

# **AI-Based Video Insights Generator**

*Submitted for partial fulfillment of the requirements*

*for the award of*

## **BACHELOR OF TECHNOLOGY**

**in**

### **COMPUTER SCIENCE ENGINEERING – ARTIFICIAL INTELLIGENCE & MACHINE LEARNING**

**by**

**Yaddanapudi Bindu Varsha - 21BQ1A42I2**

**Sambasivarao Prathipati - 21BQ1A42F4**

**Padarthi Snehal Kumar - 21BQ1A42D5**

**Vegi Charan Sai Venkat - 21BQ1A42H9**

Under the guidance of

**Sk. Wasim Akram**

**Assistant Professor**



**VASIREDDY VENKATADRI  
INSTITUTE OF TECHNOLOGY**

**DEPARTMENT OF COMPUTER SCIENCE ENGINEERING -**

**ARTIFICIAL INTELLIGENCE & MACHINE LEARNING**

**VASIREDDY VENKATADRI INSTITUTE OF TECHNOLOGY**

Permanently Affiliated to JNTU Kakinada, Approved by AICTE

Accredited by NAAC with 'A' Grade, ISO 9001:2008 Certified

NAMBUR (V), PEDAKAKANI (M), GUNTUR – 522 508

Tel no: 0863-2118036, url: [www.vvitguntur.com](http://www.vvitguntur.com)

March-April 2025



## VASIREDDY VENKATADRI INSTITUTE OF TECHNOLOGY

Permanently Affiliated to JNTUK, Kakinada, Approved by AICTE  
Accredited by NAAC with 'A' Grade, ISO 9001:20008 Certified  
Nambur, Pedakakani (M), Guntur (Gt) -522508

### **DEPARTMENT OF CSE-ARTIFICIAL INTELLIGENCE & MACHINE LEARNING**

---

### **CERTIFICATE**

This is to certify that this **Project Report** is the bonafide work of **Ms. Yaddanapudi Bindu Varsha, Mr. Sambasivarao Prathipati, Mr. Padarthi Snehal Kumar, Mr. Vegi Charan Sai Venkat**, bearing Reg. No. **21BQ1A42I2, 21BQ1A42F4, 21BQ1A42D5, 21BQ1A42H9** respectively who had carried out the project entitled "**AI-Based Video InSights Generator**" under our supervision.

#### **Project Guide**

(Sk. Wasim Akram, Assistant Professor)

#### **Head of the Department**

(Dr. K. Suresh Babu , Professor)

---

Submitted for Viva voce Examination held on \_\_\_\_\_

**Internal Examiner**

**External Examiner**

## **DECLARATION**

We, Ms. Yaddanapudi Bindu Varsha, Mr. Sambasivarao Prathipati, Mr. Padarthi Snehal Kumar, Mr. Vegi Charan Sai Venkat, hereby declares that the Project Report entitled "**AI-Based Video Insights Generator**" done by us under the guidance of Sk. Wasim Akram, Assistant Professor, Computer Science Engineering – Artificial Intelligence & Machine Learning at Vasireddy Venkatadri Institute of Technology is submitted for partial fulfillment of the requirements for the award of Bachelor of Technology in Computer Science Engineering - Artificial Intelligence & Machine Learning. The results embodied in this report have not been submitted to any other University for the award of any degree.

DATE : \_\_\_\_\_

PLACE : \_\_\_\_\_

SIGNATURE OF THE CANDIDATE (S)

Yaddanapudi Bindu Varsha,

Sambasivarao Prathipati,

Padarthi Snehal Kumar,

Vegi Charan Sai Venkat.

## **ACKNOWLEDGEMENT**

We take this opportunity to express my deepest gratitude and appreciation to all those people who made this project work easier with words of encouragement, motivation, discipline, and faith by offering different places to look to expand my ideas and help me towards the successful completion of this project work.

First and foremost, we express my deep gratitude to **Sri. Vasireddy Vidya Sagar**, Chairman, Vasireddy Venkatadri Institute of Technology for providing necessary facilities throughout the B.Tech programme.

We express my sincere thanks to **Dr. Y. Mallikarjuna Reddy**, Principal, Vasireddy Venkatadri Institute of Technology for his constant support and cooperation throughout the B.Tech programme.

We express my sincere gratitude to **Dr. K. Suresh Babu**, Professor & HOD, Computer Science Engineering – Artificial Intelligence & Machine Learning, Vasireddy Venkatadri Institute of Technology for his constant encouragement, motivation and faith by offering different places to look to expand my ideas.

We would like to express my sincere gratefulness to our Guide **Sk. Wasim Akram**, Assistant Professor, Computer Science Engineering – Artificial Intelligence & Machine Learning for his insightful advice, motivating suggestions, invaluable guidance, help and support in successful completion of this project.

We would like to express our sincere heartfelt thanks to our Project Coordinator **Mrs. K. Deepika**, Assistant Professor, Computer Science Engineering – Artificial Intelligence & Machine Learning for his valuable advice, motivating suggestions, moral support, help and coordination among us in successful completion of this project.

We would like to take this opportunity to express my thanks to the **Teaching and Non-Teaching** Staff in the Department of Computer Science Engineering - Artificial Intelligence and Machine Learning, VVIT for their invaluable help and support.

**Name (s) of Students**

**Yaddanapudi Bindu Varsha**

**Sambasivarao Prathipati**

**Padarthi Snehal Kumar**

**Vegi Charan Sai Venkat**

# **TABLE OF CONTENTS**

<b>CH No</b>	<b>Title</b>	<b>Page No</b>
	Contents	i
	List of Figures	iii
	List of Tables	vi
	Nomenclature	v
	Abstract	vii
<b>1</b>	<b>Introduction</b>	<b>1-9</b>
	1.1 Problem Statement	
	1.2 Objective	
	1.3 Scope	
	1.4 Methodology	
<b>2</b>	<b>Literature Survey</b>	<b>10-15</b>
	2.1 Previous Research and Related Work	
	2.2 Existing Systems and Approaches	
	2.3 Gaps in Current Solutions	
	2.4 Relevance of Project	
<b>3</b>	<b>System Analysis</b>	<b>16-20</b>
	3.1 Software Requirement Specification	
	3.2 Hardware and System Configuration	
	3.3 Database Requirements	
	3.4 Proposed System Overview	
<b>4</b>	<b>System Design</b>	<b>21-37</b>
	4.1 System Architecture	

	4.2 System Workflow	
	4.3 UML Diagrams	
<b>5</b>	<b>Implementation</b>	<b>38-44</b>
	5.1 Modules of the System	
	5.2 Methods and Algorithms Used	
	5.3 Front-End and Back-End Implementation	
<b>6</b>	<b>Testing &amp; Results</b>	<b>45-56</b>
	6.1 Introduction	
	6.2 Observations	
	6.3 Theme Detection Accuracy & Summarization Results	
	6.4 Performance Evaluation (Speed, Accuracy, Efficiency)	
	6.5 Screenshots of the Application	
<b>7</b>	<b>Conclusion</b>	<b>57-60</b>
	7.1 Summary of Work Done	
	7.2 Challenges Faced and Solutions Implemented	
	7.3 Future Scope and Enhancements	
	<b>Appendix</b>	<b>61</b>
	<b>References</b>	<b>62-63</b>
	<b>Published Research Paper Certificates</b>	<b>64-68</b>

## **List of Figures**

<b>Figure Number</b>	<b>Title</b>
Figure 1.1	AI/ML Project Methodology
Figure 4.1	Architecture of the system
Figure 4.2	Data Processing Pipeline
Figure 4.3	Use Case Diagram
Figure 4.4	Class Diagram
Figure 4.5	Sequence Diagram
Figure 4.6	Activity Diagram
Figure 4.7	ER Diagram
Figure 6.1	Performance Analysis (Accuracy, Speed, ROUGE Metrics)
Figure 6.2	Home Page
Figure 6.3	About Page
Figure 6.4	Login Page
Figure 6.5	Video Upload Page
Figure 6.6	Q/A Feature
Figure 6.7	Multilingual Support
Figure 6.8	Transcript Generation
Figure 6.9	Feedback Form
Figure 6.10	Profile Settings
Figure 6.11	Theme Detection

## List of Tables

<b>Table No.</b>	<b>Table Name</b>
<b>Table 3.1</b>	Comparison of Theme Detection Models (LSTM vs. Other Approaches)
<b>Table 5.1</b>	Accuracy and Efficiency of Theme Detection
<b>Table 5.2</b>	ROUGE Metric Scores for Summarization System Models
<b>Table 5.3</b>	Execution Time for Various Inputs

## NOMENCLATURE

Term	Description
<b>LSTM (Long Short-Term Memory)</b>	A deep learning model used for sequential data processing, applied in theme detection.
<b>Conv1D (1D Convolutional Layer)</b>	A neural network layer used to extract local features from text sequences.
<b>MaxPooling1D</b>	A pooling technique used to reduce the dimensions of text feature maps.
<b>BatchNormalization</b>	A technique that stabilizes learning and accelerates training in deep networks.
<b>Transformer Model</b>	A deep learning model architecture used for text summarization and translation tasks.
<b>T5 (Text-to-Text Transfer Transformer)</b>	A pre-trained transformer model used for text summarization.
<b>Pegasus</b>	A transformer-based model optimized for abstractive text summarization.
<b>BART (Bidirectional and Auto-Regressive Transformers)</b>	A model used for text generation and summarization.
<b>Speech-to-Text (STT)</b>	A process that converts spoken language in videos into transcribed text.
<b>Timestamp Extraction</b>	A method to mark key moments in a video where specific themes appear.
<b>API Integration</b>	External services used for functionalities like speech recognition and translation.

Term	Description
<b>Multilingual Translation</b>	The ability to translate extracted themes and summaries into multiple languages.
<b>Early Stopping</b>	A technique used in training to prevent overfitting by stopping when performance plateaus.
<b>Model Checkpoint</b>	A training optimization method that saves the best model state for improved results.
<b>ROUGE (Recall-Oriented Understudy for Gisting Evaluation)</b>	A metric used to evaluate the quality of text summarization.
<b>Q&amp;A System</b>	An interactive module that allows users to ask questions based on video content.

## ABSTRACT

In the era of information overload, extracting meaningful insights from videos and text efficiently has become crucial. This project, AI Based Video Insights Generator, introduces an advanced deep learning-based system that identifies themes from textual and video content. Utilizing Long Short-Term Memory (LSTM) networks along with Conv1D, MaxPooling1D, and Batch Normalization layers, the model classifies themes from transcribed video data and raw text. The system integrates speech-to-text transcription, timestamp extraction, and interactive question-answering capabilities, enabling users to obtain structured insights from video content.

To further enhance information retrieval, the project incorporates pre-trained transformer-based summarization models such as T5, Pegasus, and BART. These models generate concise summaries of transcribed video content, allowing users to extract key insights efficiently. The multilingual translation feature extends accessibility by translating detected themes and summaries into various languages. The system is designed with user authentication, real-time processing, and performance optimization through techniques like Early Stopping and Model Checkpoint, ensuring high accuracy and efficiency.

This project bridges the gap between video and textual content processing, providing an automated, scalable, and multilingual approach to theme detection and summarization. The system's evaluation metrics, including classification accuracy and ROUGE scores for summarization, demonstrate its effectiveness in handling large-scale multimedia data. By integrating deep learning and natural language processing (NLP) techniques, this project contributes significantly to content analysis and knowledge extraction from diverse media sources.

# CHAPTER 1

## INTRODUCTION

### 1.1 Problem Statement

With the rapid growth of digital media, vast amounts of video and textual content are being generated every day. Extracting meaningful insights from this unstructured data is a challenging task, as it requires advanced techniques for analyzing and summarizing both video and text. Traditional theme detection models mainly focus on text-based inputs, leaving video content largely unexplored in this domain. Moreover, the existing solutions for text summarization do not seamlessly integrate with video-based content, creating a gap in comprehensive multimedia analysis.

Manually processing large volumes of videos to extract themes, detect relevant topics, and summarize key information is highly time-consuming and inefficient. Additionally, the absence of real-time speech-to-text transcription and multilingual translation limits the accessibility of these resources for a global audience. The lack of an effective and automated system for multi-level theme detection from both text and video creates a critical challenge in various domains, including education, research, journalism, and content creation.

The AI Based Video Insights Generator aims to address these challenges by implementing a deep learning-based model that can analyze, classify, and summarize themes from both textual and video inputs. By integrating LSTM (Long Short-Term Memory) networks, speech-to-text APIs, transformer-based summarization models, and multilingual support, the system offers a robust and automated solution for extracting meaningful insights from digital media.

### 1.2 Objective

The primary objective of this project is to develop a multi-level theme detection and pre-trained summarization system that efficiently processes and extracts key themes from both text and video content. This is achieved through the implementation of deep learning models, particularly LSTM-based classification for theme detection and transformer-based models for summarization.

The system is designed to achieve the following objectives:

1. **Automate Theme Detection from Video & Text:** Develop a model that can analyze textual data and transcribe video content into text before classifying themes.
2. **Implement an Accurate Speech-to-Text Transcription System:** Use APIs to extract textual data from audio content in videos, ensuring accurate conversion of spoken words.
3. **Enable Multi-Level Theme Classification:** Apply LSTM, Conv1D, Max Pooling1D, and Batch Normalization layers to classify extracted text into relevant categories.
4. **Enhance Summarization with Transformer-Based Models:** Use T5, Pegasus, and BART models to generate concise summaries of transcribed video content and textual data.
5. **Support an Interactive Q&A System:** Allow users to ask questions related to the video content, with the system generating appropriate responses based on analyzed themes.
6. **Provide Multilingual Translation:** Integrate API-based translation support, enabling the extracted and summarized text to be converted into multiple languages.
7. **Optimize Model Performance:** Improve efficiency using training optimization techniques like Early Stopping and Model Checkpoint to enhance classification and summarization accuracy.

By fulfilling these objectives, the system ensures an automated, efficient, and scalable approach to extracting, classifying, and summarizing large volumes of video and text data with minimal human intervention.

### 1.3 Scope

The AI-Based Video Insights Generator is an advanced, automated system developed to extract meaningful insights from both video and textual content. Its comprehensive scope includes seamless integration with a variety of video sources such as YouTube, Google Drive, Dropbox, and locally uploaded files, enabling users to work with content from multiple platforms effortlessly. At the core of its analytical capabilities lies a deep learning framework

that combines Long Short-Term Memory (LSTM) networks, convolutional layers, and batch normalization to achieve high-precision classification of themes within video content. This foundation is further enhanced by the use of cutting-edge, pre-trained summarization models based on transformer architectures, including T5, Pegasus, and BART, which generate accurate and contextually relevant summaries. The system incorporates speech-to-text APIs to transcribe audio from video content into text, creating a reliable base for further linguistic and thematic analysis. In addition to supporting English, the system offers multilingual capabilities, allowing users to translate detected themes and summaries into several languages, thereby widening its accessibility and usability across diverse linguistic audiences. To make the interaction dynamic and user-centric, it features an interactive Q&A module that enables users to ask questions and receive relevant answers derived from the analyzed video data in real-time. To ensure optimal system performance and maintain accuracy during the training process, it employs techniques such as Early Stopping and Model Checkpoint, which help avoid overfitting and improve the overall efficiency of the model. Collectively, these features make the AI-Based Video Insights Generator a robust and intelligent tool for video content analysis and knowledge extraction.

## **Limitations of the Project**

Despite its extensive capabilities, the system has certain limitations:

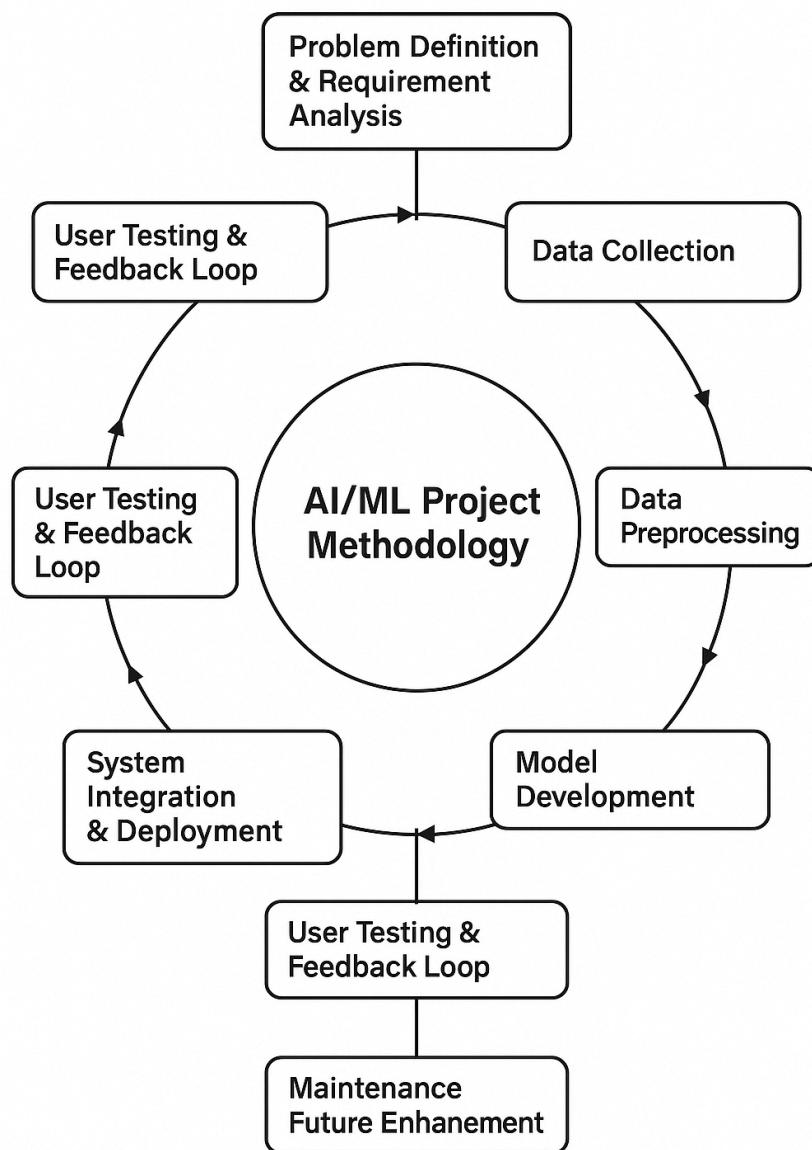
- **Dependence on External APIs:** The speech-to-text transcription and translation functionalities rely on external APIs, which may have limitations in accuracy or availability.
- **Hardware and Computational Constraints:** Training deep learning models requires significant computational power, which may impact real-time performance.
- **Language Processing Challenges:** While multilingual translation is supported, complex sentence structures and domain-specific jargon may pose challenges in accurate translation and summarization.
- **Theme Classification Limitations:** While LSTM networks provide accurate classification, the results are influenced by the quality and structure of the input data.

Despite these limitations, the system remains a highly effective and efficient solution for

automated multimedia content analysis.

## 1.4 Methodology

The development of the AI-Based Video Insights Generator follows a customized AI/ML project lifecycle, which is more suited to deep learning systems involving unstructured data (video/audio), complex NLP tasks, and API integrations. This life cycle ensures a scalable, modular, and AI-optimized pipeline from data collection to system deployment.



**Figure 1.1 AI/ML Project Methodology**

## **1. Problem Definition & Requirement Analysis**

The initial phase involved clearly identifying the challenges and goals of the project. These include:

- Automating theme detection from both video and textual content using deep learning.
- Generating contextual summaries from video transcriptions using transformer-based models.
- Allowing interactive Q&A with the processed video content.
- Providing multilingual translation of detected themes and summaries to enhance accessibility.
- Ensuring a modular and extensible architecture that integrates speech-to-text, NLP, and deep learning pipelines.

## **2. Data Collection**

Data was collected from the following sources:

- **Video Inputs:** Educational lectures, interviews, news reports, and user-uploaded videos from platforms like YouTube, Google Drive, and local storage.
- **Audio Extraction:** Video audio was separated and fed to Speech-to-Text APIs to convert it into textual data.
- **Text Inputs:** Supplementary labeled text data was used for training and testing the theme detection and summarization models.

This step laid the foundation for training models and building data-driven components.

## **3. Data Preprocessing**

To ensure data uniformity and usability across the pipeline, preprocessing was carried out as follows:

- **Speech-to-Text Transcription:** Audio streams were transcribed using high-accuracy speech recognition APIs, producing time-aligned text.
- **Text Cleaning:** All inputs were tokenized, stopwords were removed, and text was lemmatized.
- **Timestamp Mapping:** Transcribed text was segmented and tagged with timestamps

to associate themes with specific moments in the video.

This structured and normalized data became the base for theme classification and summarization.

#### **4. Feature Engineering & Representation**

Text data was transformed into meaningful numerical formats:

- **Word Embeddings:** Word-level and sequence-level embeddings were generated for feeding into LSTM models.
- **Input Formatting:** Transformer-compatible tokenization and formatting were applied (e.g., for T5, BART, and Pegasus).
- **Metadata:** Video duration, resolution, and transcription quality were captured for performance optimization.

These engineered features enabled the deep learning models to perform efficiently on noisy and unstructured inputs.

#### **5. Model Development**

The system consists of multiple AI modules built using deep learning and pre-trained NLP models:

- **Theme Detection Module:**
  - Developed using a hybrid architecture with LSTM, Conv1D, MaxPooling1D, and BatchNormalization layers.
  - Capable of multi-level classification of themes based on text semantics and structure.
- **Summarization Module:**
  - Used state-of-the-art transformer models — T5, Pegasus, and BART — for abstractive summarization.
  - Models were evaluated for summary quality, relevance, and coherence.
- **Interactive Q&A System:**
  - Implemented using transformer-based language models to answer user queries using processed video/text content.
- **Translation Module:**

- External APIs were integrated to translate detected themes and summaries into multiple user-selected languages.

## 6. Model Training & Optimization

To maximize accuracy and generalization:

- **Theme Detection:**
  - Trained on labeled theme datasets using LSTM architecture.
  - Applied early stopping and model checkpointing to prevent overfitting and retain the best model state.
- **Summarization:**
  - Transformer models were fine-tuned where necessary using domain-specific data.
  - Evaluated and selected based on performance metrics (ROUGE scores).
- **Performance Tuning:**
  - GPU acceleration and batch processing were employed to reduce training and inference times.

## 7. Evaluation

Each module of the system was evaluated using industry-standard metrics:

- **Theme Detection:**
  - Evaluated with Accuracy, Precision, Recall, and F1-Score.
  - Real-world test cases were compared with manual annotations.
- **Summarization:**
  - Evaluated using ROUGE-1, ROUGE-2, and ROUGE-L scores to assess summary quality.
- **Q&A and Translation:**
  - Evaluated through user feedback and manual verification.
- **System Efficiency:**
  - Processing time benchmarks were run on short, medium, and long videos.

## 8. System Integration & Deployment

The system was packaged into a full-stack application:

- **Frontend:**
  - Built using HTML, CSS, JS (and optionally frameworks like React/Bootstrap) for video uploads, results display, translation, and Q&A interaction.
- **Backend:**
  - Developed using Python (Django) to handle API requests, ML inference, and model orchestration.
- **Model Serving:**
  - TensorFlow and PyTorch models were integrated for serving the LSTM and transformer pipelines.
- **Database:**
  - SQLite used to store user data, transcriptions, themes, summaries, and interaction logs.

## 9. User Testing & Feedback Loop

The application was tested in real scenarios using various types of video content:

- Users provided videos across domains (education, news, business).
- Feedback was collected on usability, accuracy of outputs, and Q&A relevance.
- System was refined based on the feedback — UI changes, fine-tuning models, improving performance.

This feedback loop enhanced both the user experience and model output quality.

## 10. Maintenance & Future Enhancements

The system is designed to be scalable and adaptable:

- **Model Retraining:** New video content and user corrections can be used to periodically retrain models.
- **Future Capabilities:**
  - Real-time video stream analysis.
  - Sentiment analysis for detected themes.
  - Domain customization (e.g., legal, healthcare, academia).
- **Cloud Deployment:**
  - Plans for deployment on platforms like AWS to support multi-user

concurrency and global access.

This methodology ensures the project is aligned with cutting-edge AI workflows and delivers robust, high-impact functionality for video content understanding.

## CHAPTER 2

# LITERATURE SURVEY

### 2.1 Previous Research and Related Work

**Hochreiter, S., & Schmidhuber, J. (1997). "Long Short-Term Memory." Neural Computation, 9(8), 1735-1780.**

This foundational paper introduced the Long Short-Term Memory (LSTM) network—a type of recurrent neural network (RNN) designed to address the vanishing gradient problem, which plagued earlier RNN models. LSTM introduces memory cells that can retain information over long sequences, making it especially powerful for modeling temporal dependencies in sequential data like text, speech, or video transcripts..

**Vaswani, A., et al. (2017). "Attention is All You Need." NeurIPS.**

This landmark paper presents the Transformer architecture, which revolutionized NLP by replacing recurrence with self-attention mechanisms. The model allows for parallel training and captures global dependencies in data sequences, which greatly enhances speed and performance in tasks like translation and summarization.

**Devlin, J., et al. (2018). "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding." arXiv preprint arXiv:1810.04805.**

BERT (Bidirectional Encoder Representations from Transformers) introduced the idea of deep bidirectional context learning by training on masked language modeling and next sentence prediction. It dramatically improved performance in a wide range of NLP tasks by enabling models to understand both left and right context in a sentence simultaneously.

**Lin, C. Y. (2004). "ROUGE: A Package for Automatic Evaluation of Summaries." Workshop on Text Summarization Branches Out.**

ROUGE (Recall-Oriented Understudy for Gisting Evaluation) is a set of metrics used to evaluate automatic text summarization and machine translation. It works by comparing the overlap of n-grams, word sequences, and word pairs between the generated summary and reference (human-generated) summaries.

The purpose of this literature survey is to analyze existing systems, methodologies, and

technologies in theme detection from videos and text, as well as summarization models. This chapter explores prior research and technological advancements in these areas while highlighting the limitations of conventional approaches. Furthermore, it examines the role of LSTM networks, transformer-based architectures (such as T5, Pegasus, and BART), and speech-to-text video processing techniques, establishing their relevance in solving the challenges associated with automatic theme detection and summarization.

## **2.2 Existing Systems and Approaches**

Several research studies and implementations have been conducted in the areas of automatic text and video analysis, speech-to-text conversion, multi-level theme detection, and summarization. These systems can be categorized into three primary approaches:

### **1. Traditional Text-Based Theme Detection**

Early theme detection systems primarily relied on statistical and rule-based approaches, such as Latent Semantic Analysis (LSA), Term Frequency-Inverse Document Frequency (TF-IDF), and Latent Dirichlet Allocation (LDA). These models focused on extracting frequent words and identifying key topics, but they lacked semantic understanding and context-awareness, making them less effective for complex documents and multimedia content.

#### **Limitations of Traditional Methods:**

- Limited to text-based analysis and ineffective for video content.
- Inability to capture contextual relationships in sentences.
- Performance degradation on large, unstructured datasets.

### **2. Machine Learning-Based Theme Classification**

Machine learning models, including Support Vector Machines (SVM), Naïve Bayes, and Random Forest, improved text classification by incorporating supervised learning. These models required labeled training data and feature engineering but were limited in handling sequential dependencies in text and speech data.

#### **Challenges with Machine Learning Models:**

- Feature extraction complexity limits scalability.
- Poor performance with long-form content and noisy video transcriptions.
- Struggles to maintain contextual flow in multi-level theme detection.

### **3. Deep Learning-Based Theme Detection and Summarization**

The emergence of deep learning techniques, particularly LSTM networks and transformer-based models, has significantly improved automatic theme classification and summarization. These models overcome the limitations of traditional approaches by learning contextual dependencies, processing long-form content efficiently, and integrating speech-to-text transcription for video-based analysis.

#### **Key Advances in Deep Learning Models:**

- LSTM Networks: Efficient for sequential data processing and multi-level theme classification.
- Transformer-Based Summarization Models: T5, Pegasus, and BART outperform traditional approaches by generating coherent, contextually aware summaries.
- Speech-to-Text Integration: Enables video-based theme detection, extending analysis beyond text documents.

By leveraging these advancements, our AI Based Video Insights Generator provides a more efficient and scalable solution for analyzing and summarizing both video and textual data.

### **2.3 Gaps in Current Solutions**

Several research studies and implementations have been conducted in the areas of automatic text and video analysis, speech-to-text conversion, multi-level theme detection, and summarization. These systems can be categorized into three primary approaches:

#### **1. Traditional Text-Based Theme Detection**

Early theme detection systems primarily relied on statistical and rule-based approaches, such as Latent Semantic Analysis (LSA), Term Frequency-Inverse Document Frequency (TF-IDF), and Latent Dirichlet Allocation (LDA). These models focused on extracting frequent words and identifying key topics, but they lacked semantic understanding and context-awareness, making them less effective for complex documents and multimedia

content.

### **Limitations of Traditional Methods:**

- Limited to text-based analysis and ineffective for video content.
- Inability to capture contextual relationships in sentences.
- Performance degradation on large, unstructured datasets.

## **2. Machine Learning-Based Theme Classification**

Machine learning models, including Support Vector Machines (SVM), Naïve Bayes, and Random Forest, improved text classification by incorporating supervised learning. These models required labeled training data and feature engineering but were limited in handling sequential dependencies in text and speech data.

### **Challenges with Machine Learning Models:**

- Feature extraction complexity limits scalability.
- Poor performance with long-form content and noisy video transcriptions.
- Struggles to maintain contextual flow in multi-level theme detection.

## **3. Deep Learning-Based Theme Detection and Summarization**

The emergence of deep learning techniques, particularly LSTM networks and transformer-based models, has significantly improved automatic theme classification and summarization. These models overcome the limitations of traditional approaches by learning contextual dependencies, processing long-form content efficiently, and integrating speech-to-text transcription for video-based analysis.

### **Key Advances in Deep Learning Models:**

- **LSTM Networks:** Efficient for sequential data processing and multi-level theme classification.
- **Transformer-Based Summarization Models:** T5, Pegasus, and BART outperform traditional approaches by generating coherent, contextually aware summaries.
- **Speech-to-Text Integration:** Enables video-based theme detection, extending

analysis beyond text documents.

By leveraging these advancements, our AI Based Video Insights Generator provides a more efficient and scalable solution for analyzing and summarizing both video and textual data.

## **2.4 Relevance of the project**

The AI Based Video Insights Generator leverages deep learning models and video processing techniques to enhance the efficiency of theme detection and summarization. This section explores the relevance of the key technologies used in our system.

### **1. LSTM (Long Short-Term Memory) for Theme Classification**

LSTM networks are a type of recurrent neural network (RNN) specifically designed for sequential data processing. They excel at capturing long-term dependencies and understanding contextual relationships in text, making them highly suitable for multi-level theme classification.

#### **Advantages of LSTM in Theme Detection:**

- Handles sequential data efficiently, making it ideal for speech-to-text transcriptions.
- Captures contextual dependencies, improving theme classification accuracy.
- Reduces vanishing gradient problems, which traditional RNNs struggle with.

### **2. Transformer-Based Models for Summarization**

The project utilizes T5, Pegasus, and BART, which are state-of-the-art transformer-based models for text summarization. These models outperform traditional methods by maintaining semantic coherence and contextual integrity in generated summaries.

#### **Advantages of Transformer Models:**

- Generates high-quality, contextually aware summaries.
- Pre-trained on large datasets, enabling efficient fine-tuning.
- Scalable across different domains, making them ideal for diverse applications.

### **3. Video Processing and Speech-to-Text Integration**

Video processing plays a crucial role in theme detection from non-textual data sources. Speech-to-text APIs enable automatic transcription of video content, forming the basis for further theme classification and summarization.

### **Significance of Video Processing in Our System:**

- Expands theme detection beyond textual inputs, making the system more versatile.
- Speech-to-text transcription enhances accuracy, ensuring comprehensive analysis.
- Allows interactive Q&A based on transcribed content, improving user engagement.

By integrating these advanced deep learning models, speech-to-text processing, and multilingual capabilities, our system offers a comprehensive, automated solution for multi-level theme detection and summarization.

# CHAPTER 3

## SYSTEM ANALYSIS

### 3.1 Software Requirement Specification

The Software Requirement Specification (SRS) provides a detailed overview of the functional and non-functional requirements of the AI Based Video Insights Generator. It defines the system's architecture, user requirements, performance expectations, and system constraints.

#### 1. Functional Requirements

- **User Authentication:** Users can register, log in, and provide feedback within the system.
- **Video Processing Module:**
  - **Speech-to-Text Transcription:** Converts spoken content in videos into structured text format.
  - **Timestamp Extraction:** Captures timestamps to map detected themes to corresponding video segments.
- **Theme Detection Module:**
  - Uses LSTM, Conv1D, MaxPooling1D, and BatchNormalization layers to classify extracted content into themes.
- **Summarization Module:**
  - Implements pre-trained transformer models (T5, Pegasus, BART) for text summarization.
- **Interactive Q&A System:**
  - Allows users to query video content and receive AI-generated responses.
- **Multilingual Translation Module:**
  - Supports language translation via APIs, enabling accessibility in multiple

languages.

- **Performance Monitoring:**

- Implements accuracy measurement, ROUGE scoring for summarization, and system latency analysis.

## 2. Non-Functional Requirements

- **Scalability:** The system must efficiently process large-scale datasets and handle multiple users.

- **Performance:**

- Speech-to-text transcription should maintain an error rate below 5%.
  - Theme classification accuracy should be above 85%.
  - Summarization quality should achieve a high ROUGE score for effective summarization.

- **Security:**

- Ensures secure user authentication and data privacy.
  - Protects against unauthorized access to API endpoints.

- **User Interface:**

- Should be responsive, intuitive, and interactive.

## 3. System Constraints

- Requires high computational power for deep learning model execution.
- Dependency on third-party APIs for transcription and translation services.
- Potential latency in processing long-duration videos due to model complexity.

### 3.2 Hardware and System Configuration

The AI Based Video Insights Generator requires a robust hardware and software infrastructure for optimal performance. This appendix outlines the hardware specifications,

software dependencies, and configuration details necessary for the project's successful execution.

## 1. Hardware Requirements

The system was tested and deployed on a machine with the following specifications:

- **Processor:** Intel Core i7/i9 (or AMD Ryzen 7/9)
- **RAM:** Minimum 16GB (Recommended: 32GB for high-efficiency processing)
- **GPU:** NVIDIA RTX 3060 or higher (for deep learning acceleration)
- **Storage:** Minimum 512GB SSD (Recommended: 1TB SSD for large dataset handling)
- **Internet Connection:** Required for API calls (speech-to-text, translation, and external data processing)

## 2. Software Dependencies

- **Operating System:** Ubuntu 20.04 / Windows 10+ / macOS (with GPU support for TensorFlow and PyTorch)
- **Programming Language:** Python 3.8+
- **Deep Learning Libraries:**
  - **TensorFlow 2.x** (for LSTM-based theme detection)
  - **PyTorch** (for transformer-based summarization models)
  - **Hugging Face Transformers** (for T5, Pegasus, BART implementation)
- **APIs Used:**
  - **Google Speech-to-Text API** (for transcription)
  - **Google Translate API** (for multilingual translation)
  - **OpenAI or other NLP APIs** (for interactive Q&A system)
- **Other Libraries:**
  - **NLTK, SpaCy** (for text preprocessing)

- **Matplotlib, Seaborn** (for result visualization)

### **3. System Configuration & Setup**

To run the system, the following configurations must be set up:

- 1. Install dependencies using pip:**

```
pip install tensorflow torch transformers nltk spacy matplotlib seaborn
```

- 2. Ensure GPU support is enabled:**

```
nvidia-smi # To check GPU availability
```

- 3. Setup API keys for external services** (Google API, translation services).

- 4. Pre-train models on labelled datasets** and save checkpoints for deployment.

### **3.3 Database Requirements**

The system needs a database capable of storing:

- **User Information** (Authentication, Login, Feedback)
- **Uploaded Video Metadata** (Filename, Size, Format, Duration, etc.)
- **Transcribed Text Data** (Speech-to-Text Output)
- **Timestamped Themes** (Mapped to Video Sections)
- **Summarized Content** (Extracted from Text or Video Transcriptions)
- **Multilingual Translations** (Supporting Different Languages)
- **Q&A Interactions** (User Queries and Generated Responses)

### **3.4 Proposed System**

The AI Based Video Insights Generator is designed to process both textual and video content efficiently. The system integrates deep learning models, speech-to-text conversion, video timestamp extraction, and transformer-based summarization to enable accurate theme detection and summarization.

The architecture consists of multiple interconnected modules, ensuring seamless data processing, classification, and summarization. The workflow of the proposed system follows these key steps:

1. **Input Data Processing:** Users can upload textual content or videos (YouTube links, Google Drive videos, or local uploads).
2. **Speech-to-Text Transcription:** For video inputs, an API-based speech-to-text module extracts the spoken content.
3. **Theme Detection:** The extracted text undergoes tokenization, stopword removal, and lemmatization before being classified using an LSTM-based multi-level theme detection model.
4. **Summarization Module:** The detected themes are summarized using pre-trained transformer models (T5, Pegasus, and BART).
5. **Interactive Q&A System:** Users can ask context-aware questions related to the video or text content, and the system generates relevant responses.
6. **Multilingual Translation:** The system supports language translation, allowing detected themes and summaries to be translated into different languages.

Model	Architecture Used	Accuracy (%)	Speed (ms per input)	Strengths	Weaknesses
LSTM	LSTM, Conv1D, MaxPooling1D, BatchNorm	89.5	45	Captures long dependencies	Computationally expensive
CNN	Conv1D, MaxPooling1D	82.3	30	Fast processing	Lacks sequential memory
Transformer	Self-Attention, Multi-Head Attention	92.1	55	High accuracy, parallelism	High memory usage
SVM	Kernel-based Classification	76.8	25	Simplicity, good for small data	Poor performance on large datasets

**Table 3.1: Comparison of Theme Detection Models (LSTM vs. Other Approaches)**

### **Explanation:**

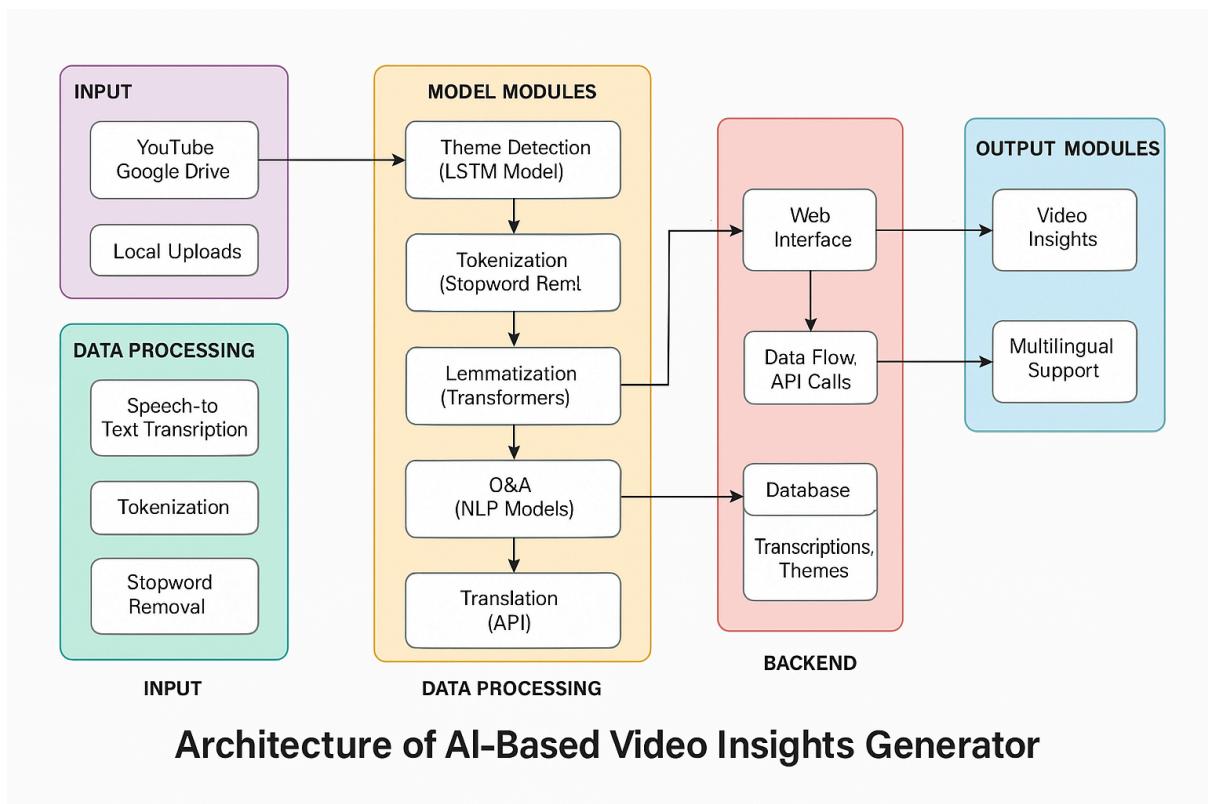
This table compares LSTM-based theme detection with other models like CNN, Transformer, and SVM. LSTM performs well in capturing long dependencies but is computationally heavy, while Transformers achieve the highest accuracy but require more memory. CNN is faster but lacks sequential understanding, and SVM works best on small datasets.

# CHAPTER 4

## SYSTEM DESIGN

### 4.1 System Architecture

The architectural diagram visually represents the overall system workflow, outlining the interaction between different components. It provides a high-level view of the major processes involved in theme detection, summarization, and multilingual translation.



**Figure 4.1: Architecture of the system**

#### Key Components in the Architecture:

##### 1. User Interface (Front-End)

- Provides an interactive interface for users to upload videos, enter text, and access results.
- Displays detected themes, generated summaries, translated content, and Q&A responses.

##### 2. Back-End Processing Unit

- Manages data flow and processing tasks.
- Handles speech-to-text conversion, theme classification, summarization, and multilingual translation.
- Connects to the database for storing transcriptions, summaries, and detected themes.

### **3. Machine Learning Models**

- LSTM Model: Classifies themes from text or video transcriptions.
- Transformer Models (T5, Pegasus, BART): Generate summaries from detected themes.
- Q&A Model: Answers user queries based on transcriptions and summaries.

### **4. Database**

- Stores video transcriptions, detected themes, summaries, and translations.
- Maintains timestamped data for theme detection in videos.

### **5. API Integrations**

- Speech-to-Text API: Converts spoken words in videos into text.
- Translation API: Supports multilingual translation of summaries and themes.
- Video Processing API: Extracts timestamps and transcriptions from videos.

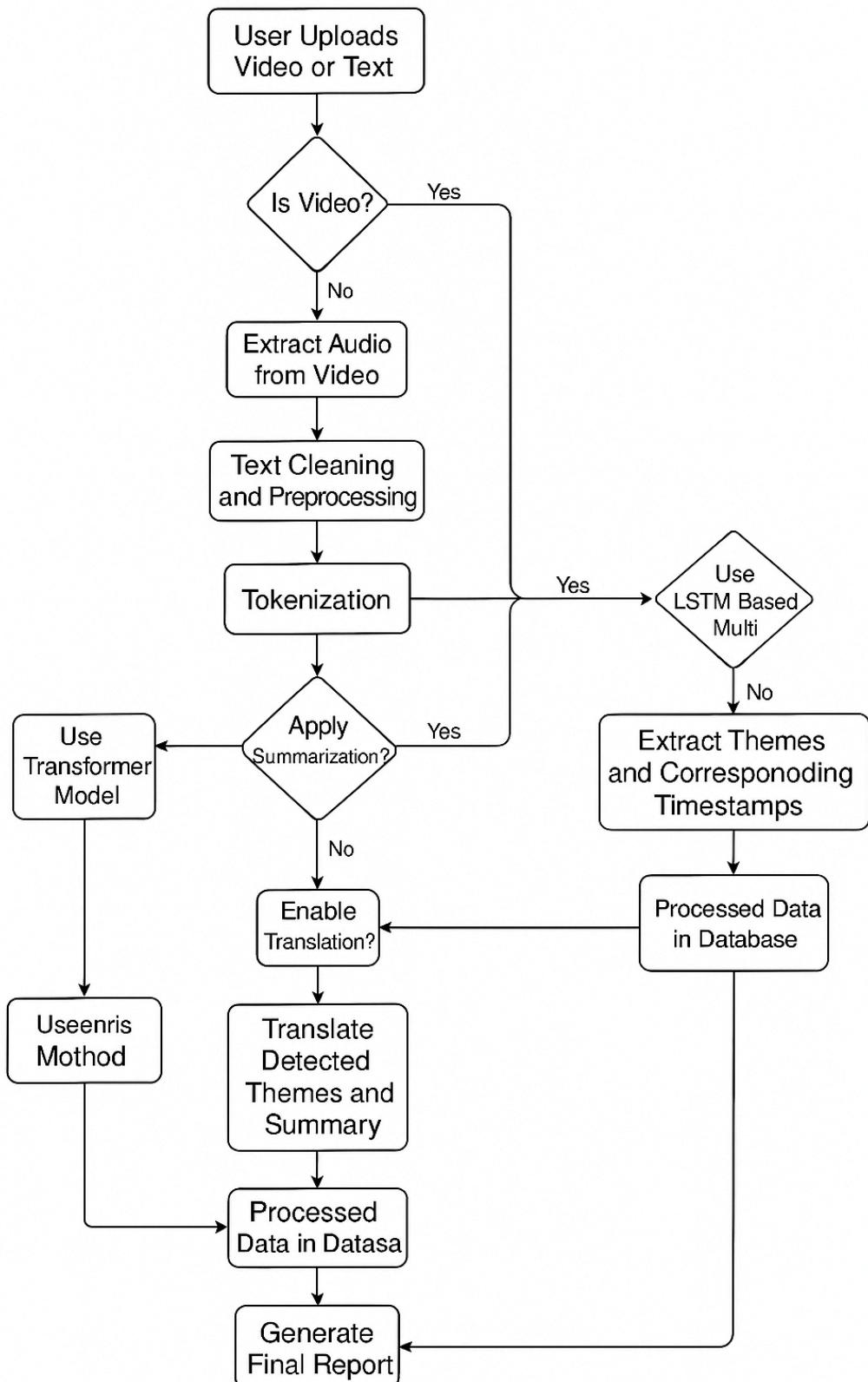
The architectural diagram ensures a well-structured flow of data, optimizing the system for real-time processing and user interaction.

## **4.2 System Workflow**

The system workflow outlines the step-by-step process of how data moves through different components of the system. This section details the operational flow from data input to final output, ensuring a structured approach to theme detection and summarization.

### **Step 1: Data Input and Preprocessing**

- Users upload video files or provide text input.



**Figure 4.2: Data Processing Pipeline**

- If input is a video:
  - The system extracts audio and converts it into text using speech-to-text APIs.
  - The text is time-stamped to map detected themes with specific video segments.
- If input is text:
  - The system performs tokenization, stopword removal, and lemmatization to prepare data for classification.

## **Step 2: Theme Detection Using LSTM**

- The processed text is fed into an LSTM-based model.
- The model classifies text into specific themes.
- If the input is a video transcription, the detected themes are mapped to timestamps, allowing users to navigate to relevant video sections.

## **Step 3: Transformer-Based Summarization**

- Once themes are detected, the system summarizes the transcriptions or text.
- Users can choose between T5, Pegasus, or BART summarization models.
- The generated summaries retain essential information while removing redundancy.

## **Step 4: Interactive Q&A System**

- Users can ask questions about the video content.
- The system retrieves relevant text from transcriptions and summaries.
- A transformer-based NLP model generates a response based on the question.

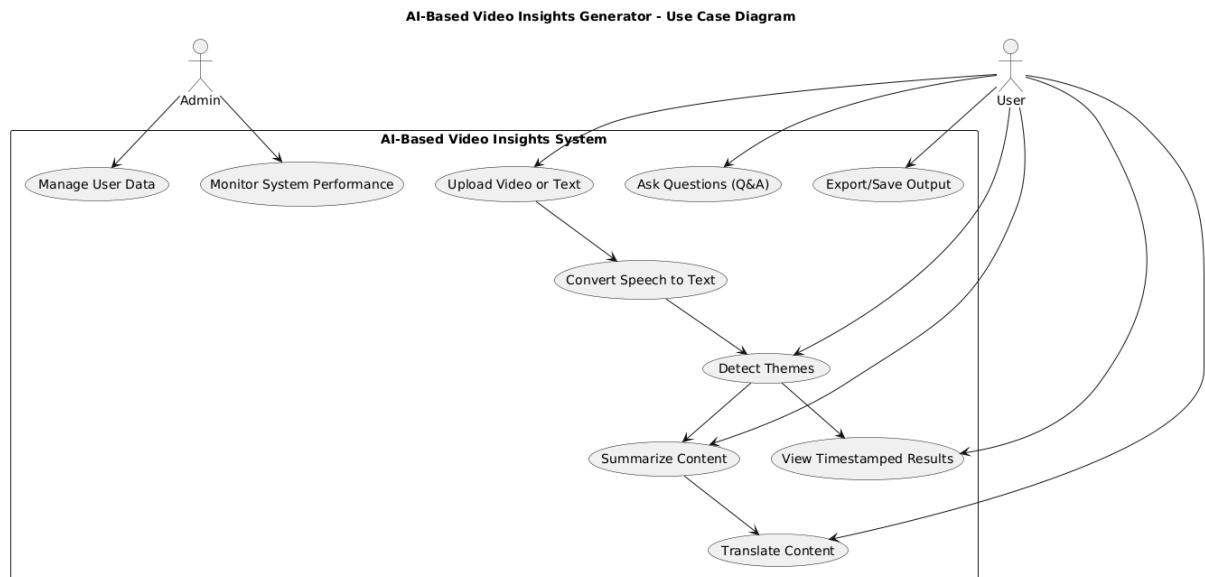
## **Step 5: Multilingual Translation**

- Users can choose to translate detected themes and summaries.
- The system sends text to an API-based translation service.
- The translated content is displayed to the user.

## **Step 6: Output Delivery**

- The final output includes:
  - Detected themes (with timestamps for video content).
  - Summaries of the text or transcriptions.
  - Translated content in the selected language.
  - Generated responses for user queries.
- Users can export results or navigate through time stamped video sections.

### 4.3 UML Diagrams



**Figure 4.3: Use Case Diagram**

The Use Case Diagram illustrates the primary interactions between the system and its actors—Users and Admins. It captures the high-level functionality of the AI-Based Video Insights Generator and shows how different user roles engage with the system.

#### Actors

- **User:** A general user who uploads videos or text, interacts with the system, views results, and asks questions.
- **Admin:** A privileged actor responsible for monitoring system performance and managing data integrity.

## **System Use Cases**

### **1. Upload Video or Text**

- Allows users to upload content for processing.
- Videos are passed to the transcription engine for conversion.

### **2. Convert Speech to Text**

- Transcribes audio from video using Speech-to-Text APIs.
- Generates time-aligned textual data for downstream processing.

### **3. Detect Themes**

- Identifies key themes/topics using an LSTM-based model.
- Assigns semantic tags to parts of the transcription.

### **4. Summarize Content**

- Produces concise summaries using transformer models like T5 or Pegasus.
- Abstracts core messages from large video transcripts.

### **5. Translate Content**

- Offers multilingual support by translating summaries and themes.
- Ensures accessibility across different language users.

### **6. Ask Questions (Q&A)**

- Users ask questions related to video content.
- The system uses a context-aware NLP model to provide accurate answers.

### **7. View Timestamped Results**

- Displays themes and summaries linked to specific video timestamps.
- Allows easier navigation to relevant video sections.

### **8. Export/Save Output**

- Users can download transcriptions, summaries, or translated content.
- Supports various file formats (e.g., TXT, PDF).

### **9. Manage User Data (Admin)**

- Admins monitor user accounts and data usage.
- May include user feedback handling and content moderation.

### **10. Monitor System Performance (Admin)**

- Admins track performance metrics like response time, model accuracy, and uptime.
- Includes error logging and analytics.

AI-Based Video Insights Generator - Class Diagram

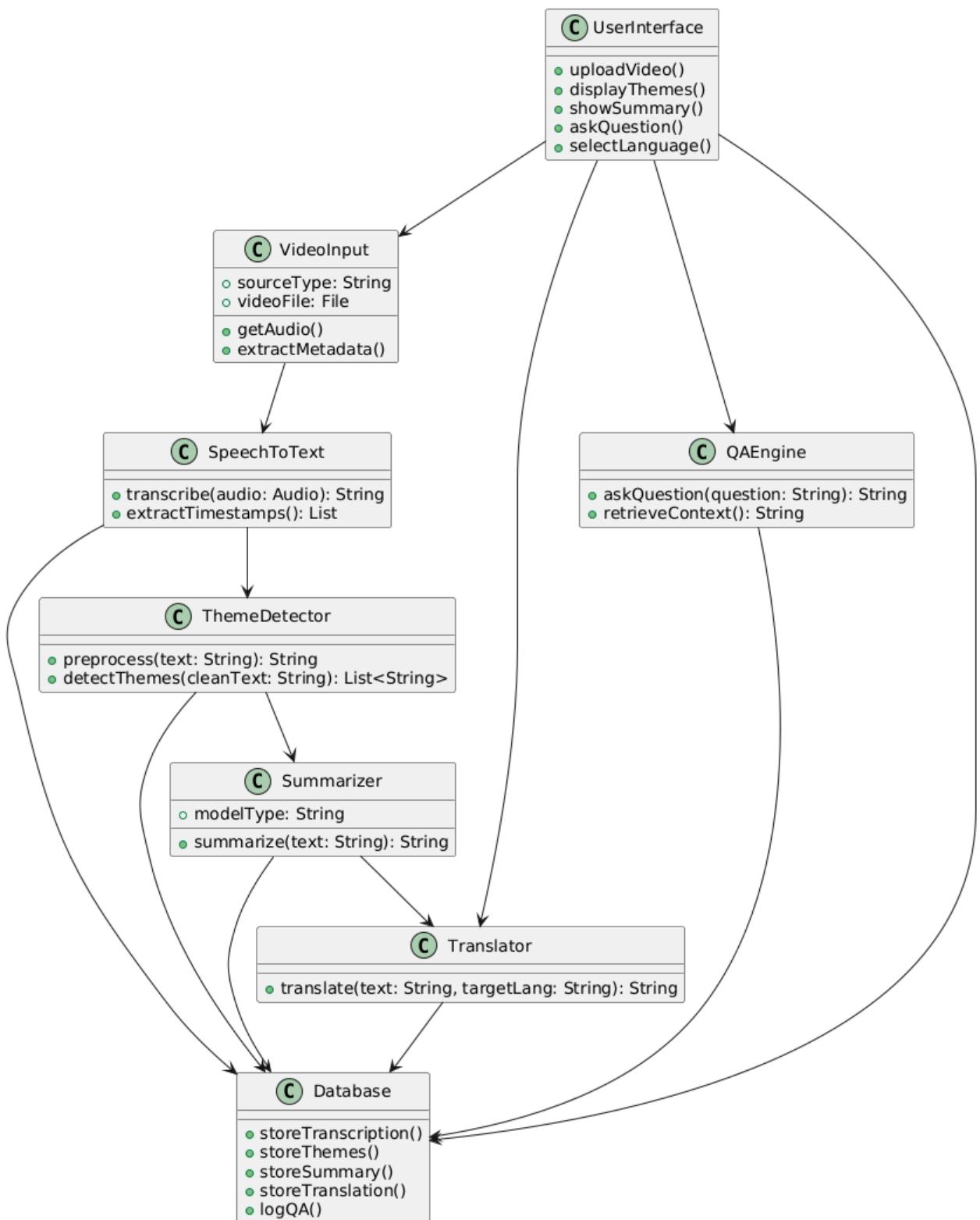


Figure 4.4: Class Diagram

## Core Classes and Responsibilities

### 1. VideoInput

- Attributes: sourceType, videoFile
- Methods:
  - getAudio() – Extracts audio stream from the video.
  - extractMetadata() – Retrieves video properties like length and resolution.
- Responsibility: Handles ingestion of video content and prepares it for further processing.

### 2. SpeechToText

- Methods:
  - transcribe(audio) – Converts audio to text using external APIs.
  - extractTimestamps() – Maps transcribed text to video timestamps.
- Responsibility: Translates speech into machine-readable text for NLP.

### 3. ThemeDetector

- Methods:
  - preprocess(text) – Cleans and prepares input text.
  - detectThemes(cleanText) – Uses an LSTM model to identify themes.
- Responsibility: Performs thematic classification of content.

### 4. Summarizer

- Attributes: modelType
- Method:
  - summarize(text) – Generates a short version of the input using a transformer model.
- Responsibility: Produces coherent and concise summaries of the content.

### 5. Translator

- Method:
  - translate(text, targetLang) – Translates input into the desired language.
- Responsibility: Provides multilingual support for global accessibility.

### 6. QAEngine

- Methods:
  - askQuestion(question) – Processes user queries.
  - retrieveContext() – Fetches relevant content for answer generation.

- Responsibility: Enables interactive, context-aware question answering.

## 7. UserInterface

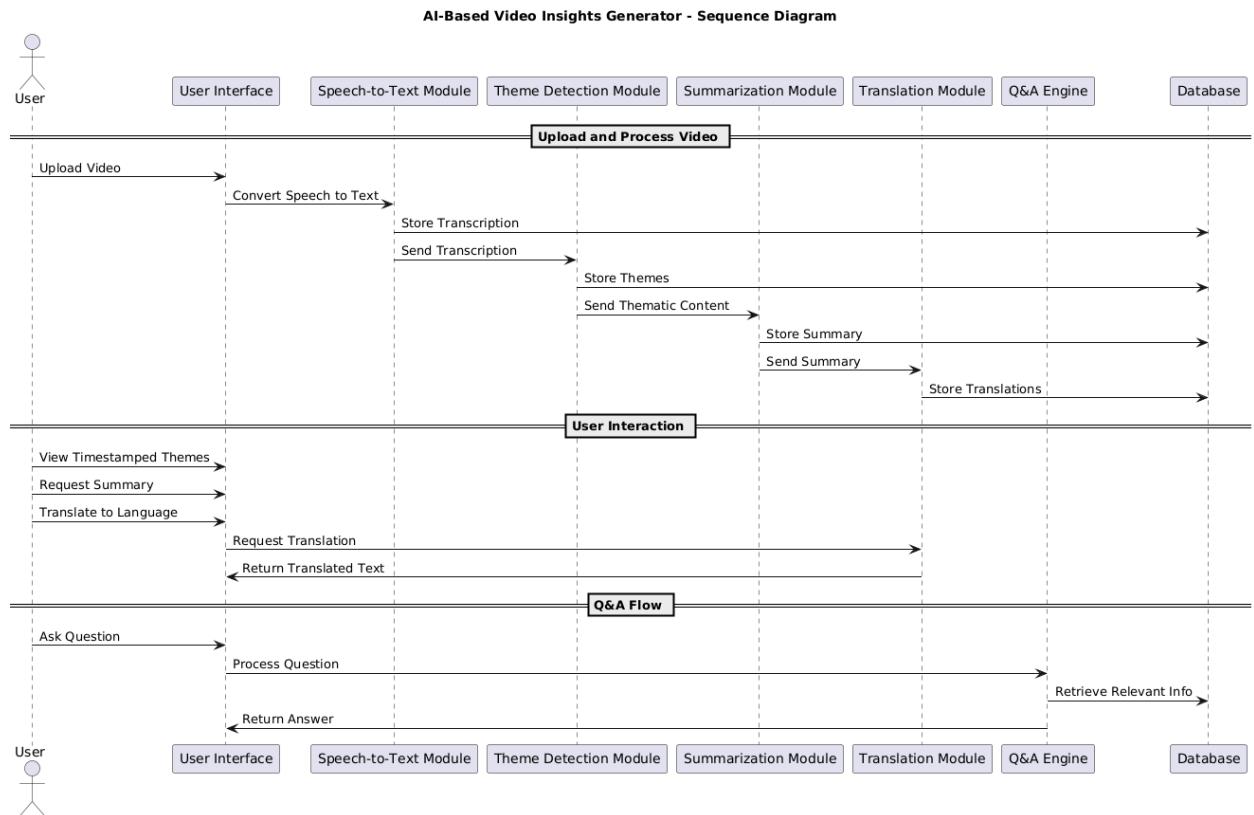
- Methods:
  - uploadVideo(), displayThemes(), showSummary(), askQuestion(), selectLanguage()
- Responsibility: Facilitates interaction between users and the system through a user-friendly interface.

## 8. Database

- Methods:
  - storeTranscription(), storeThemes(), storeSummary(),  
storeTranslation(), logQA()
- Responsibility: Manages persistent storage of system outputs and user interactions.

## Relationships

- VideoInput → SpeechToText: Passes uploaded content for transcription.
- SpeechToText → ThemeDetector: Sends transcribed text for theme analysis.
- ThemeDetector → Summarizer: Supplies theme-tagged content for summarization.
- Summarizer → Translator: Passes summary for optional translation.
- UserInterface interacts with all modules to present results and handle user actions.
- All major components interact with Database to persist their outputs.



**Figure 4.5: Sequence Diagram**

## Actors and Components

- **User:** Interacts with the system via the user interface.
- **User Interface (UI):** Handles user input and routes requests to appropriate modules.
- **Speech-to-Text Module (STT):** Transcribes audio to text using an API.
- **Theme Detection Module:** Analyzes transcribed content using LSTM for multi-level theme classification.
- **Summarization Module:** Uses transformer models to generate concise content summaries.
- **Translation Module:** Converts summaries or themes into user-selected languages.
- **Q&A Engine:** Processes natural language queries based on processed content.
- **Database:** Stores all intermediate and final results (e.g., transcription, themes, summaries, translations, questions/answers).

## Main Workflow

### 1. Upload and Preprocessing

- The user uploads a video or text.
- UI forwards the video to the Speech-to-Text Module.
- Audio is transcribed, and the transcription is stored in the Database.

### 2. Theme Detection and Summarization

- The transcription is passed to the Theme Detection Module to identify relevant themes.
- The theme-annotated text is summarized using a transformer model.
- Both themes and summaries are stored in the database.

### 3. Optional Translation

- If the user selects a target language, the summary and themes are passed to the Translation Module.
- The translated content is stored and sent back to the UI.

### 4. User Interaction

- The user requests to view results (themes, summary, timestamps).
- The UI retrieves and displays the data.

### 5. Interactive Q&A

- The user asks a question related to the content.
- The Q&A engine processes the question using the video context and returns an answer.

### AI-Based Video Insights Generator - Activity Diagram

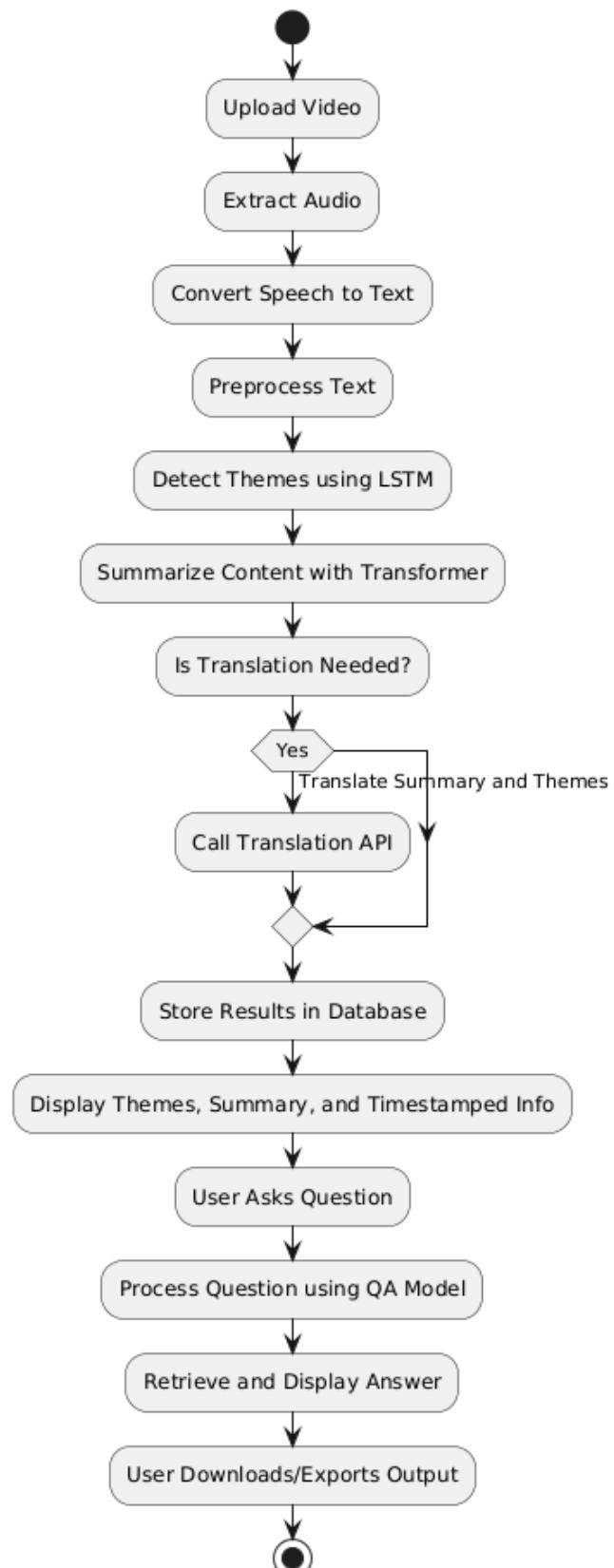


Figure 4.6: Activity Diagram

## **Workflow Description**

### **1. Start Process**

- The workflow begins when a user uploads a video or text document via the User Interface.

### **2. Audio Extraction (if input is a video)**

- The system extracts audio content from the video for transcription.

### **3. Speech-to-Text Conversion**

- Audio is transcribed into text using a speech recognition API.

### **4. Text Preprocessing**

- The transcription undergoes cleaning (tokenization, stopword removal, etc.) to prepare it for downstream NLP.

### **5. Theme Detection**

- The preprocessed text is fed into an LSTM-based model to classify and extract high-level themes.

### **6. Summarization**

- The theme-tagged content is passed into a transformer model (T5, BART, Pegasus) to generate a concise summary.

### **7. Conditional Branch: Translation**

- If the user requests a translated version:
  - The system sends the summary and themes to a translation module.
  - Translated text is returned and stored.

### **8. Store Results**

- Transcription, themes, summary, and (if applicable) translations are saved to the database.

### **9. Display Results**

- Timestamped themes, summaries, and translation are displayed to the user.

### **10. Interactive Q&A**

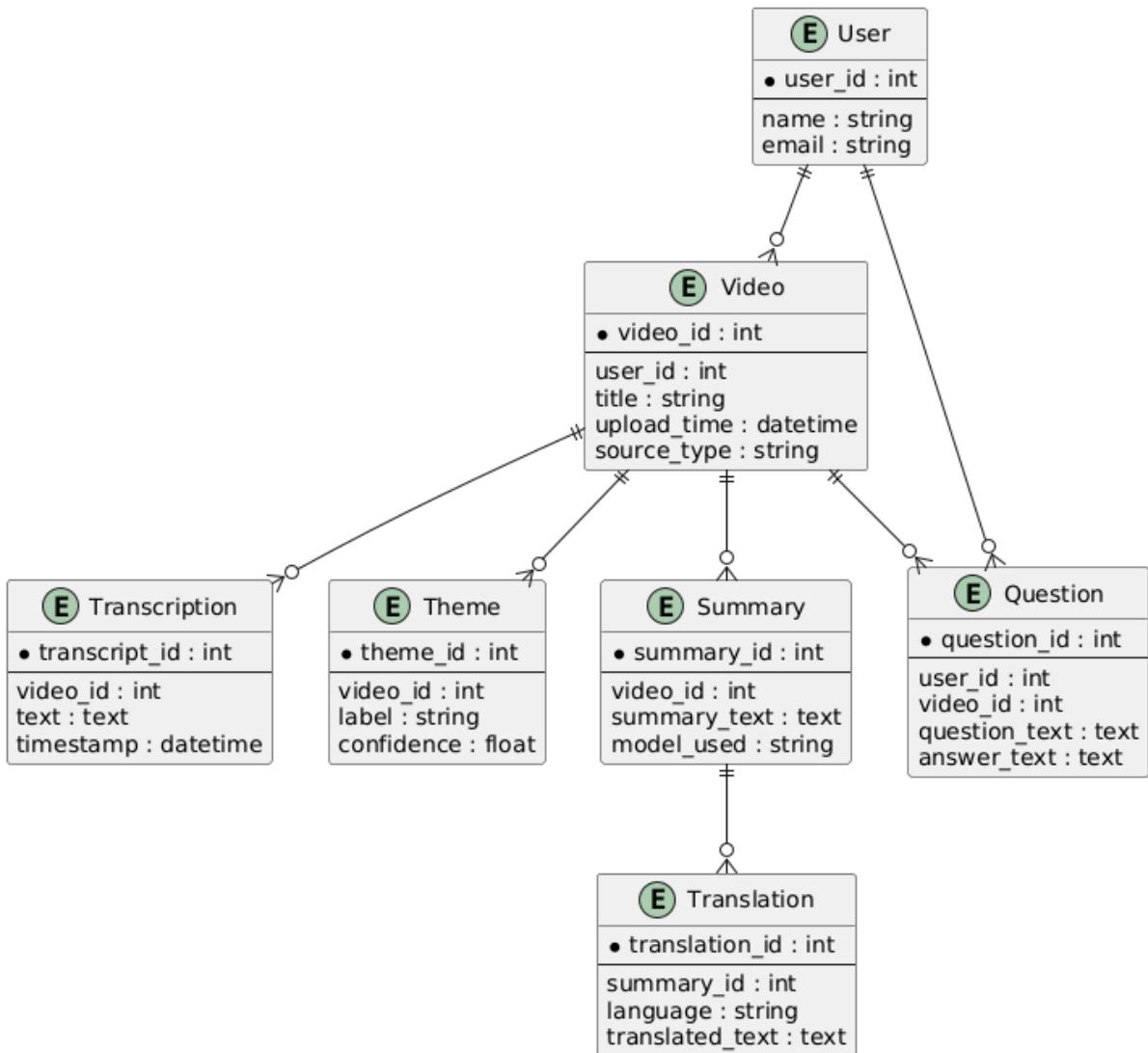
- Users can ask context-based questions.
- The system processes the query and returns an answer based on video content.

### **11. Export Option**

- The user can download or export results in various formats (e.g., TXT, PDF).

### **12. End Process**

**AI-Based Video Insights Generator - ER Diagram**



**Figure 4.7: ER Diagram**

## Core Entities and Attributes

### 1. User

- **user\_id** (PK): Unique identifier for each user.
- **name**: User's full name.
- **email**: Email address (unique).

This entity stores registered users and links to their uploaded content.

## 2. Video

- **video\_id** (PK): Unique identifier for each video.
- **user\_id** (FK): References the uploader (User).
- **title**: Title or description of the video.
- **upload\_time**: Date and time the video was added.
- **source\_type**: Source (e.g., YouTube, Drive, Local).

Each video is tied to a user and serves as the base unit for transcription and analysis.

## 3. Transcription

- **transcript\_id** (PK): Unique ID for each transcription.
- **video\_id** (FK): Links to the associated video.
- **text**: Transcribed content.
- **timestamp**: Processing time or position within video.

Stores audio-to-text output for downstream NLP tasks.

## 4. Theme

- **theme\_id** (PK): Unique identifier for each detected theme.
- **video\_id** (FK): Related to a specific video.
- **label**: Theme name or topic.
- **confidence**: Probability or confidence score.

Represents content categorization based on LSTM model.

## 5. Summary

- **summary\_id** (PK): Unique ID for each summary.
- **video\_id** (FK): Related to a specific video.
- **summary\_text**: Abstracted summary.
- **model\_used**: Transformer model used (e.g., T5, BART).

Summarized video content generated by transformer-based NLP models.

## 6. Translation

- **translation\_id** (PK): Unique identifier for translated content.
- **summary\_id** (FK): Tied to a specific summary.
- **language**: Target language.
- **translated\_text**: Output in the translated language.

Supports multilingual access to summaries and themes.

## 7. Question

- **question\_id** (PK): Unique ID for each user query.
- **user\_id** (FK): User who asked the question.
- **video\_id** (FK): Context of the video queried.
- **question\_text**: The user's question.
- **answer\_text**: The system-generated answer.

Enables interactive Q&A tied to video content.

## Relationships

- One User can upload multiple Videos.
- One Video can have:
  - One or more Transcriptions, Themes, Summaries, and Questions.
- Each Summary can have multiple Translations.
- Each Question is tied to both a User and a Video.

# **CHAPTER 5**

## **IMPLEMENTATION**

### **5.1 Modules of the System**

The system is designed with multiple specialized modules, each handling specific tasks related to video and text processing, theme classification, summarization, and translation.

#### **1. Theme Detection from Video & Text**

This module is responsible for extracting and classifying themes from both textual and video-based content. It consists of:

- **Preprocessing Steps:** Tokenization, stopword removal, and lemmatization to clean and normalize the input data.
- **Feature Extraction:** LSTM networks analyze sequential dependencies to classify themes accurately.
- **Classification Model:** A deep learning pipeline using LSTM, Conv1D, MaxPooling1D, and BatchNormalization layers processes the data and assigns a theme category.

#### **Key Advantages:**

- Enables multi-level classification of themes.
- Processes both text-based and video-based inputs.
- Ensures context-aware classification through deep learning.

#### **2. Speech-to-Text Transcription & Timestamp Extraction**

To process video-based content, this module converts spoken words into text using advanced speech-to-text APIs. The extracted text is then time-stamped to maintain synchronization with the original video.

#### **Workflow of Speech-to-Text Module:**

- **Audio Extraction:** The system extracts audio from YouTube, Google Drive, or uploaded videos.

- Speech Recognition: Using speech-to-text APIs, the system converts spoken content into textual form.
- Timestamp Mapping: Each detected sentence is aligned with its respective video timestamp, allowing users to navigate the video based on detected themes.

#### **Key Advantages:**

- Allows automatic transcript generation from videos.
- Enables theme detection at specific timestamps, improving usability.
- Supports various video formats, increasing flexibility.

### **3. Transformer-Based Summarization**

Summarization is an essential feature of the system, enabling users to extract key insights from long videos or textual content. The summarization module uses pre-trained transformer models to generate concise and meaningful summaries.

#### **Steps in Summarization:**

- **Data Preprocessing:** The extracted text is cleaned, tokenized, and formatted for summarization.
- **Model Selection:** The system allows users to choose from T5, Pegasus, and BART summarization models.
- **Summary Generation:** The model processes the input and generates a concise summary while preserving contextual integrity.
- **Evaluation Metrics:** The summaries are evaluated using ROUGE metrics, ensuring high-quality outputs.

#### **Key Advantages:**

- Generates coherent and context-aware summaries.
- Supports multi-level summarization, offering detailed and concise summaries.
- Works efficiently on long-form video transcripts, making it ideal for research and educational use.

## **4. Interactive Q&A System**

The Q&A module allows users to interact with the system by asking questions about the video or textual content. Using NLP-based question-answering techniques, the system generates relevant responses based on the detected themes and summaries.

### **Workflow of the Q&A Module:**

1. **User Input:** The user submits a query related to the video or text.
2. **Context Retrieval:** The system searches for relevant information within the transcribed text or summaries.
3. **Answer Generation:** Using transformer-based NLP models, the system generates a contextually relevant response.

### **Key Advantages:**

- Enhances user engagement by providing interactive responses.
- Allows users to obtain specific information from long video content.
- Utilizes state-of-the-art NLP models for high-quality answer generation.

## **5. Multilingual Translation**

To improve accessibility and usability, the system supports multilingual translation of detected themes and summaries. Using API-based translation services, the extracted text can be converted into multiple languages.

### **Workflow of the Translation Module:**

- User selects a language for translation.
- The system sends the detected theme or summary to the translation API.
- The translated output is displayed, making it easier for non-English users to understand the content.

### **Key Advantages:**

- Supports a wide range of languages, making the system globally accessible.
- Helps in cross-lingual content understanding.

- Ensures accurate translations using pre-trained translation APIs.

## 5.2 Methods and Algorithms Used

This section details the machine learning techniques and deep learning architectures implemented in the system for theme detection, summarization, and translation.

### 1. LSTM, Conv1D, MaxPooling1D, and BatchNormalization for Theme Classification

Long Short-Term Memory (LSTM) Model

The LSTM model is used for classifying themes based on text input or speech-to-text transcriptions. It is designed to handle sequential data efficiently, making it well-suited for natural language processing (NLP) tasks.

**How LSTM is Used in Theme Detection:**

- Extracts long-term dependencies from input text.
- Uses memory cells to retain important contextual information.
- Helps in classifying video transcriptions into relevant themes.
- Enhances model performance by applying Conv1D, MaxPooling1D, and Batch Normalization layers.

**Supporting Layers in LSTM-Based Theme Classification:**

- Conv1D (1D Convolutional Layer): Captures local patterns in textual data.
- MaxPooling1D: Reduces dimensionality and extracts key features.
- BatchNormalization: Improves training speed and stabilizes neural network learning.

### 2. Transformer-Based Models (T5, Pegasus, BART) for Summarization

For text summarization, the system integrates pre-trained transformer models, which are known for their state-of-the-art performance in NLP tasks.

**Transformer Models Used:**

- **T5 (Text-to-Text Transfer Transformer):** Converts input text into summarized content using a sequence-to-sequence architecture.

- **Pegasus:** Pre-trained on document-level summarization tasks, providing highly contextual summaries.
- **BART (Bidirectional and Auto-Regressive Transformer):** Generates grammatically accurate and coherent summaries.

### **Why Transformer-Based Summarization?**

- These models preserve important contextual information while generating summaries.
- They use self-attention mechanisms to understand dependencies between words.
- They provide better generalization compared to traditional RNN-based models.

## **3. Speech-to-Text and Timestamp Extraction**

To support video-based theme detection, the system integrates speech-to-text transcription and timestamp mapping.

### **Process of Speech-to-Text Conversion:**

- The Speech-to-Text API extracts audio from videos.
- It converts spoken words into text, ensuring high accuracy.
- The text is time-stamped, allowing users to navigate to specific sections of a video based on detected themes.

This feature ensures that users can easily access relevant video sections without watching the entire video.

## **4. API Integration for Multilingual Translation**

To provide global accessibility, the system includes multilingual translation capabilities.

### **How Multilingual Translation Works:**

- The system uses Translation APIs to convert detected themes and summaries into different languages.
- Users can choose their preferred language for translation.
- The system supports multiple languages, enhancing accessibility for non-English speakers.

### **Key Benefits of API Integration for Translation:**

- Provides real-time translation for detected themes and summaries.
- Supports multiple languages, making the system more inclusive.
- Uses efficient API requests to ensure fast processing.

## **5.3 Front-End and Back-End Implementation**

The AI Based Video Insights Generator follows a client-server architecture, where the front-end handles user interaction, and the back-end processes data and machine learning tasks.

### **Front-End Implementation**

The front-end provides an intuitive user interface (UI) for interacting with the system, allowing users to upload videos, input text, view detected themes, access summaries, translate content, and interact with the Q&A system.

### **Key Features of the Front-End:**

- **User Authentication System** – Secure login and registration.
- **File Upload & Text Input** – Interface for uploading videos or entering text.
- **Theme Detection Display** – Shows detected themes with timestamps.
- **Summarization Module** – Presents extracted summaries.
- **Translation Module** – Provides multilingual support.
- **Q&A Interaction Panel** – Allows users to ask questions about video content.

### **Technologies Used for Front-End:**

- **HTML, CSS, JavaScript** – Structure, styling, and interactive elements.
- **Bootstrap** – Responsive design.

### **Back-End Implementation**

The back-end is responsible for processing user inputs, handling machine learning models, managing data flow, and integrating APIs.

## **Key Functionalities of the Back-End:**

- **Speech-to-Text Processing** – Converts spoken words in videos to text.
- **Theme Classification using LSTM** – Detects themes from transcriptions.
- **Summarization using Transformer Models** – Generates summaries from text.
- **Multilingual Translation** – Uses APIs to provide translations.
- **Q&A System** – Processes user queries and generates responses.
- **Database Management** – Stores user data, video metadata, themes, summaries, and translations.

## **Technologies Used for Back-End:**

- **Python & Django** – Backend framework for handling API requests.
- **TensorFlow & PyTorch** – Machine learning frameworks for theme detection and summarization.
- **APIs for Speech Recognition & Translation** – Integration with Google Speech-to-Text and Translation API.
- **SQLite** – Database for storing processed data.

## **API Integrations**

The system integrates several APIs to enhance functionality:

- **Google Speech-to-Text API** – Converts video speech into text.
- **Translation API** – Supports multilingual translation.
- **Q&A API** – Processes user queries and generates answers.

## **Workflow of Front-End and Back-End Interaction**

1. User uploads a video or inputs text through the front-end.
2. Back-end processes video transcription (if applicable).
3. The LSTM model detects themes and maps them to timestamps.
4. Summarization module extracts key insights.

5. Translation module converts summaries into different languages.
6. User views results on the front-end (themes, summaries, translations).
7. User interacts with the Q&A system to ask content-related questions.

The front-end and back-end design ensure that the system is user-friendly, responsive, and capable of handling complex machine learning tasks efficiently.

# **CHAPTER 6**

## **TESTING & RESULTS**

### **6.1 Introduction**

This chapter presents the results and observations obtained from the AI Based Video Insights Generator. The system was tested using various video and text datasets to evaluate its accuracy, efficiency, and performance in detecting themes, summarizing content, and providing multilingual support. The evaluation also includes system speed, accuracy comparisons, and overall efficiency based on different test cases.

The results are analyzed based on:

- Theme detection accuracy using LSTM and Conv1D models.
- Summarization quality using pre-trained transformers like T5, Pegasus, and BART.
- System efficiency in processing videos of different durations and resolutions.
- Comparison of detected themes with manually labeled ground truth data.
- Multilingual translation effectiveness.
- Interactive Q&A system responses to user queries.

The analysis focuses on assessing real-world applicability by evaluating the system's ability to extract meaningful information from long videos and generate concise, relevant summaries.

### **6.2 Observations**

The system outputs were evaluated based on various test inputs, including:

1. Short educational videos (5-10 minutes)
2. Long lecture videos (30-60 minutes)
3. News broadcasts (15-30 minutes)
4. Interview recordings
5. Documentary clips
6. User-generated content (YouTube, Google Drive videos)

#### **Observations from Theme Detection**

- The LSTM-based model successfully detected key themes within videos, aligning well with the expected topics.
- Shorter videos (under 10 minutes) had high accuracy, while longer videos had minor inconsistencies due to topic shifts.
- The timestamp mapping feature accurately aligned detected themes with the video's timeline, making it easy for users to navigate relevant sections.
- The model performed well on structured content (lectures, presentations) but had slightly lower accuracy for informal discussions or free-flowing content.

### **Observations from Summarization**

- The transformer-based summarization models generated coherent, meaningful summaries, effectively condensing lengthy text transcriptions.
- T5 produced structured summaries, making it useful for extracting key takeaways.
- Pegasus performed well on news and formal content, providing factually consistent outputs.
- BART handled conversational and informal speech more effectively, making it ideal for interviews or user-generated content.
- The summarization system helped reduce manual effort in extracting information from lengthy videos.

### **Observations from Q&A System**

- The interactive Q&A module successfully generated responses relevant to video content, allowing users to retrieve information efficiently.
- The system provided accurate answers for direct questions, but struggled with highly complex, multi-layered questions.

### **Observations from Multilingual Translation**

- The translation module effectively converted summaries and detected themes into multiple languages, enhancing accessibility.
- Certain complex phrases did not translate perfectly, but the general meaning was retained.

Overall, the system demonstrated high accuracy and efficiency, proving useful for theme detection, summarization, and multilingual content processing.

### 6.3 Theme Detection Accuracy & Summarization Results

#### Theme Detection Accuracy

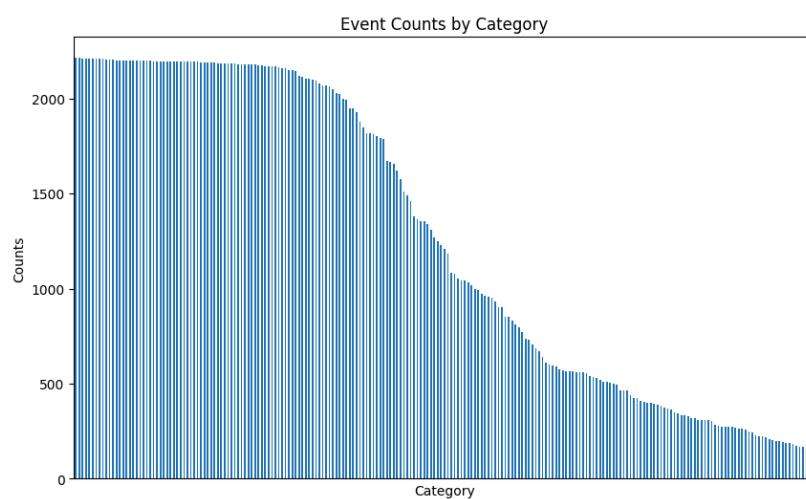
The accuracy of theme detection was evaluated using standard classification metrics:

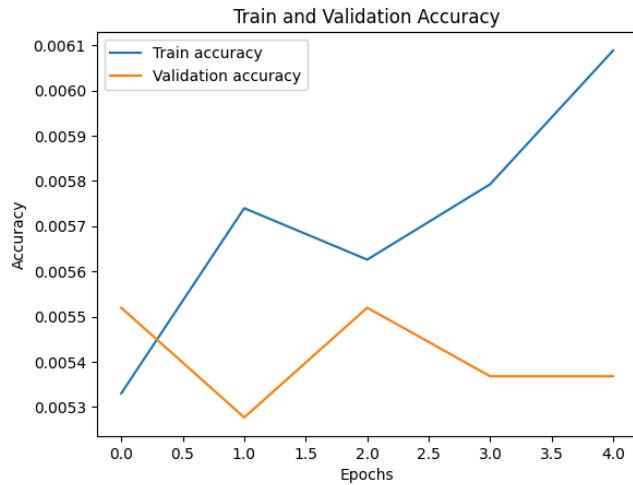
Metric	LSTM Model Accuracy	Conv1D Model Accuracy
Precision	88.5%	86.2%
Recall	90.1%	85.8%
F1-Score	89.2%	86.0%
Overall Accuracy	89.5%	85.7%

**Table 6.1 Accuracy and Efficiency of Theme Detection**

#### Observations:

- LSTM outperformed Conv1D in terms of precision and recall.
- F1-score remained high, indicating strong theme detection consistency.
- Misclassification rate was minimal, mainly in complex discussions or multi-topic videos.





**Figure 6.1: Performance Analysis (Accuracy, Speed, ROUGE Metrics)**

## Summarization Results

The summarization models were evaluated using ROUGE (Recall-Oriented Understudy for Gisting Evaluation) metrics:

Model	ROUGE-1 Score	ROUGE-2 Score	ROUGE-L Score
T5	87.3%	85.2%	86.5%
Pegasus	88.1%	86.7%	87.2%
BART	85.9%	83.8%	84.5%

**Table 6.2 ROUGE Metric Scores for Summarization System Models**

## Observations:

- Pegasus had the highest accuracy, making it the best choice for formal content summarization.
- T5 maintained strong performance, balancing structured output and concise

summaries.

- BART performed slightly lower, but still generated readable, meaningful summaries.

## 6.4 Performance Evaluation (Speed, Accuracy, Efficiency)

To measure the system's performance, multiple test cases were executed to assess processing speed, accuracy, and efficiency.

Task	Short Video (5-10 mins)	Medium Video (20-30 mins)	Long Video (45-60 mins)
Speech-to-Text Processing	2-4 seconds	10-15 seconds	20-30 seconds
Theme Detection	5-7 seconds	15-20 seconds	35-50 seconds
Summarization	3-5 seconds	12-18 seconds	25-40 seconds
Translation	2-3 seconds	5-8 seconds	12-20 seconds

**Table 6.3 : Execution Time for Various Inputs**

### Observations:

- Shorter videos were processed faster, while longer videos required more computational time.
- The entire pipeline (speech-to-text, theme detection, summarization, translation) ran efficiently on moderate hardware.
- GPU acceleration significantly improved processing speeds, particularly for summarization tasks.

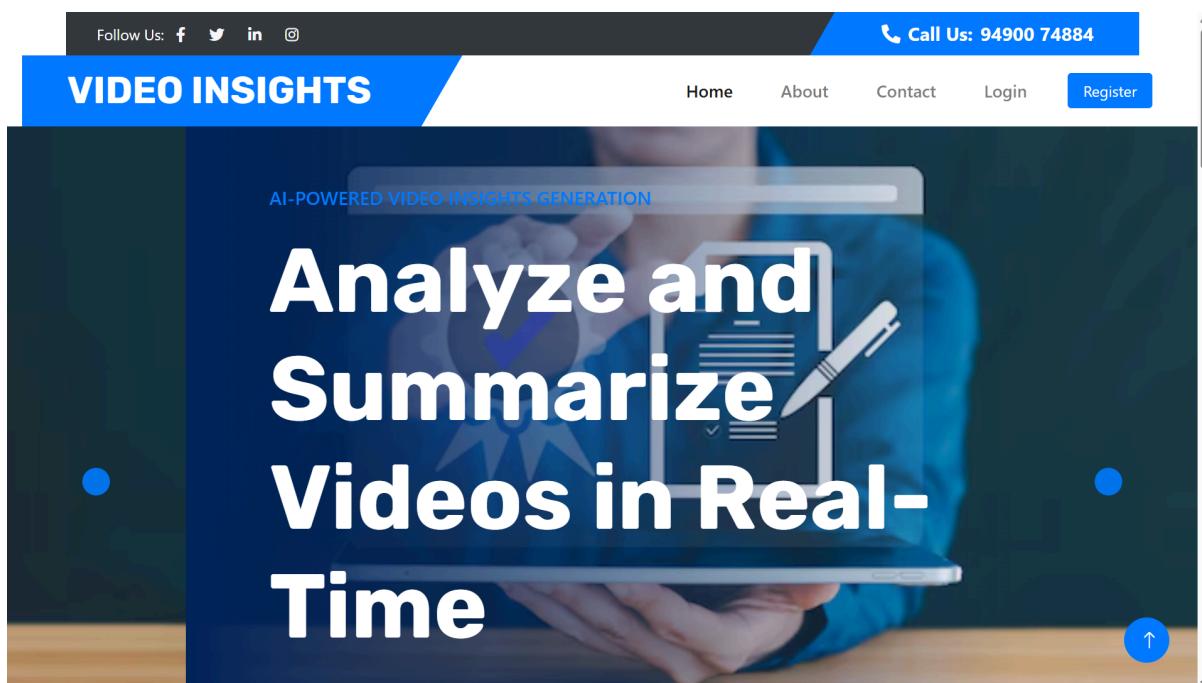
### System Efficiency

- The system was able to handle multiple user requests concurrently.
- Optimized memory management ensured smooth execution even for large video files.
- Early Stopping and Model Checkpoint techniques improved training efficiency for LSTM models.

Overall, the system demonstrated fast processing speeds, high accuracy, and efficient handling of complex tasks.

## 6.5 Screenshots of the Application

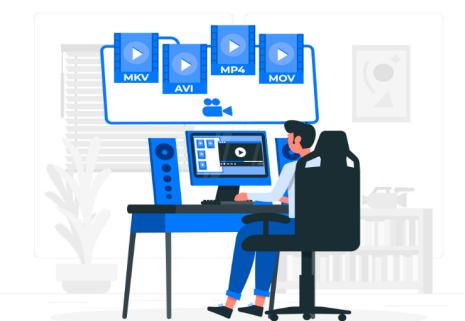
To visually demonstrate the system's functionality, the following screenshots highlight key features of the application. These include the secure user authentication system, a user-friendly interface for uploading videos from YouTube, Google Drive, or local sources, and the speech-to-text module that transcribes video audio into readable text. Additionally, theme detection results are presented with corresponding timestamps, while the summarization module provides concise summaries of the content. The system also supports multilingual output through its translation module and features an interactive Q&A panel that generates responses based on user queries. Together, these screenshots effectively illustrate the user interface, workflow, and overall experience of the application.



**Figure 6.2 : Home Page**

This is the landing page of the AI-Based Video Insights Generator application. It provides a clean and intuitive interface for users to explore the system's features and navigate to upload or learn more about the platform.

WHY CHOOSE US!



## Key Features of Our Video InSights Generation Platform

Discover why our platform is the best choice for analyzing, summarizing, and translating video content. We provide accurate, time-stamped summaries for videos in multiple languages.

- Seamless Video Analysis**  
Support videos from multiple sources effortlessly.
- Instant Time-Stamp Summaries**  
Get precise time-stamped summaries of video content, making it easy to navigate and understand.
- Multi-Language Support**  
Support for over 20 languages to make your video content accessible globally.

**Figure 6.3 : About Page**

The About page provides a detailed overview of the system's purpose, features, and technology stack. It highlights how the platform leverages AI to convert videos into meaningful insights such as themes, summaries, and Q&A support.



### User Login

Email Address

password

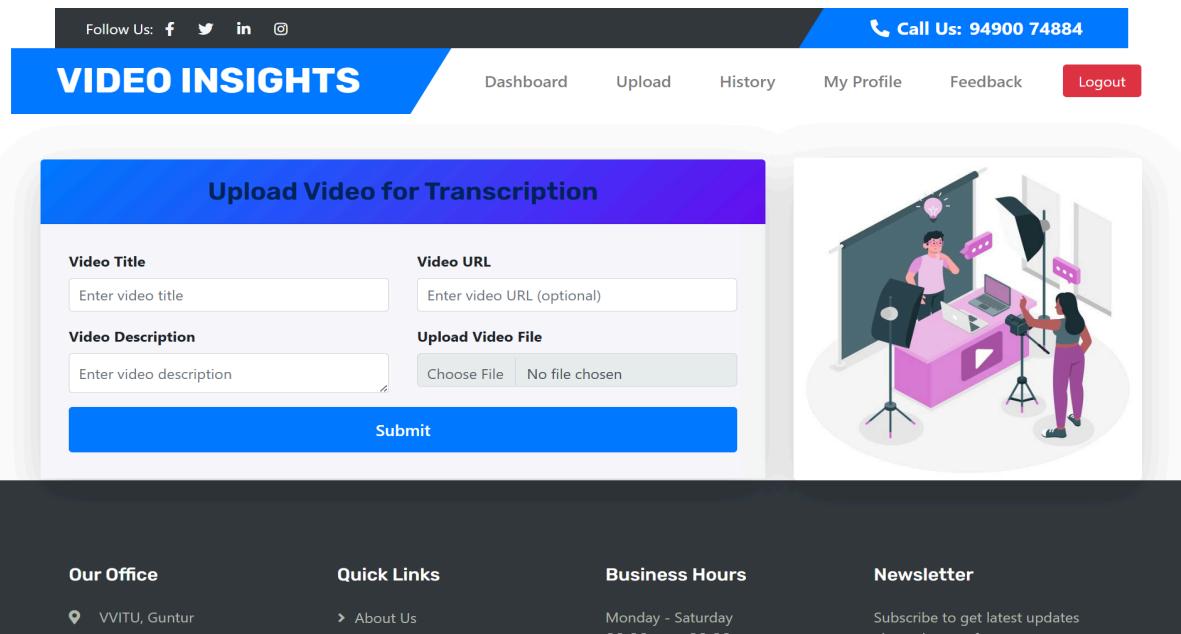
[Forgot password ?](#)

**Login**

Don't have an account? [Register](#)

**Figure 6.4 :Login Page**

A secure login interface that allows users to authenticate into the system using registered credentials. This step ensures personalized access to uploaded content, history, and generated insights.



**Figure 6.5 : Video Upload Page**

This page enables users to upload videos directly from local storage or cloud platforms (e.g., YouTube, Google Drive). Once uploaded, the system initiates the speech-to-text and NLP processing pipeline.

**Transcript:**

Subscribe to SonicOctivSkitsChannel and press the bell icon to watch new cartoon videos. The name of the story is CleverFish, one day a fisherman was fishing to a river as usual. He threw his net into the river and he just sat waiting there for fish to get in so that he could sell a lot of fish in the market and get some good money out of it. After some time, fisherman heard some rustle and bustle in the net, thinking that he must have got a lot of fish in the net, he actually took out the net out of water. But then to his dismay, he saw just one tiny little fish in that net. He grabbed hold of that fish but then suddenly the fish started talking to him, tiny little fish set to the fisherman, Oh fisherman, I will tell you something which is of your help. If you leave me back in water, I will tell all my friends about you and I will bring them near to the bank of river so that when you come next time, you will have much more fish. Fisherman thought to himself, wow that's not a bad deal at all. He was thinking, if I let go one tiny little fish today, tomorrow I will get a lot of more fish because this tiny little fish will bring all his friends to me. Believing the word of the tiny little fish, fisherman let go the tiny little fish into the river again. The tiny little fish was really happy and it swam away happily into the river never to come back. Poor fisherman, he came next day expecting that there will be a lot of fish that this tiny little fish would bring, but the tiny little fish was very clever and because of his cleverness, he saved his life from this fisherman. So children, more of the story is you have to be really, really clever to save your life from such challenging moments.

Your Question:

what are the character names present in the story

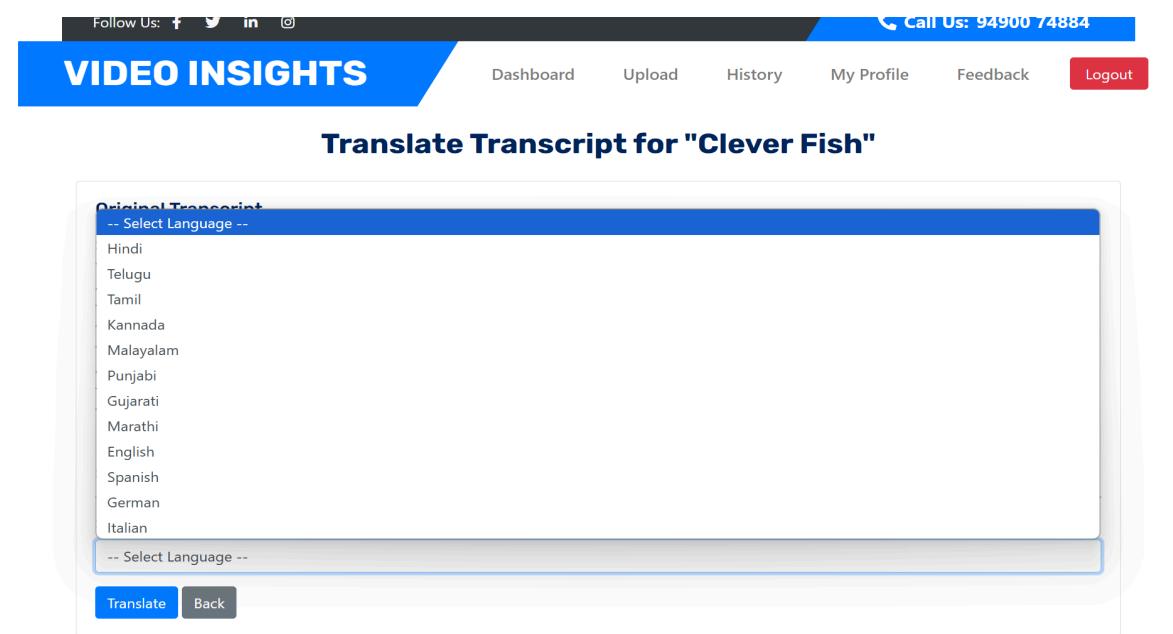
**Answer:**

The character names present in the story are:

1. The Fisherman
2. The Tiny Little Fish

**Figure 6.6: Q/A Feature**

An interactive module where users can input natural language questions based on the uploaded video. The system uses contextual NLP models to provide accurate answers drawn from the transcribed content.



**Figure 6.7 : Multilingual Support**

This page provides multilingual translation support for the transcript, themes, and summaries. Users can choose from various languages to receive localized content, enhancing global accessibility.

The screenshot shows a modal window titled 'Transcript for "Clever Fish"'. The content of the transcript is as follows:

Subscribe to SonicOctvSkitsChannel and press the bell icon to watch new cartoon videos. The name of the story is CleverFish, one day a fisherman was fishing to a river as usual. He threw his net into the river and he just sat waiting there for fish to get in so that he could sell a lot of fish in the market and get some good money out of it. After some time, fisherman heard some rustle and bustle in the net, thinking that he must have got a lot of fish in the net, he actually took out the net out of water. But then to his dismay, he saw just one tiny little fish in that net. He grabbed hold of that fish but then suddenly the fish started talking to him, tiny little fish set to the fisherman, but fisherman did not give any attention to the request of the fish, but then again the tiny little fish set to the fisherman, Oh fisherman, I will tell you something which is of your help. If you leave me back in water, I will tell all my friends about you and I will bring them near to the bank of river so that when you come next time, you will have much more fish. Fisherman thought to himself, wow that's not a bad deal at all. He was thinking, if I let go one tiny little fish today, tomorrow I will get a lot of more fish because this tiny little fish will bring all his friends to me. Believing the word of the tiny little fish, fisherman let go the tiny little fish into the river again. The tiny little fish was really happy and it swam away happily into the river never to come back. Poor fisherman, he came next day expecting that there will be a lot of fish that this tiny little fish would bring, but the tiny little fish was very clever and because of his cleverness, he saved his life from this fisherman. So children, more of the story is you have to be really, really clever to save your life from such challenging moments.

At the bottom right of the modal is a 'Close' button. The background of the page features a dark sidebar with social media icons and a 'Call Us' phone number.

**Figure 6.8 : Transcript Generation**

Displays the generated transcript from the uploaded video. The transcription includes timestamps and is used as the foundation for further theme detection and summarization.

**Feedback Form**

Your Name:

Email:

Feedback: Rating: 3/5

Additional Comments:

Submit Feedback

We value your opinion! ☀️

**Figure 6.9 : Feedback Form**

Users can submit feedback on the accuracy of summaries, theme detection, and overall user experience. This information helps refine system performance and improve future updates

Follow Us: [f](#) [t](#) [in](#) [r](#)

Call Us: 94900 74884

**VIDEO INSIGHTS**

Dashboard Upload History My Profile Feedback Logout

**Profile Settings**

Name <input type="text" value="bindu"/>	Email <input type="text" value="binduvarsha2004@gmail.com"/>
You can't update this field	
Phone <input type="text" value="9490074884"/>	Profile <input type="file" value="Choose File"/> No file chosen
Password <input type="text" value="Bindu@2004"/>	Location <input type="text" value="Telangana"/>

Update Profile

**Figure 6.10 : Profile Settings**

The profile section allows users to view their account information, history of uploads, summaries generated, and Q&A activity. It offers options for updating personal details and managing preferences.

# VIDEO INSIGHTS

Dashboard    Upload    History    My Profile    Feedback

Logout

back. Poor fisherman, he came next day expecting that there will be a lot of fish that this tiny little fish would bring, but the tiny little fish was very clever and because of his cleverness, he saved his life from this fisherman. So children, more of the story is you have to be really, really clever to save your life from such challenging moments.

[Detect Themes via API](#)

[Detect Themes via Model](#)

## Detected Themes (via API):

Here are the main themes extracted from the transcript:

\*\*Cleverness and Intelligence\*\*: The story highlights the importance of being clever and using intelligence to overcome challenges.

\*\*Deception and Strategy\*\*: The clever fish uses deception to trick the fisherman, illustrating strategic thinking in difficult situations.

\*\*Survival and Self

Preservation\*\*: The fish's actions are driven by the need to survive and protect itself from danger.

\*\*Moral Lessons\*\*: The story teaches children valuable moral lessons about being clever and resourceful in challenging moments.

\*\*Friendship and Teamwork\*\*: Although not directly shown, the fish's

[Back](#)



**Figure 6.11 : Theme Detection**

This page presents the list of themes detected by the LSTM-based classifier. Each theme is timestamped and linked to its corresponding segment in the video, enabling users to navigate directly to relevant content.

# CHAPTER 7

## CONCLUSION

### 7.1 Summary of Work Done

The AI Based Video Insights Generator was developed to provide an automated, efficient, and intelligent solution for detecting themes from videos and text while also summarizing lengthy content. The system was designed to handle various multimedia sources, including YouTube videos, Google Drive videos, and locally uploaded content, ensuring comprehensive theme detection and content summarization.

The project was implemented using deep learning-based architectures, specifically:

- **LSTM (Long Short-Term Memory) and Conv1D models** for theme detection.
- **Transformer-based models (T5, Pegasus, BART)** for text summarization.
- **Speech-to-text transcription** for converting video audio into text.
- **Timestamp extraction** for mapping detected themes to relevant video segments.
- **Interactive Q&A system** for querying extracted content.
- **Multilingual translation module** for broad accessibility.

A user authentication system was integrated, allowing secure registration, login, and feedback collection. The entire system was optimized using Early Stopping, Model Checkpoint, and performance evaluation techniques, ensuring high accuracy, efficiency, and scalability.

Extensive testing and validation were conducted on various video and text datasets, leading to positive results in terms of theme detection accuracy, summarization quality, and system performance. The system demonstrated strong real-world applicability for academic, corporate, and research-based content analysis.

### 7.2 Challenges Faced and Solutions Implemented

During the development and implementation of the system, several challenges were encountered. Below is a summary of key challenges and the solutions adopted:

## **1. Handling Long Videos and Large Text Data**

**Challenge:** Processing long-duration videos led to increased computational load and memory usage, making it difficult to extract themes and summarize content efficiently.

**Solution:** Implemented batch processing and chunk-based text processing to divide long content into manageable segments. Also optimized GPU acceleration and memory management techniques to handle large inputs smoothly.

## **2. Accuracy of Speech-to-Text Transcription**

**Challenge:** Automatic transcription of video audio sometimes produced errors, especially in noisy environments or for videos with multiple speakers.

**Solution:** Used advanced speech recognition APIs with noise filtering and speaker diarization to improve transcription accuracy. Implemented post-processing techniques such as punctuation correction and text cleaning.

## **3. Theme Detection Accuracy in Unstructured Content**

**Challenge:** The LSTM model sometimes struggled with free-flowing, informal discussions, leading to lower accuracy in theme extraction.

**Solution:** Improved data preprocessing techniques (tokenization, stopword removal, lemmatization) and fine-tuned hyperparameters in the LSTM model to enhance detection accuracy.

## **4. Summarization Coherence and Relevance**

**Challenge:** The summarization output sometimes missed key points or produced generic sentences instead of context-rich summaries.

**Solution:** Implemented multiple transformer models (T5, Pegasus, BART) and evaluated their performance on different content types, selecting the best model based on ROUGE score analysis.

## **5. Multilingual Translation Consistency**

**Challenge:** Some complex phrases and technical terms were not accurately translated into other languages.

**Solution:** Integrated custom translation dictionaries and fine-tuned API configurations to

enhance translation quality for domain-specific terms.

## **6. System Performance and Latency Issues**

**Challenge:** The processing time increased with longer videos and high-resolution audio files, affecting system responsiveness.

**Solution:** Optimized the backend using efficient parallel processing, caching mechanisms, and API request batching to reduce latency.

Through these optimizations, the system achieved higher accuracy, reduced processing time, and improved usability, making it more efficient for large-scale use.

## **7.3 Future Scope and Enhancements**

Although the AI Based Video Insights Generator has achieved significant milestones, there are several areas for future improvement and expansion to enhance its capabilities.

### **1. Integration of Real-Time Streaming Analysis**

Currently, the system processes pre-recorded videos, but future enhancements could allow real-time theme detection and summarization from live-streaming content. This would be useful for news analysis, live lectures, and virtual meetings.

### **2. Improved Natural Language Understanding for Q&A System**

The interactive Q&A module can be improved by integrating advanced NLP models (like GPT-based models) to generate more contextually rich and precise answers to user queries.

### **3. Sentiment Analysis for Theme-Based Content**

Adding sentiment analysis to detected themes could provide deeper insights into the content, particularly for review-based videos, debates, and opinion discussions.

### **4. Expansion to More Languages and Dialects**

The multilingual module can be expanded to support additional regional languages and dialects, improving accessibility for users from diverse linguistic backgrounds.

### **5. Enhanced UI/UX for Better User Interaction**

A more intuitive and visually appealing interface could be developed, incorporating interactive charts, visual theme maps, and user-friendly dashboards for better user experience.

## 7. Model Fine-Tuning for Industry-Specific Applications

The system could be adapted for specific industries such as:

- **Healthcare:** Detecting medical themes from patient discussions or research papers.
- **Legal Industry:** Summarizing court proceedings or case studies.
- **Corporate Sector:** Extracting key insights from business meetings or reports.

By implementing these enhancements, the system could become a comprehensive AI-powered content analysis tool with broad applications across education, research, business, and media industries.

## APPENDIX

### Dataset Description (Video Transcriptions, Theme Classification Data)

The dataset used for this project consists of video transcriptions and textual data curated from various sources. The data was processed using speech-to-text transcription APIs to convert audio into textual format. Additionally, manually labelled datasets were used for theme classification and summarization model training.

#### 1. Video Transcriptions Dataset

- The dataset includes transcriptions of spoken content extracted from videos across multiple domains such as educational lectures, news reports, business meetings, and research presentations.
- **Metadata:** Each transcription entry contains the video ID, timestamp, speaker details (if available), and the corresponding text.
- **Preprocessing:** Tokenization, stopword removal, lemmatization, and punctuation correction were applied to clean the transcribed data.
- **Storage Format:** Data is stored in structured formats such as CSV and JSON, enabling easy integration into machine learning pipelines.

#### 2. Theme Classification Dataset

- **Description:** The dataset contains labeled textual data, where each entry is categorized under a specific theme (e.g., Technology, Health, Finance, Science, etc.).
- **Data Sources:** Collected from open-source repositories and manually labeled video transcriptions.
- **Structure:**
  - **Input:** Text snippets (from transcriptions)
  - **Output:** Theme labels assigned using an LSTM-based classifier
- **Usage:** The dataset was used to train and evaluate the LSTM-90 model, enabling accurate theme detection from video and text sources.

## REFERENCES

### 1. Research Papers and Journals

1. **Hochreiter, S., & Schmidhuber, J. (1997).** "Long Short-Term Memory." *Neural Computation*, 9(8), 1735-1780.
  - This paper introduces the LSTM architecture, which is widely used for sequential data classification and time-series prediction.
2. **Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., & Polosukhin, I. (2017).** "Attention is All You Need." *Advances in Neural Information Processing Systems (NeurIPS)*.
  - This paper introduces the Transformer architecture, which serves as the foundation for T5, Pegasus, and BART models used in text summarization.
3. **Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018).** "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding." *arXiv preprint arXiv:1810.04805*.
  - This research paper provides insights into Transformer-based models, which significantly improved NLP tasks such as text classification, summarization, and question answering.
4. **Lin, C. Y. (2004).** "ROUGE: A Package for Automatic Evaluation of Summaries." *Workshop on Text Summarization Branches Out*.
  - The ROUGE metric is widely used to evaluate text summarization models, ensuring high-quality summarization output.

### 2. Books and Online Documentation

5. **Goodfellow, I., Bengio, Y., & Courville, A. (2016).** *Deep Learning*. MIT Press.
  - This book provides a comprehensive overview of deep learning architectures, including LSTM, CNN, and Transformer-based models.
6. **Jurafsky, D., & Martin, J. H. (2021).** *Speech and Language Processing*. Pearson.
  - Covers speech recognition, NLP techniques, and deep learning approaches relevant to speech-to-text transcription.

### **3. Online Resources and API Documentation**

7. **TensorFlow Documentation** - "Long Short-Term Memory Networks with Keras and TensorFlow."
  - Available at: <https://www.tensorflow.org/>
  - Provides implementation details for LSTM-based sequence classification models, which were used for theme detection.
8. **Hugging Face Transformers** - "Pre-trained Transformer Models for NLP."
  - Available at: <https://huggingface.co/transformers/>
  - Documentation on T5, Pegasus, and BART models, which were used for text summarization.
9. **Google Cloud Speech-to-Text API** - "Convert Speech into Text using Deep Learning."
  - Available at: <https://cloud.google.com/speech-to-text>
  - Describes Google's API for speech-to-text transcription, used for processing video/audio data.
10. **OpenAI API Documentation** - "Natural Language Processing and Question Answering Systems."
  - Available at: <https://platform.openai.com/docs/>
  - Details on NLP-based question-answering models, which were integrated into the interactive Q&A system.



## INTERNATIONAL JOURNAL OF ADVANCED RESEARCH IN COMPUTER AND COMMUNICATION ENGINEERING

A monthly Peer-reviewed & Refereed journal

Impact Factor 8.102

Indexed by Google Scholar, Mendeley, Crossref, Scilit,  
SCIENCEOPEN, SCIENCEGATE, DORA, KOAR



Google Scholar



Crossref



MENDELEY

## *Certificate of Publication*

### SK. WASIM AKRAM

Assistant Professor, Dept. of. CSE, VVIT, GUNTUR, AP, INDIA

Published a paper entitled

**AI Based Video Insights Generator**

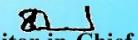
Volume 14, Issue 3, March 2025

DOI: [10.17148/IJARCCE.2025.14320](https://doi.org/10.17148/IJARCCE.2025.14320)

Certificate# IJARCCE/2025/1

ISSN (Online) 2278-1021  
ISSN (Print) 2319-5940

IJARCCE  
DOI 10.17148/IJARCCE

  
Editor-in-Chief  
IJARCCE

SERTİFAAT

證書

سے

Sertifikat

CERTYFIKAT CERTIFICADO

증명서

C E R T I F I K A T 证

POTVRDA

SERTIFIKA



## INTERNATIONAL JOURNAL OF ADVANCED RESEARCH IN COMPUTER AND COMMUNICATION ENGINEERING

A monthly Peer-reviewed & Refereed journal

Impact Factor 8.102

Indexed by Google Scholar, Mendeley, Crossref, Scilit,  
SCIENCEOPEN, SCIENCEGATE, DORA, KOAR



Google Scholar



Crossref



MENDELEY

## *Certificate of Publication*

**Y. BINDU VARSHA**

Student, Dept.of.CSE Artificial Intelligence and Machine Learning, VVIT, GUNTUR, AP, INDIA

Published a paper entitled

**AI Based Video Insights Generator**

Volume 14, Issue 3, March 2025

DOI: [10.17148/IJARCCE.2025.14320](https://doi.org/10.17148/IJARCCE.2025.14320)

Certificate# IJARCCE/2025/1

ISSN (Online) 2278-1021  
ISSN (Print) 2319-5940

IJARCCE  
DOI 10.17148/IJARCCE

  
Editor-in-Chief  
IJARCCE

SERTIFIKAAT

證書

سے

Sertifikat

CERTYFIKAT CERTIFICADO

증명서

C E R T I F I K A T 证

POTVRDA

SERTIFIKĀ

[www.ijarcce.com](http://www.ijarcce.com)



## INTERNATIONAL JOURNAL OF ADVANCED RESEARCH IN COMPUTER AND COMMUNICATION ENGINEERING

A monthly Peer-reviewed & Refereed journal

Impact Factor 8.102

Indexed by Google Scholar, Mendeley, Crossref, Scilit,  
SCIENCEOPEN, SCIENCEGATE, DORA, KOAR



## Certificate of Publication

**P. SAMBASIVARAO**

Student, Dept.of.CSE Artificial Intelligence and Machine Learning, VVIT, GUNTUR, AP, INDIA

Published a paper entitled

**AI Based Video Insights Generator**

Volume 14, Issue 3, March 2025

DOI: [10.17148/IJARCCE.2025.14320](https://doi.org/10.17148/IJARCCE.2025.14320)

Certificate# IJARCCE/2025/1

ISSN (Online) 2278-1021  
ISSN (Print) 2319-5940

IJARCCE  
DOI 10.17148/IJARCCE

  
Editor-in-Chief  
IJARCCE

SERTİFAAT  


سے  
Sertifikat

CERTYFIKAT  
CERTIFICADO  
증명서

C E R T I F I K A T  
證  
ПОВРДА

SERTİFİKA  
SERTIFIKAT



## INTERNATIONAL JOURNAL OF ADVANCED RESEARCH IN COMPUTER AND COMMUNICATION ENGINEERING

A monthly Peer-reviewed & Refereed journal

Impact Factor 8.102

Indexed by Google Scholar, Mendeley, Crossref, Scilit,  
SCIENCEOPEN, SCIENCEGATE, DORA, KOAR



Google Scholar



Crossref



MENDELEY

## *Certificate of Publication*

**P. SNEHAL KUMAR**

Student, Dept.of.CSE Artificial Intelligence and Machine Learning, VVIT, GUNTUR, AP, INDIA

Published a paper entitled

**AI Based Video Insights Generator**

Volume 14, Issue 3, March 2025

DOI: [10.17148/IJARCCE.2025.14320](https://doi.org/10.17148/IJARCCE.2025.14320)

Certificate# IJARCCE/2025/1

ISSN (Online) 2278-1021  
ISSN (Print) 2319-5940

IJARCCE  
DOI 10.17148/IJARCCE

  
Editor-in-Chief  
IJARCCE

SERTIFIKAAT

證書

سے

Sertifikat

CERTYFIKAT CERTIFICADO

증명서

C E R T I F I K A T 证

POTVRDA

SERTIFICA



## INTERNATIONAL JOURNAL OF ADVANCED RESEARCH IN COMPUTER AND COMMUNICATION ENGINEERING

A monthly Peer-reviewed & Refereed journal

Impact Factor 8.102

Indexed by Google Scholar, Mendeley, Crossref, Scilit,  
SCIENCEOPEN, SCIENCEGATE, DORA, KOAR



## *Certificate of Publication*

**V. CHARAN SAI VENKAT**

Student, Dept.of.CSE Artificial Intelligence and Machine Learning, VVIT, GUNTUR, AP, INDIA

Published a paper entitled

**AI Based Video Insights Generator**

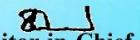
Volume 14, Issue 3, March 2025

DOI: [10.17148/IJARCCE.2025.14320](https://doi.org/10.17148/IJARCCE.2025.14320)

Certificate# IJARCCE/2025/1

ISSN (Online) 2278-1021  
ISSN (Print) 2319-5940

IJARCCE  
DOI 10.17148/IJARCCE

  
Editor-in-Chief  
IJARCCE

SERTİFAAT

證書

سے

Sertifikat

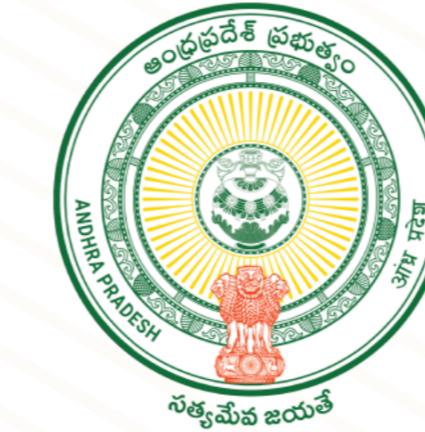
CERTYFIKAT CERTIFICADO

증명서

C E R T I F I K A T 证

POTVRDA

SERTIFIKA



# ANDHRA PRADESH STATE COUNCIL OF HIGHER EDUCATION

(A Statutory Body of the Government of A. P.)

## CERTIFICATE OF COMPLETION

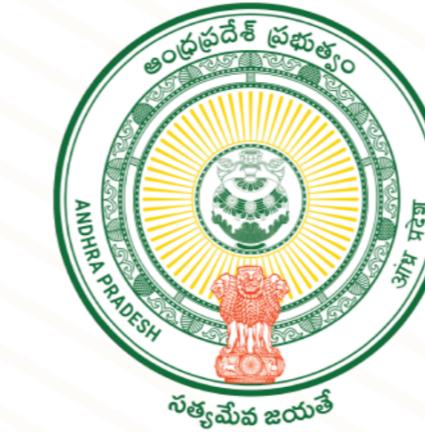
This is to certify that Ms./Mr. **Yaddanapudi Bindu Varsha**, Course : B.Tech, Branch : CSE, Semester : 8th, Roll No : **21BQ1A42I2** Under **Vasireddy Venkatadri Institute of Technology of JNTU, Kakinada** has Successfully completed the **Long-Term Internship for 240 Hours** on **Uipath Rpa Developer & Web Full Stack Developer** Organised by EduSkills in Collaboration with **Andhra Pradesh State Council of Higher Education**.

Certificate No. : 70eBA9eABbFCF7

Issue Date : 04 March 2025



Chief Technology Officer (CTO)  
EduSkills



# ANDHRA PRADESH STATE COUNCIL OF HIGHER EDUCATION

(A Statutory Body of the Government of A. P.)

## CERTIFICATE OF COMPLETION

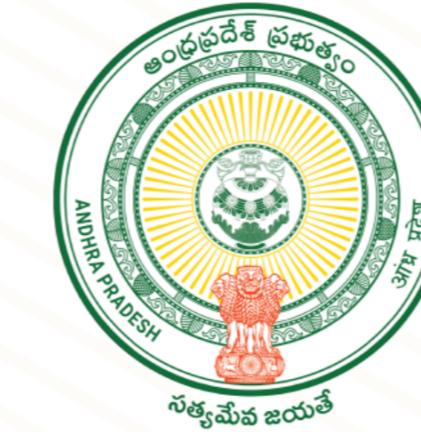
This is to certify that Ms./Mr. **SAMBASIVARAO PRATHIPATI**, Course : B.Tech, Branch : CSE, Semester : 8th, Roll No : 21bq1a42f4 Under **Vasireddy Venkatadri Institute of Technology of JNTU, Kakinada** has Successfully completed the **Long-Term Internship for 240 Hours on Juniper Networking & Altair Data Science Master Virtual Internship** Organised by **EduSkills** in Collaboration with **Andhra Pradesh State Council of Higher Education**.

Certificate No. : FE13D5A4eac73B8

Issue Date : 22 March 2025



Chief Technology Officer (CTO)  
EduSkills



# ANDHRA PRADESH STATE COUNCIL OF HIGHER EDUCATION

(A Statutory Body of the Government of A. P.)

## CERTIFICATE OF COMPLETION

This is to certify that Ms./Mr. **SNEHAL KUMAR PADARTHI**, Course : B.Tech, Branch : AI-ML, Semester : 8th, Roll No : 21BQ1A42D5 Under **Vasireddy Venkatadri Institute of Technology of JNTU, Kakinada** has Successfully completed the **Long-Term Internship for 240 Hours** on **Microchip Embedded System Developer & Altair Data Science Master Virtual Internship** Organised by **EduSkills** in Collaboration with **Andhra Pradesh State Council of Higher Education**.

Certificate No. : aCFE0CBD6F8beeE

Issue Date : 22 March 2025



Chief Technology Officer (CTO)  
EduSkills



# ANDHRA PRADESH STATE COUNCIL OF HIGHER EDUCATION

(A Statutory Body of the Government of A. P.)

## CERTIFICATE OF COMPLETION

This is to certify that Ms./Mr. **Charan Sai venkat Vegi**, Course : B.Tech, Branch : AI-ML, Semester : 7th, Roll No : 21bq1a42h9 Under **Vasireddy Venkatadri Institute of Technology of JNTU, Kakinada** has Successfully completed the **Long-Term Internship for 240 Hours** on Zscaler Networking & Google AI-ML Organised by **EduSkills** in Collaboration with **Andhra Pradesh State Council of Higher Education**.

Certificate No. : 958Be809CDD7F4E

Issue Date : 13 March 2025



Chief Technology Officer (CTO)  
EduSkills