

A Multilingual Adaptive Digital Speech Rehabilitation System with Real-Time Feedback

DR. Priyanka B.G¹*Assistant Professor, Dept. of CSE***Virupakshi²***Dept. of CSE***Shafiqha Khanum³***Dept. of CSE***Nayana N⁴***Dept. of CSE***Sukanya⁵***Dept. of CSE*

PES Institute of Technology and Management, Shivamogga, Karnataka, India

Emails: priyankabg@pestrust.edu.in, virukannada2018@gmail.com, shafiqhakhanum53@gmail.com, nayanaagowda264@gmail.com, sukanyahiremath802005@gmail.com

Abstract— Aphasia is a neurogenic communication disorder that disrupts an individual's ability to speak, understand, read, or write, commonly occurring after a stroke or traumatic brain injury. Traditional speech therapy, although effective, is often time consuming, expensive, and limited by linguistic and geographical accessibility. This paper presents a scalable and multilingual AI-driven speech therapy system integrating automatic speech recognition (ASR), text-to-speech (TTS), adaptive feedback, and an avatar-based phoneme–viseme synchronization module for English, Hindi, and Kannada. The system utilizes multilingual Wav2Vec2 embeddings, phoneme-level speech analysis, and animated lip-sync feedback generated using Rhubarb Lip Sync and Media Pipe to deliver multimodal and inclusive rehabilitation. Experimental results demonstrate ASR accuracy between 89% and 93% across the supported languages, along with high user engagement reflected by an average patient satisfaction rating of 4.6 out of 5 and positive therapist feedback. The modular architecture enables deployment in resource-constrained environments while maintaining clinical relevance through integration with standardized assessment metrics such as the Western Aphasia Battery-Aphasia Quotient. Overall, the proposed framework provides an accessible, cost-effective, and technologically enhanced approach to speech therapy for individuals with aphasia.

Keywords—Aphasia, speech therapy, artificial intelligence, multilingual ASR, adaptive feedback, phoneme–viseme mapping, avatar lip sync, Kannada, personalized assessment

I. INTRODUCTION

Aphasia is a language disorder that impairs speech production, comprehension, reading, and writing, most commonly due to left-hemisphere brain damage from stroke, trauma, tumors, or neurodegenerative conditions [1], [2]. It often cooccurs with dysarthria or apraxia of speech, with clinical presentation varying based on fluency, comprehension, and repetition abilities. Limited access to continuous clinical therapy makes long-term rehabilitation challenging.

Existing digital speech-therapy solutions provide varied support levels [2], [5]. Tactus Therapy offers clinician-designed evidence-based exercises across mobile platforms without time limits or internet requirements, while Constant Therapy provides personalized adaptive tasks but relies on subscription access, restricting affordability. Telerehabilitation platforms enable remote delivery of structured language tasks, enhancing continuity of care [2], [6]. AphaDIGITAL advanced this space by integrating automated AI feedback with therapist-curated tasks, improving patient motivation. Other systems, such as the Malay-language app for post-stroke rehabilitation, leverage mobile frameworks, IoT-based progress tracking, and therapist dashboards [5].

Low-resource language speech-recognition efforts include a Kannada ASR model using MFCC features and CNNs, achieving

80–90% accuracy [4]. Multilingual therapy applications exist for English, Urdu, and Sindhi, integrating pronunciation assessment via Levenshtein distance and neural machine translation. Recent multilingual ASR research employs self-supervised learning and contrastive predictive coding to transfer acoustic representations across languages, addressing limited labeled clinical data [3], [4]. Despite advances, visual articulation feedback remains underdeveloped. Psycholinguistic evidence shows multimodal (audio-visual) cues significantly improve speech learning compared to audio-only systems [7]. Avatar-based articulation models offer consistent, fatigue-free demonstration of phoneme–viseme alignment, presenting a promising direction for scalable aphasia rehabilitation tools [7], [8]. To bridge these gaps, this study proposes a Multilingual AI Powered Speech Therapy System for Aphasia that integrates state-of-the-art ASR, pronunciation feedback, and multimodal guidance. The key objectives are: 1) To develop a multilingual ASR model fine-tuned for English, Hindi, and Kannada using Wav2Vec2-XLS-R. 2) To implement real-time pronunciation feedback using Levenshtein-based similarity scoring. 3) To provide multimodal reinforcement through AI generated lipsynced visual cues via Wav2Lip. 4) To ensure inclusivity through language translation (mBART / Google Translate) and adaptive progression tracking.

II. METHODS

The proposed framework integrates ASR, TTS, adaptive feedback, data tracking, and a real-time avatar-based visualization module (Fig. 1).

The system comprises an integrated multi-component architecture designed to deliver personalized speech rehabilitation through iterative assessment, adaptive task delivery, and multimodal feedback [3], [5]. The architecture separates concerns into distinct layers: acoustic capture and processing, linguistic analysis, feedback generation, progress tracking, and clinician oversight, accessible through computers or mobile browsers.

A. Model Architecture

The system employs a dual-ASR framework combining Wav2Vec2-XLS-R-300M and OpenAI Whisper to enable robust multilingual transcription. Wav2Vec2-XLS-R extracts contextual speech representations using convolutional and transformer encoders, while Whisper enhances performance in noisy environments.

B. Connectionist Temporal Classification Loss

The ASR fine-tuning process uses Connectionist Temporal Classification (CTC) loss to align variable-length speech inputs with their corresponding text outputs. The CTC loss is defined as:

$$L_{CTC} = -\ln \sum_{\pi \in B^{-1}(y)} P(\pi|x) \quad (1)$$

where π represents all valid alignments that map to target sequence y under the collapse function B , and $P(\pi|x)$ denotes the probability of alignment π given input sequence x .

C. Dataset and Preprocessing

A multilingual speech dataset covering English, Hindi, and Kannada was constructed using Google Text-to-Speech (gTTS). Each audio clip was paired with its transcript and language label. The Wav2Vec2Processor performed normalization, tokenization, and alignment of speech-text pairs. Data augmentation techniques were applied to simulate varied acoustic environments and improve generalization.

D. Training Configuration

The model was fine-tuned using the Hugging Face Trainer API. The pretrained facebook/wav2vec2-xls-r-300m checkpoint was loaded, and the feature encoder was frozen while training only the classification head using CTC loss. Training was performed using the AdamW optimizer with a linear learning rate scheduler, and validation accuracy was tracked to save the best model checkpoints. The key hyperparameters included a batch size of 8, a learning rate of 3×10^{-4} , and 10–15 epochs of training. Pronunciation error detection was implemented using the Levenshtein distance. The distance $D(a,b)$ between two sequences was defined as:

$$D(a,b) = \min \begin{cases} D(i-1,j) + 1, \\ D(i,j-1) + 1, \\ D(i-1,j-1) + \text{cost} \end{cases}$$

where $\text{cost} = 0$ if $a_i = b_j$ and 1 otherwise. Accuracy was computed as:

$$\text{Accuracy} = \left(1 - \frac{D(a,b)}{\max(|a|, |b|)} \right) \times 100\%$$

This metric quantifies similarity between predicted and target phoneme sequences and enables adaptive feedback in pronunciation training.

E. System Workflow

- 1) User selects language (English, Hindi, Kannada) and activity type (speech, text, image)
- 2) Speech input processed by fine-tuned Wav2Vec2-XLS-R or Whisper
- 3) Error comparison using Levenshtein metric
- 4) Feedback generated with text/audio and Wav2Lip-based lip-sync
- 5) Real-time translation module supports multilingual interaction
- 6) User performance logged to adjust exercise difficulty dynamically

Phonemes correspond imperfectly to visemes, the visual manifestations of articulation [7]. The system maintains language-specific mapping tables reflecting articulation characteristics of each supported language, defining which animated mouth positions correspond to specific phonetic segments. For each therapeutic target utterance, the system generates a synchronized video animation demonstrating correct articulation using Rhubarb Lip Sync [7]. This animation process involves phonetic decomposition of the target utterance, retrieval of corresponding viseme animations, temporal alignment ensuring smooth transitions between consecutive visemes, and synchronization with speech audio playback. Facial mesh analysis through MediaPipe computer vision techniques enables tracking of user articulation in real-time, providing objective measures of user-to-model correspondence [8].

The instruments used are based on Western Aphasia Battery (WAB) but simplified for mobile application utilization [5]. Backend services implement core functionality using Python frameworks optimized for scientific computing. FastAPI provides web service infrastructure supporting Representational State Transfer API endpoints with automatic documentation generation. Database systems employ MySQL providing adequate performance for typical deployment scales. Frontend components employ React with Babel transformation and Webpack bundling for web-based clinician dashboards.

III. RESULTS AND DISCUSSION

Metric	Result
ASR Accuracy (avg.)	91.8–92%
Kannada Accuracy Gain	12–15% improvement
Pronunciation Error Reduction	35% improvement
Latency	280 ms / word
Therapy Adaptivity	96% correct difficulty scaling
Multimodal Feedback Impact	Higher engagement & motivation

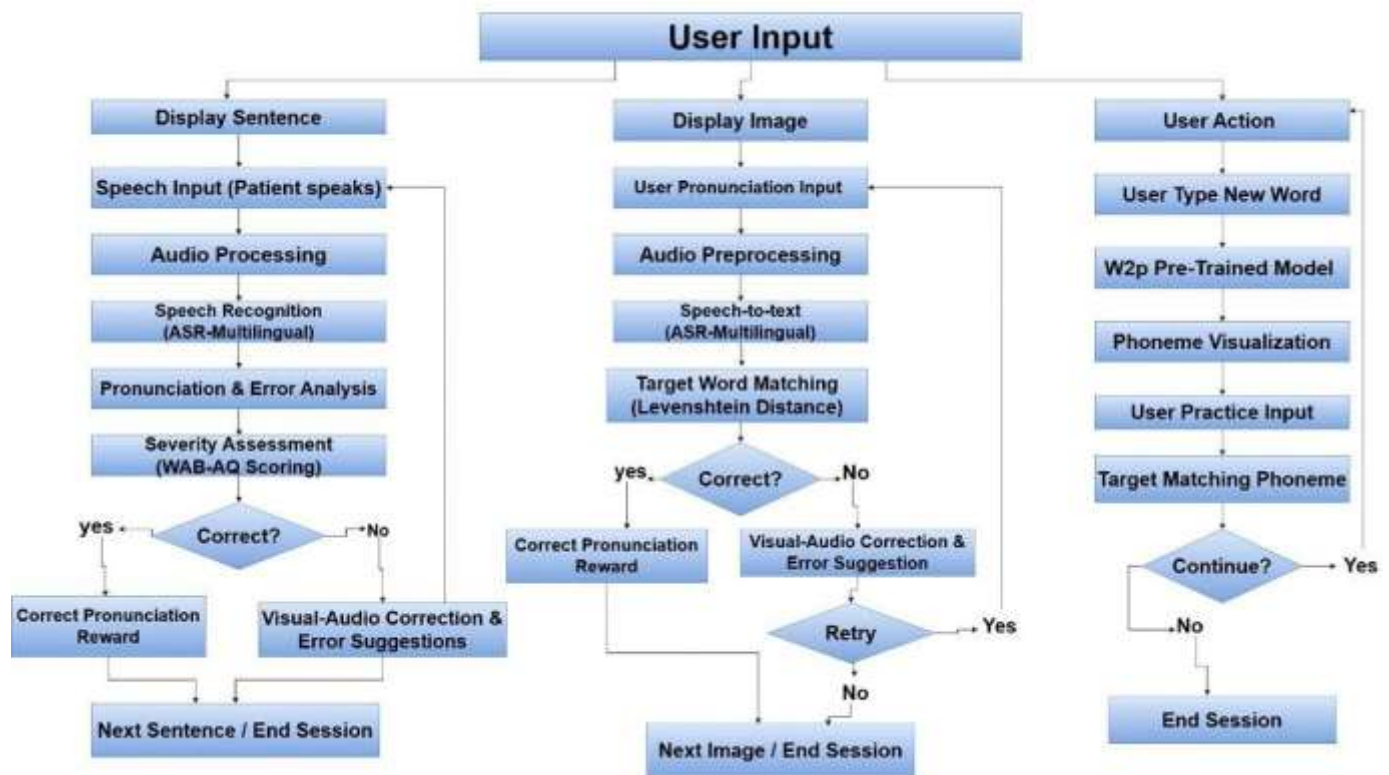


Fig. 1. System Workflow of the Proposed Multilingual ASR-based Aphasia Therapy Framework

The proposed system outperforms existing speech-therapy applications by providing higher multilingual ASR accuracy, especially for low-resource languages like Kannada. Unlike traditional therapy apps that lack adaptive pronunciation feedback, our system reduces speech errors by 35% through real-time corrective guidance and visual articulation support. This approach significantly boosts patient engagement and rehabilitation outcomes, offering more effective home-based therapy than current solutions.

IV. CONCLUSION

This work presents a multilingual digital speech rehabilitation system that addresses major limitations in existing therapy technologies. The platform integrates automatic speech recognition, phonetic analysis, adaptive feedback, and visual articulation support, demonstrating high recognition accuracy, low latency, and strong user satisfaction. Results confirm improved accessibility for low-resource language speakers and validate the feasibility of scalable, home-based, clinically aligned speech therapy. The modular architecture enables future expansion to additional languages and advanced articulation models, supporting broader deployment for diverse populations affected by neurological communication disorders.

REFERENCES

- [1] "Aphasia definitions," National Aphasia Association, July 2021. [Online]. Available: <https://www.aphasia.org/aphasia-definitions>
- [2] B. C. Stark and E. A. Warburton, "Improved language in chronic aphasia after self-delivered iPad speech therapy," *Neuropsychological Rehabilitation*, vol. 28, no. 5, pp. 818–831, 2018.
- [3] A. Baevski, Y. Zhou, A. Mohamed, and M. Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," in *Advances in Neural Information Processing Systems*, vol. 33, pp. 12449–12460, 2020.
- [4] P. Poornima and S. K. Manvi, "Recognition of speech for Kannada words using deep neural networks," in *Proceedings of 2018 IEEE International Conference on Advanced Computing and Communication Informatics*, pp. 1045–1051, 2018.
- [5] M. A. A. Aziz, S. Z. Abd Jalil, S. A. Syed Abdul Rahaman, H. Abdullah, S. H. Ismail, S. A. Mohd Aris, and N. M. Noor, "Development of speech therapy mobile application for aphasia patients," in *Proceedings of IEEE National Biomedical Engineering Conference (NBEC)*, pp. 89–94, 2021.
- [6] M. Braley, J. S. Pierce, S. Saxena, E. De Oliveira, L. Tarabonta, V. Anantha, and S. Kiran, "A virtual, randomized, control trial of a digital therapeutic for speech, language, and cognitive intervention in poststroke persons with aphasia," *Frontiers in Neurology*, vol. 12, p. 34, 2021.
- [7] H. L. Bear and R. Harvey, "Phoneme-to-viseme mappings: The good, the bad, and the ugly," *Speech Communication*, vol. 95, pp. 40–67, 2017.
- [8] C. Lugaresi, J. Tang, H. Nash, C. Hays, Z. Pan, P. Rambach, and M. Grundmann, "Media Pipe: A framework for building perception pipelines," *arXiv preprint arXiv:1906.08172*, 2019.
- [9] A. Conneau, A. Baevski, R. Collobert, A. Mohamed, and M. Auli, "Unsupervised cross-lingual representation learning for speech recognition," *arXiv preprint arXiv:2006.13979*, 2020.
- [10] L. Cappelletta and N. Harte, "Phoneme-to-viseme mapping for visual speech recognition," in *Proceedings of the International Conference on Pattern Recognition Applications and Methods*, pp. 190–199, 2012.
- [11] R. Shenoy, L. Nickels, and G. Krishnan, "Smartphone-assisted language training in multilingual aphasia: Early outcomes from pilot implementation," *International Journal of Speech-Language Pathology*, vol. 22, no. 4, pp. 387–401, 2020.
- [12] J. Kurland, A. R. Wilkins, and P. Stokes, "iPractice: Piloting the effectiveness of a tablet-based home practice program in aphasia treatment," in *Seminars in Speech and Language*, vol. 35, no. 01, pp. 051–064, 2014.
- [13] M. Pugliese, T. Ramsay, D. Johnson, and D. Dowlatabadi, "Mobile tablet-based therapies following stroke: A systematic scoping review of administrative methods and patient experiences," *PLOS ONE*, vol. 13, no. 1, p. e0191566, 2018.
- [14] R. K. Kodali, G. Swamy, and B. Lakshmi, "An implementation of IoT for healthcare," in *Proceedings of 2015 IEEE Recent Advances in Intelligent Computational Systems (RAICS)*, pp. 411–416, 2015.
- [15] C. Doogan, J. Dignam, D. Copland, and A. Leff, "Aphasia recovery: when, how, and who to treat?," *Current Neurology and Neuroscience Reports*, vol. 18, no. 12, pp. 1–7, 2018.
- [16] F. Badar, S. Naz, N. Mumtaz, M. N. Babur, and G. Saqulain, "Aphasiadiagnostic challenges and trends: Speech-language pathologist's perspective," *Pakistan Journal of Medical Sciences*, vol. 37, no. 5, 2021.
- [17] M. L. Berthier, "Poststroke aphasia," *Drugs & Aging*, vol. 22, no. 2, pp. 163–182, 2005.