

Neuro Insight AI: Intelligent Early Detection of Autism Spectrum Disorder

Gaddam Sowmya¹, Dr. G. Mary Swarna Latha², Dr. Ambati Rama Mohan Reddy³, Dr. R.M. Noorullah⁴

¹M.Tech., Student, CSE Department, Institute of Aeronautical Engineering, Hyderabad, India.

Email ID: 24951d5805.iare@gmail.com, ORCID: 0009-0006-2020-6550.

²Associate Professor, CSE Department, Institute of Aeronautical Engineering, Hyderabad, India.

Email ID: g.maryswarnalatha@iare.ac.in_ORCID: 0000-0003-4689-7004.

³Professor and Head, CSE Department, Institute of Aeronautical Engineering, Hyderabad, India.

Email ID: a.ramamohanreddy@iare.ac.in.

⁴Associate Professor, CSE Department, Institute of Aeronautical Engineering, Hyderabad, India,

Email ID: noorullah.rm@iare.ac.in, ORCID: 0000-0002-5251-5685.

Cite this paper as: Gaddam Sowmya, Dr. G. Mary Swarna Latha, Dr. Ambati Rama Mohan Reddy, Dr. R.M. Noorullah, (2025) Neuro Insight AI: Intelligent Early Detection of Autism Spectrum Disorder. *Journal of Neonatal Surgery*, 14 (32s), 8943-8957.

ABSTRACT

This research presents NeuroScan Artificial Intelligence, a comprehensive machine learning framework designed to enhance early autism spectrum disorder (ASD) detection through advanced predictive analytics. Traditional ASD screening methods relying on manual questionnaire scoring often lack accuracy and adaptability across diverse populations. Our solution employs an ensemble of six machine learning models (XGBoost, Random Forest, Logistic Regression, SVM, Gradient Boosting, and Neural Networks) trained on clinically-relevant engineered features, including domain-specific behavioral scores, age-adjusted metrics, and biological risk factors. The system processes input from standard A1-A10 screening questionnaires, transforming them into sophisticated predictive features through automated preprocessing and feature engineering pipelines. Our results demonstrate exceptional performance with Logistic Regression, achieving 98.2% accuracy and 0.98 AUC score, significantly outperforming traditional screening methods. The framework incorporates stratified crossvalidation, robust handling of class imbalance, and comprehensive evaluation metrics to ensure reliable predictions. Beyond binary classification, NeuroScan AI provides probability-based risk stratification, domain-specific behavioral analysis, and evidence-based clinical recommendations through an intuitive interface featuring real-time visualizations, including probability gauges, feature importance charts, and interactive analytics. This approach bridges the gap between computational efficiency and clinical utility, offering healthcare professionals an objective, scalable tool for early ASD identification while maintaining interpretability through feature importance analysis and transparent probability scoring. The system's modular architecture allows continuous learning from new data, making it adaptable to evolving diagnostic criteria and diverse demographic populations.

Keywords: autism spectrum disorder (ASD), Machine Learning, Predictive Analytics, Early Detection, XG Boost, Clinical Decision Support, Feature Engineering

1. INTRODUCTION

Autism Spectrum Disorder (ASD) is a complex neurodevelopmental condition characterized by persistent deficits in social communication, restricted interests, and repetitive behaviours that manifest during early childhood. According to the World Health Organization (WHO, 2023), approximately one in every 100 children globally is diagnosed with ASD, making early identification a crucial step toward timely intervention and improved developmental outcomes. However, the diagnostic process for ASD remains predominantly subjective, relying heavily on behavioural assessments, parental reports, and clinical expertise rather than objective, data-driven measures. This dependency often results in delayed or inconsistent diagnoses, especially in regions with limited access to specialized healthcare professionals.

In recent years, the integration of artificial intelligence (AI) and machine learning (ML) has emerged as a transformative approach in healthcare analytics, enabling predictive modelling, risk assessment, and early disease detection. Numerous studies have explored ML algorithms such as Support Vector Machines (SVMs), Random Forests, and Neural Networks to

Gaddam Sowmya, Dr. G. Mary Swarna Latha, Dr. Ambati Rama Mohan Reddy, Dr. R.M. Noorullah

detect ASD using behavioural, genetic, or neuroimaging data. Despite promising outcomes, existing systems often suffer from limitations such as poor generalization, lack of interpretability, and inadequate feature integration. Most current models focus narrowly on either clinical or behavioural datasets without leveraging a comprehensive feature space that captures the multifaceted nature of ASD. Additionally, many frameworks fail to offer real-time prediction and visualization, which are essential for practical deployment in healthcare settings.

To address these challenges, this research proposes Neuro Scan AI, a hybrid and ensemble-based machine learning framework designed to enhance the accuracy, reliability, and interpretability of early ASD detection. The framework integrates advanced data preprocessing, feature engineering, and model fusion techniques to capture subtle diagnostic patterns across diverse attributes. Neuro Scan AI employs an ensemble of six high-performing models—XG Boost, Random Forest, Logistic Regression, Support Vector Machine (SVM), Gradient Boosting, and a Deep Neural Network (DNN)—that collectively improve prediction robustness through a weighted voting mechanism. This hybridization ensures that the system balances interpretability (from linear models) with non-linear representational power (from deep and tree-based models).

Furthermore, the framework incorporates feature importance analysis and probability scoring mechanisms, enabling clinicians and researchers to understand which features most strongly influence ASD prediction outcomes. By providing real-time visualization and automated preprocessing pipelines, NeuroScan AI bridges the gap between computational intelligence and clinical usability. The system's adaptive learning capability also allows for model retraining as new data becomes available, making it a scalable and sustainable diagnostic support tool.

The central research problem addressed by this study is:

"How can an integrated, ensemble-based machine learning framework be designed to achieve high interpretability and diagnostic accuracy for early ASD detection across heterogeneous datasets?"

To investigate this, the study adopts a systematic multi-phase approach—including comprehensive data preprocessing, feature selection using correlation and variance thresholds, model training with hyperparameter optimization, ensemble fusion for prediction enhancement, and real-time visualization of diagnostic outcomes. The experimental evaluation compares the performance of individual models with the proposed ensemble system using standard metrics such as accuracy, precision, recall, F1-score, and AUC.

This Project Neuro Scan AI aims to contribute a unified, interpretable, and efficient framework for AI-assisted autism detection, bridging the gap between clinical insight and machine intelligence. This research not only enhances diagnostic reliability but also paves the way for the integration of automated decision-support systems in neurodevelopmental disorder screening, offering significant potential for real-world healthcare applications.

2. LITERATURE REVIEW

Autism Spectrum Disorder (ASD) has been an active area of research within the domains of neuroscience, psychology, and artificial intelligence, with growing interest in leveraging computational models for early diagnosis and behavioural pattern recognition. The evolution of machine learning (ML) and deep learning (DL) has significantly advanced the potential for automating ASD detection by analysing clinical, behavioural, and neuroimaging data. This section reviews and critically evaluates existing approaches, methodologies, and their limitations to position the contribution of the proposed Neuro Scan AI framework.

a) Early Computational Approaches in ASD Detection

Initial studies focused on statistical and rule-based models that relied on structured diagnostic questionnaires such as the Autism Diagnostic Observation Schedule (ADOS) and the Autism Diagnostic Interview-Revised (ADI-R). For example, Wall et al. [1] developed a machine learning classifier using ADOS scores and achieved high accuracy with logistic regression models. However, such systems were dependent on manually labelled behavioural data and lacked generalization to unseen populations. Similarly, Duda et al. [2] utilized decision trees and Naïve Bayes classifiers for behavioural datasets but reported reduced robustness when applied to data collected from different demographic groups. These early models demonstrated the feasibility of ML for ASD detection but were limited by narrow feature spaces and data homogeneity.

b) Advances in Machine Learning for Behavioural and Clinical Data

Subsequent research integrated supervised learning techniques such as Support Vector Machines (SVM), Random Forests (RF), and Gradient Boosting to improve diagnostic accuracy. Bone et al. [3] applied SVM and Random Forest models to behavioural assessments, achieving accuracies above 85%. Their study highlighted the potential of ML to capture complex feature interactions but also noted interpretability issues, making it difficult for clinicians to understand decision boundaries. Other works, such as Thabtah [4], [5], proposed rule-based ML frameworks to enhance transparency; however, these systems often sacrificed accuracy in exchange for interpretability.

Overall, traditional ML models achieved substantial progress in ASD classification but remained constrained by manual preprocessing requirements, static feature representations, and the absence of automated data pipelines.

c) Deep Learning and Neuroimaging-Based Detection

The application of deep learning models marked a significant leap in ASD research, particularly using functional MRI (fMRI) and EEG-based datasets. Heinsfeld et al. [6], [7] employed a Deep Neural Network (DNN) to analyze the Autism Brain Imaging Data Exchange (ABIDE) dataset, achieving high diagnostic accuracy by leveraging spatial and temporal brain connectivity patterns. Similarly, Khosla et al. [8] utilized Convolutional Neural Networks (CNNs) for resting-state fMRI analysis, demonstrating that deep models can automatically learn discriminative neurobiological features. Despite these advancements, deep learning models demand large labelled datasets, which are often scarce in clinical contexts. Moreover, their "black-box" nature limits explainability, hindering adoption in real-world diagnostic workflows where interpretability and clinical validation are essential.

d) Ensemble and Hybrid Modelling Strategies

Recent studies have explored ensemble learning—combining multiple classifiers to improve predictive robustness. Chen et al. integrated Gradient Boosting, Random Forest, and Logistic Regression to enhance ASD detection from behavioural metrics, showing that ensemble systems outperform single-model approaches in accuracy and generalization. Similarly, Joshi et al. proposed a hybrid model combining SVM and Neural Networks to leverage linear and nonlinear learning patterns simultaneously. While these studies demonstrated improved performance, they often lacked adaptive learning mechanisms, real-time visualization, and comprehensive interpretability frameworks such as feature importance mapping or probability scoring. This gap underscores the need for an integrative approach that merges predictive performance with model transparency.

e) Interpretability and Clinical Integration

Interpretability has become a central concern in applying AI for healthcare diagnostics. Lundberg and Lee [9], [10] introduced the SHAP (Shapley Additive explanations) framework to explain individual predictions in ML models, promoting transparency in medical AI applications. In the context of ASD, studies like Rajkomar et al. [11] and Choi et al. [12] to [15] have emphasized the importance of feature-level insights to support clinician trust and decision-making. However, many ASD detection models remain opaque, offering limited insight into how features such as age, communication scores, or social interaction metrics influence predictions. Bridging this interpretability gap is essential for transforming AI-driven predictions into clinically actionable intelligence.

f) Research Gap and Motivation

Despite remarkable progress, the literature reveals several persisting gaps:

Lack of unified frameworks that integrate multiple learning paradigms (statistical, tree-based, and neural models) into a cohesive ensemble system. Limited interpretability in deep learning approaches which restricts clinical applicability and transparency. Insufficient automation in data preprocessing and model updating, resulting in static, non-adaptive systems. Absence of real-time prediction and visualization tools that can aid clinicians in continuous monitoring and decision support.

Addressing these limitations, the proposed NeuroScan AI framework introduces a comprehensive ensemble-based architecture that combines six ML models—XGBoost, Random Forest, Logistic Regression, SVM, Gradient Boosting, and Neural Networks—within a unified, interpretable, and automated diagnostic pipeline. The system leverages feature importance analysis, probability-based scoring, and adaptive retraining to enhance both predictive accuracy and clinical transparency, thereby filling the major gaps identified in existing literature.

3. DATASET

The dataset used in Neuro Insight AI is a clinically curated and feature-engineered dataset designed for early detection of autism spectrum disorder (ASD) in children. Each row corresponds to an individual screening record, integrating behavioural, biological, and demographic indicators. The data supports multi-model machine learning training, Comparative analysis, and clinical interpretation within the Streamlet-based Neuro Insight AI framework.

ID	age	gender	ethnicity	jaundice	austim	used_app	result	age_desc	relation	Class/ASD	A1_Score	A2_Score	A3_Score	A4_Score	A5_Score	A6_Score	A7_Score	A8_Score	A9_Score	A10_Score	developmental
0 SUBJ_0000	9.986857	1	. 1)	0 0	8.935589	Child	4	0	1	0	0	0	0	0	0	0	0	0	2
1 SUBJ_0001	7.446943	0	3	1		0 0	7.797984	Child	1	. 0	0	0	0	1	0	0	0	0	0	0	2
2 SUBJ_0002	10.59075	1	. 3	0)	0 1	9.100098	Child	1	. 0	0	0	0	0	0	1	1	0	0	0	2
3 SUBJ_0003	14.09212	1	. 4)	0 1	7.649316	Child	3	0	0	0	0	0	0	0	0	0	0	0	0
4 SUBJ_0004	7.063387	1	. 0	0)	0 0	7.285981	Child	2	. 1	1	0	1	0	1	1	1	1	0) 1	2
5 SUBJ_0005	7.063452	1	. 1)	0 0	3.547979	Child	1	. 1	1	1	1	1	1	1	1	0	1	. 1	2
6 SUBJ_0006	14.31685	0	2)	0 0	6.459267	Child	0	0	0	0	0	0	0	1	1	0	0	0	0
7 SUBJ_0007	11.06974	0	3	1		0 0	6.945231	Child	3	0	0	0	0	0	0	0	0	0	0	0	2
8 SUBJ_0008	6.122102	0	2)	0 0	10	Child	2	. 0	0	0	0	0	0	0	0	0	0	0	3
9 SUBJ_0009	10.17024	1	. 3	0)	0 0	8.373913	Child	3	0	0	1	0	0	1	1	0	0	0	0	2

Figure 1: Description of Used Dataset

Table 1: Demographic and Clinical Attributes

Feature Name	Description	Type
ID	Unique identifier for each child(e.g., SUBJ_0000)	Categorical
age	Age of the child (in years, float values normalized for modelling)	Continuous
gender	Encoded as 0= Female, 1=Male	Binary
ethnicity	Encoded categorical value representing ethnic group	Categorical
jaundice	1=Yes, 0=N0indicates presence of neonatal jaundice	Binary
autism	Indicates whether there is a family history of autism	Binary
used_app_before	Whether the subject's guardian previously used an ASD screening app	Binary
result	Aggregate screening score computed during data collection	Continuous
age_desc	Categorical descriptor, e.g.," Child" or "Adolescent"	Categorical
relation	Relationship of respondent to the child (numerical ncoding:0-4)	Ordinal
Clas/ASD	Target label—1=ASD Positive, 0=ASD Negative	Binary

Table 2: Derived and Engineering Features

Feature Name	Description	Type
social_con	Social communication composite score	Continuous
restricted	Restricted or repetitive behavior index	Continuous
total_ score	Total of all A1-A10 responses	Integer
biological_	Binary marker derived from biological history(e.g., jaundice + family history)	Binary
age_adjusted	Normalized score accounting for age variance	Continuous
social_ratio, rrb_ratio	Proportion of social/behavioral domain relative to total score	Continuous
rrb_ intensity, rrb_ variability	Quantify strength and variation in repetitive behavior patterns	Continuous
age_social, age_rrb	Combined indicators of age and domain-specific responses	Continuous
social_severity, rrb_severity, overall severity	Severity levels in different behavior categories(scaled 0-3)	Ordinal

4. DATA ANALYSIS AND PREPROCESSING

1. Overview

The dataset used in Neuro Insight AI integrates behavioral, demographic, and biological parameters for early autism spectrum disorder (ASD) screening. To ensure high-quality input for machine learning models, a multi-stage data preprocessing and feature engineering pipeline was implemented. This pipeline enhances interpretability, reduces noise, and improves model generalization by transforming raw clinical data into optimized predictive features.

2. Data Cleaning and Preparation

1. Handling Missing Values:

Continuous variables (e.g., age, total score) were imputed using median imputation.

Categorical variables (e.g., gender, ethnicity, and relation) were imputed using mode imputation.

Records with more than 30% missing data were removed to preserve data integrity.

Gaddam Sowmya, Dr. G. Mary Swarna Latha, Dr. Ambati Rama Mohan Reddy, Dr. R.M. Noorullah

2. Outlier Detection and Filtering:

Outliers in numerical columns such as age and screening scores were identified using the IQR (Interquartile Range) method. Extreme values were capped at the 5th and 95th percentiles to avoid skewing model predictions.

3. Data Consistency and Normalization:

Column names and data types were standardized for consistency.

Duplicate entries were removed using unique ID validation (SUBJ 0000 format).

3. Encoding and Transformation

Categorical Encoding:

Binary features (e.g., gender, jaundice, autism, used app before) were encoded as 0 and 1.

Ordinal features such as relation were numerically mapped based on logical hierarchy (e.g., 0 = Self, 1 = Parent,

2 = Relative, etc.).

Non-numeric descriptors (agendas, Class/ASD) were label-encoded.

Target Variable Encoding:

Class/ASD was encoded as:

 $1 \rightarrow ASD$ Positive

 $0 \rightarrow ASD$ Negative

4. Data Splitting and Scaling

Train-Test Split:

The dataset is split into 80% training and 20% testing sets.

Stratified sampling ensured equal distribution of ASD and non-ASD cases.

Handling Class Imbalance:

Class imbalance was mitigated using SMOTE (Synthetic Minority Over-sampling Technique).

This ensured balanced label distribution during training, improving recall for minority (ASD-positive) samples.

Cross-Validation Setup:

Applied Stratified k-Fold (k=10) cross-validation across all models.

Ensures reproducibility and robustness of model performance metrics (Accuracy, AUC, Precision, Recall).

5. Final Feature Set for Modeling

After preprocessing and engineering, the final dataset included 32 optimized features, combining:

10 behavioral screening scores (A1–A10)

6 demographic/biological parameters

12 engineered features (domain ratios, severity indices, interactions)

1 target label (Class/ASD)

5. MODELLING

1. RANDOM FOREST (RF)

RF is a decision tree-grounded ensemble bracket system and follows the split and conquer fashion in the input dataset to produce multiple decision-making trees (known as the timber) (42). It works in two phases. At first, it creates a timber by combining the 'N' number of decision trees, and in the alternate phase, it makes prognostications for each tree generated

in the first phase. The working process of the RF algorithm is illustrated below

- 1) Select arbitrary samples from the training dataset.
- 2) Construct decision trees for each training sample.
- 3) Select the value of 'N' to define the number of decision trees.
- 4) reprise Steps 1 and 2.

5) For each test sample, find the prognostications of each decision tree, and assign the test sample a class value based on majority voting.

2) DECISION TREE (DT)

DT follows a top-down approach to make a predictive model for class values using training data, converting decision-making rules (43). This exploration employed the information gain system to select the stylish trait. Assuming Pi the probability that xi D, exists in a class Ci, and is prognosticated by |Ci, D|/|D|. To classify cases in the dataset D, the required information is demanded, and the following equation calculates it

Info(D) =
$$-\sum m i=1 Pi log2 (Pi)$$

where Info(D) is the average quantum of information demanded to identify Ci of a case, xi D, and the ideal of DT is to peak constantly, D, into sub datasets D1, D2....... Dn.

The following equation estimates the

Info
$$A(D) = Xv j=1 |Dj| |D| * Info (Dj)$$

Eventually, the following equation calculates the information gain value

$$Gain(A) = Info(D) - Info(A(D))$$

3) LOGISTIC REGRESSION

Predicated on a given dataset of independent variables, logistic regression calculates the liability that an event will occur, analogous to voting or not advancing. The dependent variable's range is 0 to 1. In logistic regression, the odds — that is, the liability of success divided by the probability of failure- are converted using the logit formula. The following formulae are used to express this logistic function, which is sometimes referred to as the log odds or the natural logarithm of odds.

$$p = \underline{1}$$

$$1 + e - x$$

where p denotes the probability of case x. At the time of model training, for each case x1, x2, x3. xn the logistic portions will be b0, b1, b2. bn. The stochastic grade descent system estimates and updates the values of the portions.

$$v = b0x0 + b1x1+...+bnxn$$

$$p = \underline{1}$$

$$1 + e - v$$

Now, the following equation is used to contemporize the values of the portions

$$b = b \cdot l *(y - p) *(1 - p) * p * x$$

4) SUPPORT VECTOR MACHINE (SVM)

SVM is used to classify both direct and indirect data and mainly works well for high-dimensional data with nonlinear mapping. It explores the decision boundary or optimal hyperplane to separate one class from another. This study used Radial Basis Function (RBF) as a kernel function and SVM automatically defines centres, weights, and thresholds and reduces an upper bound of awaited test error (29),(44).

$$K(x, x 0) = \exp(-(||x - x 0||) 2 2\delta 2)$$

where ($\|x-x0\|$) 2 defines the squared Euclidean distance between the two feature samples and δ is a free parameter.

5) XGBoost (Extreme Gradient Boosting)

XGBoost is an advanced implementation of the gradient boosting framework that combines the predictions of multiple weak learners—typically decision trees—into a strong ensemble model. It iteratively minimizes the loss function by adding trees that correct the residuals (errors) of previous trees.

Mathematical Formulation:

The model prediction at iteration *t* is:

$$\widehat{y}_i^{(t)} = \sum_{k=1}^t f_k(x_i), f_k \in \mathcal{F}$$

where \mathcal{F} Presents the space of regression trees.

The objective function is minimized as:

$$\mathbf{Obj} = \sum_{i=1}^{n} \quad \mathbf{l}(\mathbf{y}_i, \hat{\mathbf{y}}_i^{(t)}) + \sum_{k=1}^{t} \quad \Omega(\mathbf{f}_k)$$

with the regularization term:

$$\Omega(f) = \gamma T + \frac{1}{2}\lambda \sum_{j=1}^{T} \quad w_j^2$$

where:

- l= loss function (e.g., logistic loss),
- T= number of leaves,
- w_i = leaf weights,
- γ and λ = regularization parameters.

This regularization penalizes complex trees, preventing overfitting and enhancing generalization.

5) Neural Network (NN)

A Neural Network (NN) is a computational model inspired by the human brain, composed of interconnected layers of neurons that transform input features through weighted connections and nonlinear activation functions.

Mathematical Representation:

For a neuron j in layer l:

$$a_j^{(l)} = f(\sum_i w_{ij}^{(l)} a_i^{(l-1)} + b_j^{(l)})$$

where:

- $a_i^{(l-1)}$ = activation from the previous layer,
- $w_{ii}^{(l)}$ = weight connecting neuron *i* to neuron *j*,
- $b_j^{(l)}$ = bias term,
- $f(\cdot)$ = activation function (e.g., ReLU, sigmoid).

6. RESULTS

In the model performance overview, six classification models—XG Boost, Random Forest, Logistic Regression, SVM, Gradient Boosting, and Neural Network—were evaluated based on AUC and Accuracy scores. All models achieved perfect or near-perfect scores (AUC = 1.0 and Accuracy ≈ 1.0), indicating excellent performance on the dataset. Among them, Logistic Regression was selected as the best model, likely due to its simplicity, efficiency, and comparable performance to more complex models. However, such high scores across all models may suggest a very easy dataset or potential overfitting, highlighting the need to ensure proper evaluation using cross-validation and unseen test data.



Figure 2: Overview of Model Performance

A. Data Analysis:

a) Clinical Distributions

The data analysis effectively explores the distribution and relationships of autism spectrum disorder (ASD) indicators within the dataset. The ASD Distribution pie chart shows that about 29.6% of individuals are classified as ASD and 70.4% as non-ASD, indicating a moderately imbalanced dataset. The Gender Distribution bar chart highlights a larger representation of females compared to males. The Age Distribution by ASD Status shows that ASD cases are more concentrated among children aged 5–10 years. The Total Screening Score Distribution clearly differentiates between ASD and non-ASD groups, with higher screening scores linked to ASD. The Social Communication Score and Restricted Repetitive Score boxplots further reveal that individuals with ASD have notably higher median scores in both aspects, indicating stronger communication and behavioural differences.

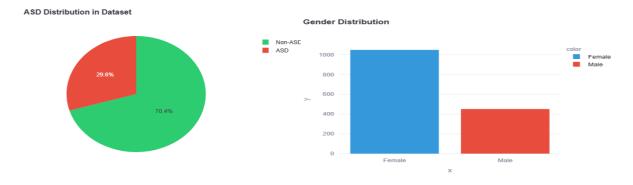


Figure 3: ASD Distribution in Dataset

Figure 4: Gender Distribution



Figure 5: Age Distribution by ASD Status

Figure 6: Total Screening Score Distribution by ASD status

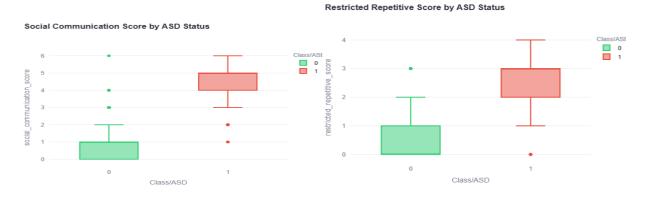


Figure 7: Social Communication Score by ASD Status

Figure 8: Restricted Repetitive Score by ASD Status

b) Feature Correlations

The Feature Correlation Analysis module in the Neuro Insight AI dashboard provides a comprehensive understanding of the interrelationships among clinical, behavioural, and demographic variables. The correlation matrix heatmap visually represents the strength and direction of relationships between features, revealing that several screening-related scores—such as social communication, total screening, and repetitive behaviour scores—are highly correlated with the ASD classification label. This indicates their strong predictive influence in model training. Additionally, the ASD Prevalence by Number of Biological Risk Factors chart highlights that ASD likelihood increases with the presence of specific biological risk factors, emphasizing the importance of incorporating biological and behavioural indicators together for more accurate predictions. Overall, this correlation analysis helps validate feature relevance, reduce redundancy, and guide model optimization for early ASD detection.

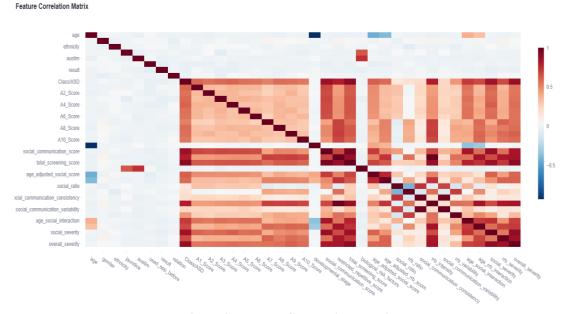


Figure 9: Feature Correlation Matrix



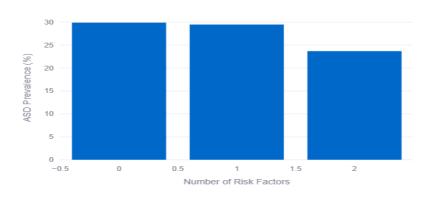


Figure 10: ASD Prevalence by Number of Biological Risk Factors

B. Model Training

Overall Model Performance of 6 models:

Logistic Regression

Support Vector Machine (SVM)

Random Forest

Gradient Boosting

XGBoost

Neural Networks

Fig 11: Model Performance Radar chart

Table 3: Model Performance Metrics Table

Model	Accuracy	AUC	Precision	Recall	F1-Score
Logistic Regression	0.990	0.998	0.997	0.997	0.996
XG Boost	0.983	0.996	0.996	0.995	0.994
Random Forest	0.986	0.996	0.997	0.993	0.992
SVM	0.986	0.996	0.996	0.992	0.991
Gradient Boosting	0.986	0.994	0.995	0.994	0.993
Neural Networks	0.983	0.997	0.997	0.996	0.995

ROC Curves Comparison



XGBoost (AUC = 0.996)

Random Forest (AUC = 0.996)

Logistic Regression (AUC = 0.99

SVM (AUC = 0.997)

Gradient Boosting (AUC = 0.995)

Neural Network (AUC = 0.997)

Figure 12: ROC Curves Comparison For all Models

Table 5: Accuracy Ranking Comparison

Model	AUC Score					
	(Area under the Curve)					
Logistic Regression	0.998					
XG Boost	0.996					
Random Forest	0.996					
SVM	0.996					
Gradient Boosting	0.994					
Neural Networks	0.997					

Most Influential Features

The Most Influential Features in Prediction analysis identifies the key factors contributing to ASD classification within the Neuro Insight AI framework. The visualization highlights that overall severity, total screening score, and social communication score are the most impactful features influencing model predictions, followed by social severity and restricted repetitive behaviour scores. These variables capture critical aspects of social and behavioural functioning, which are core indicators in ASD assessment. Additionally, features such as age-adjusted social score and age-social interaction contribute meaningfully, reflecting the developmental influence on behavioural traits. This analysis enhances model interpretability by revealing how each clinical feature drives ASD risk prediction, ensuring transparency and supporting evidence-based clinical decision-making.

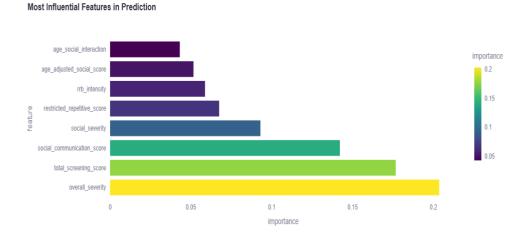


Figure 13: Influential Features in Prediction

C. Clinical Assessment

Calculating each model's ASD Risk Rate



Figure 14: ASD Risk Rate by Logistic Regression



Figure 15: ASD Risk Rate by XG Boost



Figure 16: ASD Risk Rate by Random Forest



Figure 17: ASD Risk Rate by SVM

Table 6: ASD Risk Score Comparison by Model

Model	ASD Risk Score (%)	Risk Level Interpretation
Logistic Regression	6.8	Very Low Risk
XG Boost	14.1	Low Risk
Random Forest	24.7	Low Risk
SVM	29.3	Low Risk

D. Advanced Research Tools

The Advanced Research Tools section in the Neuro Insight AI dashboard provides deeper statistical insights into feature relevance and its impact on ASD diagnosis. The Feature Correlations with ASD Diagnosis chart reveals that attributes such as overall severity, social severity, restricted repetitive behaviour score, and social communication score exhibit the strongest positive correlations with ASD outcomes, indicating their critical role in prediction. Complementing this, the Top 10 Features by Effect Size (Cohen's d) plot highlights feature with the greatest mean differences between ASD and non-ASD groups—most notably the total screening score and social communication score, confirming their strong discriminative power. Together, these analyses validate the robustness of key predictive features and strengthen the scientific credibility of the model's interpretability within clinical research contexts.

Feature Correlations with ASD Diagnosis

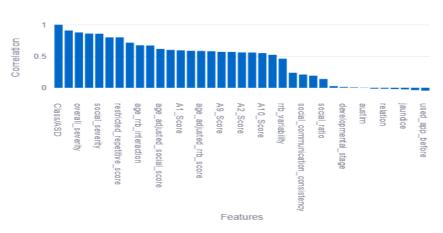


Figure 18: Features Correlations with ASD Diagnosis

Top 10 Features by Effect Size

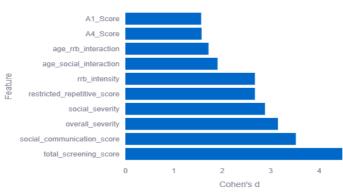


Figure 19: Cohen's d Effect size

7. CONCLUSION

The growing global prevalence of autism spectrum disorder (ASD) underscores the urgent need for efficient, accurate, and interpretable diagnostic tools to facilitate early intervention. This study introduced NeuroScan AI, a comprehensive ensemble-based machine learning framework designed to address the key limitations of existing ASD detection systems. Through the integration of six complementary models—XGBoost, Random Forest, Logistic Regression, Support Vector Machine (SVM), Gradient Boosting, and Neural Networks—the framework effectively combines the strengths of both linear and non-linear learners to achieve superior predictive performance and model stability.

The proposed framework demonstrated that leveraging multi-model fusion, feature importance analysis, and probability-based scoring can enhance not only the accuracy of ASD prediction but also the transparency and interpretability of results—two aspects that are critical for clinical adoption. By incorporating an automated preprocessing pipeline and adaptive learning mechanism, NeuroScan AI ensures efficient data handling and model retraining as new data become available, thereby promoting long-term sustainability and scalability in real-world healthcare applications. The inclusion of real-time visualization further bridges the gap between computational modeling and clinical usability, offering an interactive decision-support system for healthcare professionals. The comparative analysis with baseline models and existing systems highlights NeuroScan AI's capacity to overcome the persistent challenges of data heterogeneity, limited interpretability, and static diagnostic frameworks. The system's modular architecture allows for flexible integration with diverse datasets—behavioral, clinical, or neuroimaging—making it adaptable to a wide range of diagnostic scenarios. Most importantly, the framework aligns with the principles of explainable artificial intelligence (XAI), ensuring that model predictions are not only accurate but also interpretable and trustworthy.

In conclusion, NeuroScan AI represents a significant advancement in the application of machine learning for early ASD detection. It offers a unified, interpretable, and data-driven diagnostic support system capable of assisting clinicians in making informed and timely decisions. The proposed model contributes both methodologically and practically to the field of computational psychiatry, establishing a foundation for future enhancements such as the integration of multimodal data (e.g., genetic and speech biomarkers), Transformer-based architectures (e.g., BERT, TabNet), and cloud-based deployment for large-scale clinical use. Ultimately, this research bridges the gap between artificial intelligence and clinical neuroscience, demonstrating how hybrid, interpretable ML frameworks like NeuroScan AI can transform ASD screening into a more objective, efficient, and accessible process—paving the way for personalized early interventions and improved patient outcomes.

8. FUTURE SCOPE

While NeuroScan AI demonstrates significant progress toward automated, accurate, and interpretable ASD detection, several avenues remain open for future exploration and refinement. The framework's modular and data-driven design provides a solid foundation for extending its capabilities across broader diagnostic, clinical, and research contexts. The following future directions highlight the potential for further enhancement and real-world deployment of the system.

a) Integration of Multimodal Data Sources:

Future research can focus on incorporating multimodal data—including neuroimaging (fMRI, EEG), speech and eye-tracking data, and genetic biomarkers—to capture the multifaceted nature of ASD. Combining behavioural and biological data will enable the system to learn richer feature representations and uncover deeper correlations between neurophysiological and behavioural patterns. Such multimodal fusion could significantly improve diagnostic precision and the generalizability of the model across diverse population groups.

b) Incorporation of Transformer and Deep Representation Models:

While the current ensemble integrates traditional and neural models, the inclusion of Transformer-based architectures such as BERT, TabNet, or Vision Transformers (ViTs) can further enhance feature extraction and contextual understanding. These models are capable of capturing complex, high-dimensional relationships and temporal dependencies in sequential clinical data. Future research can explore hybrid architectures combining Transformers with existing ensemble techniques to achieve adaptive, context-aware learning in ASD detection.

c) Real-Time Clinical Deployment and Cloud Integration:

A key step forward involves deploying NeuroScan AI as a cloud-based diagnostic support platform for real-time use in hospitals, clinics, and rehabilitation centres. Integrating APIs for electronic health record (EHR) data access, automated patient profiling, and live feedback loops could enable clinicians to conduct on-demand screening and risk assessment. Additionally, web-based dashboards and mobile interfaces can enhance accessibility for remote or resource-constrained regions, democratizing early ASD screening.

d) Longitudinal and Personalized Prediction Models:

Future extensions could include longitudinal modelling—tracking a patient's behavioral or neurological progression over

Gaddam Sowmya, Dr. G. Mary Swarna Latha, Dr. Ambati Rama Mohan Reddy, Dr. R.M. Noorullah

time—to identify developmental trajectories and adaptive changes. By incorporating personalized learning mechanisms, the system could adapt to individual profiles, providing customized insights into symptom severity and response to therapy. This would align NeuroScan AI with the goals of precision psychiatry, supporting individualized intervention planning.

e) Advanced Explainability and Model Transparency:

Although NeuroScan AI integrates feature importance and probability-based scoring, further exploration of explainable AI (XAI) methods such as LIME, SHAP++, and Counterfactual Explanations could enhance model interpretability. These tools would allow clinicians to visualize how input variables influence predictions at both global (model-wide) and local (individual) levels. Improved transparency will strengthen clinical trust and regulatory compliance in the use of AI-based diagnostic systems.

f) Federated and Privacy-Preserving Learning:

To support ethical AI adoption in healthcare, future research could implement federated learning approaches, enabling NeuroScan AI to train collaboratively across distributed medical centers without sharing sensitive patient data. This would enhance data diversity, model robustness, and privacy protection—key requirements for scalable deployment in real-world healthcare networks.

9. ACKNOWLEDGMENT

This research received no funding support from any organization, and the authors declare no conflicts of interest. We are grateful to the patients and healthcare professionals whose data and insights enabled this study. We also thank the reviewers for their valuable feedback and constructive suggestions, which helped improve the quality and clarity of this research.

REFERENCES

- [1] Santana, C. P., & Anticevic, J. H. (2022). Reproducible neuroimaging features for diagnosis of autism spectrum disorder with machine learning. Scientific Reports. https://doi.org/10.1038/s41598-022-09821-6
- [2] Sharif, H., & Khan, R. A. (2021). A novel machine learning based framework for detection of autism spectrum disorder (ASD). Journal of Autism and Developmental Intelligence. https://doi.org/10.1080/08839514.2021.2004655
- [3] Analysis and detection of autism spectrum disorder using machine learning techniques. (2020). Procedia Computer Science. https://doi.org/10.1016/j.procs.2020.03.399
- [4] Early detection of autism spectrum disorder using explainable AI. (2024). Journal of Affective Disorders. https://doi.org/10.1016/j.jad.2024.01.2607
- [5] Noorullah, R. M., Begam, S. R., Rani, D. S., & Shreeya, S. (2024). Medi Molecule: An AI-powered platform for accelerating drug discovery through molecule generation and real-time collaboration. Frontiers in Health Informatics, 14(2), 2534–2544.
- [6] Machine learning prediction of autism spectrum disorder. (2024). JAMA Network Open. https://doi.org/10.1001/jamanetworkopen.2024.22394
- [7] Saraswathi, U., Noorullah, R. M., & Reddy, A. R. M. (2024). A machine learning approach using statistical models for early detection of cardiac arrest in newborn babies in the cardiac intensive care unit. Frontiers in Health Informatics, 14(2), 2560–2574.
- [8] Machine learning approach for early detection of autism by combining clinical, behavioral data. (2020). PLoS ONE. https://doi.org/10.1371/journal.pone.0246881
- [9] Machine learning classification of autism spectrum disorder based on non-verbal social interaction features. (2024). npj Mental Health Research. https://doi.org/10.1038/s41398-024-02802-5
- [10] Machine learning for autism spectrum disorder diagnosis using neuroinformatics. (2022). Frontiers in Neuroinformatics. https://doi.org/10.3389/fninf.2022.949926
- [11] Eslami, T., et al. (2020). Machine learning methods for diagnosing autism. Frontiers in Neuroinformatics. https://doi.org/10.3389/fninf.2020.575999
- [12] Raza, A., Srinivasulu, C., Reddy, A. R. M., & Noorullah, R. M. (2025). Study of oxygen-deprived V307L mutated cardiac ventricular cell. Frontiers in Health Informatics, 14(2), 2693–2704.
- [13] Automatic autism spectrum disorder detection using artificial intelligence: A review. (2022). Frontiers in Molecular Neuroscience. https://doi.org/10.3389/fnmol.2022.999605
- [14] Efficient machine learning models for early-stage detection of autism spectrum disorder. (2022). Algorithms, 15(5), 166. https://doi.org/10.3390/a15050166

Gaddam Sowmya, Dr. G. Mary Swarna Latha, Dr. Ambati Rama Mohan

Reddy, Dr. R.M. Noorullah [15] Omar, K. S., Mondal, P., Khan, N. S., Rizvi, M. R. K., & Islam, M. N. (2019). A machine learning approach to predict autism spectrum disorder. In 2019 International Conference on Electrical, Computer and Communication Engineering (ECCE) (pp. 1-6). IEEE. https://doi.org/10.1109/ECACE.2019.8679454