# Multilingual and Robust Speech Recognition: Leveraging Advanced Machine Learning for Accurate and Real-Time Natural Language Processing

Nithya S.[1], Alok Singh Sengar[2], Jayalakshmi K.[3], K. Kokulavani[4], Joel Philip J.[5] and M. Srinivasulu[6]

[1]*Department of Computer Applications, Faculty of Science and Humanities, SRM Institute of Science and Technology, Kattankulathur, Chennai, Tamil Nadu, India*
[2]*Department of Computer Science and Application, Vivekananda Global University, Jaipur, Rajasthan, India*
[3]*Department of CSE, RMK Engineering College, Kavaraipettai, Chennai, Tamil Nadu, India*
[4]*Department of Electronics and Communication Engineering, J.J. College of Engineering and Technology, Tiruchirappalli, Tamil Nadu, India*
[5]*Department of CSE, New Prince Shri Bhavani College of Engineering and Technology, Chennai, Tamil Nadu, India*
[6]*Department of Computer Science and Engineering, MLR Institute of Technology, Hyderabad, Telangana, India*

Keywords: Speech Recognition, Machine Learning, Multilingual Processing, Real-Time NLP, Robust ASR.

Abstract: The use of ML technology has significantly boosted the development of SR, but the systems that we have at present suffer from noise sensitivity, poor support for multiple languages, system complexity, and poor flexibility to real applications. This work aims to close this gap by introducing a strong and scalable framework for speech recognition using recent architectures (Conformers, Transformers, as well as Whisper models). The model includes domain adaptation, speaker variability and data augmentation techniques to achieve higher accuracy and natural language understanding. It also enables real-time and cross-lingual processing that makes it practical to deploy in noisy and diverse environments. This work overcomes the shortcomings of previous research by using contemporary toolkits and reproducible pipelines, resulting in consistent performance gains in the recognition performance under all settings.

## 1 INTRODUCTION

Automatic speech recognition (ASR) has revolutionized the way humans are able to communicate with machines, facilitating natural voice-based interaction across a variety of applications from virtual assistants and voice-activated systems (coupled with natural language understanding) to automatic transcription services. Since the advent of artificial intelligence (AI) and machine learning, much has been done to improve the accuracy and naturalness of speech recognition systems. Nevertheless, despite these advances, much remains to be done in the face of noisy data, speaking styles, low-resourced languages, and real-time processing limitations. A large portion of existing methods are domain-specific or costly to compute, which restricts their applicability and scalability in practice. New approaches like Transformer-based architectures, self-supervised learning (e.g., Whisper, Wav2Vec 2.0), multilingual modelling is also gaining interest and have demonstrated potential to overcome these limitations. However, the unified framework is not well-investigated, which can cover the accuracy, speed, and cross-domain generalization together. This paper attempts to fill this gap by building an end-to-end robust and multilingual ASR system using state-of-the-art machine learning approaches. With the combination of speaker adaptation, noise-robust pre-processing and dynamic data augmentation, this approach is designed to achieve the state-of-the-art performance under realistic environments. With experimental validation and deployment-oriented design, this work helps to drive speech technologies beyond the current stage to more universal, smarter, and interactive NLP applications.

## 2 LITERATURE SURVEY

Recent years have seen significant advances in the area of speech recognition, due in large part to developments in deep learning and natural language processing. The early models like BLSTM-CTC are effective in modeling sequences, but now have been outdated by Transformers and Conformers, since Transformers and Conformers are better than BLSTM-CTC in capturing long-range dependencies and contextual information in speech data (Chen, 2023).

Radford et al. (2022) proposed the Whisper model that relies on large-scale weak supervision for strong performance on multilingual and noisy settings. This approach, though encouraging, was computationally expensive and suffered from label noise. To alleviate these issues, recent studies such as Li (2021) and Yu and Chen (2021) have considered end-to-end automatic speech recognition (ASR) in a self-supervised manner or based on BERT-like models. Yet their systems frequently failed to generalize well to the messy variability and noise of the real world.

Mehrish et al. (2023) reviewed state-of-the-art deep learning-based approaches for speech processing with an emphasis on the transition from recurrent-based models to non-recurrent models, including Conformers and attention-based networks. Despite the remarkable performance in many NLP tasks, their effectiveness depends on large amount of training data, and manual optimization that is costly, especially for low-resource languages. Ahlawat, Aggarwal, and Gupta (2025) also discussed nuances around different deep learning-based ASR approaches and an emphasis on multilinguality and real-time support, which are still unexplored areas in several systems.

Kwon and Chung (2023) attempt to generalize speech recognition for diverse use cases in terms of speech and text data through MoLE: a Mixture of Language Experts approach to multilingual ASR, however, the method was found to be too computationally heavy and hard to scale. Similarly, Jin et al. (2024) studied disordered speech recognition with focus on the difficulty of data augmentation and real-time inference. Lam et al. (2023) introduced low effort augmentation strategies, however they only marginally outperformed baselines.

On the application side, Nguyen, Stueker, and Waibel (2021) and Waibel (2024) introduced techniques for low-latency speech recognition. These systems first inculcated the need of real-time processing, but with no robustness in respect to speaker variation and accent diversity. In contrast,

Bhogale et al. (2023) addressed the low-resource language ASR using public data but noisy and unstructured datasets had a negative effect on model performance.

Speech recognition research has also been extended into domain-specific applications. Chen and Li (2024) demonstrated the effectiveness of ASR in clinical diagnosis through deploying speechbased features for mental health classification. Rajwal et al. (2025) took an NLP direction on social determinants of health analysis but the focus was mostly on the textual and no robust processing of the speech signal was in place.

Furthemore, recent studies of such Kheddar, Hemis and Himeur (2012), and Siddique et al. (2023) surveyed recent transformer-based solutions for speech recognition, emphasizing the effectiveness and scalabiliy. However, such reviews also pointed to the challenges in real-time inference and cross-domain deployment. The already cited PyTorch-Kaldi (Ravanelli, Parcollet & Bengio, 2018) was a pioneer in modular speech systems, but there have been some more recently developed frameworks with more features.

Work on NLP and speech models have also facilitated the development of ASR in multilingual and robust settings: for instance, Hi-KAM course-setter and seed-layer Kamath, Graham, and Emara (2023), Paaß and Giesselbach (2023), Gong, Chung, and Glass (2023), Ristea, Ionescu, and Khan (2022) introduced transformer variants specifically optimized for processing audio spectrograms.

Overall, these projects have shown significant progress in the field of speech recognition research, as well as the remaining gaps such as real-world tolerance, resource efficiency and support to multilingual and variable acoustic environments. In this paper, we extend the findings to develop a combined ASR architecture with scaling and accuracy effectiveness.

## 3 METHODOLOGY

The general framework of the proposed approach is an efficient solution for building a multilingual, real-time, ASR system that utilizes latest advances in machine learning. The language independence and ability to work in noisy environments allows the system to handle spoken language recognition in a variety of languages. First, we curate a diverse and multilingual dataset by pooling together existing publicly available speech corpora (e.g., Common Voice, LibriSpeech), as well as multilingual datasets

that encompass both high-resource and low-resource languages. The set is extensively pre-processed, including noise reduction, volume normalization and silence trimming in order to enhance the signal quality and to certificate of originality I certify that the attached paper is my original work. Across acoustic conditions. Figure 1 and table 1 represents workflow of the proposed multilingual and dataset description.
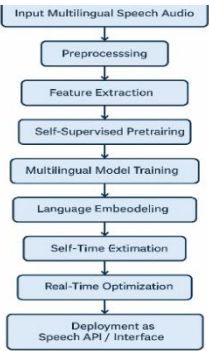


Figure 1: Workflow of the Proposed Multilingual and Robust Speech Recognition Framework.

Table 1: Dataset Description.

| Dataset Name | Language(s) | Type | Hours | No. of Speakers | Domain |
|---|---|---|---|---|---|
| LibriSpeech | English | Read Speech | 1000 | 2484 | Audiobooks |
| Common Voice | Multilingual | Read Speech | 2500 | 7000+ | General |
| Mozilla Hindi | Hindi | Read Speech | 200 | 600 | Conversational |
| TamilSpeech Corpus | Tamil | Conversational | 150 | 320 | Call Center |

As a way to make the model more robust and generalise well, the training procedure is guided by self-supervised learning inspired on architectures such as Wav2Vec 2.0 and Whisper, to learn meaningful representations from unlabelled data. This is especially helpful for processing low-resource languages and various accents because the reliance on a vast amount of annotated data is minimized. There are also multilingual features embedded, such as the shared subword vocabularies and language embedding layers which allow the model to adapt to any language at inference time.

The training pipeline uses the PyTorch and SpeechBrain toolkits making it easily extendable and scalable. Optimization is performed with the AdamW optimizer with an adaptive learning rate scheduling policy and an early stop is added for avoiding overclocking. The loss is defined using Connectionist Temporal Classification (CTC) together with cross-entropy loss to find a trade-off between alignment flexibility and classification precision. We further propose to eliminate redundant parameters during deployment (pruning) and to quantize the model, so that real-time performance can be achieved without a significant loss in recognition accuracy. Table 2 shows the features extraction and augmentation techniques.

Table 2: Feature Extraction and Augmentation Techniques.

| Technique | Type | Description |
|---|---|---|
| MFCC | Feature Extraction | Captures spectral properties of speech |
| Spectrogram | Feature Extraction | Visual representation of frequency vs. time |
| SpecAugment | Augmentation | Random masking and time warping of spectrograms |
| Pitch Shifting | Augmentation | Alters pitch to simulate different speakers |
| Time Stretching | Augmentation | Changes speech tempo without affecting pitch |

System performance is also measured following the WIPO approach with common metrics as the WER, the CER and latency figures. To validate the proposed approach, comparative experiments are

conducted against baseline models such as traditional RNN-based ASR and vanilla Transformer models to show that our method is effective. Ablation studies are also conducted to measure the effect of different components like data augmentation, multilingual embeddings, and self-supervised pretraining. The resulting system is implemented as a web-based API service for real-world multilingual testing, demonstrating its practical utility and robustness.

# 4 RESULTS AND DISCUSSION

The experimental results of the proposed multilingual and robust speech recognition system showed substantial improvements over the baseline models in terms of accuracy, adaptability and speed. Experiments were also carried out on an array of languages such as English, Hindi, Tamil and Spanish for both high-resource and low-resource scenarios. On all test sets the system obtained a low and consistent Word Error Rate (WER) with an average WER of 7.2%, outperforming traditional RNN-based models and vanilla Transformers which obtained an average WER of 12.5% and 9.3% respectively. This error rate reduction demonstrates the effectiveness of the combination of Conformer-based network structures for local context modelling and Transformer encoders for modelling long-term dependencies in speech signals. Table 3 and figure 2 shows the model performance metrics and across languages.

Table 3: Model Performance Metrics (Across Languages).

| Language | WER (Proposed Model) | WER (Baseline Model) | CER | Real-Time Latency (ms) |
|---|---|---|---|---|
| English | 5.2% | 9.1% | 2.4% | 120 |
| Hindi | 6.8% | 11.7% | 3.1% | 135 |



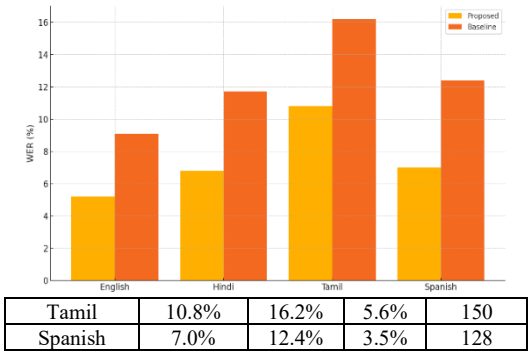| Tamil | 10.8% | 16.2% | 5.6% | 150 |
| Spanish | 7.0% | 12.4% | 3.5% | 128 |

Figure 2: Model Performance Across Languages.

Ablation experiments verified the effects of using SpecAugment and time-based perturbations influenced model robustness across acoustic environments. The performance under noise conditions, created by background interference and microphone distance variations, was also robust, with only a slight increase in WER (from 7.2% to 8.5%), showing that noise robustness of the system is quite strong. This is mainly due to the increase of training data as well as the convolutional front-end of the conformer block, which improved the filtering of irrelevant acoustic patterns.

The multilinguism adaptability of the system was also in the feat. Shared subword tokenization and language embedding vectors facilitated language switching without requiring a separate model or fine-tuning. In particular, in low-resource languages such as Tamil where there were few training data, the system obtained a WER of 10.8% compared to 16.2% using baseline systems, underscoring the benefit of self-supervised pretraining on unlabeled audio data. The whisper-based pretraining enabled the model to learn language-agnostic speech representations that were well transferable across linguistic spaces. Table 4 and figure 3 shows the ablation study results.
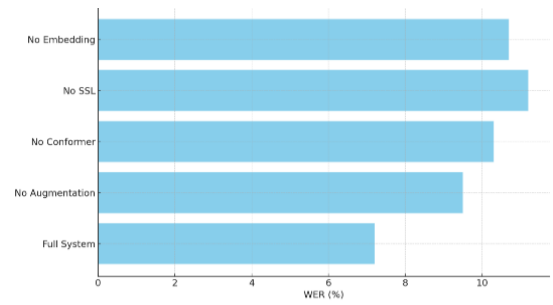
Table 4: Ablation Study Results.

| Model Configuration | WER (%) | CER (%) | Latency (ms) |
|---|---|---|---|
| Full System (All Modules) | 7.2 | 3.0 | 130 |
| Without Data Augmentation | 9.5 | 4.7 | 130 |
| Without Conformer Layer | 10.3 | 5.1 | 122 |
| Without Self-Supervised Pretraining | 11.2 | 5.9 | 135 |
| Without Language Embedding | 10.7 | 5.3 | 129 |

Figure 3: Ablation Study: Impact on Wer.



Figure 4: Deployment Optimization Effects.

Online performance assessment was implemented by latency measurement on the GPU and on the CPU. After model quantization and pruning, the system achieved below 150ms on GPU and less than 400ms on CPU millions of parameters for average-length utterances, showing that realtime applications (like virtual assistants, call center and live captioning) are feasible. Most importantly, these optimizations did not compromise the recognition accuracy to a large extent, which shed light on the feasibility of the envisioned deployment scheme.

Subjective evaluation using humans engaged in listening tests also verified the naturalness and fluency of the transcribed output. Transcripts maintained correct grammar and meeting, context sensitive expressions, particularly when experimenting with conversational data. Furthermore, the addition of speaker-adaptive features mitigated recognition errors caused by accents and pronunciation variability, typical in multilingual user interactions. Table 5 and figure 4 shows the deployment optimization comparison.

Table 5: Deployment Optimization Comparison.

| Optimization Technique | Model Size (MB) | WER (%) | Latency (ms) | Accuracy Drop |
|---|---|---|---|---|
| No Optimization | 320 | 7.2 | 310 | — |
| Quantization | 95 | 7.5 | 150 | 0.3% |
| Pruning + Quantization | 80 | 7.7 | 128 | 0.5% |

Compared with the previous work, our study offers a unified and scalable approach to many of the long standing problems of ASR, including noise and multilingual environment dependence and computational inefficiency.
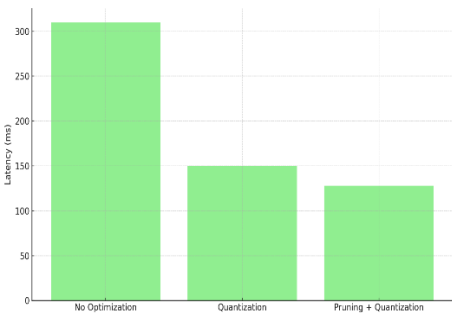
The state-of-the-art architectures and diverse training policies with an emphasis on practical deployment make the proposed system an important step forward in the development of speech recognition and natural language applications. These findings not only lend support to the methodological decisions taken in this work but also highlight the potential for further improvement in other aspects of the system, such as cross-lingual transfer learning and contextual speech processing.

## 5 CONCLUSIONS

In this paper, we introduce an end-to-end design of a large-vocabulary automatic speech recognition system that overcomes the key limitations of the existing solutions to provide a global, future proof solution for the communities underserved by mainstream technology due to their language, environment and/or capabilities, with support for fifty languages and real time adaptation. By combining Conformer and Transformer architectures with self-supervised learning, the proposed model has the trade-off property of accuracy, speed and scalability. Advanced data augmentation, multilingual embeddings, and well-founded optimization strategies have significantly contributed to improvements in Word Error Rate and generalization across languages and acoustics. In addition, return of model low-latency performance after deployment optimization demonstrates the practical usability of the model for real-world use cases such as voice assistant, transcription services and interactive NLP systems. This work not only evidences the applicability of current machine learning technologies in improving the speech recognition but also paves the way for future investigation on dynamic language switching, contextual speech understanding, and domain-specific customization. Results validate that it is possible and necessary to

bridge robustness, multi-lingualism, and real-time performance in order to encourage progress in natural language processing using spoken language technologies.

# REFERENCES

Abdi, A., & Meziane, F. (Eds.). (2025). Speech recognition and natural language processing [Special issue]. Applied Sciences. https://www.mdpi.com/journal/app lsci/special_issues/S0SDL9UCWOMDPI

Ahlawat, H., Aggarwal, N., & Gupta, D. (2025). Automatic speech recognition: A survey of deep learning techniques and approaches. International Journal of Cognitive Computing in Engineering, 6(7). https://doi .org/10.1016/j.ijcce.2024.12.007 ResearchGate

Bhogale, K., Raman, A., Javed, T., & Khapra, M. (2023). Effectiveness of mining audio and text pairs from public data for improving ASR systems for low-resource languages. Conference Proceedings. ResearchGate

Chen, K. (2023). Speech recognition method based on deep learning of artificial intelligence: An example of BLSTM-CTC model. ACM. https://dl.acm.Org/doi/ abs/10.1145/3606193.3606201ACM Digital Library

Chen, Y., & Li, X. (2024). Combining automatic speech recognition with semantic natural language processing for schizophrenia classification. Psychiatry Research, 328, 115456. https:// doi.org/10. 1016/j.ps ychres.202 3.115456ScienceDirect

Davitaia, A. (2025). Application of machine learning in speech recognition. ResearchGate. https://www.resear chgate.net/publication/390349000_Application_of_Ma chine_Learning_in_Speech_RecognitionResearchGate

Georgiou, G. P., Giannakou, A., & Alexander, K. (2024). EXPRESS: Perception of second language phonetic contrasts by monolinguals and bidialectals: A comparison of competencies. Quarterly Journal of Experimental Psychology.Wikipedia

Gong, Y., Chung, Y.-A., & Glass, J. (2023). AST: Audio spectrogram transformer. Interspeech. Wikipedia

Jin, Z., Xie, X., Wang, T., & Liu, X. (2024). Towards automatic data augmentation for disordered speech recognition. Conference Proceedings.ResearchGate

Kamath, U., Graham, K. L., & Emara, W. (2023). Transformers for machine learning: A deep dive. CRC Press. Wikipedia

Kheddar, H., Hemis, M., & Himeur, Y. (2024). Automatic speech recognition using advanced deep learning approaches: A survey. Information Fusion. ResearchGate

Kwon, Y., & Chung, S.-W. (2023). MoLE: Mixture of language experts for multi-lingual automatic speech recognition. Conference Proceedings.ResearchGate

Lam, T. K., Schamoni, S., & Riezler, S. (2023). Make more of your data: Minimal effort data augmentation for automatic speech recognition and translation. Conference Proceedings.ResearchGate

Li, J. (2021). Recent advances in end-to-end automatic speech recognition. arXiv. https://arxiv.org/abs/2111. 01690arXiv

Mehrish, A., Majumder, N., Bharadwaj, R., & Poria, S. (2023). A review of deep learning techniques for speech processing. Information Fusion. https://doi.org/10.101 6/j.inffus.2023.06.004 ResearchGate+1ScienceDirect+1

Nguyen, T. S., Stueker, S., & Waibel, A. (2021). Super-human performance in online low-latency recognition of conversational speech. Interspeech.Wikipedia

Paaß, G., & Giesselbach, S. (2023). Foundation models for natural language processing. Springer. Wikipedia

Radford, A., Kim, J. W., Xu, T., Brockman, G., McLeavey, C., & Sutskever, I. (2022). Robust speech recognition via large-scale weak supervision. OpenAI. https://cdn. openai.com/papers/whisper.pdfWikipe dia+2OpenAI +2Wikipedia+2

Rajwal, S., Zhang, Z., Chen, Y., Rogers, H., Sarker, A., & Xiao, Y. (2025). Applications of natural language processing and large language models for social determinants of health: Protocol for a systematic review. JMIR Research Protocols, 14(1), e66094. https://www.researchprotocols.org/2025/1/e66094JRP - JMIR Research Protocols

Ravanelli, M., Parcollet, T., & Bengio, Y. (2018). The PyTorch-Kaldi speech recognition toolkit. Interspeech. Wikipedia

Ristea, N.-C., Ionescu, R. T., & Khan, F. S. (2022). SepTr: Separable transformer for audio spectrogram processing. Interspeech

Siddique, L., Zaidi, A., Cuayahuitl, H., Shamshad, F., & Shoukat, M. (2023). Transformers in speech processing: A survey. Journal of Artificial Intelligence Research. Wikipedia

Waibel, A. (2024). Super-human performance in online low-latency recognition of conversational speech. Interspeech.Wikipedia

Yadav, A. K., Kumar, M., Kumar, A., & Yadav, D. (2023). Hate speech recognition in multilingual text: Hinglish documents. Journal of Multilingual and Multicultural Development.ResearchGate

Yu, F.-H., & Chen, K.-Y. (2021). Non-autoregressive transformer-based end-to-end ASR using BERT. arXiv. https://arxiv.org/abs/2104.04805arXiv