# Early Autism Detection Through Machine Learning Analysis of Vocal Patterns: A Comparative Study of Neural Network Approaches

Wei Chen Tsinghua University

Kenji Tanaka University of Tokyo Maria Rodriguez University of Barcelona

Fatima Al-Mansoori King Saud University

#### Abstract

This research investigates the efficacy of machine learning algorithms in detecting autism spectrum disorder (ASD) through vocal pattern analysis. We developed and compared three neural network architectures convolutional neural networks (CNNs), recurrent neural networks (RNNs), and hybrid CNN-RNN models—for classifying vocal samples from 450 children aged 2-6 years. The dataset comprised 15,000 audio samples collected from clinical settings across three countries. Our methodology involved preprocessing vocal data through Mel-frequency cepstral coefficients (MFCCs) extraction and implementing a novel feature selection algorithm based on information gain. Results demonstrated that the hybrid CNN-RNN model achieved superior performance with 92.3% accuracy, 89.7% sensitivity, and 94.1% specificity in ASD detection. The mathematical formulation of our feature optimization process revealed significant improvements in classification efficiency. These findings suggest that automated vocal analysis systems could serve as valuable screening tools for early ASD identification, potentially reducing diagnostic delays and improving intervention outcomes.

**Keywords:** autism spectrum disorder, machine learning, vocal analysis, neural networks, early detection, diagnostic screening

### Introduction

Autism Spectrum Disorder (ASD) represents a complex neurodevelopmental condition characterized by challenges in social communication, restricted interests, and repetitive behaviors. Early identification of ASD is crucial for initiating

timely interventions that can significantly improve long-term outcomes. However, current diagnostic procedures often rely on behavioral observations and standardized assessments, which can be subjective and time-consuming. The average age of ASD diagnosis remains around 4-5 years, despite evidence suggesting that reliable detection is possible as early as 18-24 months.

Recent advances in machine learning and artificial intelligence have opened new avenues for developing objective, automated screening tools. Vocal patterns have emerged as a promising biomarker for ASD detection, as children with autism often exhibit distinctive prosodic features, speech timing abnormalities, and atypical vocal quality. Previous research has documented differences in pitch variation, speech rate, and vocal resonance between children with ASD and typically developing peers.

This study addresses the critical need for early, objective ASD screening methods by developing and comparing multiple machine learning approaches for vocal pattern analysis. We hypothesize that neural network architectures can effectively distinguish ASD-related vocal characteristics with high accuracy, potentially providing a scalable screening solution for clinical and educational settings.

### Literature Review

The application of computational methods to autism research has evolved significantly over the past decade. Early work by Oller et al. (2010) demonstrated that automated analysis of infant vocalizations could identify developmental differences associated with ASD. Their research utilized linear discriminant analysis and achieved classification accuracies of approximately 70%.

Subsequent studies explored more sophisticated machine learning techniques. Bone et al. (2014) applied support vector machines to vocal features extracted from home recordings, reporting 85% accuracy in classifying toddlers with ASD. Their work highlighted the importance of feature selection, particularly emphasizing prosodic and spectral characteristics.

Recent research has increasingly turned to deep learning approaches. Li et al. (2018) implemented convolutional neural networks for ASD detection from speech samples, achieving 88% accuracy. Their model processed raw audio waveforms, eliminating the need for manual feature engineering. However, their dataset was limited to school-aged children, leaving open questions about applicability to younger populations.

Parallel developments in natural language processing have informed ASD research. Recurrent neural networks, particularly long short-term memory (LSTM) architectures, have shown promise in capturing temporal dependencies in speech patterns. The work of Asgari et al. (2017) demonstrated that LSTMs could effectively model conversational turn-taking patterns, a key area of

difficulty for individuals with ASD.

Despite these advances, several gaps remain in the literature. Most studies have focused on single algorithmic approaches rather than comparative analyses. Additionally, few investigations have explored hybrid architectures that combine the spatial feature extraction capabilities of CNNs with the temporal modeling strengths of RNNs. Our research addresses these gaps through systematic comparison of multiple neural network architectures applied to a large, multi-national dataset of young children's vocalizations.

### **Research Questions**

This study addresses the following research questions:

- 1. How effectively can different neural network architectures (CNN, RNN, and hybrid CNN-RNN) classify vocal samples from children with ASD versus typically developing children?
- 2. Which acoustic features contribute most significantly to accurate ASD classification across different machine learning models?
- 3. To what extent do demographic factors (age, gender, linguistic background) influence the performance of vocal-based ASD detection systems?
- 4. How does the proposed hybrid CNN-RNN model compare to existing machine learning approaches in terms of sensitivity, specificity, and computational efficiency?

## Objectives

The primary objectives of this research are:

- 1. To develop and optimize three distinct neural network architectures (CNN, RNN, and hybrid CNN-RNN) for ASD detection through vocal pattern analysis.
- 2. To implement a comprehensive feature extraction pipeline that captures both spectral and temporal characteristics of children's vocalizations.
- 3. To validate the classification performance of each model using rigorous cross-validation techniques and statistical analysis.
- 4. To identify the most discriminative vocal features for ASD detection across different age groups and demographic profiles.
- 5. To establish performance benchmarks for automated vocal analysis systems in early ASD screening contexts.

### Hypotheses to be Tested

Based on existing literature and theoretical considerations, we formulated the following hypotheses:

H1: The hybrid CNN-RNN architecture will achieve significantly higher classification accuracy compared to standalone CNN or RNN models, due to its ability to capture both spatial and temporal patterns in vocal data.

H2: Spectral features related to formant structure and harmonic-to-noise ratio will demonstrate higher discriminative power for ASD classification compared to purely prosodic features.

H3: Model performance will vary significantly across different age groups, with highest accuracy achieved in the 3-4 year age range where vocal differences between ASD and typically developing children become most pronounced.

H4: The inclusion of contextual vocal features (such as turn-taking patterns and response latency) will improve classification sensitivity without compromising specificity.

### Approach/Methodology

### Participants and Data Collection

We recruited 450 children aged 2-6 years from clinical centers in China, Spain, Japan, and Saudi Arabia. The sample included 225 children with clinically confirmed ASD diagnoses (based on ADOS-2 and ADI-R assessments) and 225 typically developing children matched for age, gender, and linguistic background. All participants provided audio recordings during structured play-based assessments conducted by trained clinicians.

### **Data Preprocessing**

Audio recordings were processed using a standardized pipeline. We applied noise reduction algorithms and normalized amplitude levels across all samples. Vocal segments were automatically extracted using voice activity detection, with manual verification by speech-language pathologists. The final dataset comprised 15,000 vocal samples averaging 2-5 seconds in duration.

#### **Feature Extraction**

We extracted 128-dimensional feature vectors from each vocal sample, including:

- Mel-frequency cepstral coefficients (MFCCs) - Spectral centroid, rolloff, and flux - Pitch and formant frequencies - Jitter, shimmer, and harmonic-to-noise ratio - Pause duration and speech rate metrics

Feature selection was performed using information gain criteria, with the optimization function defined as:

$$IG(Feature) = \sum_{c \in \{ASD, TD\}} P(c|Feature) \log_2 \frac{P(c|Feature)}{P(c)} \tag{1}$$

where P(c|Feature) represents the conditional probability of class c given the feature value, and P(c) is the prior probability of class c.

#### **Model Architectures**

We implemented three neural network architectures:

- 1. **CNN Model**: 5 convolutional layers with max pooling, followed by 3 fully connected layers. Optimized for spatial feature extraction from spectrogram representations.
- 2. **RNN Model**: Bidirectional LSTM with 128 hidden units, designed to capture temporal dependencies in sequential vocal features.
- 3. **Hybrid CNN-RNN Model**: Combined convolutional layers for feature extraction with LSTM layers for temporal modeling, implementing the following computational flow:

$$\mathbf{h}_t = LSTM(CNN(\mathbf{x}_t), \mathbf{h}_{t-1}) \tag{2}$$

where  $\mathbf{x}_t$  represents the input feature vector at time t, and  $\mathbf{h}_t$  is the hidden state.

#### Training and Evaluation

All models were trained using 5-fold cross-validation with an 80-20 train-test split. We employed Adam optimization with categorical cross-entropy loss and implemented early stopping to prevent overfitting. Performance metrics included accuracy, precision, recall, F1-score, and area under the ROC curve.

### Results

Our comparative analysis revealed significant differences in performance across the three neural network architectures. The hybrid CNN-RNN model demonstrated superior classification capabilities across all evaluation metrics.

Table 1: Performance Comparison of Neural Network Architectures for ASD Detection

Model	Accuracy	Precision	Recall	F1-Score	AUC
CNN RNN Hybrid CNN-RNN	86.7% $88.9%$ $92.3%$	86.5%	82.9% 85.1% 89.7%	83.5% $85.8%$ $90.2%$	0.923

Feature importance analysis revealed that spectral features, particularly MFCC coefficients 2-5 and formant dispersion metrics, contributed most significantly to accurate classification. The hybrid model effectively leveraged both spectral and temporal patterns, demonstrating robust performance across different age groups and demographic subsets.

Age-stratified analysis showed peak performance in the 3-4 year age group (94.1% accuracy), with slightly reduced but still substantial accuracy in younger (2-3 years: 87.3%) and older (5-6 years: 90.8%) age groups. Gender-based analysis revealed comparable performance across male and female participants, addressing concerns about potential gender bias in automated screening tools.

Computational efficiency analysis indicated that while the hybrid model required approximately 40% more training time than standalone architectures, it achieved convergence with fewer epochs and demonstrated superior generalization to unseen data.

### Discussion

Our findings support the central hypothesis that hybrid neural network architectures can significantly enhance ASD detection through vocal pattern analysis. The superior performance of the CNN-RNN model aligns with theoretical expectations, as ASD-related vocal characteristics manifest through both spectral abnormalities (captured by CNNs) and temporal patterns (modeled by RNNs).

The 92.3% overall accuracy achieved by our hybrid model represents a substantial improvement over previous automated screening approaches. This performance level approaches the inter-rater reliability of experienced clinicians using standardized diagnostic instruments, suggesting potential clinical utility as a screening tool.

The feature importance results provide insights into the acoustic markers most strongly associated with ASD. The prominence of MFCC and formant features suggests that voice quality and resonance characteristics may be particularly discriminative. This finding aligns with clinical observations of atypical vocal quality in children with ASD, though the specific acoustic correlates have been poorly characterized until now.

The age-dependent performance patterns merit careful consideration. The peak accuracy in the 3-4 year age group may reflect the developmental period when vocal differences become most pronounced, before compensatory strategies emerge in older children. The maintained high accuracy in the 5-6 year group suggests that vocal markers persist beyond early childhood, supporting the potential utility of this approach across a broader age range.

Several limitations should be acknowledged. Our sample, while multi-national, was collected in clinical settings and may not fully represent the diversity of vocal patterns in community populations. Additionally, the cross-sectional design prevents analysis of how vocal characteristics evolve over time within individuals.

### Conclusions

This research demonstrates that machine learning analysis of vocal patterns can achieve high accuracy in ASD detection, with hybrid CNN-RNN architectures outperforming single-modality approaches. The developed system shows promise as an objective, scalable screening tool that could complement existing diagnostic procedures.

The identification of specific acoustic features associated with ASD provides new insights into the vocal characteristics of autism, potentially informing both screening technologies and theoretical models of communication differences in ASD.

Future research should focus on longitudinal studies tracking vocal development in children with ASD, investigation of cross-cultural and cross-linguistic generalizability, and integration of vocal analysis with other behavioral and physiological markers. The development of mobile applications implementing these algorithms could facilitate widespread screening in diverse settings.

# Acknowledgements

We extend our gratitude to the participating families and clinical centers in Beijing, Barcelona, Tokyo, and Riyadh. This research was supported by the International Autism Research Consortium and the Global Health Innovation Fund. We thank Dr. Elena Petrov for statistical consultation and the technical staff at all participating institutions for data collection support.

99 Oller, D. K., et al. (2010). Automated vocal analysis of naturalistic recordings from children with autism, language delay, and typical development. *Proceedings of the National Academy of Sciences*, 107(30), 13354-13359.

Bone, D., et al. (2014). Applying machine learning to facilitate autism diagnostics: pitfalls and promises. *Journal of Autism and Developmental Disorders*, 44(10), 2478-2495.

Li, B., et al. (2018). A deep learning approach for autism spectrum disorder detection from speech signals. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 26(11), 2153-2161.

Asgari, M., et al. (2017). Automatic detection of autism based on conversational turn-taking and vocal intonation. *Proceedings of the Interspeech Conference*, 202-206.

Cohen, I. L. (2019). Behavioral and electrophysiological markers of autism spectrum disorder. Journal of Developmental and Behavioral Pediatrics, 40(7), 561-573.

Schuller, B. W., et al. (2020). The INTERSPEECH 2020 computational paralinguistics challenge: Elderly emotion, breathing & masks. *Proceedings of the Interspeech Conference*, 2042-2046.