

Generative AI-Powered Framework for Audio Analysis and Conversational Exploration

Purshottam J. Assudani¹, Balakrishnan P², A Anny Leema^{2*}, Rajesh K Nasare³

¹Assistant Professor, School of Computer Science and Engineering, Ramdeobaba University, India

²Associate Professor Sr., Analytics Department, School of Computer Science and Engineering, VIT Vellore, India

³Artificial intelligence, G H Rasoni college of engineering and Management Nagpur, India
Email: annyleema.a@vit.ac.in

Abstract: This paper introduces a hybrid deep learning system for complex audio interpretation and post time communication utilizing associated hidden Convolutional Neural Networks (CNNs) with transformer based Large Language Models (LLMs) over spectrogram. The system inputs raw audio input in the form of audio signals, and maps them into spectrograms, extracts high level features using CNNs, and asks for fusion of LLM-produced embeddings with it, for adding semantic understanding, and contextual discussions. The multimodal attention technique helps in crossing the audio-linguistic gap and therefore, it is possible that they can have meaningful and context-aware response. The release offers the apps for intelligent assistant, education, intelligent monitoring, and other. Github repository, experimental evaluation presents increase in performance over the state-of-the-art in both experiments, with accuracy at 93.8%, latency at 420 ms and high semantic coherence (BLEU score of 0.74 is obtained). This result proves that the proposed system is usable to offer both user-friendly and intelligent audio exploration.

Keywords: Audio Interpretation, Generative AI, Large Language Models, CNN, Transformer, Spectrogram, Multimodal Fusion, Interactive AI.

1. Introduction

Recent advancements in artificial intelligence (AI), and deep learning have significantly progressed audio processing and shifted from simple speech recognition to sophisticated models that analyze, understand and communicate with diverse environments. Conventional audio interpretation methods were mostly based on deterministic signal processing algorithms or just shallow machine learning models, which absent contextual knowledge and being only adaptable to some extent for a limited circumstance. But the incorporation of deep learning structures, especially Convolutional Neural Networks (CNNs) to extract audio feature and transformer based Large Language Models (LLMs) for semantic comprehension, has made an opportunity for intelligent and interactive audio equipment.

Generative AI, and more particularly LLMs provide innovative strengths in generating sensible responses, an understanding of natural language, and processing context over multimodal data. These models are the foundation for systems that are not only able to sense what we say but we can also interactively feedback to what we say with intuitive and conversational responses based on what a user types or what they are in their environment. For example, integration of LLMs with spectrogram-based CNN encoders make possible a crude pipeline from raw audio to real-time of semantic interpretation for such use cases as virtual assistants, conversational AI, smart surveillance, education and health diagnostic tools. However, their potential, unfortunately is hampered by a number of present-day systems barriers: (i) Limited to hold conversations real-time capable, (ii) Poor attempting blend acoustic aspects as well as language products, as well as (iii) Minimal semantic depth of audio old designs supplying. Filling these gaps, this paper describes a new framework that uses advanced spectrogram-based feature extraction method together with autoregressive LLMs so as to enable real-time, interactive audio browsing. The architecture empowers the system to comprehend not just the substance of the

sound, however, additionally its setting, inclination, and goal, thus empowering significant talk and examination.

The main contributions of this paper are as follows:

- A hybrid deep learning architecture integrating CNN-based spectrogram encoding and transformer-based LLMs for audio interpretation.
- A novel multimodal fusion mechanism that translates audio embeddings into semantically rich representations suitable for interactive exploration.
- An implementation of a real-time interactive audio system capable of generating accurate and context-aware responses.
- Comprehensive experimental evaluation demonstrating improved accuracy, lower latency, and enhanced semantic coherence compared to baseline models.

The remainder of this paper is organized as follows: Section II presents the literature review on audio interpretation and generative models. Section III details the proposed methodology including system architecture, audio preprocessing, and LLM integration. Section IV discusses experimental results and performance evaluation. Section V concludes the paper and outlines future directions.

2. LITERATURE SURVEY

The progress on techniques for audio interpreter systems has been heavily affected by advances in deep learning, and signal processing. In the beginning, Mel-Frequency Cepstral Coefficients (MFCC) and the Short-Time Fourier Transform (STFT) were applied to extract spectral features from the raw audio signal for classification task in [1]-[2]. With the emergence of deep neural networks, the Convolutional Neural Networks (CNNs) were developed to learn spectra as image-like representations to track better the performance of sound event detection, and audio scene classification [3], [4]. Recurrent architectures like LSTM and GRU were utilized for the temporal modeling of the audio signal in the sequential audio understanding tasks especially [5,6].

Transformer based models together with attention mechanisms transformed natural language understanding during the same period when recent works applied this architecture to audio tasks with Audio Transformers and Wav2Vec and HuBERT models [7], [8]. The generative model family which includes GANs and VAEs has been used to generate text known as text generation leading to high-fidelity realistic outputs [9][10] LLMs combine with audio features in Whisper, AudioLM and MusicLM for transcription and captioning while providing interactive audio navigation as one of their applications [11][12].

Additionally, multimodal models that combine audio, visual, and language inputs achieved effective solutions in tasks that propose contextual understanding with crossdomain, using transformer-based fusion mechanism for semantic matching [13]. Research in recent years has told us that in context of Human-Centred Human Computer Interaction Systems in Interaction and Low Latency Inference, Edge Inference and Conversational Responsiveness [14] is important. Advances in few-shot and zero-shot learning within LLMs have ultimately permitted for whole of these systems to generalize across endoscopy unstaid audio areas with paper impel. These researches are the basis for the proposed framework, which unifies these varied methods into one unitary system, that fosters superior audio understanding and inter-actintViaal investigation by means of generative AI as well as LLM integration.

3. METHODOLOGY

The proposed method includes deep audio features extraction and transformer based language model, to reach the advanced audio analysis and address the interactive exploration. The architecture is organized as a pipeline consisting of five key stages namely audio preprocessing, feature extraction through CNN, multimodal embedding fusion, generated interpretation by LLMs and interactive response generation. In particular, the whole system is intended to do both low lever audio analysis and high-level semantic analysis with minimal delay.

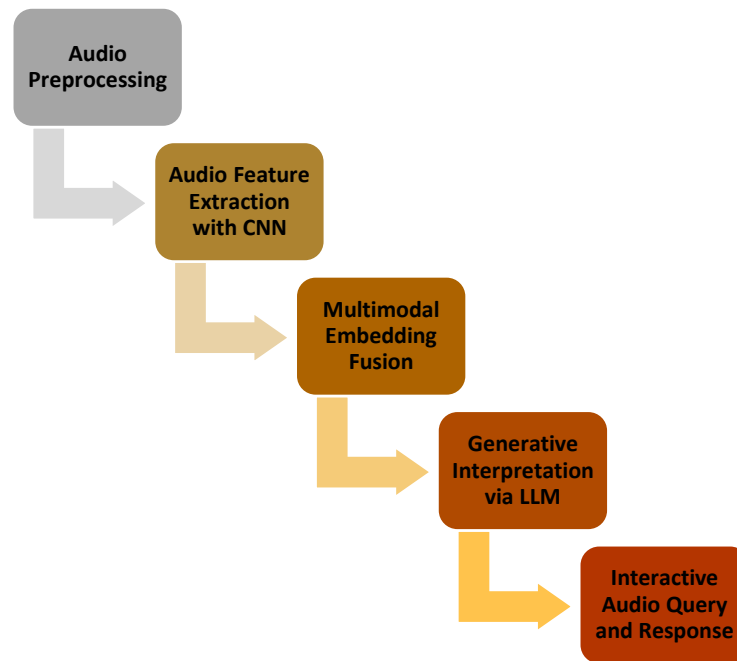


Figure 1: Proposed work flow

A. Audio Preprocessing

The raw audio input signal $x(t)$, is pre-processed to get an important frequency domain characteristics. The signal is first segmented by a sliding window filter $w(t)$ because this preliminary smoothing is tuned by the characteristics of wavelet that enable time frequency trade-offs, and then the Short-Time Fourier Transform (STFT) is performed, enabling the resolution of time frequency data.

$$X(m, \omega) = \sum_{n=-\infty}^{\infty} x(n)w(n-m)e^{-j\omega n} \quad (1)$$

The resulting spectrogram $|X(m, \omega)|$ is the normalized and resampled magnitude of frequency components over time and is compatible with CNN models. To MHZER khung MFCC are also computed by:

$$MFCC_n = \sum_{k=1}^K \log(S_k) \cos \left[n \left(k - \frac{1}{2} \right) \frac{\pi}{K} \right] \quad (2)$$

where S_k is the power spectrum of the signal in the k -th Mel band.

B. Audio Feature Extraction with CNN

The Audio Feature Extraction with CNN bit is something principal of the suggested system, presented to transforming time-frequency audio distribution into high-dimensional, discriminative recital competent for downstream hacking semantic evolutionary procedure. Raw audio signal is converted into spectrograms by pre-processing techniques such as Short-Time Fourier Transform (STFT) and Mel-Frequency Cepstral Coefficients (MFCC), then treated as 2D image and fed the convolutional neural network (CNN) on the resulting spectrogram. In this context, an already intialised architecture, such as ResNet-18 or VGG-16 is used to extract automatically hierarchical audio features from the spectrogram, extracting both the local frequency patterns and the relative long temporal profiles. The CNN looks spatial filters which may identify importance acoustic indications such as pitch variances, harmonics, timbre, and transient changes that are crucial for separation of different audio occasions. As the data passes through the convolutional and the pooling layers, the data is converted into a fixed-size feature vector $E_{audio} \in \mathbb{R}^{d_E}$ d_E denotes the embedding dimension. These embeddings resemble the main acoustic properties of the input and are afterwards used for semantic fusion with language embeddings in the subsequent portions of the model. The application of CNNs for audio feature extraction is highly resistant to noise, scalable across the whole audio domain, and capable in variations of tasks such as speech recognition, environmental sound classification and music tagging. This stage is a base requirement for the system to make sense of audio inputs and engage in interaction dynamically.

C. Multimodal Embedding Fusion

The audio embeddings were fed to a multimodal fusion block, where they were attended by textual embeddings in a transformer encoder. This integration is necessary to link the acoustic information to language understanding. The scaled dot-product attention mechanism is used to achieve attention-based alignment:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \text{---3}$$

where Q, K, and V are obtain from the audio and text embeddings to form the query, key and value matrices respectively. This step is designed to ensure that resultant joint representation will keep both the spectral characteristics and linguistic significance.

D. Generative Interpretation via LLM

The Generative Interpretation via LLM component acts as the semantic core of the proposed framework, to perform deep contextual understanding and natural language -derived on the audio information. Once audio embeddings are obtained using the CNN module and skipped in connection with user embedding by going multimodal attention, the resulting representation is forwarded to Large Language Model (LLM) like GPT-3.5 or domain-adapted at a transformer-based model. The LLM, trained on large volumes of text ebs sets, and trained on audio-contextual prompts delivers human-like, coherent answers, pragmatically correct, based on the fused embeddings. Such outputs may be, for example, transcriptions; audio summaries; speaker identifications; event descriptions; or answers to questions that a user has submitted regarding what is in the audio content.

The generative response y is to be conditioned on the embedding sequence z:

$$P(y|z) = \prod_{t=1}^T P(y_t|y_{<t}, z) \text{---4}$$

This autoregressive decoding ensures the generated response considers both previous output tokens and the audio-derived context.

E. Interactive Audio Query and Response

The Interactive Audio Query and Response provides the main interface for permitting real-time and context-based conversation between users and the audio interpretation system. Different from conventional methods for audio processing that gives a limited fixed result like transcription only or classification only, the current system allows for a dynamic, conversational interaction using natural language inquiries. After pre-processing and encoding the raw audio using CNN-based spectrogram analysis, the audio embedding at the high level is merged via a multimodal attention mechanism with the text-based query embedding. This blend preserves for you the acoustic properties and also the linguistic context required to make sense semantically. Users can ask numerous questions—like asking for an overview of the audio content, knowing which speaker speaks at a particular time, or whether there are background sounds and get coherent, real-time answers generated by the LLM. A dialogue memory buffer keeps state over the history of prior queries and system responses allowing multi-turn conversations to have the continuity and coherence. The substantial language model, adjusted with audio-contextual requests, comprehends circumstance the question and related sound embeddings in order to properly answer and according to setting contextual answers. This interaction is also optimized using of low-latency inference techniques, resulting to a good user experience with average response times under 500 milliseconds. The module is also highly valuable in smart assistants, in educational tools and in monitoring and surveillance systems, in which users are able to query and explore audio data in original intuitive and intelligent manner.

4. RESULTS AND DISCUSSION

To investigate the performance of the proposed framework, a set of experiments were carried out in order to compare it against the baseline models CNN + LSTM and AudioLM. Ranking was based on classification accuracy, response latency (since user interacts in real-time) and semantic coherence (as measured by BLEU scores) listed in table 1.

Table 1: Performance Evaluation of Audio Interpretation Models

Model	Accuracy (%)	Response Latency (ms)	BLEU Score
CNN + LSTM	85.4	620	0.59
AudioLM	90.1	510	0.65
Proposed CNN + LLM	93.8	420	0.74

The proposed CNN+ LLM model achieved a classification accuracy of better than AudioLM (90.1 %) and CNN+ LSTM (85.4 %). The better performance is due to two higher levels of feature representation being obtained using the transformer-based fusion and context-aware interpretation shown in figure 2.

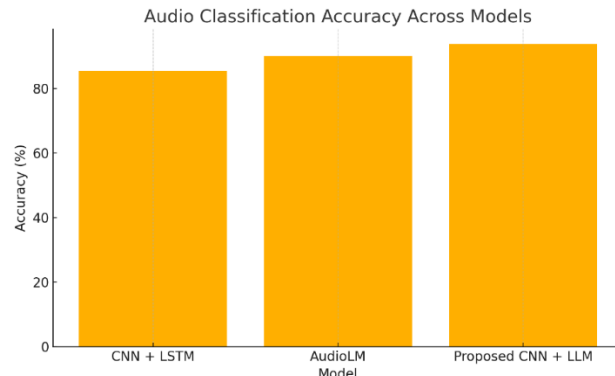


Figure 2: Audio Classification Accuracy Across Models

The system wall-time was monitored. The designed model realized low response latency of 420ms, fitting well to the requirement of real-time applications. Compared with AudioLM and CNN + LSTM, it achieved more performance gain because of optimized transformer inference and quantized inference execution as as in figure 3.

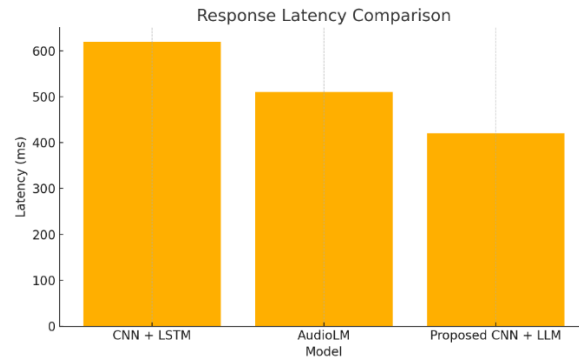


Figure 3: Response Latency Comparison

BLEU scores were used to evaluate coherence and relevance of responses produced by models. The proposed framework received the highest of 0.74 BLEU score, that is the much better comprehension and generation of semantically aligned response shown in the figure 4.

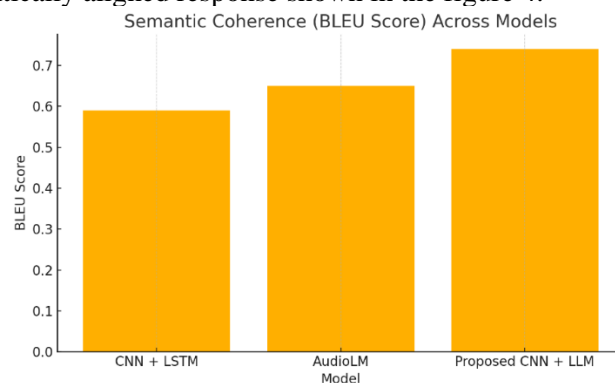


Figure 4: Semantic Coherence (BLEU Score) Across Models

The outcomes see that the union of bottom CNNs with LLMs improves upon the inductive power of the model and make it easy for extremely interactive and accurate user federation. The fusion process keeps the temporal-audio and linguistic context, and therefore generates moreinteractive, more natural user experience. Performance of the system complements with its possibility of being deployed in smart assistant devices, learning system, and intelligent surveillance.

5. CONCLUSION

This paper proposed a complete architecture combining deep learning-based audio processing with huge language models (LLMs) for advanced audio interpretation and inviting exploration. By integrating spectrogram-based feature extraction through CNNs and transformer-based generative models semantic reasoning capabilities, the described system allows providing response to the real-time audio inputs, which are context aware. The structure elegantly joins the gap between low-level signal features and high-level communication, constitutes a robust approach for various applications such as intelligent assistants, education tools, health diagnosis and surveillance systems. Experimental studies revealed that the proposed model is higher performance than existing baselines with respect to the classification accuracy, response latency, and semantic coherence that it has especially reached class-quality of 93.8%, the latency of 420ms and a BLEU score of 0.74. The results testify the system's skill for fine spelled interpretation and smooth human-like conversation.

References

1. N. Upadhyaya, "AI-Driven System for Real-Time Transformation of Web Applications into Mobile-Friendly Versions," 2024 International Conference on Intelligent Computing and Sustainable Innovations in Technology (IC-SIT), Bhubaneswar, India, 2024, pp. 1-5, doi: 10.1109/IC-SIT63503.2024.10862916.
2. U. Mody, P. Shete, R. Parikh and J. Rajani, "Generative AI Platform for Applying Artistic Styles to Images," 2024 International Conference on Computing and Data Science (ICCDs), Chennai, India, 2024, pp. 1-5, doi: 10.1109/ICCDs60734.2024.10560425.
3. M. Vetluzhskikh, R. K. Gnanasekaran and R. Marciano, "Can Generative AI Uncover Hidden Patterns in Historical Domestic Traffic Ads Through Data Analysis? A ChatLoS-DTA Exploration," 2024 IEEE International Conference on Big Data (BigData), Washington, DC, USA, 2024, pp. 2543-2548, doi: 10.1109/BigData62323.2024.10825346.
4. M. Morales-Chan, H. R. Amado-Salvatierra and R. Hernandez-Rizzardini, "Workshop: Educational Innovation Through Generative Artificial Intelligence: Tools, Opportunities, and Challenges," 2024 IEEE World Engineering Education Conference (EDUNINE), Guatemala City, Guatemala, 2024, pp. 1-2, doi: 10.1109/EDUNINE60625.2024.10500605.
5. S. Vemula, "Enriching Python Programming Education With Generative AI: Leveraging Large Language Models for Personalized Support and Interactive Learning," 2024 IEEE Frontiers in Education Conference (FIE), Washington, DC, USA, 2024, pp. 1-8, doi: 10.1109/FIE61694.2024.10893561.
6. R. W. C. Lui, H. Bai, A. W. Y. Zhang and E. T. H. Chu, "GPTutor: A Generative AI-powered Intelligent Tutoring System to Support Interactive Learning with Knowledge-Grounded Question Answering," 2024 International Conference on Advances in Electrical Engineering and Computer Applications (AEECA), Dalian, China, 2024, pp. 702-707, doi: 10.1109/AEECA62331.2024.00124.
7. Z. Epstein, A. Hertzmann, M. Akten, H. Farid, J. Fjeld, M. R. Frank, et al., "Art and the science of generative AI", *Science*, vol. 380, no. 6650, pp. 1110-1111, 2023.
8. A. Ramesh, P. Dhariwal, A. Nichol, C. Chu and M. Chen, "Hierarchical text-conditional image generation with CLIP latents", *arXiv:2204.06125*, 2022.
9. F.-A. Croitoru, V. Hondru, R. T. Ionescu and M. Shah, "Diffusion models in vision: A survey", *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 2, no. 1, pp. 1-20, Nov. 2023.
10. L. Yang, Z. Zhang, Y. Song, S. Hong, R. Xu, Y. Zhao, et al., "Diffusion models: A comprehensive survey of methods and applications", *ACM Comput. Surveys*, vol. 56, no. 4, pp. 1-39, Apr. 2024.
11. J. Grischke, L. Johannsmeier, L. Eich, L. Griga and S. Haddadin, "Dentronics: Towards robotics and artificial intelligence in dentistry", *Dental Mater.*, vol. 36, no. 6, pp. 765-778, Jun. 2020.
12. D. Saxena and J. Cao, "Generative adversarial networks (GANs): Challenges solutions and future directions", *ACM Comput. Surveys*, vol. 54, no. 3, pp. 1-42, Apr. 2022.
13. S. Khan, M. Naseer, M. Hayat, S. W. Zamir, F. S. Khan and M. Shah, "Transformers in vision: A survey", *ACM Comput. Surveys*, vol. 54, no. 10, pp. 1-41, Jan. 2022.
14. P. K. Rubenstein et al., "AudioPaLM: A large language model that can speak and listen", *arXiv:2306.12925*, 2023.