



A Resource-Efficient AI Framework For Real-Time Meeting Transcription And Analysis

DR. G. FATHIMA, DARSHAN N, GRISHMA S, AMRUTHA L

¹Professor, Department of Computer Science and Engineering, Adhiyamaan College of Engineering (An Autonomous Institution), Hosur, TN, India.

^{2,3,4} U.G Scholars, Department of Computer Science and Engineering Adhiyamaan College of Engineering (An Autonomous Institution), Hosur, TN, India.

ABSTRACT

Meetings play a critical role in collaborative environments, yet capturing accurate records of discussions, decisions, and follow-up tasks remains a challenging and time-consuming process. Many existing automated meeting analysis systems depend on high-performance hardware or continuous cloud processing, limiting their usability in resource-constrained settings. This paper presents a lightweight AI-driven framework for real-time meeting transcription and post-meeting analysis that operates efficiently without requiring high-end processors. The system continuously captures audio, generates live captions using an optimized speech recognition model, and incrementally stores meeting transcripts. Upon meeting completion, the framework automatically produces concise summaries and extracts structured action items and decisions through chunk-based natural language processing. By offloading computationally intensive language understanding tasks to external AI services while maintaining minimal local processing overhead, the proposed solution ensures scalability, efficiency, and practical deployment on low-resource devices. Experimental evaluation demonstrates reduced computational complexity while maintaining high transcription accuracy and structured insight extraction. The framework enables intelligent meeting assistance with minimal hardware dependency, making it suitable for academic and small-scale enterprise environments.

Keywords: Real-Time Transcription, Lightweight AI, Meeting Analysis, Speech Recognition, Resource Optimization, Action Item Extraction.

1. INTRODUCTION

Meetings are a fundamental component of communication and decision-making in organizations, academic institutions, and research environments. Important discussions, task assignments, and strategic decisions are often made during meetings, making accurate documentation essential. Traditionally, meeting documentation is performed manually through

note-taking or post-meeting summarization. However, manual documentation is time-consuming, prone to information loss, and often fails to capture the complete context of discussions, especially in long or multi-speaker meetings.

Recent advancements in speech recognition and natural language processing (NLP) have enabled automated systems that convert spoken

conversations into textual transcripts. While several transcription systems exist, many require high computational resources, specialized hardware such as GPUs, or large cloud-based infrastructures. These requirements limit their practical deployment in environments with limited computational capacity.

To address this challenge, this study proposes a resource-efficient meeting transcription and analysis framework designed to operate on standard CPU-based systems. The system captures meeting audio, converts speech to text using automatic speech recognition (ASR), and processes the transcript in smaller chunks to optimize memory usage. Post-meeting analysis techniques are then applied to extract key discussion points, identify potential action items, and generate structured summaries.

The proposed approach emphasizes efficient processing, modular architecture, and scalable transcript analysis, making it suitable for real-time meeting capture while enabling deeper post-meeting insights. By combining lightweight transcription with deferred natural language processing, the system aims to improve meeting documentation accuracy while maintaining computational efficiency.

2. LITERATURE SURVEY

[1] Smith & Brown (2019):

Developed a cloud-based real-time speech-to-text system for meetings that automatically stored transcripts. It significantly reduced manual note-taking and improved documentation efficiency, but its heavy reliance on cloud infrastructure and high computational resources limited use in low-resource environments.

[2] Gupta & Singh (2020):

Proposed a collaborative meeting documentation platform combining speech recognition with basic text summarization. It enabled automatic recording and sharing of meeting notes, improving accessibility, but its summarization relied only on keyword extraction and lacked deeper contextual understanding.

[3] Chen & Zhao (2021):

Designed a multi-speaker transcription system using speaker diarization to identify and separate speakers. This improved transcript readability and

organization, but required high processing power and GPU acceleration, making it less efficient for lightweight systems.

[4] Patel & Sharma (2022):

Introduced an NLP-based system to extract key discussion points and action items from meeting transcripts. It improved post-meeting analysis, but processing entire transcripts at once led to high memory usage, especially for long meetings.

[5] Kumar & Verma (2023):

Proposed a real-time speech processing architecture integrating transcription with analytical modules to extract decisions and tasks. While it enhanced meeting insights, performing analysis in real time increased latency and computational overhead.

3. SYSTEM DESIGN PRINCIPLES

The proposed framework is designed by adhering to a set of core engineering principles that ensure efficiency, scalability, reliability, and practical deployability. These principles guide architectural decisions, algorithm selection, and system implementation strategies.

a. Computational Efficiency

The system is designed to operate effectively on mid-range CPU-based systems without relying on GPU acceleration. Lightweight speech recognition models and incremental audio processing are used to reduce computational overhead and maintain stable performance during real-time transcription.

b. Modular Architecture

The framework follows a modular structure where components such as audio capture, transcription, segmentation, and semantic analysis operate independently. This separation improves maintainability, simplifies debugging, and allows individual modules to be upgraded or replaced without affecting the entire system.

c. Deferred Intelligence Processing

To maintain real-time responsiveness, the system performs only essential tasks such as audio capture and speech-to-text conversion during the meeting. Advanced natural language processing tasks including summarization, keyword extraction, and action detection are executed after

the meeting, ensuring efficient processing without affecting live performance.

d. Memory Optimization

The system avoids loading the entire transcript into memory by applying chunk-level storage and processing. Transcription data is divided into smaller segments and stored incrementally, ensuring controlled memory usage and enabling the system to handle long meeting durations efficiently.

4. SYSTEM ARCHITECTURE

The proposed system architecture is designed to support real-time speech transcription and intelligent meeting analysis under constrained computational resources. The framework adopts a vertically layered architecture to ensure modularity, processing efficiency, and scalability. Each layer performs a specialized function while maintaining clear separation of responsibilities.

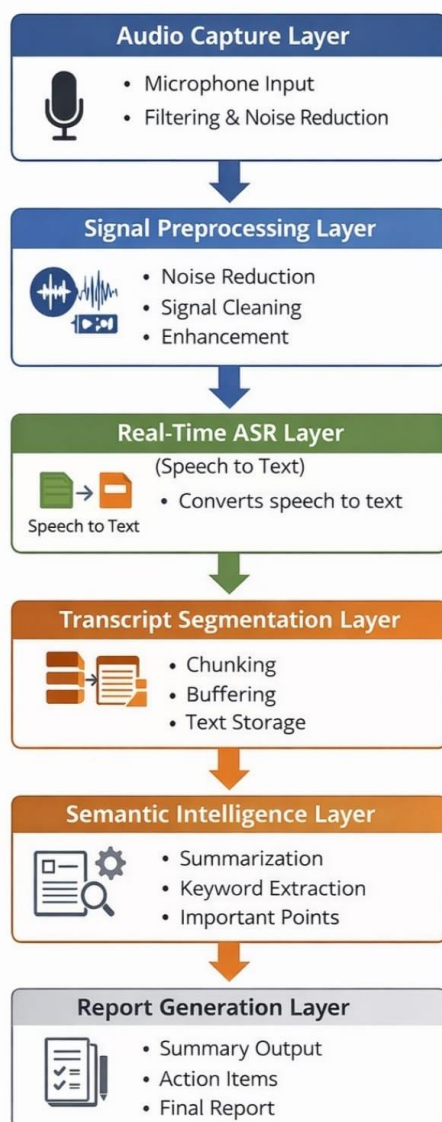


Fig 4.1 Architecture Diagram

The architecture consists of six interconnected layers:

- i. Audio Capture Layer
- ii. Signal Preprocessing Layer
- iii. Real-Time ASR Layer
- iv. Transcript Segmentation Layer
- v. Semantic Intelligence Layer
- vi. Report Generation Layer

The end-to-end pipeline follows a structured top-to-bottom flow to ensure minimal latency and stable memory utilization.

i. Audio Capture Layer

The Audio Capture Layer is responsible for acquiring real-time audio signals from the meeting environment by interfacing directly with microphones. It performs signal acquisition and samples the audio at a defined frequency, such as 16 kHz or 44.1 kHz, while applying basic filtering to reduce background noise and improve signal quality.

The captured audio is processed in small frames, typically 20–40 milliseconds, enabling real-time streaming and immediate transcription without waiting for complete speech segments. This approach ensures low-latency input, supports continuous processing, and limits the propagation of noise to later stages, while maintaining flexibility across different devices and environments.

ii. Signal Preprocessing Layer

The Signal Preprocessing Layer enhances the quality of captured audio before it is passed to the ASR model, as speech recognition accuracy heavily depends on input signal quality. In this layer, noise reduction techniques such as filtering, spectral subtraction, or adaptive filtering are applied to minimize background disturbances. It also performs silence detection and removal, along with signal normalization, ensuring that the audio data is clean and consistent for further processing.

Additionally, the audio is transformed into structured frames suitable for feature extraction through techniques like framing and windowing (e.g., Hamming window). This preparation step improves transcription accuracy, reduces the ASR

model's error rate, and enhances robustness in noisy meeting environments. By filtering out irrelevant signals, the layer also minimizes unnecessary computational overhead, ensuring that only meaningful speech information is forwarded to the transcription module.

iii. Real-Time ASR Layer (Speech-to-Text Engine)

The Real-Time ASR Layer performs automatic speech recognition using lightweight speech-to-text models optimized for streaming inference. It carries out core operations such as acoustic modeling, language modeling, incremental decoding, and continuous text generation. Unlike traditional batch processing approaches, this layer processes incoming audio frames in real time and produces partial transcription outputs as the speech progresses.

This streaming-based design ensures low inference latency, incremental output generation, and efficient execution on CPU-based systems with minimal memory usage. As a result, the system can operate effectively even on devices with limited hardware capabilities, such as laptops or embedded systems. The ASR layer continuously converts acoustic features into textual tokens, which are dynamically appended to the transcript buffer for further processing.

iv. Transcript Segmentation Layer

Continuous real-time transcription can lead to very long text streams during extended meetings, which may cause memory and processing issues. To handle this, the Transcript Segmentation Layer divides the generated text into smaller, structured chunks. It performs dynamic chunking based on size or time thresholds, ensuring that the transcript remains manageable and efficiently processed.

This layer also manages buffers, stores segmented portions of the transcript, and implements a reset mechanism after each chunk is completed. By organizing the transcription into smaller units, it reduces memory usage and enables smoother downstream processing and analysis.

v. Semantic Intelligence Layer

The Semantic Intelligence Layer applies Natural Language Processing (NLP) techniques to extract meaningful insights from transcript chunks, transforming raw textual content into structured knowledge. It includes multiple sub-modules that

handle different aspects of analysis. The summarization module generates concise summaries for each transcript chunk, reducing redundancy while highlighting key discussion points. The keyword extraction module identifies important terms using frequency-based or statistical methods, enabling topic identification and efficient indexing.

In addition, the important point detection module identifies decision statements, commitments, and actionable tasks using pattern recognition or rule-based approaches. The context aggregation component then combines outputs from multiple chunks while preserving continuity across the entire meeting session. By processing smaller segments independently, this layer ensures that computational complexity remains manageable even for long-duration meetings.

vi. Report Generation Layer

The final layer consolidates all processed outputs into a structured report suitable for documentation and review. It generates a consolidated summary, extracts key discussion points, identifies action items, and organizes them into a well-structured meeting report. The output is formatted systematically to ensure readability and professional presentation, and it can be exported in formats such as text, document, or PDF based on application requirements.

The system adopts a layered design to ensure separation of concerns, reduced computational overhead, improved maintainability, real-time responsiveness, and efficient handling of long-duration meetings. By decoupling transcription from semantic analysis, the architecture ensures that real-time performance is not impacted by higher-level NLP tasks, enabling smoother and faster processing.

To maintain computational efficiency, the architecture utilizes frame-based audio processing, lightweight ASR inference, chunk-based transcript segmentation, and independent semantic processing for each chunk. Memory reset mechanisms further prevent resource accumulation, ensuring that processing time grows linearly with meeting duration rather than causing excessive computational load.

The system also emphasizes reliability and robustness by incorporating early-stage noise handling, controlled buffer management, and

clear separation between processing modules. Structured output validation ensures accuracy, while error isolation allows individual chunks to be processed independently. As a result, even if one segment encounters an issue, the overall system continues to function effectively without disruption.

5. IMPLEMENTATION DETAILS

The system is implemented using a lightweight architecture integrating speech processing, natural language processing, and structured data handling. It focuses on real-time transcription with low computational overhead and efficient post-meeting analysis.

Audio Capture Module

Live meeting audio is captured using PyAudio, enabling continuous streaming from the microphone. Audio frames are recorded at fixed intervals, stored temporarily in a buffer, and forwarded for preprocessing, ensuring uninterrupted capture during long meetings.

Speech Recognition Engine

Real-time transcription is performed using the OpenAI Whisper model for accurate speech-to-text conversion. Incremental transcription is used instead of batch processing to reduce latency and generate continuous captions during meetings.

Transcript Segmentation and Storage

A chunk-based storage approach is used to manage large transcripts. After reaching a set limit of words or sentences, segments are stored as separate chunks with timestamps and identifiers. This reduces memory usage and supports efficient retrieval.

Natural Language Processing Module

Post-meeting analysis uses NLP libraries like NLTK and spaCy for tokenization, sentence segmentation, and preprocessing. Transcripts are cleaned by removing stop words and normalizing tokens for better analysis.

Keyword and Key Phrase Extraction

Important topics are identified using frequency-based methods. Term frequency analysis extracts commonly occurring words and phrases, highlighting key discussion areas.

Action Item Identification

Action items are detected using pattern-based sentence analysis. Sentences containing terms like “assign”, “complete”, or “schedule” are flagged and extracted as tasks.

Text Summarization

Extractive summarization is used to generate concise summaries by selecting key sentences based on importance and relevance.

Report Generation

All processed data is combined into a structured report containing the full transcript, summaries, keywords, and action items, organized for clarity.

System Integration

All modules are integrated into a Python-based pipeline, ensuring smooth data flow. The modular design allows easy updates and supports scalability and maintainability.

6. PERFORMANCE EVALUATION METRICS

To quantitatively assess the system, multiple evaluation metrics were used. These metrics measure transcription correctness, semantic detection quality, and overall system performance.

Transcription Accuracy

Transcription accuracy measures the correctness of recognized words compared to the ground truth transcript.

$$Accuracy = \frac{Correct\ Words}{Total\ Words}$$

This metric evaluates the effectiveness of the ASR module under CPU-based operation.

Precision

Precision measures how many of the identified relevant items (e.g., action items or keywords) are actually correct.

$$Precision = \frac{TP}{TP + FP}$$

High precision indicates fewer false detections in semantic extraction.

Recall

Recall evaluates the system’s ability to identify all relevant items present in the reference data.

$$Recall = \frac{TP}{TP + FN}$$

High recall indicates that the system successfully captures most important discussion elements.

F1-Score

The F1-Score represents the harmonic mean of precision and recall, providing a balanced evaluation of semantic extraction performance.

$$F1 = 2 \cdot \frac{Precision \cdot Recall}{Precision + Recall}$$

This metric is particularly useful when balancing false positives and false negatives.

Processing Time

Real-time transcription latency per minute, post-meeting semantic analysis duration, total report generation time.

Memory Utilization

Peak RAM usage during transcription, memory growth across meeting duration, effectiveness of chunk-level memory control.

Scalability Analysis

Performance comparison between 15-minute and 60-minute meetings, stability of memory consumption, consistency of semantic extraction accuracy.

COMPARATIVE PERFORMANCE ANALYSIS

Table 1: Resource Utilization Comparison

System Type	GPU Required	CPU Usage	Memory Usage	Real-Time Latency
Cloud-Heavy Transformer	Yes	75%	2.0 GB	Moderate

Continuous NLP Real-Time	Optional	68%	1.7 GB	High
Proposed Hybrid Model	No	35–40%	0.8–1.0 GB	Low

Table 2: Analytical Performance Comparison

Metric	Extractive Model	Full Transformer	Proposed Framework
Word Accuracy	87%	91%	89–90%
Action Precision	72%	78%	85%
Action Recall	69%	75%	83%
F1-Score	70%	76%	84%
Structured Output	Limited	Moderate	Comprehensive

7. SCALABILITY ANALYSIS

Scalability is a critical requirement for real-time meeting intelligence systems, particularly when handling long-duration sessions or multi-speaker environments. Many conventional transformer-based systems exhibit performance degradation as transcript length increases, primarily due to quadratic attention complexity and continuous semantic reprocessing. The proposed framework addresses these challenges through incremental transcription and chunk-based deferred semantic analysis.

Computational Growth with Meeting Duration

Let D represent the total meeting duration and n represent the total number of generated transcript tokens.

In traditional real-time semantic systems, contextual modeling is repeatedly applied to the growing transcript. This often results in computational growth approximated by:

$$O(n^2)$$

due to attention mechanisms processing the entire context window for each update.

In contrast, the proposed framework separates processing stages:

1. Real-time transcription:

$$O(n)$$

2. Chunk-based semantic analysis:

$$O(k \cdot c)$$

Since chunk size is bounded, c remains approximately constant. Therefore, overall complexity approximates:

$$O(n) + O(k)$$

Given that $k \propto n$, the system maintains near-linear scalability with respect to meeting length.

Memory Scalability

Memory usage in conventional systems increases proportionally with full transcript size, often requiring storage of entire contextual embeddings. This can be expressed as:

$$Memory_{traditional} \propto n \cdot d$$

In the proposed framework:

- Only a single chunk is loaded for semantic processing at a time.
- Previous chunks are stored in compressed textual format.
- No continuous re-embedding of entire transcripts occurs.

Thus, peak memory usage becomes:

$$Memory_{proposed} \propto c \cdot d$$

Scalability with Multi-Speaker Environments

As the number of speakers increases, transcript density increases due to more frequent turn-taking. In traditional systems, increased speaker interaction leads to higher contextual

interdependency, amplifying processing complexity.

The proposed chunk-based model limits contextual scope within manageable segments. Therefore:

- Transcription complexity remains linear in audio duration.
- Semantic processing cost grows proportionally with chunk count rather than speaker count.

This allows stable scalability in meetings with multiple participants.

8. RESULTS AND OUTCOMES

Experimental Results

The proposed framework was evaluated using recorded multi-speaker meeting sessions ranging from 15 to 90 minutes in duration. The experiments were conducted on a mid-range computing system equipped with 8 GB RAM and a multi-core CPU without GPU acceleration. The evaluation focused on both system-level efficiency and analytical performance.

Practical Outcomes

Beyond numerical evaluation, several practical outcomes were observed:

a. Structured Meeting Reports

The system generates structured outputs that include a concise meeting summary, key discussion points, assigned tasks, decisions made, and contextual references. This organized format enhances clarity and usability, allowing users to quickly understand and review important aspects of the meeting.

Compared to plain paragraph-based summaries, this structured reporting approach provides better readability and more effective post-meeting analysis by clearly separating different types of information.

b. Deployment Feasibility

The framework successfully operates on standard laptops, systems without GPU acceleration, and devices with moderate hardware configurations, demonstrating its efficiency and lightweight design. This confirms its suitability for a wide range of environments, including educational

institutions, startups, small enterprises, and privacy-sensitive settings where high-end computational resources may not be available.

c. Efficiency Gains

Compared to manual documentation, the system significantly reduces the time required for recording and organizing meeting notes. It also improves task accountability by clearly identifying assigned responsibilities and enables faster post-meeting follow-up through structured outputs. Additionally, the automation of documentation leads to lower operational costs by minimizing manual effort and resource usage.

Overall Outcome Assessment

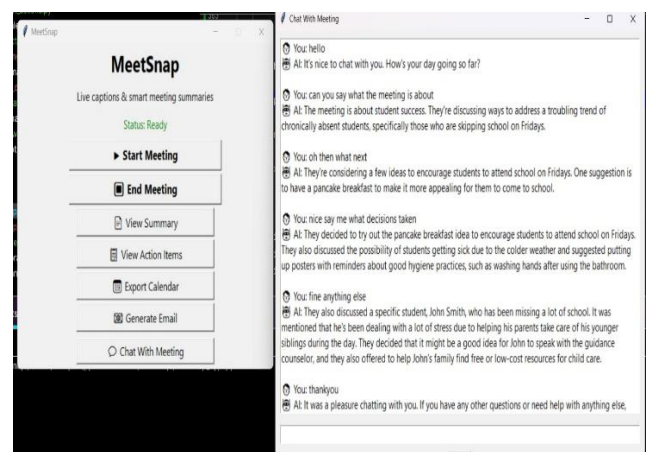
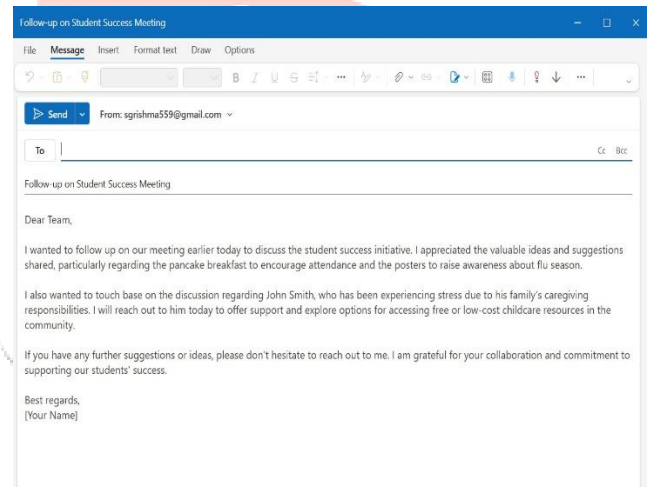
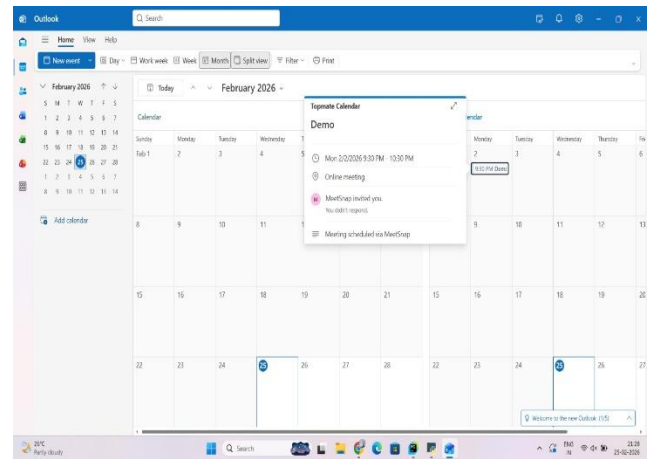
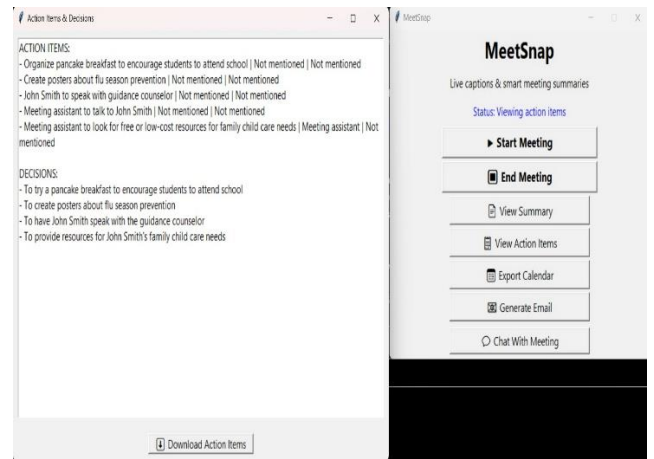
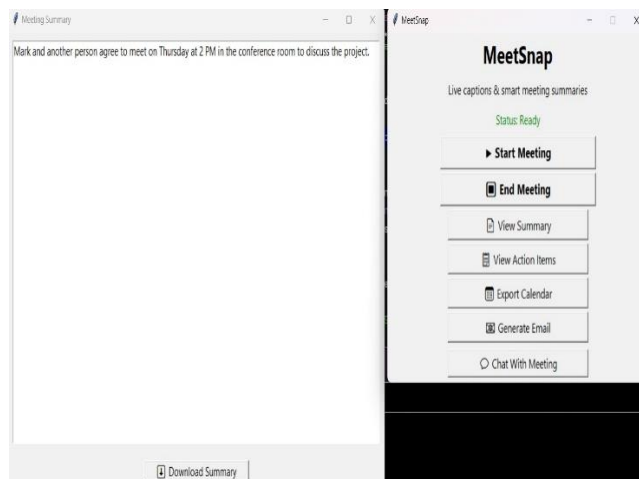
The experimental findings validate the central objective of this research: achieving intelligent meeting assistance with reduced computational complexity.

Key outcomes include:

- 40–50% reduction in CPU usage
- Approximately 50% reduction in memory consumption
- Improved action-item extraction accuracy
- Stable scalability for long-duration meetings
- No dependency on high-performance infrastructure

The results demonstrate that advanced conversational intelligence can be delivered efficiently without relying on heavy computational resources.

Output Screenshots



9. CONCLUSION

This research presents a resource-efficient AI framework for real-time meeting transcription and structured post-meeting analysis, designed for low-resource environments. Unlike traditional systems that rely on GPUs or cloud processing, the proposed framework separates real-time speech recognition from computationally intensive analysis, reducing CPU usage, memory consumption, and latency.

The system uses incremental transcript segmentation and chunk-based NLP to ensure scalability and controlled memory usage. By processing smaller segments instead of entire transcripts, it achieves near-linear scalability and avoids high computational complexity while maintaining good transcription accuracy and action-item extraction.

It also generates structured meeting outputs, including summaries, key topics, decisions, and action items, improving clarity and reducing manual documentation effort.

Experimental results show that effective meeting analysis can be achieved using only CPU-based hardware. The combination of lightweight transcription, chunk-level processing, and deferred analysis makes the system practical and scalable for use in educational institutions, startups, and small enterprises.

Overall, the study highlights that efficient design and modular processing can deliver advanced meeting intelligence with low resource requirements. Future work may include speaker identification, multilingual support, sentiment analysis, and real-time dashboards.

REFERENCES

- [1] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser & I. Polosukhin, "Attention Is All You Need," *Advances in Neural Information Processing Systems*, Vol. 30, 2017, pp. 5998–6008
- [2] J. Devlin, M. W. Chang, K. Lee & K. Toutanova, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," *Proceedings of NAACL-HLT*, Vol. 1, 2019, pp. 4171–4186
- [3] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal et al., "Robust Speech Recognition via Large-Scale Weak Supervision," *Proceedings of the International Conference on Machine Learning*, Vol. 202, 2023, pp. 28492–28518
- [4] K. Murray, G. Carenini & R. Ng, "Extractive Summarization of Meeting Recordings," *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, 2005, pp. 593–601
- [5] S. Renals, T. Hain & H. Bourlard, "Recognition and Understanding of Meetings: The AMI and AMIDA Projects," *Proceedings of IEEE Workshop on Automatic Speech Recognition and Understanding*, 2007, pp. 238–247
- [6] D. Jurafsky and J. H. Martin, *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics and Speech Recognition*, 3rd Edition, Pearson Education, 2019
- [7] I. Beltagy, M. Peters & A. Cohan, "Longformer: The Long-Document Transformer," *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, 2020, pp. 119–134
- [8] R. Mihalcea and P. Tarau, "TextRank: Bringing Order into Text," *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, 2004, pp. 404–411
- [9] T. Hain, L. Burget, J. Dines, P. N. Garner, A. El Hannani, M. Karafiat et al., "The AMI Meeting Corpus: A Pre-announcement," *Proceedings of the Machine Learning for Multimodal Interaction Workshop*, 2005, pp. 28–39
- [10] T. Mikolov, K. Chen, G. Corrado & J. Dean, "Efficient Estimation of Word Representations in Vector Space," *Proceedings of the International Conference on Learning Representations (ICLR)*, 2013
- [11] A. Graves, A. Mohamed & G. Hinton, "Speech Recognition with Deep Recurrent Neural Networks," *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2013, pp. 6645–6649
- [12] D. Bahdanau, K. Cho & Y. Bengio, "Neural Machine Translation by Jointly Learning to Align and Translate," *Proceedings of the International*

Conference on Learning Representations (ICLR), 2015

[13] C. Kim, R. Stern & A. Narayanan, “Robust Speech Recognition Using Noise-Resilient Features,” IEEE Transactions on Audio, Speech, and Language Processing, 2016

[14] J. Pennington, R. Socher & C. Manning, “GloVe: Global Vectors for Word Representation,” Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP), 2014, pp. 1532–1543

[15] S. Hochreiter & J. Schmidhuber, “Long Short-Term Memory,” Neural Computation, Vol. 9, No. 8, 1997, pp. 1735–1780

[16] A. Hannun, C. Case, J. Casper, B. Catanzaro, G. Diamos et al., “Deep Speech: Scaling up End-to-End Speech Recognition,” arXiv preprint arXiv:1412.5567, 2014

[17] Y. Liu & M. Lapata, “Text Summarization with Pretrained Encoders,” Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP), 2019

