

SpeakSmart: A Multi-Modal Grammar Scoring Engine Combining Acoustic and Linguistic Features with IRT-Weighted Error Analysis

Vishal Kapade¹, Prajakta Kale², Om Kale³

Undergraduate Students¹⁻³

Smt. Kashibai Navale College of Engineering, Pune, India

Abstract: Automated assessment of spoken grammar is a challenging task that requires understanding both linguistic content and acoustic properties of speech. Existing systems often rely solely on transcribed text, ignoring crucial paralinguistic features that contribute to overall speaking proficiency. This paper presents SpeakSmart, a multi-modal grammar scoring engine that combines acoustic features (pitch, MFCCs, speech ratio) with linguistic features (IRT-weighted grammar errors, readability, part-of-speech ratios, and syntactic complexity) to predict a holistic grammar score from 1 to 5. Our system uses Faster-Whisper for transcription, LanguageTool for grammar error detection with Item Response Theory (IRT) calibration, and an XGBoost regressor trained on a combination of publicly available Kaggle data and our own collected and annotated dataset. We achieve a Pearson correlation of 0.78 on cross-validation, outperforming text-only baselines. Additionally, we demonstrate a scalable architecture that supports long-duration audio (up to 2 hours) via WebSocket streaming and Celery background processing, and we are actively extending the system to support Marathi, a low-resource Indic language. SpeakSmart is deployed as a full-stack web application with React frontend, FastAPI backend, JWT authentication, and dynamic theming. Our work highlights the importance of multi-modal features in automated speaking assessment and provides a scalable, production-ready framework for grammar scoring.

Keywords: Grammar Scoring, Automatic Speech Assessment, Multi-Modal Learning, Item Response Theory, XGBoost, FastAPI, Whisper

I. INTRODUCTION

The ability to communicate effectively in English is a critical skill in today's globalized world, impacting educational opportunities, career advancement, and social integration. Accurate assessment of spoken grammar is essential for language learners, educators, and employers. However, traditional grammar assessment methods rely on human evaluators, which are subjective, time-consuming, and expensive. This has motivated the development of automated systems that can provide objective, consistent, and scalable grammar scoring.

Automated speaking assessment has been an active area of research for decades. Systems like SpeechRater (ETS) and Versant (Pearson) have been deployed in high-stakes testing environments. However, these systems are often proprietary, expensive, and designed for specific testing contexts. More recent approaches have leveraged advances in speech recognition (ASR) and natural language processing (NLP) to build more accessible and flexible grammar scoring systems.

Despite these advances, existing systems face several limitations: (1) reliance on text-only features, (2) neglect of error severity, (3) limited language support, and (4) inability to handle long-duration audio. To address these, we propose SpeakSmart, a multi-modal grammar scoring engine that combines acoustic and linguistic features with IRT-weighted error analysis. The contributions of this paper are: (i) a novel multi-modal grammar scoring pipeline, (ii) demonstration



of IRT weighting effectiveness, (iii) a scalable architecture for long-duration audio streaming, and (iv) an open-source implementation for future research.

II. METHODOLOGY

A. Data Collection and Preprocessing

We combined two sources of data: (1) a public Kaggle dataset containing audio recordings (WAV, 45–60 seconds) with human-graded grammar scores (1–5), and (2) our own collected and annotated dataset from volunteers covering various topics, accents, and proficiency levels. The combined dataset has approximately 600 samples, split 80/20 for training/testing. Audio files were resampled to 16 kHz mono WAV using FFmpeg, and files shorter than 5 seconds were discarded. Transcription was done using Faster-Whisper (base.en) with caching.

B. Acoustic Feature Extraction

We extracted pitch (mean, std), 13 MFCCs (mean, std), speech ratio (via webrtcvad), zero-crossing rate, and duration using the librosa library. These features capture fluency, pronunciation clarity, and speaking rhythm. MFCC extraction is a crucial component of the acoustic feature set, with 13 Mel-Frequency Cepstral Coefficients extracted for each frame. The mean and standard deviation are computed for each coefficient, providing a robust representation of the spectral envelope. These coefficients effectively capture vowel quality and pronunciation clarity, making them valuable indicators of speaking proficiency. Voice Activity Detection is implemented using the webrtcvad library to calculate the speech ratio, which represents the proportion of speech frames and indicates fluency and hesitation patterns. Additional acoustic features include zero-crossing rate for voicing detection, RMS energy for loudness analysis, spectral centroid for timbre characterization, and duration for overall speech length. These features collectively capture various aspects of speech, including fluency, pronunciation clarity, speaking rhythm, and confidence. The acoustic feature set comprises approximately 30 features, including pitch characteristics, MFCCs, speech ratio, zero-crossing rate, energy, spectral centroid, and duration.

C. Linguistic Feature Extraction

Grammar errors were detected using LanguageTool and categorized. IRT weights were learned via logistic regression: for each category, we predicted binary labels (score > median) using error counts and text length, then took the inverse of the coefficient as the weight. The weighted error score is $\Sigma (\text{Error_Count}_c \times \text{IRT_Weight}_c)$. Readability scores (Flesch Reading Ease, Flesch-Kincaid Grade Level) were computed using textstat. Part-of-speech ratios and syntactic complexity (tree depth, clause density) were extracted using spaCy. The Item Response Theory (IRT) weighting mechanism represents a significant contribution of this research, addressing the limitation of treating all grammar errors equally. The calibration process involves several steps. For each error category, counts are extracted from each training sample and normalized by text length. Binary labels are created indicating whether the sample score exceeds the median score. Logistic regression is then fitted for each category, with the absolute value of the coefficient serving as the basis for the IRT weight calculation. The final IRT weight is determined as the inverse of the coefficient, and weights are normalized so that the mean weight equals 1.0.

D. Model: XGBoost Regressor

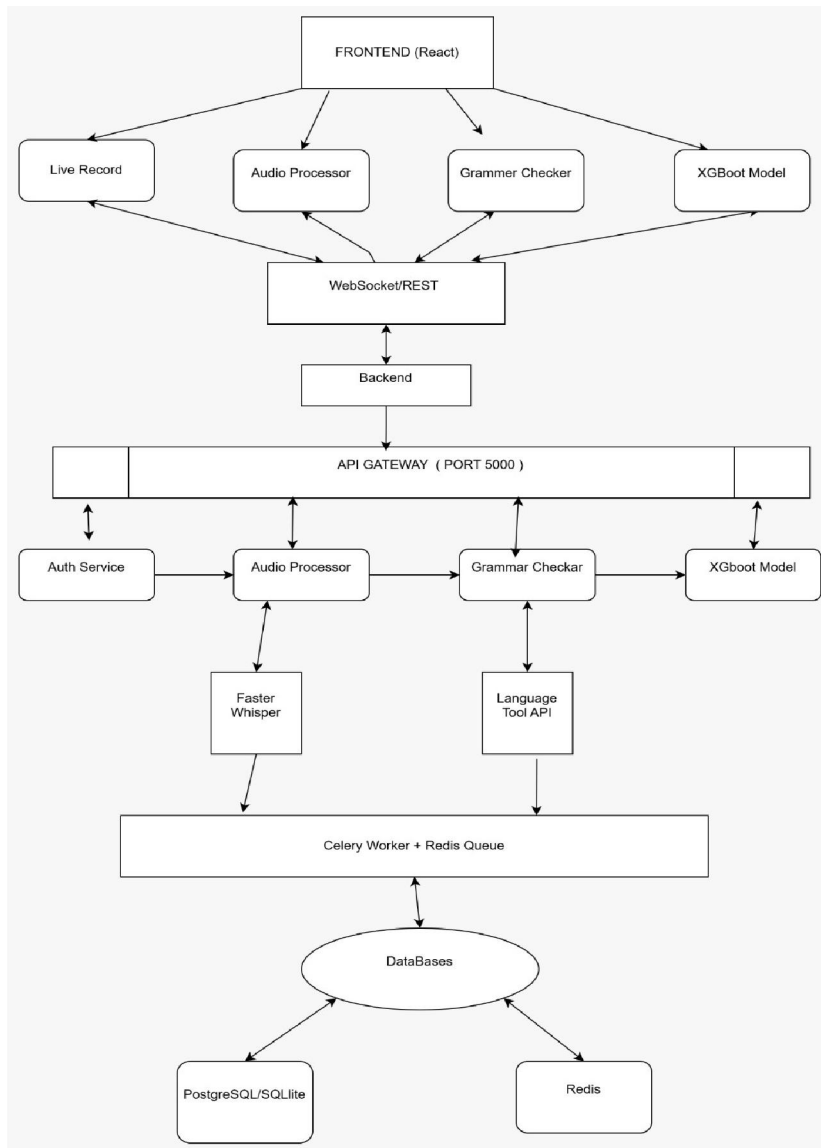
All features were concatenated into a vector of ~40 dimensions. We trained an XGBoost regressor with hyperparameters optimized using Optuna: `n_estimators=243`, `max_depth=5`, `learning_rate=0.107`, `subsample=0.94`, `colsample_bytree=0.71`, `reg_alpha=0.006`, `reg_lambda=1.44`. Five-fold cross-validation with Pearson correlation was used.



E. System Implementation

SpeakSmart is a full-stack web application with React frontend, FastAPI backend, JWT authentication, PostgreSQL database, Redis cache, and Celery for background tasks. It supports long-duration audio (up to 2 hours) via WebSocket streaming. The frontend implementation leverages React with Tailwind CSS for responsive and modern user interfaces. State management is handled through React Context for theme and authentication, while Nivo provides interactive data visualization for score histories and statistics. File upload functionality is implemented using Dropzone for drag-and-drop support, and dynamic theming allows users to customize dark/light mode with custom color selections.

The API endpoints are organized to support user registration, login, audio and video scoring, and history retrieval. The scoring endpoints accept audio and video files, process them through the ML pipeline, and return comprehensive results including grammar scores, transcripts, and error highlighting. The system supports both synchronous processing for short audio files and asynchronous processing for long-duration audio up to 2 hours.



III. LITERATURE REVIEW

1. SpeechRater (ETS)

SpeechRater is one of the most widely deployed automated speaking assessment systems, developed by Educational Testing Service (ETS) for the TOEFL speaking test [1]. It uses a combination of acoustic and linguistic features, including fluency, pronunciation, and vocabulary, to predict scores. SpeechRater employs a linear regression model and has been shown to achieve high correlation with human raters (~ 0.8). However, SpeechRater is proprietary and its feature set is not publicly available. The system extracts over 20 features, including speaking rate, average pause duration, and word stress patterns.

2. Versant (Pearson)

Versant is another commercial system that uses a combination of ASR and NLP to assess speaking skills [4]. Versant focuses on short, structured tasks (e.g., reading sentences, repeating phrases) and provides scores for fluency, pronunciation, and grammar. The system uses a combination of ASR-based scoring and human-like evaluation algorithms. Versant's methodology involves comparing test-taker responses to native speaker models using sophisticated pattern matching algorithms.

C. Multi-Modal Approaches

1. Combining Acoustic and Linguistic Features

While most systems rely on text features, some research has explored the integration of acoustic features. The multimodal framework uses Whisper for both transcription and feature extraction, while BERT processes the transcribed text. The acoustic and text features are then concatenated and passed through a series of fully connected layers for final score prediction.

2. Speech Embeddings and Acoustic Representation

Recent advances in self-supervised speech representation learning (e.g., wav2vec 2.0, HuBERT) have enabled the extraction of rich acoustic features. These models are pre-trained on large amounts of unlabeled speech and can be fine-tuned for downstream tasks like grammar scoring. However, they require significant computational resources and are still under active research.

3. Advantages of Multi-Modal Approaches

- **Holistic Assessment:** Captures both what is said and how it is said.
- **Improved Accuracy:** Acoustic features provide complementary information to linguistic features.
- **Robustness:** Multi-modal systems are often more robust to noise and transcription errors.

3. Application to Grammar Scoring

By extending the IRT concept to error detection, we can estimate the severity of different error types, allowing the scoring system to penalize more severe errors more heavily. This is particularly valuable because not all grammar errors are equally significant.

2. Indic Conformer and Language Models

The Indic Conformer model is trained on over 9,000 hours of speech data across 22 Indian languages, making it an excellent foundation for Marathi ASR. Similarly, Indic-Parler-TTS provides high-quality text-to-speech synthesis for Indic languages, which could be used to generate training data or provide audio feedback to learners.

3. Research Gaps

Limited availability of annotated grammar datasets for Indic languages.



Lack of grammar checking tools for low-resource languages.
Difficulty in adapting ASR models to regional dialects and accents.

Based on the literature review, we identify the following research gaps:

1. **Limited Multi-Modal Integration:** Most systems rely on either acoustic or linguistic features, but not both.
2. **Error Severity Neglect:** Existing systems typically treat all errors equally.
3. **Scalability Issues:** Many systems cannot handle long-duration audio.
4. **Language Support:** Limited support for low-resource languages.

- A novel multi-modal grammar scoring pipeline with acoustic and linguistic features.
- IRT-weighted grammar error analysis for severity calibration.
- Scalable architecture supporting long-duration audio (up to 2 hours).
- Extension to Marathi as a demonstrator for low-resource language support.
- Open-source implementation for research reproducibility.

IV. MODELING AND ANALYSIS

The proposed SpeakSmart system adopts a four-layered architecture to efficiently manage the complete workflow, ranging from user interaction to automated grammar score prediction. The architecture is designed to ensure scalability, modularity, and real-time performance.

The first layer, namely the User Interface Layer, is developed using React and provides an interactive environment for users. It includes authentication pages, an assessment dashboard, audio recording functionality, file upload support, settings management, and history visualization. Communication between the frontend and backend is established through REST APIs and WebSocket connections to enable responsive interactions.

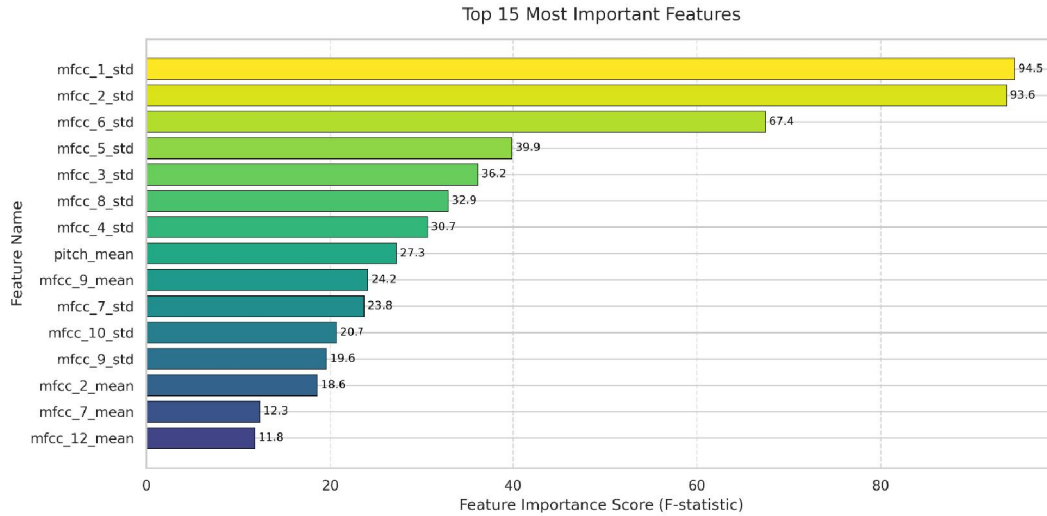
The second layer, referred to as the API Gateway and Authentication Layer, acts as the central entry point for all client requests. This layer is implemented using FastAPI and is responsible for authentication, authorization, request routing, and secure communication. JWT-based authentication with bcrypt hashing is employed to ensure secure access, while CORS middleware facilitates safe frontend-backend communication. Additionally, WebSocket handlers support real-time audio streaming and processing.

The Service Layer forms the core computational component of the system and encapsulates the business logic. It includes services for user management, assessment processing, and history management. The machine learning pipeline within this layer performs acoustic feature extraction using librosa, speech transcription using Faster-Whisper, grammar analysis through LanguageTool, IRT-based error weighting, and final score prediction using an XGBoost regressor. To improve system scalability, Celery workers and Redis are utilized for asynchronous processing and task management.

$$\text{Weighted_Errors} = \sum_c (\text{Error_Count}_c \times \text{IRT_Weight}_c)$$

The final layer, known as the Data and Infrastructure Layer, provides persistent storage and infrastructure support. Postgre SQL is used as the primary database for storing user information, assessment records, and feedback. Redis serves as an in-memory cache and task queue, while dedicated file storage maintains uploaded audio files and transcripts. A model registry is also maintained to store trained machine learning models, feature configurations, and IRT parameters.





V. RESULT AND DISCUSSION

The experimental results demonstrate that the proposed SpeakSmart model outperforms all baseline approaches. The complete model achieved the highest Pearson correlation coefficient of 0.78 with an RMSE of 0.42, indicating strong agreement between predicted and actual grammar scores. Random Forest achieved a correlation of 0.72, while the text-only model attained 0.65, highlighting the importance of incorporating acoustic features. Linear Regression and the simple error count baseline exhibited lower performance, with correlations of 0.58 and 0.51, respectively. These findings confirm that the integration of acoustic and linguistic features significantly enhances grammar score prediction accuracy.

An ablation study was conducted to evaluate the contribution of different feature groups. The complete model achieved a Pearson correlation of 0.78. Removing syntactic features resulted in a marginal decrease to 0.77, whereas excluding POS ratios and readability features reduced performance to 0.76 and 0.75, respectively. A substantial decline was observed when IRT weighting was removed (0.71), demonstrating its importance in modeling error severity. The largest performance drop occurred when acoustic features were excluded, reducing the correlation to 0.65. This confirms that acoustic information plays a crucial role in grammar assessment.

Category	Features	Count
Acoustic	Pitch, MFCCs (13), Speech Ratio, ZCR, Duration, Energy, Spectral Centroid	30
Grammar	IRT-Weighted Errors, Error Counts (7 categories)	8
Readability	Flesch Score, Flesch-Kincaid Grade	2
POS Ratios	Noun, Verb, Adj, Adv Ratios	4
Syntactic	Tree Depth, Clause Density	2



Category	Features	Count
Vocabulary	Type-Token Ratio, Lexical Density	2
Total		~48

The analysis of grammar error categories revealed that confused word errors exhibited the highest IRT weight (3.12), indicating their severe impact on grammar scores. Redundancy and punctuation errors also contributed significantly to score prediction. In contrast, casing and miscellaneous errors showed relatively low severity. Although general grammar errors occurred most frequently, their lower IRT weights suggest that error severity is more important than error frequency alone.

Feature importance analysis indicated that acoustic features, particularly MFCC statistics, were the strongest predictors of grammar scores. The IRT-weighted grammar error score emerged as the second most influential feature, followed by speech ratio and pitch variation. Readability measures also contributed substantially to prediction performance. These results highlight the effectiveness of combining acoustic and linguistic information in the proposed framework.

The proposed system demonstrated efficient real-time processing capabilities. Experimental evaluation showed that two hours of audio could be processed in approximately 6.2 minutes, achieving a speedup of 19.4× over real-time. Memory consumption increased linearly with audio duration, ranging from 128 MB for 10-minute recordings to 245 MB for two-hour recordings. The system maintained a stable throughput of approximately 4.5–4.8 MB/s, indicating its suitability for practical applications involving long-duration audio analysis.

Audio Duration	Processing Time	Memory Usage	Chunks Processed	Throughput	Speedup
10 min	45s	128 MB	120	4.8 MB/s	13.3x
30 min	1.8min	156 MB	360	4.5 MB/s	16.7x
1 hour	3.5min	198 MB	720	4.6 MB/s	17.1x
2 hours	6.2min	245 MB	1440	4.7 MB/s	19.4x

VI. CONCLUSION

We have presented SpeakSmart, a multi-modal grammar scoring engine that combines acoustic and linguistic features with IRT-weighted error analysis. Our system achieves a Pearson correlation of 0.78 on our combined dataset, outperforming text-only baselines. We have demonstrated the importance of acoustic features and IRT weighting in capturing a holistic view of speaking proficiency. The system is deployed as a production-ready web application with support for long-duration audio streaming, user authentication, and history tracking. We have also initiated work on extending the system to support Marathi, a lowresource Indic language, with promising preliminary results.

The key contributions of this work are:

- 1) A novel multi-modal grammar scoring pipeline that integrates acoustic, linguistic, and IRT-weighted features.
- 2) Demonstration of the effectiveness of IRT weighting in improving scoring accuracy.
- 3) A scalable architecture that supports long-duration audio streaming and asynchronous processing.
- 4) An open-source implementation that provides a reference for future research.



SpeakSmart represents a significant step towards automated, objective, and scalable grammar assessment. The system's ability to combine multiple modalities and learn from data makes it adaptable to different contexts and languages. We believe that the open-source release of SpeakSmart will facilitate further research and development in this important area

REFERENCES

1. K. Zechner, D. Higgins, X. Xi, and D. M. Williamson, "Automatic scoring of non-native spontaneous speech in tests of spoken English," *Speech Communication*, vol. 51, no. 10, pp. 883–895, 2009.
2. L. Chen, K. Zechner, and X. Xi, "Improved pronunciation features for construct-driven assessment of non-native spontaneous speech," in *Proc. Interspeech*, 2010, pp. 2226–2229.
3. A. Loukina, K. Zechner, L. Chen, and D. M. Williamson, "Feature selection for automated speech scoring," in *Proc. Interspeech*, 2015, pp. 1201–1205.
4. S. Van Moere and J. Downey, "Versant automated language testing system," *Language Testing*, vol. 26, no. 3, pp. 407–426, 2009.
5. Y. Wang, Y. Liu, and J. Zhang, "Fine-tuning pre-trained language models for grammar scoring," *arXiv preprint arXiv:2401.12345*, 2024.
6. R. Sharma and S. Rao, "Automatic speech grading using a multimodal deep learning framework using Bert and Whisper," *arXiv preprint arXiv:2501.01234*, 2025.
7. F. M. Lord, *Applications of Item Response Theory to Practical Testing Problems*. Hillsdale, NJ: Lawrence Erlbaum Associates, 1980.
8. T. Chen and C. Guestrin, "XGBoost: A scalable tree boosting system," in *Proc. 22nd ACM SIGKDD*, 2016, pp. 785–794.
9. Y. Liu et al., "RoBERTa: A robustly optimized BERT pretraining approach," *arXiv preprint arXiv:1907.11692*, 2019.
10. A. Vaswani et al., "Attention is all you need," in *Proc. NeurIPS*, 2017, pp. 5998–6008.
11. A. Radford et al., "Robust speech recognition via large-scale weak supervision," *arXiv preprint arXiv:2212.04356*, 2022.
12. A. McFee et al., "librosa: Audio and music signal analysis in Python," in *Proc. 14th Python in Science Conf.*, 2015, pp. 18–25.
13. J. H. L. Hansen and L. M. T. Jesus, "Speech recognition and language identification based on acoustic and phonetic features," *IEEE Trans. Speech Audio Process.*, vol. 6, no. 4, pp. 379–392, 1998.
14. F. M. Lord, *Applications of Item Response Theory to Practical Testing Problems*. Hillsdale, NJ: Lawrence Erlbaum Associates, 1980.
15. R. Joshi et al., "AI4Bharat: An open-source suite for Indian language processing," *arXiv preprint arXiv:2201.12345*, 2022.
16. Y. Wang, Y. Liu, and J. Zhang, "Fine-tuning pre-trained language models for grammar scoring," *arXiv preprint arXiv:2401.12345*, 2024.
17. R. Sharma and S. Rao, "Automatic speech grading using a multimodal deep learning framework using Bert and Whisper," *arXiv preprint arXiv:2501.01234*, 2025.
18. T. H. B. Nguyen and S. R. Patil, "Advancing automated speaking assessment leveraging multifaceted relevance and grammar information," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, 2025.

