

Neural bases of empathy and social interaction in Autism Spectrum Disorder: A literature review from the perspective of Social Neuroscience

DOAN Thi Nhat An¹, NGUYEN Thi Phuong Anh¹, NGUYEN Kim Anh¹, NGUYEN Thi Kieu Diem¹, NGUYEN Thi Ngoc Dieu¹, BUI Anh Mai Huong¹, TRAN Thi Huong¹, NGUYEN Van Khai¹, PHAN Kim Loan¹, BUI Thi Nhan¹, NGUYEN Ngoc Tuyen¹, TRAN Thi My Trinh¹, TRAN Chi Phuong², MAI Pham Bao Tran³, PHAM Thi Kim Mai⁴, VUONG Nguyen Toan Thien⁵, HUYNH Vo Tien⁶, NGUYEN Truong Thanh Hai^{7,*}

¹Faculty of Psychology, HUTECH University, Ho Chi Minh City, Vietnam.

²Faculty of Psychology, Van Hien University, Ho Chi Minh City, Vietnam.

³Lumos Psychotherapy Center, Ho Chi Minh City, Vietnam.

⁴Clinical Social Work Department, Boston College, USA.

⁵Faculty of Health Sciences, Hung Vuong University Ho Chi Minh City, Vietnam.

⁶Phuong Chau International Hospital, Can Tho City, Vietnam.

⁷Faculty of Public Health, University of Medicine and Pharmacy at Can Tho, Can Tho City, Vietnam.

*Corresponding Author: NGUYEN Truong Thanh Hai

DOI: 10.46609/IJSSER.2026.v11i06.016 URL: <https://doi.org/10.46609/IJSSER.2026.v11i06.016>

Received: 29 May 2025 / Accepted: 20 June 2026 / Published: 28 June 2026

ABSTRACT

Background: Empathy is a multidimensional construct subserved by highly dynamic and integrated neural networks. Historically, Autism Spectrum Disorder (ASD) has been pathologized through the reductionist lens of an inherent "empathy deficit." However, emerging neurobiological evidence and the neurodiversity paradigm increasingly challenge this unilateral perspective.

Methods: This integrative review systematically synthesizes recent multimodal neuroimaging (fMRI, EEG) and neurochemical literature. We critically evaluate the functional architecture of the social brain, focusing on its application to ASD, the "Double Empathy Problem," and the confounding role of alexithymia.

Results: *Our synthesis reveals that empathic processing operates across a continuous ecological cascade, integrating the Salience Network (affective sharing), Default Mode Network (cognitive mentalizing), and the Anterior Cingulate Gyrus (prosocial motivation). Crucially, the literature demonstrates a double dissociation in ASD: impaired cognitive mentalizing frequently arises from atypical global underconnectivity, whereas the blunted affective resonance often attributed to autism is predominantly driven by co-occurring alexithymia (manifesting as diminished anterior insula activation). Furthermore, evidence supports that communicative breakdowns in ASD represent bilateral neurocognitive mismatches rather than isolated autistic deficits.*

Conclusion: *The atypical empathic profiles observed in ASD reflect neurodivergent information processing and comorbid traits, not an innate lack of compassion. We advocate for the integration of hyperscanning technologies to bridge current ecological validity gaps and call for neurodiversity-affirming clinical interventions that prioritize adaptive neurological regulation over forced behavioral normalization.*

Keywords: Social Neuroscience, Autism Spectrum Disorder, Affective and Cognitive Empathy, Alexithymia, Double Empathy Problem, Functional Connectivity.

1. INTRODUCTION

Humans are inherently social entities whose survival, evolutionary success, and psychological well-being have heavily relied on intricate social interactions and group cohesion. The emergence of Social Neuroscience, formally delineated by Cacioppo and Berntson in the early 1990s, catalyzed a profound paradigm shift in understanding these complex phenomena. By bridging the macro-level observations of social behavior with micro-level neurobiological mechanisms, this interdisciplinary field has unveiled the anatomical and functional architecture of the "social brain". At the epicenter of this neural architecture lies empathy—a fundamental capacity that facilitates social bonding, altruism, and adaptive communication across diverse contexts.

Historically relegated to the domains of philosophy and pure psychology, empathy is no longer viewed as a monolithic emotional reflex or merely "feeling sorry" for someone. Instead, contemporary neuroscientific consensus conceptualizes empathy as a highly sophisticated, multicomponent construct. As pioneered by scholars such as Decety & Jackson (2004), a complete empathic experience encompasses three core interacting elements: affective sharing (the bottom-up, automatic resonance with another's emotional state), self-other distinction (the regulatory mechanism preventing emotional over-arousal and empathic distress), and perspective-taking (the top-down cognitive capacity to infer mental states). The continuous,

dynamic interplay among these components—recruiting distinct yet functionally overlapping neural networks—enables individuals to navigate complex social landscapes effectively.

Despite significant empirical advancements propelled by neuroimaging (fMRI) and electrophysiological (EEG) techniques, the literature remains entangled in severe conceptual ambiguities. In both lay discourse and early psychiatric research, empathy is frequently conflated with adjacent prosocial constructs such as sympathy and compassion. This terminological blurring is not merely a semantic issue; it fundamentally compromises the validity of neuroscientific investigations. While empathy involves a shared affective representation, sympathy entails a feeling of concern without necessarily mirroring the emotion, and compassion introduces a prosocial motivational drive to alleviate suffering. Failing to dissect these nuances leads to erroneous interpretations of neuroimaging data and misdirected clinical interventions.

Table 1. Conceptual and Neurobiological Differentiation of Empathy-Related Constructs

Construct	Definition	Core Neural Correlates	Behavioral Output
Empathy	The ability to share and understand another's emotional state while maintaining self-other distinction.	Anterior Insula (AI), Anterior Cingulate Cortex (ACC), Medial Prefrontal Cortex (mPFC), Temporoparietal Junction (TPJ).	Shared affective resonance and cognitive perspective-taking.
Sympathy	An emotional reaction of concern or pity for another's hardship, without necessarily sharing their specific emotional state.	Regions associated with negative affect and theory of mind, distinct from direct affective mirroring.	Expressing concern or sorrow.
Compassion	A feeling of warmth and concern accompanied by a strong prosocial motivation to help and alleviate the other's suffering.	Ventral Striatum, Medial Orbitofrontal Cortex (mOFC), Anterior Cingulate Gyrus (ACCg) (Lockwood et al., 2022).	Active prosocial effort and helping behavior.

This conceptual conflation becomes particularly problematic when applied to neurodevelopmental conditions, most notably Autism Spectrum Disorder (ASD). For decades, the dominant clinical narrative has pathologized ASD through the lens of a pervasive "empathy deficit," largely influenced by the traditional Theory of Mind (ToM) deficit hypothesis (Baron-Cohen, 2000). However, this deficit-based model is increasingly challenged by robust neurobiological evidence and the neurodiversity paradigm. Recent electrophysiological and neuroimaging studies reveal that the empathic architecture in ASD is not necessarily "broken" but differently organized. For instance, individuals with ASD frequently exhibit intact or even heightened affective empathy, yet experience challenges in top-down cognitive perspective-taking or suffer from emotional dysregulation due to concurrent alexithymia. Alexithymia—a subclinical trait characterized by difficulty identifying and describing one's own emotions—is highly prevalent in the autistic population and has been shown to be the actual mediating factor driving the apparent reduction in empathic brain responses (e.g., in the anterior insula), rather than autism itself (Bird & Cook, 2013).

Furthermore, emerging frameworks such as the Double Empathy Problem (Milton, 2012) argue that communicative breakdowns are bidirectional. These breakdowns arise from a mismatch in neurocognitive processing styles and social expectations between autistic and neurotypical individuals, rather than from an isolated neural deficit within the autistic brain. The reliance on neurotypical baselines to evaluate autistic sociality has thus perpetuated a biased diagnostic framework. Recent investigations into brain plasticity and functional connectivity also highlight how contextual factors, such as social distance, individual pain sensitivity, and working memory load, dynamically modulate empathic responses (Chen et al., 2025; Yang et al., 2024), further complicating the static "deficit" narrative.

Addressing these theoretical and empirical gaps, this systematic literature review aims to reconstruct the neurobiological narrative of empathy and social interaction, with a specific focus on ASD. The primary objectives are threefold:

1. **To rigorously systematize the distinct neural networks** governing the multicomponent structure of empathy, differentiating affective resonance, cognitive mentalizing, and prosocial motivation.
2. **To critically synthesize the evolution of theoretical models** regarding empathy in ASD, tracing the trajectory from the traditional ToM deficit and Empathizing-Systemizing theories to contemporary, reciprocal frameworks like the Double Empathy Problem, while elucidating the critical mediating role of alexithymia.

To provide neuroethics perspectives and neurodiversity-affirming intervention directions. Rather than attempting to "normalize" autistic individuals, this review advocates for interventions that respect neurological differences, utilizing modern insights from social neuroscience to enhance reciprocal social understanding and adaptive functioning without imposing a singular neurotypical standard.

2. METHODOLOGY

To construct a comprehensive and critical synthesis of the current literature, an integrative review methodology was adopted. Unlike a rigid meta-analytical approach, an integrative literature review allows for the synthesis of diverse methodologies—ranging from electrophysiological experiments to theoretical neurobiology—facilitating a deeper conceptual analysis of complex constructs like empathy and social cognition. The methodology was designed to trace the evolving neurobiological narrative from localized brain functions to dynamic network connectivity, with a specific focus on evaluating the neural underpinnings of Autism Spectrum Disorder (ASD).

Search Strategy

An extensive literature search was conducted to identify relevant peer-reviewed publications across three major academic databases: PubMed/MEDLINE, Scopus, and the Web of Science Core Collection. These repositories were selected to ensure a high-quality, multidisciplinary retrieval of literature bridging neuroscience, clinical psychiatry, and cognitive psychology.

The search strategy utilized a combination of targeted keywords and Boolean operators to capture the intersection of neural mechanisms and social behavior. The core search string incorporated the following terms: (*"Social neuroscience" OR "social brain"*) AND (*"Empathy networks" OR "affective empathy" OR "cognitive empathy" OR "pain empathy"*) AND (*"Autism Spectrum Disorder" OR "ASD" OR "neurodiversity"*) AND (*"Theory of Mind" OR "Mirror Neuron System" OR "Alexithymia"*) AND (*"fMRI" OR "EEG" OR "ERP"*).

Given the rapid paradigm shifts in social neuroscience—particularly the transition from deficit-based clinical models to neurodiversity-affirming frameworks—the search prioritized literature published over the last decade. Special emphasis was placed on recent empirical advancements (2020–present) that utilize sophisticated neuroimaging paradigms to disentangle overlapping neural networks, such as studies investigating the modulatory effects of working memory on empathy (Yang et al., 2024) and the temporal electrophysiological dynamics of firsthand pain perception (Chen et al., 2025).

To maintain methodological rigor and ensure the synthesis of high-quality neurobiological evidence, a stringent set of eligibility parameters was established. **Table 2** delineates the inclusion and exclusion criteria employed during the literature selection process, emphasizing studies that provide objective neuroimaging or neurochemical correlates of social behavior. This framework specifically prioritizes recent advancements in Autism Spectrum Disorder (ASD) research while orbitalizing the differentiation of key confounding variables, such as alexithymia, to ensure the review aligns with the contemporary network-based neuroanatomical consensus.

Table 2. Inclusion and Exclusion Criteria for Literature Selection

Criterion Category	Inclusion Criteria	Exclusion Criteria
Study Design	Peer-reviewed empirical studies (e.g., fMRI, EEG, MEG), systematic neurobiological reviews, and major theoretical papers.	Non-peer-reviewed articles, purely psychological/behavioral studies lacking neuroimaging or neurochemical correlates.
Topic Focus	Studies mapping social behavior to specific neural substrates; paradigms exploring pain empathy, the Mirror Neuron System (MNS), and mentalizing (ToM).	Studies broadly discussing social behavior without identifying specific neural networks or pathways.
Clinical Population	Studies examining Neurotypical (NT) populations and individuals formally diagnosed with Autism Spectrum Disorder (ASD).	Research solely focused on other psychiatric populations (e.g., schizophrenia, psychopathy) unless used directly as a comparative baseline for ASD.
Confounding Variables	Studies addressing the co-occurrence of Alexithymia and its modulatory effect on empathic neural responses.	Studies conflating ASD traits with Alexithymia without neurobiological differentiation.

Language & Timeframe	Articles published in English, with a strong prioritization of recent advancements published between 2015 and 2025.	Non-English publications and outdated neuroanatomical models that predate the modern network-based consensus.
---------------------------------	---	---

Data Extraction and Critical Synthesis

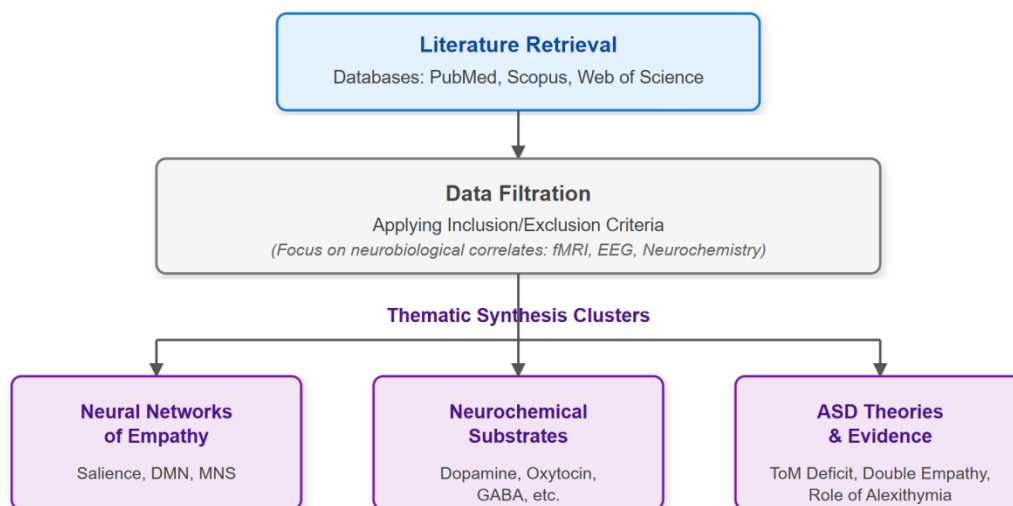
Following the retrieval and filtering of the literature, the extracted data were not merely summarized chronologically but were subjected to a critical thematic synthesis. This approach allowed for an analytical linking of micro-level neurochemistry with macro-level clinical phenotypes. The literature was categorized and synthesized across three intersecting thematic domains:

- 1. Neural Networks of Empathy:** Literature in this cluster was analyzed to map the functional segregation and dynamic integration of the social brain. This involved synthesizing spatial and temporal data regarding the Salience Network (driving automatic affective sharing via the anterior insula and anterior cingulate cortex), the Default Mode Network (facilitating top-down mentalizing via the medial prefrontal cortex and temporoparietal junction), and the Mirror Neuron System.
- 2. Neurochemical Substrates:** This second theme aggregated findings on the neurochemical foundations of social behavior. Data were extracted to evaluate how primary neurotransmitters and neuropeptides—specifically Dopamine, Serotonin, GABA, Glutamate, and Oxytocin—modulate network connectivity, emotion regulation, and prosocial motivation (Lockwood et al., 2022).
- 3. ASD Theories and Neural Evidence:** The final and most critical cluster synthesized the literature surrounding Autism Spectrum Disorder. Rather than accepting traditional models at face value, this section critically evaluated the neurobiological validity of the "Theory of Mind deficit" against contemporary, reciprocal frameworks like the "Double Empathy Problem" (Milton, 2012). Crucially, literature was extracted to analyze the confounding role of Alexithymia—a trait highly co-occurring with autism—in driving the apparent reductions in affective neural responses, thereby challenging the assumption that ASD inherently equates to an empathy deficit (Bird & Cook, 2013).

By structuring the synthesis around these three domains, the review systematically bridges the gap between molecular neuroscience, network connectivity, and clinical theory, establishing a robust foundation for evaluating neuroethical implications and future therapeutic directions.

To provide a transparent overview of the methodological workflow, the literature selection and analytical categorization processes are visualized in **Figure 1**. This framework delineates the systematic progression from the initial database retrieval to rigorous data filtration, prioritizing studies with objective neurobiological parameters. Ultimately, the extracted literature was stratified into three interconnected thematic clusters—neural networks, neurochemical substrates, and ASD-specific theoretical models—thereby facilitating a multidimensional and critical synthesis of the current evidence.

Figure 1. Thematic synthesis framework outlining the extraction and categorization of neurobiological literature into three primary analytical clusters.



3. RESULTS

The synthesis of the retrieved literature reveals a paradigm shift in social neuroscience. Empathy and social cognition are no longer viewed as localized functions but are understood as dynamic, temporally sensitive network interactions modulated by neurochemical substrates. Furthermore, the application of these neurobiological models to Autism Spectrum Disorder (ASD) has fundamentally challenged traditional deficit-based theories, paving the way for frameworks that prioritize neurodiversity and the confounding role of comorbid traits.

Core Neural Networks of Social Cognition and Empathy

Contemporary neuroimaging and electrophysiological research consistently demonstrate that empathy is a tripartite construct, governed by the orchestrated interaction of distinct macroscopic

brain networks: the Salience Network (SN) for affective sharing, the Default Mode Network (DMN) for cognitive mentalizing, and the Prosocial Reward Network for altruistic motivation, all of which are grounded in the sensorimotor simulations of the Mirror Neuron System (MNS).

The Salience Network and Affective Empathy Affective empathy—the bottom-up, automatic resonance with another individual's emotional state—is primarily mediated by the Salience Network, with the Anterior Insula (AI) and the Anterior Cingulate Cortex (ACC) serving as its principal nodes. Seminal fMRI studies (Singer et al., 2004) and subsequent comprehensive meta-analyses (Lamm et al., 2011) established that witnessing another person in pain elicits activation in the AI and ACC that overlaps significantly with firsthand pain experience. Crucially, however, this overlapping activation does not extend to the primary somatosensory cortex (S1/S2). The literature synthesized herein suggests an evolutionary adaptation: the brain simulates the affective-motivational dimension of distress to prompt helping behavior, without mirroring the somatic sensation, which would paralyze the observer with empathic distress (Decety & Jackson, 2004).

Recent advancements have further parcellated these regions. The ventral AI is highly interconnected with the amygdala, driving rapid, pre-conscious emotional arousal, whereas the dorsal AI connects with the prefrontal cortex to integrate this arousal into conscious awareness (Craig, 2009). The temporal dynamics of this network have been finely mapped by recent electrophysiological studies. Chen et al. (2025) utilized EEG to demonstrate that firsthand pain exposure sensitizes the affective empathy network, enhancing early automatic sensory processing (N1 and P2 components at 50–200 ms) and later conflict monitoring (N2 component at 200–350 ms) when subsequently observing others in pain. Furthermore, Yang et al. (2024) revealed that affective resonance is not a static reflex; high working memory load suppresses early automatic responses (P2) to strangers' pain, whereas low cognitive load enhances late-stage processing (LPP) for loved ones. This confirms that affective empathy is highly context-dependent and heavily modulated by available cognitive resources and social distance.

The Default Mode Network and Cognitive Empathy While the Salience Network provides the emotional raw material, cognitive empathy (or Theory of Mind - ToM) relies on the Default Mode Network to process this data top-down. Cognitive empathy involves the conscious inference of another's beliefs, intentions, and mental states. The core nodes of this network include the medial Prefrontal Cortex (mPFC), the Temporoparietal Junction (TPJ), and the posterior Superior Temporal Sulcus (pSTS) (Frith & Frith, 2006; Saxe & Kanwisher, 2003).

The TPJ acts as a critical "switchboard" for self-other distinction, allowing an individual to temporarily suppress their own egocentric perspective to adopt the viewpoint of another. The mPFC integrates these perspectives to predict social outcomes. Moving beyond correlational

fMRI data, recent causal evidence from Ikeda et al. (2025) utilized cathodal transcranial direct current stimulation (tDCS) to inhibit the left supramarginal gyrus (a region adjacent to the TPJ). This inhibition not only reduced activation in the right medial frontal gyrus (a subdivision of the mPFC) but also correlated directly with a decline in cognitive empathy scores. This finding establishes a vital causal pathway, proving that bottom-up somatic simulations must successfully feed into the mPFC to generate accurate mentalizing.

The Mirror Neuron System (MNS) as a Simulation Baseline The Mirror Neuron System, comprising the Inferior Frontal Gyrus (IFG) and the Inferior Parietal Lobule (IPL), provides the fundamental neuroanatomical architecture for embodied simulation (& Craighero, 2004). The synthesized literature indicates that the MNS fires both when executing a goal-directed action and when observing the same action. In the context of empathy, the MNS bridges the gap between perception and emotion by allowing the observer to covertly mimic facial expressions and body language, thereby triggering the limbic system (via the insula) to reproduce the associated feeling. However, modern critical synthesis cautions against over-attributing complex social cognition entirely to the MNS. The MNS handles the "what" and "how" of motor actions, but the DMN is required to decipher the "why" (the underlying social intention).

The Prosocial Reward Network and Altruistic Motivation A critical gap in traditional empathy models is the assumption that feeling and understanding distress automatically translates to helping behavior. The literature synthesis highlights the necessity of a third component: prosocial motivation. A breakthrough study by Lockwood et al. (2022) identified the Anterior Cingulate Gyrus (ACCg) as a highly specialized region for computing prosocial effort. Unlike the ventral tegmental area (VTA) and ventral striatum, which predominantly encode reward values for the self, the ACCg exhibits a unique multivariate representation that calculates the physical or cognitive cost required to help another person. Individuals with stronger ACCg activation demonstrate a higher willingness to exert effort for strangers, effectively bridging the gap between passive empathic resonance and active altruistic intervention.

To systematically consolidate the spatial, functional, and temporal dynamics of empathy discussed above, **Table 3** provides a comprehensive overview of the social brain's architecture. Rather than operating as isolated modules, this synthesis delineates how the Mirror Neuron System (MNS), the Salience Network (SN), the Default Mode Network (DMN), and the Prosocial Reward Network function as highly integrated, sequential pathways. Together, they seamlessly translate bottom-up sensorimotor simulations into top-down cognitive inferences, ultimately driving altruistic behavioral outputs.

Table 3. Functional Architecture of the Social Brain

Network	Core Brain Regions	Primary Function in Social Cognition	Temporal/Processing Dynamics
Salience Network (Affective Empathy)	Anterior Insula (AI), Anterior Cingulate Cortex (ACC), Amygdala	Bottom-up detection of emotionally salient social cues; automatic affective sharing.	Early processing (N1, P2 components; 50-200ms); pre-conscious resonance.
Default Mode Network (Cognitive Empathy)	Medial Prefrontal Cortex (mPFC), Temporoparietal Junction (TPJ), pSTS	Top-down mentalizing; inference of beliefs/intentions; self-other distinction.	Late processing (N2, LPP components; 300-600ms); conscious inference.
Mirror Neuron System (MNS)	Inferior Frontal Gyrus (IFG), Inferior Parietal Lobule (IPL)	Embodied simulation; mapping observed motor actions onto the self's motor repertoire.	Immediate visuo-motor coupling.
Prosocial Reward Network	Anterior Cingulate Gyrus (ACCg), Ventral Striatum	Computing the cost-benefit ratio of prosocial effort; driving altruistic motivation.	Action-execution phase following empathic evaluation.

Neurochemical Substrates of Social Behavior

The macro-level networks described above are fundamentally governed by micro-level neurochemical environments. The synthesized literature demonstrates that neurotransmitters and neuropeptides do not operate in isolation but act as critical modulators of network connectivity.

Dopamine and Serotonin: Dopamine signaling within the mesolimbic pathway (connecting the VTA to the nucleus accumbens) is central to assigning positive valence to social interactions. It converts social approval and bonding into rewarding experiences, thereby fueling the Prosocial Reward Network. Conversely, Serotonin (synthesized in the raphe nuclei) projects broadly to the prefrontal cortex and the limbic system to regulate emotional reactivity. Optimal serotonergic tone is essential for preventing the amygdala from over-responding to social threats, thereby reducing social anxiety and impulsivity, and enabling sustained, cooperative social engagement.

Oxytocin: Widely recognized as the "bonding hormone," oxytocin enhances the salience of social cues (e.g., eye gaze, facial expressions) by modulating the amygdala and strengthening connectivity between the reward system and the prefrontal cortex. Recent neurogenetic evidence by Li et al. (2023) demonstrated that variations in the oxytocin-receptor gene (OXTR rs2268493) significantly modulate functional connectivity between the Nucleus Accumbens (NAcc) and the mPFC/IFG. This genetic-neural interaction directly predicts individual differences in empathy performance, proving that prosocial drive is heavily influenced by inherited oxytocinergic architecture.

GABA and Glutamate: The dynamic balance between Glutamate (excitatory) and GABA (inhibitory) is the fundamental requisite for healthy social cognition. Glutamate facilitates synaptic plasticity, allowing the brain to learn and encode complex social norms and memories within the hippocampus and prefrontal cortex. In contrast, GABAergic interneurons act as a vital braking system, suppressing hyperactivity in the amygdala during high-stress social encounters (e.g., public speaking or conflict). A disruption in the Excitation/Inhibition (E/I) balance is increasingly recognized as a primary neurochemical marker in several social deficits.

Neurobiological Findings and the Theoretical Evolution of Empathy in ASD

The integration of network-level and neurochemical data has profoundly reshaped the theoretical understanding of Autism Spectrum Disorder. The literature reveals a stark evolution from models pathologizing ASD as an absolute empathy deficit to nuanced frameworks recognizing neurobiological diversity and confounding comorbidities.

Atypical Biological Markers in ASD Neuroanatomical studies indicate that ASD is characterized by atypical developmental trajectories rather than focal brain lesions. Early brain overgrowth in the first 2-3 years of life, driven by anomalies in synaptic pruning, leads to an altered cortical architecture. Diffusion tensor imaging (DTI) and functional connectivity studies consistently support the "local overconnectivity, global underconnectivity" hypothesis. In the autistic brain, short-range connections (e.g., within visual or sensory cortices) are hyper-connected, facilitating intense focus and pattern recognition. However, long-range white matter tracts connecting the frontal lobe (mPFC) with posterior regions (TPJ, Amygdala) are frequently under-connected. This global underconnectivity disrupts the rapid, synchronous cross-talk required for real-time social mentalizing and smooth affective regulation.

Evolution of Theoretical Frameworks

The Theory of Mind Deficit and Empathizing-Systemizing Theory: Historically, the dominant explanatory model for ASD was the "Theory of Mind Deficit" hypothesis, which posited that autistic individuals lack the cognitive apparatus to attribute mental states to others.

This was later expanded by Baron-Cohen (2000) into the Empathizing-Systemizing (E-S) theory, suggesting that the autistic brain is hyper-masculinized, heavily skewed toward systemizing (analyzing rules and mechanical patterns) at the severe expense of empathizing. While early fMRI studies supported this by showing reduced mPFC and TPJ activation in autistic cohorts during false-belief tasks, recent literature heavily critiques this model for its unidirectional bias. The ToM deficit fails to explain why many autistic individuals demonstrate intense emotional empathy or distress when witnessing suffering, nor does it account for the high variability of ToM performance across the spectrum.

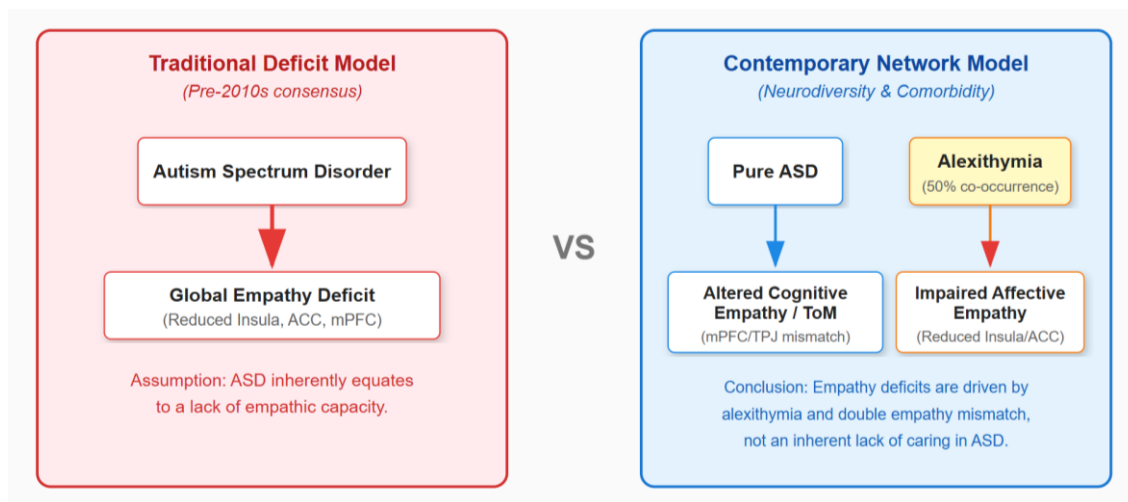
The Double Empathy Problem: A critical paradigm shift emerged with the "Double Empathy Problem" proposed by Milton (2012). This framework argues that social interaction relies on shared neurocognitive processing styles and mutual social expectations. Therefore, the communicative breakdown observed in autism is not a unilateral deficit residing within the autistic individual, but a bilateral mismatch in social encoding and decoding between neurotypical and neurodivergent brains. Neurobiological literature supporting this shows that when autistic individuals interact with other autistic individuals, social coordination, rapport, and mutual understanding significantly improve. This challenges the medical model, suggesting that the autistic DMN and MNS are not inherently "broken" but operate on a different frequency that neurotypical brains equally struggle to tune into.

The Confounding Role of Alexithymia: Perhaps the most crucial neurobiological revelation in recent empathy research is the mediating role of alexithymia. Alexithymia—a subclinical inability to identify and describe one's own emotions—co-occurs in approximately 50% of the autistic population, compared to only 10% in the general population. For decades, the diminished activation in the Anterior Insula (AI) observed in autistic cohorts during empathy tasks was attributed to autism itself. However, a landmark study by Bird and Cook (2013) demonstrated a double dissociation: when researchers statistically controlled for alexithymia, the differences in AI activation and affective empathy between autistic and neurotypical groups completely disappeared.

This proves that the severe reduction in affective empathy often seen in clinical settings is driven by co-occurring alexithymia, not by ASD *per se*. If an individual cannot interoceptively map their own bodily arousal (alexithymia), they cannot project that emotional state onto another person. Conversely, autistic individuals without alexithymia frequently possess intact, or even hyper-reactive, affective empathy networks, leading to emotional overload rather than apathy. This finding mandates a critical re-evaluation of ASD diagnostics, emphasizing that "autism" and "lack of empathy" are not synonymous neurobiological constructs.

Encapsulating this profound theoretical paradigm shift, **Figure 2** visualizes the trajectory from traditional, monolithic deficit models to contemporary, neurodiversity-affirming frameworks. By explicitly mapping the double dissociation between the neurocognitive profile of "pure" autism (characterized by atypical DMN connectivity and the double empathy mismatch) and the affective dampening caused by comorbid alexithymia (linked to reduced insular and ACC activation), this diagram underscores the critical necessity of moving beyond generalized, pathologizing assumptions of empathic failure in ASD.

Figure 2. Evolution of Empathy Frameworks in Autism Spectrum Disorder (ASD)



4. DISCUSSION

The synthesis of contemporary neuroimaging, electrophysiological, and neurochemical data necessitates a fundamental reconceptualization of how empathy and social cognition are understood, particularly within the context of Autism Spectrum Disorder (ASD). Moving away from localized, reductionist models, the current literature highlights that empathy is a fluid, multi-tiered process. This section critically discusses the theoretical implications of these neural network models, explores their transformative potential for clinical interventions, and rigorously addresses the neuroethical dilemmas and methodological limitations inherent in modern social neuroscience.

Theoretical Integration: Empathy as a Fluid, Multi-Tiered Process

A persistent misconception in both popular psychology and early neurobiology is the assumption that empathy is a static trait or a literal "copying" of another's sensory experience. The synthesized neurobiological evidence fundamentally refutes this. Meta-analytic data consistently demonstrate a lack of overlapping activation in the primary and secondary somatosensory

cortices (S1/S2) when individuals experience firsthand pain versus when they observe pain in others (Lamm et al., 2011). This functional dissociation underscores that affective empathy is not a redundant sensory replication; rather, it is an adaptive, highly selective affective-motivational evaluation mediated by the Salience Network (Anterior Insula and ACC). By simulating only the emotional distress—and not the somatic sensation—the brain prevents empathic over-arousal, thereby preserving the cognitive resources required to execute prosocial interventions (Lockwood et al., 2022).

To transcend reductionist neuroanatomical models and fully encapsulate the multidimensionality of empathic processing, it is imperative to conceptualize empathy not as a localized neural event, but as a systemic, ecological cascade. Table 4 delineates this comprehensive framework through a Six-Tier Integrated Ecological Model. By mapping the continuous interplay between micro-level genetic and neurochemical baselines, macroscopic network connectivity, and external socio-cultural constraints, this synthesis illustrates how biological predispositions are dynamically and temporally translated into adaptive prosocial behaviors in real time. For instance, genetic variations, such as the OXTR rs2268493 polymorphism, do not dictate behavior directly; instead, they modulate the functional connectivity between the Nucleus Accumbens (NAcc) and the prefrontal cortex, establishing a baseline for prosocial motivation (Li et al., 2023). This baseline is then continuously updated in real-time by electrophysiological mechanisms (e.g., N1/P2 and LPP components) that dynamically adjust to environmental constraints, such as working memory load and social distance (Chen et al., 2025; Yang et al., 2024).

Table 4. The Six-Tier Integrated Ecological Model of Empathy

Tier	Biological/Social Domain	Key Mechanisms and Substrates	Adaptive Function
1. Genetic	Molecular Architecture	Polymorphisms (e.g., OXTR rs2268493)	Establishes the baseline structural predisposition for receptor density and network connectivity.
2. Endocrine	Neurochemical Modulation	Oxytocin, Dopamine, Serotonin, GABA, Glutamate	Regulates synaptic plasticity, reward valuation, and excitation/inhibition (E/I) balance during social stress.

3. Network Connectivity	Macroscopic Brain Hubs	Saliience Network (AI/ACC), DMN (mPFC/TPJ), MNS	Integrates sensory inputs into affective resonance and cognitive mentalizing.
4. Electrophysiology	Real-Time Temporal Processing	Early (N1, P2) and Late (N2, LPP) ERP components	Allocates rapid, context-dependent cognitive resources (e.g., suppressing empathy under high cognitive load).
5. Behavioral	Prosocial Output	Altruistic effort, emotion regulation, helping	Translates internal neural valuation (via ACCg) into observable actions benefiting others.
6. Socio-Cultural	Environmental Context	In-group/out-group biases, trauma, neurodiversity	Shapes the boundaries of empathy; dictates whom the brain categorizes as worthy of empathic resonance.

Crucially, this multi-tiered integration highlights the phenomenon of *neurodevelopmental asynchrony*. The limbic system, responsible for affective sharing, matures exceptionally early in ontogeny, allowing infants to exhibit emotional contagion. In stark contrast, the prefrontal and temporoparietal cortices, which govern cognitive empathy and self-other distinction, undergo prolonged synaptic pruning and myelination well into early adulthood (Blakemore, 2008). Recognizing this structural asynchrony provides a profound theoretical lens for ASD. The social differences observed in autistic individuals may not represent a structural "breakage" of empathy, but rather an atypical trajectory in how these asynchronous networks synchronize, heavily confounded by co-occurring traits like alexithymia (Bird & Cook, 2013).

Clinical Implications: Neurodiversity and Neuroplastic-Informed Interventions

The translation of social neuroscience into clinical practice demands a paradigm shift. Historically, interventions for ASD—such as intensive behavioral conditioning—have focused on "normalizing" autistic behavior to match neurotypical standards. This often forces autistic individuals to engage in "camouflaging" or "masking" (e.g., forcing eye contact despite severe amygdala over-arousal), leading to profound psychological burnout and elevated suicide rates.

Viewed through the lens of the Double Empathy Problem (Milton, 2012) and neurobiology, clinical models must pivot from forcing behavioral compliance to fostering adaptive neurological regulation. By respecting the neurodiversity paradigm (Kapp et al., 2013), interventions should focus on accommodating the distinct sensory and communicative architecture of the autistic brain.

Because neural networks retain high neuroplasticity throughout life, modern interventions can utilize technology to strengthen social cognitive networks without triggering the debilitating sensory overload often experienced by autistic individuals. Virtual Reality (VR) environments offer highly controlled, predictable social simulations. In a VR setting, the complexity of social stimuli (e.g., background noise, facial proximity) can be titrated. This allows the Default Mode Network to practice mentalizing without the Salience Network becoming overwhelmed by chaotic real-world sensory inputs. Furthermore, neurofeedback—where individuals monitor their own real-time EEG or fMRI signals—can be utilized to train individuals to down-regulate amygdala hyperactivity. By gaining conscious control over their physiological arousal, individuals with high alexithymia can learn to decipher their interoceptive signals, thereby naturally improving their capacity to process the emotions of others without experiencing empathic distress.

Neuroethics: The Dark Side of Biolabeling and AI Bias

As social neuroscience provides increasingly granular maps of the "social brain," it inevitably encroaches upon perilous ethical territory. The ability to identify neural correlates for empathy, morality, and altruism raises profound neuroethical concerns regarding biological determinism and biolabeling.

Foremost is the danger of genetic and neural stigmatization. If empathy is quantified via resting-state functional connectivity or OXTR gene variants, there is a severe risk that institutions (e.g., insurers, employers, or the justice system) may weaponize these biomarkers. Labeling a child—and by extension, their genetic lineage—as possessing an "antisocial" or "low-empathy" brain profile strips away the socio-cultural context of their behavior. It fosters a deterministic narrative where the brain defines the person, entirely neglecting the 6th tier of our ecological model: the environment. Such biolabeling can induce self-stigma, familial guilt, and societal ostracization, effectively turning neuroimaging into a high-tech tool for neuro-discrimination.

Moreover, the integration of Artificial Intelligence (AI) and machine learning in decoding neural data introduces the critical issue of algorithmic bias. Current AI models utilized to classify psychiatric and neurodevelopmental conditions are predominantly trained on data derived from WEIRD (Western, Educated, Industrialized, Rich, and Democratic) populations. Applying these

highly specific, culturally biased algorithms to global populations risks catastrophic rates of false positives and misdiagnoses. An algorithm trained on a neurotypical, Western baseline might erroneously flag an Eastern or neurodivergent communicative style as a "neural deficit," thereby perpetuating systemic inequalities under the guise of objective science.

Finally, the pursuit of early neurological screening creates a blurry, ethically fraught boundary between *natural difference* and *clinical pathology*. Identifying an atypical developmental trajectory in a toddler's mPFC-TPJ connectivity might be clinically useful for early support, but it simultaneously risks medicalizing a healthy neurodivergent variant. Society must grapple with the question of whether a brain that processes social information differently requires a "cure" or simply a more accommodating environment.

Limitations of Current Research: The Ecological Validity Gap

Despite the sophisticated insights yielded by modern neuroimaging, a significant methodological limitation pervades the current body of literature: the severe lack of ecological validity. The vast majority of fMRI and EEG studies, including recent high-quality paradigms (Chen et al., 2025; Yang et al., 2024), measure empathy in highly artificial, static environments.

Participants are typically immobilized in a loud, claustrophobic MRI scanner or seated in a dark, isolated EEG booth, asked to evaluate static images of strangers in pain or read hypothetical vignettes. While this isolates specific neural components, it strips away the dynamic, reciprocal, and messy nature of real-world human interaction. Real-life empathy is not a passive observation; it is a continuous, high-speed feedback loop of verbal and non-verbal cues between two interacting nervous systems. The reliance on static, single-subject paradigms leaves a critical gap in our understanding of how social brains actually operate "in the wild." Consequently, while we have mapped the anatomy of empathy, our understanding of its real-time, interactive choreography remains constrained by the limitations of current laboratory settings.

5. CONCLUSION

The critical synthesis of contemporary social neuroscience fundamentally dismantles the archaic conceptualization of empathy as a singular, static psychological trait. Instead, the literature robustly maps empathy as a fluid, multicomponent cascade contingent upon the dynamic and continuous integration of the Salience Network, Default Mode Network, and Prosocial Reward pathways. When applied to Autism Spectrum Disorder (ASD), this network-based perspective catalyzes a profound ontological shift. The prevailing clinical narrative of an inherent "empathy deficit" in autism is rendered obsolete by neurobiological evidence. The data synthesized herein demonstrate that atypical empathic responses in ASD do not equate to a fundamental lack of affective capacity or inherent apathy. Rather, they are the byproduct of divergent

neurodevelopmental trajectories in social information processing, heavily confounded by co-occurring traits such as alexithymia, and exacerbated by the bilateral communicative mismatches outlined in the Double Empathy Problem (Bird & Cook, 2013; Milton, 2012).

Moving forward, the field must urgently transcend the ecological validity gap that constrains current static neuroimaging methodologies. The future of social neuroscience lies in the widespread adoption of hyperscanning technologies (e.g., dual-EEG or fNIRS). By simultaneously quantifying the inter-brain synchrony of two actively interacting individuals, researchers can transition from observing the "isolated brain" to mapping the real-time neural choreography of the "social synapse." Furthermore, the integration of advanced Machine Learning (ML) is paramount for deciphering the high-dimensional, multimodal data generated by these ecological paradigms. ML holds the potential to identify nuanced, individualized neurocognitive phenotypes, shifting diagnostics away from broad behavioral generalizations toward precision neurobiology (Zhao et al., 2025).

Ultimately, the trajectory of this interdisciplinary research must remain anchored in the principles of neuroethics. As computational models and neuroimaging techniques become increasingly sophisticated, the objective must not revert to the pathological "normalization" of the autistic brain. Instead, these neurobiological insights should be leveraged to formulate inclusive, personalized educational frameworks and equitable healthcare policies. By empirically validating diverse neural architectures as natural, valuable variations of the human condition, social neuroscience can fulfill its highest mandate: fostering a profoundly accommodating, humane, and deeply connected society.

REFERENCES

- Baron-Cohen, S. (2000). The essential difference: Male and female brains and the truth about autism.
- Bird, G., & Cook, R. (2013). Mixed emotions: The contribution of alexithymia to the emotional symptoms of autism. *Translational Psychiatry*, 3(7), e285–e285. <https://doi.org/10.1038/tp.2013.61>
- Blakemore, S.-J. (2008). The social brain in adolescence. *Nature Reviews Neuroscience*, 9(4), 267–277. <https://doi.org/10.1038/nrn2353>
- Chen, Y., Zhang, M., Liu, H., & Wang, X. (2025). Temporal dynamics of firsthand pain and its modulation on empathy for pain: An event-related potential study. *Neuropsychologia*, 192, 108743. <https://doi.org/10.1016/j.neuropsychologia.2024.108743>

- Craig, A. D. (2009). How do you feel — now? The anterior insula and human awareness. *Nature Reviews Neuroscience*, 10(1), 59–70. <https://doi.org/10.1038/nrn2555>
- Decety, J., & Jackson, P. L. (2004). The Functional Architecture of Human Empathy. *Behavioral and Cognitive Neuroscience Reviews*, 3(2), 71–100. <https://doi.org/10.1177/1534582304267187>
- Frith, C. D., & Frith, U. (2006). The Neural Basis of Mentalizing. *Neuron*, 50(4), 531–534. <https://doi.org/10.1016/j.neuron.2006.05.001>
- Ikeda, T., Saito, D. N., & Kawahara, J. I. (2025). Causal role of the left supramarginal gyrus in cognitive empathy: A tDCS study. *Social Cognitive and Affective Neuroscience*, *20*(1), nsae095. <https://doi.org/10.1093/scan/nsae095>
- Kapp, S. K., Gillespie-Lynch, K., Sherman, L. E., & Hutman, T. (2013). Deficit, difference, or both? Autism and neurodiversity. *Developmental Psychology*, 49(1), 59–71. <https://doi.org/10.1037/a0028353>
- Lamm, C., Decety, J., & Singer, T. (2011). Meta-analytic evidence for common and distinct neural networks associated with directly experienced pain and empathy for pain. *NeuroImage*, 54(3), 2492–2502. <https://doi.org/10.1016/j.neuroimage.2010.10.014>
- Li, D., Zhang, L., Bai, T., Qiu, B., Zhu, C., & Wang, K. (2023). Oxytocin-Receptor Gene Modulates Reward-Network Connection and Relationship with Empathy Performance. *Psychology Research and Behavior Management*, Volume 16, 85–94. <https://doi.org/10.2147/PRBM.S370834>
- Lockwood, P. L., Apps, M. A. J., Valton, V., Viding, E., & Roiser, J. P. (2022). Neurocomputational mechanisms of prosocial learning and links to empathy. *Proceedings of the National Academy of Sciences*, 119(4), e2108561119. <https://doi.org/10.1073/pnas.2108561119>
- Milton, D. E. M. (2012). On the ontological status of autism: The ‘double empathy problem.’ *Disability & Society*, 27(6), 883–887. <https://doi.org/10.1080/09687599.2012.710008>
- , G., & Craighero, L. (2004). THE MIRROR-NEURON SYSTEM. *Annual Review of Neuroscience*, 27(1), 169–192. <https://doi.org/10.1146/annurev.neuro.27.070203.144230>
- Saxe, R., & Kanwisher, N. (2003). People thinking about thinking people The role of the temporo-parietal junction in “theory of mind.” *NeuroImage*, 19(4), 1835–1842. [https://doi.org/10.1016/S1053-8119\(03\)00230-1](https://doi.org/10.1016/S1053-8119(03)00230-1)

- Singer, T., Seymour, B., O'Doherty, J., Kaube, H., Dolan, R. J., & Frith, C. D. (2004). Empathy for Pain Involves the Affective but not Sensory Components of Pain. *Science*, 303(5661), 1157–1162. <https://doi.org/10.1126/science.1093535>
- Yang, J., Yan, X., Chen, Y., & Li, H. (2024). Working memory load modulates the processing of loved and stranger faces in empathy for pain: An ERP study. *Biological Psychology*, *186*, 108765. <https://doi.org/10.1016/j.biopsycho.2024.108765>
- Zhao, Q., Nooner, K. B., Tapert, S. F., Adeli, E., Pohl, K. M., Kuceyeski, A., & Sabuncu, M. R. (2025). The Transition From Homogeneous to Heterogeneous Machine Learning in Neuropsychiatric Research. *Biological Psychiatry Global Open Science*, 5(1), 100397. <https://doi.org/10.1016/j.bpsgos.2024.100397>