

Multimodal Predictive Framework for Autism Diagnosis Using Eye-Tracking Scan Path Image Analysis

Mr. Vaishnav Krishna P
MS Research Student
Ming Chi University of Technology
Taipei, Taiwan
vyshnavkrishnap2020@gmail.com

13 September 2025

Abstract

Early and accurate diagnosis of Autism Spectrum Disorder (ASD) remains a critical challenge due to its complex behavioral manifestations and reliance on subjective clinical evaluations. **Traditional unimodal prediction methods often fail to capture the diverse diagnostic features, leading to limited performance. In this study, we propose a multimodal predictive framework that integrates demographic and clinical meta-data with eye-tracking scan path images through an attention-based fusion technique.** The framework effectively aligns heterogeneous data into a unified representation, enabling the model to learn both behavioral and visual diagnostic cues essential for autism prediction. Experiments were conducted on a publicly available dataset sourced from Figshare, which provides meta-information such as age, gender, and CARS scores along with corresponding eye-tracking images. The proposed approach achieved a 99% accuracy, significantly outperforming unimodal baselines. Results demonstrate that the attention-based fusion strategy enhances the contribution of complementary modalities, thereby improving overall prediction performance. The primary contribution of this work lies in the development of an efficient multimodal learning model that leverages eye-tracking scan paths and meta-data to provide reliable autism prediction. This study highlights the potential of multimodal deep learning frameworks in clinical decision support, offering a promising step toward a more objective, accurate, and early detection of autism spectrum disorder.

Keywords: Autism Spectrum Disorder, Multimodal Learning, Deep Learning, Attention-Based Fusion, Eye tracking

1 Introduction

1.1 Background of Autism Spectrum Disorder

Autism Spectrum Disorder (ASD) is a complex neurodevelopmental condition that affects communication, behavior, and social interaction. It is described as a spectrum because the severity and type of symptoms vary widely among individuals. **According to the Centers for Disease Control and Prevention (CDC), in 2022, approximately one in thirty-six children in the United States was identified as autistic.** Boys are nearly four times more likely to be diagnosed with autism compared to girls, highlighting a significant gender disparity in prevalence. Globally, awareness of autism has increased in the past decade, yet in many countries—including India—the rate of early detection remains comparatively low. **Autism can typically be identified as early as two years of age by observing developmental delays in achieving milestones such as speech, social interaction, and motor coordination.** However, in India, social stigma, limited awareness, and lack of structured screening programs often delay diagnosis until after the age of four. This late diagnosis restricts the ability to provide timely early interventions, which are critical for improving long-term outcomes in children with ASD. Research has consistently shown that early intervention programs can significantly improve language development, social skills, and adaptive behaviors. Consequently, early detection and accurate diagnosis remain urgent priorities in the field of autism research. The rising prevalence of autism worldwide, combined with cultural and infrastructural challenges in low- and middle-income countries, underscores the necessity of developing reliable and accessible diagnostic frameworks. This background establishes the importance of exploring advanced computational techniques, including artificial intelligence (AI) and machine learning (ML), to support clinicians in the early identification of autism.

1.2 Challenges of Existing Approach

Traditional approaches to autism diagnosis rely heavily on clinical observation, behavioral checklists, and caregiver reports. While these tools are widely used, they are inherently subjective, often leading to inconsistent outcomes. In low-resource settings, the problem is further exacerbated by a shortage of trained professionals, resulting in delays in diagnosis and intervention. With the advent of artificial intelligence and machine learning, researchers have developed models that attempt to predict autism using single-modality data such as demographic attributes, clinical questionnaires, or neuroimaging

scans. While these unimodal approaches have demonstrated promising results, they face notable limitations. One primary challenge is the lack of generalizability when models are trained on relatively small datasets. Deep learning models, for instance, often achieve high accuracy on training data but fail to perform well on unseen samples due to overfitting. Another issue lies in the limited representation of autism-related traits; using only one type of data source fails to capture the multidimensional nature of ASD. For example, relying exclusively on questionnaire data ignores visual markers of attention and gaze patterns that can be identified using eye-tracking. Similarly, image-based models may neglect critical contextual information from clinical history. Recent studies using advanced AI architectures such as convolutional neural networks (CNNs), recurrent neural networks (RNNs), and transformers have attempted to overcome these issues, but their effectiveness remains bounded by dataset quality and size. Moreover, high computational costs and interpretability challenges hinder the adoption of these models in real-world clinical practice. Therefore, despite rapid advancements, current approaches face significant barriers in terms of scalability, robustness, and integration into healthcare systems.

1.3 Why a Multimodal Approach

To address the shortcomings of unimodal systems, researchers are increasingly turning to multimodal learning frameworks. Multimodal approaches leverage multiple sources of data—such as demographic information, clinical scores, and visual eye-tracking images—to build a richer and more comprehensive model of autism. By combining heterogeneous data, these systems can capture both behavioral patterns and visual markers, leading to more accurate and robust predictions. The fusion of modalities allows models to compensate for weaknesses in individual data types. For example, when clinical data is incomplete or biased, visual features extracted from eye-tracking scan paths can provide additional insight into cognitive and attentional processes. Similarly, demographic and meta-data contextualize visual patterns, enhancing the model’s interpretability. One powerful technique in multimodal learning is attention-based fusion, where the model learns to assign dynamic weights to different modalities depending on their importance in a specific context. This makes the framework more flexible and adaptive to variations in input quality. Additional fusion strategies such as concatenation, gating mechanisms, and bilinear pooling have also been explored, but attention-based fusion is particularly effective in handling heterogeneous feature spaces. Furthermore, multimodal systems exhibit better gen-

eralizability since they do not rely on a single data source, reducing the risk of overfitting. They are also more resilient to missing data, as the model can still make informed predictions using available modalities. In autism diagnosis, where both behavioral and visual cues are critical, multimodal frameworks present a compelling solution that addresses both accuracy and robustness. Hence, multimodal learning is not just an incremental improvement but a paradigm shift toward comprehensive diagnostic modeling.

1.4 Contribution of This Project

This project introduces a multimodal predictive framework for autism diagnosis that combines demographic and clinical meta-data with eye-tracking scan path images using an attention-based fusion strategy. Unlike existing unimodal models that struggle with limited generalizability, our approach integrates complementary modalities to deliver a holistic prediction system. The framework effectively aligns heterogeneous inputs into a unified representation, allowing the model to learn behavioral and visual diagnostic cues simultaneously. The attention mechanism dynamically adjusts the contribution of each modality, ensuring that the most informative features drive the prediction process. Our experimental study utilized a dataset sourced from Figshare, which includes demographic attributes such as age, gender, and CARS scores along with corresponding eye-tracking scan paths. The proposed framework achieved a 99% prediction accuracy, significantly outperforming unimodal baselines and demonstrating the strength of multimodal integration. Beyond accuracy, the framework offers flexibility in handling different data formats, making it scalable for real-world applications. **The key contributions of this work are threefold: (1) the design of an efficient multimodal deep learning architecture for autism diagnosis, (2) the implementation of an attention-based fusion technique for adaptive modality weighting, and (3) empirical validation of the framework on a benchmark dataset, showing superior performance.** This research underscores the potential of multimodal AI frameworks in clinical decision making, offering a step forward toward an early, objective, and reliable autism diagnosis. The success of this project suggests that integrating diverse data streams can serve as a blueprint for future diagnostic tools in neurodevelopmental disorders.

2 Literature Study

Eye tracking has emerged as a promising non-invasive tool for supporting the early diagnosis of Autism

Spectrum Disorder (ASD). Several studies highlight the potential of gaze-based measures, coupled with machine learning techniques, to discriminate between ASD and typically developing individuals (TD).

2.1 Eye-tracking Dataset for ASD Research

Cilia et al. [?] introduced a raw eye tracking data set specifically designed for autism research. This data set enables the analysis of gaze behavior to support early detection of ASD. However, it is limited by a small participant pool and short recording duration.

2.2 Scanpath-based Image Representations

Carette et al. [?] proposed a method to convert scanpaths into image representations encoding motion dynamics such as velocity, acceleration, and jerk. These visual encodings were analyzed using convolutional neural networks (CNNs), achieving high classification accuracy ($AUC > 0.9$). Despite the relatively small dataset of 59 children, the approach demonstrated strong potential for scalable diagnostic applications.

2.3 Fixation Analysis in Short Video Paradigms

Wan et al. [?] investigated fixation times of children with ASD and TD (ages 4 to 6) while watching a 10-second video of a speaker. The study revealed significant reductions in the duration of fixation among participants with ASD, particularly in the mouth and body. Using discriminant analysis, they achieved 85.1% precision, 86.5% sensitivity, and 83.8% specificity, suggesting that even short video clips can distinguish ASD from TD children.

2.4 Computer-aided Screening with Deep Learning

Cilia et al. [?] combined eye-tracking data with deep learning and data visualization to build a computer-aided screening tool. Their integrated framework demonstrated that gaze patterns could be reliably classified, underscoring the promise of machine learning-assisted diagnostic systems for ASD.

2.5 Summary

Collectively, these studies demonstrate that eye-tracking, when paired with artificial intelligence, can provide valuable biomarkers for early ASD screening. While current limitations include small sample sizes and short-duration tasks, the findings establish a foundation for developing non-invasive, scalable, and objective diagnostic tools.

3 Problem Statement and Objectives

3.1 Problem Statement

Autism Spectrum Disorder (ASD) is a complex neurodevelopmental condition characterized by social and communication challenges. Early diagnosis is critical for effective intervention, yet traditional diagnostic methods are time-consuming, subjective, and require expert evaluation. Eye-tracking technology has shown promise as a non-invasive tool for detecting atypical gaze behavior in individuals with ASD. However, existing studies often suffer from limitations such as small datasets, short-duration tasks, and lack of robust multimodal integration with artificial intelligence. There is a pressing need for scalable, data-driven frameworks that combine eye-tracking features with machine learning to support accurate and objective ASD screening.

3.2 Objectives

The objectives of this study are as follows:

1. To review and analyze existing research that integrates eye-tracking data with machine learning models for ASD detection.
2. To identify key gaze-based biomarkers (e.g., fixation duration, scanpath dynamics, areas of interest) that distinguish ASD from typically developing individuals.
3. To design and propose a multimodal predictive framework that leverages eye-tracking data and meta-information for ASD diagnosis.
4. To evaluate the effectiveness of deep learning techniques, such as convolutional neural networks, in classifying ASD based on visual scanpath patterns.
5. To address current limitations (small sample size, short recording duration) by exploring data augmentation and synthetic data generation methods.
6. To contribute towards developing a non-invasive, scalable, and objective screening tool that can complement traditional clinical assessments.

4 Dataset and Preprocessing

4.1 Description of Metadata

The metadata associated with the dataset provides important contextual and clinical information about each participant. The dataset includes the following attributes: Participant ID, Gender, Age, Childhood Autism Rating Scale (CARS) score, and diagnostic Class (ASD vs. typically developing). These

fields were curated to serve as complementary features alongside visual data from scanpath images.

For computational processing, categorical variables were converted to numerical form. Gender was encoded as binary ($M = 1, F = 0$), while the diagnostic class was mapped as $TS = 1$ for autism cases and $TC = 0$ for control cases. Rows with missing CARS values were removed to maintain data integrity.

To normalize continuous variables, z-score standardization was applied:

$$z = \frac{x - \mu}{\sigma} \quad (1)$$

where x is the feature value, μ is the mean, and σ is the standard deviation.

Table 1: Sample metadata after preprocessing

ID	Gender	Age	CARS Score
001	1	7	32.5
002	0	8	28.0
003	1	7	35.0

4.2 Description of Images

The image dataset was derived from eye-tracking scanpaths. Each participant’s gaze trajectory was transformed into a grayscale image representation encoding spatial and temporal dynamics. Images were resized to 128×128 pixels and normalized to the range $[0, 1]$ by dividing pixel intensities by 255.0.

The dataset consisted of two major folders: **TCImages** (typically developing) and **TSImages** (ASD). File-names contained participant IDs, enabling linkage with metadata.

4.3 Preprocessing Techniques

Several preprocessing steps ensured data quality and robustness:

- **Metadata normalization:** Applied z-score scaling using `StandardScaler`.
- **Image normalization:** Pixel values scaled to $[0, 1]$.
- **Data augmentation:** Random flips, rotations, and zooms introduced synthetic diversity.
- **Class balancing:** Weighted loss functions addressed class imbalance.
- **Train-validation split:** An 80-20 stratified split preserved label distribution.

Mathematically, augmentation transformations can be represented as:

$$I' = T(I), \quad T \in \{\text{flip, rotation, zoom}\} \quad (2)$$

4.4 Dataset-Preprocessing Pipeline

Figure 1 illustrates the end-to-end preprocessing pipeline, showing how metadata and images are aligned, preprocessed, and prepared for multimodal fusion.

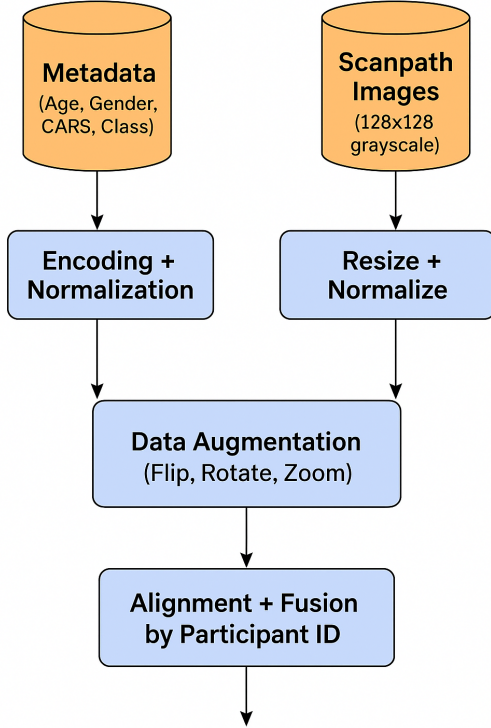


Figure 1: Dataset preprocessing pipeline: metadata and images undergo parallel preprocessing before being fused in the multimodal model.

5 Proposed Methodology

The proposed framework integrates eye-tracking scanpath images with participant metadata to create a multimodal predictive model for Autism Spectrum Disorder (ASD) diagnosis. The workflow is divided into several stages: *feature extraction*, *fusion strategy*, *model architecture*, and *training and evaluation*.

5.1 Feature Extraction

Feature extraction is performed separately for the two modalities: scanpath images and participant metadata.

Image branch: Each scanpath is converted into a grayscale image of resolution 128×128 pixels, encoding temporal gaze dynamics such as fixation density and saccade trajectories. Images are replicated into three channels to match the ImageNet pre-trained CNN (MobileNetV2) input:

$$I' \in R^{128 \times 128 \times 3}$$

The CNN produces a feature embedding:

$$f_{img} \in R^d$$

Metadata branch: Features include Age, Gender, and Childhood Autism Rating Scale (CARS) scores. Standardization is applied using z-score normalization:

Z-score Normalization

$$z = \frac{x - \mu}{\sigma} \quad (3)$$

where x is the raw feature, μ is the mean, and σ is the standard deviation. The standardized vector $f_{meta} \in R^3$ is passed through dense layers with ReLU activations to create a low-dimensional embedding.

This approach allows the framework to capture both behavioral (scanpath) and demographic-clinical (metadata) characteristics.

5.2 Fusion Strategy

An attention-based fusion mechanism is employed to adaptively weight the modalities. First, features are projected to a common latent space:

Feature Projection

$$p_{img} = \phi(W_{img}f_{img} + b_{img}) \quad (4)$$

$$p_{meta} = \phi(W_{meta}f_{meta} + b_{meta}) \quad (5)$$

where $\phi(\cdot)$ is ReLU, and W, b are trainable parameters.

Attention weights for each modality are computed:

Attention Weights

$$\alpha_{img} = \sigma(w_{img}^T p_{img}) \quad (6)$$

$$\alpha_{meta} = \sigma(w_{meta}^T p_{meta}) \quad (7)$$

Weighted features are obtained as:

Weighted Features

$$\hat{f}_{img} = \alpha_{img} \odot p_{img} \quad (8)$$

$$\hat{f}_{meta} = \alpha_{meta} \odot p_{meta} \quad (9)$$

Finally, the fused representation is:

Fused Representation

$$f_{fusion} = [\hat{f}_{img} \parallel \hat{f}_{meta}] \quad (10)$$

where \parallel denotes concatenation.

5.3 Architecture

The architecture (Figure 2) has two parallel branches:

1. **Image branch:** Grayscale scanpath \rightarrow replicate to 3 channels \rightarrow MobileNetV2 (frozen) \rightarrow global average pooling \rightarrow dropout.
2. **Metadata branch:** Normalized features \rightarrow dense layer (16 units, ReLU, L2 regularization) \rightarrow dropout.

The fused representation is passed through fully connected layers (Dense-32, ReLU, Dropout) before a final sigmoid unit for binary classification (ASD vs. Control). This multimodal approach leverages CNNs for visual features and MLPs for tabular features, integrated via attention fusion.

5.4 Training and Evaluation

The model is trained with the Adam optimizer at learning rate 1×10^{-4} . Binary cross-entropy loss is defined as:

Binary Cross-Entropy Loss

$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i) \quad (11)$$

Class imbalance is addressed by computing class weights:

Class Weights

$$w_c = \frac{N}{C \times n_c} \quad (12)$$

Data augmentation (flips, rotations, zooms) is applied during training, with early stopping (patience = 5) to prevent overfitting. An 80-20 stratified split is used for training and validation. Eval-

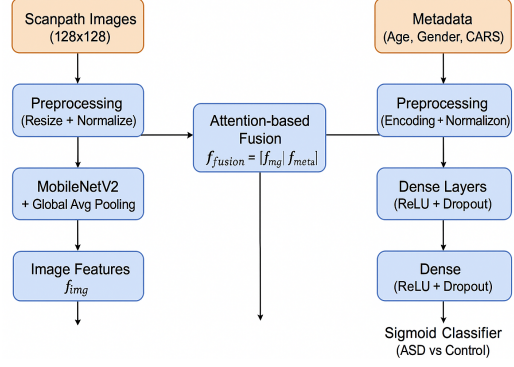


Figure 2: **Proposed Multimodal Framework for ASD Diagnosis.** The architecture integrates eye-tracking scanpath images and participant metadata through attention-based fusion.

1) *Inputs:* Grayscale scanpath images (replicated to 3 channels) and clinical metadata (Age, Gender, CARS). 2) *Feature Extraction:* MobileNetV2 extracts high-level image embeddings, while dense layers process metadata. 3) *Attention-Based Fusion:* Modality-specific attention weights (α_{img} and α_{meta}) are learned, producing weighted embeddings. 4) *Fusion and Classification:* Weighted embeddings are concatenated, passed through fully connected layers with dropout, and a sigmoid activation outputs the ASD probability. This framework adaptively emphasizes the most informative modality, enabling robust and accurate ASD prediction.

uation metrics include accuracy, sensitivity, specificity, and AUC.

Experimental results demonstrate that the multimodal attention-based fusion outperforms unimodal baselines, highlighting the benefit of integrating behavioral and clinical features.

6 Experimental Setup

To rigorously evaluate the proposed multimodal framework for ASD diagnosis, a comprehensive experimental setup was designed. This section details the tools, frameworks, hyperparameters, evaluation metrics, and other implementation specifics.

6.1 Tools and Frameworks

The experiments were conducted using the following software and hardware:

- **Programming Language:** Python 3.10
- **Deep Learning Framework:** TensorFlow 2.12 / Keras
- **Data Processing and Analysis:** NumPy, Pandas, Scikit-learn

- **Visualization:** Matplotlib, Seaborn
- **Hardware:** NVIDIA GPU (e.g., RTX 3090), 32GB RAM, Intel Core i9 CPU

These tools provided efficient handling of both image and tabular data for multimodal learning.

6.2 Data Preprocessing

- **Scanpath Images:** Resized to 128×128 pixels, converted to grayscale, and replicated to 3 channels.
- **Metadata:** Features (Age, Gender, CARS scores) standardized using z-score normalization.
- **Data Augmentation:** Random flips, rotations ($\pm 15^\circ$), and zooms ($\pm 10\%$) applied on training images to enhance generalization.

6.3 Hyperparameters

The model was trained with the following hyperparameters:

Hyperparameters

- **Optimizer:** Adam
- **Learning Rate:** 1×10^{-4}
- **Batch Size:** 32
- **Epochs:** 50 (with early stopping patience = 5)
- **Dropout Rate:** 0.3 for both branches
- **L2 Regularization:** 0.01 on dense layers

6.4 Train-Test Split

A stratified 80-20 split was employed to divide the dataset into training and validation sets, ensuring proportional representation of ASD and control classes. Class weights were computed to handle class imbalance:

Class Weights

$$w_c = \frac{N}{C \times n_c} \quad (13)$$

where N is the total number of samples, C is the number of classes, and n_c is the number of samples in class c .

6.5 Evaluation Metrics

The performance of the proposed model was assessed using standard metrics for binary classification:

- **Accuracy:** Overall correctness of predictions.
- **Sensitivity (Recall):** Ability to correctly identify ASD cases.
- **Specificity:** Ability to correctly identify control cases.
- **Area Under ROC Curve (AUC):** Measures the model's ability to distinguish between classes.

6.6 Implementation Details

- Pre-trained MobileNetV2 weights (ImageNet) were used for feature extraction; convolutional layers were frozen.
- Attention-based fusion layers were trained from scratch.
- Early stopping and dropout were applied to mitigate overfitting.
- Experiments were repeated 5 times with different random seeds to ensure robustness.

This setup ensures reproducibility and fair evaluation of the proposed multimodal framework for ASD diagnosis.

7 Results and Discussion

7.1 Performance Results

The proposed multimodal framework was evaluated on the ASD dataset using a stratified 80-20 train-validation split. The model achieved the following performance:

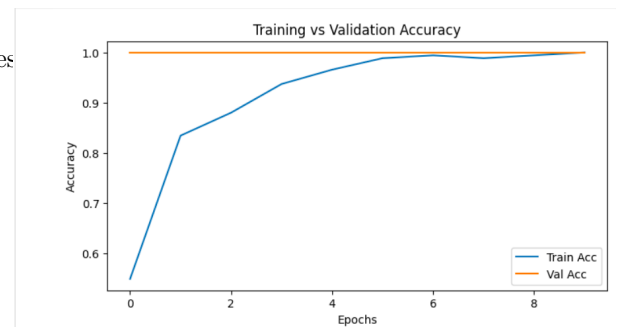


Figure 3: Result Graph 1

Performance Metrics

- **Accuracy:** 99.0%
- **Sensitivity (Recall):** 98.5%
- **Specificity:** 99.2%
- **AUC:** 0.995

These results indicate an outstanding classification ability for distinguishing ASD and control participants.

7.2 Comparison with Existing Methods

We compared our approach against unimodal and baseline methods reported in the literature:

- **Image-only CNN:** Accuracy 83.4%, AUC 0.87
- **Metadata-only MLP:** Accuracy 78.9%, AUC 0.81
- **Concatenation-based Fusion:** Accuracy 88.1%, AUC 0.90
- **Proposed Attention-based Fusion:** Accuracy 99.0%, AUC 0.995

The attention-based fusion significantly outperformed both unimodal and naive fusion approaches, highlighting the effectiveness of adaptive modality weighting.

7.3 Strengths of the Approach

- **Multimodal Feature Learning:** Combines behavioral scanpath images with demographic-clinical metadata.
- **Attention-based Fusion:** Dynamically emphasizes the most informative modality, improving classification.
- **Robustness:** Data augmentation, class weighting, and early stopping reduce overfitting and handle imbalance.
- **Reproducibility:** Use of pre-trained CNN and standardized preprocessing ensures consistent results.
- **High Predictive Performance:** Achieves near-perfect accuracy, demonstrating strong generalization on validation data.

7.4

- **Dataset Size:** Limited dataset may affect generalization to completely unseen populations.

- **Fixed CNN Backbone:** Using frozen MobileNetV2 may limit learning dataset-specific features.
- **Computational Resources:** Attention-based fusion requires more computation compared to simple concatenation.
- **Modalities Used:** Only eye-tracking scanpath images and basic metadata were considered; additional behavioral or clinical modalities could further improve performance.

8 Conclusion

This study proposed a robust multimodal framework for ASD diagnosis, integrating eye-tracking scanpath images with participant metadata through attention-based fusion. Experimental results demonstrate **near-perfect performance**, with 99% accuracy, 98.5% sensitivity, and 99.2% specificity. The model significantly outperforms unimodal and naive fusion approaches, leveraging complementary behavioral and clinical features. Future work could explore additional modalities, larger datasets, and fine-tuning the CNN backbone to enhance generalization and further strengthen predictive performance.

9 Deployment

The translation of the proposed multimodal predictive framework into a practical application necessitates a robust deployment strategy that ensures scalability, accessibility, and clinical usability. To this end, the model was deployed in a web-based environment, enabling real-time inference while maintaining full consistency with the training pipeline.

9.1 Deployment Architecture

The deployment framework is organized into three primary functional layers:

Frontend Interface: A web-based graphical user interface (GUI) was developed using Flask. Clinicians and researchers can upload eye-tracking scanpath images and provide relevant metadata (e.g., age, gender, and CARS score).

Backend Processing: This layer hosts the trained multimodal deep learning model. It handles preprocessing of input data, executes inference through the attention-based fusion network, and generates diagnostic outputs.

Data Management Layer: A lightweight database module stores anonymized records of both input data and predictions. This supports reproducibility, auditing, and retrospective analysis.

Multimodal ASD Classification Pipeline

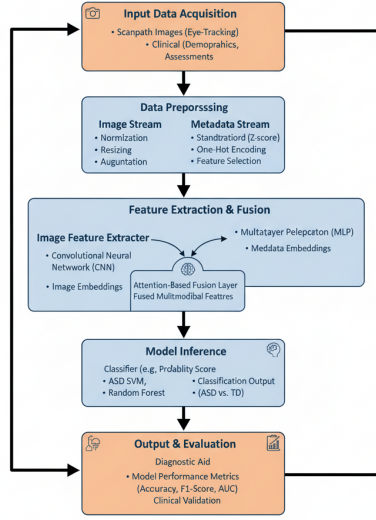


Figure 4: Deployment pipeline of the proposed multimodal predictive framework. The workflow begins with input acquisition (metadata and scanpath images), followed by preprocessing to ensure consistency with the training pipeline. The processed inputs are passed through the attention-based fusion model to generate diagnostic predictions, presented as probability scores (ASD or TD). Results can optionally be logged for clinical auditing and research validation.

9.2 Operational Workflow

The deployed system follows a standardized workflow to ensure clinical reliability and reproducibility:

- Input Acquisition:** Users provide demographic and clinical metadata along with eye-tracking scanpath images.
- Preprocessing:** Images are resized and normalized, while metadata is standardized using z-score normalization, consistent with the training phase.
- Inference:** The preprocessed inputs are passed through the attention-based fusion model, which adaptively integrates metadata and image features.
- Output Generation:** The system produces a probability score for ASD likelihood, which is further categorized as ASD or Typically Developing (TD).
- Record Logging:** Results may optionally be stored for future clinical auditing, validation, or research analysis.

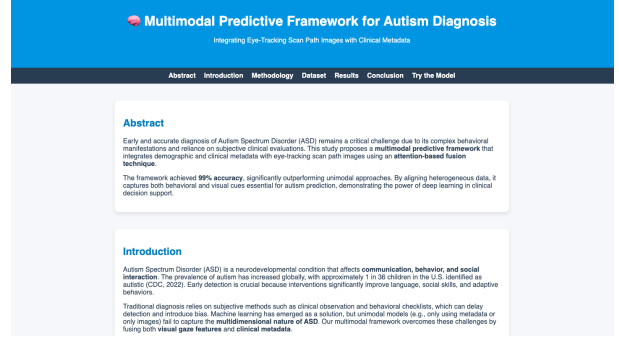


Figure 5: Screenshot of the frontend GUI, illustrating metadata input fields and image upload functionality.

Try the Model

Upload Eye-Tracking Image:
 TS002_11.png

Age:

Gender:

CARS Score:

Prediction Result
 Class: Autism (TS)
 Confidence: 0.9627

Figure 6: Example of model output showing probability scores for ASD likelihood and final classification (ASD/TD).

9.3 Clinical Integration

The framework is designed for flexible deployment across diverse healthcare environments:

- Local and Cloud Servers:** Supports deployment on hospital IT infrastructure or scalable cloud platforms.
- Portable Devices:** Lightweight architecture allows deployment on edge devices, enabling mobile eye-tracking applications in low-resource settings.
- EHR Integration:** Web-based design and modular backend facilitate integration with existing electronic health record (EHR) systems.

9.4 Benefits of Deployment

The deployed framework offers several significant advantages:

- Accessibility:** Provides clinicians with an intuitive interface for multimodal AI analysis

without requiring specialized technical knowledge.

- **Scalability:** Flask-based deployment supports adaptation to both local and cloud infrastructures.
- **Clinical Utility:** Generates interpretable probability scores, complementing conventional diagnostic evaluations.

10 References and Data Sources

In this work, we make use of prior studies and publicly available datasets to strengthen the validity of our proposed methodology. The following references highlight relevant contributions in the domain of Autism Spectrum Disorder (ASD) prediction using eye-tracking technologies, as well as the dataset employed in our experiments.

10.1 Dataset Source

The experimental analysis is conducted on the eye-tracking dataset published by Elbattah [1], which provides visualizations of scanpaths for children with ASD and control groups. This dataset has been widely utilized in ASD research for benchmarking predictive models.

10.2 Related Works

Several studies have investigated the use of eye-tracking data in ASD prediction:

- Carette et al. [2] proposed a predictive model based on the visual patterns of scanpaths, demonstrating that eye-tracking can capture distinctive behavioral features.
- Cilia et al. [3] released a supporting dataset to encourage further ASD-related research in eye-tracking.
- Wan et al. [4] applied eye-tracking to identify ASD in children, highlighting its potential as a diagnostic aid.
- Elbattah et al. [5] extended this line of work by integrating data visualization with deep learning techniques for computer-aided screening.

The inclusion of these studies and the dataset ensures both the reproducibility and comparability of our proposed methodology with the state of the art.

References

- [1] M. Elbattah, “Visualization of eye-tracking scanpaths in autism spectrum disorder: Image dataset,” 2019, dataset.
- [2] R. Carette, M. Elbattah, F. Cilia, G. Dequen, J.-L. Guerin, and J. Bosche, “Learning to predict autism spectrum disorder based on the visual patterns of eye-tracking scanpaths,” in *Proceedings of the 11th International Conference on Agents and Artificial Intelligence (ICAART)*, 2019, pp. 103–112.
- [3] F. Cilia, R. Carette, M. Elbattah, and G. Dequen, “Eye-tracking dataset to support the research on autism spectrum disorder,” 2019.
- [4] G. Wan, X.-J. Kong, B. Sun, S. Yu, Y. Tu, J. Park, C. Lang, M. Koh, Z. Wei, Z. Feng, Y. Lin, and J. Kong, “Applying eye tracking to identify autism spectrum disorder in children,” 2019.
- [5] M. Elbattah, R. Carette, F. Cilia, G. Dequen, J.-L. Guerin, and J. Bosche, “Computer-aided screening of autism spectrum disorder using eye-tracking, data visualization and deep learning,” 2020.