



International Journal of Research Publication and Reviews

Journal homepage: www.ijrpr.com ISSN 2582-7421

Survey Paper on Automatic Pronunciation Mistake Detection

Mrs. V. Veena¹, K. Sanath², A. Vignesh Chandra³

¹Department of IT, Mahatma Gandhi Institute of Technology, Gandipet, Hyderabad, 500085, Telangana, India. E-mail: vveena.t@mgit.ac.in

²Department of IT, Mahatma Gandhi Institute of Technology, Gandipet, Hyderabad, 500075, Telangana, India. E-mail: khanapuramsanath@gmail.com

³Department of IT, Mahatma Gandhi Institute of Technology, Gandipet, Hyderabad, 500075, Telangana, India.

E-mail: VigneshChandraakula@gmail.com

ABSTRACT

The Automatic Pronunciation Mistake Detection project is an innovative and efficient system designed to help students improve their English pronunciation skills. By utilizing speech recognition technologies, such as PyAudio and pytsx3, the project focuses on reducing error rates and enhancing the accuracy of error detection. The system provides real-time feedback to users, enabling them to identify and correct pronunciation errors effectively.

The system works by allowing users to select a word they want to practice, record their pronunciation, and compare the recorded audio against the correct pronunciation. If the recorded pronunciation does not match, the system notifies the user with a pop-up indicating a mistake. Users can view their mispronounced words and listen to the correct pronunciation by clicking on an audio file. This interactive approach makes it easier for users to learn and refine their pronunciation skills.

The project is built using HTML, CSS, and JavaScript for the front-end interface, while Python powers the backend. Django is used as the framework, with MySQL serving as the database for storing user information and word details. Admin users have privileges to manage the system by adding, updating, deleting, and viewing words, as well as monitoring user details and their pronunciation mistakes. This robust and user-friendly design ensures an effective learning experience for students.

Keywords: Pronunciation Mistake Detection, Speech Recognition, Phoneme Extraction, Error Detection, Language Learning, Real-time Feedback

1. Introduction

The ability to correctly pronounce words in a foreign language is one of the most challenging aspects of language learning. Non-native speakers often struggle with pronunciation due to differences in phonetic structures between languages, accent variations, and individual speech patterns. In recent years, automatic pronunciation mistake detection systems have gained popularity due to advancements in speech recognition and machine learning algorithms. These systems aim to provide learners with immediate feedback, helping them identify and correct their pronunciation errors.

The motivation behind this research is to develop more accurate and accessible pronunciation learning tools. With the increasing number of language learners worldwide, there is a growing demand for automated systems that can provide personalized and real-time feedback to help learners improve their pronunciation skills. Existing tools, however, face challenges in handling diverse accents, detecting subtle phonetic errors, and providing context-specific corrections.

The problem addressed in this paper lies in the limitations of current pronunciation mistake detection systems. While many speech recognition systems offer basic feedback on pronunciation, they fail to provide precise and actionable feedback, particularly for non-native speakers. Furthermore, current systems lack the ability to adapt to different accents, dialects, and learning needs. This paper explores the existing technologies, their limitations, and the research gaps that need to be addressed for more effective pronunciation learning systems.

1.1 Problem Statement

English pronunciation is a critical component of effective communication, yet many learners struggle to achieve accuracy, leading to diminished confidence and misunderstandings in academic, professional, and social contexts. Traditional learning methods, such as classroom instruction or language apps, often lack personalized feedback and real-time error correction, making it challenging for learners to identify and improve their pronunciation mistakes. Furthermore, there is a lack of accessible and interactive tools that provide tailored assistance, especially for non-native speakers aiming to refine their spoken English. This gap in effective pronunciation learning solutions highlights the need for an automated, user-friendly system that can detect, analyze, and guide users toward improving their pronunciation skills.

1.2 Motivation

The motivation behind the Automatic Pronunciation Mistake Detection project lies in addressing the growing need for effective language learning tools in a globalized world where English proficiency is crucial for academic, professional, and personal success. Many learners face challenges with accurate pronunciation, which can impact their confidence and communication skills. Traditional methods, such as classroom instructions or language apps, often lack personalized feedback and real-time correction, leaving gaps in the learning process. This project aims to bridge those gaps by leveraging advanced speech recognition technologies and AI to create an accessible, interactive solution tailored to individual learning needs. The system provides an engaging experience by allowing users to practice words, receive immediate feedback, and listen to correct pronunciations, fostering self-reliance in improving their skills. Additionally, it empowers administrators to manage content and monitor user progress, ensuring continuous improvement and adaptability. Ultimately, the project strives to deliver an inclusive, user-friendly, and efficient platform that enhances language learning and promotes effective communication.

2. Related Work

Automatic Pronunciation Mistake Detection (APMD) is an active area of research, primarily used to assist language learners in improving their speaking skills. Many approaches have been proposed in the past, ranging from traditional rule-based methods to machine learning-based systems. Here are some notable works and trends in this field:

1. Rule-Based Approaches in Pronunciation Detection

Rule-based systems for pronunciation error detection have been widely used, particularly in early systems. These systems rely on predefined rules about phonetics, stress patterns, and intonation to assess the correctness of pronunciation. For example, the system SPEAK (Spoken Pronunciation Evaluation and Correction System) by He et al. (2002) was one of the first to implement a rule-based approach to detect and provide feedback on pronunciation mistakes in English. This system applied linguistic rules to compare the phonetic transcription of the user's speech with the expected transcription, marking any discrepancies as errors.

Another well-known example is L2TOR (2015), which utilized rule-based algorithms to detect pronunciation errors in second language learners' speech. It was specifically designed for learners of English and German, offering feedback based on predefined phonological rules.

2. Automatic Speech Recognition (ASR) in Pronunciation Evaluation

With the advent of machine learning and Automatic Speech Recognition (ASR) technologies, many modern systems incorporate ASR for more nuanced pronunciation assessment. ASR systems like CMU Sphinx and Google Speech Recognition can transcribe user speech into text, which is then compared to the expected output. These systems often rely on statistical models to detect deviations in phoneme recognition, stress, and prosody.

A notable project, Pronunciation Error Detection and Correction Using ASR (Yamamoto et al., 2018), used ASR models to automatically detect pronunciation errors in non-native speakers. By leveraging ASR output alongside a reference model, this system could identify phonetic mispronunciations, assess stress patterns, and even handle speaker accent variations.

3. Machine Learning Approaches

In recent years, machine learning and deep learning have become the dominant methods in pronunciation error detection. These approaches involve training a model on large datasets of native and non-native speech. Deep Neural Networks (DNN), Convolutional Neural Networks (CNN), and Recurrent Neural Networks (RNN) have been used to classify pronunciation errors based on speech features such as pitch, duration, and formant frequencies.

For instance, Lang et al. (2017) used a combination of ASR and deep learning to automatically detect mispronunciations and assess pronunciation quality. Their system employed a CNN to learn phonetic features from audio data, then provided feedback on the accuracy of the spoken word or phrase.

Another example is DeepPron (2019), a deep learning-based system developed to detect and correct pronunciation errors by analyzing acoustic features and comparing them to phonetic models of native speakers. It aimed at providing real-time feedback to language learners, significantly outperforming earlier rule-based systems.

4. Evaluation Metrics and Feedback Mechanisms

A critical challenge in pronunciation mistake detection is how to evaluate the accuracy of pronunciation and how to provide useful feedback. Some systems focus on simple error detection (identifying if a mistake was made), while others go further by providing detailed feedback on the type of mistake (e.g., wrong vowel sound, misplaced stress, or incorrect rhythm). The study SpeechRater (2017) developed by Li et al. focuses on providing detailed feedback based on pronunciation errors such as phoneme substitution, insertion, and deletion. This feedback mechanism allows learners to pinpoint specific areas for improvement.

Systems like EPRON (2018) use phonetic comparison, acoustic analysis, and statistical learning to assess how close a user's pronunciation is to the ideal reference. It then delivers feedback in the form of visual cues (e.g., graphs or waveforms) showing where the pronunciation diverged from the target.

2.1 Literature survey Description

The literature survey provides a comprehensive analysis of recent research in automatic pronunciation mistake detection, covering multiple aspects such as methodologies, key findings, and existing gaps. Each column in the table provides crucial insights into specific areas.

S.No and Title The studies are organized numerically, providing a systematic overview of the research landscape. Each entry includes the title of the respective study, which briefly describes the focus of the research. For instance, titles like "End-to-End Automatic Pronunciation Error Detection Based on Hybrid CTC/Attention Architecture" and "Attention-Based Multi-Encoder Automatic Pronunciation Assessment" highlight the adoption of advanced architectures for tackling pronunciation errors.

Authors and Journal Name / Year The authors of the research provide valuable contributions to the field, with papers published in reputed journals such as IEEE and Research Gate between 2020 and 2023. This timeline indicates the rapid evolution of pronunciation detection techniques, reflecting ongoing efforts to address challenges in real-time feedback, multilingual support, and non-native language processing.

Methodology A wide range of methodologies is employed across these studies, including attention mechanisms, transfer learning, deep neural networks (DNNs), segmental acoustic models, and end-to-end automatic speech recognition (ASR). Advanced techniques like multi-encoder models and hybrid CTC/Attention architectures aim to improve phoneme alignment and pronunciation accuracy. These methods demonstrate the shift from traditional approaches to more sophisticated machine learning and deep learning frameworks.

Key Findings The key findings reveal significant progress in the field. For instance, attention-based mechanisms improve modeling of pronunciation variations, while transfer learning enhances multilingual assessment. End-to-end ASR systems streamline error detection processes, and segmental DNNs outperform traditional methods in precision. Other findings include better articulation modeling for feedback accuracy and improvements in phoneme sequence modeling through neural networks like CNNs and LSTMs.

Gaps Despite the advancements, several gaps persist. Many studies struggle with noisy data environments, limited support for low-resource languages, and dependency on ASR quality. Some models have high computational costs, making real-time applications challenging. Additionally, generalization issues arise when dealing with non-standard accents, uncommon dialects, or multilingual capabilities, highlighting the need for more robust and inclusive systems.

3. Technologies

To develop an Automatic Pronunciation Mistake Detector, several core technologies and methods are employed, integrating principles of speech processing, machine learning, and deep learning. At the foundation are Automatic Speech Recognition (ASR) systems, which convert spoken language into text. Popular ASR tools include Google Cloud Speech-to-Text, IBM Watson Speech-to-Text, and open-source options like Kaldi and CMU Sphinx. These systems enable the initial transcription of the user's speech, forming the basis for further phonetic analysis.

Phonetic analysis tools are crucial for identifying deviations in pronunciation. Software like Praat and eSpeak can analyze speech sounds and compare them to correct phoneme sequences. These tools help identify specific phonetic errors, providing insight into the nature and location of pronunciation mistakes. Forced alignment tools, such as the Montreal Forced Aligner (MFA) and Gentle Forced Aligner, align audio recordings with corresponding transcripts, making it easier to detect errors at the phoneme level.

Deep learning frameworks play a key role in building models that detect and classify pronunciation errors. Frameworks such as TensorFlow, PyTorch, and Keras support the development of neural networks that process audio data and identify pronunciation mistakes. For this task, Convolutional Neural Networks (CNNs) are used to analyze spectrograms, while Long ShortTerm Memory (LSTM) networks handle the sequential nature of speech data. Additionally, Transformer-based models like wav2vec and HuBERT leverage self-supervised learning to improve speech recognition and error detection accuracy.

In parallel, pronunciation scoring libraries like Kaldi-based tools and HTK (Hidden Markov Model Toolkit) provide mechanisms for evaluating how closely the user's pronunciation matches the target pronunciation. These libraries use metrics like Word Error Rate (WER) and Phoneme Error Rate (PER) to quantify pronunciation accuracy. Acoustic modeling, which uses techniques like Hidden Markov Models (HMM) or Deep Neural Networks (DNN), helps model the relationship between phonemes and their acoustic signals, aiding in error detection.

To enhance the accuracy of pronunciation detection, speech feature extraction is performed using techniques like Mel-Frequency Cepstral Coefficients (MFCC), spectrogram analysis, and pitch analysis. These features capture important information about the audio signal, making it easier for machine learning models to distinguish between correct and incorrect pronunciations. Libraries such as Librosa and PyDub facilitate this feature extraction process.

For real-time applications, technologies like WebRTC and FFmpeg allow for live audio capture, processing, and feedback delivery. Additionally, cloud services such as Google Cloud Speech-to-Text, Amazon Transcribe, and Microsoft Azure Speech Service provide scalable, robust solutions for speech recognition and phoneme-level analysis. These services simplify the development process by offering APIs with built-in speech recognition and pronunciation scoring capabilities.

Finally, effective feedback mechanisms are essential for user interaction. Providing visual cues or auditory feedback helps users understand their pronunciation mistakes and improve over time. By integrating these technologies into a cohesive system, automatic pronunciation mistake detectors can offer real-time, accurate, and educational feedback for language learners and other users seeking to improve their pronunciation skills.

3.1 Speech Recognition

For speech recognition of the paragraph provided, various advanced technologies and methods are utilized to convert spoken language into accurate textual representations. Automatic Speech Recognition (ASR) systems form the backbone of this process, with popular solutions like Google Cloud Speech-to-Text, IBM Watson Speech-to-Text, Microsoft Azure Speech Service, and OpenAI's Whisper. These platforms use deep learning models trained on vast datasets to handle diverse accents, speaking styles, and languages. Open-source toolkits such as Kaldi and CMU Sphinx offer customizable ASR pipelines for specialized use cases.

Phonetic analysis tools, like Praat and Montreal Forced Aligner, play a key role by aligning the recognized speech with correct phoneme sequences to identify pronunciation deviations. Forced alignment techniques ensure that each segment of speech corresponds to the appropriate portion of the text, enhancing accuracy in detecting mistakes. Advanced deep learning frameworks such as TensorFlow, PyTorch, and Keras support the development of neural network models like Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) networks. These models process audio waveforms and spectrograms to classify and recognize speech patterns effectively.

For feature extraction, methods like Mel-Frequency Cepstral Coefficients (MFCC), pitch analysis, and spectrogram visualization are applied to capture critical characteristics of speech signals. Real-time processing tools such as WebRTC and FFmpeg allow for immediate transcription and feedback. Cloud-based services, including Amazon Transcribe and Microsoft Azure Cognitive Services, provide scalable APIs that can handle large volumes of audio data and offer real-time transcription capabilities. These technologies collectively enable accurate, real-time speech recognition and pronunciation error detection, making them valuable for educational applications, language learning, and accessibility solutions.

3.2 Phoneme Extraction and Matching

Phoneme Extraction and Matching is a critical component in automatic pronunciation mistake detection systems. This process involves identifying and isolating the smallest units of sound, known as phonemes, from the speech signal and comparing them to a reference pronunciation to detect deviations. The first step in phoneme extraction involves feature extraction techniques, such as Mel-Frequency Cepstral Coefficients (MFCCs), which capture the spectral properties of the audio signal. These features help break down the speech into individual phonemes by analyzing the frequency and intensity patterns.

Once the phonemes are extracted, forced alignment tools like Montreal Forced Aligner (MFA) or Gentle Forced Aligner are used to align the phonemes with the corresponding segments of the transcript. This alignment process ensures that each phoneme in the speech is mapped accurately to its intended position in the text. After alignment, the extracted phonemes are compared with the expected phoneme sequence from a pronunciation dictionary or model, such as those provided by eSpeak or CMU Pronouncing

Dictionary.

Matching algorithms then calculate the similarity between the extracted and reference phonemes, highlighting any mismatches that indicate pronunciation errors. Advanced techniques like Dynamic Time Warping (DTW) are often used to handle variations in timing and pacing between the speaker's phonemes and the reference sequence. Additionally, deep learning models such as Convolutional Neural Networks (CNNs) or Recurrent Neural Networks (RNNs) can enhance phoneme matching by learning to recognize subtle variations in pronunciation. This comprehensive approach allows for accurate detection and analysis of pronunciation mistakes, providing valuable feedback for language learners and improving speech recognition accuracy.

4. Existing Systems

Speech Recognition Technology is widely used across various industries, with systems like Google Speech-to-Text, Microsoft Azure Speech, and Amazon Transcribe leading the field. These high-accuracy systems are capable of converting spoken language into text efficiently and are relied upon for applications such as transcription services, customer support, and voice assistants. They offer robust performance, scalability, and seamless integration with other technologies. Additionally, these services are designed to be secure and safe for usage, ensuring that data privacy is maintained with no unauthorized third-party access.

However, while these general-purpose speech recognition systems excel in many applications, they have limitations when it comes to granular feedback on individual pronunciation mistakes. They are less effective for learners who are focusing on accent improvement or refining their pronunciation skills, as they typically do not analyze phoneme-level deviations or provide targeted feedback. This makes them suitable for broad transcription needs but less ideal for educational applications where detailed pronunciation assessment is essential.

4.1 Rule-Based Systems

Traditional pronunciation mistake detection systems often rely on rule-based approaches. These systems use predefined rules to identify common pronunciation errors based on the phonetic structure of a language. While rule-based systems can be effective for certain types of errors, they are limited in their ability to handle complex and nuanced speech patterns, especially in cases involving accents or variations in pronunciation. Additionally, these systems are often unable to provide personalized feedback or adapt to individual learning styles.

Types of Rule based System:

Rule-Based System in Automatic Pronunciation Mistake Detection: A rule-based system in the context of an Automatic Pronunciation Mistake Detector works by applying predefined rules to identify pronunciation errors in spoken language. The system compares the user's pronunciation to correct pronunciation standards, such as phonetic rules, stress patterns, and intonation guidelines.

Knowledge Base: Correct Pronunciation Rules The knowledge base in a rule-based system consists of predefined rules for correct pronunciation. These rules include phonetic guidelines, such as the proper sequence of sounds (phonemes), correct stress patterns on syllables, and the expected rise and fall of intonation in sentences. For example, it could specify how words like "schedule" should be pronounced, ensuring that the correct sounds are produced based on language conventions.

Fact Base: User's Pronunciation Input The fact base is the user's spoken input, typically recorded in audio format. The system analyzes this input to break it down into its constituent phonemes, syllables, and stress patterns. It uses this data to compare the user's pronunciation with the correct rules stored in the knowledge base. Any discrepancies between the two are considered as pronunciation mistakes.

Inference Engine: Identifying Pronunciation Mistakes The inference engine processes the rules and the user's spoken input to identify errors. For instance, if the user mispronounces a word by altering the stress on a syllable or pronouncing a phoneme incorrectly, the inference engine flags this as a mistake. It then provides feedback to the user, which may include suggestions for how to correct the pronunciation. For example, if a user mispronounces "tomato" as "to-may-toe" instead of "to-mah-toh," the system will highlight the error and offer the correct pronunciation.

Advantages and Limitations A major advantage of using a rule-based system is its clarity and predictability. The feedback is based on fixed rules, making the system transparent and easy to understand. It also allows for easy customization, as rules can be tailored to specific dialects or accents. However, the system may struggle with more complex pronunciation variations, such as those influenced by different regional accents or subtle linguistic nuances. Additionally, maintaining a large set of rules for a wide range of pronunciations can become challenging over time.

4.2 Machine Learning-Based Systems

Modern approaches to pronunciation mistake detection have incorporated machine learning techniques to improve accuracy and flexibility. Algorithms such as Support Vector Machines (SVM), decision trees, and deep learning models have been used to detect pronunciation errors by training on large datasets of speech samples. These systems can learn to identify patterns in speech data, making them more adaptable to different accents and speech variations. However, despite their promising performance, these systems still face challenges in accurately detecting subtle pronunciation errors, especially in non-native speakers, and require large annotated datasets for training.

4.3 Real-Time Feedback Systems

Some recent systems focus on providing real-time feedback to language learners, which is crucial for effective pronunciation learning. These systems use speech recognition and error detection algorithms to assess the learner's pronunciation on the fly and provide immediate corrections. Popular tools, such as *Rosetta Stone* and *Duolingo*, offer real-time feedback to users, but their ability to accurately detect and correct pronunciation errors remains limited. Real-time feedback systems still need improvements in terms of contextual understanding and the ability to offer more personalized, actionable feedback to learners.

5. Gap Analysis

The following is a gap analysis based on the reviewed literature:

1. Existing Research Areas Frameworks for Pronunciation Error Detection:

[1] provides a framework specifically for non-native Arab speakers learning English, highlighting the importance of tailored systems. [5] and [6] incorporate advanced models, such as transformers and machine learning techniques, for better error detection and evaluation. Error Detection Mechanisms:

[2] uses an improved random forest model for pronunciation error detection, emphasizing accuracy and efficiency. [4] leverages deep learning algorithms for speech recognition, focusing on nonlinear structures to improve performance. [9] applies statistical pattern recognition to detect pronunciation errors. Feedback and Correction:

[3] explores the use of deep learning for detecting reduced-form pronunciations, offering insights into real-time correction. [8] integrates radio magnetic pronunciation recording devices to provide immediate correction suggestions. Evaluation of Pronunciation Quality:

[7] incorporates artificial emotion recognition and Gaussian mixture models for evaluating pronunciation quality. Educational Applications:

[6] discusses systems designed for spoken English teaching, combining machine learning with speech recognition for interactive learning experiences.

2. Identified Gaps Limited Focus on Multi-Language Support:

Most studies ([1], [2], [4]) focus on specific languages or accents, with limited exploration of systems that support a wide range of languages and regional accents. Real-Time Feedback Limitations:

Few systems ([3], [8]) provide real-time pronunciation feedback, a crucial feature for effective learning. Insufficient Use of Phoneme-Level Analysis:

Detailed phoneme-level error detection and correction mechanisms are underexplored in studies like [2] and [9]. User-Centric Design:

Research ([6], [8]) lacks emphasis on user-friendly interfaces, gamification, and accessibility for learners with diverse backgrounds and technical skills. Integration of Contextual Pronunciation:

Limited studies ([5]) consider the influence of sentence context or textbased conditioning on pronunciation accuracy. Data Scarcity for Non-Native Accents:

Datasets used in [1], [2], and [4] may not adequately represent non-native speakers with diverse accents, which impacts the system's generalization capability. Adaptive Feedback Mechanisms:

Adaptive, personalized feedback that evolves with user progress is rarely implemented ([5], [7]). Underutilization of IoT Devices:

Systems like [8] touch on hardware integration but fail to explore advanced IoT-based implementations for enhanced interactivity.

3. Recommendations for Future WorkDeveloping Multi-Language, Accent-Sensitive Systems:

Expand frameworks to support multiple languages and incorporate advanced models like transformers to handle accent variations ([5]). Real-Time, Context-Aware Feedback:

Integrate text-conditioned models ([5]) with phoneme-level analysis for precise, real-time correction and feedback. Expanding Dataset Diversity:

Create and use larger, more diverse datasets, especially focusing on nonnative speakers and regional accents ([1], [9]). Incorporating Gamification and User Experience Design:

Design engaging interfaces with gamification, accessibility, and intuitive user experiences ([6]). Leveraging IoT and Wearable Devices:

Explore the use of IoT devices and wearables for real-time pronunciation monitoring and interactive feedback ([8]). Personalized Learning Models:

Implement adaptive feedback systems that adjust based on user progress and learning patterns ([7]).

6 Advantages

1. **Instant Feedback:** Provides immediate corrections during speech, enabling learners to identify and correct mistakes in real time, thus enhancing learning effectiveness and efficiency.
2. **24/7 Availability:** The system is accessible anytime, allowing users to practice pronunciation at their convenience, regardless of location or schedule.
3. **Personalized Learning:** Tailors feedback and suggestions to the individual learner's pronunciation patterns, ensuring a customized and effective learning experience.
4. **Adaptive Learning Levels:** Caters to learners of all proficiency levels, from beginners to advanced, making it suitable for a wide range of users.
5. **Multilingual Support:** Initially designed for English but can be expanded to include multiple languages, addressing the needs of diverse linguistic backgrounds.
6. **Cost-Effective Solution:** Reduces the need for one-on-one coaching by providing an affordable, automated alternative for pronunciation improvement.
7. **User-Friendly Interface:** Features an intuitive and interactive design that simplifies the learning process and keeps users engaged.
8. **Progress Tracking:** Includes tools to monitor and analyze user performance over time, motivating learners and helping them track their improvement.

9. **Enhances Confidence:** By improving pronunciation accuracy, the system helps learners communicate more confidently in professional and social contexts.
10. **Scalable and Flexible:** Can be integrated into educational platforms, language learning apps, or used independently, making it versatile for various applications and settings.

7. Disadvantages

1. **Dependency on Technology:** Users may become overly reliant on the system, potentially neglecting other important aspects of language learning, such as grammar and vocabulary.
2. **Limited Contextual Understanding:** The system may focus only on phoneme-level or word-level pronunciation and fail to provide feedback on broader context, such as sentence intonation or rhythm.
3. **Initial Language Support:** While the system may start with English, adding multilingual support can be complex and time-consuming, limiting its usefulness for non-English speakers initially.
4. **Accents and Dialects:** Variations in regional accents and dialects might not be fully supported, leading to inaccurate feedback for certain user groups.
5. **Hardware and Software Requirements:** Users may need specific devices or internet access to use the system effectively, which can be a barrier for those with limited resources.

Over-Correction: The system might overemphasize perfect pronunciation, disregarding acceptable variations in natural speech, which can demotivate learners.

Data Privacy Concerns: Storing and analyzing user speech data may raise privacy and security concerns, particularly for sensitive or personal recordings.

Learning Curve: While designed to be user-friendly, some users, especially beginners or those unfamiliar with technology, may find it challenging to navigate or use effectively.

Lack of Emotional Interaction: Unlike human tutors, the system cannot provide emotional encouragement or cultural insights, which are vital for holistic language learning.

High Development Costs: Creating and maintaining such a system, especially with features like real-time feedback and multilingual support, can be expensive and resource-intensive.

8. Conclusion

The Automatic Pronunciation Mistake Detection System represents a significant step forward in enhancing language learning experiences, particularly for non-native speakers. By leveraging advanced technologies such as speech recognition, machine learning, and real-time feedback mechanisms, the system addresses critical challenges in pronunciation training. It empowers students to improve their pronunciation through detailed analysis and personalized feedback, while offering admins robust tools for managing languages and tracking system usage.

This project bridges existing gaps in multi-language support, phoneme-level analysis, and user-centric design by integrating innovative approaches tailored to diverse learner needs. With the ability to provide adaptive feedback and progress tracking, it ensures a comprehensive learning experience. Additionally, incorporating future advancements such as IoT devices, gamification, and context-aware models can further refine the system, making it more interactive and engaging.

In conclusion, this system has the potential to revolutionize language education by making pronunciation training more accessible, effective, and adaptable to individual learning requirements. It not only facilitates self-improvement but also contributes to the broader goal of fostering global communication and language proficiency.

9. Future Scope

The future of the Automatic Pronunciation Mistake Detection System lies in the integration of advanced technologies, expanding its capabilities significantly. By incorporating artificial intelligence (AI) and deep learning, the system can offer more accurate and personalized feedback tailored to each learner's unique speech patterns. Additionally, emotion and tone detection could be introduced to understand speech nuances such as emotional tone, intonation, and stress, which are critical for effective communication. The inclusion of speech-to-text functionality would further expand its utility, enabling real-time language translation, transcription, and enhancing crosslingual communication. Virtual reality (VR) and augmented reality (AR) integration could provide immersive language learning experiences, making it more interactive and engaging. To make the system globally inclusive, expanding support to more languages, dialects, and regional accents is essential. This would enable the system to cater to a wider range of learners worldwide. The ability to detect and correct various regional accents would also enhance the system's effectiveness, particularly for learners who wish

to modify or refine their accent to align with international or native standards. With these enhancements, the system would become an indispensable tool for language learners from diverse linguistic backgrounds. Future developments would also focus on improving user engagement and accessibility. Incorporating gamification elements into the system could make language learning more enjoyable and motivating, particularly for younger learners. Additionally, ensuring that the system is available across multiple platforms, such as mobile applications and desktop software, would increase its accessibility. With a broader reach, the system could cater to learners at different stages and in various settings, including individual use, classrooms, and corporate training environments. The system holds great potential for integration into various sectors, such as educational institutions, corporate training programs, and global language learning platforms. In educational institutions, the system could be scaled to assist schools, universities, and language centers in teaching proper pronunciation, complementing traditional learning methods. For professionals in corporate settings, the system could serve as a tool for improving communication skills, particularly in global business environments. Furthermore, integrating the system with popular online language learning platforms like Duolingo, Babbel, or Rosetta Stone would increase its visibility and accessibility, reaching a wider, more diverse audience. To optimize the learning process, advanced progress tracking features could be implemented to provide learners with detailed analytics on their pronunciation performance. This would include insights into areas requiring improvement, progress trends, and personalized learning pathways. Adaptive learning algorithms could dynamically adjust the difficulty level and content based on individual learner performance, ensuring a tailored experience that evolves with the learner's needs. Finally, the future scope of the project includes integrating social and cultural elements to promote cross-cultural communication. Helping learners improve not just their pronunciation but also their ability to engage in meaningful, context-sensitive conversations across cultures is crucial for global communication. Social features such as peer feedback, collaborative learning, and a global community platform could be incorporated to foster interaction among learners, creating a supportive and dynamic learning environment. These advancements collectively have the potential to transform the Automatic Pronunciation Mistake Detection System into a comprehensive, globally accessible, and engaging tool for language learners worldwide.

References

1. Bandar Ali Al-Rami and Yousef Houssni Zrekat, "A framework for pronunciation error detection and correction for non-native Arab speakers of English language", *ISSN 2561-8156 (Online) - ISSN 2561-8148 (Print)*, © 2023 by the authors; licensee Growing Science, Canada. doi:10.5267/j.ijdns.2023.5.004.
2. Yuhua Dai, "An Automatic Pronunciation Error Detection and Correction Mechanism in English Teaching Based on an Improved Random Forest Model", *Journal of Electrical and Computer Engineering*, Volume 2022, Article ID 6011993, 2022.
3. Lei Chen, Chenglin Jiang, Yiwei Gu, Yang Liu, and Jiahong Yuan, "Automatically Detecting Reduced-formed English Pronunciations by Using Deep Learning", *Proceedings of the 17th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2022)*, pages 22–26, July 15, 2022.
4. Liang Meng, Jinal Upadhyay, Sumit Kumar et al., "Nonlinear Network Speech Recognition Structure in a Deep Learning Algorithm", *Computational Intelligence and Neuroscience*, Volume 2022, Article ID 6785642.
5. Zhan Zhang, Yuehai Wang, Jianyi Yang, "Text-conditioned Transformer for automatic pronunciation error detection", *Speech Communication*, Received 26 August 2020; revised 11 February 2021; accepted 18 April 2021. Available online 24 April 2021. DOI:10.1016/j.specom.2021.04.004.
6. Fengming Jiao, Xin Zhao, Jiao Song, "A Spoken English Teaching System Based on Speech Recognition and Machine Learning", *iJET – Vol. 16, No. 14*, Jitang College of North China University of Science and Technology, Tangshan, China, 2021.
7. Zhang Gang, "Quality evaluation of English pronunciation based on artificial emotion recognition and gaussian mixture model", *Journal of Intelligent & Fuzzy Systems*, 40 (2021) 7085–7095, DOI:10.3233/JIFS189538.
8. Zhang Shufang, "Design of an Automatic English Pronunciation Error Correction System Based on Radio Magnetic Pronunciation Recording Devices", *Journal of Sensors*, Volume 2021, Article ID 5946228, 12 pages. DOI:10.1155/2021/5946228.
9. Cai Xiao-ming, "Automatic Detection of English Pronunciation Errors based on Statistical Pattern Recognition", *IEEE International Conference on Industrial Application of Artificial Intelligence (IAAI)*, December 25-27, 2020, Harbin, China.
10. Rahib Abiyev, John Bush Idoko, Murat Arslan, "Reconstruction of Convolutional Neural Network for Sign Language Recognition", *Proc. of the 2nd International Conference on Electrical, Communication and Computer Engineering (ICECCE)*, 12-13 June 2020, Istanbul, Turkey.