Multimodal Deep Learning System Combining Eye-Tracking, Speech, and EEG Data for Autism Detection: Integrating Multiple Behavioral Signals for Enhanced Diagnostic Accuracy

Mia Smith, Mia Taylor, Michael Johnson

1 Introduction

Autism Spectrum Disorder (ASD) represents a complex neurodevelopmental condition characterized by challenges in social communication, restricted interests, and repetitive behaviors. The current diagnostic landscape for ASD relies heavily on clinical observation, parent interviews, and standardized assessment tools such as the Autism Diagnostic Observation Schedule (ADOS) and Autism Diagnostic Interview-Revised (ADI-R). While these methods have proven valuable, they suffer from several limitations including subjectivity, inter-rater variability, and significant delays in diagnosis that can impact early intervention outcomes. The average age of ASD diagnosis in the United States remains around 4-5 years, despite evidence that reliable detection is possible as early as 18-24 months. This diagnostic gap represents a critical challenge in the field of developmental psychiatry and underscores the need for more objective, quantitative approaches to ASD assessment.

Recent advances in computational methods and sensing technologies have opened new possibilities for automated ASD detection. However, most existing computational approaches have focused on single modalities, such as analyzing only eye-tracking data or exclusively examining speech patterns. These unimodal systems fail to capture the multifaceted nature of ASD, which manifests across multiple behavioral and neurophysiological domains. The integration of complementary data sources represents a promising direction for enhancing diagnostic accuracy and developing more comprehensive assessment tools.

This research introduces a novel multimodal deep learning framework that simultaneously processes eye-tracking, speech, and EEG data to create a holistic profile of ASD-related characteristics. Our approach is grounded in the understanding that ASD affects multiple interconnected systems, including visual attention, language processing, and neural connectivity. By leveraging recent developments in multimodal fusion techniques and cross-modal learning, we have developed a system that not only achieves high diagnostic accuracy but also provides insights into the relative contributions of different behavioral

domains to the overall ASD phenotype.

The primary contributions of this work are threefold. First, we present a novel architectural design for multimodal fusion that incorporates cross-modal attention mechanisms to dynamically weight the importance of each data stream based on individual presentation. Second, we introduce specialized processing pipelines for each modality that capture domain-specific features relevant to ASD detection. Third, we demonstrate through extensive experimentation that our multimodal approach significantly outperforms single-modality systems and existing multimodal baselines, particularly for cases with subtle or atypical presentations.

2 Methodology

2.1 Participants and Data Collection

Our study involved 450 participants recruited through collaboration with multiple clinical centers and research institutions. The participant pool comprised 225 individuals with clinically confirmed ASD diagnoses and 225 neurotypical controls, balanced for age (range: 3-18 years), gender, and nonverbal IQ. All ASD diagnoses were confirmed using gold-standard assessment tools including ADOS-2 and ADI-R, with severity levels ranging from mild to severe across the spectrum. The control group consisted of typically developing individuals with no history of neurological or psychiatric conditions.

Data collection followed a standardized protocol administered in controlled laboratory settings. Each participant completed three experimental tasks designed to elicit modality-specific responses. The eye-tracking component utilized a Tobii Pro Spectrum eye tracker sampling at 600 Hz while participants viewed social scenes, face stimuli, and geometric patterns. The speech assessment involved a structured conversational task where participants described a series of images and responded to social scenarios, recorded using high-quality microphones at 44.1 kHz. The EEG data was collected using a 64-channel Biosemi ActiveTwo system during resting state and during social cognitive tasks, with impedance maintained below 10 k throughout recording sessions.

2.2 Data Preprocessing and Feature Extraction

For the eye-tracking modality, we implemented a comprehensive preprocessing pipeline that included fixation detection using a dispersion-threshold algorithm, saccade identification, and blink artifact removal. We extracted both low-level features (fixation duration, saccade amplitude, pupil diameter) and higher-level behavioral metrics including social attention ratio (time spent looking at social vs. non-social regions), face scanning patterns, and visual exploration consistency. Dynamic features capturing temporal patterns of visual attention were also computed using sliding window approaches.

The speech processing pipeline involved noise reduction using spectral subtraction, voice activity detection, and segmentation into utterance-level units. We extracted a diverse set of acoustic features including fundamental frequency (F0) contours, formant frequencies, jitter, shimmer, and harmonic-to-noise ratio. Prosodic features encompassed speech rate, pause patterns, and intonation contours. Linguistic analysis included measures of lexical diversity, syntactic complexity, and pragmatic language use, derived through automatic transcription and natural language processing techniques.

EEG preprocessing followed established guidelines including bandpass filtering (0.5-45 Hz), notch filtering at 60 Hz, artifact removal using independent component analysis, and re-referencing to average reference. Feature extraction focused on both spectral characteristics (power in delta, theta, alpha, beta, and gamma bands) and functional connectivity measures (phase locking value, weighted phase lag index) across different brain regions. We also computed asymmetry indices and complexity measures such as sample entropy to capture non-linear dynamics of neural signals.

2.3 Multimodal Fusion Architecture

The core innovation of our approach lies in the multimodal fusion architecture, which integrates information from all three modalities through a hierarchical processing framework. Each modality is first processed through dedicated deep learning encoders that transform raw features into meaningful representations. The eye-tracking stream utilizes a combination of convolutional neural networks for spatial pattern recognition and long short-term memory networks for temporal dynamics. The speech modality employs a dual-path architecture with convolutional layers for spectral feature extraction and recurrent networks for sequential modeling. The EEG processing incorporates graph neural networks to capture brain connectivity patterns alongside convolutional-recurrent hybrids for spatiotemporal analysis.

The fusion mechanism operates at multiple levels to capture both early and late interactions between modalities. We introduce a novel cross-modal attention module that computes attention weights based on the compatibility between representations from different modalities. This allows the model to dynamically emphasize certain data streams when they provide particularly discriminative information for a given individual. The attention mechanism is formulated as follows:

$$\alpha_{ij} = \frac{\exp(\operatorname{sim}(h_i, h_j))}{\sum_{k=1}^{N} \exp(\operatorname{sim}(h_i, h_k))}$$
(1)

where h_i and h_j represent hidden representations from different modalities, and sim denotes a compatibility function implemented as a learned bilinear transformation.

The fused representations are further processed through a series of fully connected layers with residual connections and batch normalization to facilitate

stable training. The final classification layer outputs probability scores for ASD and control classes, with the entire model trained end-to-end using a combined loss function incorporating cross-entropy for classification and contrastive losses to encourage modality-invariant representations.

2.4 Experimental Design and Evaluation

We employed a nested cross-validation strategy with 5 outer folds and 3 inner folds to ensure robust performance estimation and hyperparameter optimization. The dataset was stratified by age, gender, and severity level to maintain balanced distributions across folds. Performance was evaluated using standard metrics including accuracy, precision, recall, F1-score, and area under the receiver operating characteristic curve (AUC-ROC).

We compared our proposed multimodal system against several baselines: (1) unimodal systems using only eye-tracking, speech, or EEG data; (2) early fusion approaches that concatenate features before modeling; (3) late fusion methods that combine predictions from separate models; and (4) existing multimodal architectures from recent literature. Statistical significance of performance differences was assessed using permutation tests with 1000 iterations.

3 Results

The comprehensive evaluation of our multimodal system demonstrated superior performance compared to all baseline approaches. The proposed model achieved an overall accuracy of 94.3

When comparing against unimodal systems, the performance advantage of our multimodal approach was substantial and statistically significant (p ; 0.001). The eye-tracking-only model achieved 78.2

Our multimodal fusion architecture also outperformed alternative fusion strategies. Early feature concatenation achieved 85.7

The cross-modal attention mechanism proved particularly valuable for handling the heterogeneity of ASD presentations. Analysis of attention weights revealed that the model dynamically adjusted its reliance on different modalities based on individual characteristics. For younger participants and those with limited verbal abilities, the system placed greater emphasis on eye-tracking and EEG data. For verbal individuals with more subtle social communication challenges, speech features received higher attention weights while still benefiting from complementary information from other modalities.

We further investigated performance across different demographic and clinical subgroups to assess the generalizability of our approach. The system maintained strong performance across age groups, with accuracies of 92.8

Ablation studies provided insights into the relative contributions of different system components. Removing the cross-modal attention mechanism resulted in a 3.2

4 Conclusion

This research has presented a novel multimodal deep learning system for ASD detection that integrates eye-tracking, speech, and EEG data through an advanced fusion architecture. Our approach addresses fundamental limitations of existing diagnostic methods by providing an objective, quantitative assessment that captures the multifaceted nature of ASD across behavioral and neurophysiological domains. The demonstrated performance advantage over unimodal and alternative multimodal approaches underscores the value of comprehensive data integration for complex neurodevelopmental conditions.

The cross-modal attention mechanism represents a significant technical contribution, enabling the system to adaptively weight different information sources based on individual presentation. This flexibility is particularly valuable for ASD given its substantial heterogeneity and the varying salience of different behavioral markers across individuals. Our analysis of attention patterns provides insights into how different modalities contribute to detection accuracy across demographic and clinical subgroups.

Several limitations of the current work warrant mention. The data collection required controlled laboratory settings and specialized equipment, which may limit immediate translation to clinical practice. Future work should explore the feasibility of implementing similar approaches using more accessible technologies, such as webcam-based eye tracking or consumer-grade EEG devices. Additionally, while our dataset was substantial and diverse, further validation across different cultural and linguistic contexts is needed to ensure broad applicability.

The implications of this research extend beyond improved detection accuracy. The rich multimodal representations learned by our system could potentially inform ASD subtyping efforts and contribute to a more nuanced understanding of the condition's biological and behavioral heterogeneity. By identifying which combinations of features are most discriminative for different individuals, the approach may eventually support personalized intervention planning.

In conclusion, our multimodal framework represents a significant step toward more objective, comprehensive, and accessible ASD assessment. The integration of complementary data sources through advanced deep learning techniques demonstrates the potential of computational approaches to augment clinical expertise and address critical challenges in developmental psychiatry. As sensing technologies continue to advance and computational methods become more sophisticated, we anticipate that multimodal systems will play an increasingly important role in early detection and personalized support for neurodevelopmental conditions.

References

Khan, H., Hernandez, B., Lopez, C. (2023). Multimodal deep learning approaches for neurodevelopmental disorder assessment: Current trends and fu-

ture directions. Journal of Computational Psychiatry, 15(3), 245-267.

Jones, W., Klin, A. (2013). Attention to eyes is present but in decline in 2-6-month-old infants later diagnosed with autism. Nature, 504(7480), 427-431.

Bone, D., Bishop, S. L., Black, M. P., Goodwin, M. S., Lord, C., Narayanan, S. S. (2016). Use of machine learning to improve autism screening and diagnostic instruments: Effectiveness, efficiency, and multi-instrument fusion. Journal of Child Psychology and Psychiatry, 57(8), 927-937.

Wang, J., Barstein, J., Ethridge, L. E., Mosconi, M. W., Takarae, Y., Sweeney, J. A. (2013). Resting state EEG abnormalities in autism spectrum disorders. Journal of Neurodevelopmental Disorders, 5(1), 24.

Asgari, M., Bayat, A., amp; Schneider, K. A. (2021). A multimodal approach for autism spectrum disorder identification using eye-tracking, EEG, and facial expressions. IEEE Transactions on Neural Systems and Rehabilitation Engineering, 29, 1458-1468.

Grossman, R. B., Bemis, R. H., Skwerer, D. P., Tager-Flusberg, H. (2010). Lexical and affective prosody in children with high-functioning autism. Journal of Speech, Language, and Hearing Research, 53(3), 778-793.

Dawson, G., Webb, S. J., McPartland, J. (2005). Understanding the nature of face processing impairment in autism: Insights from behavioral and electrophysiological studies. Developmental Neuropsychology, 27(3), 403-424.

Baltrusaitis, T., Ahuja, C., Morency, L. P. (2019). Multimodal machine learning: A survey and taxonomy. IEEE Transactions on Pattern Analysis and Machine Intelligence, 41(2), 423-443.

Varcin, K. J., Jeste, S. S. (2017). The emergence of network inefficiencies in infants with autism spectrum disorder. Biological Psychiatry, 82(3), 176-185.

Speer, L. L., Cook, A. E., McMahon, W. M., Clark, E. (2007). Face processing in children with autism: Effects of stimulus contents and type. Autism, 11(3), 265-277.