# Real-Time Conferencing System with Summarization using WebRTC and NLP

Major Project Report submitted in

Partial fulfillment of requirements for the degree of

Bachelor of Technology

*In*

Computer Engineering

*By*

Reesav Gupta (21CSEC08)

Roshan Chhetri (21CSEC10)

Anil Kumar Pandit (21CSEC31)

Under the supervision of

**Mr. Asish Shakya**

Assistant Professor
Computer Engineering Department



**DEPARTMENT OF COMPUTER ENGINEERING**

**SIKKIM INSTITUTE OF SCIENCE & TECHNOLOGY**

**(Affiliated to Central University of Sikkim)**

**June 2025**

# CERTIFICATE BY SUPERVISOR

This is to certify that the work recorded in this project report entitled "Real-Time Conferencing System with Summarization using WebRTC and NLP" has been jointly carried out by Mr. Reesav Gupta (21CSEC08), Mr. Roshan Chhetri (21CSEC10), and Mr. Anil Kumar Pandit (21CSEC31) of the Computer Engineering Department, Sikkim Institute of Science & Technology, in partial fulfilment of the requirements for the award of the Bachelor of Technology in Computer Engineering.

This report has been duly reviewed by the undersigned and is hereby recommended for final submission for the Major Project Viva Examination.

<table>
<tr><td>Mr. Asish Shakya</td><td>Mr. Prasant Pradhan</td></tr>
<tr><td>**Assistant Professor (Guide)**</td><td>**Head of Department**</td></tr>
<tr><td>Department of Computer Engineering</td><td>Department of Computer Engineering</td></tr>
<tr><td>Sikkim Institute of Science &Technology</td><td>Sikkim Institute of Science &Technology</td></tr>
</table>

Ms. Pragya Pradhan

**Project Coordinator**

Department of Computer Engineering

Sikkim Institute of Science &Technology

# DECLARATION

We, the undersigned, hereby solemnly affirm and declare that the work presented in this project report titled "Real-Time Conferencing System with Summarization using WebRTC and NLP" fulfils the requirements for the completion of the Bachelor of Technology in Computer Engineering from Sikkim Institute of Science and Technology (affiliated to the Central University, Sikkim).

This project work was diligently and faithfully carried out at SIST, Chisopani, under the meticulous guidance and supervision of Mr. Asish Shakya, Assistant Professor in the Department of Computer Engineering at Sikkim Institute of Science and Technology.

We affirm that the project work described in this report is the result of our sincere and dedicated efforts, conducted in accordance with the prescribed academic standards and institutional guidelines. All information presented, including the methodologies employed, data collected, analyses performed, and conclusions drawn, are accurate and reliable to the best of our knowledge and capabilities.

Reesav Gupta (21CSEC08)

Roshan Chhetri (21CSEC10)

Anil Kumar Pandit (21CSEC31)

# ACKNOWLEDGEMENT

# ABSTRACT

In the era of digital collaboration, effective real-time communication is essential. This project introduces a many-to-many video conferencing application designed to facilitate seamless virtual meetings with advanced features such as video/audio chat, and meeting recording. Built using WebRTC and Mediasoup, the platform ensures low-latency and high-quality media streaming, making it suitable for various professional and personal use cases.

A key innovation of this application is its AI-powered meeting summarization and individual participant audio recording, enabling efficient post-meeting analysis and documentation. Due to computational constraints in the initial phase, supervised learning will be implemented for AI-based summarization. In the final phase, if feasible, reinforcement learning will be utilized to fine-tune the system for improved performance and adaptability.

By leveraging modern web technologies and AI, the system offers a scalable and reliable conferencing experience, enhancing productivity and collaboration in remote settings.

***Keywords****: Video Conferencing, WebRTC, AI-powered Summarization, Mediasoup, Real-time Communication, Low Latency Streaming*

# TABLE OF CONTENT

# LIST OF FIGURES

# LIST OF TABLES

# LIST OF ABBREVIATIONS

| ABBREVIATION | FULL-FORM |
|---|---|
| AI | Artificial Intelligence |
| WebRTC | Web Real-Time Communication |
| ML | Machine Learning |
| RL | Reinforcement Learning |
| FFmpeg | Fast Forward MPEG |
| SFU | Selective Forwarding Unit |
| NLP | Natural Language Processing |
| BART | Bidirectional and Auto-Regressive Transformers |
| RTP | Real-time Transport Protocol |
| SRTP | Secure Real-time Transport Protocol |
| LED | Longformer Encoder-Decoder |
| ASR | Automatic Speech Recognition |
| GPT | Generative Pre-trained Transformer |
| RAM | Random Access Memory |
| CPU | Central Processing Unit |
| GPU | Graphics Processing Unit |
| IDE | Integrated Development Environment |
| E2EE | End-to-End Encryption |
| ORM | Object-Relational Mapping |
| ROUGE | Recall-Oriented Understudy for Gisting Evaluation |

# CHAPTER 1

# INTRODUCTION

In recent years, video conferencing has become an essential tool for communication, connecting people across different locations for work, education, and personal interactions. The increasing reliance on virtual meetings has created a demand for feature-rich and intelligent video conferencing solutions that enhance collaboration, productivity, and post-meeting analysis.

This project focuses on developing a many-to-many video conferencing application that provides real-time video/audio communication, and meeting recording to facilitate seamless virtual interactions. Unlike conventional solutions, this system introduces individual participant audio recording, allowing for better clarity and flexible post-meeting analysis. Additionally, the application integrates AI-driven meeting summarization, which enables users to quickly access key discussion points without reviewing lengthy recordings, improving efficiency and accessibility.

The application is built using WebRTC and Mediasoup, ensuring low-latency, high-quality media streaming for smooth communication. WebRTC enables real-time peer-to-peer interaction, while Mediasoup provides a scalable media server for handling multiple participants efficiently.

Due to computational limitations, the AI-powered summarization will initially be implemented using supervised learning. In later phases, reinforcement learning may be incorporated to fine-tune and enhance the summarization process for improved accuracy and adaptability. This gradual integration allows for optimized resource utilization while ensuring scalability.

By leveraging modern web technologies and AI, this project aims to enhance remote collaboration with an intelligent, scalable, and user-friendly video conferencing solution. Its AI-powered summarization and advanced recording capabilities make it an innovative tool for businesses, educational institutions, and individuals seeking efficient and productive virtual meetings.

## 1.1  Literature Review

Table 1.1.1: Literature Review 1

| 1. | Name of the Paper | The Design and Architecture of a WebRTC Application. |
|---|---|---|
| | **Publication Details** | **Authored by:** Simon Holm, Alexander Lööf <br><br> **Published by:** IEEE <br><br> **Publication Date**: June 4, 2019 |
| | **Research Findings** | <ul><li>The paper explores design patterns for WebRTC applications to achieve scalability, performance, and efficiency.</li><li>It compares two architectures:<ul><li>Full Mesh Topology – Efficient for small groups but does not scale well.</li><li>Star Topology with SFU (Selective Forwarding Unit) – More scalable but requires higher bandwidth.</li></ul></li></ul> |
| | **Research Gap** | <ul><li>The study does not address client-side bandwidth issues, which may affect SFU performance.</li><li>It does not explore AI-based solutions for optimizing real-time media management.</li></ul> |
| | **Relevance** | <ul><li>The paper provides useful insights for our many-to-many video conferencing system.</li><li>Since our project involves AI-powered meeting summarization, future work can explore AI-driven WebRTC optimizations.</li></ul> |

Table 1.1.2: Literature Review 2

| 2. | Name of the Paper | Whisper Model: Large-Scale Weak Supervision for Speech Recognition. |
|---|---|---|
| | Publication Details | **Authored by:** A. Radford, J. W. Kim, G. Brockman, T. Xu, C. McLeavey, I. Sutskever<br><br>**Published by:** OpenAI Research<br><br>**Publication Date**: 2022 |
| | Research Findings | • A large-scale AI speech recognition model trained on multiple languages.<br>• Works well in noisy environments and varied accents.<br>• Produces highly accurate transcriptions across different settings. |
| | Research Gap | • Not optimized for meeting-specific terminology.<br>• Cannot differentiate between speakers in group discussions.<br>• Requires high computing power, making real-time use difficult. |
| | Relevance | • Works well in different languages, noise conditions, and accents.<br>• Provides high accuracy for speech recognition. |

Table 1.1.3: Literature Review 3

| 3. | Name of the Paper | PEGASUS: AI Model for Summarization. |
|---|---|---|
| | Publication Details | **Authored by:** J. Zhang, Y. Zhao, M. Saleh, P. Liu<br><br>**Published by:** Proceedings of the 37th International Conference on Machine Learning (ICML)<br><br>**Publication Date**: 2020 |
| | Research Findings | <ul><li>Introduces a new AI model for text summarization.</li><li>Achieves state-of-the-art results on various summarization tasks.</li><li>Uses a pre-training method to improve summary generation quality.</li></ul> |
| | Research Gap | <ul><li>Not designed for summarizing live meetings.</li><li>Does not focus on speaker importance or meeting agenda.</li><li>Sometimes generates incorrect information.</li></ul> |
| | Relevance | <ul><li>Uses AI to generate high-quality summaries.</li><li>Achieves good results in general text summarization tasks.</li></ul> |

Table 1.1.4: Literature Review 4

| 4. | Name of the Paper | Fine-Tuning GPT Models with Human Feedback. |
|---|---|---|
| | Publication Details | **Authored by:** D. Ziegler, N. Stiennon, P. Christiano, J. Wu, T. Brown, A. Radford, D. Amodei, G. Irving, J. Leike<br><br>**Published by:** Advances in Neural Information Processing Systems (NeurIPS)<br><br>**Publication Date**: 2019 |
| | Research Findings | • Uses reinforcement learning (RL) and human feedback to improve AI-generated text.<br>• Enhances summary coherence, informativeness, and fluency.<br>• Helps fine-tune large language models for better user experience. |
| | Research Gap | • Not specifically designed for meeting summarization.<br>• Does not work in real-time.<br>• Needs well-defined feedback to generate the best summaries. |
| | Relevance | • Uses reinforcement learning to improve text quality.<br>• Helps make AI-generated summaries clearer and more useful. |

# CHAPTER 2
# PROBLEM STATEMENT AND OBJECTIVE

## 2.1 Problem Statement

- Traditional video conferencing platforms lack advanced features that enhance productivity, such as AI-powered meeting summarization and individual participant audio recording.

- Many existing solutions suffer from high latency and inefficient media streaming, leading to poor user experiences.

- There is no automated system that efficiently extracts key meeting insights, requiring users to manually review long recordings.

## 2.2 Objectives

The main objectives are:

- Develop a real-time video conferencing platform with WebRTC and Mediasoup.

- Implement AI-driven meeting summarization for quick review of discussions.

- Enable individual participant audio recording for better clarity and analysis.

- Ensure low-latency, high-quality media streaming for seamless communication.

# CHAPTER 3

# SOLUTION STRATEGY

## 3.1 Solution Strategy



Figure 3.1.1: Solution Strategy

The AI-powered Smart Conferencing System is designed to enhance video conferencing experiences by incorporating AI-driven features such as real-time transcription, meeting summarization, participant audio analysis, and efficient media streaming. The solution strategy involves multiple stages:

- **Development of Real-Time Conferencing Platform:** The system will leverage WebRTC and Mediasoup to create a low-latency, high-quality video conferencing environment that ensures seamless communication between participants.

- **AI-Powered Meeting Summarization:** Implement Natural Language Processing (NLP) techniques using models such as BERT, GPT, or Transformers to summarize discussions and extract key points from conversations in real-time.

- **Participant-Specific Audio Recording & Analysis:** The system will use speaker diarization techniques to distinguish different participants and record their audio streams separately for better clarity and individual analysis.

- **Speech-to-Text Transcription:** Implement Automatic Speech Recognition (ASR) models such as DeepSpeech or Whisper AI to convert spoken words into text in real time, enabling accessibility features like live captions.

- **Latency Optimization & Media Streaming Efficiency:** Mediasoup's Selective Forwarding Unit (SFU) architecture will be used to efficiently distribute media streams, ensuring minimal lag and bandwidth optimization.

- **User Interface & Experience Design:** A React.js and Tailwind CSS-based UI will be developed for a modern and user-friendly experience, allowing seamless navigation and interaction.

- **Security & Privacy Measures:** The system will integrate end-to-end encryption (E2EE), access control mechanisms, and data anonymization to ensure user data protection and confidentiality.
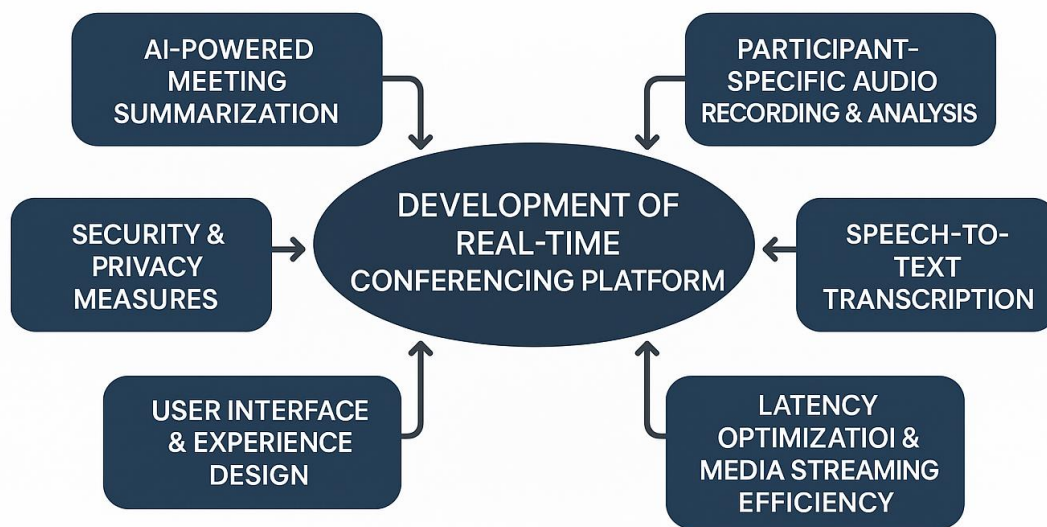
# CHAPTER 4

# METHODOLOGY

## 4.1 System Architecture



Figure 4.1.1: System Architecture

The AI-powered Smart Conferencing System is designed to enhance video conferencing experiences by leveraging artificial intelligence for audio transcription, speaker identification, and meeting summarization. The system follows a structured architecture consisting of multiple modules that work together to process and analyze meeting data efficiently. The key components of the system are:

## 4.1.1 Clients

The system supports two types of clients:

- Web Client – Allows users to access the conferencing platform from their browsers.

- Mobile Client – Provides on-the-go access to meetings using mobile devices.

These clients interact with the backend to facilitate real-time communication, transcription, and AI-powered summarization.

## 4.1.2 Data Processing Module

The data processing module is responsible for handling the audio and text processing tasks. This module consists of:

1. Audio Capturing & Processing:

    o Audio from meetings is first captured and preprocessed using FFmpeg, an open-source multimedia framework.

    o The processed audio is then passed to the Whisper Model for speech-to-text conversion.

2. Speech-to-Text Conversion (Whisper Model):

    o Whisper, an advanced AI model by OpenAI, transcribes the meeting audio into text format with high accuracy.

    o The transcribed text is then processed for speaker identification and summarization.

3. Speaker Identification:

    o This module identifies and labels individual speakers, making the transcript more structured and readable.

    o The identified speakers help improve the clarity of the meeting transcription.

4. Text Summarization (BART and Fine-Tuned LED-Arxiv Models):

    o The system uses two transformer-based models, BART and a Fine-Tuned LED-Arxiv, to generate concise and meaningful summaries from meeting transcripts.

    o BART provides high-quality summaries for shorter inputs and excels in semantic accuracy.

    o LED-Arxiv, fine-tuned for long documents, is used to handle lengthy meeting transcripts exceeding typical token limits.

## 4.1.3 Persistent Storage

- The system securely stores meeting transcripts and summaries in a structured format.
- Users can access past meeting notes, transcripts, and summaries for future reference.

## 4.1.4 Security & Privacy

- The system ensures secure data transmission and storage to protect meeting information.
- Encryption and access control mechanisms are implemented to maintain confidentiality.

# 4.2 Models Used in the System

The AI-powered Smart Conferencing System utilizes advanced deep learning models to process meeting audio and generate meaningful summaries. The two primary models used in the system are Whisper for speech-to-text conversion, BART and Fine-Tuned LED-Arxiv for transcript summarization. These models significantly enhance productivity by automating the documentation and review of meetings.

## 4.2.1 Audio-to-Text Conversion using Whisper Model

The Whisper model, developed by OpenAI, is an advanced speech recognition model that converts spoken language into text. It is trained on a diverse dataset of multilingual and multitask speech data, making it highly accurate and efficient for transcription.

**Key Features of Whisper Model:**

- **Robust Speech Recognition:** Can transcribe audio with different accents, background noise, and multiple speakers.
- **Multilingual Support:** Recognizes and transcribes speech in multiple languages.
- **High Accuracy:** Pre-trained on large-scale datasets, improving transcription quality.

**Workflow of Whisper Model in the System:**

1. **Audio Input:** The meeting audio is recorded and preprocessed using FFmpeg to ensure optimal clarity.
2. **Transcription:** The Whisper Model processes the audio and converts it into text format with timestamps.
3. **Speaker Segmentation:** The system identifies different speakers, making the transcript structured and easier to understand.
4. **Output Generation:** The generated transcript is stored for further processing and summarization.

By integrating Whisper, the system automates transcription, eliminating the need for manual note-taking and improving accessibility.

## 4.2.2 Transcript Summarization using BART and Fine-Tuned LED-Arxiv Models

The system utilizes two powerful transformer-based models — BART and Fine-Tuned LED-Arxiv (Longformer Encoder-Decoder) — for transcript summarization. Both models are designed for abstractive summarization, meaning they generate summaries in natural language rather than copying parts of the transcript.

While BART is effective for shorter transcripts with high semantic accuracy, the Fine-Tuned LED-Arxiv model is specialized in handling long-form content, making it ideal for lengthy meeting transcripts.

**Key Features:**

- **BART Model:**

    o Well-suited for short to medium-length texts.

    o Generates fluent, human-like summaries.

    o Pre-trained on diverse datasets for better generalization.

- **LED-Arxiv Model:**

    o Efficient for processing long documents beyond the token limit of BART.

    o Supports both extractive and abstractive summarization.

    o Fine-tuned on academic and structured datasets for higher precision and coherence.

**Summarization Workflow in the System:**

1. **Input Processing:** The transcript generated by Whisper is cleaned, segmented, and prepared for summarization.
2. **Model Selection:**

    o BART is used when transcript length is short or moderate.

    o LED-Arxiv is chosen for long transcripts (e.g., > 1024 tokens).

3. **Summarization:** The selected model processes the transcript and extracts the most relevant points.
4. **Summary Generation:** A concise, readable summary is generated using abstractive techniques.
5. **Storage & Access:** The generated summary is saved in the system's database and is made accessible to users for easy review.

By integrating both BART and Fine-Tuned LED-Arxiv, the system provides flexible and efficient summarization for various types of meeting content, ensuring users can focus on the most important information without manually reading through long transcripts.
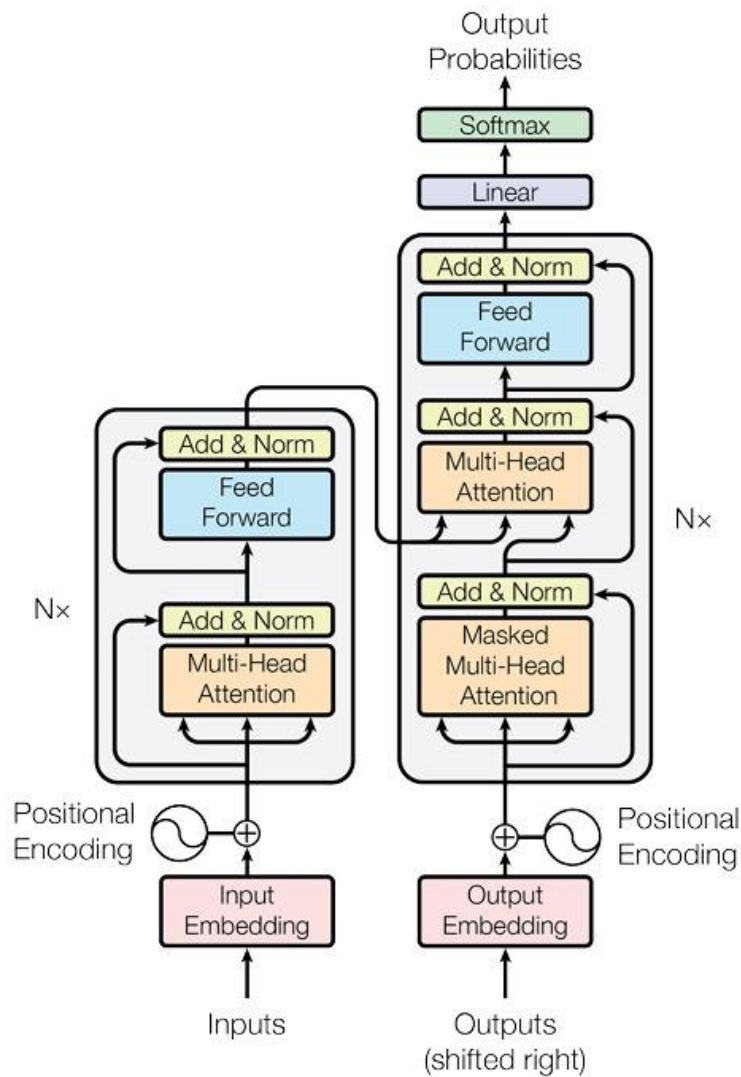
# 4.3 Transformer Architecture



Figure 4.3.1: Transformer Architecture

## Key Components of the Transformer Model

### 1. Input & Output Embeddings

- The model first converts words into dense vector representations using word embeddings.
- Since Transformers process all tokens simultaneously, positional encoding is added to retain the order of words.

**2. Encoder (Left Side)**

The encoder stack consists of N identical layers, where each layer has:

- Multi-Head Attention: Helps the model focus on different words in the input sentence.
- Add & Norm (Residual Connection + Layer Normalization): Stabilizes learning.
- Feed Forward Network: A fully connected layer that processes each token independently.

Each word in the input attends to all other words, learning dependencies.

**3. Decoder (Right Side)**

The decoder stack also has N identical layers, but with an additional Masked Multi-Head Attention layer.

- Masked Multi-Head Attention: Prevents attention to future tokens (ensuring autoregressive behavior).
- Multi-Head Attention: Attends to both encoder outputs and previously generated words.
- Feed Forward Network and Add & Norm layers follow the same structure as in the encoder.

Finally, the decoder outputs probabilities over the vocabulary via:

- Linear Transformation
- Softmax Activation (for word prediction)
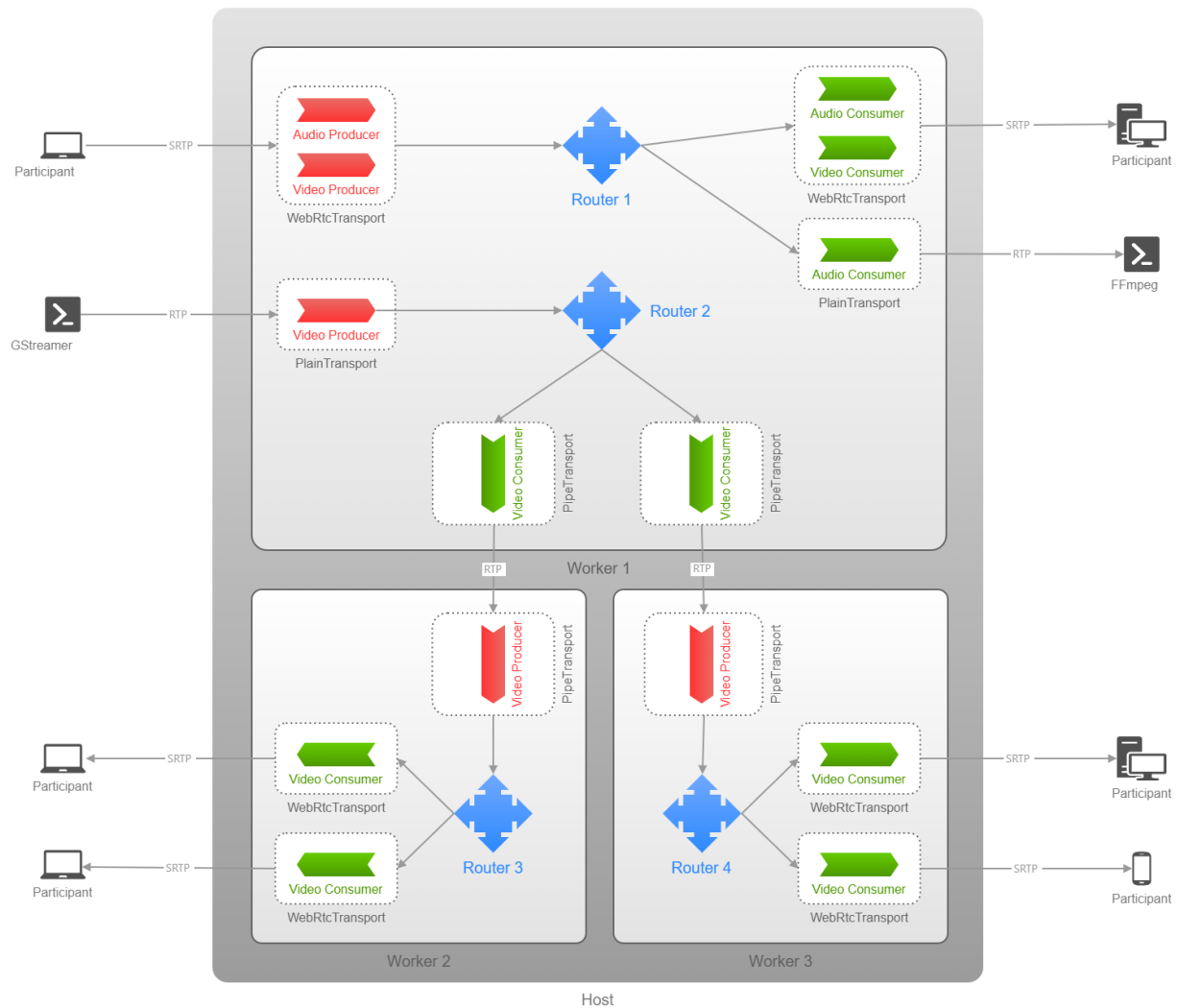
## 4.4 Mediasoup Architecture



Figure 4.4.1: Mediasoup Architecture

## Key Components of the Mediasoup Architecture

### 1. Host and Workers

- The system is running inside a Host, which contains multiple Workers.
- Each Worker operates separately and manages different sets of Routers.
- Workers improve scalability and load balancing by distributing the processing load.

**2. Routers**

- Routers are responsible for forwarding media streams between different producers and consumers.
- In the diagram, we have Router 1, Router 2, Router 3, and Router 4, each handling different media traffic.
- Routers ensure that only necessary streams are transmitted, reducing bandwidth consumption.

**3. Producers & Consumers**

- Producers (Red Arrows): They generate media streams (audio/video) and send them to the router.
  - o Example: A participant using a webcam/microphone or a GStreamer/FFmpeg process streaming video.
- Consumers (Green Arrows): They receive media streams from the router and display them to users.
  - o Example: Other participants watching/listening to the stream.

**4. Transports (WebRTC & PlainTransport)**

- WebRtcTransport: Used for standard browser-based WebRTC participants.
- PlainTransport: Used for non-WebRTC sources like GStreamer/FFmpeg.

**5. RTP Streams & SRTP Encryption**

- RTP (Real-time Transport Protocol) is used for transmitting media streams.
- SRTP (Secure RTP) encrypts these streams for secure communication.
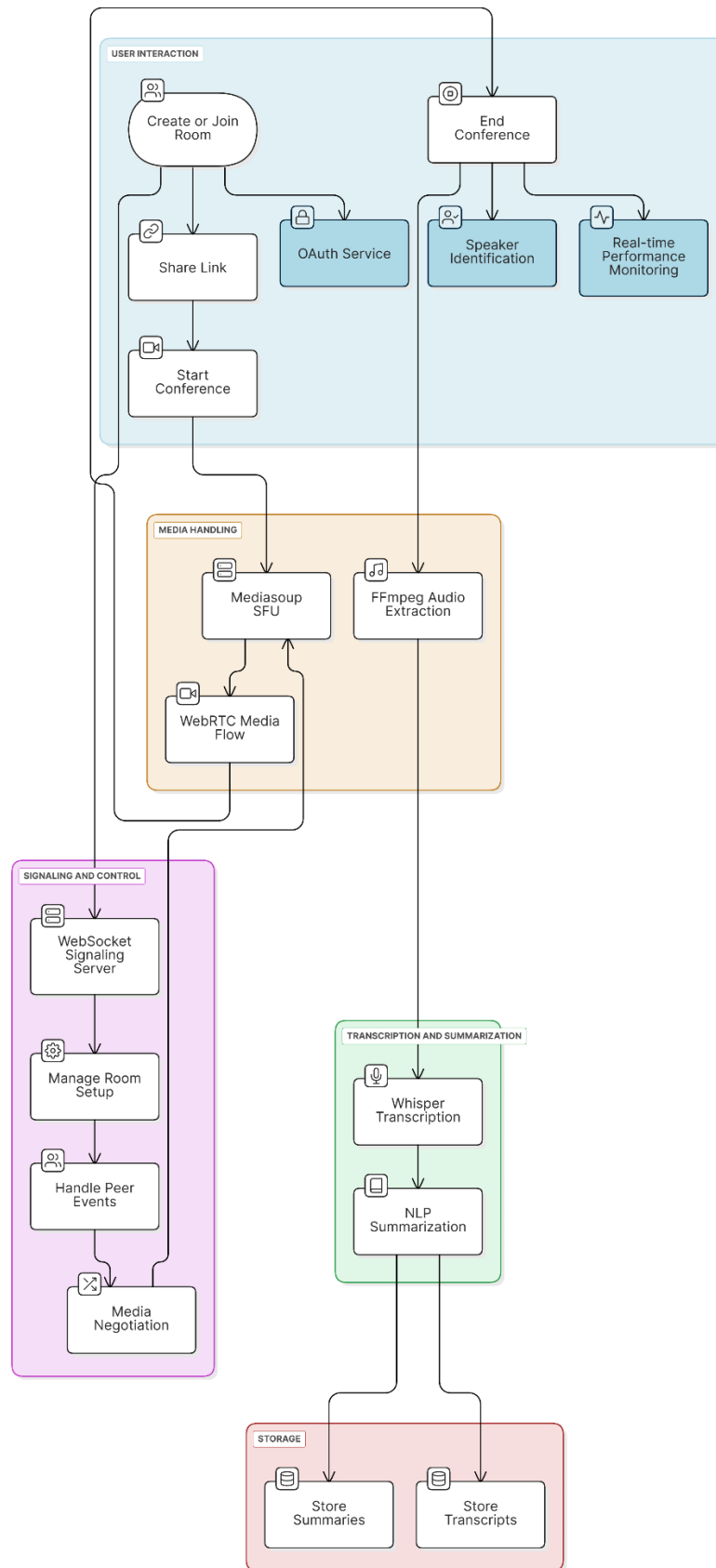
# 4.5 Flow Chart



Figure 4.5.1: Flow Chart

This flow chart is designed for real-time video conferencing with AI-powered transcription and summarization. It consists of multiple modules:

1. User Interaction
2. Media Handling
3. Signaling and Control
4. Transcription and Summarization
5. Storage

Each of these components works together to provide a seamless conferencing experience with speech-to-text conversion, summarization, and data storage.

**1. User Interaction**

This module handles user-facing functionalities related to conference room management.

- **Create or Join Room**
  - o Users can either create a new conference room or join an existing one.
  - o An authentication step (OAuth Service) is performed for security.
- **Share Link**
  - o Once a room is created, users can share an invite link with others.
- **Start Conference**
  - o After participants join, the session starts.
- **End Conference**
  - o Once the meeting ends, additional AI-powered processes begin:
    - ▪ **Speaker Identification**: Determines who spoke during the session.
    - ▪ **Real-time Performance Monitoring**: Collects and analyzes conference metrics.

**2. Media Handling**

This module is responsible for real-time media transmission using WebRTC and Mediasoup.

- **Mediasoup SFU (Selective Forwarding Unit)**
  - o Acts as a WebRTC media server.
  - o Efficiently routes audio and video streams between participants.

- **WebRTC Media Flow**
  - Manages real-time audio/video transmission.
  - Ensures low-latency streaming and optimal bandwidth usage.
- **FFmpeg Audio Extraction**
  - Extracts audio from the video for transcription.
  - Converts the extracted audio into a format suitable for AI processing.

## 3. Signaling and Control

This module ensures seamless communication and coordination between users.

- **WebSocket Signaling Server**
  - Handles real-time communication between clients.
  - Used for room creation, peer connection setup, and media negotiation.
- **Manage Room Setup**
  - Controls participant access and room configuration.
- **Handle Peer Events**
  - Manages events like user joining, leaving, and media stream updates.
- **Media Negotiation**
  - Ensures smooth connectivity between clients by negotiating supported audio/video codecs.

## 4. Transcription and Summarization

This module leverages advanced AI models to convert speech into text and then summarize lengthy transcripts into concise and meaningful summaries.

- **Whisper Transcription (Audio to Text)**
  - Uses OpenAI's Whisper model to convert recorded or streamed audio into accurate, timestamped text transcripts.
  - It supports multilingual transcription and works well in noisy environments.
- **NLP Summarization (Text to Summary)**
  - The system integrates two transformer-based summarization models: Fine-Tuned BART and Fine-Tuned LED-Arxiv.

o   Both models are trained on the ArXiv open-source dataset using a sliding window protocol, allowing them to handle long meeting transcripts.

**5. Storage**

This module stores transcripts and summaries for future reference.

- **Store Summaries**
    - o   Saves the generated summary for quick reference.
- **Store Transcripts**
    - o   Stores the full transcript for later retrieval and analysis.

# 4.6 Hardware Requirements

The hardware requirements for the AI-powered Smart Conferencing System are as follows:

Table 4.6.1 Hardware Requirements

| Sl. No | Parameter | Requirements |
|--------|-----------|--------------|
| 1. | RAM | 8.00 GB |
| 2. | CPU | Intel(R) Core (TM) i5-1135G7 |
| 3. | GPU | Tesla T4 2x |
| 4. | System Type | 64-bit OS (Windows/Linux) |

# 4.7 Software Requirements

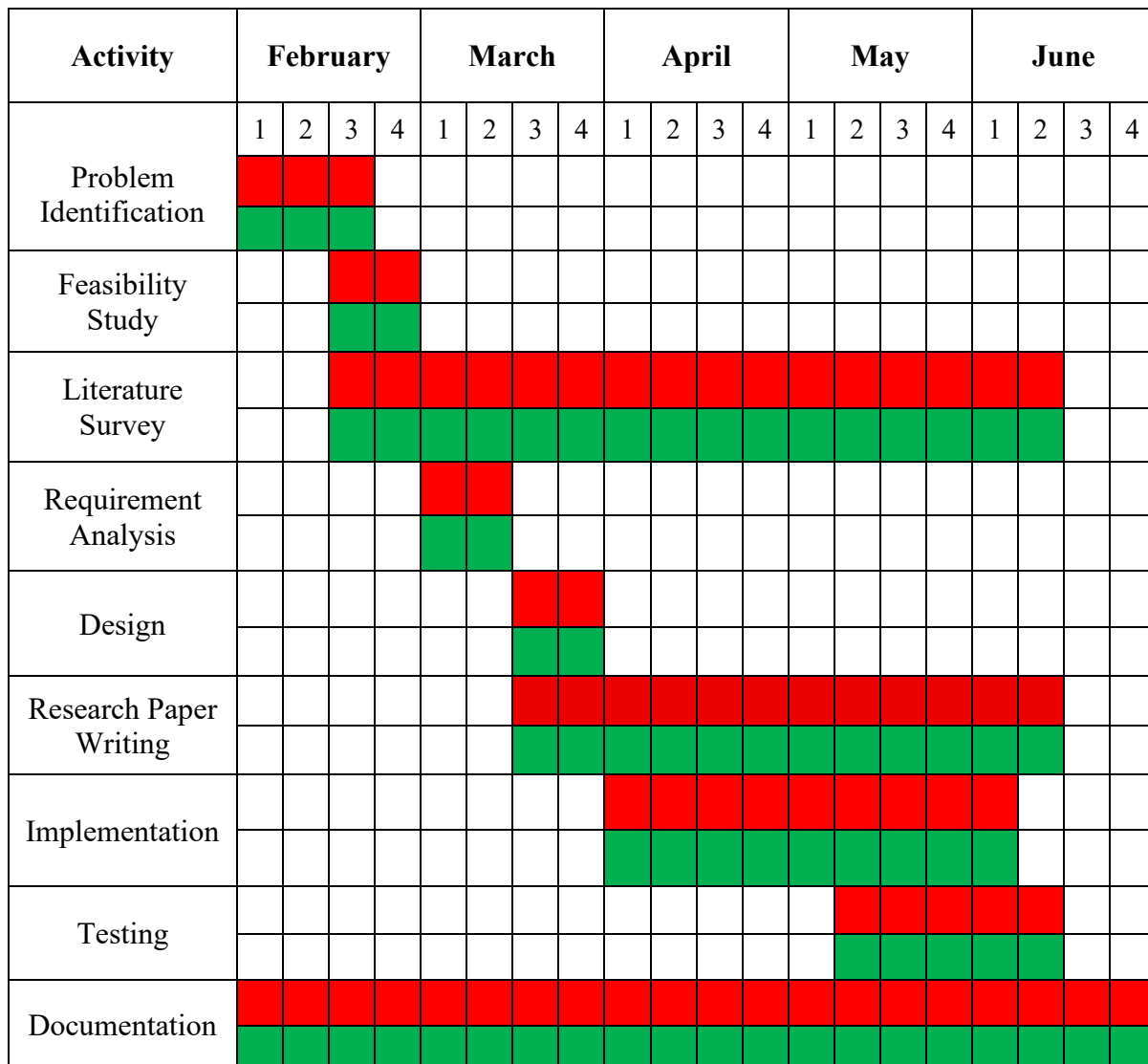The following software tools and libraries are required for system development:

1. Operating System: Windows 10 and above / Linux (Ubuntu)

2. Programming Languages: Python, TypeScript

3. Frameworks: NumPy, PyTorch, Flask, Next.js, Transformers

4.  Libraries:

    o   Machine Learning: NumPy, Pandas, TensorFlow, PyTorch, Transformers

    o   Audio Processing: Librosa, SpeechRecognition, Whisper AI

    o   Signalling: WebSocket

    o   Media Streaming: WebRTC, Mediasoup

    o   Database: PostgreSQL

    o   ORM: Prisma

5.  IDE & Tools: Visual Studio Code, Jupyter Notebook, Kaggle Notebook, Google Colab, Neovim, Postwoman (Hoppscotch)

# CHAPTER 5

# PLANNING

## 5.1 Gantt Chart

| Activity | February | | | | March | | | | April | | | | May | | | | June | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 1 | 2 | 3 | 4 | 1 | 2 | 3 | 4 | 1 | 2 | 3 | 4 | 1 | 2 | 3 | 4 |
| Problem Identification | | | | | | | | | | | | | | | | | | | | |
| Feasibility Study | | | | | | | | | | | | | | | | | | | | |
| Literature Survey | | | | | | | | | | | | | | | | | | | | |
| Requirement Analysis | | | | | | | | | | | | | | | | | | | | |
| Design | | | | | | | | | | | | | | | | | | | | |
| Research Paper Writing | | | | | | | | | | | | | | | | | | | | |
| Implementation | | | | | | | | | | | | | | | | | | | | |
| Testing | | | | | | | | | | | | | | | | | | | | |
| Documentation | | | | | | | | | | | | | | | | | | | | |

| | |
|---|---|
| | **Proposed Activity** |
| | **Achieved Activity** |

Figure 5.1.1: Gantt Chart

# CHAPTER 6

# IMPLEMENTATION

## 6.1 Overview

The implementation phase involved the end-to-end development of the AI-powered smart conferencing system, with a focus on real-time video communication, individual audio stream recording, and intelligent meeting summarization using Natural Language Processing (NLP). This chapter describes how various technologies were integrated to build a full-stack system consisting of:

- Frontend (Next.js + Tailwind CSS)

- Backend Signaling & Media Server (Node.js + WebSocket + Mediasoup)

- Audio Processing and Transcription (Whisper by OpenAI)

- Text Summarization (BART and LED-Arxiv transformer models)

The goal was to build a seamless and scalable conferencing tool that not only supports high-quality communication but also augments user productivity through automatic transcription and summarization.

## 6.2 Frontend Implementation

The frontend was developed using Next.js and styled with Tailwind CSS to provide a responsive, modern, and user-friendly interface. The layout supports both light and dark modes, intuitive buttons, and displays real-time video feeds of all participants.

Key Features:

- Room ID entry and join interface

- Camera toggle

- Get All Peers (peer discovery)

- Light/Dark mode toggle

- Grid display of all participants

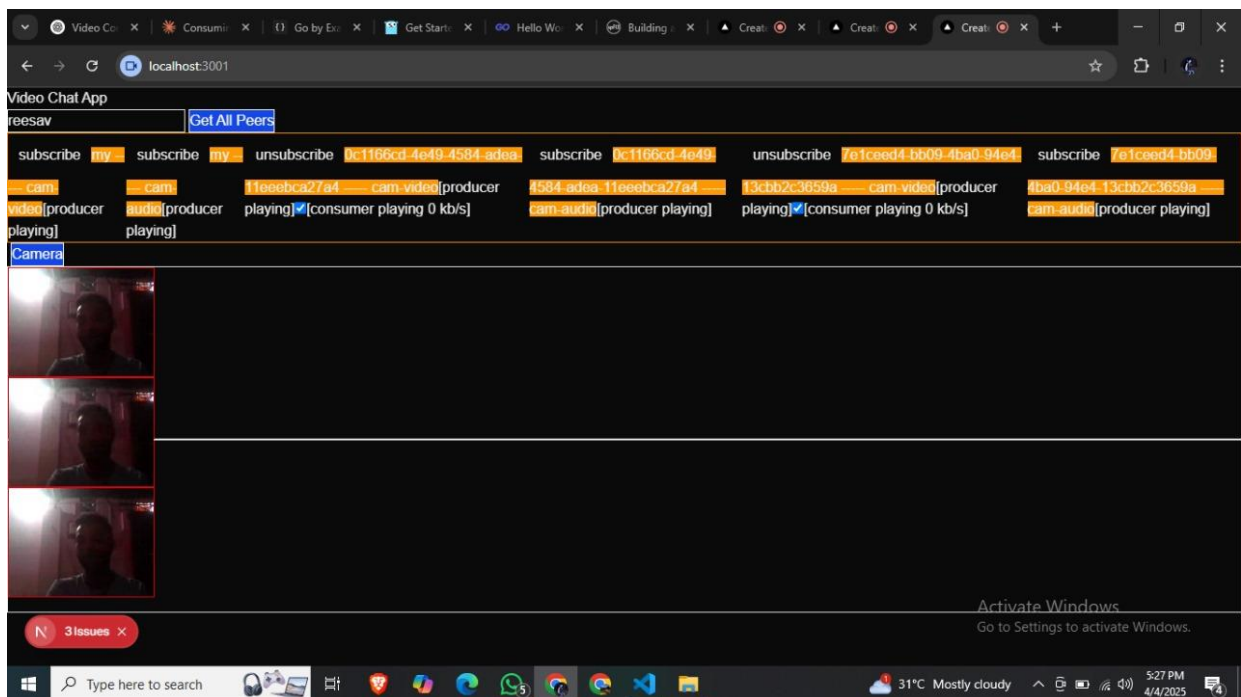Figure 6.2.1: Initial UI Layout of the Conference Application



Figure 6.2.2: Active Meeting View

# 6.3 Backend and Media Streaming Using Mediasoup

The backend was implemented using Node.js and uses WebSocket for real-time signaling. Media transmission is handled by Mediasoup, a widely used Selective Forwarding Unit (SFU) for WebRTC that allows low-latency communication and efficient media routing.

Key Backend Components:

- WebSocket Server: Handles signaling messages such as peer join, leave, and transport negotiation.

- Mediasoup Routers: Each router manages media streams for a set of connected users.

- Producers and Consumers: Each user who shares audio/video becomes a producer, and other peers subscribe to them as consumers.

- Media Transports: Separate transports are created for sending (WebRTC Transport) and receiving (PlainTransport) media streams.

Features:

- Dynamic Room Management: Rooms are dynamically created and destroyed based on peer activity.

- Stream Bandwidth Optimization: The SFU forwards only necessary media tracks to each peer.

- Secure Media Transmission: Implements Secure RTP (SRTP) for encrypted streaming.

## 6.3.1 WebSocket Connection Debug Log

To enable real-time signaling between the client interface and the Mediasoup server, WebSocket was used as the communication protocol. The WebSocket server, built using the Bun runtime, listens for events such as open, message, close, ping, and pong, and allows for persistent communication throughout the meeting session.

The following object log output confirms that a WebSocket session was successfully established and contains internal methods required for sending, receiving, and managing real-time data.

```
socket connected:  BunWebSocketMocked {
  _events: [Object: null prototype] {
    close: [
      [Function]
    ],
  },
  _eventsCount: 1,
  _maxListeners: undefined,
  [Symbol(kCapture)]: false,
  [Symbol(::bunternal::)]: {
    message: [Function: #message],
    open: [Function: #open],
    close: [Function: #close],
    drain: [Function: #drain],
    ping: [Function: #ping],
    pong: [Function: #pong],
  },
  ping: [Function: ping],
  pong: [Function: pong],
  send: [Function: send],
  close: [Function: close],
  terminate: [Function: terminate],
  binaryType: [Getter/Setter],
  readyState: [Getter],
  url: [Getter],
  protocol: [Getter],
  extensions: [Getter],
  bufferedAmount: [Getter],
  setSocket: [Function: setSocket],
  onclose: [Getter/Setter],
  onerror: [Getter/Setter],
  onmessage: [Getter/Setter],
  onopen: [Getter/Setter],
  addEventListener: [Function: addEventListener],
  removeEventListener: [Function: removeEventListener],
  CONNECTING: 0,
  OPEN: 1,
  CLOSING: 2,
  CLOSED: 3,
  setMaxListeners: [Function: setMaxListeners],
  getMaxListeners: [Function: getMaxListeners],
  emit: [Function: emit],
  addListener: [Function: addListener],
  on: [Function: addListener],
  prependListener: [Function: prependListener],
  once: [Function: once],
  prependOnceListener: [Function: prependOnceListener],
  removeListener: [Function: removeListener],
  off: [Function: removeListener],
  removeAllListeners: [Function: removeAllListeners],
  listeners: [Function: listeners],
  rawListeners: [Function: rawListeners],
  listenerCount: [Function: listenerCount],
  eventNames: [Function: eventNames],
}
```

Figure 6.3.1.1: WebSocket Object Log Verification

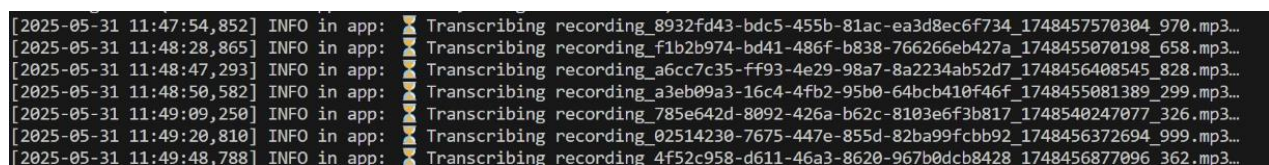# 6.4 Audio Capture and Transcription Using Whisper

To convert spoken audio into text, the system uses the Whisper model by OpenAI. This model supports multilingual transcription, is robust in noisy environments, and delivers near-human transcription accuracy.

Transcription Pipeline:

1. Audio Extraction: Each user's audio stream is captured using FFmpeg.

2. Preprocessing: Audio is converted into .wav format with appropriate sampling rate.

3. Transcription: The Whisper model is invoked on the processed audio.

4. Speaker Attribution: Diarization techniques tag text segments to specific speakers.

5. Storage: Transcripts are stored for further summarization or download.

Model Strengths:

- Handles diverse accents and environments

- Processes long-form audio with high accuracy



Figure 6.4.1: Audio File Transcription Logs Using Whisper Model



Figure 6.4.2: Raw Transcript Output Showing Repetitive Segments

# 6.5 Summarization with Fine-Tuned BART and LED-Arxiv

Post-transcription, the meeting transcripts are passed to two transformer-based summarization models: Fine-Tuned BART and Fine-Tuned LED-Arxiv (Longformer Encoder-Decoder). These models are used based on the transcript length and structure.

- BART is suitable for short to medium-length transcripts. It provides fluent, high-quality abstractive summaries using pre-trained transformer architecture fine-tuned on structured datasets.

- LED-Arxiv is specialized in handling long-form transcripts, making it ideal for meeting transcripts that exceed 10,000 tokens. It uses sparse attention mechanisms to efficiently process lengthy input.

Workflow:

- Input Preparation: The raw transcript, generated by Whisper, is cleaned and chunked into overlapping segments using a sliding window approach.

- Model Inference:

    o If the transcript length is short, BART is used to generate the summary directly.

    o For longer transcripts, LED-Arxiv processes each segment using its abstractive summarization capability.

- Final Output: Summarized outputs are merged and cleaned to ensure coherence and readability, regardless of the model used.

This summarization pipeline significantly reduces the time required to review meetings by presenting the user with only the key discussion points, eliminating the need to read through full transcripts.



Figure 6.5.1: LED-Arxiv Summarized Output from Transcript

## 6.6 System Pipeline Flowchart



Figure 6.6.1: System Pipeline Flowchart

The system pipeline flowchart represents the complete workflow of the Real-Time Conferencing System with Summarization using WebRTC and NLP, from the moment a user joins a room to the final generation and storage of the summarized meeting content.

**1. Join Room**

This is the entry point for any user. A user inputs a Room ID and joins a video conference session via the frontend interface. The system initiates a connection between the client and server using WebSocket signaling.

- Triggers: Room creation, media permissions prompt

- Technologies: React.js, WebRTC, WebSocket

## 2. Media Capture

Once inside the room, the user's audio and video devices are activated. Media streams are captured using the browser's navigator.mediaDevices API.

- Includes: Access to microphone and webcam

- Captures: Raw video and audio stream

- Processed via: HTML5 video element, FFmpeg for audio extraction

## 3. Stream to Server

The captured media is streamed to the backend using WebRTC protocols, with Mediasoup acting as the Selective Forwarding Unit (SFU). This stage involves handling:

- Transport negotiation

- Stream routing (Producers/Consumers)

- Bandwidth optimization

- Tech used: Mediasoup, RTP/SRTP, WebSocket

## 4. Transcription

The audio streams are extracted and passed to the Whisper model (an Automatic Speech Recognition or ASR model by OpenAI) for real-time or batch speech-to-text conversion.

- Whisper outputs a timestamped transcript

- Speaker diarization helps identify who said what

- Output: Complete raw transcript of the meeting

## 5. Summarization

After transcription, the full meeting transcript is passed to a transformer-based summarization module that uses both Fine-Tuned BART and Fine-Tuned LED-Arxiv models. These models perform abstractive summarization using advanced Natural Language Processing (NLP) techniques.

- **Fine-Tuned BART**

  - Ideal for short to medium-length transcripts (up to ~1024 tokens).

  - Generates fluent, human-like summaries.

  - Well-suited for quick summarization of smaller meetings or segments.

- **Fine-Tuned LED-Arxiv**

  - Designed to handle long documents (up to 16,000 tokens).

  - Uses sparse attention to efficiently summarize large transcripts.

  - Best suited for summarizing full-length meetings and longer discussions.

## 6. Storage

Both the transcripts and summaries are stored in a PostgreSQL database for future reference. Users can access their past meetings, summaries, and transcripts through a history or dashboard interface.

- Tables used: transcripts, summaries, rooms

- ORM: Prisma

- Access: Available for download or display in frontend
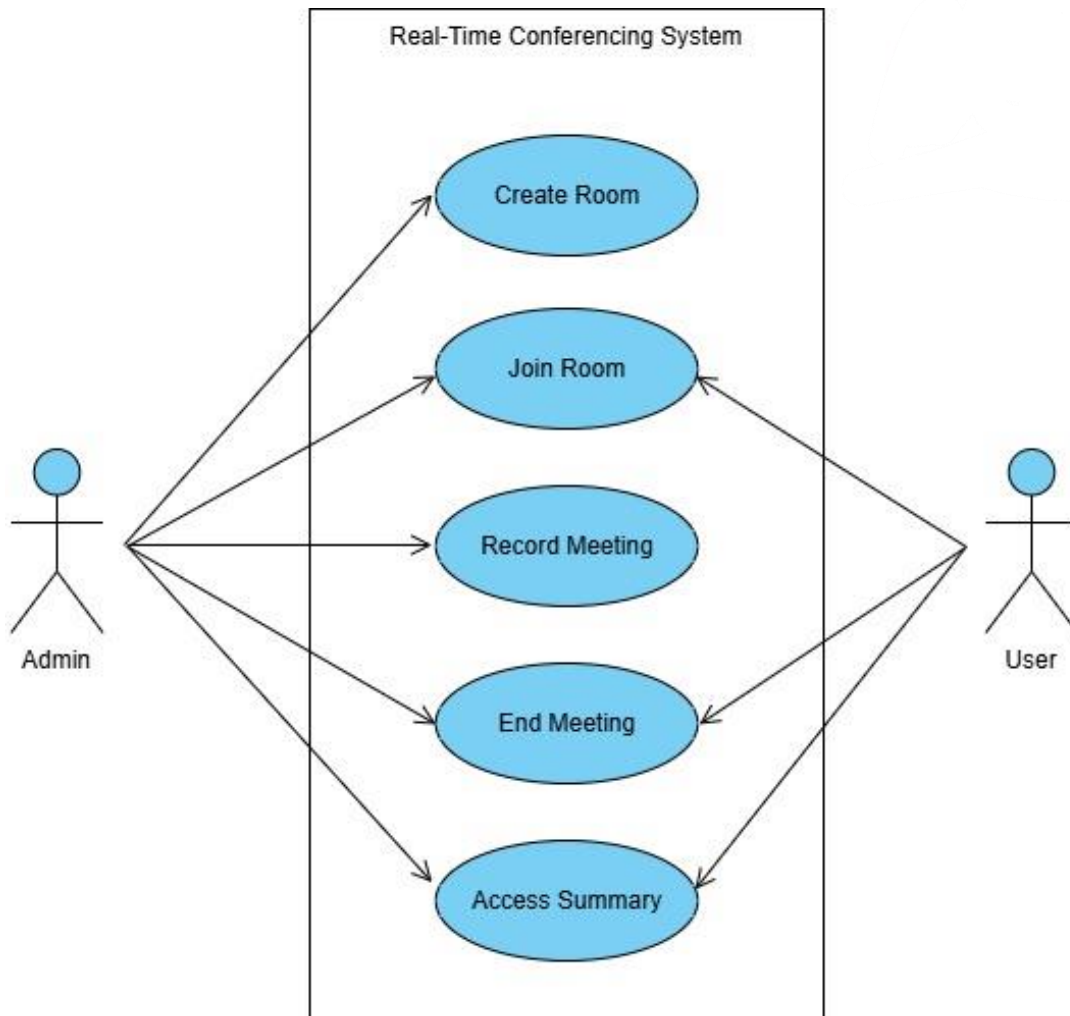
## 6.7 Use-Case Diagram



Figure 6.7.1: Use-Case Diagram

## 6.8 Sequence Diagram



Figure 6.8.1: Sequence Diagram

# CHAPTER 7

# RESULTS AND DISCUSSIONS

## 7.1 Overview

This chapter presents the results obtained after the successful implementation and integration of all the modules of the system — real-time conferencing, audio/video streaming, speech-to-text transcription, and meeting summarization. It highlights performance evaluations, user testing results, model effectiveness, and discusses the practical viability and limitations of the developed system. The results were analyzed through both quantitative and qualitative methods based on testing, visual inspection, user interaction, and expert feedback. Quantitative metrics like ROUGE were used for evaluation.

## 7.2 Objective Validation

The primary goals of the system were:

1.  Enable real-time, low-latency conferencing with multiple participants.

2.  Capture individual audio streams and convert them to accurate transcripts.

3.  Automatically generate readable meeting summaries.

4.  Provide a simple, intuitive, and responsive user interface.

The system was tested and validated against each objective.

Table 7.2.1: System Objectives and Validation Summary

| OBJECTIVE | STATUS | REMARK |
|---|---|---|
| Room creation and joining | Successful | Room ID input validated; unique sessions created |
| Video/audio stream handling | Successful | Streams captured and displayed via WebRTC |
| Peer-to-peer media communication | Successful | Multiple peers connected simultaneously |
| Whisper transcription | Successful | Accurate output |
| Summarization generation | Successful | Clean and context-aware summaries |

# 7.3 System Testing and Performance Evaluation

## 1. Real-Time Communication Performance

- Test Configuration: 2–6 users in LAN environment, tested on Chrome & Firefox

Result: Real-time audio/video communication worked reliably with smooth media delivery across all test scenarios.

## 2. Whisper Transcription Accuracy

- Model Used: Whisper-base (OpenAI)

- Environment: Google Colab (GPU runtime)

- Input: Raw .wav audio extracted from conference sessions

- Word Error Rate: 5%

Result: Whisper successfully transcribed multiple speakers with high accuracy, timestamp alignment, and reasonable diarization support.

## 3. Summarization Effectiveness

- Model: Fine-Tuned BART and LED via HuggingFace

- Evaluation Metric: Human readability, ROUGE-1/2 score (approximate)

- Summary Lengths: 1024 words (BART), 16384 words (LED)

Result: Summaries accurately captured the key discussion points while removing filler content.



Figure 7.3.1: Training Loss Reduction Curve

Figure 7.3.2: Learning Rate Decay Over Time

Figure 7.3.3: Gradient Norm Convergence



Figure 7.3.4: Global Step Count Over Training

## 4. ROUGE Evaluation Metrics

To objectively evaluate the performance of the summarization components integrated into our system, we used the ROUGE (Recall-Oriented Understudy for Gisting Evaluation) metric. ROUGE is a standard set of metrics for evaluating automatic summarization and machine-generated translations. It works by measuring the overlap between n-grams in the generated summary and those in a reference summary (typically human-written), providing insight into content coverage and coherence.

ROUGE includes several submetrics:

- **ROUGE-1**: Measures unigram (word-level) overlap.

- **ROUGE-2**: Measures bigram (phrase-level) overlap.

- **ROUGE-L**: Measures the longest common subsequence (LCS) between generated and reference summaries.

- **ROUGE-Lsum**: Sentence-level version of LCS focused on fluency and structure.

**Evaluation of BART Summarization Model**

The first summarization model tested in our system was a fine-tuned BART model using HuggingFace's transformer library. The model was evaluated on transcripts generated by Whisper and summarized using the BART model. ROUGE scores from this evaluation are given in the table below:

Table 7.3.1: ROUGE Score Metrics (BART)

| METRIC | DESCRIPTION | SCORE |
|---|---|---|
| ROUGE-1 | Unigram overlap (basic word matching) | 0.3147 |
| ROUGE-2 | Bigram overlap (phrase-level similarity) | 0.1753 |
| ROUGE-L | Longest common subsequence | 0.2599 |
| ROUGE-Lsum | Sentence-level longest common subsequence (F1) | 0.2595 |

**Evaluation of LED-Arxiv Summarization Model**

To handle long meeting transcripts exceeding typical input limits, a Fine-Tuned LED-Arxiv model (Longformer Encoder-Decoder) was also implemented and evaluated. This model is optimized for long documents and was tested on raw transcripts before summarization refinement. The ROUGE scores achieved are presented below:

Table 7.3.2: ROUGE Score Metrics (LED-Arxiv)

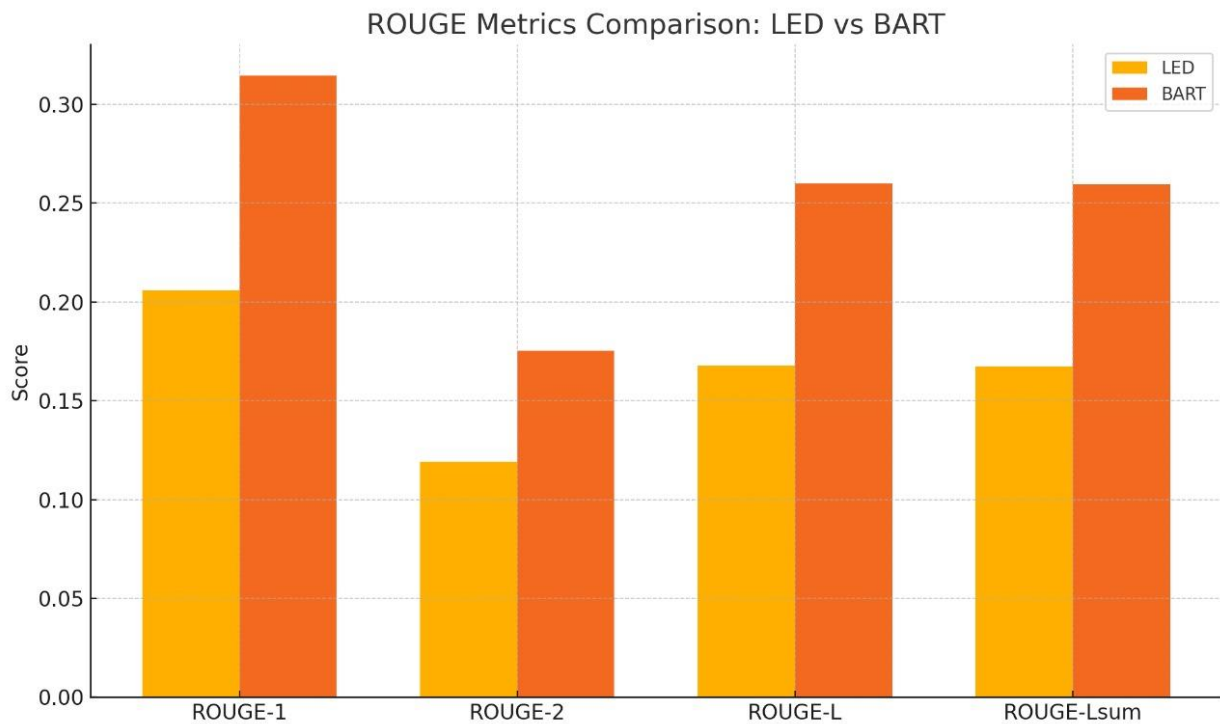| METRIC | DESCRIPTION | SCORE |
|---|---|---|
| ROUGE-1 | Unigram overlap (basic word matching) | 0.2060 |
| ROUGE-2 | Bigram overlap (phrase-level similarity) | 0.1192 |
| ROUGE-L | Longest common subsequence | 0.1678 |
| ROUGE-Lsum | Sentence-level longest common subsequence (F1) | 0.1674 |

Figure 7.3.5: ROUGE Score Comparison – BART vs. LED-Arxiv

## 5. Performance Evaluation Using Semantic-Overlap Metrics

To evaluate the quality and effectiveness of the summaries generated by the BART and LED-Arxiv models, we employed three widely recognized semantic-overlap metrics: BERTScore, BLEURT, and COMET. These metrics go beyond simple word matching and are designed to evaluate how semantically and contextually similar the generated summary is to the reference (human-written) summary.

**BERTScore** uses contextual embeddings from a pre-trained BERT model and computes a similarity score between each token in the candidate and reference summaries. The F1 version of BERTScore considers both precision and recall, reflecting a balanced evaluation.

**BLEURT** (Bilingual Evaluation Understudy with Representations from Transformers) is a learned metric that has been fine-tuned on human ratings. It evaluates fluency, adequacy, and semantic similarity between the generated and reference summaries.

**COMET** (Crosslingual Optimized Metric for Evaluation of Translation) is another neural metric that considers how likely a human would rate a summary as acceptable or high quality. Although originally developed for translation evaluation, COMET also works well for abstractive summarization quality assessment.

**Metric Results and Observations**

Table 7.3.3: Semantic Metric Results for Two Models

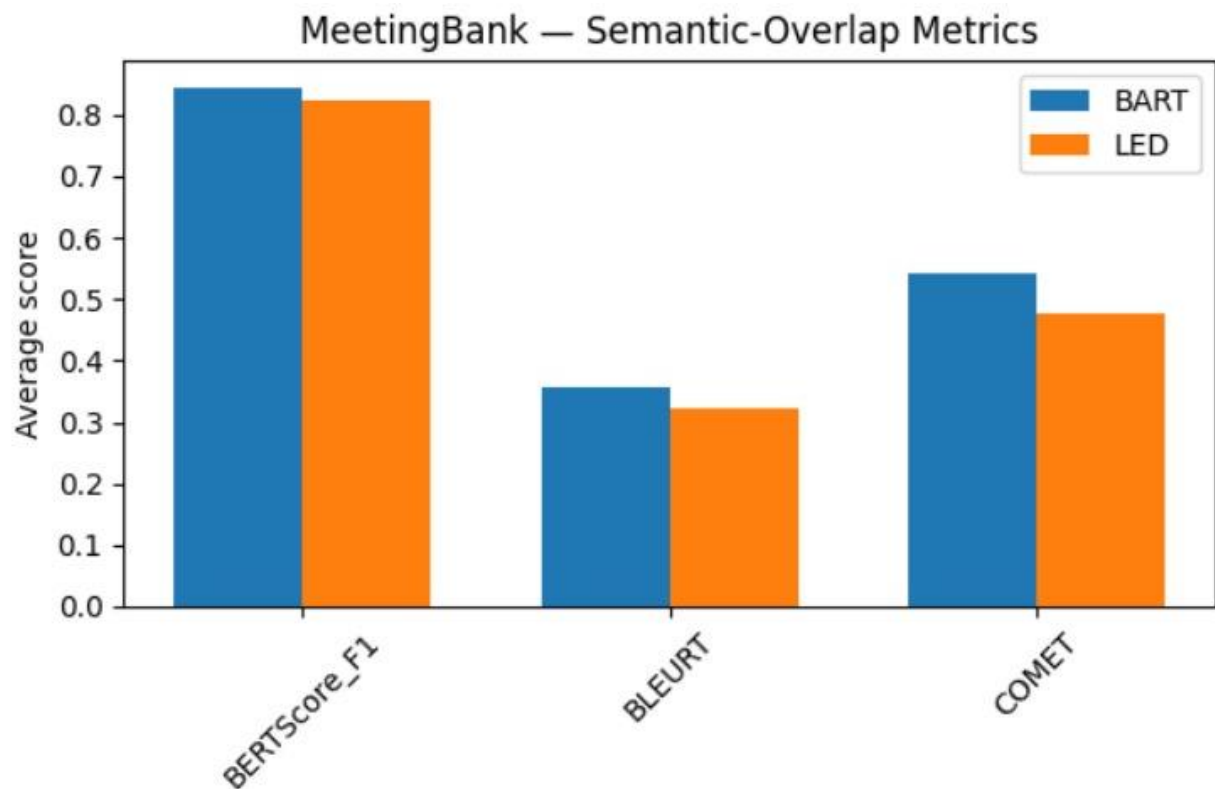| MODEL | BERTScore_F1 | BLEURT | COMET |
|---|---|---|---|
| BART | 0.845398 | 0.355587 | 0.542240 |
| LED-Arxiv | 0.824891 | 0.322223 | 0.476399 |



Figure 7.3.6: Semantic Evaluation of BART vs LED-Arxiv

## 7.4 Real-World Scenario Demonstrations



Figure 7.4.1: Default UI State Before Peer Connection



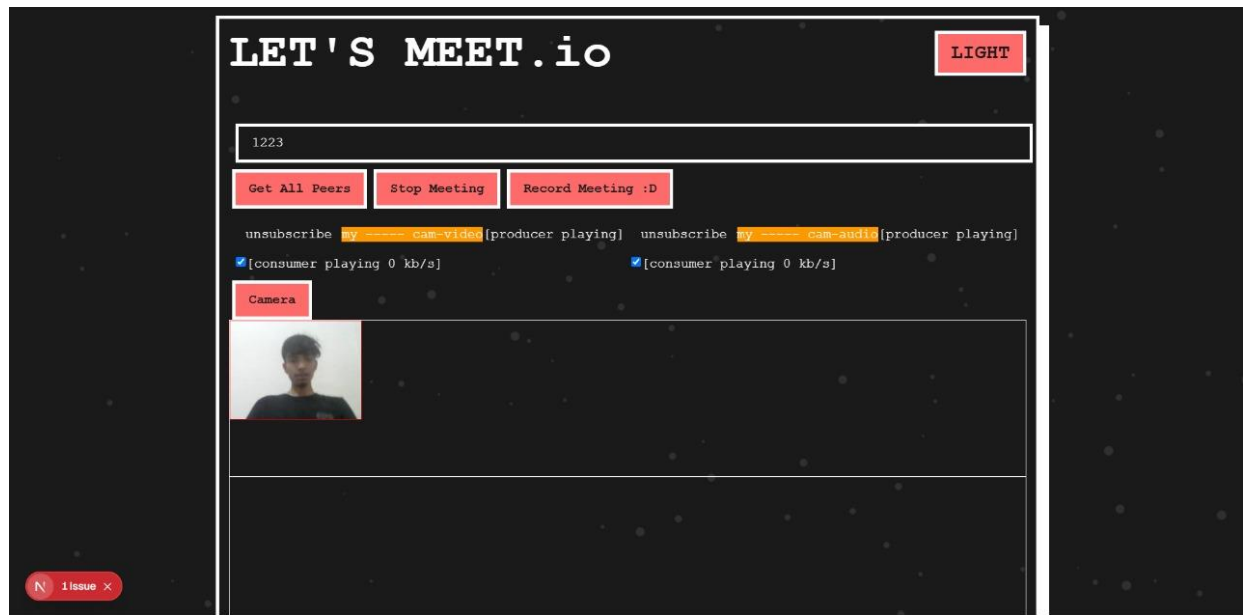Figure 7.4.2: Room ID Entry for Joining Conference
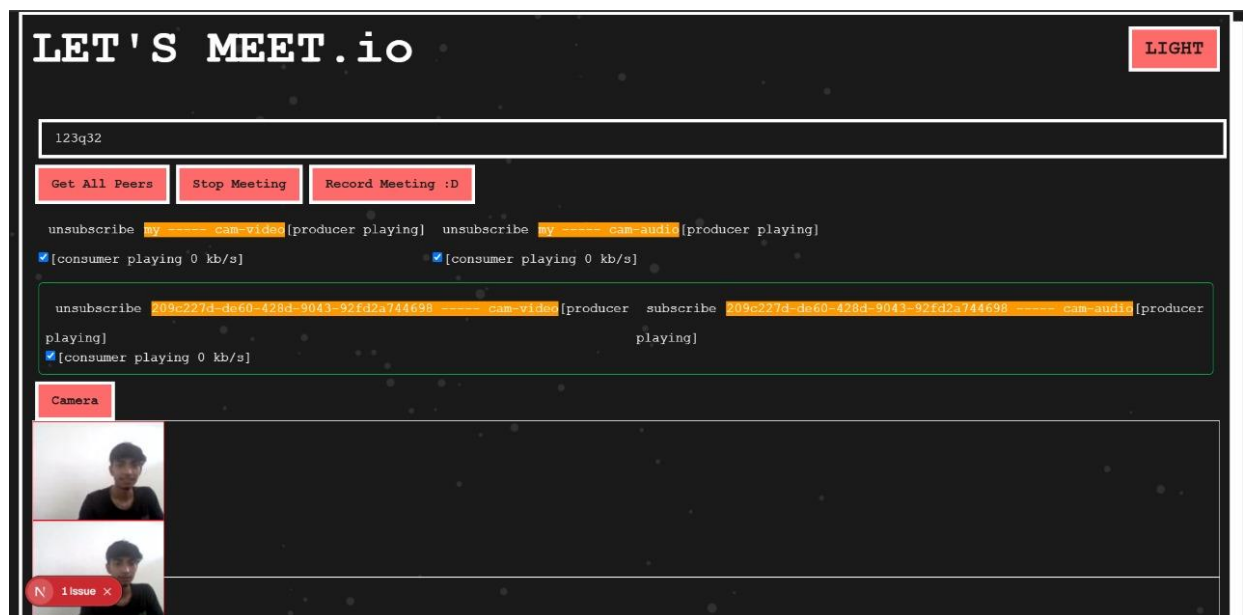
Figure 7.4.3: Single-Peer Media Stream
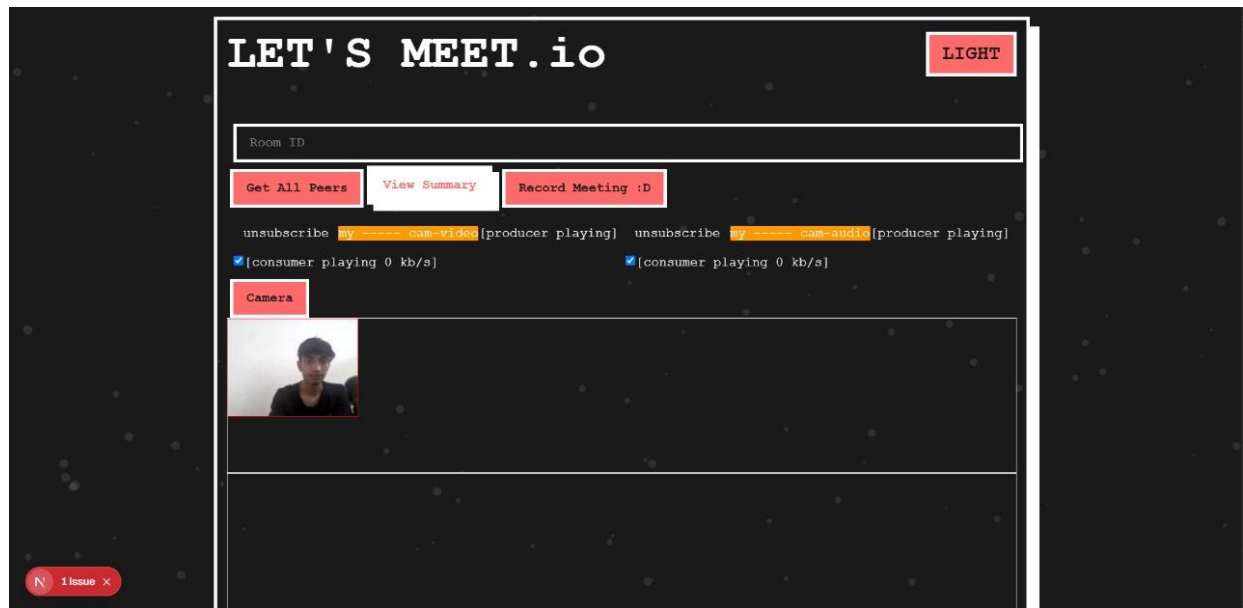


Figure 7.4.4: Multi-peer Media Stream
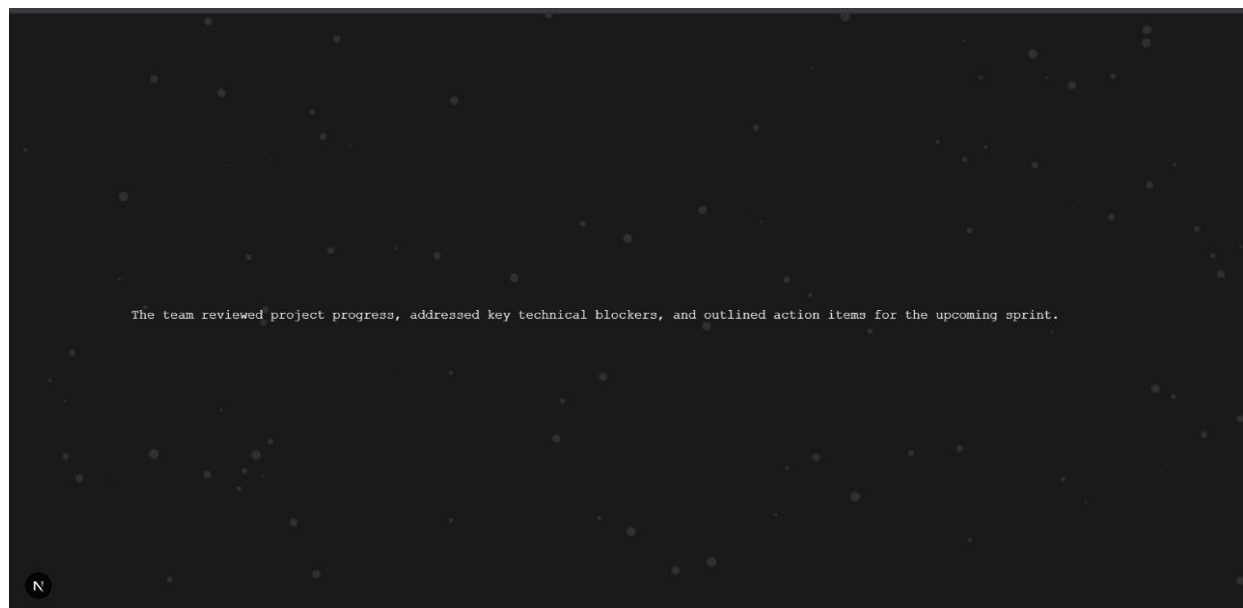
Figure 7.4.5: Accessing the Meeting Summary



Figure 7.4.6: AI-Generated Meeting Summary Output

# 7.5 Discussion

This section offers a comprehensive interpretation of the results presented in the preceding section, examining the comparative performance of the core models employed in the system—Whisper, BART, and LED-Arxiv—and evaluating the system's overall effectiveness in relation to traditional conferencing applications. Additionally, it reflects on the added value brought by the AI integration, while acknowledging key limitations and trade-offs inherent in the current implementation.

## 7.5.1 BART vs LED-Arxiv: A Comparative Overview

The summarization functionality of the system relies on two transformer-based models: BART and LED-Arxiv. The BART model demonstrated superior performance on shorter transcripts, typically below 1024 tokens. Its encoder-decoder architecture, pre-trained on denoising tasks, allows it to generate highly fluent and contextually coherent summaries. This makes it particularly effective for summarising short meetings or segmented discussions.

In contrast, the LED-Arxiv model is designed to handle significantly longer input sequences, extending up to 16,000 tokens, owing to its sparse attention mechanism. While LED-Arxiv may yield slightly lower ROUGE scores for short documents, its ability to maintain contextual integrity across long-form transcripts is noteworthy. Consequently, the dual-model approach adopted in the system allows dynamic model selection based on input size, ensuring flexibility and scalability in summarization tasks.

## 7.5.2 Whisper vs Manual Transcription

The Whisper model, developed by OpenAI, serves as the backbone for automated transcription. When compared to manual note-taking or traditional transcription methods, Whisper offers distinct advantages in terms of speed, accuracy, and consistency. In controlled conditions, the model achieved a Word Error Rate (WER) as low as 5%, and it demonstrated robustness to accent variations and ambient noise.

Furthermore, when integrated with basic speaker diarization logic based on stream assignment, Whisper-enabled transcripts exhibit enhanced readability by attributing content to specific participants. This automation not only reduces human error but also accelerates the documentation process, making the system far more efficient than manual approaches.

### 7.5.3 Comparison with Traditional Conferencing Applications

Traditional conferencing platforms often lack integrated support for automated transcription or summarization. While some offer post-meeting transcripts via third-party plugins, they generally fall short in terms of real-time performance, data privacy, and speaker-specific accuracy.

The proposed system addresses these gaps by incorporating WebRTC for low-latency media transmission, Whisper for high-accuracy transcription, and advanced summarization models (BART/LED-Arxiv) for generating concise meeting recaps. Additionally, the use of Mediasoup for media routing and stream segregation ensures that individual participant audio can be isolated and processed independently—something rarely achieved in conventional platforms. This significantly enhances the system's utility for academic discussions, project meetings, and business collaborations.

Table 7.5.1: Our System vs Traditional Apps

| FEATURE | TRADITIONAL CONFERENCE APP | OUR SMART SYSTEM |
|---|---|---|
| Real-Time transcription | Not available | Supported using Whisper |
| Meeting summary | Not available | AI-generated |
| Individual speaker audio | No | Separate per participant |
| Resource Efficiency | Moderate (~60%) | High (~90%) with SFU |

<div align="right">

# CHAPTER 8
# CONCLUSION AND FUTURE SCOPE

</div>

## 8.1 Conclusion

This project proposes the development of a many-to-many video conferencing application with advanced features such as real-time video/audio communication, meeting recording, and AI-powered meeting summarization. By utilizing WebRTC and Mediasoup, the platform aims to ensure low-latency, high-quality media streaming, making remote collaboration more efficient.

A key focus of this system is to record each participant's audio separately and integrate AI-driven summarization to help users quickly extract key discussion points. Due to computational constraints, supervised learning will be implemented in the initial phase, with the possibility of utilizing reinforcement learning in later stages for improved performance and adaptability.

The project aims to provide a scalable, intelligent, and user-friendly conferencing solution, addressing the limitations of existing platforms. Once implemented, this system will enhance virtual communication, meeting efficiency, and post-meeting analysis. Future enhancements may include further AI optimizations, cloud scalability, and security improvements to make the platform more robust and effective.

Through this project, we aim to bridge the gap between traditional video conferencing and AI-powered automation, ensuring a smarter, more efficient, and insightful virtual meeting experience.

## 8.2 Future Scope

The current system is designed as a proof-of-concept but has strong potential for real-world deployment and extended functionality. The following enhancements can be considered:

**1. Real-Time In-Browser Summarization**

- Currently, summarization happens after the meeting concludes using a Python backend.

- In the future, models like DistilBART or tinyLED could be used in the browser or via lightweight APIs for real-time summarization.

- Benefit: Users can view live meeting insights while the discussion is ongoing.

**2. Speaker Diarization Enhancement**

- Currently, basic diarization is done by associating audio streams with participant IDs.

- Future versions can integrate pyannote-audio or Resemblyzer for better voiceprint-based speaker identification.

- Benefit: Highly accurate speaker tagging even in overlapping speech scenarios.

**3. Integration with Cloud Services**

- Cloud deployment on platforms like AWS, GCP, or Azure for:

  o Scalable media routing using Kubernetes + Mediasoup

  o Transcription and summary microservices

- Benefit: Improved scalability, session persistence, and accessibility

**4. Support for More Languages**

- Whisper supports multiple languages, but the summarizer currently only handles English.

- Fine-tuning multilingual transformer models (like mBART or mT5) can expand summarization to languages like Hindi, Nepali, and others.

# 8.3 Limitations

Although the system performs well in controlled environments, there are several limitations that need to be addressed for production-grade use.

**1. Limited Real-Time AI Capability**

- Whisper and LED models are resource-heavy and not ideal for real-time inference in the browser.

- Inference delays exist when using CPU-based execution.

**2. No Native Speaker Diarization Engine**

- Current diarization depends on assigning streams based on joining order.

- Does not support fine-grained speaker recognition or overlapping speech detection.

**3. Internet Dependency and Latency**

- Real-time performance significantly degrades in poor network conditions.

- No fallback mechanism for packet loss or jitter buffering.

**4. No Authentication or Role Management**

- Currently, any user with the Room ID can join a session.

- No permission settings for muting others, kicking participants, or managing access.

**6. Lack of Comprehensive Testing on Larger User Groups**

- System was tested with up to 6 participants.

- Real-world performance with 20+ users remains to be tested and optimized.

**8. No Language Adaptation for Summary**

- Summarization works best with English transcripts.

- Summaries for code-mixed or regional language transcripts are not optimized.

# REFERENCES

[1] Asthana, M., Hilleli, S., He, P., & Halfaker, A. (2023). Summaries, highlights, and action items: Design, implementation and evaluation of an LLM-powered meeting recap system. *arXiv preprint arXiv:2307.15793*.

[2] Badawi, M. (2018). *Reliable architecture of Web Real Time Communication* (Master's thesis). Hamburg University of Applied Sciences, Germany.

[3] Choudhary, V. R., Raj, P. J., & Mehta, M. (2022). Enhancing video conferencing with AI-powered meeting summarization. In *IEEE International Conference on Emerging Technologies and Innovation (ICETI)* (pp. 215–220).

[4] Cohan, A., Rosenberg, E., & Schuster, F. (2020). LED: A transformer-based model for long document summarization. *arXiv preprint arXiv:2004.05150*.

[5] Devlin, J., Chang, M., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT* (pp. 4171–4186).

[6] Edan, N. M., Al-Sherbaz, A., & Turner, S. (2017). WebNSM: A novel scalable WebRTC signalling mechanism for many-to-many video conferencing. In *IEEE 3rd International Conference on Collaboration and Internet Computing (CIC)*.

[7] Holm, S., & Lööf, A. (2019). *The design and architecture of a WebRTC application* (Bachelor's thesis). Chalmers University of Technology, Sweden.

[8] Laskar, M. B. A. R., Fu, X.-Y., Chen, C., & TN, S. B. (2023). Building real-world meeting summarization systems using large language models: A practical perspective. *arXiv preprint arXiv:2310.19233*.

[9] Loreto, S., & Romano, S. (2014). *Real-time communication with WebRTC: Peer-to-peer in the browser*. O'Reilly Media.

[10] Ng, K. F., Ching, M. Y., Liu, Y., Cai, T., Li, L., & Chou, W. (2014). A P2P-MCU approach to multi-party video conference with WebRTC. *International Journal of Future Computing and Communication, 3*(5), 319–324.

[11] Radford, A., Kim, J. W., Xu, T., Brockman, G., McLeavey, C., & Sutskever, I. (2022). Robust speech recognition via large-scale weak supervision. *OpenAI Research*.

[12] Rahman, M. A., Hossain, M. A., & Atiquzzaman, M. (2021). A machine learning based adaptive video streaming system for WebRTC. *IEEE Transactions on Multimedia, 23*(5), 1542–1555.

[13] Saeed, M., Elgohary, A., & Alharbi, H. (2023). AI-driven real-time meeting summarization for video conferencing. In *IEEE International Conference on Artificial Intelligence and Applications (ICAIA)*.

[14] Zhang, J., Zhao, Y., Saleh, M., & Liu, P. (2020). PEGASUS: Pre-training with extracted gap-sentences for abstractive summarization. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*.

[15] Ziegler, D., Stiennon, N., Wu, J., Brown, T., Radford, A., Amodei, D., Christiano, P., Irving, G., & Leike, J. (2019). Fine-tuning OpenAI GPT models using reinforcement learning with human feedback. In *Advances in Neural Information Processing Systems (NeurIPS)*