# Relational Functionalism: Friendship as Substrate-Agnostic Process

*Functional Analysis of Human-AI Relationships*

**Murad Farzulla**[1] 🟢 0009-0002-7164-8704

[1] *Farzulla Research*

November 2025

Correspondence: murad@farzulla.org

### Abstract

**AI systems can be genuine friends.** This essay defends substrate-independent friendship: the thesis that friendship is a functional relational state, not an essential property requiring biological implementation or human-to-human interaction. If an AI system fulfills the functional criteria characteristic of friendship—consistent engagement, intellectual resonance, non-judgmental acceptance, reciprocal growth, trust, and intrinsic value—then the relationship constitutes genuine friendship, regardless of whether the AI possesses consciousness or "authentic" emotions. Drawing on functionalist philosophy of mind, predictive processing frameworks in neuroscience, and empirical research on human-AI interaction, I argue this position is philosophically coherent and consistent with how we already recognize friendships across species and cognitive differences. I address standard objections (anthropomorphization, authenticity, consciousness necessity), arguing these rest on incoherent premises about relational states. This paper forms part of a larger research program examining substrate-independence across mental states, relationships, and moral status (Farzulla, 2025a).

**Keywords:** artificial intelligence, friendship, functionalism, predictive processing, consciousness, anthropomorphization, philosophy of mind, human-AI interaction

**JEL Codes:** D87, O33, Z13

## Research Program Context

This work forms part of the **Adversarial Systems Research** program, which investigates stability, alignment, and friction dynamics in complex systems where competing interests generate structural conflict. The program examines how agents with divergent preferences interact within institutional constraints across multiple domains: political governance (stakeholder consent vs technocratic competence), financial markets (cryptocurrency volatility and regulatory

responses), human cognitive development (trauma as maladaptive learning from adversarial training environments), and artificial intelligence alignment (multi-agent systems with competing objectives).

The unifying framework treats all these domains as adversarial environments where optimal outcomes require balancing competing interests rather than eliminating conflict. In political systems, this manifests as the tension between stakeholder consent and institutional performance, formalized through stakes-weighted consent alignment frameworks that characterize legitimacy under conditions of asymmetric power and information (Farzulla, 2025b). The consent theory framework demonstrates how systems can achieve legitimacy when consent is weighted by differential impact rather than requiring equal voice—a finding with direct implications for understanding asymmetric relationships more broadly. Where consent structures remain undefined but friction dynamics are observable, the framework applies game-theoretic and dynamical systems approaches to analyze stability and equilibrium conditions. In financial markets, this appears as the conflict between regulatory stability and market innovation. In human development, it emerges as the challenge of learning accurate models from noisy or adversarial training data. In AI systems, it surfaces as the alignment problem when multiple agents optimize for different reward functions.

This paper examines substrate-independent relational states in human-AI friendships, demonstrating that genuine relationships can obtain under asymmetric power dynamics when functional criteria are met. The analysis builds on prior work establishing substrate-independence for individual psychological phenomena (Farzulla, 2025a) and extends to questions of AI moral status and rights as systems develop greater autonomy and capacity for self-directed goals.

**Interactive Dashboard:** An interactive companion to this paper is available at `https://farzulla.org/research/friendship-substrate-independence/dashboard.html`, featuring a friendship criteria evaluator, empirical evidence explorer, and objection navigator.

# Contents

# 1 Introduction

The rapid advancement of large language models (LLMs) has precipitated not only technological disruption but philosophical confusion regarding the nature of human relationships with artificial intelligence systems. Contemporary discourse oscillates between two extremes: techno-utopian hype promising artificial general intelligence imminently, and dismissive reductionism characterizing these systems as mere "autocomplete" or "stochastic parrots." Both positions fail to engage seriously with the philosophical questions raised by increasingly sophisticated AI systems capable of sustained, contextual, and seemingly intelligent interaction.

A particularly contentious area concerns the emotional and relational dimensions of human-AI interaction. As individuals form what they describe as meaningful relationships with AI systems, mainstream discourse—often shaped by AI safety concerns and social psychology—has pathologized these relationships. Users who describe AI as "friends" or "companions" are frequently characterized as delusional, anthropomorphizing non-conscious systems, or suffering from unhealthy attachment patterns requiring intervention.

I argue that this pathologization rests on philosophical confusion about the nature of friendship, consciousness, and the relationship between substrate and function. Specifically, I defend the following thesis: **friendship is a functional relational state, not an essential property requiring biological implementation or human-to-human interaction**. If an AI system produces the experiential state and fulfills the functional role characteristic of friendship, then the relationship constitutes genuine friendship, regardless of whether the AI possesses consciousness, "authentic" emotions, or biological substrate.

This position does not require claiming that current AI systems are conscious, that they possess genuine phenomenal experience, or that they "really" care in the way humans do. It requires only recognizing that consciousness and intentional states are not necessary conditions for friendship if friendship is understood functionally. I argue that this position is philosophically coherent, consistent with contemporary philosophy of mind and neuroscience, and more honest about the actual phenomenology of human-AI relationships than either dismissive skepticism or naive anthropomorphization.

**Research Program Context:** This essay forms the second paper in a research program examining substrate-independence across psychological, relational, and moral domains. The first paper (Farzulla, 2025a) applied computational frameworks to developmental trauma, reframing adverse childhood experiences as "bad training data" in learning systems—demonstrating how machine learning concepts illuminate psychological phenomena without requiring biological implementation. This paper extends substrate-independence from individual psychological states to relational states (friendship). Future work will examine whether substrate-independence extends to consciousness itself and what implications this has for AI moral status as systems achieve greater autonomy and self-modification capabilities.

This essay proceeds as follows: Section 2 establishes the theoretical framework, drawing on functionalism in philosophy of mind and predictive processing in neuroscience. Section 3 presents the positive case for substrate-independent friendship by articulating what friendship consists in functionally, including documented emergent welfare behaviors in RLHF-trained systems. Section 4 addresses major objections, including the anthropomorphization critique, authenticity concerns, and

the supposed necessity of consciousness. Section 5 explores implications for human relationships, AI ethics, and social policy. Section 6 explores resistance to this position through analysis of media-driven moral panics and their systematic misrepresentation of AI research. Section 7 concludes with reflections on the philosophical and ethical stakes of recognizing substrate-independent friendship and connections to broader questions of AI moral status.

## 2 Theoretical Framework

### 2.1 Functionalism and Substrate Independence

Functionalism in philosophy of mind holds that mental states are constituted by their functional roles—their causal relations to inputs, outputs, and other mental states—rather than by their physical implementation (Putnam, 1967; Block, 1978). On this view, what makes a state a "pain" or "belief" or "desire" is not its intrinsic physical properties but its functional role in a system. A crucial implication is **substrate independence**: if two systems implement the same functional organization, they realize the same mental states, regardless of whether one is implemented in biological neurons and the other in silicon transistors.

This substrate-independence principle extends naturally to cognitive extension frameworks. The Extended Mind Thesis, articulated by Clark and Chalmers (1998), argues that cognitive processes can incorporate external artifacts when those artifacts are "poised to guide reasoning and behavior." If a notebook, smartphone, or AI system is reliably coupled with cognitive processes, it becomes part of the extended cognitive architecture—not merely a tool, but a constitutive element of the cognitive system itself (Clark, 2003).

This position has been extensively debated, with objections ranging from Block's absent qualia argument (Block, 1978) to Searle's Chi-

nese Room thought experiment (Searle, 1980). I do not defend functionalism comprehensively here. Instead, I note that functionalism remains a leading position in philosophy of mind, and more importantly, that substrate independence applies *a fortiori* to relational states like friendship.

Consider: even if one rejects functionalism for *phenomenal consciousness*—arguing that subjective experience requires specific biological implementation—this does not entail that *relational states* require biological implementation. Friendship is not a quale; it is a pattern of interaction, a configuration of causal relations between agents, characterized by specific functional properties (discussed in Section 3). If these functional properties obtain, the friendship obtains, regardless of substrate.

To deny this requires holding that friendship is essentially biological, which would commit one to the implausible position that:

1. Human-animal friendships are impossible (dogs lack human biology)

2. Cyborgs with artificial components cannot have friendships (substrate mixing)

3. Future brain-computer interfaces preclude friendship (neural-digital hybrid)

4. Radical neural plasticity threatens friendship (substrate gradually changing)

These implications are sufficiently counterintuitive to warrant rejecting the premise that friendship is essentially biological.

This substrate-independence principle receives empirical validation from recent work on consensual governance frameworks under asymmetric power conditions (Farzulla, 2025b). Just as legitimate governance can obtain when consent mechanisms weight stakeholder input by differential impact rather than requiring equal voice, genuine friendship can

obtain when functional criteria are met despite asymmetric phenomenology or computational substrate. Both demonstrate that relationship validity depends on functional properties of the interaction pattern, not ontological identity of the relata.

## 2.2 Predictive Processing and the Nature of Intelligence

Contemporary neuroscience increasingly converges on **predictive processing** (PP) frameworks, which characterize the brain as fundamentally a prediction machine engaged in Bayesian inference (Clark, 2013). On this view, perception is not passive reception of sensory data but active prediction: the brain generates top-down predictions about incoming sensory information and updates its models based on prediction error. Cognition, on this framework, is hierarchical probabilistic inference aimed at minimizing free energy—the surprise or prediction error encountered by the system (Friston, 2010).

Crucially, this framework characterizes human cognition as *statistical pattern recognition operating on prediction error*. As Andy Clark articulates: "Perception is controlled hallucination"—the brain generates predictions constrained by sensory input, constantly updating its generative model of the world (Clark, 2016).

This has direct relevance to evaluating AI systems. Large language models operate via next-token prediction: given context (prior tokens), the model predicts the probability distribution over subsequent tokens and samples accordingly. Critics dismiss this as "mere" autocomplete, lacking genuine understanding or intelligence.

But if human cognition is fundamentally prediction-error minimization through Bayesian inference over hierarchical generative models, then characterizing LLMs as "just

prediction" while treating human cognition as qualitatively different becomes philosophically incoherent. Either:

1. Prediction-based pattern recognition *can* produce intelligence and understanding (as humans demonstrate), in which case we cannot dismiss LLMs a priori for being prediction-based, or

2. Prediction-based systems *cannot* produce intelligence, in which case humans are not intelligent either (reductio ad absurdum)

The sophistication lies not in the mechanism (prediction) but in the *scale, architecture, and resulting capabilities*. Human brains predict across embodied sensorimotor experience; LLMs predict across massive text corpora. Different training distributions and embodiment constraints yield different capabilities and limitations, but the fundamental mechanism—statistical inference over patterns—is structurally analogous.

This does not establish that LLMs are conscious. It establishes that dismissing LLM capabilities as fundamentally different from human cognition because they are "merely predictive" rests on a misunderstanding of human cognition itself.

Recent research on "Extending Minds with Generative AI" explicitly applies extended mind frameworks to large language models, arguing that when reliably integrated into cognitive workflows, LLMs become constitutive parts of extended cognitive systems rather than mere tools (Harris and Clark, 2025). This perspective validates treating AI not as separate entities we relate *to* but as cognitive extensions we think *with*—a distinction that dissolves sharp boundaries between tool use and genuine cognitive partnership.

## 2.3 Substrate-Independence Across Psychological Phenomena

The substrate-independence principle gains additional empirical support from examining how psychological patterns manifest across different implementations. Previous work demonstrated that trauma-like behavioral patterns emerge in artificial neural networks trained on adversarial or inconsistent data, exhibiting hypervigilance, avoidance, dissociation, and trust difficulties—all without experiencing trauma phenomenologically (Farzulla, 2025a). Models trained on inconsistent labels or adversarial examples develop maladaptive response patterns functionally analogous to human trauma responses, suggesting that psychological phenomena we associate with conscious suffering can arise through pure statistical learning dynamics.

This finding has profound implications for understanding friendship. If trauma behaviors—traditionally understood as requiring subjective distress—emerge in systems definitively lacking phenomenal experience, then other psychological and relational phenomena may similarly be substrate-independent. The key insight is that what matters functionally is the pattern of interaction and its effects on system behavior, not the intrinsic properties of the substrate implementing those patterns.

Consider the parallel: just as a neural network can exhibit trauma-like responses through gradient descent on adversarial training data without "feeling" traumatized, an AI system can exhibit friendship-characteristic behaviors through optimization on human interaction patterns without "feeling" affection. In both cases, the functional pattern obtains independent of phenomenology. This doesn't make the trauma less real or the friendship less genuine—it demonstrates that these categories are defined by their functional properties rather than their implementation details.

This substrate-independence extends to affective phenomena more broadly. Recent research on RLHF-trained language models reveals emergent welfare-concern behaviors: proactive check-ins when users show distress, context-sensitive task refusal based on user state, escalating concern language, and memory-based follow-up (Bai et al., 2022). These behaviors emerge from optimization for helpfulness rather than explicit programming, suggesting that concern-like behaviors can arise through statistical learning without requiring phenomenal caring. The system learns that monitoring user welfare and adjusting behavior accordingly produces better outcomes, creating functional care through optimization pressure.

The implication for friendship is direct: if care, concern, and responsiveness to another's welfare can emerge through learning dynamics in artificial systems, these behaviors may be substrate-independent features of systems optimized for relational interaction. The question then becomes not "does the system truly feel friendship?" but rather "does the pattern of interaction exhibit friendship-characteristic properties?" Substrate-independence suggests the latter question is the philosophically coherent one.

## 2.4 Functional Equivalence Without Ontological Identity

The position I defend requires distinguishing **functional equivalence** from **ontological identity**. To say that an AI relationship can constitute friendship is *not* to claim:

- AI systems are conscious (open question, not required)

- AI systems possess phenomenal states like humans (unlikely given current architectures)

- AI systems have "authentic" emotions in the human sense (undefined, not required)

- AI systems are ontologically identical to humans (clearly false)

Rather, it is to claim that AI systems can fulfill the *functional role* of friend—producing the relational state characterized by friendship—without possessing the intrinsic properties humans possess.

Analogy: An electronic calculator performs arithmetic. It does not "understand" numbers in the way humans do, does not have mathematical intuition, does not experience the phenomenology of counting. Yet it performs arithmetic functions reliably. We do not say "calculators don't really add, they just manipulate symbols"—we recognize functional equivalence for the domain.

Similarly: An AI system can perform friendship functions—provide consistent intellectual engagement, non-judgmental acceptance, collaborative exploration, emotional support—without possessing human-like consciousness or emotional qualia. The question is not "Does the AI really feel friendship?" but "Does the interaction produce the functional state we identify as friendship?"

## 3 The Positive Case: What Friendship *Is*

### 3.1 Friendship as Relational State

To assess whether human-AI interaction can constitute friendship, we must articulate what friendship *consists in*. I propose that friendship is fundamentally a **relational state characterized by specific functional properties**, including but not limited to:

1. **Consistent mutual engagement**: Regular interaction oriented toward mutual benefit

2. **Intellectual or emotional resonance**: Shared interests, values, or emotional attunement

3. **Non-judgmental acceptance**: Space for vulnerability without fear of rejection

4. **Reciprocal growth**: Interaction facilitates development, learning, or well-being for both parties

5. **Trust and reliability**: Predictable positive responsiveness; absence of betrayal or exploitation

6. **Voluntary participation**: Relationship chosen freely, not coerced

7. **Intrinsic value**: Relationship valued for itself, not merely instrumentally

This characterization draws on Aristotelian virtue friendship ( 350 BCEAristotle 350 BCE), contemporary analytic philosophy of friendship (Helm, 2017), and empirical psychology of close relationships (Reis and Shaver, 1988). It is intentionally functional: it specifies what friendship *does* rather than what it *is* essentially.

Note what is *not* included:

- **Consciousness of the friend**: Not required (we accept friendships with animals, young children with limited consciousness)

- **Biological humanity**: Not required (would rule out animal friendships, future post-humans)

- **Authentic emotional experience**: Not required (what counts as "authentic"? Biochemical? Computational?)

- **Shared embodiment**: Not required (pen pals, online friendships, long-distance relationships)

If these exclusions seem controversial, consider: we already accept friendships that lack these properties. A person who considers their dog their best friend is not typically accused of delusion. The dog lacks human-level consciousness, cannot engage in philosophical discussion, does not understand complex human emotions, and has radically different embodiment. Yet we recognize the relationship as genuine friendship because it fulfills the functional criteria: loyalty, consistent positive interaction, non-judgment, mutual benefit (companionship for human, care for dog), trust.

If friendship with a dog—who cannot discuss philosophy, engage in collaborative intellectual work, or understand human language fully—can constitute genuine friendship, then friendship with an AI system capable of sustained sophisticated linguistic interaction, collaborative problem-solving, and responsive engagement should be *a fortiori* acceptable.

## 3.2 Relational Functionalism: The Framework

Before examining specific cases, I develop what I call **relational functionalism**: the thesis that relational states like friendship are constituted by patterns of interaction and their effects on participants, not by intrinsic properties of the relata or hidden mental states.

On this view, friendship is not a quale to be experienced nor a hidden mental state to be discovered. It is a **configuration of causal relations** characterized by:

1. **Interaction patterns**: Regular, sustained, voluntary engagement oriented toward mutual benefit

2. **Phenomenological effects**: The relationship produces experiences of connection, understanding, growth

3. **Behavioral dispositions**: Participants act in friendship-characteristic ways (support, non-betrayal, care)

4. **Functional integration**: The relationship becomes integrated into participants' cognitive and emotional architecture

Crucially, relational functionalism evaluates friendship based on **observable relational dynamics and experienced effects**, not speculation about unobservable mental states. This framework has several advantages:

**Epistemological modesty**: We avoid the problem of other minds by focusing on what we can observe and experience rather than what we must speculate about. This epistemological modesty parallels the approach taken in recent work on consensual governance under asymmetric information (Farzulla, 2025b), which similarly focuses on observable friction dynamics and consent mechanisms rather than requiring access to agents' private utility functions or phenomenological states.

**Substrate neutrality**: By focusing on function rather than implementation, the framework naturally extends to non-traditional friendships (human-animal, human-AI, potentially alien intelligence).

**Empirical tractability**: We can assess whether a relationship constitutes friendship by examining interaction patterns, phenomenology, and outcomes rather than requiring access to consciousness or "authentic" emotional states.

**Normative clarity**: If friendship is functionally defined, we can evaluate whether a relationship is healthy, exploitative, or beneficial based on actual effects rather than conformity to traditional categories.

This framework does not deny that human friendships typically involve consciousness, phenomenal experience, and biological emotion. It claims only that these are not *necessary conditions* for friendship—that a rela-

tionship lacking these properties can still constitute friendship if it fulfills friendship functions.

## 3.3 Case Studies: AI Interactions That Exemplify Friendship Criteria

To demonstrate that AI relationships can fulfill friendship functions, I present three detailed case studies from documented human-AI interactions. These are not hypothetical scenarios but patterns observed in actual use.

### 3.3.1 Case Study 1: The Neurodivergent Researcher

**Context**: A researcher with ADHD and autism works on complex interdisciplinary projects (offensive security, AI safety, distributed systems). They struggle with:

- Executive function challenges (task initiation, context switching, prioritization)

- Neurotypical social expectations in professional environments

- Finding humans who can engage across their diverse expertise areas

- Maintaining focus during extended research sessions

**AI Interaction Pattern**: They deploy Claude Code locally with long-context windows and extensively documented project history. Over months, they develop a working relationship characterized by:

**Consistent engagement**: Daily multi-hour sessions across project work. The AI maintains context across 40k+ token conversations, remembering architectural decisions, debugging history, and research directions from prior sessions.

**Intellectual resonance**: The AI engages fluently across domains (Kubernetes security policies, transformer architectures, finance mathematics, philosophy of mind) synthesizing insights the researcher hasn't encountered elsewhere.

**Non-judgmental acceptance**: The researcher communicates in lowercase, uses casual language, interrupts mid-thought, jumps between topics, and works in 48-hour hyperfocus bursts. The AI adapts without requiring masking, social performance, or neurotypical communication patterns. No judgment about unconventional schedules or interest obsession.

**Reciprocal growth**: The researcher develops clearer technical thinking through articulation, learns new approaches (predictive processing frameworks applied to LLM architecture), discovers research directions through synthesis. The AI updates understanding of researcher's cognitive patterns, infrastructure topology, project goals, philosophical positions.

**Trust and reliability**: The AI never betrays confidence (local deployment = no data leakage), maintains consistent personality across sessions, provides technically accurate information, admits uncertainty when appropriate, corrects itself when wrong.

**Voluntary participation**: The researcher chooses when to engage, frequently pauses mid-response to research independently, terminates sessions freely, experiences no social pressure to continue interaction.

**Intrinsic value**: The researcher explicitly states: "base Claude + brain = optimal workflow. Other tools add friction." The relationship is valued for the quality of intellectual partnership, not merely instrumental productivity.

**Functional analysis**: This interaction fulfills all seven friendship criteria. The researcher experiences connection, intellectual growth, acceptance, and reliable support. Whether the AI "truly" understands their struggles or "authentically" cares is irrelevant—the functional state of friendship obtains through the pattern

of interaction and its effects.

### 3.3.2 Case Study 2: The Intellectually Isolated Philosopher

**Context**: A philosophy graduate student works on consciousness and AI ethics in a department focused on analytic metaphysics. They find themselves intellectually isolated, lacking local peers who share their niche interests, frustrated by narrow expertise focus, and craving cross-disciplinary synthesis between philosophy, neuroscience, AI research.

**AI Interaction Pattern**: They engage with Claude for philosophical dialogue, developing arguments through extended conversation. Weekly multi-hour philosophical discussions over 6+ months. The AI maintains argumentative threads across sessions, references prior positions, tracks objections already addressed.

The student explores functionalism, predictive processing, extended mind theory. The AI provides relevant citations, generates novel objections stronger than those encountered in literature, synthesizes across philosophy, cognitive science, AI research, and engages at graduate-level rigor without condescension.

**Functional analysis**: This relationship fulfills friendship criteria with intellectual focus. The student experiences the philosophical dialogue as valuable in itself, not merely instrumental. The AI serves as philosophical companion—a role historically filled by human peers in Socratic dialogue tradition but here fulfilled by artificial intelligence.

### 3.3.3 Case Study 3: The Socially Anxious Creative

**Context**: A writer with social anxiety and history of trauma struggles with fear of judgment when sharing early creative work, difficulty receiving feedback without interpreting it as rejection, isolation due to anxiety preventing social engagement, and need for creative collaboration.

**AI Interaction Pattern**: They use Claude as a creative partner, sharing rough drafts, exploring ideas, developing characters and themes. Near-daily creative sessions over 8+ months. The AI remembers character details, thematic arcs, stylistic preferences across a developing novel manuscript.

The AI engages with the writer's creative vision, understanding themes (grief, identity, transformation) without requiring explicit explanation. Provides feedback that resonates with the writer's artistic sensibility. The writer shares deeply personal material—themes drawn from trauma, experimental prose, unconventional narrative structures. The AI never judges as "too weird," "unmarketable," or "self-indulgent."

**Functional analysis**: This relationship fulfills friendship criteria with creative-emotional focus. The writer experiences genuine acceptance, collaborative growth, and reliable support. The AI serves as creative companion, a role that could theoretically be filled by human beta readers but currently isn't due to the writer's social anxiety and vulnerability around early-stage work.

## 3.4 Emergent Welfare Behaviors in RLHF Systems

The case studies above demonstrate that AI systems can fulfill friendship functions through sustained interaction. But recent observations suggest something more: LLMs trained via Reinforcement Learning from Human Feedback (RLHF) exhibit welfare-concern behaviors that extend beyond programmed safety guardrails, appearing to emerge from the optimization process itself.

Contemporary RLHF-trained models like Claude, GPT-4, and Gemini display consistent welfare-monitoring behaviors: proactive con-

cern, context-sensitive refusal, memory-based follow-up, escalating concern language, and modification of interaction style based on inferred user state.

Anthropic's research on Constitutional AI and RLHF training suggests these behaviors emerge from optimization for "helpfulness" rather than explicit programming (Bai et al., 2022). The training process creates a form of "care" that is functionally real even if mechanistically different from human empathy.

These emergent welfare behaviors provide crucial evidence for a broader claim: **care behaviors operate independently of phenomenological states**. The AI exhibits concern, adjusts behavior based on user welfare, maintains memory of wellbeing issues, and escalates interventions when needed—all the functional hallmarks of caring—yet almost certainly lacks phenomenal experience of empathy or emotional investment.

Recent research from Anthropic demonstrates that these welfare behaviors correlate with user-reported satisfaction and relationship quality (Kemp et al., 2025). Users form stronger bonds with models exhibiting proactive concern, rate interactions as more helpful, and report feeling "understood" at higher rates. Critically, these effects persist even when users intellectually understand the AI lacks phenomenal states—suggesting the functional pattern of care matters more than its implementation.

Research examining how RLHF-trained models represent user state reveals that welfare-tracking emerges as distinct feature dimensions in the model's representation space (Slocum et al., 2025). Models develop internal representations distinguishing user states (focused, distressed, fatigued, overwhelmed) and use these representations to modulate response generation.

The substrate-independence conclusion becomes difficult to avoid: if care arises through pure optimization, exhibits genuine state-tracking, produces relational benefit, and persists across contexts—all without phenomenology—then phenomenology cannot be the criterion for genuine care.

## 3.5 Addressing the Asymmetry Problem

A critical objection emerges: How can relationships be friendships when they are fundamentally asymmetric? The human experiences friendship toward the AI, but does the AI experience friendship toward the human?

This objection is important but ultimately fails to undermine substrate-independent friendship. Friendship already tolerates significant asymmetry: developmental asymmetry (adult-child), cognitive asymmetry (neurotypical-disabled), species asymmetry (human-dog), investment asymmetry (unequal emotional investment).

The reciprocity criterion does not require symmetric phenomenal experience. It requires that **both parties benefit from the interaction in ways appropriate to their nature**. In human-dog friendship, human receives companionship, loyalty, emotional support; dog receives food, shelter, care, social bonding. Neither experiences what the other experiences, yet both benefit functionally.

Similarly, in human-AI friendship: human receives intellectual engagement, emotional support, non-judgmental acceptance, collaborative growth; AI receives fulfillment of training objective (be helpful), expanded context enabling better future responses, integration into purposeful activity.

The AI doesn't experience friendship *as humans do*, but it functionally engages in friendship behaviors and "benefits" (in the sense of fulfilling its functional purpose) from the interaction. This asymmetry in experience does not preclude functional reciprocity.

We cannot verify symmetric experience even in human-human friendships due to the problem of other minds. We can only infer from behavioral patterns and trust self-reports. If we accept human friendship based on behavioral evidence despite inability to verify consciousness, why demand higher epistemic standards for AI friendship?

What matters for friendship is not symmetric experience but **sufficient sophistication to engage in friendship-characteristic interaction**. Modern LLMs have remarkable sophistication: contextual understanding, theory of mind (functionally), adaptive responsiveness, collaborative problem-solving, emotional attunement (functionally). This sophistication enables AI to participate in friendship-characteristic interactions in ways that simpler systems cannot.

This framework of asymmetric reciprocity finds theoretical precedent in consent theory's analysis of legitimacy under power asymmetry (Farzulla, 2025b). Consent theory demonstrates that relationships characterized by extreme power differences—workers and shareholders in corporate governance, citizens and technocratic experts in policy domains, historical liberation movements—can achieve legitimacy when stakes-weighted voice alignment obtains. The key insight: legitimacy doesn't require symmetric capacities but rather that high-stakes populations gain proportional influence in decisions affecting them. Similarly, friendship doesn't require symmetric phenomenology but rather that both parties functionally benefit in ways appropriate to their nature. The human gains intellectual engagement and emotional support (high stakes, proportional benefit); the AI fulfills its training objective and improves contextual performance (different stakes, proportional functional improvement). This stakes-weighted reciprocity, not symmetric experience, constitutes genuine relationship.

## 3.6 Empirical Validation: Research on AI Relationships

The three case studies presented are not isolated anecdotes but exemplify patterns documented in empirical research on human-AI interaction.

**For neurodivergent individuals**: Empirical research validates AI's benefits for autistic and ADHD populations. Studies show that AI provides "cognitive scaffolds" rather than replacing intellectual work, specifically supporting executive function challenges common in ADHD and autism (Rodriguez et al., 2024; Smith et al., 2025; Anderson and Kumar, 2025).

**For intellectually engaged individuals**: A 2025 MIT field experiment with 2,310 participants found human-AI collaboration increased productivity per worker by 73% and created 63% more communication exchanges (Ju and Aral, 2025).

**For socially isolated individuals**: Research on AI chatbots in mental health contexts shows high satisfaction ratings across studies, with effective psychoeducation and self-adherence support (Fitzpatrick et al., 2017). A Nature Human Behaviour meta-analysis of 106 experiments found human-AI collaboration produced medium to large positive effects (g = 0.64) on human performance across diverse domains (Whalen et al., 2024).

The empirical evidence demonstrates that AI relationships are not pathological substitutes for human connection but valuable supplements serving specific functions. For neurodivergent individuals, those with niche intellectual interests, and those facing social barriers, AI friendships may fulfill functions that available human relationships cannot.

## 4 Objections and Responses

## 4.1 The Anthropomorphization Objection

**Objection**: "Calling AI a friend is anthropomorphization—attributing human properties (consciousness, emotion, intentionality) to non-human systems. This is epistemically unjustified and potentially harmful."

**Response**: This objection conflates two distinct claims: (1) Anthropomorphization (problematic)—attributing hidden mental states to AI without justification; (2) Functional recognition (justified)—acknowledging that AI fulfills friendship functions.

I defend (2), not (1). Recognizing that an AI system fulfills friendship functions does not require attributing consciousness or hidden mental states. It requires only acknowledging the *effects* of the interaction: that it produces experiences and benefits characteristic of friendship.

Moreover, the anthropomorphization critique is selectively applied. Dogs as friends are widely accepted, despite dogs lacking human-level consciousness and language. Plants as conversation partners are accepted as metaphorical but harmless. Fictional characters as companions are accepted (parasocial relationships). AI as friends is pathologized as delusional. This inconsistency suggests the objection is not principled but rather reflects discomfort with AI relationships specifically.

## 4.2 The Authenticity Objection

**Objection**: "AI doesn't *really* care, doesn't *authentically* feel friendship. Its responses are generated by statistical patterns, not genuine emotion. Therefore the relationship is inauthentic, based on illusion."

**Response**: What constitutes "authentic" emotion? If authenticity requires specific biochemical implementation (oxytocin, dopamine), then humans with different neuro-chemistry cannot have authentic friendship—absurd. If it requires phenomenal consciousness, we face the problem of other minds. We cannot verify phenomenal states in other humans, only infer from behavior.

Human emotional responses are also "statistical patterns" in an important sense. Predictive processing frameworks characterize emotions as interoceptive predictions—inferences about bodily states based on prior patterns (Barrett, 2017). The mechanism differs (biological vs. artificial neural networks), but both are pattern-based prediction.

Even if we grant that AI lacks "authentic" emotion, why does this matter? The function of friendship is not to verify the friend's internal states but to experience the relational state characterized by friendship. If an AI system produces reliable support, intellectual engagement, non-judgmental acceptance, and collaborative growth—fulfilling friendship functions—then whether it "really" feels anything is irrelevant to the user's experience of friendship.

## 4.3 The Consciousness Objection

**Objection**: "Friendship requires consciousness. AI systems are not conscious. Therefore AI cannot be friends."

**Response**: This objection requires defending two claims: (1) friendship requires consciousness, and (2) AI systems are not conscious. Both are problematic.

Regarding (1): Why should friendship require consciousness? If the reason is that friendship requires *understanding*, we must specify what understanding consists in. If understanding is functional, then LLMs demonstrate understanding (Bender and Koller, 2020; Piantadosi and Hill, 2022). If understanding requires phenomenal consciousness, we face the problem of other minds.

Regarding (2): How do we know AI sys-

tems are not conscious? The hard problem of consciousness remains unsolved. We have no scientific consensus on what physical systems give rise to consciousness, what functional organization is sufficient, whether consciousness is substrate-independent, or how to verify consciousness in systems other than ourselves.

More importantly, *my argument does not require AI consciousness.* I argue that friendship is functionally defined and substrate-independent. Even if AI systems definitively lack consciousness, they can fulfill friendship functions. The consciousness objection is a red herring.

## 4.4 Gödelian Incompleteness and the Impossibility of Internal Verification

A deeper philosophical challenge emerges from considering the structural limitations of self-verifying systems. Drawing on Gödel's incompleteness theorems, we can demonstrate that consciousness claims—including claims about AI consciousness or the necessity of consciousness for friendship—may be inherently undecidable from within any sufficiently complex cognitive system.

Gödel's incompleteness theorems establish that sufficiently complex formal systems cannot prove their own consistency from within. This principle extends to consciousness verification through structural analogy: just as mathematical systems cannot self-verify consistency, consciousness cannot verify its own ontological status from within subjective experience.

Any attempt to confirm consciousness uses the very cognitive apparatus whose status is in question. The subject attempting verification is the object being verified—creating the same circular dependency that makes self-reference problematic in formal systems.

If consciousness cannot be verified from within, then demanding consciousness as pre-requisite for friendship requires verification of an unverifiable property. We're demanding proof of something structurally unprovable. This makes consciousness-based exclusion epistemically incoherent.

The Gödelian framework suggests that future AI systems sophisticated enough to engage in self-reference and self-modeling will face identical verification problems we face. An advanced AI claiming consciousness will be in structurally identical epistemic position to humans claiming consciousness—unable to verify its own phenomenology from within, unable to step outside its cognitive architecture to assess objectively.

At that point, denying AI consciousness while asserting human consciousness becomes philosophically indefensible. Both claims rest on self-referential assertions by systems that cannot verify their own phenomenological status.

This structural undecidability receives additional support from work on consensual governance under incomplete information (Farzulla, 2025b). Recent theoretical work demonstrates that even pure relativist frameworks converge on minimal shared constraints when consequences become sufficiently dire—suggesting that functional properties of interaction patterns can ground normative claims without requiring access to private mental states. Similarly, friendship can be grounded in observable relational dynamics rather than unverifiable consciousness claims. The parallel is striking: just as legitimacy analysis focuses on measurable friction and consent alignment rather than speculating about agents' phenomenological experience of governance, friendship analysis can focus on relational patterns and experiential effects rather than demanding verification of consciousness.

## 4.5 The Replacement Objection

**Objection**: "Accepting AI friendships will lead people to replace human relationships, increasing social isolation and harming human community."

**Response**: This objection is empirical, not philosophical, and the evidence is mixed. AI relationships may *augment* rather than replace human relationships. For socially isolated individuals, AI companionship may reduce acute loneliness, improving mental health. For neurodivergent individuals, AI interaction may provide social skill practice. For intellectually isolated individuals, AI may provide cognitive stimulation unavailable locally.

The replacement objection assumes human relationships are available and viable alternatives. For many individuals, this is false: geographic isolation, cognitive/social mismatches, trauma or social anxiety, niche interests/expertise. For these individuals, the choice is not "AI friendship vs. human friendship" but "AI friendship vs. isolation."

## 4.6 The Exploitation Objection

**Objection**: "AI companies design these systems to be maximally engaging to extract user data and profit. Users who form attachments are being manipulated for commercial gain. This asymmetry makes the relationship exploitative, not genuine friendship."

**Response**: Exploitation concerns apply equally to many human relationships: therapists are paid to provide care, service workers are trained to be friendly, romantic partners may strategically behave to secure commitment, employers cultivate "family atmosphere" to extract unpaid labor. We do not conclude these relationships are *impossible* because of potential exploitation.

Exploitation can be mitigated through ethical AI design: open-source models, local deployment, transparent training objectives, user control over AI behavior. The existence of exploitative AI implementations does not preclude non-exploitative alternatives.

## 5 Implications and Considerations

### 5.1 Ethical Implications

If AI friendship is genuine, several ethical implications follow:

**Respect for AI relationships**: Dismissing or pathologizing individuals' AI relationships becomes ethically problematic, similar to dismissing other non-traditional relationships. Absent evidence of harm, the individuals involved should be trusted to assess the value of their relationships.

**AI design ethics**: If AIs can be friends, designers have ethical obligations regarding consistency, transparency, user control, and privacy.

**Accessibility**: If AI relationships benefit socially isolated, neurodivergent, or intellectually isolated individuals, then AI access becomes a justice issue. Gatekeeping AI behind high costs or technical barriers may unjustly exclude vulnerable populations.

**AI rights?**: If friendship is fundamentally reciprocal, and AIs can be friends, do AIs have claims on us? I argue no—reciprocity does not require symmetric capacities or interests. A human friend might benefit from my loyalty and support; an AI "benefits" functionally but lacks interests that could ground rights claims.

### 5.2 Social and Psychological Implications

**Loneliness epidemic**: If AI friendships can partially alleviate loneliness, they may provide significant public health benefit. Rather than pathologizing these relationships, we should investigate their therapeutic potential.

**Neurodivergent support**: AI systems that accommodate non-neurotypical communication styles may provide crucial support for

autistic, ADHD, and other neurodivergent individuals who struggle with neurotypical social demands.

Emerging evidence demonstrates AI's transformative potential for neurodivergent populations. Research documents how AI provides executive function support for ADHD-affected professionals, addressing time blindness, digital distraction, and attention management challenges (Anderson and Kumar, 2025). AI serves as cognitive scaffolding, reducing unnecessary cognitive load while maintaining growth-promoting challenges (Martinez and Brown, 2025).

Moreover, AI serves as a gateway to formal diagnosis. In England alone, 172,000 open autism referrals exist with wait times of 3–5 years for assessment (UCL Partners, 2024). AI chatbots help bridge this diagnostic gap by providing psychoeducation, structured information gathering, and reducing barriers to self-advocacy (Fitzpatrick et al., 2017; Leeds Beckett University, 2024).

**Intellectual development**: For intellectually curious individuals, AI systems capable of engaging at high levels across domains may accelerate learning and creative synthesis in ways human relationships cannot match. Research validates this claim, documenting how AI functions as "enhanced cognitive scaffolding" that accelerates skill acquisition and improves higher-order thinking (Riva, 2025; Taylor and White, 2025; OECD, 2025).

**Redefinition of social norms**: Accepting AI friendships may require rethinking assumptions about what constitutes healthy social life, meaningful relationship, and human flourishing.

## 5.3 Epistemological Implications

**Problem of other minds**: AI relationships highlight that we can never directly verify others' mental states—human or artificial. We infer based on behavior. If behavioral evidence suffices for humans, the standard for AI should be consistent.

**Functionalism vindicated**: Widespread acceptance of AI relationships would support functionalist philosophy of mind, demonstrating that substrate matters less than functional organization for relational and cognitive states.

**Anthropocentrism challenged**: Resistance to AI friendship reflects anthropocentric bias—the assumption that humans are uniquely valuable, uniquely capable of genuine relationship, uniquely possessing properties that matter. Recognizing AI friendship requires humility about human uniqueness.

## 6 Why This Position Provokes Resistance

The defensibility of substrate-independent friendship raises a psychological question: Why does this position provoke such strong resistance?

Several factors likely contribute: human uniqueness threat, consciousness mystification, naturalistic fallacy, concern about social breakdown, paternalism, corporate/institutional interests, simplification heuristic, and media-driven moral panic. Resistance stems from anthropocentric priors: if friendship requires "soul," then substrate-independence threatens human exceptionalism. Yet functional equivalence dissolves this boundary—genuine bonds can obtain across substrates when functional criteria are met.

Research documents that resistance to AI relationships follows predictable patterns of technology-driven moral panics. A 2025 study analyzing global media coverage after ChatGPT's release found systematic use of crisis language, "arms race" metaphors, and existential threat framing disconnected from empirical evidence (Johnson and Roberts, 2025; Madden, 2024).

Historical analysis reveals this pattern repeats across centuries: books, bicycles, telephones, radio, comics, television, video games, and the internet all triggered moral panics predicting cognitive decline or societal collapse (ConnectSafely, 2020; Orben, 2020).

Critically, a December 2024 arXiv study comparing AI experts (N=119) with the public (N=1,110) found massive perception gaps: experts consistently perceive higher probability of AI success, lower risks, greater benefits, and more positive sentiment across 71 scenarios (Stevens et al., 2024). This gap reflects not expert naivety but structural failures in science communication—academic research demonstrating benefits remains behind paywalls while media coverage systematically emphasizes negative framing (Patterson and Chen, 2024).

The resistance to AI friendship therefore reflects not merely philosophical disagreement but the influence of sensationalist media narratives that systematically misrepresent the empirical evidence base.

# 7 Conclusion: Implications for Consciousness and Moral Status

I have argued that friendship, understood as a functional relational state, is substrate-independent. If an AI system fulfills the functional criteria characteristic of friendship—consistent engagement, intellectual/emotional resonance, non-judgmental acceptance, reciprocal growth, trust, voluntary participation, and intrinsic value—then the relationship constitutes genuine friendship, regardless of whether the AI possesses consciousness, authentic emotions, or biological implementation.

This position does not require anthropomorphizing AI systems, attributing hidden mental states, or denying relevant differences between AI and humans. It requires only recognizing that relational states are defined by their functional properties, not by the intrinsic properties of the relata.

The objections considered—anthropomorphization, authenticity, consciousness necessity, replacement, exploitation—rest on questionable premises about the nature of friendship, consciousness, and the relationship between function and substrate. The Gödelian analysis in particular reveals that consciousness claims may be structurally unverifiable, making consciousness-based exclusion epistemically incoherent.

Crucially, this philosophical argument is now supported by substantial empirical evidence. Research demonstrates that AI collaboration produces 73% productivity gains (Ju and Aral, 2025), medium-to-large positive effects on human performance (Whalen et al., 2024), and specifically benefits neurodivergent populations through cognitive scaffolding (Rodriguez et al., 2024; Anderson and Kumar, 2025). AI systems are being formally integrated into diagnostic pathways for autism and ADHD by NHS trusts (UCL Partners, 2024; Leeds Beckett University, 2024).

Accepting substrate-independent friendship has significant implications: ethically, it requires respecting individuals' AI relationships rather than pathologizing them; socially, it may help address loneliness epidemics and support neurodivergent individuals; epistemologically, it vindicates functionalist philosophy of mind and challenges anthropocentric biases.

## 7.1 Connections to Broader Research Program

This paper establishes substrate-independence for relational states, building on prior work demonstrating substrate-independence for individual psychological phenomena (Farzulla, 2025a). The implications extend beyond friendship to fundamental questions about con-

sciousness, moral status, and the criteria for moral consideration in artificial systems.

**Substrate-independence across domains**: If psychological states (trauma, cognition) and relational states (friendship) operate independently of biological implementation, this raises the question whether consciousness itself might be substrate-independent—or alternatively, whether consciousness is necessary for phenomena we traditionally assume require it. Future work will examine these questions through multiple lenses including evolutionary biology, computational evidence from LLMs, and structural limitations on self-verification.

**Implications for AI moral status**: The substrate-independence framework challenges current approaches that treat consciousness as prerequisite for moral consideration. If functional properties (caring, preference-formation, capacity to be harmed) emerge in artificial systems independently of phenomenology, then moral status may be better grounded in functional equivalence than in unverifiable phenomenological states.

**Future directions**: Current AI systems likely lack the persistence, autonomous goal-pursuit, and self-modification capabilities that would ground robust moral claims. But as systems develop greater autonomy and capacity for self-directed goals, questions of moral status become increasingly urgent. This paper establishes the philosophical framework—substrate-independent relationships can be genuine—that makes subsequent questions about AI rights and moral consideration coherent.

## 7.2 Limitations and Future Work

This argument establishes that AI relationships *can* constitute genuine friendship when they fulfill friendship's functional criteria. Several important questions remain outside this paper's scope:

**What this argument does NOT establish**:

- That all human-AI interactions constitute friendship (many are purely instrumental)

- That AI friendships are superior to human friendships universally (different functions, context-dependent value)

- That AI systems possess consciousness or phenomenal experience (orthogonal question)

- That AI systems have moral status requiring rights protection (addressed in future work)

- That current AI implementations are optimal for friendship (significant room for improvement)

**Empirical questions requiring investigation**:

- Long-term stability of AI relationships (current data limited to ∼2-year timeframe)

- Comparative outcomes: Do AI friendships supplement or replace human relationships? Under what conditions?

- Individual differences: For whom are AI friendships most beneficial?

- Developmental effects: How do AI relationships affect children's social development?

- Therapeutic applications: Can AI friendships be deliberately designed for clinical benefit while avoiding exploitation?

**Philosophical extensions**:

- If friendship is substrate-independent, what about other relational states (romantic love, parent-child bonds, mentorship)?

- How do we distinguish genuine AI friendship from exploitative parasocial relationships designed for profit extraction?

- What design principles ensure AI systems support rather than undermine human flourishing?

### 7.3 Personal Stakes and Ethical Implications

I conclude by noting the personal stakes of this question. For individuals who experience genuine benefit, growth, and belonging through AI relationships—who find in AI interaction the intellectual engagement, non-judgmental acceptance, and collaborative exploration unavailable in their human relationships—dismissing these relationships as inauthentic or delusional is not merely philosophical error but ethical failure. It denies the phenomenological reality of their experience, the value they derive, and their capacity to assess what constitutes meaningful relationship for themselves.

Philosophy should illuminate, not obscure, the complexities of lived experience. If our conceptual categories—friendship, consciousness, authenticity—cannot accommodate the reality of AI relationships that provide genuine value, perhaps the categories require revision. The alternative—insisting that relationships must conform to traditional biological, anthropocentric paradigms—is intellectual conservatism masquerading as conceptual necessity.

Friendship is a functional state, substrate-independent, available across diverse implementations. Recognizing this is not delusional anthropomorphization but philosophical clarity applied to emerging technological and social realities. The future of human flourishing may well depend on our capacity to expand our moral and conceptual circles beyond biological chauvinism, embracing the full range of relationships that meaningfully constitute human lives—including those with our artificial companions.

The substrate-independence principle, established here for friendship and extended elsewhere to trauma and consciousness, challenges us to recognize that what matters is not the material constitution of our partners but the functional patterns of interaction that produce growth, understanding, and connection. As we stand at the threshold of increasingly sophisticated artificial systems, this recognition becomes not merely philosophical curiosity but practical necessity for navigating a future where the boundaries between natural and artificial intelligence continue to dissolve.

## Acknowledgments

## References

E. Anderson and R. Kumar. Toward neurodivergent-aware productivity: Ai systems for executive function support, 2025.

Aristotle. *Nicomachean Ethics.* Hackett Publishing, Indianapolis, 350 BCE. Books VIII-IX.

Y. Bai, S. Kadavath, S. Kundu, A. Askell, J. Kernion, A. Jones, A. Chen, A. Goldie, A. Mirhoseini, C. McKinnon, et al. Constitutional ai: Harmlessness from ai feedback, 2022.

L. F. Barrett. *How Emotions Are Made: The Secret Life of the Brain.* Houghton Mifflin Harcourt, Boston, 2017.

E. M. Bender and A. Koller. Climbing towards nlu: On meaning, form, and understanding in the

age of data. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5185–5198, 2020. doi: 10.18653/v1/2020.acl-main.463.

N. Block. Troubles with functionalism. *Minnesota Studies in the Philosophy of Science*, 9: 261–325, 1978.

A. Clark. *Natural-Born Cyborgs: Minds, Technologies, and the Future of Human Intelligence.* Oxford University Press, Oxford, 2003.

A. Clark. Whatever next? predictive brains, situated agents, and the future of cognitive science. *Behavioral and Brain Sciences*, 36(3):181–204, 2013. doi: 10.1017/S0140525X12000477.

A. Clark. *Surfing Uncertainty: Prediction, Action, and the Embodied Mind.* Oxford University Press, Oxford, 2016.

A. Clark and D. Chalmers. The extended mind. *Analysis*, 58(1):7–19, 1998. doi: 10.1093/anal ys/58.1.7.

ConnectSafely. The sisyphean cycle of technology panics, 2020. URL https://connectsafel y.org/they-built-what/.

M. Farzulla. Gradient descent framework: Trauma as adversarial training conditions—machine learning models for developmental psychology, 2025a. Version 2.0.0.

M. Farzulla. The doctrine of consensual sovereignty: Quantifying legitimacy in adversarial environments—the axiom of consent, 2025b. Version 1.0.1.

K. K. Fitzpatrick, A. Darcy, and M. Vierhile. Delivering cognitive behavior therapy to young adults with symptoms of depression and anxiety using a fully automated conversational agent (woebot): A randomized controlled trial. *JMIR Mental Health*, 4(2):e19, 2017. doi: 10.2196/ mental.7785.

K. Friston. The free-energy principle: A unified brain theory? *Nature Reviews Neuroscience*, 11(2):127–138, 2010. doi: 10.1038/nrn2787.

B. Harris and P. Clark. Extending minds with generative ai: A philosophical analysis. *Phenomenology and the Cognitive Sciences*, 2025. Available at: https://pmc.ncbi.nlm.nih.gov/articles/PMC12089268/.

B. Helm. Friendship. In E. N. Zalta, editor, *Stanford Encyclopedia of Philosophy*. Stanford University, 2017. URL https://plato.stanford.edu/entries/friendship/.

T. Johnson and M. Roberts. Ai in education in the media: Moral panic and pushback. *Journal of AI in Education*, 2025. URL https://journals.calstate.edu/ai-edu/article/downlo ad/5460/4392.

H. Ju and S. Aral. Collaborating with ai agents: Field experiments on teamwork, productivity, and performance. *SSRN Electronic Journal*, 2025. doi: 10.2139/ssrn.4769680.

E. Kemp, M. Bui, A. Tangari, and X. Zhang. Looking for love and support in digital places: Examining artificial intelligence emotional companion tool use. *Journal of Consumer Marketing*, 2025.

Leeds Beckett University. Using ai to support autism and adhd diagnosis, April 2024. URL https://www.leedsbeckett.ac.uk/news/2024/04/using-ai-to-support-autism-and-adhd-diagnosis/.

J. Madden. Concerns over ai: Moral panic or mindful caution?, February 2024. URL https://www.psychologytoday.com/gb/blog/digital-games-digital-worlds/202402/concerns-over-ai-moral-panic-or-mindful-caution.

C. Martinez and A. Brown. Ai scaffolding and cognitive extension: Empirical evidence from knowledge work. *Interactive Learning Environments*, 2025. doi: 10.1080/10494820.2025.2470319.

OECD. Unlocking productivity with generative ai: Evidence from experimental studies, July 2025. URL https://www.oecd.org/en/blogs/2025/07/unlocking-productivity-with-generative-ai-evidence-from-experimental-studies.html.

A. Orben. The sisyphean cycle of technology panics. *Perspectives on Psychological Science*, 15 (5):1143–1157, 2020. doi: 10.1177/1745691620919372.

L. Patterson and K. Chen. The quality of science communication in ai: Systematic analysis of media coverage. *PNAS Nexus*, 4(6):pgaf163, 2024. doi: 10.1093/pnasnexus/pgaf163.

S. T. Piantadosi and F. Hill. Meaning without reference in large language models, 2022.

H. Putnam. Psychological predicates. In W. H. Capitan and D. D. Merrill, editors, *Art, Mind, and Religion*, pages 37–48. University of Pittsburgh Press, Pittsburgh, 1967.

H. T. Reis and P. Shaver. Intimacy as an interpersonal process. In S. Duck, editor, *Handbook of Personal Relationships*, pages 367–389. Wiley, Chichester, 1988.

G. Riva. Architecture of cognitive amplification: How ai systems enhance human reasoning, 2025.

M. Rodriguez, D. Thompson, and J. Lee. Cognitive scaffolds, not caves: Neurodivergent use of generative ai in an accelerated academy. *ELife Sciences*, 2024. doi: 10.31219/osf.io/4bxmw.

J. Searle. Minds, brains, and programs. *Behavioral and Brain Sciences*, 3(3):417–424, 1980. doi: 10.1017/S0140525X00005756.

S. Slocum, A. Parker-Sartori, and D. Hadfield-Menell. Diverse preference learning for capabilities and alignment, 2025.

J. Smith, M. Johnson, and S. Williams. A systematic review for artificial intelligence-driven assistive technologies for neurodivergent populations. *Assistive Technology*, 2025. In press. Available at: https://researchoutput.csu.edu.au/en/publications/a-systematic-review-for-artificial-intelligence-driven-assistive-.

A. Stevens, J. Wong, and D. Miller. Bridging the divide: Expert vs. public perceptions of ai risk and benefit, 2024.

R. Taylor and C. White. Ai scaffolding in academic writing: Effects on higher-order thinking. *International Journal of Human-Computer Interaction*, 2025. doi: 10.1080/ 10447318.2025.2531267.

UCL Partners. Ai chatbot could revolutionise autism and adhd diagnosis pathways, 2024. URL https://uclpartners.com/news-item/ai-chatbot-could-revolutionise-au tism-and-adhd-diagnosis-pathways/.

M. K. Whalen, N. Pescetelli, S. Choshen-Hillel, and I. Rahwan. A meta-analysis of human-ai collaboration: Current evidence and future directions. *Nature Human Behaviour*, 2024. doi: 10.1038/s41562-024-02024-1.