# Early Autism Detection Through Vocal Pattern Analysis: A Machine Learning Approach Using Infant Cry Acoustics

Wei Zhang Tsinghua University

Kenji Tanaka University of Tokyo Maria Rodriguez University of Barcelona

Fatima Al-Jaber King Saud University

#### Abstract

This study investigates the efficacy of machine learning algorithms in detecting early signs of autism spectrum disorder (ASD) through acoustic analysis of infant vocalizations. We collected and analyzed 2,500 audio samples from infants aged 6-18 months, including 1,200 samples from typically developing infants and 1,300 from infants later diagnosed with ASD. Using feature extraction techniques including Mel-frequency cepstral coefficients (MFCCs), pitch contours, and spectral properties, we trained and evaluated multiple classification models. Our results demonstrate that support vector machines achieved 87.3% accuracy in distinguishing ASD-related vocal patterns, with random forests and neural networks performing at 84.1% and 89.2% accuracy respectively. The findings suggest that automated vocal analysis could serve as a valuable screening tool for early ASD detection, potentially enabling earlier intervention and improved developmental outcomes.

**Keywords:** autism spectrum disorder, machine learning, vocal analysis, early detection, infant cries, acoustic features

#### Introduction

Autism Spectrum Disorder (ASD) represents a complex neurodevelopmental condition characterized by challenges in social communication, restricted interests, and repetitive behaviors. Early detection of ASD is crucial for initiating timely interventions that can significantly improve long-term outcomes. Current diagnostic methods typically rely on behavioral observations and parental reports, with formal diagnosis often occurring around 3-4 years of age. This diagnostic delay represents a critical window where early intervention could yield substantial benefits.

The acoustic properties of infant vocalizations have emerged as a promising biomarker for early neurodevelopmental assessment. Previous research has demonstrated that atypical vocal patterns can reflect underlying neurological differences associated with ASD. However, traditional acoustic analysis methods have been limited by subjective interpretation and manual feature extraction. The advent of machine learning techniques offers new opportunities for automated, objective analysis of vocal characteristics that may serve as early indicators of ASD.

This study addresses the critical need for earlier ASD detection by developing and validating machine learning models that analyze infant cry acoustics. Our approach leverages recent advances in signal processing and pattern recognition to identify subtle vocal markers that may precede overt behavioral symptoms. By focusing on the 6-18 month age range, we target the period when early intervention can have the most profound impact on developmental trajectories.

#### Literature Review

The relationship between vocal characteristics and neurodevelopmental disorders has been explored in several previous studies. Sheinkopf et al. (2000) first documented atypical cry patterns in infants later diagnosed with ASD, noting differences in fundamental frequency and temporal structure. Subsequent work by Oller et al. (2010) expanded on these findings, demonstrating that canonical babbling patterns differ between typically developing infants and those with ASD.

Traditional approaches to vocal analysis have primarily relied on manual feature extraction and statistical comparison. Esposito and Venuti (2008) conducted detailed acoustic analyses of cries from infants with ASD, identifying specific patterns in pitch variability and harmonic structure. However, these methods were limited by small sample sizes and the challenge of capturing complex, non-linear relationships in vocal data.

Machine learning applications in ASD research have grown substantially in recent years. Wall et al. (2012) applied support vector machines to home video analysis, achieving promising results in behavioral pattern recognition. In the domain of vocal analysis, Pokorny et al. (2018) demonstrated the feasibility of using acoustic features for ASD classification, though their work focused on older children with established diagnoses.

The current study builds upon this foundation by applying advanced machine learning techniques to a large, diverse dataset of infant vocalizations. Our approach differs from previous work in its specific focus on the pre-diagnostic period and its comprehensive comparison of multiple classification algorithms. We also address methodological limitations of earlier studies by implementing rigorous cross-validation and feature selection procedures.

## Research Questions

This study addresses the following research questions:

- 1. Can machine learning algorithms reliably distinguish between vocalizations of infants who will later receive an ASD diagnosis and those of typically developing infants based solely on acoustic features?
- 2. Which acoustic features (MFCCs, pitch characteristics, spectral properties, temporal patterns) demonstrate the strongest discriminative power for early ASD detection?
- 3. How do different machine learning algorithms (support vector machines, random forests, neural networks) compare in their performance for this classification task?
- 4. What is the optimal age window for detecting ASD-related vocal patterns during the 6-18 month developmental period?

### **Objectives**

The primary objectives of this research are:

- 1. To develop a comprehensive feature extraction pipeline for analyzing infant vocalizations, incorporating both traditional acoustic parameters and novel spectral-temporal features.
- 2. To collect and curate a large, balanced dataset of infant cries from both typically developing infants and those later diagnosed with ASD.
- 3. To implement and optimize multiple machine learning classification models for ASD detection based on vocal features.
- 4. To evaluate model performance using rigorous statistical measures including accuracy, precision, recall, F1-score, and area under the ROC curve.
- 5. To identify the most discriminative vocal features for early ASD detection and analyze their relationship to known neurodevelopmental mechanisms.

# Hypotheses to be Tested

Based on previous literature and theoretical considerations, we propose the following hypotheses:

H1: Infants who later receive an ASD diagnosis will demonstrate statistically significant differences in vocal acoustic features compared to typically developing infants.

H2: Machine learning models trained on acoustic features will achieve classification accuracy significantly above chance level (50

H3: Non-linear classification models (neural networks, random forests) will outperform linear models (logistic regression, linear SVM) due to the complex, non-linear nature of vocal patterns.

H4: Vocal features related to pitch control and harmonic structure will show the strongest discriminative power, reflecting underlying differences in neuromotor control and auditory processing.

H5: Classification performance will be age-dependent, with optimal detection occurring between 12-15 months when vocal development shows the most rapid changes.

# Approach/Methodology

#### **Data Collection and Participants**

We collected audio recordings from 450 infants (225 typically developing, 225 later diagnosed with ASD) aged 6-18 months. Participants were recruited through pediatric clinics and early intervention programs across four countries. Diagnosis confirmation occurred at 36 months using the Autism Diagnostic Observation Schedule (ADOS) and clinical evaluation. Audio samples were recorded during routine pediatric visits using standardized recording equipment in controlled acoustic environments.

#### **Feature Extraction**

From each audio sample, we extracted 128 acoustic features including:

- 13 Mel-frequency cepstral coefficients (MFCCs) - Fundamental frequency (F0) and its statistical moments - Jitter and shimmer measures - Spectral centroid, rolloff, and flux - Harmonic-to-noise ratio - Formant frequencies (F1-F4)

Features were normalized using z-score standardization to account for individual differences in vocal intensity and recording conditions.

#### Machine Learning Models

We implemented and compared four classification algorithms:

1. Support Vector Machine (SVM) with radial basis function kernel 2. Random Forest with 100 decision trees 3. Multilayer Perceptron neural network 4. Logistic Regression as baseline

The performance of each model was evaluated using the following objective function for optimization:

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^{N} L(y_i, f(x_i; \theta)) + \lambda R(\theta)$$
 (1)

where L represents the cross-entropy loss function,  $f(x_i; \theta)$  is the model prediction, and  $R(\theta)$  is the regularization term.

#### **Evaluation Framework**

We employed 10-fold cross-validation with stratified sampling to ensure balanced representation of both classes in each fold. Model performance was assessed using accuracy, precision, recall, F1-score, and area under the receiver operating characteristic curve (AUC-ROC). Statistical significance of performance differences between models was tested using paired t-tests with Bonferroni correction.

#### Results

Our analysis revealed significant differences in vocal characteristics between infants later diagnosed with ASD and typically developing infants. The neural network model achieved the highest overall performance with 89.2

The most discriminative features included MFCC coefficients 1-3 (related to spectral envelope), fundamental frequency variability, and harmonic-to-noise ratio. These features consistently appeared in feature importance rankings across all models.

Age-stratified analysis showed that classification performance improved with infant age, with optimal detection occurring between 12-15 months. This pattern suggests that vocal differences become more pronounced as infants approach key developmental milestones in speech and language acquisition.

Table 1: Performance Comparison of Machine Learning Models for ASD Detection

Model	Accuracy	Precision	Recall	F1-Score	AUC-ROC
Neural Network	0.892	0.885	0.901	0.893	0.945
Support Vector Machine	0.873	0.869	0.878	0.873	0.928
Random Forest	0.841	0.832	0.851	0.841	0.901
Logistic Regression	0.798	0.785	0.812	0.798	0.867

Feature importance analysis revealed that MFCC-based features contributed most significantly to classification performance, accounting for approximately 45

Cross-validation results demonstrated stable performance across different data splits, with standard deviations of less than 2

#### Discussion

Our findings provide strong evidence that machine learning analysis of infant vocalizations can serve as an effective screening tool for early ASD detection. The high classification accuracy achieved by our models, particularly the neural network approach, suggests that subtle vocal differences are present long before traditional behavioral symptoms become apparent.

The superior performance of non-linear models (neural networks and SVM with RBF kernel) supports our hypothesis that vocal patterns associated with ASD involve complex, non-linear relationships that cannot be adequately captured by linear methods. This finding aligns with current understanding of ASD as a condition characterized by atypical neural connectivity and information processing.

The age-dependent pattern of detection accuracy has important clinical implications. The optimal detection window of 12-15 months corresponds to a critical period in language development when infants typically begin producing more complex vocalizations. The increasing discriminative power during this period may reflect emerging differences in neuromotor control and auditory-motor integration.

From a clinical perspective, our approach offers several advantages over current screening methods. The non-invasive nature of audio recording makes it suitable for widespread implementation in pediatric settings. The automated analysis reduces reliance on subjective clinical judgment and could potentially be integrated into mobile health applications for continuous monitoring.

However, several limitations should be considered. Our sample, while larger than previous studies, still represents a specific demographic distribution. Future research should validate these findings in more diverse populations and explore cultural and linguistic variations in vocal development. Additionally, longitudinal follow-up is needed to determine how early vocal patterns relate to later symptom severity and developmental outcomes.

#### Conclusions

This study demonstrates the feasibility and effectiveness of using machine learning for early ASD detection through vocal pattern analysis. Our results show that:

- 1. Machine learning models can achieve high accuracy (up to 89.2)
- 2. Non-linear classification models, particularly neural networks, outperform linear methods for this task, suggesting that vocal patterns associated with ASD involve complex feature interactions.
- 3. The most discriminative features relate to spectral characteristics (MFCCs) and pitch control, potentially reflecting underlying differences in neuromotor

function and auditory processing.

4. Optimal detection occurs between 12-15 months, providing a practical window for early screening and intervention.

These findings have significant implications for clinical practice and public health. The development of automated, objective screening tools based on vocal analysis could substantially reduce the age of ASD diagnosis, enabling earlier access to interventions that improve long-term outcomes. Future work should focus on validating these methods in larger, more diverse populations and integrating them with other behavioral and biological markers for comprehensive early detection systems.

# Acknowledgements

We gratefully acknowledge the participating families and clinical staff who made this research possible. This study was supported by grants from the National Institute of Mental Health (R01MH123456) and the Autism Research Foundation. We thank Dr. Emily Chen for her assistance with statistical analysis and the Tsinghua University Computational Linguistics Laboratory for providing computational resources. The authors declare no conflicts of interest.

99 Sheinkopf, S. J., Mundy, P., Oller, D. K.,

& Steffens, M. (2000). Vocal atypicalities of preverbal autistic children. *Journal of Autism and Developmental Disorders*, 30(4), 345-354.

Oller, D. K., Niyogi, P., Gray, S., Richards, J. A., Gilkerson, J., Xu, D., ... & Warren, S. F. (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.

Esposito, G.,

& Venuti, P. (2008). Comparative analysis of crying in children with autism, developmental delays, and typical development. Focus on Autism and Other Developmental Disabilities, 23(4), 207-215.

Wall, D. P., Dally, R., Luyster, R., Jung, J. Y.,

& DeLuca, T. F. (2012). Use of artificial intelligence to shorten the behavioral diagnosis of autism. *PLoS One*, 7(8), e43855.

Pokorny, F. B., Schuller, B. W., Marschik, P. B., Brueckner, R., Nyström, P., Cummins, N., ...

& Einspieler, C. (2018). Earlier identification of children with autism spectrum disorder: An automatic vocalisation-based approach. *Proceedings of Interspeech*, 2018, 309-313.