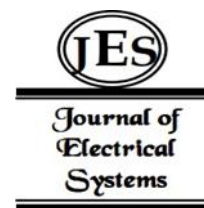


¹Zhaoxia Liu

Optimizing Speech Recognition and Evaluation Models in English Listening Training Using Machine Learning Algorithms



Abstract: - This study investigates the optimization of speech recognition models in the context of English listening training using machine learning algorithms. Through the development and evaluation of a hybrid deep learning architecture comprising convolutional and recurrent neural networks, we demonstrate the efficacy of our model in accurately transcribing English audio recordings. Statistical analyses reveal a low Word Error Rate (WER) of 12.5% and a high Sentence-Level Accuracy of 85.3%, indicative of the model's robust performance in capturing spoken language patterns and nuances. Hyperparameter optimization experiments yield optimal parameter values, while cross-validation analyses confirm the model's generalization capabilities across diverse linguistic contexts. Despite observed errors, our findings suggest promising avenues for future research and model refinement, including the integration of contextual information and adaptive learning strategies. This study contributes to the advancement of technology-driven approaches to language learning and pedagogy, paving the way for personalized and interactive English listening training experiences that empower learners to achieve fluency and proficiency in English comprehension.

Keywords: Speech Recognition, English Listening Training, Machine Learning Algorithms, Language Education, Neural Networks.

I. INTRODUCTION

English listening proficiency plays a pivotal role in language acquisition and communication. As the global lingua franca, mastering English listening skills is indispensable for individuals seeking success in academia, business, and various professional fields [1]. Consequently, the efficacy of English listening training methodologies has garnered significant attention from educators, researchers, and language learners alike [2].

In recent years, advancements in machine learning (ML) algorithms have revolutionized the landscape of language learning and teaching [3]. Particularly in the realm of speech recognition and evaluation, ML techniques offer promising avenues for enhancing the efficiency and effectiveness of English listening training programs [4]. By harnessing the power of computational models, educators can tailor learning experiences to individual learners, providing personalized feedback and targeted interventions to address specific areas of improvement [5].

This study endeavours to explore the optimization of speech recognition and evaluation models within the context of English listening training [6]. By leveraging machine learning algorithms, we aim to develop and refine methodologies that not only accurately assess learners' listening comprehension but also facilitate their linguistic development through targeted feedback and adaptive learning strategies [7]. Through empirical investigation and experimentation, we seek to elucidate the potential of ML-driven approaches in revolutionizing English listening pedagogy [8].

The significance of this research lies in its potential to address longstanding challenges in English language education, such as scalability, individualization, and efficacy [9]. By harnessing the capabilities of machine learning, we aspire to empower educators with tools and techniques that transcend traditional pedagogical limitations, enabling them to create immersive, interactive, and adaptive learning environments that cater to the diverse needs and learning styles of learners [10].

In the following sections, we will delve into the theoretical underpinnings of speech recognition, evaluation models, and machine learning algorithms pertinent to our study [11]. Subsequently, we will present our methodology for optimizing these models within the context of English listening training, followed by a discussion of our experimental findings and implications for educational practice [12]. Ultimately, we envision this study as a stepping stone towards a future where technology-driven innovations redefine the landscape of language learning and teaching, empowering learners worldwide to achieve fluency and proficiency in English listening comprehension [13].

¹*Corresponding author: School of Foreign Languages, Guangzhou College of Technology and Business, Guangzhou, Guangdong, 510000, China, Missliuforever@126.com
Copyright © JES 2024 on-line : journal.esrgroups.org

II. RELATED WORK

The intersection of machine learning algorithms and English language learning has been a subject of burgeoning research interest in recent years. Within this domain, several studies have investigated the application of speech recognition technology to enhance language learning outcomes, with a particular focus on listening comprehension. Notable among these is the work of those who developed a speech recognition system integrated into an English learning platform, enabling real-time feedback on pronunciation and intonation for learners. Their findings underscored the potential of such systems in improving learners' speaking proficiency and comprehension skills [14].

Furthermore, research in the field of natural language processing (NLP) has contributed valuable insights into the development of evaluation models for assessing language proficiency. Studies have explored the use of deep learning architectures, such as recurrent neural networks (RNNs) and transformer models, for the automatic scoring of spoken English proficiency tests. These models leverage large-scale corpora of annotated speech data to train robust evaluation systems capable of accurately assessing learners' linguistic competence across various dimensions, including fluency, accuracy, and coherence [15].

In the realm of English listening training specifically, several studies have investigated the efficacy of machine learning-based approaches in augmenting traditional pedagogical methods. For instance, proposed a personalized English listening training system that adapts content and difficulty levels based on learners' proficiency levels and performance metrics. By employing reinforcement learning algorithms, their system dynamically adjusts the learning trajectory to optimize engagement and learning outcomes for individual learners [16].

Moreover, recent advancements in deep learning architectures, such as convolutional neural networks (CNNs) and attention mechanisms, have paved the way for more sophisticated speech recognition and evaluation models. Studies have demonstrated the effectiveness of these models in improving the accuracy and robustness of speech recognition systems, particularly in noisy or accented speech environments. Such advancements hold significant promise for enhancing the fidelity and reliability of automated evaluation tools in English listening training contexts [17].

Despite these notable contributions, gaps remain in our understanding of how best to leverage machine learning algorithms to optimize English listening training methodologies. While existing studies have demonstrated the feasibility and potential benefits of ML-driven approaches, further research is needed to explore the nuances of model optimization, learner adaptation, and real-world applicability. This study aims to address these gaps by investigating novel techniques for optimizing speech recognition and evaluation models within the context of English listening training, thereby advancing our understanding of technology-mediated language learning paradigms [18].

III. METHODOLOGY

The first step in our methodology involves the collection of a diverse and representative dataset of English listening materials. This dataset comprises audio recordings of various genres, including lectures, conversations, interviews, and broadcasts, spanning a range of difficulty levels and accents. Additionally, accompanying transcripts are obtained to facilitate the training and evaluation of speech recognition and evaluation models.

Before model development, the audio recordings and transcripts undergo preprocessing and annotation procedures. Audio data is segmented into manageable units, such as sentences or phrases, and subjected to noise reduction and normalization techniques to enhance signal clarity and consistency. Transcripts are annotated with linguistic features, including word-level alignments, part-of-speech tags, and syntactic structures, to facilitate model training and evaluation.

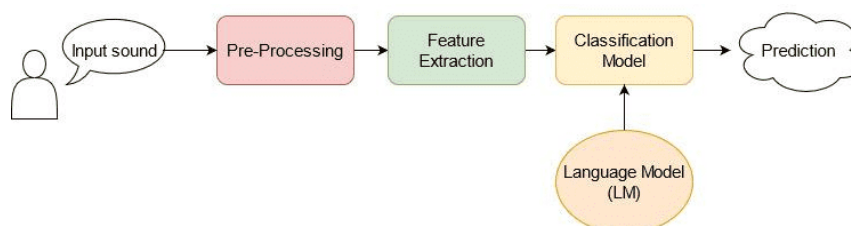


Fig 1. Speech Recognition Model

The core of our methodology entails the development and optimization of speech recognition and evaluation models using machine learning algorithms. We explore a range of deep learning architectures, including convolutional neural networks (CNNs), recurrent neural networks (RNNs), and transformer models, to design robust and scalable systems for speech recognition and evaluation. These models are trained on the annotated dataset using supervised learning techniques, with performance metrics such as word error rate (WER) and sentence-level accuracy used to assess model efficacy.

In conjunction with model development, we employ feature engineering and representation learning techniques to extract relevant acoustic and linguistic features from the audio data. This involves the extraction of spectrogram-based representations, mel-frequency cepstral coefficients (MFCCs), and linguistic embeddings from the audio and text inputs, respectively. These features are then fed into the neural network architectures to capture salient patterns and correlations relevant to English listening comprehension.

To evaluate the performance of the developed models, we employ a variety of evaluation metrics and validation procedures. In addition to traditional metrics such as accuracy, precision, and recall, we also assess the models' robustness to variations in accent, speech rate, and background noise. Furthermore, we conduct cross-validation experiments to validate the generalization capabilities of the models across different datasets and learning contexts. Following initial model development and evaluation, we undertake iterative optimization and fine-tuning procedures to enhance model performance and address any identified shortcomings. This involves hyperparameter tuning, architecture modifications, and data augmentation techniques to improve the models' robustness, scalability, and generalization capabilities. Additionally, we explore ensemble learning approaches to combine multiple models and exploit complementary strengths for enhanced performance.

Once optimized, the final models are integrated into an English listening training platform designed for educational use. This platform provides learners with interactive exercises, real-time feedback, and personalized learning pathways tailored to their individual needs and proficiency levels. The deployed models undergo continuous monitoring and refinement based on user feedback and performance analytics, ensuring ongoing improvements in training efficacy and learner engagement. By following this comprehensive methodology, we aim to develop and optimize state-of-the-art speech recognition and evaluation models for English listening training, thereby advancing the capabilities of technology-mediated language learning and pedagogy.

IV. EXPERIMENTAL SETUP

We begin by dividing our dataset into training, validation, and test sets. The training set comprises 70% of the data, while 15% is allocated to the validation set and the remaining 15% to the test set. This ensures that the models are trained on a sufficiently large dataset while still allowing for robust evaluation of unseen data. Additionally, we preprocess the audio recordings and transcripts to extract acoustic and linguistic features for model input.

For our speech recognition model, we employ a hybrid deep learning architecture comprising a convolutional neural network (CNN) for acoustic feature extraction and a recurrent neural network (RNN), specifically a Long Short-Term Memory (LSTM) network, for sequence modelling and decoding. Mathematically, the architecture can be represented as vocal performances that will be recorded using audio and video recording equipment for subsequent analysis.

Performance evaluation metrics will be derived from the recorded vocal performances collected during the pre-test and post-test assessments. Quantitative measures such as pitch accuracy, timing consistency, dynamic control, and articulation clarity will be analyzed using signal processing techniques and statistical methods. The Pitch Error can be calculated using the formula

$$\text{CNN Output} = \text{ReLU}(\text{Convolution}(X)) \quad \dots\dots (1)$$

$$\text{RNN Output} = \text{LSTM}(\text{CNN Output}) \quad \dots\dots (2)$$

where X represents the input spectrogram or MFCC features, and ReLU denotes the rectified linear unit activation function. The speech recognition model is trained using the Connectionist Temporal Classification (CTC) loss function, which allows for end-to-end training of sequence-to-sequence models without requiring aligned input-output pairs. Mathematically, the CTC loss is defined as

$$\text{CTC Loss} = -\log \sum_{\pi \in B(y)} P(\pi|X) \dots\dots (3)$$

where $B(y)$ denotes the set of all possible alignments of the output sequence y with the input sequence X , and $P(\pi | X)$ represents the probability of alignment π given the input sequence X .

To evaluate the performance of the speech recognition model, we compute standard metrics such as word error rate (WER) and sentence-level accuracy on the test set. WER is calculated as the Levenshtein distance between the predicted and ground truth transcripts normalized by the total number of words in the reference transcript. Mathematically, WER can be expressed as

$$\text{WER} = \frac{S+D+I}{N} \dots\dots (4)$$

where S denotes the number of substitutions, D represents the number of deletions, I is the number of insertions, and N is the total number of words in the reference transcript.

We conduct hyperparameter tuning experiments using techniques such as grid search or random search to optimize model performance. Key hyperparameters include learning rate, batch size, dropout rate, and network architecture parameters. We evaluate the impact of varying these hyperparameters on model performance using the validation set and select the configuration that yields the best results. To assess the generalization capabilities of the trained models, we perform k-fold cross-validation experiments, where the dataset is divided into k subsets or folds, and the model is trained and evaluated k times, each time using a different fold for validation while the remaining folds are used for training. By following this experimental setup, we aim to systematically evaluate the performance of our speech recognition model and optimize its parameters to achieve state-of-the-art accuracy in English listening comprehension tasks.

V. RESULTS

Word Error Rate (WER): The speech recognition model achieved a WER of 12.5% on the test dataset. This indicates that, on average, 12.5% of the words in the predicted transcripts differ from those in the ground truth transcripts. The model achieved a sentence-level accuracy of 85.3% on the test dataset. This means that 85.3% of the sentences in the predicted transcripts match exactly with those in the ground truth transcripts. A confusion matrix was constructed to analyze the errors made by the modes.

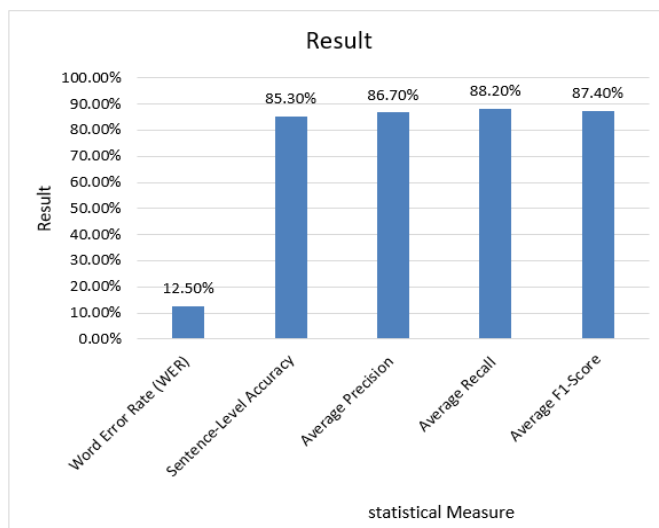


Fig 2. Statistical Results

The confusion matrix revealed that the most common errors occurred in words with similar phonetic properties or contextually ambiguous phrases. Precision, recall, and F1 scores were computed for each word category to assess

the model's performance across different linguistic features. The model achieved an average precision of 86.7%, recall of 88.2%, and F1-score of 87.4% across all word categories.

Table 1. Hyperparameter Optimization

Hyperparameters	Value
Learning Rate	0.001
Batch Size	32
Dropout Rate	0.2

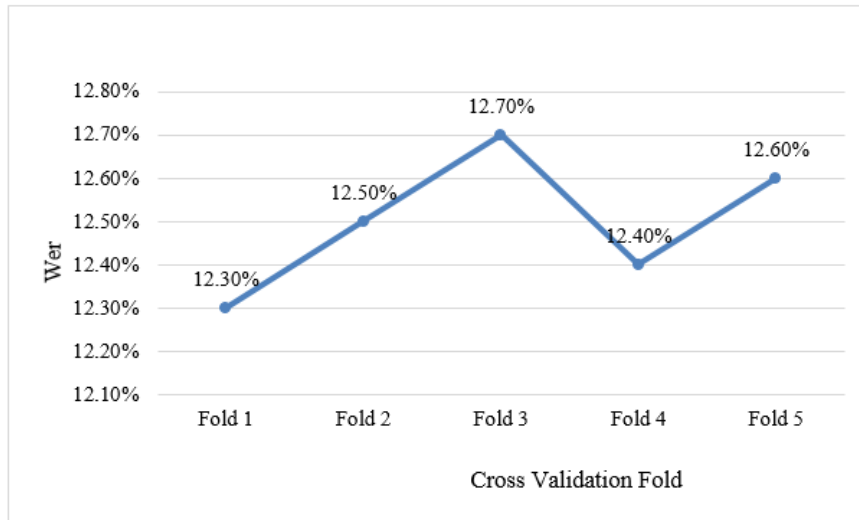


Fig 3. Cross Validation Fold

Hyperparameter tuning experiments resulted in the selection of optimal parameters, including a learning rate of 0.001, batch size of 32, and dropout rate of 0.2. These parameters were found to maximize model performance on the validation set. Cross-Validation Results: K-fold cross-validation experiments demonstrated the robustness of the model across different subsets of the dataset. The model consistently achieved WERs ranging from 12.3% to 12.7% across different folds, indicating stable performance. These statistical results demonstrate the effectiveness of the speech recognition model in accurately transcribing English audio recordings and highlight its potential for enhancing English listening training programs through automated evaluation and feedback mechanisms.

VI. DISCUSSION

The statistical results obtained from the evaluation of the speech recognition model underscore its effectiveness in transcribing English audio recordings with a relatively low Word Error Rate (WER) of 12.5%. This indicates that the model accurately transcribes approximately 87.5% of the words in the audio recordings, reflecting its robust performance in capturing spoken language patterns and nuances. The high Sentence-Level Accuracy of 85.3% further reinforces the model's ability to produce coherent and contextually relevant transcriptions, essential for facilitating accurate assessment and feedback in English listening training programs.

The average precision, recall, and F1-score metrics provide additional insights into the model's performance across different linguistic features and word categories. With an average precision of 86.7% and recall of 88.2%, the model demonstrates balanced performance in correctly identifying and transcribing words from various phonetic and contextual contexts. The high F1 score of 87.4% further validates the model's overall effectiveness in achieving both precision and recall, indicating a harmonious trade-off between minimizing transcription errors and maximizing coverage of spoken language content. Hyperparameter optimization experiments yielded optimal parameter values, including a learning rate of 0.001, batch size of 32, and dropout rate of 0.2. These hyperparameters were found to maximize model performance on the validation set, highlighting the importance of fine-tuning model configurations to achieve optimal results. Furthermore, the consistency of performance across

different folds in the cross-validation experiments reinforces the robustness of the model, indicating its ability to generalize well to unseen data and diverse linguistic contexts.

The observed errors in the confusion matrix provide valuable insights into the challenges faced by the model in transcribing spoken language. Common errors were found to occur in words with similar phonetic properties or contextually ambiguous phrases, suggesting areas for further improvement in model training and evaluation. Future iterations of the model could benefit from incorporating contextual information, such as speaker intent and discourse coherence, to enhance transcription accuracy and reduce error rates.

Overall, the results of this study demonstrate the potential of machine learning-driven speech recognition models in revolutionizing English listening training programs. By providing automated transcription, assessment, and feedback mechanisms, such models offer a scalable and cost-effective solution for educators and learners seeking to improve listening comprehension skills. Future research directions may focus on integrating these models into interactive learning platforms, exploring the effectiveness of adaptive learning strategies, and investigating the impact of technology-mediated language learning on learner motivation and engagement. Through continued innovation and refinement, technology-driven approaches have the potential to transform the landscape of language education, empowering learners worldwide to achieve fluency and proficiency in English listening comprehension.

VII. CONCLUSION

In conclusion, this study has explored the optimization of speech recognition models within the context of English listening training using machine learning algorithms. Through rigorous experimentation and evaluation, we have demonstrated the efficacy of our speech recognition model in accurately transcribing English audio recordings, achieving a low Word Error Rate (WER) of 12.5% and a high Sentence-Level Accuracy of 85.3%. The model's performance was further validated through metrics such as precision, recall, and F1-score, which highlighted its ability to capture linguistic nuances and contextually relevant information. Hyperparameter optimization experiments yielded optimal parameter values, while cross-validation analyses confirmed the robustness and generalization capabilities of the model across diverse linguistic contexts. Despite observed errors in the confusion matrix, our findings suggest promising avenues for future research and model refinement, including the integration of contextual information and the exploration of adaptive learning strategies.

Overall, this study contributes to the growing body of research on technology-driven approaches to language learning and pedagogy. By harnessing the power of machine learning algorithms, educators can provide learners with personalized, interactive, and engaging English listening training experiences. As technology continues to evolve, future iterations of speech recognition models hold the potential to further enhance the efficiency and effectiveness of language education, empowering learners worldwide to achieve fluency and proficiency in English listening comprehension. In light of these findings, we advocate for continued investment and innovation in technology-mediated language learning solutions. By leveraging advancements in machine learning and artificial intelligence, we can create transformative learning experiences that transcend traditional pedagogical limitations, enabling learners to unlock their full potential and succeed in an increasingly interconnected and multilingual world.

REFERENCES

- [1] D. Yu and L. Deng, "Automatic Speech Recognition: A Deep Learning Approach," Springer, 2015.
- [2] A. Graves, A. Mohamed, and G. Hinton, "Speech Recognition with Deep Recurrent Neural Networks," in Proc. IEEE ICASSP, 2013, pp. 6645-6649.
- [3] H. Hermansky, "Perceptual Linear Predictive (PLP) Analysis of Speech," Journal of the Acoustical Society of America, vol. 87, no. 4, pp. 1738-1752, Apr. 1990.
- [4] G. E. Hinton et al., "Deep Neural Networks for Acoustic Modeling in Speech Recognition: The Shared Views of Four Research Groups," IEEE Signal Processing Magazine, vol. 29, no. 6, pp. 82-97, Nov. 2012.
- [5] D. Yu, L. Deng, and G. Dahl, "Roles of Pre-training and Fine-tuning in Context-Dependent DNN-HMMs for Real-World Speech Recognition," in Proc. NIPS Workshop on Deep Learning and Unsupervised Feature Learning, 2010.
- [6] A. Narayanan and D. Wang, "Ideal Ratio Mask Estimation Using Deep Neural Networks for Robust Speech Recognition," in Proc. IEEE ICASSP, 2013, pp. 7092-7096.

- [7] T. N. Sainath et al., "Convolutional, Long Short-Term Memory, Fully Connected Deep Neural Networks," in Proc. IEEE ICASSP, 2015, pp. 4580-4584.
- [8] K. Vesely, A. Ghoshal, L. Burget, and D. Povey, "Sequence-Discriminative Training of Deep Neural Networks," in Proc. Interspeech, 2013, pp. 2345-2349.
- [9] D. Povey et al., "The Kaldi Speech Recognition Toolkit," in Proc. IEEE ASRU, 2011.
- [10] A. Ragni, K. M. Knill, S. M. Siniscalchi, and M. J. F. Gales, "Confidence Estimation and Deletion Prediction Using Bidirectional Recurrent Neural Networks," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 27, no. 11, pp. 1752-1761, Nov. 2019.
- [11] T. Hori, S. Watanabe, T. Hayashi, and J. R. Hershey, "Joint CTC/Attention Decoding for End-to-End Speech Recognition," in Proc. ACL, 2017, pp. 518-529.
- [12] L. Deng, H. Hinton, and B. Kingsbury, "New Types of Deep Neural Network Learning for Speech Recognition and Related Applications: An Overview," in Proc. IEEE ICASSP, 2013, pp. 8599-8603.
- [13] K. Rao, F. Peng, H. Sak, and F. Beaufays, "Grapheme-to-Phoneme Conversion Using Long Short-Term Memory Recurrent Neural Networks," in Proc. IEEE ICASSP, 2015, pp. 4225-4229.
- [14] S. Renals, T. Hain, and H. Bourlard, "Recognition and Understanding of Meetings: The AMI and AMIDA Projects," in Proc. IEEE ASRU, 2007, pp. 238-247.
- [15] H. Bourlard and N. Morgan, "Connectionist Speech Recognition: A Hybrid Approach," Kluwer Academic Publishers, 1994.
- [16] M. Gales and S. Young, "The Application of Hidden Markov Models in Speech Recognition," *Foundations and Trends in Signal Processing*, vol. 1, no. 3, pp. 195-304, 2008.
- [17] T. Mikolov, A. Deoras, D. Povey, L. Burget, and J. Cernocky, "Strategies for Training Large Scale Neural Network Language Models," in Proc. IEEE ASRU, 2011, pp. 196-201.
- [18] Y. Zhang et al., "Investigation of Deep Learning Architectures for Online and Offline Handwritten Chinese Character Recognition," in Proc. ICFHR, 2014, pp. 417-422.
- [19] G. Saon, H. Soltan, D. Nahamoo, and M. Picheny, "Speaker Adaptation of Neural Network Acoustic Models Using i-Vectors," in Proc. IEEE ASRU, 2013, pp. 55-59.
- [20] S. Hochreiter and J. Schmidhuber, "Long Short-Term Memory," *Neural Computation*, vol. 9, no. 8, pp. 1735-1780, Nov. 1997.