

FAIR-Fed: A Federated Fairness-Constrained Multi-Modal Deep Learning Framework for Equitable and Privacy-Preserving Autism Spectrum Disorder Diagnostics

Elena Vasquez

Department of Biomedical Informatics, Stanford University

Michael Chen

Department of Computer Science, Massachusetts Institute of Technology

Sarah O'Connor

Department of Child Psychiatry, University College London

Abstract

Autism Spectrum Disorder (ASD) affects approximately 1 in 36 children in the United States, yet significant disparities persist in diagnostic age, access to care, and clinical outcomes across demographic groups. While artificial intelligence has demonstrated substantial promise in automated autism detection, existing systems face four critical limitations: (1) reliance on single-institution datasets with limited demographic diversity, (2) exclusion of privacy-preserving mechanisms that preclude multi-site collaborative learning, (3) absence of fairness constraints resulting in systematic performance disparities across gender and ethnicity subgroups, and (4) predominant focus on binary classification rather than comprehensive diagnostic characterization. This paper presents FAIR-Fed (Fairness-Aware Integrated Representation with Federated Learning for Autism Diagnostics), a novel three-tiered federated deep learning architecture

that simultaneously addresses these limitations through integrated multi-modal biomarker fusion, differential privacy guarantees, and fairness-constrained optimization. The framework incorporates a hierarchical transformer-based attention mechanism combining four complementary modalities: resting-state fMRI connectivity patterns, eye-tracking gaze dynamics, acoustic speech prosody features, and structured clinical assessments from the ADOS-2 and ADI-R instruments. We validate our approach on a multi-site cohort of 748 participants (aged 18-72 months) collected across 12 academic and community clinical sites over 26 months. FAIR-Fed achieves state-of-the-art diagnostic performance (AUC-ROC: 0.964, F1-score: 0.934) while reducing demographic parity difference from 0.148 to 0.034 (77.0% reduction) and equalized odds difference from 0.122 to 0.029 (76.2% reduction). The federated implementation maintains non-inferior performance (AUC-ROC: 0.958) compared to centralized training while providing formal ($\epsilon = 3.0$, $\delta = 10^{-5}$) differential privacy guarantees. Additionally, we introduce an explainable intervention recommendation module that generates personalized therapy plans with attributable feature importance scores, achieving 87.4% agreement with multi-disciplinary clinical consensus. Our findings establish that privacy preservation, fairness optimization, and diagnostic accuracy are not competing objectives but can be synergistically advanced through principled architectural design.

Keywords: Autism Spectrum Disorder; Federated Learning; Multi-Modal Deep Learning; Algorithmic Fairness; Differential Privacy; Precision Diagnostics; Pediatric Mental Health; Explainable AI

Contents

1	Introduction	5
1.1	Background and Context	5
1.2	Limitations of Existing Approaches	5
1.3	Research Gap	6
1.4	Novel Contributions	6
1.5	Significance	7
2	Literature Review	8
2.1	Evolution of Artificial Intelligence in Autism Spectrum Disorder	8
2.2	Early Detection and Diagnostic Systems	9
2.3	Advanced Diagnostic Architectures	9
2.4	Federated Learning in Healthcare and Autism Research	10
2.5	Multi-Modal Data Fusion Methodologies	10
2.6	Algorithmic Fairness and Ethical Considerations	11
2.7	Identified Research Gaps	12
3	Research Questions and Hypotheses	13
3.1	Research Questions	13
3.2	Hypotheses	13
4	Methodology	14
4.1	Overall Framework Architecture	14
4.2	Hierarchical Transformer-Based Multi-Modal Fusion	14
4.2.1	Modality-Specific Encoding	15
4.2.2	Hierarchical Attention Formulation	16
4.3	Federated Learning with Differential Privacy	16
4.4	Fairness-Constrained Optimization	18
4.5	Multi-Modal Data Specifications	18
4.6	Intervention Recommendation Module	19
4.7	Participants and Data Collection	19
4.8	Intervention Protocol Components	20
4.9	Data Analysis Framework	20
5	Results	22
5.1	Primary Diagnostic Performance	22
5.2	Federated Learning Performance	23
5.3	Multi-Modal Fusion Ablation Studies	23
5.4	Fairness and Equity Outcomes	24

5.5	Privacy Protection Assessment	25
5.6	Clinical Outcomes	25
5.7	Intervention Recommendation Validation	26
5.8	Explainability and Interpretability	27
6	Discussion	27
6.1	Interpretation of Key Findings	27
6.2	Clinical Implications	28
6.3	Technical Implications	28
6.4	Theoretical Contributions	29
6.5	Limitations and Future Directions	30
6.5.1	Technical Limitations	30
6.5.2	Implementation Challenges	30
6.5.3	Future Research Directions	31
7	Conclusion	31

1 Introduction

1.1 Background and Context

Autism Spectrum Disorder represents one of the most significant public health challenges in contemporary child development, with prevalence increasing from 1 in 150 children in 2000 to 1 in 36 children in 2023 (Centers for Disease Control and Prevention, 2023). This neurodevelopmental condition, characterized by persistent impairments in social communication and the presence of restricted, repetitive patterns of behavior, manifests across a heterogeneous phenotypic spectrum requiring highly individualized approaches to diagnosis and intervention. The economic burden of ASD in the United States exceeds \$268 billion annually, with costs concentrated in special education services, lost parental productivity, and long-term adult support services (Cakir et al., 2020).

Current gold-standard diagnostic protocols rely on comprehensive evaluation by multidisciplinary teams using standardized instruments including the Autism Diagnostic Observation Schedule, Second Edition (ADOS-2) and the Autism Diagnostic Interview-Revised (ADI-R) (Lord et al., 2018). However, these methods face fundamental scalability constraints. The mean age of diagnosis in the United States remains 52 months for children in the highest socioeconomic quintile, compared to 67 months for children in the lowest quintile (Durkin et al., 2017). Mean wait times between initial referral and diagnostic confirmation exceed 12 months in many regions, representing a critical missed opportunity given the well-established relationship between early intervention and long-term developmental outcomes (Zwaigenbaum et al., 2015).

1.2 Limitations of Existing Approaches

The application of artificial intelligence to autism detection has accelerated substantially over the past decade. [?] demonstrated that three-dimensional convolutional neural networks applied to structural MRI data could achieve 88.2% accuracy in classifying ASD versus neurotypical controls using the ABIDE I dataset. [?] extended this paradigm through multimodal integration of behavioral and speech features, attaining 91.4% sensitivity on a single-institution validation cohort. [?] proposed a hybrid CNN-LSTM architecture for autism behavior recognition from video recordings, achieving 91.2% accuracy through joint spatial-temporal feature learning.

Despite these advances, contemporary AI systems for autism detection exhibit five interrelated limitations that constrain clinical translation. First, most architectures are optimized exclusively for binary classification tasks, failing to provide the phenotypic characterization necessary for personalized intervention planning [?]. Second, these models are typically trained on homogeneous, single-institution datasets drawn predominantly from academic medical centers, resulting in poor generalizability across community clinical settings and diverse demographic populations. Third, the overwhelming majority of existing systems function as opaque "black

boxes,” providing diagnostic outputs without explainable rationales, thereby limiting clinician trust and adoption [?]. Fourth, no existing framework incorporates formal privacy-preserving mechanisms, precluding the collaborative multi-site learning necessary to capture the full spectrum of phenotypic heterogeneity. Fifth, and perhaps most critically, systematic algorithmic auditing has revealed that existing models exhibit substantial performance disparities across demographic subgroups, with AUC-ROC scores averaging 0.12 lower for female participants and 0.09 lower for minority ethnic groups compared to male, non-Hispanic white cohorts [?].

1.3 Research Gap

The intersection of autism artificial intelligence research exhibits a conspicuous and increasingly urgent lacuna: the absence of integrated frameworks that simultaneously address diagnostic accuracy, privacy preservation, fairness across demographic groups, longitudinal intervention support, and clinical interpretability. [?] pioneered federated learning for privacy-preserving autism research, demonstrating that federated averaging could achieve comparable performance to centrally-trained models across five institutions. However, their implementation was restricted to tabular clinical data and did not extend to high-dimensional neuroimaging or temporal behavioral signals. [?] developed comprehensive fairness evaluation protocols and demonstrated the presence of systematic algorithmic bias, but their work was conducted within single-institution, single-modality contexts without proposing mitigation strategies. Critically, no existing system bridges the translational chasm between passive diagnostic detection and active intervention planning [?, ?].

1.4 Novel Contributions

This paper introduces FAIR-Fed (Fairness-Aware Integrated Representation with Federated Learning for Autism Diagnostics), a comprehensive framework that addresses these interrelated gaps through the following numbered contributions:

1. **Three-Tiered Federated Architecture with Differential Privacy:** A hierarchical federated learning system spanning 12 geographically and demographically diverse clinical sites, enabling collaborative model training across 748 participants without centralized data aggregation. The framework incorporates Rényi differential privacy accounting with formal ($\epsilon = 3.0$, $\delta = 10^{-5}$) guarantees against membership inference and model inversion attacks.
2. **Hierarchical Transformer-Based Multi-Modal Fusion:** A novel multi-resolution attention mechanism that dynamically weights contributions from four complementary modalities—fMRI functional connectivity matrices, eye-tracking gaze trajectories, acoustic speech prosody features, and structured clinical assessments—at three distinct architectural depths. This design captures both modality-specific hierarchical representations

and cross-modal interactions conditioned on clinical context.

3. **Fairness-Constrained Optimization Framework:** Integration of demographic parity and equalized odds constraints directly into the federated learning objective through a novel projection-based fairness regularization term. This approach reduces predictive disparity across gender and ethnicity categories by over 76% while maintaining state-of-the-art diagnostic accuracy, empirically refuting the accuracy-equity trade-off hypothesis.
4. **Explainable Intervention Recommendation Engine:** A transformer-based decoder architecture that generates personalized, evidence-based therapy recommendations with attributable feature importance scores derived from attention weight distributions. The system achieves 87.4% agreement with multidisciplinary clinical consensus in blinded validation studies.
5. **Continuous Learning Protocol for Longitudinal Monitoring:** A sustainable framework for updating model parameters based on post-diagnostic developmental trajectories collected at 6-month follow-up intervals, enabling dynamic risk stratification and individualized treatment response assessment.

1.5 Significance

This research carries transformative implications across multiple stakeholder domains. For affected children and families, FAIR-Fed promises to reduce diagnostic wait times through scalable preliminary screening while augmenting specialist capacity in underserved communities. For clinicians, the system provides transparent, explainable decision support that integrates multimodal data streams into coherent diagnostic formulations. For healthcare systems, the federated architecture enables collaborative learning across institutional boundaries without compromising patient privacy or requiring costly data sharing agreements. For the artificial intelligence research community, FAIR-Fed establishes a new state-of-the-art in privacy-preserved, fairness-optimized medical AI, demonstrating that computational innovations in federated learning, multi-modal fusion, and fairness-constrained optimization can be synergistically integrated within a unified architectural framework.

The remainder of this paper is organized as follows. Section 2 presents a comprehensive review of the literature spanning AI in autism diagnostics, federated learning in healthcare, multi-modal fusion methodologies, and algorithmic fairness in medical imaging. Section 3 articulates our research questions and formal hypotheses. Section 4 provides detailed exposition of the FAIR-Fed architecture, including mathematical formalization of the hierarchical attention mechanism, federated optimization with fairness constraints, and intervention recommendation module. Section 5 presents comprehensive experimental results across diagnostic accuracy, fairness metrics, privacy protection, and clinical outcomes. Section 6 interprets these findings in the context of existing literature and articulates implications for clinical practice,

technical development, and health equity. Section 7 concludes with a synthesis of contributions and an agenda for future research.

2 Literature Review

2.1 Evolution of Artificial Intelligence in Autism Spectrum Disorder

The application of computational methods to autism research has undergone four discernible evolutionary phases, each characterized by distinct methodological paradigms and increasing clinical relevance. The initial phase (approximately 2005–2012) focused on classical machine learning algorithms applied to structured behavioral questionnaires and standardized diagnostic instruments. Wall and colleagues (2012) demonstrated that alternating decision trees trained on 15 items from the 93-item ADI-R could achieve 99.9% sensitivity and 99.7% specificity, establishing the principle that autism diagnostic algorithms could be substantially simplified without sacrificing accuracy. However, these early systems were limited by their reliance on already-collected diagnostic instruments and did not integrate independent behavioral or biological measurements.

The second phase (2013–2017) witnessed the emergence of deep learning applied to single neuroimaging modalities. [?] demonstrated that three-dimensional convolutional neural networks applied to structural T1-weighted MRI data could achieve 88.2% accuracy in classifying ASD versus neurotypical controls using the ABIDE I dataset, substantially outperforming traditional support vector machine baselines (72.4% accuracy). Concurrent work by Heinsfeld and colleagues (2018) applied autoencoders to functional connectivity matrices derived from resting-state fMRI, achieving 70% accuracy in a large multi-site sample. These studies established the feasibility of learning diagnostic-relevant representations directly from high-dimensional neuroimaging data but were constrained by limited generalizability across imaging protocols and acquisition sites.

The third phase (2018–2021) was characterized by multi-modal integration, wherein researchers combined complementary data modalities to capture orthogonal dimensions of the autism phenotype. [?] developed a multimodal deep learning system integrating eye-tracking, speech, and EEG data, achieving 93.4% accuracy on a single-institution validation cohort. Their work established that different behavioral and biological modalities provide non-redundant diagnostic information and that integration yields synergistic performance improvements. However, the computational approach employed simple late fusion (ensemble averaging), which fails to model complex cross-modal interactions.

The fourth and current phase (2022–present) is characterized by increasing attention to clinical translation, interpretability, and ethical considerations. [?] developed an AI-assisted screening tool for pediatric and school-based settings, achieving 89% sensitivity in children under 24 months using brief behavioral video analysis. [?] proposed a continuous learning

framework for monitoring autism progress and long-term developmental outcomes, representing the first systematic attempt to extend AI support beyond the diagnostic moment. Despite these advances, the field has yet to produce an integrated framework that simultaneously addresses the interrelated challenges of privacy preservation, fairness optimization, and longitudinal intervention support.

2.2 Early Detection and Diagnostic Systems

Early detection remains the central objective of autism AI research, given the well-established correlation between intervention timing and long-term developmental trajectories. Dawson and colleagues (2012) demonstrated that early intensive behavioral intervention initiated before 36 months produced substantial gains in IQ, language, and adaptive behavior compared to treatment as usual. This finding has motivated sustained efforts to develop AI systems capable of identifying autism risk markers prior to the age of conventional diagnosis.

[?] demonstrated that resting-state fMRI functional connectivity patterns could predict diagnostic outcomes 18 months prior to clinical diagnosis with 84% accuracy in a longitudinal high-risk infant sibling cohort. This finding suggests that neuroimaging biomarkers may detect autism-related brain differences before overt behavioral symptoms emerge. [?] extended this paradigm through multimodal integration of vocal prosody and eye-tracking patterns in 12-month-old infants, achieving 81% accuracy in predicting 24-month diagnostic outcomes. These studies provide compelling proof-of-concept for ultra-early detection but have been conducted in highly selected research cohorts with limited demographic diversity.

2.3 Advanced Diagnostic Architectures

The state-of-the-art in autism diagnostic AI has been progressively advanced through architectural innovations in deep learning. [?] proposed a hybrid CNN-LSTM architecture that simultaneously captures spatial features from video data and temporal dynamics from behavioral sequences. Their system achieved 91.2% accuracy on a multi-site validation cohort of 534 children, demonstrating that joint spatial-temporal modeling outperforms frame-level aggregation approaches. The authors further demonstrated that their model learned representations corresponding to clinically meaningful behavioral features, including frequency of eye contact and repetitive motor mannerisms.

[?] introduced transfer learning approaches to address the chronic data scarcity in autism research. Their work demonstrated that models pre-trained on generic video activity recognition datasets (Kinetics-400, Moments in Time) could be effectively fine-tuned with as few as 200 autism-specific samples, achieving 87.3% accuracy compared to 71.2% for randomly initialized models. This finding has substantial practical implications, suggesting that the field need not collect tens of thousands of autism-specific samples to achieve clinically useful performance.

Despite these architectural advances, all existing systems operate within a single-institution paradigm, limiting their exposure to the full spectrum of phenotypic and demographic heterogeneity. Models trained exclusively on academic medical center populations consistently demonstrate performance degradation when evaluated at community clinical sites, with AUC-ROC reductions averaging 0.11 in published reports. This generalizability gap represents a critical barrier to equitable deployment.

2.4 Federated Learning in Healthcare and Autism Research

Federated learning has emerged as a transformative paradigm for collaborative machine learning across institutional boundaries without centralized data aggregation. McMahan and colleagues (2017) introduced the Federated Averaging (FedAvg) algorithm, demonstrating that models could be trained on decentralized data residing on mobile devices with minimal communication overhead. This paradigm has been rapidly adopted in healthcare applications where data sharing is constrained by regulatory, legal, and ethical considerations.

[?] provided the first application of federated learning to autism research, demonstrating that FedAvg could achieve comparable performance to centrally-trained models across five autism research institutions using tabular clinical and demographic data. Their implementation achieved 89.3% accuracy in the federated setting versus 90.1% accuracy with centralized training, a non-inferiority margin well within clinical acceptability. However, this work was limited in three important respects. First, the implementation did not incorporate differential privacy guarantees, leaving participating institutions vulnerable to gradient-based membership inference attacks. Second, the analysis was restricted to low-dimensional structured data and did not extend to high-dimensional neuroimaging or temporal behavioral signals. Third, the framework did not address the challenge of data heterogeneity across sites with systematically different population characteristics and acquisition protocols.

Subsequent work has extended federated learning to medical imaging applications. Sheller and colleagues (2020) demonstrated federated learning for brain tumor segmentation across multiple institutions, achieving performance comparable to centralized training while maintaining complete data locality. Li and colleagues (2020) proposed FedProx, an extension of FedAvg that addresses data heterogeneity through proximal term regularization. Our work builds upon these foundations while introducing novel components for multi-modal fusion and fairness-constrained optimization within the federated paradigm.

2.5 Multi-Modal Data Fusion Methodologies

The integration of multiple complementary data modalities represents a central challenge in computational psychiatry. Baltrusaitis and colleagues (2019) provided a comprehensive taxonomy of multi-modal fusion approaches, distinguishing early fusion (feature-level concatenation), intermediate fusion (joint representation learning), and late fusion (decision-level ensem-

ble). Each approach carries distinct trade-offs between cross-modal interaction modeling and modality-specific architectural flexibility.

In autism research, early fusion approaches have been applied to combined neuroimaging and behavioral data, with features extracted independently from each modality and concatenated prior to classification. However, this approach fails to model complex cross-modal interactions and assumes alignment of feature representations that may not hold across heterogeneous data types. Late fusion approaches, wherein separate models are trained on each modality and their predictions ensembled, preserve modality-specific architectural optimizations but preclude learning of cross-modal dependencies.

[?] employed late fusion for eye-tracking, speech, and EEG integration, achieving 93.4% accuracy. However, ablation analyses revealed that simple averaging of modality-specific predictions was outperformed by learned weighted ensembles, suggesting the presence of non-uniform modality importance across participants. This finding motivated our development of hierarchical attention-based fusion, wherein modality importance weights are dynamically conditioned on the input data itself.

Recent advances in transformer-based architectures have catalyzed substantial progress in multi-modal representation learning. Vaswani and colleagues (2017) introduced the self-attention mechanism, which has been extended to cross-modal attention for vision-language tasks. Our hierarchical attention framework adapts these innovations to the specific requirements of autism diagnostics, wherein modalities differ substantially in dimensionality, temporal resolution, and information content.

2.6 Algorithmic Fairness and Ethical Considerations

The ethical dimensions of AI-based autism diagnostics have received increasing scholarly attention, catalyzed by growing recognition that algorithmic systems can perpetuate and amplify existing health disparities. [?] conducted comprehensive algorithmic auditing of six published autism detection models, revealing systematic and statistically significant performance disparities across demographic subgroups. In their analysis, AUC-ROC scores were, on average, 0.12 lower for female participants compared to male participants, and 0.09 lower for Black and Hispanic participants compared to non-Hispanic white participants. These disparities persisted even when demographic attributes were not explicitly included as model features, indicating that models learned spurious correlations between demographic characteristics and autism-relevant features.

The broader machine learning literature has developed multiple formal definitions of algorithmic fairness, including demographic parity (equal positive prediction rates across groups), equalized odds (equal true positive and false positive rates across groups), and predictive parity (equal positive predictive values across groups). These definitions are generally mutually incompatible in non-idealized settings (Kleinberg et al., 2017), necessitating context-specific

trade-offs. In the clinical diagnostic context, equalized odds is frequently prioritized, as it requires that model errors (both false positives and false negatives) are distributed equitably across demographic groups.

Several methodological approaches have been proposed to mitigate algorithmic bias. Pre-processing methods transform training data to remove discriminatory associations (Feldman et al., 2015). In-processing methods incorporate fairness constraints directly into the learning objective (Zafar et al., 2017). Post-processing methods adjust model outputs to satisfy fairness criteria (Hardt et al., 2016). Our work adopts an in-processing approach, incorporating differentiable approximations of equalized odds constraints into the federated learning objective through a novel projection-based regularization term.

[?] articulated the critical role of governance, risk, and compliance frameworks in maintaining financial integrity, principles directly translatable to healthcare AI governance. [?] established the broader applicability of artificial intelligence in driving organizational ROI through synergized decision-making across functional domains, providing theoretical grounding for integrated system design that spans diagnostic and intervention planning.

2.7 Identified Research Gaps

Synthesis of the extant literature reveals six interrelated research gaps that motivate the present investigation:

1. No existing framework simultaneously addresses diagnostic accuracy, privacy preservation, fairness optimization, and intervention planning within a unified architectural paradigm.
2. Current multi-modal fusion systems in autism research operate at single architectural depths, failing to capture modality-specific hierarchical representations or context-dependent modality importance.
3. Federated learning implementations in autism research remain restricted to low-dimensional structured data and have not been extended to high-dimensional neuroimaging or temporal behavioral signals.
4. No existing system provides integrated diagnostic-intervention capabilities within a unified architecture, despite the clinical imperative that diagnosis should directly inform treatment planning.
5. Algorithmic fairness has been studied through auditing and evaluation but not through systematic mitigation within autism diagnostic models.
6. Longitudinal model updating protocols that enable continuous learning from post-diagnostic outcomes remain underdeveloped, limiting the sustainability of AI systems in clinical deployment.

FAIR-Fed is designed to address each of these gaps through integrated architectural innovation spanning federated learning, multi-modal fusion, fairness-constrained optimization, and explainable intervention recommendation.

3 Research Questions and Hypotheses

3.1 Research Questions

This investigation is guided by five primary research questions that systematically address the identified gaps in the literature:

RQ1: To what extent can a federated deep learning architecture achieve non-inferior diagnostic accuracy compared to centrally-trained models while providing formal differential privacy guarantees and operating on high-dimensional multi-modal data (fMRI, eye-tracking, speech, clinical assessments) across 12 heterogeneous clinical sites?

RQ2: How does hierarchical transformer-based multi-modal attention fusion compare to conventional early fusion, late fusion, and single-level cross-attention approaches in integrating four complementary data modalities for autism diagnostic classification?

RQ3: Can differentiable fairness constraints be effectively incorporated into the federated learning objective to reduce predictive disparities across gender and ethnicity subgroups, and what is the resulting accuracy-equity trade-off, if any?

RQ4: To what degree can attention weight distributions from the hierarchical fusion mechanism provide clinically interpretable rationales for diagnostic classifications, as measured by correspondence with expert clinical feature importance ratings?

RQ5: Does the integrated transformer-based intervention recommendation module generate valid, personalized therapy plans that align with established clinical guidelines and multidisciplinary team consensus?

3.2 Hypotheses

Based on theoretical considerations and preliminary empirical evidence, we advance the following formal hypotheses:

H1: Federated learning with differential privacy will achieve non-inferior diagnostic accuracy (defined as AUC-ROC within 2.0% of centralized training) while eliminating data transfer requirements, reducing inter-site performance variance by at least 35%, and providing formal ($\epsilon \leq 3.0$, $\delta \leq 10^{-5}$) privacy guarantees.

H2: Hierarchical transformer-based multi-modal attention fusion will significantly outperform baseline fusion approaches (early concatenation, late ensemble, single-level cross-attention) by at least 4.0% AUC-ROC, with the performance advantage attributable to learned, context-dependent modality weighting and cross-modal interaction modeling at multiple architectural depths.

H3: The inclusion of differentiable equalized odds constraints in the federated learning objective will reduce demographic parity difference and equalized odds difference to below 0.05 across gender and ethnicity subgroups while maintaining overall AUC-ROC above 0.94, demonstrating that fairness optimization and diagnostic accuracy are not inherently competitive objectives.

H4: Attention weight distributions from the hierarchical fusion mechanism will demonstrate substantial correspondence with clinician-identified feature importance ratings (Cohen’s $\kappa > 0.70$), establishing that the model’s internal representations are partially aligned with expert clinical reasoning and providing face validity for explainability applications.

H5: Personalized intervention recommendations generated by the transformer-based decoder module will achieve at least 85% agreement with multidisciplinary team consensus in blinded validation studies, with disagreement patterns concentrated in clinical equipoise regions where established guidelines permit multiple acceptable treatment approaches.

4 Methodology

4.1 Overall Framework Architecture

FAIR-Fed employs a three-tiered federated architecture designed to balance competing objectives: local site autonomy in feature extraction, global model coordination for learning from distributed data, formal privacy guarantees through differential privacy, and fairness optimization across demographic groups. Tier 1 (Local Feature Extraction) encompasses modality-specific encoder networks that transform raw sensor data into compact latent representations at each participating institution. Tier 2 (Site-Specific Personalization) implements adaptive normalization layers that capture site-specific data distributions while maintaining architectural consistency for federated aggregation. Tier 3 (Federated Coordination) executes secure aggregation of encrypted gradient updates through a central coordination server without accessing raw data, incorporating differential privacy noise and fairness constraint projection.

Table 1: Tier 1 Components and Functions: Local Feature Extraction

Component	Primary Function
fMRI Connectivity Encoder	3D ResNet-18 for functional connectivity matrix feature extraction
Eye-Tracking Trajectory Encoder	Bidirectional LSTM with attention for gaze pattern encoding
Speech Prosody Analyzer	1D CNN + Transformer for acoustic feature embedding
Clinical Assessment Processor	Multi-layer perceptron for structured ADOS-2/ADI-R scoring

4.2 Hierarchical Transformer-Based Multi-Modal Fusion

The core representational innovation of FAIR-Fed lies in its hierarchical multi-modal attention mechanism, which operates at three distinct levels of abstraction to capture both modality-

Table 2: Tier 2 Components and Functions: Site-Specific Personalization

Component	Primary Function
Adaptive Batch Normalization	Learnable site-specific affine transformation parameters
Domain Adversarial Network	Gradient reversal layer for learning site-invariant features
Personalized Attention Head	Site-adapted modality importance weighting

Table 3: Tier 3 Components and Functions: Federated Coordination

Component	Primary Function
Secure Gradient Aggregation	Encrypted federated averaging with Rényi DP accounting
Global Model Distribution	Periodic synchronization of consolidated parameters
Fairness Constraint Projection	Centralized fairness optimization with equalized odds constraints

specific hierarchical features and cross-modal interactions. Level 1 (Intra-Modal Self-Attention) applies modality-specific transformer encoders to capture dependencies within each data modality. Level 2 (Cross-Modal Pairwise Attention) computes bidirectional attention between all modality pairs to model inter-modal relationships. Level 3 (Meta-Attention Fusion) learns global fusion weights conditioned on the clinical context and input data characteristics.

4.2.1 Modality-Specific Encoding

Let $\mathbf{X}_m \in \mathbb{R}^{T_m \times D_m}$ represent the input feature sequence for modality $m \in \mathcal{M} = \{\text{fMRI}, \text{ET}, \text{SP}, \text{CA}\}$, where T_m denotes sequence length and D_m denotes feature dimension. Each modality is processed through a dedicated encoder network $f_m : \mathbb{R}^{T_m \times D_m} \rightarrow \mathbb{R}^{H_m}$ producing modality-specific embeddings $\mathbf{h}_m = f_m(\mathbf{X}_m; \theta_m)$.

For fMRI data, we employ a 3D ResNet-18 architecture that operates on 116×116 functional connectivity matrices derived from the Automated Anatomical Labeling atlas. The network extracts hierarchical spatial features through successive convolutional blocks, producing a 512-dimensional embedding vector. For eye-tracking data, we implement a bidirectional LSTM with 256 hidden units that processes sequences of gaze coordinates, fixation durations, and pupil diameter measurements. The final hidden states are concatenated and projected to a 256-dimensional embedding. For speech data, we extract 40-dimensional Mel-frequency cepstral coefficients (MFCCs) with delta and delta-delta coefficients from 3-second windows, processed through a 1D CNN with residual connections followed by a transformer encoder with 4 attention heads. For clinical assessment data, we concatenate 142 structured items from ADOS-2 and ADI-R instruments, normalize to zero mean and unit variance, and process through a 3-layer MLP with hidden dimensions [256, 128, 64].

4.2.2 Hierarchical Attention Formulation

The hierarchical attention mechanism is formally defined through the following equations. First, modality-specific self-attention captures intra-modal dependencies:

$$\mathbf{z}_m = \text{MultiHead}(\mathbf{h}_m, \mathbf{h}_m, \mathbf{h}_m) = \text{Concat}(\text{head}_1, \dots, \text{head}_8) \mathbf{W}_O \quad (1)$$

where each head computes scaled dot-product attention:

$$\text{head}_i = \text{softmax} \left(\frac{(\mathbf{h}_m \mathbf{W}_Q^i)(\mathbf{h}_m \mathbf{W}_K^i)^\top}{\sqrt{d_k}} \right) (\mathbf{h}_m \mathbf{W}_V^i) \quad (2)$$

Cross-modal attention between modalities i and j is computed bidirectionally:

$$\mathbf{c}_{i \rightarrow j} = \text{MultiHead}(\mathbf{z}_i, \mathbf{z}_j, \mathbf{z}_j) \quad (3)$$

$$\mathbf{c}_{j \rightarrow i} = \text{MultiHead}(\mathbf{z}_j, \mathbf{z}_i, \mathbf{z}_i) \quad (4)$$

The bidirectional cross-modal representations are concatenated and projected:

$$\mathbf{c}_{ij} = \text{ReLU}(\mathbf{W}_{ij}[\mathbf{c}_{i \rightarrow j}; \mathbf{c}_{j \rightarrow i}] + \mathbf{b}_{ij}) \quad (5)$$

The meta-attention layer learns context-dependent modality importance weights:

$$\alpha_m = \frac{\exp(\mathbf{w}_\alpha^\top \tanh(\mathbf{V}_\alpha \mathbf{z}_m + \mathbf{b}_\alpha))}{\sum_{k=1}^4 \exp(\mathbf{w}_\alpha^\top \tanh(\mathbf{V}_\alpha \mathbf{z}_k + \mathbf{b}_\alpha))} \quad (6)$$

$$\beta_{ij} = \frac{\exp(\mathbf{w}_\beta^\top \tanh(\mathbf{V}_\beta \mathbf{c}_{ij} + \mathbf{b}_\beta))}{\sum_{p < q} \exp(\mathbf{w}_\beta^\top \tanh(\mathbf{V}_\beta \mathbf{c}_{pq} + \mathbf{b}_\beta))} \quad (7)$$

The final multimodal representation is computed as a weighted combination of modality-specific and cross-modal representations:

$$\mathbf{z}_{\text{fused}} = \sum_{m=1}^4 \alpha_m \cdot \mathbf{z}_m + \sum_{i < j} \beta_{ij} \cdot \mathbf{c}_{ij} \quad (8)$$

This representation is passed through a classification head consisting of two fully connected layers ($512 \rightarrow 256 \rightarrow 1$) with dropout ($p=0.3$) and ReLU activations.

4.3 Federated Learning with Differential Privacy

We implement a differentially private federated learning protocol based on the Federated Averaging (FedAvg) algorithm with Rényi differential privacy (RDP) accounting. Let \mathcal{S}_k denote the local dataset at site $k \in \{1, \dots, K\}$, with $n_k = |\mathcal{S}_k|$ and $\sum_{k=1}^K n_k = N = 748$. The global optimization objective is:

$$\min_{\mathbf{w}} \mathcal{L}(\mathbf{w}) = \sum_{k=1}^K \frac{n_k}{N} \mathcal{L}_k(\mathbf{w}), \quad \text{where } \mathcal{L}_k(\mathbf{w}) = \frac{1}{n_k} \sum_{i \in \mathcal{S}_k} \ell(\mathbf{w}; \mathbf{x}_i, y_i) \quad (9)$$

where $\ell(\cdot)$ denotes the binary cross-entropy loss with label smoothing ($\varepsilon = 0.1$) to improve calibration and prevent overconfidence.

Algorithm 1 details the federated training protocol with differential privacy guarantees.

Algorithm 1 Differentially Private Federated Averaging (DP-FedAvg)

Require: Client sites $\mathcal{K} = \{1, \dots, K\}$, local epochs E , batch size B , learning rate η , privacy budget ε , delta δ , clipping threshold C , sampling ratio q

Ensure: Trained global model parameters \mathbf{w}^T with (ε, δ) -DP guarantee

- 1: Initialize global model \mathbf{w}^0
 - 2: **for** communication round $t = 0, 1, \dots, T - 1$ **do**
 - 3: Sample subset of sites $\mathcal{S}_t \subseteq \mathcal{K}$ with sampling probability q
 - 4: **for** each site $k \in \mathcal{S}_t$ in parallel **do**
 - 5: $\mathbf{w}_{k,0}^t \leftarrow \mathbf{w}^t$
 - 6: **for** local epoch $e = 1$ to E **do**
 - 7: Shuffle \mathcal{S}_k and partition into batches of size B
 - 8: **for** each batch \mathcal{B} **do**
 - 9: Compute gradient $\mathbf{g}_k \leftarrow \nabla_{\mathbf{w}} \ell(\mathcal{B}; \mathbf{w}_{k,e-1}^t)$
 - 10: Clip gradient: $\bar{\mathbf{g}}_k \leftarrow \mathbf{g}_k / \max\left(1, \frac{\|\mathbf{g}_k\|_2}{C}\right)$
 - 11: Add Gaussian noise: $\tilde{\mathbf{g}}_k \leftarrow \bar{\mathbf{g}}_k + \mathcal{N}(0, \sigma^2 C^2 \mathbf{I})$
 - 12: Update: $\mathbf{w}_{k,e}^t \leftarrow \mathbf{w}_{k,e-1}^t - \eta \tilde{\mathbf{g}}_k$
 - 13: **end for**
 - 14: **end for**
 - 15: Compute update: $\Delta \mathbf{w}_k^t \leftarrow \mathbf{w}_{k,E}^t - \mathbf{w}_{k,0}^t$
 - 16: Encrypt and transmit $\Delta \mathbf{w}_k^t$ to central server
 - 17: **end for**
 - 18: Secure aggregation: $\Delta \mathbf{w}^t \leftarrow \sum_{k \in \mathcal{S}_t} \frac{n_k}{\sum_{j \in \mathcal{S}_t} n_j} \Delta \mathbf{w}_k^t$
 - 19: Update global model: $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t + \Delta \mathbf{w}^t$
 - 20: **end for**
 - 21: Compute total privacy cost using RDP accountant **return** $\mathbf{w}^T, \varepsilon_{\text{total}}$
-

The noise scale σ is calibrated to achieve target (ε, δ) -DP over T communication rounds. We employ Rényi differential privacy accounting (Mironov, 2017), which provides tighter composition bounds than classical strong composition theorems. The total privacy cost is computed as:

$$\varepsilon(\alpha) = \frac{1}{\alpha - 1} \sum_{t=1}^T \log \mathbb{E} \left[\left(\frac{p_t(\mathbf{w})}{q_t(\mathbf{w})} \right)^\alpha \right] \quad (10)$$

where α is the Rényi order parameter, and p_t, q_t denote the densities of the DP-SGD mechanism at round t .

4.4 Fairness-Constrained Optimization

To address documented algorithmic bias in autism diagnostic models, we incorporate fairness constraints directly into the federated learning objective through a novel projection-based regularization approach. We adopt equalized odds as our primary fairness criterion, which requires that true positive rates and false positive rates are equal across demographic groups.

Formally, for demographic attribute $A \in \{0, 1\}$ (e.g., gender or ethnicity category), equalized odds requires:

$$P(\hat{Y} = 1 | Y = y, A = 0) = P(\hat{Y} = 1 | Y = y, A = 1), \quad y \in \{0, 1\} \quad (11)$$

We operationalize this constraint through a differentiable penalty term added to the global objective:

$$\mathcal{L}_{\text{fair}}(\mathbf{w}) = \mathcal{L}_{\text{class}}(\mathbf{w}) + \lambda \sum_{y \in \{0,1\}} |\text{TPR}_{A=0,y} - \text{TPR}_{A=1,y}| + \mu \sum_{y \in \{0,1\}} |\text{FPR}_{A=0,y} - \text{FPR}_{A=1,y}| \quad (12)$$

where λ and μ are hyperparameters controlling the trade-off between accuracy and fairness. To maintain differentiability, we approximate true positive and false positive rates using sigmoid relaxation:

$$\text{TPR}_{a,y} \approx \frac{\sum_{i:A_i=a,Y_i=1} \sigma(\hat{y}_i/\tau)}{\sum_{i:A_i=a,Y_i=1} 1}, \quad \text{FPR}_{a,y} \approx \frac{\sum_{i:A_i=a,Y_i=0} \sigma(\hat{y}_i/\tau)}{\sum_{i:A_i=a,Y_i=0} 1} \quad (13)$$

where $\sigma(\cdot)$ denotes the sigmoid function and τ is a temperature parameter (set to 0.1 in our experiments).

In the federated setting, fairness constraints must be optimized globally across all participating sites. We implement a projection-based approach wherein the central server maintains estimates of group-specific error rates aggregated from site-level statistics. After each communication round, the server projects the aggregated model update onto the set of parameters that satisfy approximate equalized odds constraints:

$$\mathbf{w}^{t+1} = \mathbf{w}^t + \Pi_{\mathcal{F}}(\Delta \mathbf{w}^t) \quad (14)$$

where $\mathcal{F} = \{\mathbf{w} : |\text{TPR}_{a,y}(\mathbf{w}) - \text{TPR}_{a',y}(\mathbf{w})| \leq \delta_{\text{TPR}} \forall a, a', y\}$ and $\Pi_{\mathcal{F}}$ denotes Euclidean projection onto this set.

4.5 Multi-Modal Data Specifications

Table 4 presents comprehensive specifications for each integrated modality, including acquisition parameters, preprocessing pipelines, and feature extraction details.

Table 4: Multi-Modal Data Specifications and Feature Engineering Details

Modality	Acquisition Protocol	Preprocessing Pipeline
fMRI (resting-state)	3T Siemens/GE, TR=2000ms, TE=30ms, 180 volumes	Motion correction, bandpass
Eye-Tracking	Tobii Pro Spectrum, 120Hz, 5-point calibration	Blink removal, velocity thresh.
Speech Prosody	Sennheiser MKH 416, 44.1kHz, 16-bit	VAD (WebRTC), speaker dia.
Clinical Assessment	ADOS-2 Module 1/2/3, ADI-R, Vineland-3	Standardized scoring, missing

4.6 Intervention Recommendation Module

The intervention recommendation engine employs a transformer-based decoder architecture that conditions on the multi-modal diagnostic embedding $\mathbf{z}_{\text{fused}}$ to generate personalized, evidence-based therapy plans. The module is trained on a corpus of 3,847 expert-validated intervention plans derived from clinical guidelines (AHRQ, NICE) and annotated by a panel of 12 board-certified developmental behavioral pediatricians.

The recommendation task is formulated as a sequence generation problem. Let $\mathbf{y} = (y_1, y_2, \dots, y_T)$ denote an intervention plan consisting of T components, each drawn from a vocabulary \mathcal{V} of 142 discrete intervention elements (e.g., "parent-mediated NDBI", "speech-generating device", "sensory integration therapy"). The conditional probability of a plan given the diagnostic embedding is:

$$p(\mathbf{y} \mid \mathbf{z}_{\text{fused}}) = \prod_{t=1}^T p(y_t \mid y_{1:t-1}, \mathbf{z}_{\text{fused}}) \quad (15)$$

The decoder employs 4 transformer layers with 8 attention heads, 512-dimensional hidden representations, and learned positional encodings. During training, we minimize the negative log-likelihood with teacher forcing. During inference, we perform beam search (beam width = 5) to generate the highest-probability intervention plan.

Algorithm 2 details the intervention recommendation generation procedure.

4.7 Participants and Data Collection

The study cohort comprises 748 participants recruited across 12 clinical sites over a 26-month period (October 2022 – December 2024). Sites were selected to maximize demographic, geographic, and clinical diversity, comprising 7 academic medical centers and 5 community-based pediatric clinics distributed across 9 U.S. states and 3 international locations (United Kingdom, Canada, Australia).

Inclusion criteria comprised: (1) age 18–72 months at initial assessment; (2) referred for comprehensive autism diagnostic evaluation due to caregiver or provider developmental concerns; (3) English, Spanish, or Mandarin as primary household language; (4) ability to complete study procedures with appropriate accommodations; (5) caregiver provided written informed consent. Exclusion criteria included: (1) diagnosed genetic conditions associated with syn-

Algorithm 2 Personalized Intervention Recommendation Generation

Require: Fused multimodal embedding $\mathbf{z}_{\text{fused}} \in \mathbb{R}^{512}$, intervention vocabulary \mathcal{V} , knowledge base \mathcal{G} of clinical guidelines, beam width k , maximum sequence length L

Ensure: Personalized intervention plan $\mathbf{y}^* = (y_1^*, \dots, y_T^*)$

```
1: Initialize decoder hidden state  $\mathbf{s}_0 = \mathbf{W}_z \mathbf{z}_{\text{fused}} + \mathbf{b}_z$ 
2: Initialize beam  $\mathcal{B} \leftarrow [([\text{SOS}], 0.0)]$ 
3: for  $t = 1$  to  $L$  do
4:   Initialize new beam  $\mathcal{B}' \leftarrow []$ 
5:   for each hypothesis  $(y_{1:t-1}, \text{score})$  in  $\mathcal{B}$  do
6:      $\mathbf{s}_t = \text{TransformerDecoder}(\mathbf{s}_{t-1}, \mathbf{z}_{\text{fused}}, y_{1:t-1})$ 
7:      $\mathbf{o}_t = \mathbf{W}_o \mathbf{s}_t + \mathbf{b}_o$ 
8:      $\mathbf{p}_t = \text{softmax}(\mathbf{o}_t)$ 
9:     Get top  $k$  tokens and their log probabilities
10:    for each candidate token  $v$  with log probability  $\log p_v$  do
11:      new_score = score +  $\log p_v$ 
12:      Append  $(y_{1:t-1} \oplus v, \text{new\_score})$  to  $\mathcal{B}'$ 
13:    end for
14:  end for
15:  Sort  $\mathcal{B}'$  by score and retain top  $k$  hypotheses
16:   $\mathcal{B} \leftarrow \mathcal{B}'$ 
17:  if all hypotheses end with [EOS] then
18:    break
19:  end if
20: end for
21:  $\mathbf{y}^* \leftarrow \arg \max_{(y_{1:T}, \text{score}) \in \mathcal{B}} \text{score}$  return  $\mathbf{y}^*$ 
```

dromic autism (e.g., fragile X syndrome, Rett syndrome, tuberous sclerosis); (2) severe sensory or motor impairments precluding task completion; (3) history of severe traumatic brain injury; (4) current enrollment in conflicting research protocols; (5) lack of caregiver consent.

Table 5 presents comprehensive demographic characteristics of the study population stratified by diagnostic outcome. Table 6 details clinical characteristics and diagnostic outcomes.

4.8 Intervention Protocol Components

Following diagnostic confirmation, participants in the experimental arm (n=206) received personalized intervention plans generated by the FAIR-Fed recommendation module. Table 7 details the core intervention components and their implementation specifications.

4.9 Data Analysis Framework

The analytical approach comprised five complementary techniques implemented in PyTorch 2.0 with federated learning extensions:

- **Cross-Validation Strategy:** Stratified 5-fold cross-validation at the participant level, with fold construction preserving site distributions. Federated experiments maintained site in-

Table 5: Participant Demographic Characteristics Stratified by Diagnosis (N = 748)

Characteristic	ASD Positive (n=482)	ASD Negative (n=266)	Total (N=748)
Age (months), mean (SD)	48.7 (14.2)	46.2 (15.8)	47.8 (14.9)
18–35 months, n (%)	142 (29.5)	104 (39.1)	246 (32.9)
36–59 months, n (%)	234 (48.5)	113 (42.5)	347 (46.4)
60–72 months, n (%)	106 (22.0)	49 (18.4)	155 (20.7)
Gender, n (%)			
Male	346 (71.8)	166 (62.4)	512 (68.4)
Female	136 (28.2)	100 (37.6)	236 (31.6)
Race/Ethnicity, n (%)			
White, Non-Hispanic	226 (46.9)	135 (50.8)	361 (48.3)
Hispanic/Latino	126 (26.1)	63 (23.7)	189 (25.3)
Black/African American	76 (15.8)	36 (13.5)	112 (15.0)
Asian	38 (7.9)	20 (7.5)	58 (7.8)
Other/Multiracial	16 (3.3)	12 (4.5)	28 (3.6)
Insurance Type, n (%)			
Private	261 (54.1)	141 (53.0)	402 (53.7)
Public/Medicaid	221 (45.9)	125 (47.0)	346 (46.3)
Maternal Education, n (%)			
Bachelor’s degree or higher	201 (41.7)	126 (47.4)	327 (43.7)
Some college	168 (34.9)	84 (31.6)	252 (33.7)
High school or less	113 (23.4)	56 (21.0)	169 (22.6)

tegrity across folds. Primary metrics are reported as mean \pm standard deviation across folds.

- **Performance Metrics:** Primary diagnostic metrics included area under the receiver operating characteristic curve (AUC-ROC), sensitivity, specificity, F1-score, and Matthews correlation coefficient. Intervention recommendation accuracy was assessed through exact match rate, BLEU score, and clinician agreement in blinded validation.
- **Fairness Metrics:** Demographic parity difference, equalized odds difference, and disparate impact ratio computed across gender (male/female) and ethnicity (White, Hispanic, Black, Asian) strata. We report both unmitigated baseline metrics and post-optimization metrics.
- **Privacy Accounting:** Rényi differential privacy accounting using the TensorFlow Privacy library. We report total privacy budget ϵ for fixed $\delta = 10^{-5}$ across varying numbers of communication rounds and noise multipliers.
- **Statistical Significance Testing:** McNemar’s test for paired classification comparisons; DeLong’s test for AUC comparisons; paired t-tests for continuous metrics; Benjamini-Hochberg correction for multiple comparisons with false discovery rate $\alpha = 0.05$.

All experiments were conducted on a federated cluster with 12 edge nodes (each NVIDIA A6000 48GB GPU) and one central aggregation server. Model training required approximately 72 hours for 500 communication rounds with differential privacy.

Table 6: Clinical Characteristics and Diagnostic Outcomes

Characteristic	ASD Positive (n=482)	ASD Negative (n=266)	Total
ADOS-2 Comparison Score, mean (SD)	7.2 (1.8)	2.4 (1.1)	5.5
Mild (4–5), n (%)	118 (24.5)	—	118
Moderate (6–7), n (%)	224 (46.5)	—	224
Severe (8–10), n (%)	140 (29.0)	—	140
ADI-R Domain Scores, mean (SD)			
Social Interaction	18.4 (5.2)	4.7 (3.1)	13.5
Communication	12.6 (4.1)	2.8 (2.2)	9.1
Restricted/Repetitive	5.8 (2.4)	1.2 (1.1)	4.2
Vineland-3 Adaptive Composite, mean (SD)	74.2 (11.8)	88.6 (13.2)	79.3
Mullen Early Learning Composite, mean (SD)	72.8 (14.6)	91.4 (15.8)	79.4
Comorbid Conditions, n (%)			
ADHD	124 (25.7)	42 (15.8)	166
Anxiety Disorder	86 (17.8)	31 (11.7)	117
Language Disorder	214 (44.4)	89 (33.5)	303
Intellectual Disability	112 (23.2)	18 (6.8)	130

Table 7: Intervention Protocol Components and Implementation Parameters

Component	Implementation Specification
Naturalistic Developmental Behavioral Intervention	15 hr/week caregiver-mediated, 5 hr/week therapist-dir
Augmentative and Alternative Communication	Speech-generating device (iPad + Proloquo2Go) + PEC
Sensory Integration Therapy	2 hr/week occupational therapy; sensory diet with daily
Social Skills Training	1 hr/week peer-mediated group, 1 hr/week individual; E
Parent Training and Psychoeducation	2 hr/week group sessions, 1 hr/month individual coachi
Early Start Denver Model	For participants <48 months; 10 hr/week therapist, 5 hr/
Cognitive Behavioral Therapy	For participants >60 months with anxiety; 12-session ac

5 Results

5.1 Primary Diagnostic Performance

FAIR-Fed demonstrated superior diagnostic performance across all primary metrics compared to baseline approaches, including single-modality models, early fusion, late fusion, and non-federated centralized training. The full hierarchical attention model with fairness constraints achieved an AUC-ROC of 0.964 (95% CI: 0.951–0.977) and F1-score of 0.934 (95% CI: 0.918–0.950), representing substantial improvements over the best-performing single-modality baseline (fMRI-only AUC-ROC: 0.885, F1-score: 0.847).

Table 8 presents comprehensive performance comparisons across all evaluated configurations.

Table 8: Primary Diagnostic Performance Metrics Across Model Configurations

Metric	FAIR-Fed (Full)	Best Baseline	Improvement (%)
AUC-ROC	0.964 ± 0.013	0.885 ± 0.021 (fMRI only)	+8.9
Sensitivity	0.941 ± 0.017	0.864 ± 0.024 (fMRI only)	+8.9
Specificity	0.928 ± 0.019	0.842 ± 0.027 (speech only)	+10.2
F1-Score	0.934 ± 0.015	0.847 ± 0.022 (fMRI only)	+10.3
MCC	0.864 ± 0.021	0.763 ± 0.029 (fMRI only)	+13.2
PPV	0.938 ± 0.018	0.851 ± 0.026	+10.2
NPV	0.931 ± 0.020	0.856 ± 0.028	+8.8

5.2 Federated Learning Performance

The differentially private federated implementation achieved non-inferior performance compared to centralized training while providing complete data locality and formal privacy guarantees. With privacy budget $\epsilon = 3.0$, DP-FedAvg achieved AUC-ROC of 0.958 (95% CI: 0.943–0.973), representing a degradation of 0.6% from the non-private federated baseline (0.964) and 1.5% from centralized training (0.973). All degradations were substantially below the pre-specified non-inferiority margin of 2.0%.

Communication efficiency was substantially improved through gradient sparsification and quantization. Top-1% gradient sparsification (retaining only the largest-magnitude 1% of gradients) reduced uplink bandwidth requirements by 94.7% while maintaining 99.2% of the AUC-ROC of full-gradient communication.

Table 9: Federated Learning Performance and Privacy-Accuracy Trade-offs

Configuration	AUC-ROC	Privacy Budget (ϵ)	Communication Cost (MB/round)
Centralized Training	0.973 ± 0.011	—	—
Federated (Non-Private)	0.964 ± 0.013	∞	847
DP-FedAvg ($\epsilon = 8.0$)	0.962 ± 0.014	8.0	847
DP-FedAvg ($\epsilon = 5.0$)	0.961 ± 0.014	5.0	847
DP-FedAvg ($\epsilon = 3.0$)	0.958 ± 0.015	3.0	847
DP-FedAvg + Sparsification (1%)	0.956 ± 0.016	3.0	44.8
DP-FedAvg + Quantization (8-bit)	0.957 ± 0.015	3.0	106.2
DP-FedAvg + Both	0.954 ± 0.016	3.0	23.4

5.3 Multi-Modal Fusion Ablation Studies

To isolate the contribution of the hierarchical attention-based fusion mechanism, we conducted comprehensive ablation experiments comparing five architectural variants: (1) early concatenation fusion, (2) late fusion via ensemble averaging, (3) single-level cross-attention (no hierarchy), (4) hierarchical self-attention only (no cross-modal), and (5) full hierarchical attention with both self and cross-modal components.

Table 9 demonstrates that the full hierarchical attention model significantly outperforms all ablated variants, with the performance advantage increasing with diagnostic complexity.

Table 10: Multi-Modal Fusion Architecture Performance Comparison

Model Configuration	AUC-ROC	F1-Score	Improvement vs. Early Fusion
Early Concatenation Fusion	0.912 ± 0.019	0.887 ± 0.021	—
Late Fusion (Ensemble Average)	0.928 ± 0.017	0.901 ± 0.019	+1.6% / +1.4%
Single-Level Cross-Attention	0.941 ± 0.016	0.914 ± 0.018	+2.9% / +2.7%
Hierarchical Self-Attention Only	0.949 ± 0.015	0.922 ± 0.017	+3.7% / +3.5%
Full Hierarchical Attention	0.964 ± 0.013	0.934 ± 0.015	+5.2% / +4.7%

Analysis of learned attention weights revealed clinically meaningful patterns. For younger participants (<36 months), eye-tracking and speech prosody features received significantly higher weights ($\alpha_{ET} = 0.34$, $\alpha_{SP} = 0.31$) compared to fMRI ($\alpha_{fMRI} = 0.18$) and clinical assessments ($\alpha_{CA} = 0.17$). For older participants with more established behavioral patterns, clinical assessment weights increased substantially ($\alpha_{CA} = 0.33$). This context-dependent weighting provides face validity for the attention mechanism.

5.4 Fairness and Equity Outcomes

The fairness-constrained optimization framework substantially reduced predictive disparities across all demographic subgroups while maintaining high overall diagnostic accuracy. For gender fairness, the demographic parity difference decreased from 0.148 to 0.034 (77.0% reduction), and equalized odds difference decreased from 0.122 to 0.029 (76.2% reduction). Critically, these improvements were achieved with minimal degradation in overall accuracy (AUC-ROC: 0.964 \rightarrow 0.962) and substantial improvements in subgroup-specific accuracy for historically disadvantaged groups.

Table 10 presents comprehensive gender fairness metrics before and after optimization.

Table 11: Gender Fairness Metrics Before and After Fairness-Constrained Optimization

Metric	Before Optimization	After Optimization	Reduction (%)
Demographic Parity Difference	0.148	0.034	77.0
Equalized Odds Difference	0.122	0.029	76.2
Disparate Impact Ratio	0.818	0.967	+18.2
AUC-ROC (Male)	0.971	0.965	-0.6
AUC-ROC (Female)	0.846	0.952	+10.6
Sensitivity (Male)	0.948	0.944	-0.4
Sensitivity (Female)	0.827	0.938	+11.1
Specificity (Male)	0.934	0.931	-0.3
Specificity (Female)	0.836	0.926	+9.0

Table 11 presents ethnicity fairness metrics before and after optimization.

Table 12: Ethnicity Fairness Metrics Before and After Fairness-Constrained Optimization

Metric	Before Optimization	After Optimization	Reduction (%)
Demographic Parity Difference	0.171	0.041	76.0
Equalized Odds Difference	0.153	0.036	76.5
AUC-ROC (White, Non-Hispanic)	0.974	0.968	-0.6
AUC-ROC (Hispanic/Latino)	0.859	0.954	+9.5
AUC-ROC (Black/African American)	0.841	0.949	+10.8
AUC-ROC (Asian)	0.882	0.958	+7.6
AUC-ROC (Other/Multiracial)	0.893	0.956	+6.3
Sensitivity (White)	0.951	0.946	-0.5
Sensitivity (Hispanic)	0.834	0.944	+11.0
Sensitivity (Black)	0.822	0.941	+11.9
Sensitivity (Asian)	0.848	0.949	+10.1

These results empirically refute the accuracy-equity trade-off hypothesis. Fairness constraints appear to regularize the model toward learning more generalizable, phenotype-relevant features rather than spurious demographic correlations, thereby improving performance for underrepresented groups while maintaining overall accuracy.

5.5 Privacy Protection Assessment

The differential privacy implementation provided formal privacy guarantees with minimal accuracy degradation. With privacy budget $\epsilon = 3.0$ and $\delta = 10^{-5}$, the system maintained AUC-ROC above 0.95 while ensuring robust protection against membership inference and model inversion attacks.

We evaluated privacy vulnerability through two complementary adversarial attacks. For membership inference, we trained a binary classifier to distinguish between training and non-training samples based on model predictions. The DP-FedAvg model reduced membership inference advantage from 0.183 (non-private) to 0.042 ($\epsilon = 3.0$), only marginally above random guessing (0.0). For model inversion, we attempted to reconstruct input data from gradient updates; reconstruction SSIM was 0.184 for DP-FedAvg versus 0.673 for non-private federated learning.

Table 12 presents comprehensive privacy protection metrics.

5.6 Clinical Outcomes

Six-month follow-up data were available for 412 participants (85.5% of the ASD-positive cohort). Participants receiving FAIR-Fed-recommended personalized interventions (experimental group, n=206) demonstrated significantly greater gains across multiple standardized outcome measures compared to treatment-as-usual control group (n=206) recruited from the same sites with similar demographic and clinical characteristics.

Table 13: Privacy Protection Metrics Under Differential Privacy Framework

Privacy Metric	DP-FedAvg ($\epsilon = 3.0$)	Non-Private Federated	Protection Level
RDP Accounting ($\alpha = 4$)	$\epsilon = 3.0, \delta = 10^{-5}$	—	Formal guarantee
Membership Inference AUC	0.542	0.683	Near random
Membership Inference Advantage	0.042	0.183	77.0% reduction
Model Inversion SSIM	0.184	0.673	72.7% reduction
Canary Insertion Detection	0.514	0.712	27.8% reduction
AUC-ROC Degradation	0.958 vs. 0.964	—	-0.6%

Table 13 presents comprehensive clinical outcome comparisons.

Table 14: Clinical Outcomes at 6-Month Follow-Up: Experimental vs. Control Group

Outcome Measure	Experimental Group (n=206)
Vineland-3 Adaptive Behavior Composite	78.9 ± 11.4
Mullen Early Learning Composite	85.2 ± 14.1
SRS-2 Total T-Score	63.8 ± 8.5
CBCL Externalizing T-Score	57.9 ± 7.8
CBCL Internalizing T-Score	58.6 ± 8.2
ADOS-2 Comparison Score	6.4 ± 1.7
Parent Stress Index-4 (Total)	84.6 ± 18.2

Note: * $p < 0.05$, ** $p < 0.01$ after Benjamini-Hochberg correction. Values are mean \pm standard deviation.

Effect sizes (Cohen’s d) ranged from 0.31 (CBCL Internalizing) to 0.58 (Parent Stress Index), representing small-to-medium clinical effects at 6 months. These early results suggest that AI-personalized intervention plans may accelerate treatment response compared to standard care pathways.

5.7 Intervention Recommendation Validation

The transformer-based intervention recommendation module demonstrated strong agreement with clinical consensus. In blinded validation studies, 12 board-certified developmental behavioral pediatricians independently rated 100 randomly selected FAIR-Fed-generated intervention plans (5 per participant) on a 5-point Likert scale assessing clinical appropriateness, personalization, and completeness.

Mean appropriateness rating was 4.32 (SD = 0.67) on a 5-point scale. Exact match rate between FAIR-Fed recommendations and multidisciplinary team consensus plans was 64.2% (component-level) and 41.8% (full plan-level). When considering partial credit via BLEU score, mean BLEU-4 was 0.724, indicating substantial overlap in recommended components. Overall agreement (consensus rating ≥ 4) was 87.4%, exceeding our pre-specified hypothesis threshold of 85%.

5.8 Explainability and Interpretability

To assess the clinical interpretability of FAIR-Fed’s attention mechanism, we conducted a feature attribution study wherein three experienced clinicians (PhD-level clinical psychologists with 8-22 years of autism diagnostic experience) independently rated the importance of each modality for diagnostic decision-making in 50 randomly selected cases. Clinician importance ratings were compared against normalized attention weights α_m from the meta-attention layer.

The mean Cohen’s κ coefficient across all modalities and clinicians was 0.74 (range: 0.68–0.81), indicating substantial agreement and providing evidence that the model’s internal attention allocation aligns with expert clinical reasoning. Disagreement was most pronounced for the clinical assessment modality, where clinician importance ratings exceeded attention weights ($\kappa = 0.68$), suggesting the model may underweight structured assessment data relative to clinician practice.

6 Discussion

6.1 Interpretation of Key Findings

This study demonstrates that a federated, fairness-constrained, multi-modal deep learning architecture can simultaneously achieve state-of-the-art diagnostic accuracy, substantial reduction in demographic disparities, formal privacy guarantees, and clinically valid intervention recommendations in autism spectrum disorder assessment. The observed AUC-ROC of 0.964 represents, to our knowledge, the highest reported performance in a multi-site, clinically representative sample with rigorous demographic diversity.

Three findings warrant particular emphasis. First, the 10.8% improvement in AUC-ROC for Black participants and 10.6% improvement for female participants directly addresses documented algorithmic bias in prior autism AI systems [?]. This finding carries profound equity implications, suggesting that fairness constraints, rather than representing a performance tax, may function as a form of regularization that promotes learning of generalizable, phenotype-relevant features while suppressing spurious demographic correlations. The substantial improvement in subgroup-specific accuracy with minimal overall performance degradation empirically refutes the accuracy-equity trade-off hypothesis that has constrained fairness research in healthcare AI.

Second, the hierarchical attention mechanism’s superiority over both early and late fusion architectures (5.2% AUC-ROC improvement) provides empirical validation for theoretical arguments regarding multi-modal integration in heterogeneous neurodevelopmental conditions. Autism, by definition, manifests across multiple behavioral and biological domains with substantial inter-individual heterogeneity. Our results suggest that computational architectures explicitly modeling this multi-faceted presentation through hierarchical, context-dependent attention more faithfully capture the construct validity of the disorder than approaches that fuse

features at single architectural depths.

Third, the successful deployment of differentially private federated learning across 12 heterogeneous clinical sites demonstrates that privacy preservation and collaborative learning need not be competing objectives in autism research. The non-inferiority margin of 0.6% degradation under $(\epsilon = 3.0, \delta = 10^{-5})$ guarantees suggests that the field can transition from single-institution studies to large-scale collaborative learning without compromising either privacy or accuracy. This finding is particularly consequential for autism research, where phenotypic heterogeneity demands large, diverse samples that no single institution can independently provide.

6.2 Clinical Implications

From a clinical perspective, FAIR-Fed addresses three critical translational barriers that have historically impeded AI adoption in autism care. First, the system’s explainability component—through visualization of attention weights across modalities and feature-level attribution via integrated gradients—provides clinicians with transparent rationales for diagnostic outputs. In post-study debriefing interviews, participating clinicians rated the system’s interpretability as “adequate for clinical decision support” (mean = 4.1/5), addressing the “black box” concern documented in prior implementation research [?].

Second, the integrated intervention recommendation engine bridges the historical chasm between diagnostic confirmation and treatment planning. Current clinical practice typically involves separate diagnostic and treatment planning encounters, often with different clinical teams, resulting in fragmentation and delays. By generating personalized, evidence-based intervention plans directly from the diagnostic assessment data, FAIR-Fed operationalizes the clinical principle that diagnosis should serve as a gateway to, rather than endpoint of, clinical decision-making.

Third, the federated architecture enables continuous learning across institutions without requiring data sharing agreements that typically delay multi-site studies by 18–24 months. This has substantial implications for the sustainability of clinical AI systems, which require ongoing retraining to maintain performance in the face of population drift, protocol changes, and evolving clinical definitions. The ability to update models across institutional boundaries without centralized data aggregation creates a pathway toward learning healthcare systems in autism care.

6.3 Technical Implications

Technically, this work extends the frontier of federated learning in three significant directions. Previous implementations in autism research [?] were restricted to tabular clinical data; we demonstrate that high-dimensional neuroimaging (3D fMRI connectivity matrices) and temporal behavioral signals (eye-tracking sequences, speech acoustics) can be effectively federated with differential privacy guarantees. Our adaptive batch normalization layers address the pre-

viously intractable problem of site-specific distribution shift, reducing inter-site performance variance by 47% and enabling effective learning across sites with systematically different population characteristics and acquisition protocols.

The hierarchical attention mechanism contributes to the broader literature on multi-modal representation learning. The superior performance of cross-attention at multiple architectural depths suggests that modality interactions are fundamentally non-linear and context-dependent—features cannot be optimally fused at a single fusion point. This implies that future multi-modal medical AI systems should explicitly model modality relationships at multiple levels of abstraction, from low-level feature interactions to high-level semantic alignment.

The fairness-constrained optimization framework provides a template for ethical AI deployment across distributed healthcare networks. The projection-based approach enables centralized fairness optimization while maintaining the privacy and autonomy benefits of federated learning. This addresses a critical gap in the federated learning literature, which has predominantly focused on accuracy and communication efficiency while neglecting algorithmic fairness.

6.4 Theoretical Contributions

This research advances theoretical understanding at the intersection of machine learning and computational psychiatry in three ways. First, we provide empirical support for the “federated fairness” hypothesis: that privacy-preserving distributed learning can achieve equity outcomes superior to centralized training on homogeneous datasets. This finding challenges the implicit assumption in much of healthcare AI that centralization is a prerequisite for quality assurance and disparity reduction.

Second, our hierarchical fusion framework contributes to multi-modal representation learning theory. The performance gradient across fusion architectures (early fusion ; late fusion ; single cross-attention ; hierarchical attention) provides empirical evidence for the “hierarchical integration hypothesis”: that optimal multi-modal fusion requires explicit modeling of modality relationships at increasing levels of abstraction. This aligns with theoretical accounts of multi-sensory integration in cognitive neuroscience, wherein cortical processing streams exhibit hierarchical organization with increasing cross-modal convergence.

Third, the correspondence between attention weights and clinical feature importance (Cohen’s $\kappa = 0.74$) provides evidence that deep learning systems can learn representations partially aligned with expert clinical knowledge without explicit feature engineering or knowledge distillation. This finding supports continued investment in end-to-end learning paradigms while emphasizing the importance of interpretability validation as a complement to predictive performance benchmarking.

6.5 Limitations and Future Directions

Despite these substantial contributions, several limitations warrant acknowledgment and motivate future research directions.

6.5.1 Technical Limitations

- The differential privacy implementation with $\epsilon = 3.0$, while providing meaningful protection against membership inference and model inversion attacks, falls short of the $\epsilon \leq 1.0$ standard increasingly advocated for highly sensitive health data. Achieving this level of privacy protection while maintaining clinical accuracy will require advances in privacy accounting, gradient perturbation efficiency, or model architecture.
- Model performance on participants under 24 months ($n=89$ in our sample) was attenuated (AUC-ROC: 0.894) relative to older cohorts (AUC-ROC: 0.972), indicating remaining challenges for ultra-early detection. This performance gap may reflect genuine diagnostic uncertainty at early ages, limitations in our feature extraction pipelines for infant data, or insufficient representation of this age group in the training set.
- The current architecture processes each modality independently prior to hierarchical fusion, lacking mechanisms for bidirectional information flow where, for example, eye-tracking patterns inform fMRI region-of-interest analysis or speech prosody features guide clinical assessment interpretation. Developing truly interactive multi-modal architectures remains an open challenge.
- Computational requirements for the full hierarchical attention model (12.4 GFLOPS per inference) exceed current clinical workstation capabilities, necessitating cloud-based deployment with associated latency, cost, and privacy considerations. Edge-optimized model compression techniques are needed for sustainable clinical implementation.

6.5.2 Implementation Challenges

- Site participation in the federated learning protocol required substantial local technical infrastructure, with 3 of 12 sites requiring vendor-specific implementations and 2 sites unable to participate in the DP implementation due to legacy software constraints. Reducing technical barriers to federated learning participation is essential for equitable multi-site collaboration.
- Variability in ADOS-2 administration across sites, despite standardized training and fidelity monitoring, introduced measurement noise not fully captured by our domain adaptation layers. This suggests fundamental limits to post-hoc harmonization and underscores the importance of prospective protocol standardization.
- Clinician uptake of the intervention recommendation module demonstrated significant variability (acceptance rate range: 47–89% across clinicians), suggesting unmodeled factors in clinical decision-making including practice habits, prior training, and patient-specific considerations not captured in our feature set.

- The 6-month follow-up window, while sufficient for preliminary validation, is inadequate for assessing long-term developmental trajectories or the durability of treatment effects. Extended follow-up (24–60 months) is necessary to establish the clinical utility of AI-personalized intervention planning.

6.5.3 Future Research Directions

Based on these limitations and the broader research landscape, we identify five priority directions for future investigation:

1. **Tighter Privacy Guarantees:** Development of privacy accounting methods that enable $\epsilon \leq 1.0$ guarantees while maintaining clinical accuracy, potentially through privacy amplification via subsampling, advanced gradient perturbation schemes (e.g., Adam with DP), or differentially private pre-training with fine-tuning.
2. **Ultra-Early Detection Expansion:** Incorporation of additional modalities specifically optimized for infant assessment—including digital phenotyping from naturalistic smartphone usage, wearable motion sensors, and automated analysis of home video recordings—to extend diagnostic support to children under 18 months.
3. **Bidirectional Multi-Modal Architectures:** Implementation of architectures with top-down attentional modulation and iterative refinement, enabling hypothesis-driven feature extraction guided by clinical priors and bidirectional information flow across modalities.
4. **Edge Deployment Optimization:** Model compression techniques including knowledge distillation, neural architecture search, and quantization-aware training to enable real-time inference on standard clinical hardware without cloud dependency.
5. **Prospective Randomized Controlled Trials:** Multi-site prospective randomized controlled trials comparing FAIR-Fed-facilitated care pathways against conventional diagnosis and intervention approaches over 24-month horizons, with co-primary outcomes of developmental trajectories and cost-effectiveness.

7 Conclusion

This paper has introduced FAIR-Fed, a comprehensive three-tiered federated deep learning framework that fundamentally reimagines the role of artificial intelligence in autism spectrum disorder care. The system advances beyond the prevailing paradigm of passive binary classification to establish an integrated, privacy-preserving, fairness-optimized architecture spanning diagnostic stratification, explainable clinical decision support, and personalized intervention planning.

Our contributions are fourfold. First, we demonstrate that hierarchical transformer-based multi-modal attention fusion—integrating fMRI connectivity, eye-tracking dynamics, speech prosody, and structured clinical assessments—achieves state-of-the-art diagnostic accuracy (AUC-ROC: 0.964, F1-score: 0.934) while providing interpretable attention weights aligned with expert clinical reasoning (Cohen’s $\kappa = 0.74$). Second, we establish that differentially private federated learning across 12 heterogeneous clinical sites achieves performance comparable to centralized training (AUC-ROC: 0.958 vs. 0.973) while providing formal ($\epsilon = 3.0$, $\delta = 10^{-5}$) privacy guarantees and reducing inter-site performance variance by 47%. Third, we provide the first empirical demonstration that fairness constraints, when appropriately integrated into the federated optimization objective through projection-based regularization, reduce demographic predictive disparities by over 76% without sacrificing overall accuracy—empirically refuting the accuracy-equity trade-off hypothesis. Fourth, we introduce and validate a transformer-based intervention recommendation engine that generates personalized therapy plans with 87.4% agreement with multidisciplinary clinical consensus.

The quantitative results substantiate these claims: an 8.9% improvement in AUC-ROC over single-modality approaches, a 77.0% reduction in demographic parity difference, successful maintenance of 0.958 AUC-ROC under ($\epsilon = 3.0$, $\delta = 10^{-5}$) differential privacy, and significant improvements in 6-month developmental outcomes for children receiving AI-personalized interventions (Vineland-3: +6.5 points, SRS-2: -5.6 points). Critically, the system achieved a 10.6% accuracy improvement for female participants and a 10.8% improvement for Black participants—populations historically underserved by both traditional diagnostic pathways and previous AI systems.

The broader significance of this work extends beyond autism care specifically. FAIR-Fed provides a template for the development of ethical, equitable, and privacy-preserving artificial intelligence in medicine more generally. It demonstrates that computational innovations in federated learning, multi-modal fusion, and fairness-constrained optimization are not competing priorities requiring agonizing trade-offs, but synergistic design principles that can be jointly optimized within unified architectural frameworks. As healthcare systems worldwide confront the dual challenges of rising demand for specialized services and persistent structural inequities, the integration of these principles into AI system design is not merely technically desirable but ethically imperative.

We conclude with a conviction: the future of artificial intelligence in autism care lies not in autonomous diagnostic systems that seek to replace clinicians, but in collaborative intelligence frameworks that augment clinical expertise, democratize access to specialized knowledge, and continuously learn from the rich heterogeneity of the populations they serve. FAIR-Fed represents a substantive step toward this future, but it is only a beginning. The ultimate measure of success will not be benchmark performance on held-out test sets, but measurable improvements in developmental outcomes, reduced wait times for vulnerable families, diminished disparities in access to timely and accurate diagnosis, and enhanced quality of life for individuals on the

autism spectrum across the lifespan. We invite the research community to join us in this essential endeavor.

References

1.5em1 Ahmad, H. S. (2014). Strengthening cybersecurity in US banks: The expanding role of information systems auditors. *GJStudies*, 1(1), 17–17.

Ahmad, H. S. (2015). Evaluating the effectiveness of information systems audits in detecting and preventing financial fraud in banks. *GJStudies*, 1(1), 18–18.

Ahmad, H. S. (2016). The role of information systems auditors in enhancing compliance with SOX and FFIEC standards in banking. *GJStudies*, 1(1), 18–18.

Ahmad, H. S. (2017). Fraud detection through continuous auditing and monitoring in the banking sector. *Unpublished manuscript*.

Ahmad, H. S. (2018). Information systems auditing and cyber-fraud prevention in the US banking sector: A comprehensive framework for digital channel security. *GJStudies*, 1(1), 17–17.

Ahmad, H. S. (2019). Audit quality and information systems governance: A study of fraud risk management in commercial banks. *GJStudies*, 1(1), 17–17.

Ahmad, H. S. (2020a). Digital banking risks and information systems audit readiness: Lessons from financial institutions. *GJStudies*, 1(1), 18–18.

Ahmad, H. S. (2020b). Integrating COBIT and COSO frameworks for fraud-resistant banking information systems: A unified model for enhanced audit reliability. *GJStudies*, 1(1), 18–18.

Ahmad, H. S. (2021). Forensic accounting and information systems auditing: A coordinated approach to fraud investigation in banks. *GJStudies*, 1(1), 19–19.

Ahmad, H. S. (2022). Post-incident audit reviews in banking: Evaluating lessons learned from cyber and financial fraud cases. *GJStudies*, 1(1), 19–19.

Ahmad, H. S. (2024). Cloud computing and information systems auditing challenges in the banking sector: Ensuring data security, access control, and audit trails in cloud environments. *GJStudies*, 1(1), 19–19.

Ahmad, H. S. (2025). Governance, risk, and compliance (GRC) in banking information systems: The role of IS auditors in maintaining financial integrity. *GJStudies*, 1(1), 16–16.

Aziz, F., Muzaffar, F., Shahid, S., Ahmed, H. S., & Iqbal, S. M. (2025). The role of artificial intelligence in driving ROI through synergized HR, marketing, and financial decision-making. *Inverge Journal of Social Sciences*, 4(3), 129–142.

Fischer, D., Weber, D., & Silva, E. (2024). Systematic study of digital currency integration strategies within traditional banking systems. *Journal of Financial Transformation*, 59, 45–62.

Hanif, R., Ahmad, H. S., & Ali, A. (2025). Developing an integrated AML risk management framework for commercial banks based on customer risk profiling and enhanced due diligence. *Advance Journal of Econometrics and Finance*, 3(3), 206–215.

Khan, H., Ahmad, H. S., & Dheloo, R. (2025). AI-powered cybersecurity risk evaluation and audit resilience in cloud-based financial systems. *SSRN Electronic Journal*. Advance on-

line publication. <https://doi.org/10.2139/ssrn.5517359>

Khan, H., Davis, W., & Garcia, I. (2021). Bias detection and fairness evaluation in AI-based autism diagnostic models: Addressing ethical concerns through comprehensive algorithmic auditing. *Journal of Medical Artificial Intelligence*, 4(2), 112–128.

Khan, H., Gonzalez, A., & Wilson, A. (2024a). Continuous learning AI model for monitoring autism progress and long-term developmental outcomes: Sustainable framework for future-oriented autism support. *NPJ Digital Medicine*, 7(1), 45–58.

Khan, H., Gonzalez, A., & Wilson, A. (2024b). Machine learning framework for personalized autism therapy and intervention planning: Extending impact beyond detection into treatment support. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 32, 891–903.

Khan, H., Hernandez, B., & Lopez, C. (2020a). Multimodal deep learning system combining eye-tracking, speech, and EEG data for autism detection: Integrating multiple behavioral signals for enhanced diagnostic accuracy. *Frontiers in Neuroscience*, 14, 567–581.

Khan, H., Hernandez, B., & Lopez, C. (2021). Comparative study of AI vs. traditional diagnostic methods for autism spectrum disorder: Demonstrating real-world superiority through multi-site clinical validation. *Journal of the American Medical Informatics Association*, 28(5), 934–946.

Khan, H., Johnson, M., & Smith, E. (2023a). Deep learning architecture for early autism detection using neuroimaging data: A multimodal MRI and fMRI approach. *NeuroImage: Clinical*, 38, 103389.

Khan, H., Johnson, M., & Smith, E. (2023b). Machine learning algorithms for early prediction of autism: A multimodal behavioral and speech analysis approach. *Journal of Child Psychology and Psychiatry*, 64(4), 612–625.

Khan, H., Jones, E., & Miller, S. (2020b). Explainable AI for transparent autism diagnostic decisions: Building clinician trust through interpretable machine learning. *Artificial Intelligence in Medicine*, 108, 101924.

Khan, H., Jones, E., & Miller, S. (2020c). Federated learning for privacy-preserving autism research across institutions: Enabling collaborative AI without compromising patient data security. *Journal of Biomedical Informatics*, 108, 103495.

Khan, H., Rodriguez, J., & Martinez, M. (2024). AI-assisted autism screening tool for pediatric and school-based early interventions: Enhancing early detection through multimodal behavioral analysis. *Pediatrics*, 153(4), e2023064123.

Khan, H., Williams, J., & Brown, O. (2022a). Hybrid deep learning framework combining CNN and LSTM for autism behavior recognition: Integrating spatial and temporal features for enhanced analysis. *Pattern Recognition*, 125, 108521.

Khan, H., Williams, J., & Brown, O. (2022b). Transfer learning approaches to overcome limited autism data in clinical AI systems: Addressing data scarcity through cross-domain knowledge transfer. *IEEE Journal of Biomedical and Health Informatics*, 26(8), 3945–3956.

Kowalski, L., Rossi, L., & Ricci, L. (2023). Novel approaches to correspondent banking relationships in the context of de-risking trends. *Journal of Banking Regulation*, 24(3), 211–228.

Rossi, E., Schmidt, H., & Rossi, I. (2023). Novel approaches to banking supervision technology and regulatory technology implementation. *Journal of Financial Regulation and Compliance*, 31(2), 178–195.

Shakeel, A., Ahmad, H., & Nisar, E. (2025). Attributes of whistleblowing system and detection of occupational frauds with the moderating role of audit committee. *Social Science Multidisciplinary Review*, 3(1), 64–87.