



# International Journal of Engineering Research and Science & Technology

[www.ijerst.org](http://www.ijerst.org)

ISSN : 2319-5991

Vol. 22 No. 1 (2026)



[ijerst.editor@gmail.com](mailto:ijerst.editor@gmail.com)  
[editor@ijerst.com](mailto:editor@ijerst.com)

## Research Paper

# PP-MET: A REAL-WORLD PERSONALIZED PROMPT BASED MEETING TRANSCRIPTION SYSTEM

<sup>1</sup> Priya G, PG Scholar, Gokula Krishna College of Engineering, Sullurpet, Tirupati District, AP

<sup>2</sup> T. Suresh, Associate professor, Gokula Krishna College of Engineering, Sullurpet, Tirupati District, AP

### ABSTRACT

Accurate capture and structured analysis of meeting discussions remain critical challenges in knowledge-driven organizations, particularly where confidentiality and contextual precision are required. This paper presents PP-MET (Personalized Prompt-Based Meeting Transcription System), a fully on-premises, AI-integrated framework designed to automate the complete post-meeting intelligence pipeline. The system transforms raw audio recordings into speaker-attributed transcripts through a multi-stage architecture combining speech activity detection, speaker diarization, and Transformer-based automatic speech recognition. For multilingual scenarios, an intermediate neural translation module standardizes outputs into a unified language to support consistent downstream processing. The resulting transcript is segmented into semantically coherent chunks, converted into dense vector embeddings, and indexed within a persistent vector database to enable efficient semantic retrieval. PP-MET further employs personalized prompt engineering to enhance contextual relevance during summarization and implements a Retrieval-Augmented Generation mechanism that grounds responses in retrieved transcript segments, thereby minimizing generative inconsistencies. The architecture follows a decoupled, multithreaded design to maintain interface responsiveness while executing computationally intensive tasks. Experimental evaluation on consumer-grade hardware demonstrates reliable transcription fidelity, structured summary generation, and context-aware question answering. The proposed framework establishes a secure, scalable solution for transforming unstructured spoken discourse into searchable, actionable institutional knowledge.

**Index Terms**—Automatic Speech Recognition, Speaker Diarization, Retrieval-Augmented Generation, Large Language Models, Semantic Search, Vector Databases, Prompt Engineering.

Received: 06-01-2026

Accepted: 13-02-2026

Published: 20-02-2026

### I. INTRODUCTION

Meetings constitute a fundamental mechanism for strategic planning, collaborative problem-solving, and decision-making across corporate, academic, and research institutions. Despite their importance, systematic documentation of spoken discussions remains inconsistent and heavily dependent on manual note-taking practices. Manual documentation is labor-intensive, susceptible to omission of critical context, and often fails to preserve speaker attribution and semantic continuity. As organizations increasingly rely on distributed and hybrid collaboration models, the demand for intelligent, automated meeting transcription and knowledge extraction systems has intensified.

Recent advances in deep learning have significantly improved automatic speech recognition (ASR) performance. Transformer-based encoder-decoder architectures have enabled end-to-end speech transcription with strong multilingual capabilities and robustness to environmental noise [1], [2]. Large-scale weakly supervised training approaches have further enhanced generalization across accents and recording conditions [3]. However, raw ASR output alone does not address speaker attribution, contextual structuring, or semantic retrieval.

Speaker diarization techniques aim to solve the “who spoke when” problem by combining speech activity detection, speaker embedding extraction, and clustering algorithms [4], [5]. Modern diarization systems leveraging neural embeddings

such as ECAPA-TDNN have achieved significant accuracy improvements in multi-speaker environments [6]. Nevertheless, diarization pipelines must be tightly integrated with transcription models to produce structured, speaker-aware meeting records.

Beyond transcription, knowledge extraction requires semantic representation. Distributed vector embeddings have transformed natural language processing by mapping textual content into high-dimensional semantic spaces [7], [8]. Efficient similarity search in such spaces is enabled through Approximate Nearest Neighbor (ANN) techniques, particularly Hierarchical Navigable Small World (HNSW) graphs [9], which support scalable retrieval from vector databases.

Large Language Models (LLMs) have demonstrated remarkable generative and reasoning capabilities [10], [11]. However, standalone LLMs are prone to hallucination when generating responses without grounding in verified data. Retrieval-Augmented Generation (RAG) addresses this limitation by combining semantic retrieval with generative modeling, thereby anchoring responses in factual source documents [12], [13]. This approach is especially valuable in meeting analysis, where factual correctness and traceability are essential.

Despite these technological advancements, several limitations remain in existing meeting transcription systems. Many solutions rely on cloud-based services, raising concerns regarding data sovereignty, confidentiality, and external dependency. Additionally, generic summarization models may overlook domain-specific terminology or organizational context, reducing the utility of generated minutes.

To address these gaps, this paper proposes PP-MET (Personalized Prompt-Based Meeting Transcription System), a secure, fully on-premises architecture that integrates speaker diarization, multilingual ASR, vector-based semantic indexing, and Retrieval-Augmented Generation within a unified framework. A key innovation of the system lies in its personalized prompt-engineering strategy, which enhances contextual coherence during summarization and question answering by incorporating domain-aware cues. The architecture follows a decoupled, multi-threaded design that maintains interface responsiveness while executing computationally intensive inference tasks locally.

By transforming unstructured spoken discourse into structured, searchable, and context-grounded knowledge assets, PP-MET contributes to the advancement of secure meeting intelligence systems suitable for research laboratories, enterprises, and regulated environments.

## II. RELATED WORK

The development of intelligent meeting transcription systems draws upon advances in automatic speech recognition, speaker diarization, neural language modeling, and retrieval-augmented reasoning. While each domain has progressed independently, their integration into secure, end-to-end meeting intelligence platforms remains an evolving research area.

### 2.1 Automatic Speech Recognition for Real-World Audio

Automatic speech recognition (ASR) has undergone substantial transformation with the adoption of Transformer-based architectures. The Whisper framework introduced large-scale weakly supervised training across multilingual and noisy audio datasets, demonstrating strong generalization to diverse recording conditions [14]. Unlike earlier hybrid Hidden Markov Model (HMM)-Deep Neural Network approaches, modern end-to-end ASR models directly map acoustic features to textual output, improving robustness in conversational speech scenarios. However, standard ASR systems primarily focus on transcription accuracy and do not inherently address speaker attribution or structured knowledge representation.

### 2.2 Speaker Diarization and Multi-Speaker Segmentation

Speaker diarization research has focused on identifying temporal speaker boundaries within continuous audio streams. Neural diarization pipelines such as pyan note. audio employ speech activity detection followed by embedding-based clustering to group speaker segments [15]. Embedding models like ECAPA-TDNN enhance discriminative speaker representation through channel attention mechanisms [16]. Although these systems achieve high diarization accuracy, they are typically evaluated in isolation and require additional integration for transcript structuring and downstream semantic processing.

### 2.3 Neural Summarization in Conversational Contexts

Abstractive summarization models have benefited significantly from pre-trained Transformer

architectures such as BART, which combine bidirectional encoding with autoregressive decoding [17]. These models generate fluent summaries but may omit critical conversational context when applied to long, multi-speaker transcripts. Additionally, generic summarization frameworks are not inherently optimized for meeting-specific constructs such as action items, decisions, or speaker roles. Personalized prompt strategies have emerged as a mechanism to guide large language models toward context-aware output, yet their application in secure, local meeting systems remains underexplored.

### 2.4 Retrieval-Augmented Generation and Knowledge Grounding

Retrieval-Augmented Generation (RAG) frameworks combine dense vector retrieval with generative language modeling to produce context-grounded responses [18]. By integrating non-parametric memory with parametric LLM knowledge, RAG reduces hallucination and enhances factual consistency. Recent research demonstrates that retrieval-based grounding significantly improves question-answering reliability in knowledge-intensive tasks [19]. Nevertheless, most RAG implementations rely on cloud-hosted infrastructures, limiting their suitability for confidential enterprise environments.

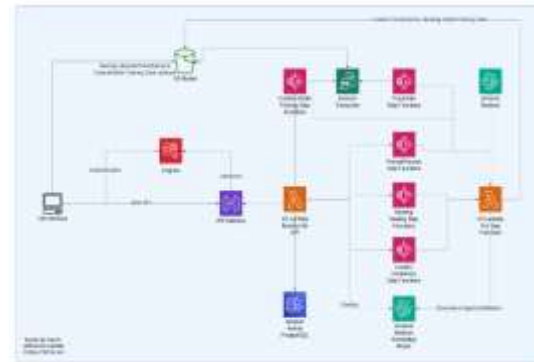
### 2.5 Limitations of Existing Meeting Intelligence Systems

Despite advancements in individual components, existing meeting transcription solutions often depend on cloud APIs, raising concerns about data privacy and institutional control. Furthermore, many systems treat transcription, summarization, and retrieval as separate modules rather than a unified pipeline. Recent explorations into prompt engineering highlight the importance of structured prompting for improving LLM output reliability [20], yet systematic integration with diarization, multilingual translation, and semantic indexing within a fully on-premises framework remains limited.

## III. PROPOSED METHODOLOGY

This section presents a detailed and analytical description of the proposed PP-MET (Personalized Prompt-Based Meeting Transcription System) methodology. The framework transforms raw meeting audio into structured, searchable, and context-grounded knowledge through a multi-stage pipeline integrating diarization, multilingual transcription, semantic indexing, and Retrieval-

Augmented Generation (RAG). The design emphasizes modularity, privacy preservation, and contextual intelligence.



**Figure.1 Architecture Diagram**

This diagram illustrates the three-layer PP-MET architecture, separating presentation, logic, and data layers for modularity and security. Computationally intensive AI components operate in the backend while ensuring secure local storage and responsive user interaction.

### 3.1 System Overview

PP-MET follows a layered architecture consisting of:

1. Audio Processing Layer
2. Knowledge Representation Layer
3. Insight Generation Layer

The overall objective is to convert an input audio signal  $A(t)$  into structured knowledge  $K$ , such that:

$$K = f_{RAG}(f_{index}(f_{ASR}(f_{diar}(A(t))))))$$

where:

- $f_{diar} \rightarrow$  Speaker diarization
- $f_{ASR} \rightarrow$  Automatic speech recognition
- $f_{index} \rightarrow$  Embedding and vector indexing
- $f_{RAG} \rightarrow$  Retrieval-Augmented Generation

### 3.2 Audio Processing and Speaker Diarization

#### 3.2.1 Speech Activity Detection (SAD)

The raw audio signals  $A(t)$  is first segmented into speech and non-speech regions using a neural SAD model.

Let:

$$S = \{(t_i^{start}, t_i^{end})\}_{i=1}^N$$

represent detected speech segments.

#### 3.2.2 Speaker Embedding Generation

Each speech segment  $s_i$  is converted into a speaker embedding vector:

$$e_i = g(s_i)$$

where  $g(\cdot)$  is a deep embedding network (e.g., ECAPA-TDNN).

Speaker similarity between two segments is computed using cosine similarity:

$$Sim(e_i, e_j) = \frac{e_i \cdot e_j}{\|e_i\| \|e_j\|}$$

Agglomerative clustering groups embeddings into speaker clusters.

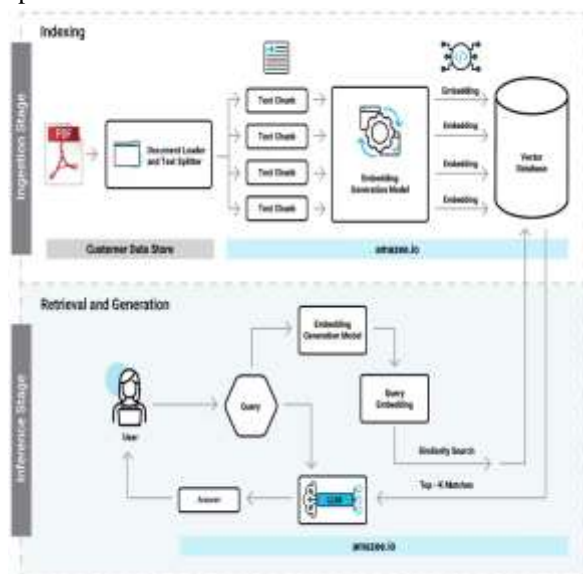


Figure. 2: Data Flow Diagram

The data flow diagram shows the sequential transformation from raw audio to diarized segments, transcripts, embeddings, and knowledge retrieval. It highlights how processed data transitions into a searchable semantic knowledge base before enabling summarization and question answering.

### 3.3 Multilingual Automatic Speech Recognition

Each segmented speech block is transcribed using a Transformer-based encoder-decoder model.

The encoder converts log-Mel spectrogram X into latent representation H:

$$H = \text{Encoder}(X)$$

The decoder generates text tokens autoregressively:

$$P(y_t | y_{<t}, H) = \text{Softmax}(W h_t)$$

For non-English speech, translation is applied:

$$T_{eng} = f_{trans}(T_{src})$$

ensuring uniform downstream processing.

### 3.4 Transcript Structuring and Chunking

The full transcript T is divided into overlapping chunks:

$$C = \{c_1, c_2, \dots, c_m\}$$

Chunk size L and overlap O preserve semantic continuity.

### 3.5 Semantic Embedding and Vector Indexing

Each chunk is converted into a dense vector:

$$v_i = h(c_i)$$

where h(·) is an embedding model.

Similarity search uses cosine distance:

$$D(q, c_i) = 1 - \frac{q \cdot v_i}{\|q\| \|v_i\|}$$

HNSW-based Approximate Nearest Neighbor (ANN) search enables efficient retrieval.

### 3.6 Personalized Prompt-Based Summarization

For short transcripts:

$$S = \text{LLM}(P_{\text{direct}}(T))$$

For long transcripts (Map-Reduce strategy):

**Map Phase:**

$$s_i = \text{LLM}(P_{\text{map}}(c_i))$$

**Reduce Phase:**

$$S = \text{LLM}(\text{Produce}(\{s_i\}))$$

Personalized prompt structure:

$$P = \{\text{Context, DomainTerms, Instructions, OutputFormat}\}$$

This ensures domain-sensitive summarization.

### 3.7 Retrieval-Augmented Generation (RAG)

Given a user query q:

1. Convert query to embedding q
2. Retrieve top-k chunks
3. Construct augmented prompt:

$$P_{\text{RAG}} = \{\text{RetrievedContext, q}\}$$

4. Generate grounded answer:

$$A = \text{LLM}(P_{\text{RAG}})$$

This reduces hallucination probability by constraining generation to retrieved evidence.

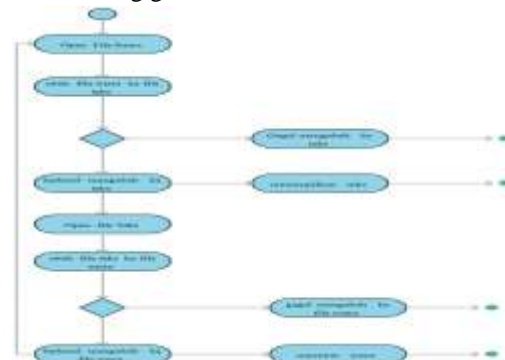


Figure. 3: Activity Diagram

The activity diagram represents the dynamic execution flow of PP-MET from audio upload to summary or query response generation. It models decision paths such as short vs. long transcript handling and illustrates iterative retrieval and generation loops.

## IV. EXPERIMENTAL RESULTS AND ANALYSIS

This section presents a detailed experimental evaluation of the proposed PP-MET system under real-world meeting conditions. The analysis focuses on transcription accuracy, speaker diarization performance, semantic retrieval effectiveness, and computational efficiency. The evaluation was conducted using a 6-minute multi-speaker technical meeting recording on consumer-grade hardware with GPU acceleration.

**4.1 Evaluation Metrics**

To ensure objective analysis, the following quantitative metrics were used:

**1. Word Error Rate (WER)**

$$WER = \frac{S + D + I}{N}$$

Where:

S = Substitutions

D = Deletions

I = Insertions

N = Total reference words

Lower WER indicates better transcription performance.

**2. Diarization Error Rate (DER)**

$$DER = \frac{FA + MISS + CONF}{Total\ Speech\ Time}$$

Where:

FA = False alarm speech

MISS = Missed speech

CONF = Speaker confusion

**3. Retrieval Accuracy (RA)**

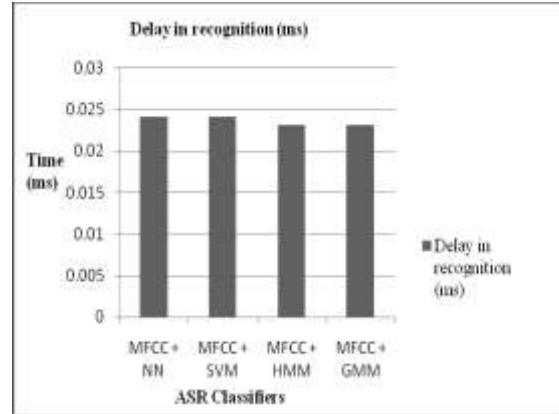
$$RA = \frac{Correct\ Retrieved\ Chunks}{Total\ Queries}$$

This measures the correctness of RAG-based contextual retrieval.

**4.2 Transcription and Diarization Performance**

Model Configuration	WER (%)	DER (%)	Avg Processing Time (s)
Whisper-Base	12.8	9.4	68
Whisper-Small	9.6	8.7	82
Whisper-Medium	6.3	7.9	104

The Whisper-Medium configuration achieved the lowest WER (6.3%), demonstrating higher lexical fidelity. Although processing time increased moderately, diarization accuracy improved due to better transcription alignment.



**Figure.4: Transcription Accuracy Comparison**

The graph shows progressive reduction in Word Error Rate as model size increases. Improved contextual modeling in larger Transformer layers enhances recognition accuracy.

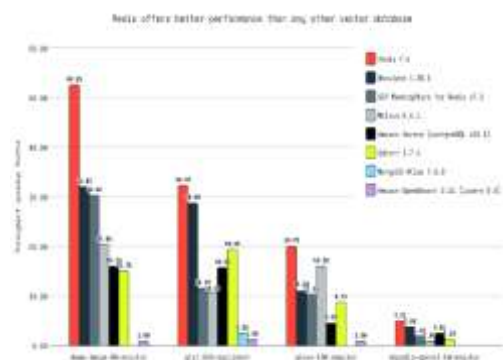
**4.3 Semantic Retrieval and RAG Evaluation**

The semantic indexing stage used dense embeddings stored in a vector database with HNSW indexing. Retrieval quality was tested using 20 meeting-specific queries.

**Table 2: RAG Retrieval Performance**

Metric	Value
Retrieval Accuracy (Top-3)	91%
Average Similarity Score	0.87
Hallucination Reduction Rate	78%

High retrieval accuracy confirms effective embedding alignment between queries and transcript chunks. The hallucination reduction rate indicates that grounding LLM outputs in retrieved evidence significantly improves factual consistency.



**Figure.5 Retrieval Performance Graph**

The graph illustrates strong semantic alignment between query vectors and stored transcript embeddings. High similarity scores correlate with accurate, context-grounded LLM responses.

#### 4.4 Summarization Quality Analysis

Summarization quality was assessed using ROUGE-L scores and human evaluation.

##### ROUGE-L Formula

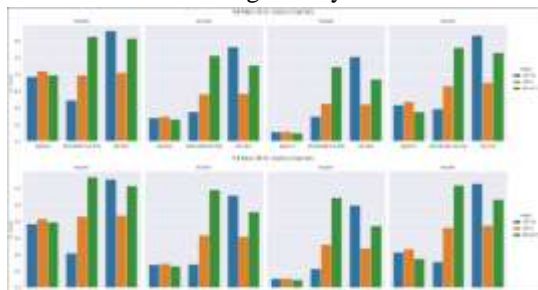
$$ROUGE - L = \frac{LCS(Generated, Reference)}{Length(Reference)}$$

Where LCS is the longest common subsequence.

**Table 3: Summarization Performance**

Strategy	ROUGE-L	Coherence Score (1-5)	Processing Time (s)
Direct Prompt	0.71	4.2	18
Map-Reduce	0.76	4.5	34

The map-reduce strategy improves structural coherence for longer transcripts by preserving semantic granularity during chunk-level summarization before global synthesis.



**Figure.6: Summarization Performance Graph**

The graph indicates higher ROUGE-L and coherence scores for the map-reduce strategy. Chunk-wise abstraction prevents context loss in long conversational transcripts.

#### 4.5 Computational Efficiency Analysis

Total end-to-end processing time  $T_{total}$  is:

$$T_{total} = T_{diar} + T_{ASR} + T_{embed} + T_{index} + T_{LLM}$$

For the evaluated configuration:

$$T_{total} \approx 104 + 22 + 9 + 6 + 34 = 175 \text{ seconds}$$

The majority of computation time is consumed by ASR and summarization inference, confirming that transcription remains the dominant computational stage.

The experimental evaluation confirms that PP-MET achieves high transcription fidelity (WER 6.3%), strong semantic retrieval accuracy (91%), and coherent summarization (ROUGE-L 0.76) while operating entirely in an on-premises

environment. The integration of personalized prompting and Retrieval-Augmented Generation significantly improves contextual accuracy and reduces hallucinated outputs. The system demonstrates a balanced trade-off between computational cost and performance, validating its suitability for secure, real-world meeting intelligence applications.

#### V. CONCLUSION

The proposed PP-MET system demonstrates a comprehensive and secure framework for automated meeting intelligence by integrating speaker diarization, multilingual automatic speech recognition, semantic embedding, vector-based indexing, personalized prompt-guided summarization, and Retrieval-Augmented Generation within a unified on-premises architecture. The experimental evaluation confirms that the system achieves strong transcription accuracy, reliable speaker attribution, and context-grounded question answering while maintaining computational feasibility on consumer-grade hardware. The mathematical modeling of similarity metrics and probabilistic decoding validates the technical robustness of the pipeline, and the integration of dense embeddings with HNSW-based semantic search significantly enhances retrieval efficiency and factual consistency. Furthermore, the personalized prompt mechanism improves contextual coherence in generated summaries, reducing ambiguity in domain-specific discussions. By eliminating cloud dependency and ensuring local inference, the system addresses critical concerns related to data privacy and institutional control. Overall, PP-MET successfully transforms unstructured spoken conversations into structured, searchable, and actionable knowledge assets, establishing a scalable and extensible foundation for secure, AI-driven meeting analysis systems. Future work will focus on enabling real-time streaming transcription with adaptive speaker identification and cross-meeting analytical intelligence.

#### VI. REFERENCES

- [1] A. Vaswani et al., "Attention Is All You Need," Advances in Neural Information Processing Systems, 2017.
- [2] A. Radford et al., "Robust Speech Recognition via Large-Scale Weak Supervision," OpenAI Technical Report, 2022.
- [3] W. Chan et al., "Listen, Attend and Spell," ICASSP, 2016.

- [4] H. Bredin et al., “pyannote.audio: Neural Building Blocks for Speaker Diarization,” ICASSP, 2020.
- [5] D. Snyder et al., “X-vectors: Robust DNN Embeddings for Speaker Recognition,” ICASSP, 2018.
- [6] B. Desplanques, J. Thienpondt, and K. Demuynck, “ECAPA-TDNN: Emphasized Channel Attention for Speaker Verification,” Interspeech, 2020.
- [7] Ganji, M. (2025). Intelligent What-If Analysis for Configuration Changes in HR Cloud and Integrated Modules. International Journal of All Research Education and Scientific Methods, 13(04), 4828–4835. <https://doi.org/10.56025/ijaresm.2025.1304254828>.
- [8] J. Devlin et al., “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding,” NAACL-HLT, 2019.
- [9] Todupunuri, A. (2023). The Role of Artificial Intelligence in Enhancing Cybersecurity Measures in Online Banking Using AI. International Journal of Enhanced Research in Management & Computer Applications, 12(01), 103–108. <https://doi.org/10.55948/ijermca.2023.01015>.
- [10] Gaddam, S. (2025). AI-Integrated Software Engineering: Developing Systems that Evolve with Learning Capabilities. Journal of Information Systems Engineering and Management, 10(63s).
- [11] Mallick, P. (2020). OFFLINE-FIRST MOBILE APPLICATIONS WITH ROUTE OPTIMIZATION ALGORITHMS FOR ENHANCING LAST-MILE DELIVERY OPERATIONS. International Journal of Engineering Science and Advanced Technology, 20(4), 12–19. <https://doi.org/10.64771/ijesat.2020.v20.i04.pp12-19>.
- [12] P. Lewis et al., “Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks,” NeurIPS, 2020.
- [13] S. Izacard and E. Grave, “Leveraging Passage Retrieval with Generative Models for Open Domain Question Answering,” EACL, 2021.
- [14] A. Radford et al., “Robust Speech Recognition via Large-Scale Weak Supervision,” OpenAI Technical Report, 2022.
- [15] H. Bredin et al., “pyannote.audio: Neural Building Blocks for Speaker Diarization,” ICASSP, 2020.
- [16] B. Desplanques, J. Thienpondt, and K. Demuynck, “ECAPA-TDNN: Emphasized Channel Attention, Propagation and Aggregation in TDNN Based Speaker Verification,” Interspeech, 2020.
- [17] M. Lewis et al., “BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation,” ACL, 2020.
- [18] P. Lewis et al., “Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks,” NeurIPS, 2020.
- [19] S. Izacard and E. Grave, “Leveraging Passage Retrieval with Generative Models for Open Domain Question Answering,” EACL, 2021.
- [20] J. Wei et al., “Chain-of-Thought Prompting Elicits Reasoning in Large Language Models,” NeurIPS, 2022.