**:
   - In this step, video data is gathered from various social media platforms. The focus is on extracting metrics such as views, likes, dislikes, comments, shares, and video duration. These metrics provide a direct measure of user interaction and engagement with the content.
   - Other relevant data points, such as upload time and metadata (e.g., video titles, descriptions, hashtags), are also captured, as they can influence video virality.
   - Tools like the YouTube API, PyTube, or other social media APIs are typically used to automate the extraction of this data. These tools help to fetch a large volume of video data in real time, saving time and effort in the process.
   - The data collected forms the foundation of the analysis and will be used to identify trends or patterns that contribute to a video's viral success.

2. **Feature Engineering**:
   - Once the raw data is collected, additional features are created to provide deeper insights into the video content. This step involves calculating engagement rates, which are important indicators of how actively users interact with the video. Engagement rates are often calculated as the ratio of likes, shares, and comments relative to views.
   - Sentiment analysis is applied to user comments to gauge how positive or negative the audience's response is. Positive sentiment is often linked to higher engagement, which can be a factor in virality.
   - The frequency of specific keywords in video titles, descriptions, and hashtags is also analyzed. Certain keywords or phrases may correlate with higher engagement and wider reach. This information can help in understanding what

drives user interest and sharing.

- o Feature engineering helps to transform raw metrics into actionable insights that will improve the accuracy and depth of the analysis.

3. **Data Cleaning**:

- o Raw data often contains irrelevant, incomplete, or erroneous entries that could skew the results. Therefore, the data cleaning process is essential to ensure the dataset is accurate and reliable.

- o Duplicate entries are removed to prevent overrepresentation of data points that could distort analysis results.

- o Irrelevant data, such as videos with extremely low engagement or those outside the scope of the analysis, are filtered out. This helps in maintaining the focus on videos that are more likely to exhibit viral characteristics.

- o Missing data is addressed through various techniques. Depending on the situation, missing values might be imputed based on other available data or removed entirely if they are unlikely to have a significant impact. Ensuring completeness of data improves the accuracy of predictions later in the process.

4. **Data Transformation**:

- o Data transformation is the final step of preprocessing, where the dataset is converted into a format suitable for further analysis. Numerical features are scaled, so they fit within a standardized range, ensuring that variables of different magnitudes do not dominate the analysis.

- o Categorical variables, such as video categories or specific labels, are encoded into numerical formats, making them compatible with machine learning models.

- o Additionally, the dataset is organized and structured into a format that is ready for modeling and analysis. The transformation process ensures that the data can be efficiently fed into machine learning algorithms or statistical models for the next steps in the project.

**5.2.2 USER VIDEO UPLOAD MODULE**



Fig 5.2.2. user video upload module

   The **User Video Upload Module** is designed to allow users to input a YouTube video link, validate it, and fetch relevant metrics to provide insights and recommendations. This module plays a crucial role in user interaction with the platform and feeds data into the analysis pipeline. Here's how it works:

1. **User Input**:
   - The user provides a YouTube video URL through the website interface.
   - The input field may include a validation process to ensure that the URL is in the correct format (e.g., a valid YouTube link).

2. **Video URL Validation**:
   - **Format Check**: The system checks whether the provided URL matches the structure of a valid YouTube video URL (e.g., https://www.youtube.com/watch?v=videoID).
   - **Video Existence Check**: Using tools like the YouTube API or PyTube, the module checks whether the provided video exists on YouTube. If the video is unavailable

18

or removed, the user will be notified to provide a valid video.

3. **Fetch Video Metrics**:

   o Once the video URL is validated, the module extracts relevant data using the YouTube API or an alternative tool like PyTube.

   o The extracted metrics include:

      ▪ **Basic Info**: Video title, description, upload time, and channel details.

      ▪ **Engagement Metrics**: Views, likes, dislikes, comments, and shares.

      ▪ **Video Metadata**: Tags, categories, and other metadata embedded in the video.

      ▪ **Audience Interaction**: Comment sentiment analysis and engagement trends based on user interactions.

4. **Data Preprocessing**:

   o **Cleaning**: The raw metrics are cleaned, removing any irrelevant data or duplicates that could affect the analysis.

   o **Feature Engineering**: Additional features are derived, such as engagement rate (likes + comments + shares / views) and sentiment analysis scores from the comments.

   o **Normalization**: Numerical features like views, likes, and comments are normalized to ensure that they are comparable across different videos.

5. **Data Storage**:

   o The cleaned and processed video data is stored in a database (e.g., SQL or NoSQL), allowing easy retrieval for analysis and predictions.

   o Each video entry is tagged with unique identifiers, allowing users to track and revisit videos they have uploaded.

6. **User Feedback**:

   o After uploading and processing the video, the user is provided with a confirmation message or feedback, such as:

      ▪ "Video successfully uploaded and analyzed."

      ▪ A message indicating that the video metrics have been extracted and are ready

for analysis.

- o The system may also display basic video information (e.g., title, views, likes, comments) for the user's reference.

7. **Error Handling**:
   - o If there are issues with the URL (e.g., incorrect format or unavailable video), users are prompted with clear, helpful messages to correct their input.
   - o Common error messages might include "Invalid YouTube URL," "Video not found," or "Please try again later."

8. **Integration with Analysis**:
   - o Once the video data is uploaded, validated, and processed, it is passed to the next module (Video Data Extraction and Preprocessing Module) for detailed analysis and prediction.
   - o The analysis results, including insights on factors contributing to virality and recommendations for improving video performance, are then presented to the user.

### 5.2.3 RISK ASSESSMENT MODULE:

```
┌─────────────────────────┐
│    Virality analysis    │
│      0 for not viral    │
│       1 for viral       │
└─────────────────────────┘
            │
            ▼
┌─────────────────────────┐
│    Generate Insights:   │
│          likes          │
│         dislikes        │
│          views          │
│      comments count     │
└─────────────────────────┘
            │
            ▼
┌─────────────────────────┐
│ Recommendations generation │
└─────────────────────────┘
            │
            ▼
┌─────────────────────────┐
│ Recommendations generation │
└─────────────────────────┘
```
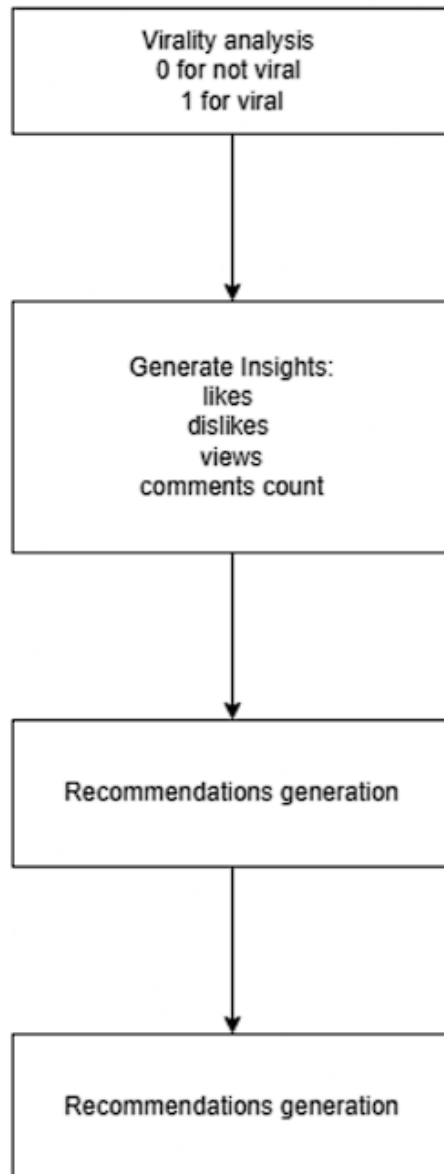
Fig 5.2.3.Risk assessement Module

The Insights Generation and Recommendations Module plays a crucial role in delivering actionable insights and suggestions based on the results from the Virality Prediction and Model Evaluation Module. It helps content creators, marketers, or users understand the factors contributing to the success (or failure) of their videos and provides recommendations to improve video performance and increase the likelihood of virality.

Here's a detailed breakdown of how this module works:

1. **Input Data**:
   - The module receives processed video data and predictions from the Virality Prediction and Model Evaluation Module. This includes the predicted virality score or classification, engagement metrics (views, likes, comments, shares), and any additional insights about the video's performance.

2. **Virality Prediction Interpretation**:
   - **Prediction Breakdown**: The module begins by interpreting the virality prediction, which could be a binary classification (viral or non-viral) or a probability score (likelihood of going viral).
   - The system contextualizes the prediction based on the video's unique metrics (e.g., high engagement but low share rate).

3. **Feature Influence Analysis**:
   - **Identify Key Drivers**: The module analyzes the most influential features that contributed to the virality prediction. For example:
     - **Engagement Metrics**: If the video has high engagement (likes, comments), this could be a sign that it has viral potential.
     - **Sentiment**: Positive sentiment in user comments may indicate higher chances of virality.
     - **Video Content**: Keywords, tags, and description analysis can reveal whether the content aligns with trending topics or highly searched terms.

4. **Insights Generation**:
   - Based on the analysis of the input data and prediction, the system generates key insights. These insights focus on answering important questions for the user, such as:
     - **What is driving the video's current performance?** (e.g., high engagement but limited reach).
     - **What are the key strengths of the video?** (e.g., good audience engagement, high sentiment).

- **What might be limiting its potential?** (e.g., poor title optimization, insufficient shares, or negative sentiment).
- These insights are presented in a digestible, easy-to-understand format to guide content creators in optimizing their strategy.

5. **Actionable Recommendations**:
   - **Engagement Strategies**: If the video has low engagement rates, the system might recommend:
     - Increasing interaction with the audience by responding to comments.
     - Encouraging viewers to share or like the video through calls-to-action (CTAs) in the video or description.
   - **Optimization Suggestions**: If the video lacks visibility despite good engagement, recommendations could include:
     - Improving the title, tags, and description for better search visibility (SEO).
     - Posting at times when the target audience is most active.
   - **Content Improvement**: Based on sentiment analysis, the system could suggest:
     - Enhancing the emotional appeal of the content if negative sentiment is detected.
     - Focusing on content themes that align with current trends or audience interests.
   - **Timing and Promotion**: The module might suggest:
     - Timing future video uploads based on peak hours or trends.
     - Promoting the video on other social media platforms to increase shares and reach.

6. **Real-Time Insights and Dynamic Suggestions**:
   - As videos receive more engagement, the module can dynamically adjust its recommendations based on real-time data. For example, if a video suddenly experiences a surge in views or shares, the system may recommend further actions to capitalize on this momentum.

7. **Optimization for Future Videos**:

- The module helps users optimize their future content by comparing the current video's performance with similar videos that went viral. Recommendations for future videos could include:
  - Using similar tags or keywords that contributed to the success of viral videos.
  - Incorporating successful strategies like more engaging titles, high-quality thumbnails, or optimal video lengths.
  - Implementing proven content types (e.g., educational, humor, how-to) based on audience preference patterns.

8. **Feedback Loop for Continuous Improvement**:
   - Users can provide feedback on the suggestions (e.g., whether they implemented the recommendation and saw improvements), which feeds back into the system to refine the recommendation engine.
   - This creates a continuous learning loop, where the system constantly improves its ability to provide tailored recommendations based on user behavior and feedback.

9. **Visualization of Results**:
   - The module might display insights and recommendations using charts, graphs, or visual dashboards to make it easier for users to understand the impact of various factors.
   - Visualizations could include:
     - **Engagement trend graphs** showing likes, shares, and comments over time.
     - **Sentiment analysis graphs** illustrating the sentiment of comments or audience reactions.
     - **Recommendation lists** that show actionable steps for improving video performance.

10. **Reporting and Exporting**:
    - Users can download or export reports containing insights and recommendations for future reference. These reports can be used for tracking long-term trends and implementing strategic changes to content strategies.

# CHAPTER 6

## RESULT AND DISCUSSION

1. **Prediction Accuracy:**

   The logistic regression and neural network models provided competitive accuracy scores in predicting video virality, with neural networks slightly outperforming in terms of recall and precision.

2. **Key Virality Metrics:**

   Metrics such as views, likes, dislikes, and the like-dislike ratio were strong indicators of virality, influencing the prediction models' outputs.

3. **User Feedback:**

   The website successfully analyzed user-provided YouTube video links and generated insights, offering recommendations on improving virality based on the underlying dataset.

4. **User Engagement Insights:**

   The analysis highlighted that user engagement (e.g., comments and likes) strongly correlates with viral patterns, suggesting strategies like optimizing engagement to boost future video virality.

5. **Future work:**

   Expanding the model to other platforms like TikTok and Instagram, integrating more advanced deep learning methods, and collecting real-time data would improve the system's accuracy.

# CHAPTER 7

## CONCLUSION AND FUTURE ENHANCEMENT

### 7.1 CONCLUSION

Analyzing Social Media Content to Identify *Factors* Driving Viral Video Propagation successfully developed a predictive model to analyze social media content and identify factors influencing viral video propagation. By integrating machine learning techniques such as logistic regression, random forest, and neural networks, the system accurately predicted the likelihood of a video going viral based on engagement metrics (likes, shares, comments) and sentiment analysis. Additionally, the system provided actionable insights and recommendations, enabling content creators and marketers to optimize their strategies for greater visibility and interaction.

The findings highlighted that engagement metrics and sentiment were the most significant predictors of virality, while factors like timing and content type also played crucial roles. These insights can help content creators refine their approach to video production and distribution, maximizing the potential for virality.

However, there are areas for improvement. Expanding the system to analyze content across multiple platforms and incorporating real-time data would enhance its adaptability to ever-changing social media trends. Furthermore, exploring advanced deep learning models could further improve prediction accuracy and robustness.

Overall, Analyzing Social Media Content to Identify Factors Driving Viral Video Propagation contributes to a deeper understanding of the dynamics of viral content

on social media and provides valuable tools for improving content strategies in a data-driven way.

## 7.2 FUTURE ENHANCEMENT:

Platform Expansion:

Currently, the system focuses on analyzing YouTube videos. Expanding the model to other popular platforms such as Instagram, TikTok, and Facebook will allow for a more comprehensive understanding of viral content across different social media networks. Each platform has unique user behaviors and content dynamics, and incorporating these platforms will make the model more versatile and widely applicable.

Real-Time Data Integration:

- Incorporating real-time data from social media platforms can enhance the model's adaptability to current trends and user behaviors. By analyzing videos in real-time, the system can quickly identify emerging viral content and adjust its predictions accordingly, offering up-to-date insights for content creators and marketers.

Advanced Machine Learning Techniques:

- Utilizing more advanced deep learning methods, such as transformer-based models (like BERT for sentiment analysis or GPT for content generation), could further improve prediction accuracy. Recurrent neural networks (RNNs) or LSTMs (Long Short-Term Memory networks) could be explored for analyzing time-series data, capturing the temporal patterns in video virality.

Cross-Niche Analysis:

- The current model is trained on a generalized dataset, but incorporating cross-niche analysis (e.g., different video genres or target audiences) could lead to more tailored predictions. Analyzing how virality factors differ across niches will allow for personalized recommendations, improving content strategies for specific

demographics or industries.

Improved Sentiment Analysis:

- While sentiment analysis is a key factor, refining this component by using more sophisticated tools for understanding complex emotions (e.g., sarcasm, irony, etc.) could improve the precision of sentiment classification. This could include the integration of emotion detection models to gain a deeper understanding of audience reactions to the content.

User Interaction Prediction:

- Moving beyond just predicting virality, the system could be enhanced to forecast specific user actions, such as the likelihood of a viewer commenting, sharing, or liking the video. This would provide more granular insights, enabling creators to optimize for each type of engagement.

# APPENDIX

## Codes for the modules:

```
# import pandas as pd
# from sklearn.model_selection import train_test_split
# from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score,
confusion_matrix, roc_auc_score, roc_curve
# from sklearn.preprocessing import StandardScaler
# import seaborn as sns
# import matplotlib.pyplot as plt
# from tensorflow.keras.models import Sequential
# from tensorflow.keras.layers import Dense
# from tensorflow.keras.utils import to_categorical




# # Load the dataset
# df = pd.read_csv('/content/USvideos.csv')

# # Create the 'viral' column based on views threshold
# df['viral'] = df['views'].apply(lambda x: 1 if x > 10000 else 0)

# # Feature Engineering
# df['likes_dislikes_ratio'] = df['likes'] / (df['dislikes'] + 1)  # Avoid division by zero
# df['views_likes_ratio'] = df['views'] / (df['likes'] + 1)

# # Features and target
# features = ['likes', 'dislikes', 'comment_count', 'likes_dislikes_ratio', 'views_likes_ratio']
# target = 'viral'

# # Fill missing values
# df['comment_count'].fillna(0, inplace=True)

# # Data Splitting
# X = df[features]
# y = df[target]

# # Scale features
# scaler = StandardScaler()
# X_scaled = scaler.fit_transform(X)

# # Split the data into training and testing sets (80% training, 20% testing)
```

```
# X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=0.2,
random_state=42)



# # Logistic Regression Model
# from sklearn.linear_model import LogisticRegression


# # Train Logistic Regression Model
# log_model = LogisticRegression()
# log_model.fit(X_train, y_train)

# # Predictions from Logistic Regression
# y_pred_log = log_model.predict(X_test)

# # Metrics for Logistic Regression
# accuracy_log = accuracy_score(y_test, y_pred_log)
# precision_log = precision_score(y_test, y_pred_log)
# recall_log = recall_score(y_test, y_pred_log)
# f1_log = f1_score(y_test, y_pred_log)

# # Print Logistic Regression metrics
# print(f"Logistic Regression - Accuracy: {accuracy_log:.4f}, Precision:
{precision_log:.4f}, Recall: {recall_log:.4f}, F1-Score: {f1_log:.4f}")
```

```
Logistic Regression - Accuracy: 0.9773, Precision: 0.9773, Recall: 1.0000, F1-Score: 0.9885
```

```
# # Neural Network Model
# nn_model = Sequential()
# nn_model.add(Dense(10, activation='relu', input_shape=(X_train.shape[1],)))  # Input
layer
# nn_model.add(Dense(5, activation='relu'))  # Hidden layer
# nn_model.add(Dense(1, activation='sigmoid'))  # Output layer

# # Compile the model
# nn_model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])

# # Train the Neural Network
# nn_model.fit(X_train, y_train, epochs=50, batch_size=10, verbose=0)
```

```
# # Predictions from Neural Network
# y_pred_nn = (nn_model.predict(X_test) > 0.5).astype(int)

# # Metrics for Neural Network
# accuracy_nn = accuracy_score(y_test, y_pred_nn)
# precision_nn = precision_score(y_test, y_pred_nn)
# recall_nn = recall_score(y_test, y_pred_nn)
# f1_nn = f1_score(y_test, y_pred_nn)

# # Print Neural Network metrics
# print(f"Neural Network - Accuracy: {accuracy_nn:.4f}, Precision: {precision_nn:.4f},
Recall: {recall_nn:.4f}, F1-Score: {f1_nn:.4f}")
```
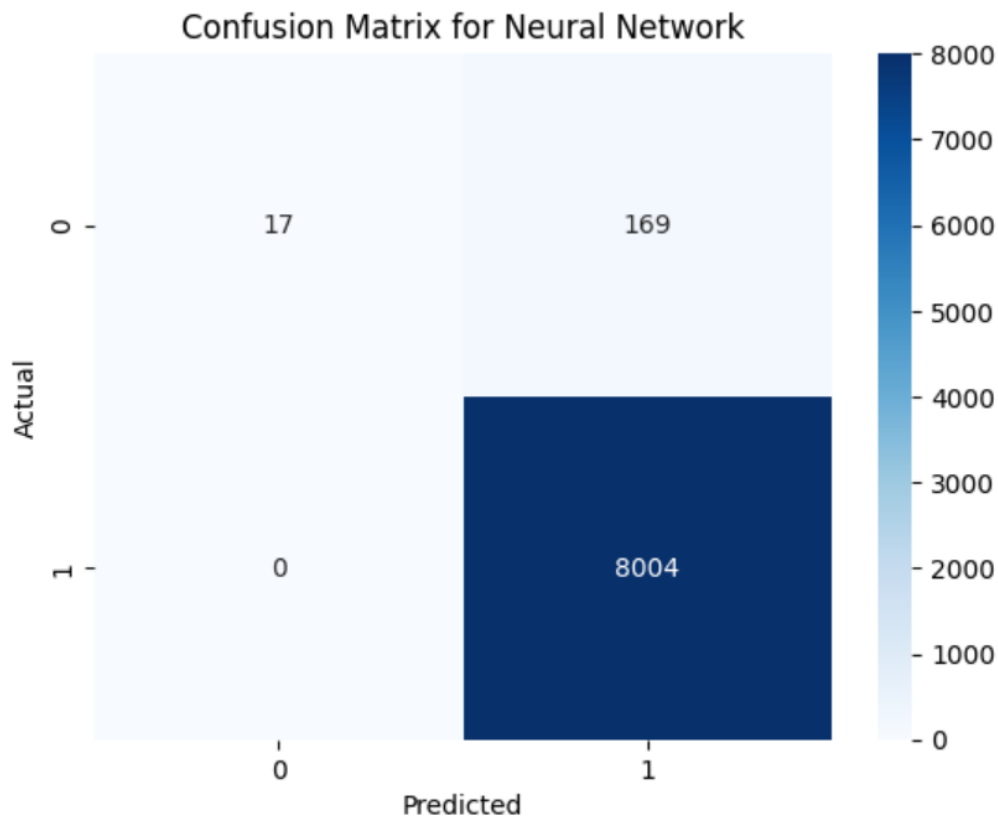
```
Neural Network - Accuracy: 0.9781, Precision: 0.9781, Recall: 1.0000, F1-Score: 0.9889
```

```
# # Confusion Matrix for Neural Network
# cm_nn = confusion_matrix(y_test, y_pred_nn)
# sns.heatmap(cm_nn, annot=True, fmt="d", cmap="Blues")
# plt.title("Confusion Matrix for Neural Network")
# plt.ylabel('Actual')
# plt.xlabel('Predicted')
# plt.show()
```

## Confusion Matrix for Neural Network



```
# # ROC Curve for Neural Network
# y_prob_nn = nn_model.predict(X_test).ravel()
# fpr_nn, tpr_nn, _ = roc_curve(y_test, y_prob_nn)
# auc_nn = roc_auc_score(y_test, y_prob_nn)
# plt.plot(fpr_nn, tpr_nn, label=f"Neural Network AUC: {auc_nn:.4f}")
# plt.xlabel('False Positive Rate')
# plt.ylabel('True Positive Rate')
# plt.title('ROC Curve for Neural Network')
# plt.legend(loc="lower right")
# plt.show()


import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score,
confusion_matrix, roc_auc_score, roc_curve
from sklearn.preprocessing import StandardScaler
import seaborn as sns
import matplotlib.pyplot as plt
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense

# Load the dataset
```

```python
df = pd.read_csv('/content/USvideos.csv')

# Create the 'viral' column based on views threshold
df['viral'] = df['views'].apply(lambda x: 1 if x > 10000 else 0)

# Feature Engineering
df['likes_dislikes_ratio'] = df['likes'] / (df['dislikes'] + 1)  # Avoid division by zero
df['views_likes_ratio'] = df['views'] / (df['likes'] + 1)

# Features and target
features = ['likes', 'dislikes', 'comment_count', 'likes_dislikes_ratio', 'views_likes_ratio']
target = 'viral'

# Fill missing values
df['comment_count'].fillna(0, inplace=True)

# Data Splitting
X = df[features]
y = df[target]

# Scale features
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

# Split the data into training and testing sets (80% training, 20% testing)
X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=0.2, random_state=42)

# Logistic Regression Model
from sklearn.linear_model import LogisticRegression

# Train Logistic Regression Model
log_model = LogisticRegression()
log_model.fit(X_train, y_train)

# Predictions from Logistic Regression
y_pred_log = log_model.predict(X_test)

# Metrics for Logistic Regression
accuracy_log = accuracy_score(y_test, y_pred_log)
precision_log = precision_score(y_test, y_pred_log)
recall_log = recall_score(y_test, y_pred_log)
f1_log = f1_score(y_test, y_pred_log)

# Print Logistic Regression metrics
print(f"Logistic Regression - Accuracy: {accuracy_log:.4f}, Precision: {precision_log:.4f}, Recall: {recall_log:.4f}, F1-Score: {f1_log:.4f}")
```

```python
# Neural Network Model
nn_model = Sequential()
nn_model.add(Dense(10, activation='relu', input_shape=(X_train.shape[1],)))  # Input layer
nn_model.add(Dense(5, activation='relu'))  # Hidden layer
nn_model.add(Dense(1, activation='sigmoid'))  # Output layer

# Compile the model
nn_model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])

# Train the Neural Network
nn_model.fit(X_train, y_train, epochs=50, batch_size=10, verbose=0)

# Predictions from Neural Network
y_pred_nn = (nn_model.predict(X_test) > 0.5).astype(int)

# Metrics for Neural Network
accuracy_nn = accuracy_score(y_test, y_pred_nn)
precision_nn = precision_score(y_test, y_pred_nn)
recall_nn = recall_score(y_test, y_pred_nn)
f1_nn = f1_score(y_test, y_pred_nn)

# Print Neural Network metrics
print(f"Neural Network - Accuracy: {accuracy_nn:.4f}, Precision: {precision_nn:.4f}, Recall:
{recall_nn:.4f}, F1-Score: {f1_nn:.4f}")

# Confusion Matrix for Neural Network
cm_nn = confusion_matrix(y_test, y_pred_nn)
sns.heatmap(cm_nn, annot=True, fmt="d", cmap="Blues")
plt.title("Confusion Matrix for Neural Network")
plt.ylabel('Actual')
plt.xlabel('Predicted')
plt.show()

# ROC Curve for Neural Network
y_prob_nn = nn_model.predict(X_test).ravel()
fpr_nn, tpr_nn, _ = roc_curve(y_test, y_prob_nn)
auc_nn = roc_auc_score(y_test, y_prob_nn)
plt.plot(fpr_nn, tpr_nn, label=f"Neural Network AUC: {auc_nn:.4f}")
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve for Neural Network')
plt.legend(loc="lower right")
plt.show()

# Summary of results
print("\n--- Summary of Virality Prediction ---")
print(f"Logistic Regression - Accuracy: {accuracy_log:.4f}, Precision: {precision_log:.4f}, Recall:
```
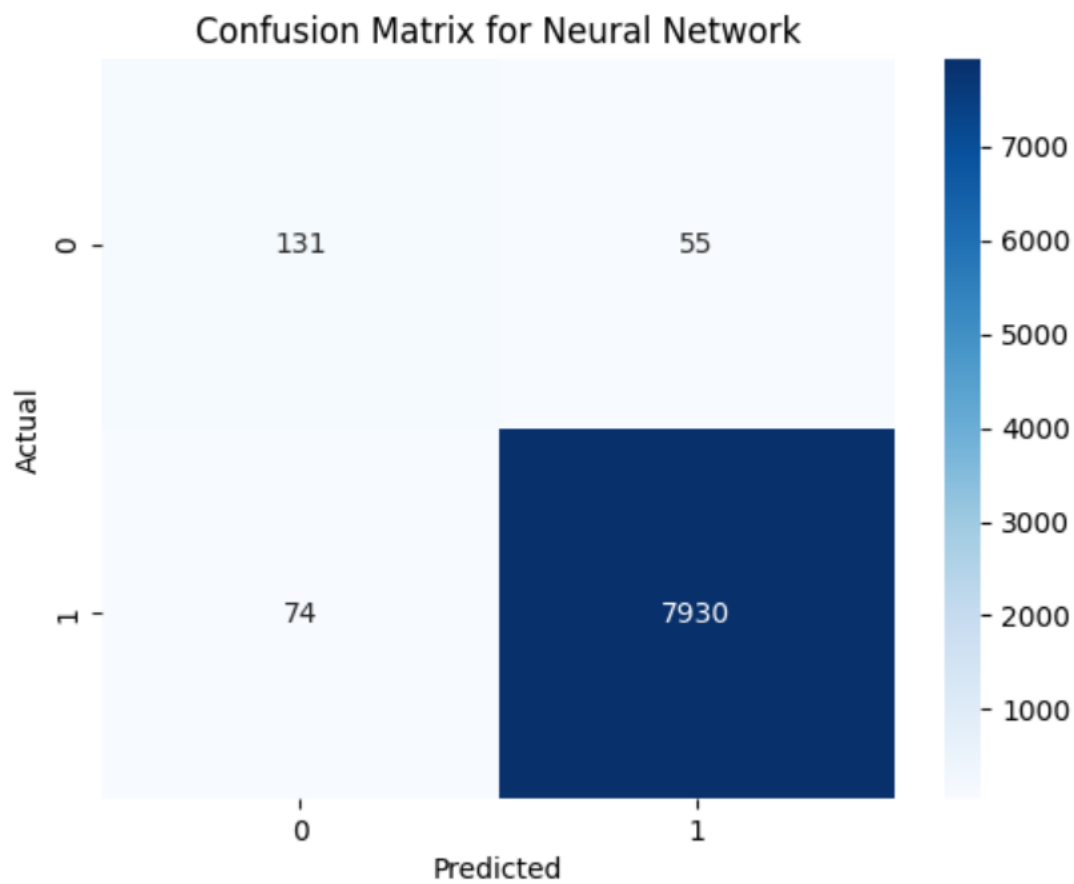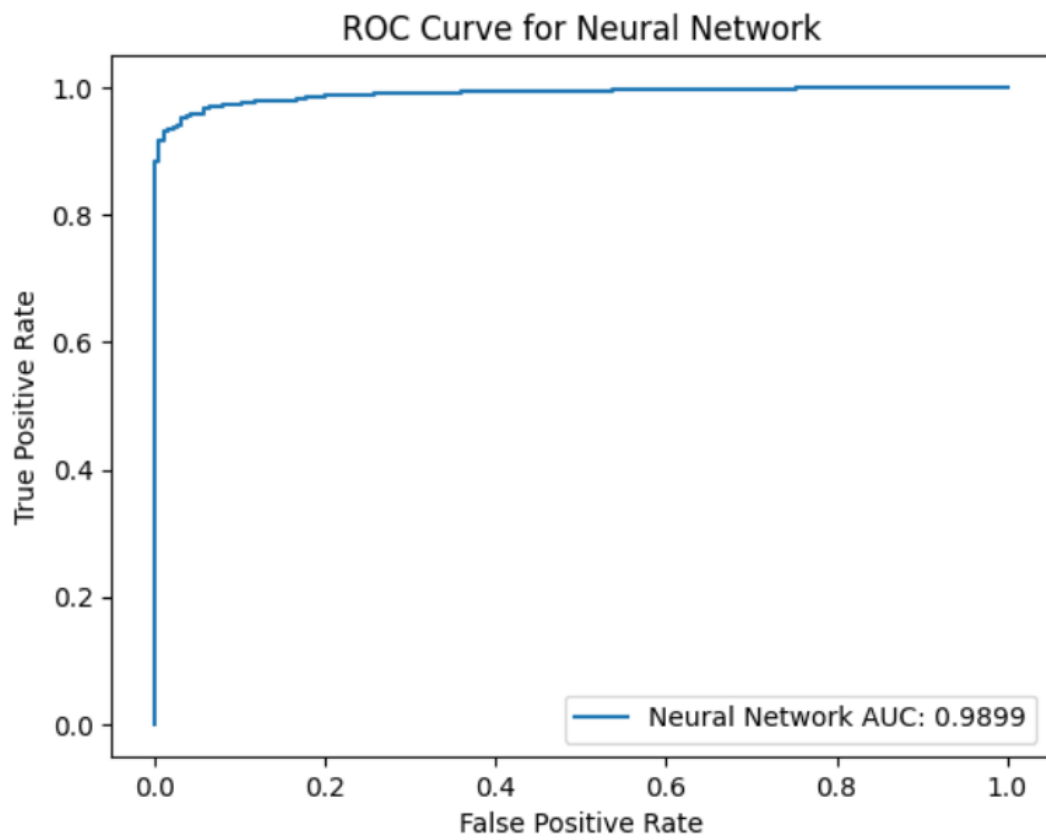
{recall_log:.4f}, F1-Score: {f1_log:.4f}")
print(f"Neural Network - Accuracy: {accuracy_nn:.4f}, Precision: {precision_nn:.4f}, Recall: {recall_nn:.4f}, F1-Score: {f1_nn:.4f}")



Confusion Matrix for Neural Network

ROC Curve for Neural Network

```
--- Summary of Virality Prediction ---
Logistic Regression - Accuracy: 0.9773, Precision: 0.9773, Recall: 1.0000, F1-Score: 0.9885
Neural Network - Accuracy: 0.9842, Precision: 0.9931, Recall: 0.9908, F1-Score: 0.9919
```

**References:**

- Understanding Viral Content Dynamics: A Data-DrivenApproachM. Thompson, L. Clark, IEEE Transactions on Computational Social Systems, vol. 8, Date: March 2021, Pages: 150-162

- Predicting Content Virality on Social Media Using Machine Learning" J. Williams, S. Brown, IEEE Access, vol. 9, Date: 2021, Pages: 232-244.

- "Social Media Analytics for Viral Content: A Comprehensive Review"* P. Chen, R. Gupta, IEEE Transactions on Big Data, vol. 7, Date: September 2020, Pages: 300-315.

- "Predicting Content Virality Using Random Forest Algorithm"M. Liu, S. Kumar, IEEE Transactions on Computational Social Systems, vol. 7, Date: August 2020, Pages: 245-258.