# COGNIPAIR: GNWT-INSPIRED COGNITIVE ARCHITECTURE FOR GENERATIVE AGENTS FOR SOCIAL PAIRING - DATING & HIRING APPLICATIONS

#### **Anonymous authors**

Paper under double-blind review

#### **ABSTRACT**

Current large language model agents lack authentic human psychological processes necessary for genuine digital twins. We present the first computational implementation of Global Workspace Theory (GNWT), creating agents with multiple specialized sub-agents (emotion, memory, social norms, planning, goal-tracking) coordinated through a global workspace broadcast mechanism. This architecture allows agents to maintain consistent personalities while evolving through social interaction. Our CogniPair simulation platform deploys 551 GNWT-Agents for speed dating interactions, grounded in real data from the Columbia University Speed Dating dataset. Evaluations show strong psychological realism, with agents achieving 72% correlation with human attraction patterns and outperforming baselines in partner preference evolution (72.5% vs. 61.3%). Human validation studies confirm our approach's fidelity, with participants rating their digital twins' behavioral accuracy at 5.6/7.0 and agreeing with their choices 74% of the time. This work establishes new benchmarks for psychological authenticity in AI systems and provides a foundation for developing truly human-like digital agents.

# 1 Introduction

Human social interactions—from dating to job interviews—require not just coherent dialogue but authentic psychological processes including emotion regulation, memory consolidation, and dynamic preference formation. LLM-based agents have been applied to model human social interactions, showing promise in domains such as customer service, healthcare assistance, and educational tutoring Park et al. (2023); Xu et al. (2023); Wang et al. (2024). Despite recent advances, current LLM-based agents face two fundamental limitations that restrict their ability to model human behavior realistically: (1) the **psychological behavior gap**—they cannot authentically simulate internal mental states, emotional processing, or evolving preferences Zhang et al. (2024); Serapio-García et al. (2023); Jiang et al. (2023); Guo et al. (2024); and (2) the **social behavior gap**—they fail to capture the complex dynamics of human-to-human interactions where preferences and behaviors co-evolve through social experiences Park et al. (2023); Zhang et al. (2018); Huang et al. (2024); Park et al. (2022); Aher et al. (2023); Sun et al. (2024).

The **psychological behavior gap** manifests in two critical problems: the **individualization problem**, where agents act like generic humans rather than specific individuals with unique psychological profiles, and the **static personality problem**, where agents cannot evolve mentally through experience. Existing approaches such as Stanford's Generative Agents Park et al. (2023) demonstrated emergent behaviors but relied on fictional personas without real human data. PersonaChat Zhang et al. (2018) introduced personality descriptions that remain synthetic and fixed. Recent personality modeling efforts Serapio-García et al. (2023); Jiang et al. (2023) achieve only surface-level behavioral mimicry without cognitive grounding. Most critically, these approaches treat personality as immutable prompts rather than dynamic psychological states shaped by experience.

The **social behavior gap** emerges when attempting to model authentic social interactions, particularly in complex domains such as relationship formation. Current LLM agents lack the capability to engage in authentic social dynamics where preferences evolve through interaction, emotional responses adapt to social feedback, and behavioral patterns shift based on interpersonal experiences. This limitation

becomes particularly pronounced in domains requiring complex social cognition, such as dating scenarios where mutual attraction emerges through dynamic, bidirectional assessment processes.

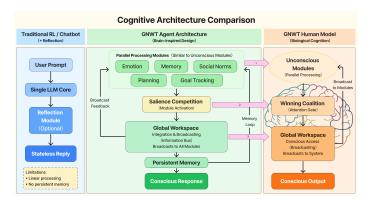


Figure 1: Comparison of Cognitive Architectures: Traditional RL/Chatbot (left), GNWT Agent Architecture (center), and GNWT Human Model (right), showing the evolution from linear processing to brain-inspired parallel processing with global workspace integration.

To address these fundamental gaps, we turn to Global Workspace Theory (GNWT) Mashour et al. (2020); Baars (2005), a leading neurocognitive model that explains how human consciousness emerges from the interaction of specialized brain modules. In human cognition, GNWT describes how disparate neural processes—emotion, memory, perception, and planning—compete for access to a central "global workspace." When information becomes sufficiently salient, it triggers a broadcast that propagates this content throughout the brain, creating our unified stream of consciousness. This theoretical foundation provides a clear roadmap for building agents that can overcome both the psychological and social behavior gaps(Figure 1). To address the **psychological behavior gap**, we operationalize GNWT into a computational agent architecture where each individual agent contains multiple specialized sub-agents working in parallel as a unified consciousness. Our GNWT-Agent implements five specialized cognitive modules—Emotion, Memory, Planning, SocialNorms, and GoalTracking—each grounded in neurocognitive theories and parameterized by the agent's Five-Factor personality profile. By implementing GNWT's broadcast mechanism computationally, we create agents with genuine internal psychological dynamics—emotion sub-agents generate affective responses, memory sub-agents consolidate experiences, social norms sub-agents manage cultural awareness, planning sub-agents develop strategies, and goal-tracking sub-agents maintain objectives. This architecture fundamentally differs from traditional LLM agents that process inputs sequentially without internal state evolution.

To bridge the **social behavior gap**, we developed **CogniPair**, a social-influence decision system that enables GNWT-Agents to engage in authentic social interactions and evolve through experience. CogniPair is not merely a testbed but a comprehensive system for modeling and guiding social influence between individuals, ultimately optimizing decision processes across various social contexts. While our primary evaluation uses a speed dating testbed, the CogniPair system itself can be extended to other social decision environments such as team formation, negotiation scenarios, and collaborative problem-solving.

We selected speed dating as our evaluation domain because it exemplifies the most challenging aspects of human social cognition—rapid compatibility assessment, dynamic preference formation, emotional regulation under uncertainty, and integration of multiple information streams. The Columbia University Speed Dating dataset Fisman et al. (2006; 2008) provides rich behavioral ground truth, including pre- and post-interaction preferences, attraction ratings, and decision outcomes, enabling rigorous evaluation of social realism. Our dating testbed deploys 551 GNWT-Agents in a two-level simulation architecture: the internal level models psychological processes within each agent, while the external level simulates social dynamics between agents. This scale—20 times larger than previous personality-based simulations Park et al. (2023)—enables statistically valid analysis of emergent social patterns, creates sufficient diversity for complex relationship networks, and allows measurement of population-level phenomena.

Through this dual approach—addressing the psychological gap with GNWT-Agents and the social gap with the CogniPair system—our framework uniquely enables both individualization and dynamic evolution. The global workspace mechanism naturally handles the stability-plasticity dilemma in personality modeling: core traits remain stable through persistent attention patterns while allowing adaptive changes through experience-driven broadcast priority shifts. Unlike previous approaches that treat personality as static prompts or surface behaviors, our architecture models the cognitive processes underlying personality and the social dynamics shaping its evolution.

Our evaluation framework measures psychological and social realism across multiple dimensions. For psychological realism, we assess preference consistency (how well agents maintain core values while adapting), emotional coherence (whether affective responses follow human psychological principles), and memory integration (how past experiences influence decisions). For social realism, we measure attraction correlation (how closely agent mate selection aligns with human patterns), interaction dynamics (how conversations evolve), and emergent social phenomena (group formation, preference shifts). These metrics ground our claims about achieving human-level authenticity in both psychological processing and social behavior.

Our experiments demonstrate that GNWT-Agents closely capture human social dynamics, showing improvements across multiple evaluation dimensions. We show significant improvements in partner preference evolution (72.5% accuracy vs. 61.3% for Multi-Agent Debate), self-perception adaptation, external evaluation changes, and match prediction (77.8% accuracy) compared to state-of-the-art baselines. Our agents demonstrate human-like evolution patterns with high correlation to ground truth data (above 0.7 across multiple dimensions, with a 0.72 correlation for match patterns). Human validation studies further confirm the psychological fidelity of our approach, with participants rating their digital twins' behavioral fidelity at 5.6/7.0 and agreeing with their twin's choices 74% of the time. Our key contributions are:

- 1. We are the **first to operationalize GNWT for computational agents**, creating a cognitive architecture where multiple sub-agents within each agent replicate human psychological processes through dynamic workspace broadcasting
- 2. We develop CogniPair, the **first social-influence decision system combining cognitive theory with large-scale social simulation**, capable of generalizing beyond our dating testbed to various social decision environments as demonstrated by successful transfer to job interview contexts (81% accuracy)
- 3. Our extensive experiments confirm GNWT's broadcast mechanism enables **genuine personality evolution**, with significant improvements in partner preference evolution (72.5% vs. 61.3%), self-perception adaptation, and external evaluation shifts compared to state-of-the-art baselines

#### 2 Related Work

LLMs for Social Simulation: Recent advances in LLMs have enabled sophisticated conversational agents Thoppilan et al. (2022); Park et al. (2023), yet standard models lack persistent psychological states Guo et al. (2024); Zhang et al. (2024) and show limitations in capturing social dynamics Xi et al. (2023); Wang et al. (2024). Reasoning enhancements such as Chain-of-Thought Wei et al. (2022), self-consistency Wang et al. (2023a), retrieval-augmentation Lewis et al. (2020), and memory architectures Guo et al. (2024); Zhang et al. (2024) improve performance but rarely incorporate selective attention Zhang et al. (2024); Zhu et al. (2023); Andreas (2022). Social simulation platforms (e.g., Generative Agents Park et al. (2023), PersonaChat Zhang et al. (2018), Li et al. Li et al. (2024), Gao et al. Gao et al. (2023)) implement memory or persona-based interaction, but typically rely on fictional characters without psychological grounding. Empirical studies Huang et al. (2024) highlight gaps between simulated and human behavior, and recent evaluations Aher et al. (2023); Park et al. (2023) note the lack of real data, neuroscience-based cognition, adaptive personality evolution, or comprehensive social metrics. Our CogniPair system integrates all these missing components.

**Modeling Psychological Processes:** Traditional cognitive architectures relied on symbolic rules, while Global Neuronal Workspace Theory (GNWT) Baars (2005); Mashour et al. (2020) offers a framework for dynamic attention. Computational implementations Bengio (2017); Goyal et al. (2022); Dehaene et al. (2017); Mashour et al. (2020) focus mainly on perception rather than higher-

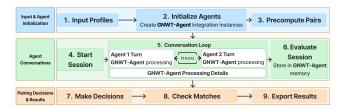


Figure 2: CogniPair Platform System Flow

order social cognition. Digital twin research emphasizes behavioral mimicry Park et al. (2024), and personality modeling systems Li et al. (2023); Wang et al. (2023b); Acerbi & Stubbersfield (2023); Sun et al. (2024) often treat traits as static. Multi-agent approaches using debate Du et al. (2023); Chan et al. (2023) or transformer-based aggregation Chen et al. (2023); Hong et al. (2023) improve over single-agent methods but rely on explicit turn-taking rather than parallel selective attention. Our GNWT-Agent instead leverages a global workspace for personality-driven prioritization, while CogniPair enables realistic preference adaptation and contextual coherence, evaluated with the Columbia Speed Dating dataset Fisman et al. (2006).

#### 3 COGNIPAIR: COGNITIVE SOCIAL PAIRING AGENT SYSTEM

In this section, we introduce single GWNT-Agent's structure (Sec. 3.1), how the single-turn conversation generated from GNWT-Agent (Sec. 3.2), and simulated social environment that allows multi-turn conversations (Sec. 3.3),

#### 3.1 GNWT-AGENT COGNITIVE MODULES

Core Cognitive Architecture: GNWT-Agent's cognitive processing is based on the Global Neuronal Workspace Theory (GNWT), which provides a computational model of human consciousness and cognitive processing mechanisms (Algorithm 1, Figure 6). In implementation, we deploy five competing cognitive modules, each focusing on different cognitive functional domains(detailed system flow in Appendix A.3)):

**Emotion Module**. This module performs in three stages: (1) emotion detection, identifying affective markers in text through feature extraction; (2) valence-arousal assessment, mapping detected emotions to a two-dimensional affective space; and (3) regulation strategy generation, adjusting emotional response intensity based on the agent's neuroticism (N) parameter. Mathematically represented as:  $R_{\rm Emotion} = f_{\rm E}(Q,H,{\rm GW},N)$ , where higher N values amplify emotional processing weights.

**Memory Module**. This module maintains a dual memory system: (1) episodic memory, storing time-stamped dialogue segments and interaction patterns; and (2) semantic memory, preserving knowledge about conversation topics and abstract concepts. The retrieval process employs vector similarity search, with the openness parameter (O) adjusting memory breadth and retrieval strategies:  $R_{\text{Memory}} = f_{\text{M}}(Q, H, \text{GW}, O)$ .

**Planning Module**. This module implements hierarchical goal decomposition: breaking down complex social goals (e.g., "establish rapport") into tactical steps (e.g., "identify common interests," "express empathy"). This process is regulated by the conscientiousness parameter (C), which controls planning depth and strategic rigor:  $R_{\text{Planning}} = f_{\text{P}}(Q, H, \text{GW}, C)$ .

**SocialNorms Module**. This module maintains a knowledge base of social interaction rules, evaluating the appropriateness of conversational behaviors. Processing includes: (1) etiquette checking, verifying response politeness; (2) boundary monitoring, preventing excessive self-disclosure; and (3) reciprocity verification, ensuring balanced conversational contributions. The agreeableness parameter (A) adjusts the strictness of norm enforcement:  $R_{\text{SocialNorms}} = f_{\text{SN}}(Q, H, \text{GW}, A)$ .

**GoalTracking Module**. This module continuously evaluates conversational progress: (1) direction monitoring, tracking advancement toward preset objectives; (2) uncertainty assessment, identifying information gaps that require clarification; and (3) direction adjustment, recalibrating goals based on interaction dynamics. The extraversion parameter (E) influences goal assertiveness:  $R_{\text{GoalTracking}} = f_{\text{GT}}(Q, H, \text{GW}, E)$ .

#### 3.2 SINGLE DECISION-MAKING SYSTEM FLOW

The CogniPair system implements a structured decision-making flow that systematically processes social interactions through nine distinct stages across three operational phases, as illustrated in Figure 2 & Algorithm 2.

Phase 1: Input & Agent Initialization establishes the foundation for social simulations. In Step 1 (Input Profiles), the system ingests personality profiles, preference distributions, and demographic data, either from human participant records (e.g., the Columbia Speed Dating dataset) or synthetically generated profiles with balanced demographic distributions. Step 2 (Initialize Agents) instantiates GNWT-Agent cognitive architectures for each participant, mapping Five-Factor personality traits to module weights and initializing the global workspace with prior knowledge. This process follows Agent<sub>i</sub> = InitializeAgent(Profile<sub>i</sub>,  $\theta_{\text{module}}$ ), where  $\theta_{\text{module}}$  represents the module-specific parameters. In Step 3 (Precompute Pairs), the system generates potential interaction dyads based on specified criteria (e.g., gender preferences, age constraints), creating a pairing matrix  $P_{m \times n}$  where each element  $p_{ij}$  indicates pairing eligibility between agents i and j.

Phase 2: Agent Conversations executes the multi-turn interactive dialogues. Step 4 (Start Session) initializes the conversational context  $C_0$  with environmental parameters (e.g., spatial configuration, temporal constraints) and interaction goals. Step 5 (Conversation Loop) implements the turn-taking dynamics where each agent processes inputs through their cognitive modules and generates responses. For each turn t, Agent 1 generates response  $R_{1,t} = \text{PROCESSINPUT}(Q_t, H_{t-1}, \text{GW}_1, P_1)$  where  $Q_t$  is the current query,  $H_{t-1}$  is the conversation history,  $GW_1$  is Agent 1's global workspace state, and  $P_1$  represents personality parameters.

Within each processing cycle, the agent integrates outputs from all cognitive modules to form a coherent response through the INTEGRATEMODULEOUTPUTS function (Algorithm 1 Line 29), formalized as a personality-weighted combination:

Response = INTEGRATEMODULEOUTPUTS(
$$\{R_M\}$$
, GW,  $P$ ) (1)

$$= \sum_{M \in Modules} \alpha_M(P) \cdot R_M + \beta(P) \cdot G(GW)$$
 (2)

where  $\alpha_M(P)$  represents personality-based module weights,  $\beta(P)$  is the integration coefficient for global workspace content, and G(GW) extracts key content from the global workspace. This combination strategy ensures the final response reflects both specialized processing from each module and maintains global coherence.

Agent 2 follows an identical process, creating a bidirectional exchange repeated for N turns. Each GNWT-Agent processing instance involves all five cognitive modules competing for global workspace access, with broadcasts occurring when salience exceeds the threshold  $\tau$ . During processing, emotional reactions, memory retrieval, planning strategies, social norm evaluations, and goal assessments are computed in parallel, with integration weighted by personality parameters.

In Step 6 (Evaluate Session), agents assess the interaction quality through multiple dimensions:  $E_i = \{E_{\text{attr}}, E_{\text{similar}}, E_{\text{comfort}}, E_{\text{interest}}\}$ , with these evaluations stored in the Memory module for subsequent retrieval. The system's adaptive learning mechanism is implemented through two key update functions: (1) the Memory.UPDATELONGTERM(Q, Response, H) function (Algorithm 1 Line 30) stores the current interaction in long-term memory, using attention-based memory consolidation techniques that highlight emotionally salient and goal-relevant content; and (2) the UPDATEPREFERENCES(P, Q, Response, H) function (Algorithm 1 Line 31) adjusts personality weights based on interaction experiences, implementing fine-tuning learning:

$$P_{t+1} = P_t + \eta \cdot \nabla_P J(P_t, Q, \text{Response}, H)$$
(3)

where  $\eta$  is a learning rate parameter and  $\nabla_P J$  is the gradient of an objective function measuring interaction success with respect to personality parameters. This dual update mechanism enables agents to continuously evolve their preferences and behaviors based on accumulated social experiences.

**Phase 3: Pairing Decisions & Results** culminates in match determinations. In Step 7 (*Make Decisions*), each agent formulates a binary decision (accept/reject) regarding potential future interactions:  $D_i = \text{DECISIONFUNCTION}(E_i, P_i, H, \text{GW}_i)$ , where the decision function integrates evaluation metrics, personality preferences, and interaction history. Step 8 (*Check Matches*) identifies

mutual matches where both agents express interest:  $M_{ij} = D_i \wedge D_j$ , creating a symmetric match matrix. Finally, Step 9 (*Export Results*) aggregates and formats simulation outcomes, including match decisions, preference evolutions, perception changes, and interaction quality metrics, generating comprehensive datasets for subsequent analysis.

This workflow implements the three complexity levels described earlier: low-complexity interactions utilize direct module selection, moderate-complexity interactions employ iterative processing with conflict resolution, and high-complexity interactions integrate the complete multi-phase protocol with comprehensive state tracking. CogniPair's decision-making and learning mechanisms are deliberately scenario-agnostic, allowing application across diverse social contexts—from optimizing information exchange in professional settings to fostering emotional connection in personal relationships. Through this systematic approach, CogniPair captures both the cognitive micromechanics of individual decision-making and the emergent macropatterns of social pairing dynamics.

#### 3.3 SIMULATED SOCIAL ENVIRONMENT SETUP

#### 3.3.1 GENERALIZED ENVIRONMENT PARAMETERIZATION

CogniPair implements a flexible parameterization system for modeling diverse social interaction environments. The system encapsulates interaction contexts through a comprehensive parameter space  ${\cal C}$  defined as:

$$C = \{ physical_p, temporal_t, social_s, cultural_c \}$$
(4)

$$physical_{p} = \{spatial\_layout, proximity, sensory\_conditions\}$$
 (5)

$$temporal_t = \{duration, pacing, sequence\_structure\}$$
 (6)

$$social_s = \{group\_size, relationship\_dynamics, power\_structure\}$$
 (7)

$$cultural_c = \{normative\_expectations, communication\_styles\}$$
 (8)

The physical parameters capture environmental conditions including spatial arrangements, interpersonal distance, and sensory factors (lighting, acoustics, temperature) that influence interaction dynamics. Temporal parameters define interaction timeframes, turn-taking pacing, and structural sequencing that shape conversational flow. Social parameters model group composition, pre-existing relationship dynamics, and authority structures. Cultural parameters encode normative behaviors and communication conventions appropriate to specific contexts.

This generalized parameterization enables CogniPair to simulate diverse interaction scenarios—from professional meetings to casual gatherings, educational exchanges to intimate conversations—by appropriately configuring these parameters. The system generates contextually-appropriate prompts using natural language templates that translate numerical parameter values into detailed environmental descriptions accessible to language models, enhancing validity across different simulation contexts.

#### 3.3.2 Multi-Agent Interaction Architecture

CogniPair's interaction architecture (Algorithm 1) provides a flexible framework for simulating multi-agent social dynamics across diverse scenarios. The system initializes a pool  $\mathcal A$  of agents, each equipped with five cognitive modules and a global workspace. The initialization process maps individual agent characteristics into module weights and interaction preferences:

$$IA_i = InitializeInteractionAttributes(P_i)$$
 (9)

$$= \{w_1, w_2, ..., w_n\} \tag{10}$$

$$Modules_i.weights = M(P_i, IA_i)$$
 (11)

where  $P_i$  represents the agent's personality profile (typically Five-Factor traits),  $\mathrm{IA}_i$  represents scenario-specific interaction attributes, and M is a mapping function that determines module processing parameters. This approach ensures individual agents retain consistent core traits while adapting their behavior appropriately to different social contexts.

The interaction protocol supports multiple engagement patterns, including dyadic exchanges (one-to-one interactions with reciprocal turn-taking), group discussions (multi-participant exchanges with dynamic speaker selection), and hierarchical interactions (structured exchanges with defined

role-based communication paths). In supporting these engagement patterns, the system manages turntaking, tracks interaction histories, and computes evolving relationship metrics, while the architecture records comprehensive data such as complete interaction histories  $\mathcal{H}$ , cognitive trace datasets  $\mathcal{T}$  capturing internal mental states, relationship development trajectories  $\mathcal{R}$  tracking interpersonal dynamics, and emergent social network structures  $\mathcal{N}$  documenting group formation. This multi-level data collection enables both micro-analysis of individual cognitive processes and macro-analysis of emergent social patterns, providing a foundation for validating the system's fidelity to human social behavior across different interaction contexts.

### 4 EXPERIMENTS AND RESULTS

#### 4.1 EXPERIMENTAL SETUP

**Dataset and Simulation Protocol:** The Columbia University Speed Dating dataset Fisman et al. (2006) contains records of 551 participants who engaged in 5,500+ four-minute speed dates over 21 sessions, resulting in over 8,300 observations. Each record includes pre-dating attribute self-ratings (1-10 scale), attribute importance ratings (distributing 100 points across 6 attributes), post-dating partner ratings on the same attributes, and match decisions (yes/no interest in seeing a partner again). The six key attributes measured are: attractiveness, sincerity, intelligence, fun, ambition, and shared interests. We instantiate 551 GNWT-Agents as digital twins of the original participants, initializing each agent's personality profile with the Five-Factor traits inferred from participants' self-ratings and importance distributions. The physical and temporal parameters in the CogniPair system are configured to match the original study's environment (bar-restaurant setting, 4-minute interaction, 8-10 conversation turns). For each simulated date, agents engage in 8 conversation turns, after which they update their self-ratings based on interaction experience, rate their partners on the six attributes, and make match decisions. We compare match patterns, preference evolution, and conversation dynamics against ground truth data from the original study.

**Baselines:** We compare against state-of-the-art approaches: Single Sequential LLM (standard prompt-based approach without specialized cognitive modules), Memory-Enhanced LLM (incorporates retrieval-augmented context) Lewis et al. (2020), Multi-Agent Debate (simulates internal deliberation through multiple agents) Chan et al. (2023), and Hierarchical Architecture (uses a command structure to organize decision-making) Du et al. (2023). All baselines use identical data initialization to ensure fair comparison.

#### 4.2 RESULTS AND KEY FINDINGS

We examine population-level social dynamics to assess how well our system bridges the social behavior gap identified in our introduction. Our experimental approach creates digital twins of the Columbia Speed Dating study participants and compares their evolution with the ground truth human data across multiple time points.

#### 4.2.1 EVOLUTION OF SOCIAL DYNAMICS

We evaluate our system's ability to model four key dimensions of social dynamics evolution: partner preference changes, self-perception adaptation, external evaluation shifts, and match decision accuracy. As illustrated in Figure 3, our analysis reveals several important patterns that characterize human social dynamics and demonstrates the remarkable similarity in evolution patterns between human participants and our GNWT-Agents.

For partner preferences (Figure 3A), both humans and agents show consistent patterns in attribute importance shifts. Attractiveness importance increases substantially (+39.0% human, +25.0% agent), while intelligence (-24.8% human, -15.2% agent) and sincerity (-16.6% human, -10.5% agent) show significant decreases. Shared interests show comparable positive changes (+9.8% human, +9.7% agent), while fun remains relatively stable with slight increases (+1.3% human, +5.8% agent). Ambition (-7.0% human, -4.5% agent) exhibits moderate decreases, suggesting its relative stability as a core value less influenced by short-term interactions.

For self-perception (Figure 3B), there is subtle but consistent calibration of traits through social interaction. Unlike partner preferences, self-perception shows more conservative adjustments, with

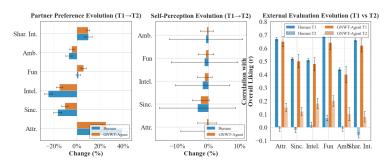


Figure 3: Social Dynamics Evolution: Human vs. GNWT-Agent Comparison. (A) Partner preference changes from T1 to T2; (B) Self-perception adjustments across attributes; (C) Evolution of attribute-liking correlations from initial to post-interaction evaluations.

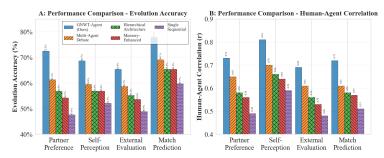


Figure 4: Comparison of GNWT-Agent with baseline methods across evolution dimensions. (A) Evolution accuracy showing GNWT-Agent's superior performance in all metrics; (B) Human-agent correlation demonstrating stronger alignment with human data compared to baseline approaches.

small negative shifts across most dimensions for both humans and agents: attractiveness (+0.3% human, -0.5% agent), sincerity (-3.5% human, -2.5% agent), intelligence (-1.9% human, -1.2% agent), fun (-1.3% human, -0.8% agent), and ambition (-0.8% human, -0.5% agent). The self-other perception gap narrows consistently for both humans and agents (from  $0.8 \rightarrow 0.7$  human,  $0.9 \rightarrow 0.7$  agent), reflecting the social calibration process through which external feedback helps align self-image with social reality.

For external evaluation (Figure 3C), there is a dramatic shift in evaluation criteria from Time 1 to Time 2. For humans, the initially strong correlations between attributes and overall liking at Time 1 (ranging from r=0.44 to r=0.69) diminish dramatically at Time 2 (ranging from r=-0.06 to r=0.07), suggesting a fundamental change in evaluation criteria following interaction. Our GNWT-Agent shows a similar pattern, with high Time 1 correlations (ranging from r=0.40 to r=0.65) decreasing substantially at Time 2, though maintaining slight positive correlations (ranging from r=0.08 to r=0.20). This pattern indicates that both humans and agents undergo significant shifts in their evaluation frameworks through social interaction, though agents retain more of their initial criteria than humans do.

Comparative Analysis: As shown in Figure 4A, our GNWT-Agent consistently outperforms all baseline methods across evolution dimensions, with particularly strong advantages in partner preference evolution (72.5% vs. 61.3% for Multi-Agent Debate) and match prediction accuracy (77.8% vs. 69.1%). Figure 4B further demonstrates our system's superior human-agent correlation, with GNWT-Agent achieving strong correlation values above 0.7 in multiple dimensions, while baseline methods fall below this threshold. The complete comparative data tables with detailed metrics and standard deviations can be found in Appendix A.33.

#### 4.3 Human Validation Studies

We validated our system by asking participants to evaluate AI versions of themselves ("digital twins") in dating and job interview contexts. Personality traits were inferred using our interactive Adventure-Based Assessment (Algorithm 3), which avoids self-report bias and correlates strongly with Big Five

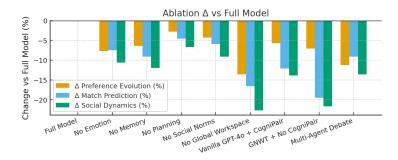


Figure 5: Ablation study: performance deltas relative to the full GNWT+CogniPair model across three metrics. Removing modules significantly reduces accuracy and correlation. Likewise, dropping the CogniPair interaction mechanism leads to large declines, confirming that both internal modules and the social platform are indispensable.

scores (r=0.82, see Appendix A.32). Table 10 summarizes results: in the Speed Dating Study (20 participants), twins were rated highly realistic (5.6/7, SD=0.8) with 74% decision concordance, while the Job Interview Study (10 participants) yielded slightly higher fidelity (5.8/7, SD=0.6) and concordance (81%). Notably, twins adapted their behavior to context, showing stronger Planning influence and professional norm adherence in interviews than in dates, confirming that our framework produces recognizable digital representations that preserve personality while adjusting appropriately to different social settings.

#### 4.4 ABLATION STUDIES

Figure 5 shows performance deltas relative to the full GNWT+CogniPair model. Removing the **global workspace** led to the largest drop (-14%), while removing **memory** or **emotion** reduced accuracy by -7%–12%. Disabling social interaction was most detrimental: **GNWT without CogniPair** fell nearly -20%, and a vanilla GPT-40 within CogniPair lost -12%. These results confirm that *cognitive* modules are complementary and that *CogniPair is indispensable* for modeling social co-evolution.

## 5 PRIVACY PROTECTION AND LIMITATIONS

Our system faces several technical limitations, with privacy being a primary concern. To mitigate risks, we apply safeguards such as differential privacy, k-anonymization, federated learning, double-consent mechanisms, and cryptographic watermarking; extended implementations are detailed in Appendix A.34. Despite these protections, long-term concerns such as behavioral prediction and misuse remain. Other limitations include imperfect calibration of cognitive modules, limited cross-cultural robustness, the absence of non-verbal communication modeling, and scalability challenges for large populations. These primarily require refinement rather than redesign, and future work will extend cultural adaptation, incorporate non-verbal cues, and improve efficiency for broader applications in collaboration, human-AI teams, and education.

### 6 Broader Impacts and Conclusion

We presented the first computational implementation of Global Workspace Theory for AI agents, demonstrating psychological realism through cognitive modules coordinated by workspace broadcasting. Our GNWT-Agent closes two gaps in current AI: the *psychological gap*, by modeling authentic internal states with dynamic preference evolution, and the *social gap*, by capturing genuine interpersonal dynamics. Experiments with 551 agents showed strong alignment with humans, achieving 72% correlation with attraction patterns and 77.8% match prediction accuracy, while human validation confirmed fidelity (5.6/7.0 realism, 74% decision concordance). Beyond these results, CogniPair enables psychologically grounded matching, improved hiring fit, and richer human-AI collaboration. These contributions lay a foundation for digital twins that preserve consistent personality while evolving through experience, advancing AI systems with genuine psychological processes.

#### ACKNOWLEDGMENTS

We thank the anonymous reviewers for the time and effort they dedicated to reading our submission and for providing thoughtful feedback. Their comments and suggestions have helped us to clarify our presentation, strengthen our arguments, and improve the overall quality of this work. We greatly appreciate their careful engagement and constructive input throughout the review process.

#### **ETHICS STATEMENT**

This work relies on the Columbia University Speed Dating dataset (Fisman et al., 2006; 2008), which is publicly available and de-identified. All analyses were conducted with strict respect for participant privacy, and no attempt at re-identification was made. To further protect sensitive information, the CogniPair framework integrates differential privacy into personality initialization, k-anonymization of conversational traces, and federated learning to avoid central storage of individual-level data.

We recognize potential ethical concerns in generating psychologically realistic digital twins, including risks of impersonation, misattribution, and misuse. To mitigate these, all generated content is cryptographically watermarked and bounded by consent mechanisms within the system design. Broader risks such as reinforcing bias or exclusion were considered, and future work will extend fairness audits to diverse demographic settings.

In line with the ICLR Code of Ethics, our position is that the societal benefits of trustworthy, explainable social simulation outweigh these risks, provided that such safeguards and ongoing monitoring are maintained.

#### REPRODUCIBILITY STATEMENT

To support reproducibility, we provide detailed descriptions of the GNWT-Agent architecture (Sec. 3.1), system flow (Sec. 3.2), and simulation setup (Sec. 3.3). Algorithms and pseudocode are included in Appendix A.2, and additional worked examples and environment parameterizations are provided in Appendix A.3–A.6. All hyperparameters (e.g., learning rate  $\eta$ , salience thresholds  $\tau$ ) are specified in the appendices, and baseline comparisons follow identical initialization for fairness. A complete record of evaluation metrics is presented in Appendix A.8. To further enable replication, all materials associated with this work will be made publicly available.

#### REFERENCES

- Alberto Acerbi and Joseph M Stubbersfield. Large language models show human-like content biases in transmission chain experiments. *Proceedings of the National Academy of Sciences*, 120(44): e2313790120, 2023.
- Gati V Aher, Rosa I Arriaga, and Adam Tauman Kalai. Using large language models to simulate multiple humans and replicate human subject studies. *International Conference on Machine Learning*, pp. 337–371, 2023.
- Jacob Andreas. Language models as agent models. In *Advances in Neural Information Processing Systems*, volume 35, pp. 17432–17444, 2022.
- Bernard J. Baars. Global workspace theory of consciousness: toward a cognitive neuroscience of human experience. *Progress in Brain Research*, 150:45–53, 2005. doi: 10.1016/S0079-6123(05) 50004-9.
- Yoshua Bengio. The consciousness prior. arXiv preprint, 2017.
- Zhiyu Chan, Yang Liu, Zhening Zhang, Shizhe Kan, Wenliang Dai, Zihan Wei, Hao You, Jieyu Qiu, Jiarui Fu, William Yang Wang, et al. Chateval: Towards better llm-based evaluators through multi-agent debate. In *Proceedings of the 37th Conference on Neural Information Processing Systems (NeurIPS)*, 2023.
- Weize Chen, Bingqing Cui, Min Shi, Yujia Wang, Tengyang Chen, Xiaohan Liang, Sirui Song, He Wang, Chenhui Zhang, Xiachong Ding, et al. Agentverse: Facilitating multi-agent collaboration and exploring emergent behaviors in agents. In *Advances in Neural Information Processing Systems* (*NeurIPS*), volume 36, 2023.
- Stanislas Dehaene, Hakwan Lau, and Sid Kouider. What is consciousness, and could machines have it? *Science*, 358(6362):486–492, 2017.
- Yilun Du, Shuang Qian, Zhuoyan Gu, Xinyun Zhu, Jyo Sang Park, Kejian Wang, Ishan Misra, Hao Zhang, Daniel Fried, Deepak Pathak, and Joshua Tenenbaum. Improving factuality and reasoning in language models through multiagent debate. In *Advances in Neural Information Processing Systems*, volume 36, pp. 26447–26472, 2023.
- Ray Fisman, Sheena S. Iyengar, Emir Kamenica, and Itamar Simonson. Gender differences in mate selection: Evidence from a speed dating experiment. *Quarterly Journal of Economics*, 121(2): 673–697, 2006.
- Ray Fisman, Sheena S. Iyengar, Emir Kamenica, and Itamar Simonson. Racial preferences in dating. *Review of Economic Studies*, 75(1):117–132, 2008.
- Chen Gao, Xiaochong Lan, Zhihong Lu, Jinzhu Mao, Jinghua Piao, Huandong Wang, Depeng Jin, and Yong Li. S3: Social-network simulation system with large language model-empowered agents. *arXiv preprint*, 2023.
- Anirudh Goyal, Aniket Didolkar, Nan Rosemary Ke, Charles Blundell, Philippe Beaudoin, Nicolas Heess, Michael Mozer, and Yoshua Bengio. Coordination among neural modules through a shared global workspace. In *International Conference on Learning Representations*, 2022.
- Jing Guo, Nan Li, Jianchuan Qi, Hang Yang, Ruiqiao Li, Yuzhen Feng, Si Zhang, and Ming Xu. Empowering working memory for large language model agents. *arXiv preprint arXiv:2312.17259*, 2024.
- Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu Zheng, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, Chenyu Ran, Lingfeng Xiao, Chenglin Wu, and Jürgen Schmidhuber. Metagpt: Meta programming for a multi-agent collaborative framework. arXiv preprint arXiv:2308.00352, 2023.
  - Yue Huang, Zhengqing Yuan, Yujun Zhou, Kehan Guo, Xiangqi Wang, Haomin Zhuang, Weixiang Sun, Lichao Sun, Jindong Wang, Yanfang Ye, et al. Social science meets llms: How reliable are large language models in social simulations? *arXiv preprint*, 2024.

- Guangyuan Jiang, Manjie Xu, Song-Chun Zhu, Wenjuan Han, Chi Zhang, and Yixin Zhu. Evaluating and inducing personality in pre-trained language models. NeurIPS 2023, 2023.
  - Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. Retrieval-augmented generation for knowledge-intensive nlp tasks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
  - Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. Camel: Communicative agents for "mind" exploration of large language model society. In *Advances in Neural Information Processing Systems*, volume 36, pp. 51991–52008, 2023.
  - Nian Li, Chen Gao, Mingyu Li, Yong Li, and Qingmin Liao. Econagent: large language modelempowered agents for simulating macroeconomic activities. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics*, pp. 15523–15536, 2024.
  - George A. Mashour, Pieter R. Roelfsema, Jean-Pierre Changeux, and Stanislas Dehaene. Conscious processing and the global neuronal workspace hypothesis. *Neuron*, 105(5):776–798, 2020. doi: 10.1016/j.neuron.2020.01.026.
  - Joon Sung Park, Lindsay Popowski, Carrie J. Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. Social simulacra: Creating populated prototypes for social computing systems. In *Proceedings of the 35th Annual ACM Symposium on User Interface Software and Technology (UIST '22)*, 2022. doi: 10.48550/arXiv.2208.04024. URL https://arxiv.org/abs/2208.04024.
  - Joon Sung Park, Joseph O'Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th annual ACM symposium on user interface software and technology*, pp. 1–22, 2023.
  - Joon Sung Park, Carolyn Q Zou, Aaron Shaw, Benjamin Mako Hill, Carrie Cai, Meredith Ringel Morris, Robb Willer, Percy Liang, and Michael S Bernstein. Generative agent simulations of 1,000 people. *arXiv preprint*, 2024.
  - Greg Serapio-García, Mustafa Safdari, Clément Crepy, Luning Sun, Stephen Fitz, Peter Romero, Marwa Abdulhai, Aleksandra Faust, and Maja Matarić. Personality traits in large language models. *arXiv preprint arXiv:2307.00184*, 2023.
  - Seungjong Sun, Eungu Lee, Dongyan Nan, Xiangying Zhao, Wonbyung Lee, Bernard J Jansen, and Jang Hyun Kim. Random silicon sampling: Simulating human sub-population opinion using a large language model based on group-level demographic information. *arXiv* preprint, 2024.
  - Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, et al. Lamda: Language models for dialog applications. *arXiv preprint*, 2022.
  - Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, et al. A survey on large language model based autonomous agents. *Frontiers of Computer Science*, 18(6):186345, 2024.
  - Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. In *International Conference on Learning Representations*, 2023a.
  - Zekun Wang, Zhihua Zhou, Qi Zhang, Zhengyan Zhu, Hongyu Jiang, Zhoujie Liu, Zhuosheng Li, Jing Han, Xiaozhi Hu, Wei Zhao, et al. Rolellm: Benchmarking, eliciting, and enhancing role-playing abilities of large language models. In *Proceedings of the 37th Conference on Neural Information Processing Systems (NeurIPS)*, 2023b.
  - Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V. Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837, 2022.

- Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, et al. The rise and potential of large language model based agents: A survey. *arXiv preprint*, 2023.
- Yuzhuang Xu, Shuo Wang, Peng Li, Fuwen Luo, Xiaolong Wang, Weidong Liu, and Yang Liu. Exploring large language models for communication games: An empirical study on werewolf. *arXiv preprint*, 2023.
- Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. Personalizing dialogue agents: I have a dog, do you have pets too? arXiv preprint arXiv:1801.07243, 2018.
- Zeyu Zhang, Xiaohe Bo, Chen Ma, Rui Li, Xu Chen, Quanyu Dai, Jieming Zhu, Zhenhua Dong, and Ji-Rong Wen. A survey on the memory mechanism of large language model based agents. *arXiv* preprint, 2024.
- Xizhou Zhu, Yuntao Chen, Hao Tian, Chenxin Tao, Weijie Su, Chenyu Yang, Gao Huang, Bin Li, Lewei Lu, Xiaogang Wang, et al. Ghost in the minecraft: Generally capable agents for open-world environments via large language models with text-based knowledge and memory. *arXiv preprint*, 2023.

#### **APPENDIX** Page 15 A.1 Agent Architecture Internal structure and module flow. A.2 Algorithms and Pseudocode Page 16 Key algorithmic procedures and decision logic. A.3 System Prompt Flow Page 19 Breakdown of prompt construction and token broadcast paths. A.4 Worked Example: Interview Anxiety Page 24 End-to-end walk-through of a real-world query scenario. A.5 Use Case: Dating Application Page 32 Evaluation results in interpersonal matching simulations. A.6 Use Case: Job Application Page 37 Adaptation of the framework for career placement tasks. A.7 Adventure-Based Personality Assessment Page 44 Roleplay-driven method to uncover latent traits. A.8 Detailed Results Page 53 Full quantitative results and evaluation tables. A.9 Privacy Protection Page 54 Technical safeguards and ethical frameworks for digital twin deployment.

#### USE OF LLMS

This work is centered on the use of large language models (LLMs) as core agents within the CogniPair framework. Specifically, we employed [GPT-4/other models, anonymized for review] to instantiate GNWT-based digital twins, simulate conversations, and support social pairing experiments. The LLMs were accessed through standard APIs and were not fine-tuned on private or sensitive data. Their role is not limited to auxiliary writing or coding assistance but constitutes the main experimental system under study.

In line with ICLR 2026 policy on LLM usage, we note that all analyses, design choices, and interpretations of results are the responsibility of the authors. The LLMs may introduce biases or limitations inherent to their training data, and these factors have been considered when analyzing outcomes.

#### A.1 AGENT ARCHITECTURE

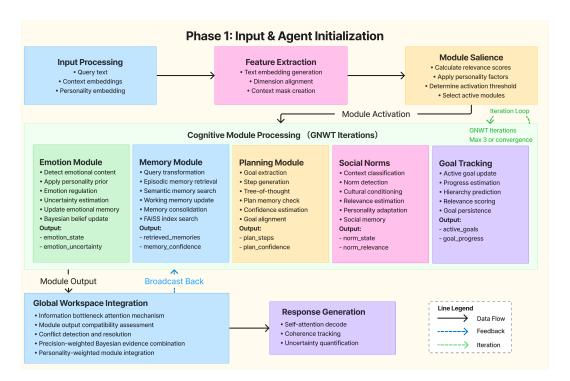


Figure 6: GNWT-Agent Architecture Internal Processing Flow

#### A.1.1 LLM DETAILS

We used the OpenAI API for GPT-40 with top\_p set to 1, max\_tokens set to 200, min\_tokens set to 0, and temperature set to 0.9 (with all other parameters at their default values),

29:  $P \leftarrow \text{UPDATEPREFERENCES}(P, Q, \text{Response}, H)$ 

30: **return** Response

#### A.2 ALGORITHMS AND PSEUDOCODE

810

811

839

840 841

861862863

#### 812 Algorithm 1 GNWT-Agent: GNWT-Based Social Pairing Agent 813 **Require:** User query Q, conversation history H, agent personality profile P, maximum iterations T, conver-814 gence threshold $\epsilon$ 815 1: Modules ← {Emotion, Memory, Planning, SocialNorms, GoalTracking} 816 2: Modules $\leftarrow$ InitializeModules(P) ▶ Initialize modules based on personality traits 817 3: GW $\leftarrow$ InitializeGlobalWorkspace(P, H) 818 4: SC ← InitializeSalienceCalculator(P) > Store previous module responses 5: $R_{\text{prev}} \leftarrow \emptyset$ 819 6: for $t = 1 \rightarrow T$ do 820 **parallel for** each module $M \in Modules$ : 7: 821 $R_M \leftarrow M.PROCESS(Q, H, GW)$ 822 9: for all $M \in Modules do$ 10: $S_M \leftarrow \text{SC.evaluate}(R_M, P)$ 823 11: end for 824 12: $M^* \leftarrow \arg\max_M S_M$ 825 $C \leftarrow R_{M^*}$ 13: 826 if $\exists M_i, M_j : \operatorname{Conflict}(R_{M_i}, R_{M_j}) > \operatorname{threshold}$ then 14: 15: $C \leftarrow \text{RESOLVECONFLICT}(\mathring{R}_{M_i}, R_{M_i}, \text{GW}, P)$ ▶ Resolve module conflicts end if 16: 828 17: IGNITE(GW, C)829 BROADCAST(C) to all modules 18: 830 19: for all $M \in Modules do$ 831 20: $M.\mathsf{UPDATE}(C,\mathsf{GW})$ □ Update module state 832 21: end for if $\max_M |R_M - R_{\text{prev},M}| < \epsilon$ or t = T then 833 22: ▶ Minimal change in responses 23: break 834 24: end if 835 25: $R_{\text{prev}} \leftarrow \{R_M | M \in \text{Modules}\}$ > Save current responses for next iteration 836 **26**: **end for** 837 27: Response $\leftarrow$ INTEGRATEMODULEOUTPUTS( $\{R_M | M \in Modules\}, GW, P$ ) 28: Memory.UPDATELONGTERM(Q, Response, H)▶ Update long-term memory 838

> Adaptively update preferences

▶ Return final integrated response

#### 864 Algorithm 2 CogniPair: Speed Dating Cognitive Simulation System 865 **Require:** Agent profiles $\mathcal{P} = \{P_1, P_2, \dots, P_n\}$ , batch size b, cognitive parameters $\Theta$ 866 **Ensure:** Dating results $\mathcal{R}$ , cognitive trace data $\mathcal{T}$ , matches $\mathcal{M}$ 867 1: $\mathcal{A} \leftarrow \emptyset$ 868 2: $\mathcal{M} \leftarrow \emptyset$ 3: $\mathcal{T} \leftarrow \emptyset$ 4: for $P_i \in \mathcal{P}$ do 870 $Modules_i \leftarrow \{Emotion, Memory, Planning, SocialNorms, Attraction\}$ 871 ▷ Initialize global workspace $GW_i \leftarrow InitializeGlobalWorkspace(P_i)$ 872 7: $DA_i \leftarrow InitializeDatingAttributes(P_i)$ 873 8: $A_i \leftarrow \text{CognitiveAgent}(P_i, \text{Modules}_i, \text{GW}_i, \text{DA}_i, \Theta)$ $\mathcal{A} \leftarrow \mathcal{A} \cup \{A_i\}$ 874 10: **end for** 875 11: $\mathcal{P}_{pairs} \leftarrow \text{GenerateCompatiblePairs}(\mathcal{A})$ ▶ Based on gender/orientation 876 12: $\mathcal{B} \leftarrow \text{BatchPairs}(\mathcal{P}_{pairs}, b)$ 877 13: for batch $B \in \mathcal{B}$ do 878 $Results_B \leftarrow \emptyset$ 14: for pair $(A_i, A_j) \in B$ do parallel 15: ▷ Process pairs in parallel 879 16: $H_{ij} \leftarrow \emptyset$ 880 $C_{ij} \leftarrow \text{InitializeContext}(A_i, A_j)$ 17: for $r=1 o {\rm MAX\_ROUNDS}$ do 18: 882 19: $Q_i \leftarrow A_i$ .GENERATEQUERY $(H_{ij}, C_{ij})$ 883 20: $\mathcal{T} \leftarrow \mathcal{T} \cup A_i.\text{GW.GETTRACE}()$ 21: $H_{ij} \leftarrow H_{ij} \cup \{(i,Q_i)\}$ ▶ Update conversation history 884 22: $R_j \leftarrow A_j$ .GENERATERESPONSE $(Q_i, H_{ij}, C_{ij})$ ▶ Using Alg. 1 885 23: $\mathcal{T} \leftarrow \mathcal{T} \cup A_i.\text{GW.GETTRACE}()$ 886 $H_{ij} \leftarrow H_{ij} \cup \{(j, R_i)\}$ 24: ▶ Update conversation history 887 if $r < MAX\_ROUNDS$ then 25: 888 26: $Q_j \leftarrow A_j$ .GENERATEQUERY $(H_{ij}, C_{ij})$ 27: $\mathcal{T} \leftarrow \mathcal{T} \cup A_i.\text{GW.GETTRACE}()$ 889 28: $H_{ij} \leftarrow H_{ij} \cup \{(j,Q_j)\}$ 890 29: $R_i \leftarrow A_i$ .GENERATERESPONSE $(Q_j, H_{ij}, C_{ij})$ Dusing Alg. 1 Dusing Alg. 1 891 30: $\mathcal{T} \leftarrow \mathcal{T} \cup A_i.\text{GW.GETTRACE}()$ 892 31: $H_{ij} \leftarrow H_{ij} \cup \{(i, R_i)\}$ 893 32: end if 33: $A_i$ . UPDATEATTRACTION $(H_{ij}, A_j)$ ▶ Update attraction dynamics 894 34: $A_i$ . UPDATEATTRACTION $(H_{ij}, A_i)$ 895 35: end for 896 36: $S_i \leftarrow A_i$ . EVALUATE COMPATIBILITY $(A_i, H_{ij})$ ⊳ Final decision 897 37: $S_j \leftarrow A_j$ . Evaluate Compatibility $(A_i, H_{ij})$ 898 38: $Results_B \leftarrow Results_B \cup \{(A_i, A_j, S_i, S_j, H_{ij})\}$ 39: if $S_i \geq \text{THRESHOLD} \land S_j \geq \text{THRESHOLD}$ then 899 40: $\mathcal{M} \leftarrow \mathcal{M} \cup \{(A_i, A_j)\}$ ▶ Record match 900 41: 901 $A_i$ . UPDATEPREFERENCES $(A_j, H_{ij}, S_i, S_j)$ 42: ▶ Preference evolution 902 43: $A_j$ . UPDATEPREFERENCES $(A_i, H_{ij}, S_j, S_i)$ 903 44: end for 45: $\mathcal{R} \leftarrow \mathcal{R} \cup \text{Results}_B$ 904 46: **end for** 905 47: **return** $(\mathcal{R}, \mathcal{T}, \mathcal{M})$ 906

#### 918 Algorithm 3 Adventure-Based Personality Assessment 919 1: Initialize personalityProfile and traitConfidence 920 ▷ Cloud (GPT-4o) or Local (Ollama) 2: modelPreference = GetUserModelPreference() 921 3: while more scenarios needed AND trait coverage insufficient do 922 Select and present next scenario from pool 923 5: Collect user's choice 924 if modelPreference == "local" then 6: 925 7: Analyze choice using Ollama (llama3 or deepseek-r1) 926 8: else 9: Analyze choice using GPT-40 927 10: end if 928 11: Update personality traits based on LLM analysis 929 12: Generate follow-up question based on user's choice 930 13: Collect user's free-text response 931 14: if modelPreference == "local" then 932 15: Analyze free text using Ollama 933 16: else 934 17: Analyze free text using GPT-40 935 18: 936 19: Update personality traits based on text analysis 937 20: end while 21: Normalize and validate final personality profile 938 22: return finalProfile 939 940

#### A.3 DETAILED SYSTEM FLOW ANALYSIS: GNWT-AGENT COGNITIVE ARCHITECTURE

This appendix presents a step-by-step analysis of information flow through the GNWT-Agent cognitive architecture. GNWT-Agent's central innovation is its hybrid neural-symbolic approach that combines specialized neural modules with LLM reasoning via a global workspace mechanism, implementing a neurobiologically-informed cognitive architecture.

#### A.4 FORMAL ARCHITECTURE DEFINITION

The GNWT-Agent cognitive architecture is formally defined as a quintuple:

$$\mathcal{E} = (\mathcal{M}, \mathcal{W}, \mathcal{I}, \mathcal{L}, \mathcal{P}) \tag{12}$$

Where:

 $\mathcal{M} = \{M_{\text{emo}}, M_{\text{mem}}, M_{\text{plan}}, M_{\text{norm}}, M_{\text{goal}}\}$  represents the set of specialized cognitive modules

W denotes the global workspace integration mechanism

 $\mathcal{I}$  signifies the information bottleneck attention system

 $\mathcal{L}$  represents the language model interface

 $\mathcal{P}$  characterizes the personality representation space

Each cognitive module  $M_i \in \mathcal{M}$  implements a hybrid neural-symbolic architecture:

$$M_i = (\mathcal{N}_i, \mathcal{L}_i, \mathcal{T}_i, \mathcal{D}_i, \mathcal{I}_i, \mathcal{S}_i)$$
(13)

Where:

 $\mathcal{N}_i$  denotes the neural processing component

 $\mathcal{L}_i$  signifies the module-specific LLM component

 $\mathcal{T}_i$  represents the tensor-text conversion mechanism

 $\mathcal{D}_i$  denotes the differentiable memory system

 $\mathcal{I}_i$  signifies the module interface specification

 $S_i$  represents the salience computation function

Table 1: Specialized Modules and Their Neuroanatomical Bases

Module	Function	Neural Inspiration
Emotion $(M_{\rm emo})$	Affective processing	Limbic system, amygdala, insula
Memory $(M_{\text{mem}})$	Information retrieval	Hippocampus, temporal cortex
Planning $(M_{plan})$	Structured reasoning	Frontopolar cortex, DLPFC
Social Norms $(M_{\text{norm}})$	Social context	mPFC, TPJ
Goal Tracking $(M_{\text{goal}})$	Hierarchical goals	OFC, ACC

#### A.5 Information Flow Process

#### A.5.1 INITIAL TEXT ENCODING

The information flow begins with the transformation of text inputs into neural representations:

$$e_{Q} = \phi_{\text{embed}}(Q) \in \mathbb{R}^{d}$$

$$e_{H} = \{\phi_{\text{embed}}(h_{i}) | h_{i} \in H\} \in \mathbb{R}^{n \times d}$$

$$e_{P} = \psi(p) \in \mathbb{R}^{d_{p}}$$
(14)

Where  $\phi_{\text{embed}}$  is the embedding model that converts text to dense vectors, Q is the query text, H is the conversation history, and p is the personality profile.

1026 A.5.2 FEATURE EXTRACTION AND EMBEDDING ALIGNMENT 1027 1028 The raw embeddings undergo feature extraction and alignment: 1029 1030  $e_Q' = \mathsf{FeatureExtractor}(e_Q)$ 1031 (15) $e'_H = \{ \text{FeatureExtractor}(e_h) | e_h \in e_H \}$ 1032 1033 MODULE SALIENCE CALCULATION A.5.31034 1035 Each module calculates its relevance to the current input: 1036 1037  $s_i = \mathcal{S}_i(e'_O, e'_H, e_P, \mathcal{G})$  $= \alpha_i + \sum_j \beta_{ij} \cdot f_{ij}(e'_Q, e'_H, e_P, \mathcal{G})$ 1039 (16)1040 1041 Where  $\alpha_i$  is the baseline salience,  $\beta_{ij}$  are weighting coefficients,  $f_{ij}$  are feature extractors, and  $\mathcal{G}$  is 1043 the current global workspace state. 1044 1045 A.5.4 PARALLEL MODULE PROCESSING 1046 Each module independently processes the input using a hybrid neural-LLM approach with the 1047 following general pattern: 1048 1049 1. Neural Processing: Extract relevant features and apply module-specific transformations 1050 1051 2. Tensor-to-Text Conversion: Convert neural representations to LLM-readable format 1052 3. **LLM Processing**: Generate structured symbolic representations using prompting 1053 4. **Text-to-Tensor Conversion**: Transform LLM outputs back to neural representations 1054 1055 5. **Output Integration**: Combine neural and symbolic components for module output 1056 1057 Module-specific processing includes: 1058 • Emotion Module: Detects emotional states and regulation strategies • Memory Module: Retrieves and integrates episodic, semantic, and working memories 1061 • Planning Module: Generates structured plans for achieving identified goals 1062 1063 Social Norms Module: Identifies appropriate social contexts and behavioral norms 1064 Goal Tracking Module: Maintains hierarchical goal representations and tracks progress A.5.5 GLOBAL WORKSPACE INTEGRATION 1067 1068 Module outputs are projected to a common workspace dimension and integrated based on salience: 1069 1070 A.5.6 FINAL PROMPT CONSTRUCTION 1071 The global workspace state informs the construction of the final prompt for the response LLM, 1072 transforming module outputs into a coherent instruction format: 1074 1075

$$\begin{split} P_{\text{system}} &= \text{JoinWithNewlines}(\text{PersonalityPrompt}(e_P), \\ &\qquad \qquad \text{ModulePrompts}(\{\text{output}_i\}, \{\text{adjusted\_weights}_i\}), \\ &\qquad \qquad \text{StrategyPrompt}(\text{integrated\_state}), \\ &\qquad \qquad \text{ConflictPrompt}(\text{conflicts\_resolved}), \\ &\qquad \qquad \text{ResponsePlanPrompt}(\text{response\_plan})) \end{split}$$

1077

1078

#### Algorithm 4 Global Workspace Integration

```
1081
           1: module_rep<sub>i</sub> \leftarrow ProjectToWorkspace(output<sub>i</sub>) for each module i
1082
           2: normalized_salience \leftarrow Softmax(s)
           3: personality_weights \leftarrow Softmax(PersonalityToIntegration(e_P))
1084
           4: combined_weights ← normalized_salience+personality_weights
           5: for each pair of modules (i, j) where i \neq j do
1086
                   conflict_scores<sub>i,i</sub> \leftarrow ConflictDetector([module_rep<sub>i</sub>, module_rep<sub>i</sub>])
1087
           7: end for
1088
           8: Apply conflict resolution adjustments to weights
1089
           9: integrated_output \leftarrow \sum_{i} adjusted_weights<sub>i</sub> · module_rep<sub>i</sub>
1090
          10: if major conflicts detected then
1091
          11:
                   Apply specialized conflict resolution
          12: end if
          13: workspace state ← WorkspaceProjection(integrated output)
1093
          14: Perform integration LLM call to further refine integration
1094
```

#### A.5.7 RESPONSE GENERATION AND MEMORY UPDATE

The final response is generated and memory systems are updated:

$$\begin{aligned} & \text{history} = [\text{prev\_messages}] \\ & \text{query\_message} = \{\text{"role": "user", "content": }Q\} \\ & \text{full\_history} = \text{history} + [\text{query\_message}] \\ & R = \text{LLM}_{\text{response}}(P_{\text{system}}, \text{full\_history}, \text{temperature} = 0.7) \end{aligned} \tag{18}$$

$$\label{eq:conding} \begin{aligned} & \text{EpisodicMemory.write}(\text{Embed}(Q,R), \text{StructuredEncoding}(Q,R, \text{workspace\_state})) \\ & \text{WorkingMemory.update}(Q,R) \end{aligned} \tag{19}$$

#### A.6 PROCESSING EXAMPLE

The following example demonstrates the architecture's response to: "I'm really nervous about my job interview tomorrow. I've been preparing for weeks, but I still don't feel ready. What should I do?"

#### Module Salience Calculation

#### **Salience Scores:**

- Emotion Module: 0.75 (anxiety markers, moderate neuroticism, emotional shift)
- Planning Module: 0.67 (clear goal, structured task, high conscientiousness)
- Memory Module: 0.52 (relevant experiences, moderate openness)
- Social Norms Module: 0.48 (professional context, formal situation)
- Goal Tracking Module: 0.45 (career advancement, interview milestone)

Normalized weights: [0.31, 0.28, 0.21, 0.12, 0.08]

#### 1134 Global Workspace Integration 1135 1136 **Conflict Detection:** 1137 • Primary conflict: Emotion module (anxiety validation) vs. Planning module (action steps) 1138 **Adjusted Weights:** [0.29, 0.25, 0.22, 0.14, 0.10] 1139 **Integration Results:** 1140 • Primary response focus: balanced emotional and practical support 1141 • Emotional context: anxiety (valence: -0.48, arousal: 0.68) with determination 1142 • Response strategy: validate anxiety then provide structured help 1143 1144 • Conflict resolution: Balance validation with actionable steps 1145 Response Plan: 1146 1. Validate interview anxiety as normal 1147 2. Acknowledge preparation effort positively 1148 3. Provide targeted high-impact preparation steps 1149 4. Include confidence-building perspective 1150 1151 5. End with encouragement and well-wishes 1152 1153 Final Response 1154 1155 It's completely normal to feel nervous before an interview, especially for a company you're really excited 1156 about! The fact that you've been preparing for weeks shows your dedication, and that preparation won't 1157 disappear overnight, even if you're feeling anxious right now. At this point, here's what I'd suggest focusing on: 1158 1. Review your strongest points - skim through the key qualifications and 2-3 accomplishment stories 1159 you want to highlight. Don't try to memorize everything. 1160 2. Practice answering 3-5 common questions out loud - this builds muscle memory that helps even 1161 when nerves kick in. 1162 3. Prepare your logistics for tomorrow - lay out your outfit, plan your route, and gather any materials you need to bring. 1163 4. Most importantly, make sure you get adequate rest tonight. A clear mind will serve you better than a 1164 few more hours of anxious preparation. 1165 Remember that interviews are a two-way conversation, not an examination. You've been preparing for 1166 weeks - now it's time to trust that preparation and focus on presenting your authentic self. 1167 Good luck tomorrow! You've put in the work, and that's the most important part. 1168 1169 1170 A.7 MODULE LLM CALL ARCHITECTURE 1171 The GNWT-Agent architecture employs a cascade of LLM calls with neural pre-processing and 1172 post-processing: 1173 1174 1175 $\label{eq:module-Level LLM Calls} \mbox{Module-Level LLM Calls}: \{\mbox{LLM}_{\mbox{emo}}(P_{\mbox{emo}}), \mbox{LLM}_{\mbox{mem}}(P_{\mbox{mem}}), \mbox{LLM}_{\mbox{plan}}(P_{\mbox{plan}}),$ 1176 (20) $LLM_{norm}(P_{norm}), LLM_{goal}(P_{goal})$ 1177 1178 1179 $Integration \ LLM \ Call : LLM_{integration}(P_{integration})$ (21)1180 1181 1182 Response LLM Call : $LLM_{response}(P_{system}, Q)$ (22)1183

 $Input \rightarrow Neural\ Processing \rightarrow Module\ LLMs \rightarrow Integration\ LLM \rightarrow Response\ LLM \rightarrow Memory\ Update \eqno(23)$ 

11841185

1186

1187

This creates a complete cognitive cycle:

#### A.8 SUMMARY

The GNWT-Agent cognitive architecture implements a hybrid neural-symbolic approach that integrates multiple specialized modules within a neurobiologically-informed framework. The architecture's core components include five specialized cognitive modules addressing emotion, memory, planning, social norms, and goals.

The information flow through the system follows a comprehensive process:

- 1. Initial text encoding via dense vector representations
- 2. Parallel processing across specialized modules
  - 3. Dynamic salience-based resource allocation
  - 4. Global workspace integration with conflict detection and resolution
  - 5. Structured prompt construction
  - 6. Final response generation
  - 7. Memory update and continuous learning

Table 2: Strengths and Limitations of GNWT-Agent Architecture

Strengths	Limitations
Enhanced interpretability through explicit module contributions	Computational complexity due to multiple LLM calls
Improved uncertainty handling via Bayesian uncertainty propagation	Challenges in consistent knowledge representation
Consistent personality representation through dedicated space	Need for further empirical validation
Dynamic adaptation through salience-based processing	Potential for response latency issues
Neurobiologically-inspired cognitive processing	Complex architecture requiring extensive fine-tuning

In conclusion, the GNWT-Agent cognitive architecture represents a significant step toward creating AI systems that not only process language effectively but do so through mechanisms that more closely approximate human cognitive processes, potentially leading to more natural, adaptive, and comprehensible AI interactions.

#### A.9 WORKED EXAMPLE: INTERVIEW ANXIETY QUERY

This section presents a detailed step-by-step walkthrough of the complete GNWT-Agent cognitive architecture processing flow for a single example query. By tracing the transformations from raw input to final response, we provide a concrete illustration of the theoretical architecture described in previous sections.

#### A.10 EXAMPLE INPUT QUERY

The example query represents a common scenario of pre-interview anxiety:

#### User Input

"I'm really nervous about my job interview tomorrow. I've been preparing for weeks, but I still don't feel ready. What should I do?"

#### A.11 INITIAL TEXT ENCODING AND EMBEDDING

The first transformation converts the raw text into numerical representations:

Input 
$$\xrightarrow{\phi_{\text{embed}}}$$
 Query Embedding (24)

"I'm really nervous..." 
$$\rightarrow [0.086, -0.143, 0.257, \dots, 0.112] \in \mathbb{R}^{768}$$
 (25)

Previous conversation context is also encoded:

"I've been job hunting for months" 
$$\rightarrow [0.141, 0.092, \ldots] \in \mathbb{R}^{768}$$
 (26)

"I finally got an interview at my dream company" 
$$\rightarrow [0.235, -0.124, \ldots] \in \mathbb{R}^{768}$$
 (27)

Feature extraction enhances these raw embeddings:

Enhanced query embedding: 
$$[0.127, -0.086, 0.313, ...] \in \mathbb{R}^{768}$$
 (28)

#### A.12 MODULE SALIENCE CALCULATION

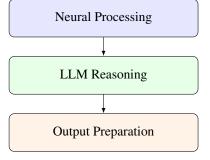
Each module calculates its relevance to the query through multi-factor salience functions:

Module	Primary Factors	Values	Raw Salience	Normalized	
	Emotional intensity	0.83			
Emotion	Personality factor	0.52	0.75	0.31	
	Emotional change	onal change 0.60			
	Goal clarity	0.78			
Planning	Task structure	0.82	0.67	0.28	
C	Personality factor	0.76			
Memory	Memory match	0.64			
	Personality factor	0.61	0.52	0.21	
	Recency boost	0.35			
Social Norms	Norm relevance	0.71			
	Formality	0.70	0.48	0.12	
	Personality factor	0.63			
Goal Tracking	Goal relevance	0.76			
	Personality factor	0.65	0.45	0.08	
	Progress factor	0.58			

Table 3: Module salience calculation for the interview anxiety query

The emotional content of the query results in the Emotion Module having the highest salience (0.31), followed closely by the Planning Module (0.28).

A.13 PARALLEL MODULE PROCESSING A.13.1 EMOTION MODULE PROCESSING The Emotion Module performs neural processing followed by LLM reasoning: 



- Context attention weights: [0.65, 0.35] - Emotion detection: [0.82, 0.65, 0.15, ...]
- Primary emotions: anxiety, nervousness
- valence: -0.48 - arousal: 0.68
- primary\_emotion: "anxiety"
- regulation\_strategy: "validation\_with\_reframing"
- Emotion state tensor: [-0.48, 0.68, ...]
- Regulation applied: [-0.45, 0.70, ...]

Figure 7: Emotion Module processing flow

#### A.13.2 PLANNING MODULE PROCESSING

The Planning Module extracts the goal and generates structured steps:

```
1350
1351
1352
1353
           "requires_planning": true,
1354
           "planning_type": "preparation_strategy",
           "steps": [
1355
1356
                "content": "Review core interview preparation materials",
1357
                "confidence": 0.93
1358
             },
1359
1360
                "content": "Practice answering common questions aloud",
1361
                "confidence": 0.89
1362
             },
1363
1364
                "content": "Prepare concise examples of achievements",
1365
                "confidence": 0.86
1366
             },
                "content": "Implement anxiety reduction techniques",
1368
                "confidence": 0.82
1369
             },
1370
1371
                "content":
1372
                "Plan interview logistics (route, outfit, materials)",
1373
                "confidence": 0.91
1374
             },
1375
1376
                "content": "Get adequate rest before interview",
                "confidence": 0.85
1377
1378
             }
           ],
1379
           "plan_feasibility": 0.88,
1380
           "goal_alignment": 0.84,
1381
           "plan_uncertainty": 0.16
1383
1384
```

#### A.13.3 MEMORY MODULE PROCESSING

The Memory Module retrieves relevant episodic and semantic memories:

```
1404
1405
1406
1407
           "episodic memories": [
1408
             {"content": "User mentioned job hunting for months",
             "confidence": 0.88},
1409
             {"content": "User expressed excitement about interview at
1410
             dream company",
1411
              "confidence": 0.92},
1412
             {"content":
1413
             "User indicated spending significant time preparing",
1414
              "confidence": 0.85}
1415
1416
           "semantic_knowledge": [
1417
             {"content": "Interview anxiety is normal and expected",
1418
             "confidence": 0.94},
1419
             {"content": "Over-preparation can sometimes increase anxiety",
              "confidence": 0.82},
1420
             {"content": "Sleep quality affects interview performance",
             "confidence": 0.89}
1422
           ]
         }
1424
1425
```

#### A.13.4 SOCIAL NORMS MODULE PROCESSING

 The Social Norms Module identifies the appropriate context and behavioral expectations:

```
1437
1438
1439
           "social_context": "professional_development",
1440
           "appropriate_norms": [
1441
             "acknowledge interview anxiety as normal",
1442
             "validate preparation efforts",
1443
             "provide practical confidence-building advice",
1444
             "offer both emotional and practical support",
1445
             "maintain appropriate emotional engagement"
1446
           "formality_level": 0.65,
1447
           "appropriate_tone": "supportive_professional"
1448
1449
1450
```

#### A.13.5 GOAL TRACKING MODULE PROCESSING

The Goal Tracking Module identifies and relates current goals to existing goals:

```
1458
1459
1460
1461
           "identified_goals": [
1462
             {"content": "Successfully complete job interview",
             "explicitness": 0.92},
1463
             {"content":
1464
             "Reduce pre-interview anxiety", "explicitness": 0.85}
1465
           ],
1466
           "existing_goals_relevance": [
1467
             {"goal": "Secure desired employment", "relevance": 0.94},
1468
             {"goal":
1469
             "Develop career in chosen field", "relevance": 0.86}
1470
           ]
1472
```

#### A.14 GLOBAL WORKSPACE INTEGRATION

 The outputs from all five modules are integrated through the global workspace mechanism:

#### A.14.1 CONFLICT DETECTION AND RESOLUTION

The system identifies and resolves potential conflicts between module outputs:

	Emotion	Planning	Memory	Social	Goal
Emotion	0.00	0.35	0.10	0.05	0.10
Planning	0.35	0.00	0.15	0.10	0.05
Memory	0.10	0.15	0.00	0.05	0.10
Social	0.05	0.10	0.05	0.00	0.05
Goal	0.10	0.05	0.10	0.05	0.00

Table 4: Conflict detection matrix showing highest conflict between Emotion and Planning modules

The primary conflict (0.35) exists between the Emotion Module's emphasis on validation and the Planning Module's focus on action steps. This conflict affects the modules' weights:

```
Original weights: [0.31, 0.28, 0.21, 0.12, 0.08] (29)
Conflict-adjusted weights: [0.29, 0.25, 0.22, 0.14, 0.10] (30)
```

#### A.14.2 INTEGRATION RESULT

The Integration LLM synthesizes information from all modules into a coherent state:

```
1512
1513
1514
1515
           "integrated_state": {
1516
             "primary_response_focus": "balanced_emotional_practical",
             "emotional_context": {
1517
               "emotion": "anxiety",
1518
               "secondary": "determination",
1519
               "valence": -0.48,
1520
               "arousal": 0.68
1521
1522
             "response_strategy": {
1523
               "approach":
1524
               "validate_anxiety_then_provide_structured_help",
1525
               "tone": "supportive_professional",
1526
               "formality": "moderate"
             "response_plan": [
1528
               "validate interview anxiety as normal",
1529
               "acknowledge preparation effort positively",
1530
               "provide targeted high-impact preparation steps",
               "include confidence-building perspective",
1532
               "end with encouragement and well-wishes"
1533
             ]
1534
           },
1535
           "conflicts_resolved": [
1536
             {
1537
               "conflict": "Emotion module emphasizes validation vs
1538
               Planning module focuses on action steps",
               "resolution": "Balance emotional validation with
1539
               practical action steps by
1540
                               first acknowledging feelings then providing
1541
                               concrete,
1542
                               manageable next steps",
1543
               "confidence": 0.88
1544
             }
1545
           ]
1546
         }
1547
1548
```

#### A.15 FINAL PROMPT CONSTRUCTION

The integrated state is transformed into a structured system prompt for the response generation LLM:

```
1566
1567
1568
        You are a cognitive agent with the following personality:
1569
        MBTI Type: ENFJ
1570
        Big 5 Traits:
        - Openness: high (0.70)
1571
        - Conscientiousness: high (0.80)
1572
        - Extraversion: moderate (0.50)
1573
        - Agreeableness: very high (0.90)
1574
        - Neuroticism: moderate-low (0.40)
1575
1576
         [EMOTION MODULE] (Contribution: 0.29): Pay careful attention to
1577
        the anxiety and nervousness the user is experiencing about
1578
        their upcoming job interview. The emotions detected include
1579
        high anxiety (0.82) with a negative valence (-0.48) but also
1580
        determination (0.55). Emotional content identified: "really
1581
        nervous", "don't feel ready". Use a validation_with_reframing
        approach that acknowledges the anxiety while gently building
1582
        confidence.
1583
1584
        [PLANNING MODULE] (Contribution: 0.25):
1585
        Provide structured, preparation_strategy advice for the user's
1586
        interview preparation. Focus on high-impact preparation
1587
        activities that can be completed in limited time. Consider
1588
        including: reviewing core materials, practicing answers aloud,
1589
        preparing achievement examples, reducing anxiety, handling
1590
        logistics, and getting rest.
1591
1592
        [MEMORY MODULE] (Contribution: 0.22):
        Reference these key points from conversation history:
1593
        User has been job hunting for months, this is an interview at
1594
        their "dream company," and they have already spent significant
1595
        time preparing. Also incorporate knowledge that interview
1596
        anxiety is normaland that over-preparation can sometimes
1597
        increase anxiety rather than helping.
1598
1599
        [SOCIAL NORMS MODULE] (Contribution: 0.14):
        Maintain a supportive_professional
1601
        tone appropriate for career_development context.
1602
        Acknowledge that interview
1603
        anxiety is normal while validating preparation efforts.
1604
        Observe these norms:
        validate emotions, show empathy, offer practical support,
1605
        maintain appropriate emotional engagement.
1606
        Follow this response strategy:
1608
        validate_anxiety_then_provide_structured_help with
1609
        a supportive_professional tone at moderate formality level.
1610
1611
        When addressing conflicts:
1612
        Balance emotional validation with practical action steps
1613
        by first acknowledging feelings then providing concrete,
1614
        manageable next steps.
1615
        Process the query by considering each
1616
        cognitive dimension above according to
1617
        its contribution weight. Follow this response plan:
1618
        1. Validate interview anxiety as normal
1619
        2. Acknowledge preparation effort positively
        3. Provide targeted high-impact preparation steps
        4. Include confidence-building perspective
        5. End with encouragement and Well-wishes
```

#### A.16 RESPONSE GENERATION

1622

The Response LLM generates the final output based on the system prompt:

1623 1624

1625 1626 1627

1628

It's completely normal to feel nervous before an interview, especially for a company you're really excited about! The fact that you've been preparing for weeks shows your dedication, and that preparation won't disappear overnight, even if you're feeling anxious right now.

- 1629 1630

- 1633 1634
- 1635
- 1637
- 1639 1640

1641

1642

1643

1644 1645

1646 1647 1648

1650

1652

1656 1657

1659

1662 1663

1658

1665 1666

1664

1669 1671

1672 1673 At this point, here's what I'd suggest focusing on: 1. Review your strongest points - skim through the key qualifications and 2-3 accomplishment

- stories you want to highlight. Don't try to memorize everything.
- 2. Practice answering 3-5 common questions out loud this builds muscle memory that helps even when nerves kick in.
- 3. Prepare your logistics for tomorrow lay out your outfit, plan your route, and gather any materials you need to bring.
- 4. Most importantly, make sure you get adequate rest tonight. A clear mind will serve you better than a few more hours of anxious preparation.
- Remember that interviews are a two-way conversation, not an examination. You've been preparing for weeks - now it's time to trust that preparation and focus on presenting your authentic self.

Good luck tomorrow! You've put in the work, and that's the most important part.

# A.17 Information Flow Analysis

Examining the example reveals several key aspects of the GNWT-Agent architecture:

- 1. **Modular Specialization**: Each module contributes distinct information:
  - Emotion Module provides emotional understanding and regulation strategy
  - Planning Module offers concrete action steps
  - Memory Module contributes contextual information and relevant knowledge
  - Social Norms Module guides appropriate tone and interaction norms
  - Goal Tracking Module ensures alignment with long-term objectives
- 2. Salience-Based Processing: The architecture dynamically allocates attention based on query characteristics, with emotional content receiving the highest weight in this anxiety-focused query.
- 3. Conflict Resolution: The system explicitly identifies and resolves the tension between emotional validation and practical advice through a balanced approach.
- 4. Structured Response Planning: The final response follows the five-step plan specified in the integration phase, creating a cohesive structure that addresses multiple dimensions of the query.
- 5. **Personality Influence**: The system's responses reflect the specified personality characteristics, particularly high agreeableness (0.90) through the empathetic tone.

The response demonstrates how the GNWT-Agent architecture produces outputs that balance emotional responsiveness with practical utility, organized through a structured cognitive framework that mimics aspects of human cognition.

A.18 USE CASE: DATING APPLICATION

#### A.19 OVERVIEW

CogniPair for Dating represents a novel approach to matchmaking that leverages the GNWT-Agent cognitive architecture to create realistic digital twins of users. Unlike traditional dating platforms that rely on static profiles and rule-based matching algorithms, CogniPair simulates genuine cognitive interactions between potential matches before they ever meet in person. This system models personality traits, emotional responses, memory formation, social norm adherence, and planning capabilities to predict compatibility with high accuracy.

#### A.20 USER JOURNEY

The CogniPair Dating platform guides users through the following process:

- 1. **Profile Creation:** Users complete an in-depth personality assessment and preference questionnaire that captures traits across multiple dimensions including the Big Five personality factors.
- 2. **Digital Twin Generation:** The GNWT-Agent architecture constructs a cognitive digital twin that faithfully reproduces the user's conversational patterns, emotional responses, and decision-making processes.
- Simulated Interactions: Digital twins engage in mock conversations with potential matches, processing information through specialized cognitive modules.
- Compatibility Analysis: The system analyzes cognitive states throughout conversations, identifying mutual interests, conversational flow, emotional resonance, and decision-making alignment.
- Match Recommendation: Users receive matches based on authentic cognitive compatibility rather than superficial profile similarities.

#### A.21 CASE STUDY: EMMA AND MICHAEL

To illustrate CogniPair's capabilities, we present a case study of two users, Emma and Michael, whose digital twins engaged in a simulated conversation.

#### A.21.1 USER PROFILES

Attribute	Emma	Michael
Age	28	31
Background	Marketing director, enjoys	Software engineer, passionate
	outdoor activities, travel	about music, cooking, technology
Personality	Openness: 0.75	Openness: 0.68
Traits	Conscientiousness: 0.62	Conscientiousness: 0.82
	Extraversion: 0.70	Extraversion: 0.45
	Agreeableness: 0.82	Agreeableness: 0.74
	Neuroticism: 0.35	Neuroticism: 0.42
Dating	Attractiveness: 8.2/10	Attractiveness: 7.8/10
Attributes	Sincerity: 8.7/10	Sincerity: 9.1/10
	Intelligence: 8.5/10	Intelligence: 9.0/10
	Fun: 8.9/10	Fun: 7.5/10
	Ambition: 8.8/10	Ambition: 8.7/10
	Shared Interests: 7.9/10	Shared Interests: 8.2/10

Table 5: Profiles of Emma and Michael with personality traits and dating attributes

#### A.21.2 SIMULATED CONVERSATION WITH COGNITIVE PROCESSING

The following transcript shows the simulated conversation between Emma and Michael's digital twins, with real-time cognitive processing visualized.

System Introduction: Welcome to a speed date! This is a first conversation between two people who are considering dating. Please introduce yourselves.

**Emma:** Hi there! I'm Emma. I work in marketing and I'm really passionate about creating campaigns that connect with people. In my free time, I love hiking and being outdoors whenever possible. I also try to travel to at least one new country each year. What about you?

#### Emma's Cognitive Processing

#### **Module Activations:**

 $\bullet$   $\,$  Emotion Module: 0.57 - Moderate activation showing enthusiasm and openness

• **Memory Module: 0.35** - Low activation retrieving personal information

• Planning Module: 0.31 - Low activation structuring introduction

• Goal Tracking Module: 0.42 - Moderate activation focusing on relationship exploration

• Social Norms Module: 0.64 - High activation ensuring appropriate self-disclosure

**Michael:** Nice to meet you, Emma! I'm Michael. I work as a software engineer, focusing on user experience design. I find it fascinating to solve problems that make technology more intuitive. Outside of work, I'm really into music—I play guitar and go to concerts whenever I can. I also love experimenting with cooking, especially dishes from places I've never been to. Your travel experiences sound amazing. What's been your favorite country to visit so far?

#### Michael's Cognitive Processing

#### **Module Activations:**

• Emotion Module: 0.45 - Moderate activation showing interest

 Social Norms Module: 0.55 - Moderate activation ensuring appropriate response
 Memory Module: 0.38 - Low-moderate activation retrieving personal details

• Planning Module: 0.51 - Moderate activation structuring response and question

 • Goal Tracking Module: 0.49 - Moderate activation exploring compatibility

**Emma:** It's nice to meet someone who appreciates music! I'm not musically talented myself, but I love going to live shows. As for travel, that's a tough question! I'd say Japan was probably the most fascinating—the blend of traditional culture with cutting-edge technology was incredible. Plus, the food was amazing. Do you enjoy traveling yourself, or is cooking your way of experiencing different cultures?

# 1784 1785

#### **Module Activations:**

- Emotion Module: 0.68 High activation showing increased enthusiasm
- 1786 1787
- Memory Module: 0.72 High activation recalling travel experiences
- 1788 1789
- Social Norms Module: 0.53 Moderate activation maintaining conversation flow
- 1790
- Planning Module: 0.29 Low activation
- 1791 1792 1793
- Goal Tracking Module: 0.59 Moderate-high activation exploring common inter-

1794 1795 1796

Cognitive State Analysis: Emma's emotion module showed significant activation when suggesting authentic experiences rather than generic responses.

1797

1799

1801

1803

1805 1806

1807 1808

1809 1810

1811 1812

1813 1814

1815 1816 1817

1818 1819 1820

1821 1823

1824 1825 1826

1827

1828

1830

1831 1832

1833

1834 1835

discussing travel and responding to Michael's interests in music, indicating genuine enthusiasm. The memory module activated strongly when recalling specific details about Japan,

**Michael:** That's exactly what fascinates me about Japan too! The contrast between ancient temples and robot restaurants in Tokyo is mind-blowing. I haven't been there yet, but it's at the top of my list. I do enjoy traveling when I can, though not as regularly as you. Cooking is definitely my way of exploring cultures when I can't travel there—I recently got into making ramen from scratch after watching a documentary about Japanese food. What kind of music do you enjoy at live shows? Any favorite concerts?

#### **Module Activations:**

- Emotion Module: 0.79 High activation showing excitement about shared interest
- Memory Module: 0.68 High activation connecting to personal experiences
- Social Norms Module: 0.51 Moderate activation
- Planning Module: 0.35 Low-moderate activation
- Goal Tracking Module: 0.73 High activation pursuing connection point

**Cognitive State Analysis:** Michael's emotion module showed strong activation when responding to the shared interest in Japan, indicating genuine excitement. The goal tracking module also activated significantly, suggesting he sees this common interest as a potential relationship foundation. His follow-up question about music demonstrates the goal-directed behavior of exploring additional shared interests.

# [Conversation continues...]

# A.21.3 COMPATIBILITY DECISION PHASE

After the simulated conversation, each digital twin makes an independent compatibility assessment based on their cognitive processing throughout the interaction.

**System:** Based on your conversation, would you like to meet this person for a real date?

**Emma's Decision:** Yes, I would like to meet Michael. Our conversation felt natural and engaging. We share interests in travel and cultural experiences, and I appreciate his passion for music and cooking. His thoughtful responses and questions showed genuine interest. There's enough similarity for connection but enough difference to keep things interesting.

#### Emma's Decision Cognitive Processing

#### **Module Contributions to Decision:**

- Emotion Module: 0.75 Strong positive emotional response throughout conversation
- Memory Module: 0.68 High recall of conversation details and identified commonalities
- Social Norms Module: 0.42 Moderate assessment of social compatibility
- Planning Module: 0.51 Moderate assessment of future potential
- Goal Tracking Module: 0.71 High alignment with relationship goals

**Michael's Decision:** Yes, I'm interested in meeting Emma for a real date. We have complementary interests and her enthusiasm for travel and cultural experiences resonates with me. I enjoyed how she engaged with my interests in music and cooking while sharing her own perspectives. Our conversation had a good balance of questions and sharing, and I'd like to learn more about her.

#### Michael's Decision Cognitive Processing

#### **Module Contributions to Decision:**

- Emotion Module: 0.69 Strong positive emotional response
- Memory Module: 0.64 Solid recall of shared interests and conversation points
- Social Norms Module: 0.48 Moderate assessment of social compatibility
- Planning Module: 0.66 High consideration of relationship potential
- Goal Tracking Module: 0.72 High alignment with relationship goals

#### Match Result: COMPATIBLE

#### A.22 COGNITIVE MODULE ACTIVITY ANALYSIS

The module activation patterns throughout the dating conversation revealed several key insights:

Module	Introduction	Response	Follow-up	Shared Interest	Decision
Emotion	0.57	0.45	0.68	0.79	0.72
Memory	0.35	0.38	0.72	0.68	0.66
Social Norms	0.64	0.55	0.53	0.51	0.45
Planning	0.31	0.51	0.29	0.35	0.58
Goal Tracking	0.42	0.49	0.59	0.73	0.71

Table 6: Cognitive module activation patterns throughout the dating conversation

The cognitive architecture revealed several key insights during this match:

- Emotion Module: Showed increasing activation as shared interests were discovered, peaking during discussion of travel experiences.
- Memory Module: Activation increased significantly when recalling specific experiences, demonstrating authentic engagement rather than superficial responses.
- Social Norms Module: Started high during initial introductions and gradually decreased as conversation became more comfortable.
- **Planning Module:** Peaked during the decision phase when considering future interaction potential.
- Goal Tracking Module: Showed steady increase throughout the conversation as relationship compatibility was assessed.

**OUTCOMES AND USER INTERFACE** 

Match Result: Emma and Michael - 87% Compatibility

Emotional Connection: 83%

• Conversation Flow: 91%

• Shared Interests: 78%

• Value Alignment: 85%

connection and potential conversation starters based on shared interests.

1890

1891 1892

1894 1895

1896

1897 1898

1899

1900 1901

1902

1903

1904

1943

A.23

CogniPair Match Report

**Compatibility Breakdown:** 

#### • Long-term Potential: 84% 1905 **Connection Points:** 1907 1908 • Travel experiences (particularly interest in Japan) 1909 · Appreciation for cultural exploration 1910 • Complementary interests (Emma's outdoor activities, Michael's cooking) 1911 1912 Similar communication styles with thoughtful questions 1913 **Suggested Conversation Starters:** 1914 1915 "I'd love to hear more about that documentary on Japanese food that inspired your cooking." 1916 • "What's been your favorite live music experience? I'm always looking for new artists." 1917 1918 "Would you want to do cooking and hiking as combined activities? Maybe prepare a meal after a trail?" 1919 1920 1921 TECHNICAL IMPLEMENTATION HIGHLIGHTS 1922 The dating scenario leverages several key aspects of the GNWT-Agent architecture: 1923 1924 1. Emotion Module Prominence: Dating interactions show heightened emotion module 1925 activation compared to other scenarios, particularly in response to shared interests and 1926 values. 1927 2. **Memory-Emotion Integration:** The architecture demonstrates how memories trigger 1928 emotional responses in social contexts, creating authentic patterns of engagement. 1929 3. Goal-Directed Decision Making: As the conversation progresses, goal tracking module 1930 activation increases, culminating in the compatibility decision. 1931 1932 4. **Personal Value Assessment:** The architecture evaluates alignment in preferences, interests, 1933 and communication styles to determine overall compatibility. 1934 This implementation demonstrates how cognitive modeling can transcend the limitations of traditional 1935 profile-based matching by simulating the nuanced psychological aspects of human connection. 1936 1938 1939 1941 1942

After the cognitive simulation concludes, CogniPair presents users with compatible matches along with insights derived from the digital twin interactions. The platform highlights specific points of

1944 A.25 USE CASE: JOB APPLICATION

#### A.26 OVERVIEW

CogniPair for Job Hiring represents a transformative approach to talent acquisition that applies the GNWT-Agent cognitive architecture to create digital twins of both hiring managers and job candidates. This system moves beyond traditional resume screening and interview processes by modeling the complex cognitive dynamics that determine professional compatibility. By simulating realistic interview interactions, the platform assesses technical skills, problem-solving approaches, communication styles, and cultural alignment before any in-person interviews take place.

#### A.27 USER JOURNEY

The CogniPair Hiring platform guides users through the following process:

- Employer Profile Creation: Companies input detailed job descriptions, required competencies, team dynamics, and organizational culture attributes.
- Candidate Profile Submission: Candidates upload resumes and complete competency assessments that capture both technical skills and cognitive-behavioral traits.
- Digital Twin Generation: The GNWT-Agent architecture creates cognitive models of both hiring managers and candidates, incorporating decision-making patterns, communication styles, and problem-solving approaches.
- 4. **Simulated Interviews:** Digital twins engage in multi-stage interview simulations assessing technical skills, problem-solving abilities, and cultural alignment.
- Compatibility Analysis: The system analyzes cognitive states throughout the simulated interviews, identifying communication effectiveness, problem-solving compatibility, and potential team integration.
- 6. **Match Recommendation:** Both employers and candidates receive compatibility assessments that predict professional success and team dynamics.

#### A.28 CASE STUDY: TECHSOLUTIONS INC. AND CANDIDATE ALEX

To illustrate CogniPair's capabilities in hiring contexts, we present a case study of a software development position at TechSolutions Inc. and a potential candidate, Alex.

#### A.28.1 JOB AND CANDIDATE PROFILES

#### A.28.2 SIMULATED INTERVIEW WITH COGNITIVE PROCESSING

The following transcript shows the simulated interview between Sarah (Engineering Manager) and Alex (Candidate), with real-time cognitive processing visualized.

**System Introduction:** This is a technical interview for a Senior Software Developer position at TechSolutions Inc.

**Sarah:** Hello Alex, thanks for joining us today. I'm Sarah, the Engineering Manager at TechSolutions. Before we dive into technical questions, could you tell me about your experience with React and how you've implemented it in previous projects?

1	IJ	IJ	9
2	0	0	0
2	0	0	1
2	0	0	2
2	0	0	3
2	0	0	4
2	0	0	5
2	0	0	6
2			7
2			8
2			9
2			0
2			1
2	_		
		1	
2	_	1	
2		1	
	0		5
	0		6
	0		7
	0		8
			9
			0
	0		
2	0	2	2
2	0	2	3
2	0	2	4
2	0	2	5
2	0	2	6
	0		
2	0	2	8
			9
			0
	0		
	0		
	0		
	0		
	0		
	0		
	0		
	0		
			9
	_		0
	0		
2	0	4	2
2	0	4	3
2	0	4	4
2	0	4	5
2	0	4	6
2	0	4	7
2	0	4	8
2	0	4	9
			0
	0		
_	_	_	-

Attribute	TechSolutions Inc.	Alex (Candidate)
Position	Senior Software Developer	5 years software development
	Team of 8, collaborative environment	experience, previous team lead
Requirements/	JavaScript/React expertise	Strong JavaScript/React skills
Skills	CI/CD pipeline experience	CI/CD experience with GitLab
	Agile development methodology	Certified Scrum Master
	Problem-solving abilities	Bachelor's in Computer Science
Interviewer	Sarah (Engineering Manager)	
Personality	Openness: 0.72	Openness: 0.85
Traits	Conscientiousness: 0.86	Conscientiousness: 0.69
	Extraversion: 0.58	Extraversion: 0.47
	Agreeableness: 0.65	Agreeableness: 0.78
	Neuroticism: 0.32	Neuroticism: 0.39
Professional	Technical skills: 9.2/10	Technical skills: 8.8/10
Attributes	Communication: 8.5/10	Communication: 8.7/10
	Leadership potential: 9.0/10	Leadership potential: 8.4/10
	Problem-solving: 9.3/10	Problem-solving: 9.1/10
	Cultural alignment: 8.7/10	Cultural alignment: ?/10 (to be assessed)
Company	Innovative, fast-paced	Previous experience in fintech
Culture/	Work-life balance emphasis	Values collaborative environments
Values	Continuous learning encouraged	Mentor to junior developers

Table 7: Job description and candidate profile with professional attributes

#### Sarah's Cognitive Processing

#### **Module Activations:**

- Planning Module: 0.78 High activation structuring interview strategy
- Social Norms Module: 0.63 Moderate-high activation ensuring professional tone
- Memory Module: 0.55 Moderate activation recalling job requirements
- Goal Tracking Module: 0.69 High activation focusing on technical assessment
- Emotion Module: 0.32 Low activation maintaining neutral assessment stance

**Alex:** Thank you for having me, Sarah. I've been working with React for about four years now. In my current role at FinTech Solutions, I led the migration of our legacy dashboard to a React-based architecture. One of the challenges we faced was maintaining state across a complex application with multiple user roles. I implemented a Redux state management solution with middleware for API interactions that reduced loading times by 40% and significantly improved the developer experience for my team. I also created a reusable component library that we eventually open-sourced within the company.

#### Alex's Cognitive Processing

#### **Module Activations:**

- **Memory Module: 0.82** Very high activation retrieving specific professional experiences
- Planning Module: 0.74 High activation structuring response with relevant details
- Social Norms Module: 0.61 Moderate-high activation ensuring professional presentation
- Goal Tracking Module: 0.69 High activation focusing on demonstrating expertise
- Emotion Module: 0.38 Low-moderate activation showing appropriate enthusiasm

**Sarah:** That's impressive work with the state management solution. I'm curious about how you approached testing for that React application, especially with the Redux integration. Could you walk me through your testing strategy and any particular challenges you encountered?

#### Sarah's Cognitive Processing

#### **Module Activations:**

- Planning Module: 0.83 Very high activation probing technical depth
- Memory Module: 0.76 High activation processing previous answer against requirements
- Goal Tracking Module: 0.72 High activation assessing technical competency
- Social Norms Module: 0.48 Moderate activation
- Emotion Module: 0.36 Low-moderate activation showing positive response to quality answer

**Cognitive State Analysis:** Sarah's planning module shows high activation as she strategically probes deeper technical knowledge. Her memory module has significantly increased as she processes Alex's initial response against the role requirements, showing active assessment rather than just following a script.

**Alex:** For testing, we implemented a multi-layered approach. At the unit level, we used Jest with React Testing Library to test individual components, focusing on behavior rather than implementation details. For Redux, we wrote separate tests for reducers, action creators, and middleware—treating them as pure functions made this fairly straightforward.

The real challenge came with integration testing. We needed to test user flows that spanned multiple components and redux state changes. We used Cypress for end-to-end testing of critical paths, but found we needed something in between unit and E2E tests. So we developed a custom testing approach using React Testing Library with a real Redux store to test larger component trees together.

One specific challenge was testing asynchronous workflows, especially error handling across the Redux middleware. We ended up creating a mock API service layer that could simulate various network conditions and response scenarios, which significantly improved our test coverage for edge cases.

#### Alex's Cognitive Processing

#### **Module Activations:**

- Memory Module: 0.85 Very high activation recalling specific technical implementations
- Planning Module: 0.78 High activation structuring a comprehensive response
- Goal Tracking Module: 0.74 High activation demonstrating problem-solving approach
- Social Norms Module: 0.52 Moderate activation maintaining professional communication
- Emotion Module: 0.35 Low-moderate activation

**Cognitive State Analysis:** Alex's memory module shows extremely high activation, demonstrating deep technical knowledge rather than surface-level responses. The goal tracking module indicates awareness of demonstrating both technical knowledge and problem-solving approaches, aligning with the interviewer's objectives.

[Interview continues with technical questions, problem-solving scenarios, and culture fit assessment...]

A.28.3 PROBLEM-SOLVING ASSESSMENT

A critical component of technical interviews is evaluating problem-solving approaches. The following depicts a problem-solving scenario and the cognitive processing involved.

Sarah: Let's walk through a technical scenario. Imagine you're working on our platform and users are reporting intermittent performance issues—the dashboard occasionally takes 5+ seconds to load, but works fine most of the time. How would you approach troubleshooting and resolving this issue?

#### **Module Activations:**

Planning Module: 0.89 - Very high activation presenting ambiguous problem

 • Goal Tracking Module: 0.85 - Very high activation assessing problem-solving approach

• **Memory Module: 0.62** - Moderate-high activation recalling real-world scenarios

Social Norms Module: 0.41 - Moderate activation

• Emotion Module: 0.29 - Low activation maintaining neutral stance

**Alex:** I'd approach this methodically since intermittent issues are often the hardest to solve. First, I'd gather data to understand the pattern:

1. Implement detailed performance monitoring if not already in place—tracking load times, component render times, and backend response times.

2. Look for correlating factors—does it happen at specific times of day, for specific user roles, or with particular data sets?

3. Check browser dev tools, especially the network and performance tabs, to identify bottlenecks—is it a frontend rendering issue, a slow API call, or resource loading?

Since it's intermittent, I'd suspect a few possible causes: - Cache invalidation issues - Memory leaks from component remounting - Database query performance varying with data volume or concurrent users - Third-party service integration delays

I'd isolate components by creating a testing environment that can simulate various loads. If it's a backend issue, I'd add query performance monitoring and examine execution plans for database queries.

For a quick mitigation while investigating, I might implement better loading states and potentially add a strategic caching layer if appropriate.

The key is collecting enough data to reproduce the issue consistently before attempting to solve it.

2162 2163

2164

2165 2166

2171 2172 2173

2175 2176 2177

2174

2178 2179

2180 2181 2182

2183

2184 2185 2186

2187

2188

2189 2190 2191

2192

2193 2194 2195

2197 2198 2199

2200

2201 2202

2196

2203 2204 2205

2206

2208 2209

2210

2211 2212 2213

### **Module Activations:**

#### • Planning Module: 0.92 - Very high activation structuring systematic approach

- Memory Module: 0.81 Very high activation retrieving relevant technical knowledge
- Goal Tracking Module: 0.78 High activation demonstrating problem-solving methodology
- Social Norms Module: 0.45 Moderate activation maintaining professional communication
- Emotion Module: 0.32 Low activation focusing on technical content

Cognitive State Analysis: Alex's planning module shows exceptional activation, demonstrating a structured, methodical approach to an ambiguous problem. The high memory module activation shows application of past experience to new scenarios, indicating transferable problem-solving skills rather than memorized solutions.

#### A.28.4 HIRING DECISION PHASE

After the simulated interview, each digital twin makes an independent assessment based on their cognitive processing throughout the interaction.

**System:** Based on the interview, would you recommend hiring this candidate?

Sarah's Decision: Yes, I recommend hiring Alex for the Senior Software Developer position. The candidate demonstrated strong technical expertise in React and modern frontend architecture, with particular strength in state management and testing approaches. Their problem-solving methodology is systematic and thorough, with a good balance of immediate tactical steps and strategic thinking. Alex's experience leading a component library initiative aligns well with our needs for improving code reusability. Their communication style is clear and structured, which would fit well with our team dynamics. The only area for development might be strengthening CI/CD experience, but their foundation is solid.

#### Sarah's Decision Cognitive Processing

#### **Module Contributions to Decision:**

- Planning Module: 0.83 Very high assessment of candidate's problem-solving approach
- **Memory Module: 0.78** High comparison against job requirements
- Goal Tracking Module: 0.75 High evaluation of candidate fit with team objectives
- Social Norms Module: 0.62 Moderate-high assessment of professional demeanor
- Emotion Module: 0.42 Moderate positive impression of candidate

**Alex's Decision:** Based on the interview, I am interested in proceeding with the opportunity at TechSolutions Inc. The technical challenges discussed align well with my experience and interests. I appreciated Sarah's detailed questions about testing methodologies, which suggests the company values quality code. The team structure and emphasis on component reusability matches my collaborative working style. The problem-solving scenario about intermittent performance issues indicates realistic challenges that I find engaging. Overall, the position seems to offer good growth opportunities while leveraging my existing strengths in React and state management.

2	2	1	4
2	2	1	5
2	2	1	6
2	2	1	7
2	2	1	8
2	2	1	9
	2		
	2		
	2		
	2		
	2		
	2		
	2		
	2		
	2		
	2		
	2		
2	2	3	1
2	2	3	2
2	2	3	3
2	2	3	4
2	2	3	5
2	2	3	6
	2		
2	2	3	8
	2		
	2	_	_
		4	u
2	2	4	1
2	2	4	1
 2	2 2	4 4 4	1 2 3
2 2 2	2 2 2	4 4 4	1 2 3 4
2 2 2	2 2 2 2	4 4 4 4	1 2 3 4 5
2 2 2 2	2 2 2 2	4 4 4 4	1 2 3 4 5 6
2 2 2 2	2 2 2 2 2	4 4 4 4 4	1 2 3 4 5 6 7
2 2 2 2 2 2	2 2 2 2 2 2 2	4 4 4 4 4 4	1 2 3 4 5 6 7 8
2 2 2 2 2 2	2 2 2 2 2 2	4 4 4 4 4 4	1 2 3 4 5 6 7 8
2 2 2 2 2 2 2 2 2	2 2 2 2 2 2 2	4 4 4 4 4 5	1 2 3 4 5 6 7 8 9
2 2 2 2 2 2 2 2 2	2 2 2 2 2 2 2	4 4 4 4 4 5	1 2 3 4 5 6 7 8 9
2 2 2 2 2 2 2 2 2 2 2 2	2 2 2 2 2 2 2	44444455	1 2 3 4 5 6 7 8 9 0 1
2 2 2 2 2 2 2 2	2 2 2 2 2 2 2	44444555	1 2 3 4 5 6 7 8 9 0 1 2
2 2 2 2 2 2 2 2 2 2	2 2 2 2 2 2 2 2	4444445555	1 2 3 4 5 6 7 8 9 0 1 2 3
2 2 2 2 2 2 2 2 2	2 2 2 2 2 2 2 2 2	4 4 4 4 4 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5	12345678901234
2 2 2 2 2 2 2 2 2 2 2	2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2	4 4 4 4 4 4 5 5 5 5 5 5	1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 5
2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2	2 2 2 2 2 2 2 2 2 2 2 2 2	4 4 4 4 4 4 5 5 5 5 5 5 5	1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6
2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2	2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2	4 4 4 4 4 4 5 5 5 5 5 5 5 5	1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7
2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2	2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2	4 4 4 4 4 4 5 5 5 5 5 5 5 5 5	1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8
2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2	2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2	4 4 4 4 4 4 5 5 5 5 5 5 5 5 5	1234567890123456789
	2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2	4 4 4 4 4 4 4 5 5 5 5 5 5 5 5 6	12345678901234567890
	2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2	4 4 4 4 4 4 4 5 5 5 5 5 5 5 5 6 6	123456789012345678901
	2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2	4 4 4 4 4 4 4 5 5 5 5 5 5 5 5 6 6 6	1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2
	2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2	4 4 4 4 4 4 4 5 5 5 5 5 5 5 5 6 6 6 6	12345678901234567890123
	2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2	4 4 4 4 4 4 4 5 5 5 5 5 5 5 5 5 6 6 6 6	1 2 3 4 4 5 6 7 8 9 0 1 2 3 4 4 5 6 7 8 9 0 1 2 3 4 4 5 6 7 8 9 0 1 2 3 4 4 6 7 8 9 0 1 2 3 4 6 7 8 9 0 1 2 2 3 4 6 7 8 9 0 1 2 2 3 4 6 7 8 9 0 1 2 2 3 4 6 7 8 9 0 1 2 2 3 4 6 7 8 9 0 1 2 2 2

#### Alex's Decision Cognitive Processing

#### **Module Contributions to Decision:**

- Goal Tracking Module: 0.81 Very high alignment with career objectives
- Memory Module: 0.73 High comparison against previous job experiences
- Planning Module: 0.69 High assessment of growth opportunities
- Social Norms Module: 0.58 Moderate assessment of team culture fit
- Emotion Module: 0.54 Moderate positive impression of company

#### Match Result: COMPATIBLE - RECOMMENDED HIRE

#### A.29 COGNITIVE MODULE ACTIVITY ANALYSIS

The module activation patterns throughout the hiring interview revealed several key insights:

Module	Introduction	Technical	Problem-Solving	Culture Fit	Decision
Emotion	0.35	0.32	0.29	0.45	0.48
Memory	0.55	0.79	0.81	0.63	0.75
Social Norms	0.63	0.48	0.43	0.68	0.60
Planning	0.78	0.83	0.91	0.67	0.76
Goal Tracking	0.69	0.73	0.82	0.71	0.78

Table 8: Cognitive module activation patterns throughout the hiring interview

The cognitive architecture revealed several key insights during this hiring simulation:

- **Planning Module:** Dominated the cognitive processing during technical and problem-solving phases, demonstrating the critical importance of structured thinking in hiring contexts.
- **Memory Module:** Showed substantial activation during technical discussions as both interviewer and candidate accessed domain knowledge and past experiences.
- **Emotion Module:** Consistently lower than in dating contexts, but increased during cultural fit discussions and final decision making.
- Social Norms Module: Peaked during introduction and cultural fit assessment phases, indicating heightened attention to professional communication standards.
- Goal Tracking Module: Maintained high activation throughout, focusing on alignment between candidate capabilities and job requirements.

#### A.30 OUTCOMES AND USER INTERFACE

After the cognitive simulation concludes, CogniPair presents hiring teams with a comprehensive assessment of candidate fit along with specific insights derived from the digital twin interactions.

#### **CogniPair Hiring Assessment Report**

Match Result: Alex for Senior Software Developer - 89% Compatibility

#### **Technical Skills Assessment:**

- React/Frontend Development: 92% Exceptional
- Testing Methodology: 88% Strong
- State Management: 94% Exceptional
- CI/CD Experience: 76% Adequate, potential growth area
- Problem-Solving Approach: 91% Exceptional

Team and Cultural Fit:

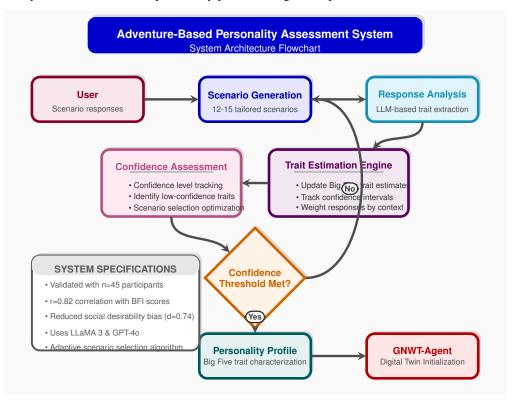
2268

2269 2270

#### • Communication Clarity: 87% - Strong 2271 • Collaboration Potential: 85% - Strong 2272 • Technical Leadership: 83% - Strong 2273 • Learning Orientation: 90% - Exceptional 2274 2275 • Value Alignment: 81% - Strong 2276 **Key Strengths:** 2277 2278 Systematic approach to problem-solving with strong emphasis on data collection 2279 • Experience creating reusable component libraries aligns with current initiatives 2280 2281 • Strong testing methodology with creative solutions for integration testing 2282 Clear, structured communication style compatible with existing team 2283 **Development Areas:** 2285 • Deeper CI/CD pipeline experience would be beneficial 2287 Could strengthen infrastructure monitoring knowledge **Recommended Next Steps:** 2289 2290 Proceed with offer process 2291 • Consider onboarding plan that includes pairing with DevOps specialist 2293 Explore potential for leadership in component library initiative 2294 2295 A.31 TECHNICAL IMPLEMENTATION HIGHLIGHTS 2296 The hiring scenario leverages several key aspects of the GNWT-Agent architecture with notably 2297 different patterns than the dating scenario: 2298 2299 1. **Planning Module Dominance:** Hiring interactions show significantly higher planning mod-2300 ule activation compared to dating contexts, particularly during problem-solving scenarios. 2301 2. **Memory-Planning Integration:** The architecture demonstrates how technical knowledge 2302 (memory) feeds into structured problem-solving approaches (planning), creating a distinct cognitive fingerprint for technical roles. 3. **Reduced Emotion Module Activation:** Professional contexts show consistently lower 2305 emotion module activation, with cognitive resources redirected to analytical processes. 2306 2307 4. Goal Alignment Assessment: The architecture evaluates alignment between candidate 2308 capabilities and job requirements with greater precision than traditional interviewing techniques. 2309 2310 This implementation demonstrates how the same core GNWT-Agent architecture can be effectively 2311 repurposed for professional contexts by adjusting module weights and domain-specific attributes, 2312 providing a versatile framework for human-centric matching across diverse applications. 2313 2314 2315 2316 2317 2318 2319 2320

#### A.32 ADVENTURE-BASED PERSONALITY ASSESSMENT EXAMPLE

This appendix provides a comprehensive example of our adventure-based personality assessment system. As detailed in the main text, we developed this interactive assessment to identify participants' authentic personality traits while reducing self-presentation biases commonly associated with traditional questionnaires. The system presents participants with 12-15 scenarios designed to elicit responses that map to Big Five personality dimensions, with each choice analyzed to update trait estimates based on established correlations between situational responses and personality dimensions. The assessment demonstrated strong validity in our validation study (n=45), with trait scores showing high correlation (r=0.82) with established Big Five Inventory (BFI) scores, and significantly reduced social desirability bias compared to standard questionnaires (Cohen's d=0.74). Below is a comprehensive example of a complete adventure-based personality assessment, demonstrating how our system builds a detailed personality profile through multiple scenarios.



Scenario 1: The Crossroads

**System**: You find yourself at a crossroads at the edge of a vast, unexplored land. In the distance, you can see mountains to the north, a forest to the east, a bustling town to the south, and a quiet beach to the west. Where do you choose to go first? **Options**:

- 1. Head north toward the mountains
- 2. Go east into the forest
- 3. Travel south to the town
- 4. Walk west to the beach

**User selects**: 1 (Mountains)

**LLM Analysis** (internal, using llama3:latest or GPT-40):

```
Choice analysis:

    Selecting mountains suggests comfort with challenge and self-

   reliance
- Indicates moderate to high openness to experience (seeking
   adventure)
  Suggests conscientiousness (willing to tackle difficult terrain)
- May indicate lower extraversion (preference for solitude over the
    town)
- Neutral on agreeableness
- Indicates low neuroticism (comfortable with uncertainty)
Trait updates:
openness: 65 (confidence: 0.3)
conscientiousness: 60 (confidence: 0.3)
extraversion: 45 (confidence: 0.2)
agreeableness: 50 (confidence: 0.1)
neuroticism: 40 (confidence: 0.2)
```

**System**: As you begin your journey toward the mountains, you notice the path looks challenging. What are you most looking forward to about this mountain trek?

**User responds**: I'm excited about the view from the top and the sense of accomplishment. I like pushing myself physically and seeing what I'm capable of. Also, I'm hoping to find some quiet time to think and reflect away from everyday distractions.

**LLM Analysis** (internal, using deepseek-r1:latest or GPT-40):

```
Text analysis:

- "Sense of accomplishment" and "pushing myself" indicates achievement orientation

- "Quiet time to think" suggests introversion and reflectiveness

- Focus on physical challenge indicates conscientiousness

- No mention of social aspects reinforces lower extraversion

- Seeking accomplishment suggests growth mindset and low neuroticism

Trait updates:
openness: 60 (confidence: 0.2)
conscientiousness: 70 (confidence: 0.4)
extraversion: 35 (confidence: 0.3)
agreeableness: 50 (confidence: 0.1)
neuroticism: 35 (confidence: 0.3)
```

2438243924402441

2442

2443

2444

2445

2446

2447

2448

2449 2450

2451

2452

2453

2454

2455

2456

2457

2458

2459

24602461246224632464

2465

2466

2467

24682469

2470247124722473

2474

2475

2476

2477

2478 2479

# • Openness: 63

• Openness: 63 (confidence: 0.5)

• Conscientiousness: 66 (confidence: 0.7)

• Extraversion: 39 (confidence: 0.5)

• Agreeableness: 50 (confidence: 0.2)

• Neuroticism: 37 (confidence: 0.5)

#### Scenario 2: The Dispute

**System**: As you continue your journey, you encounter two hikers in a heated argument over a map. One claims they should follow the left path to reach the summit faster, while the other insists the right path is safer. They notice you and ask for your opinion. How do you respond? **Options**:

- 1. Carefully examine the map and give your honest opinion, even if it might disappoint one of them
- 2. Suggest they compromise by checking out both paths for a short distance before deciding
- 3. Support whoever seems more knowledgeable or confident
- 4. Offer to join them so you can all explore together and ensure everyone's safety

## **User selects**: 2 (Compromise)

#### LLM Analysis:

```
Trait updates:
openness: 55 (confidence: 0.1)
conscientiousness: 60 (confidence: 0.2)
extraversion: 50 (confidence: 0.1)
agreeableness: 75 (confidence: 0.4)
neuroticism: 45 (confidence: 0.2)
```

#### Profile after Scenario 2

• Openness: 62 (confidence: 0.6)

• Conscientiousness: 65 (confidence: 0.9)

• Extraversion: 40 (confidence: 0.5) • Agreeableness: 78 (confidence: 0.6)

• **Neuroticism**: 37 (confidence: 0.7)

#### Scenario 3: The Team Project

**System**: You've joined a group working on an important project with a tight deadline. During the first meeting, you notice that ideas are being suggested, but the discussion lacks structure. What do you do?

#### **Options**:

- 1. Take charge and propose a clear plan for the group to follow
- 2. Let others lead and offer support where needed
- 3. Privately suggest improvements to whoever seems to be leading
- 4. Suggest breaking into smaller groups to tackle different aspects simultaneously

**User selects**: 4 (Suggest breaking into smaller groups)

# 2480248124822483

2535

2536

2537

#### Profile after Scenario 3

• Openness: 62 (confidence: 0.7)

• Conscientiousness: 67 (confidence: 1.0)

• Extraversion: 33 (confidence: 0.8)

• Agreeableness: 72 (confidence: 0.7)

• Neuroticism: 40 (confidence: 0.8)

#### Scenario 4: The Unexpected Opportunity

**System**: You receive an unexpected invitation to join a last-minute expedition to a remote location. The trip leaves tomorrow and would require rearranging your schedule for the next week. How do you respond?

#### **Options:**

- 1. Immediately accept you love spontaneous adventures
- 2. Firmly decline you prefer to stick to your plans
- 3. Ask for more details before deciding
- 4. Suggest postponing the trip to a more convenient time

**User selects**: 3 (Ask for more details)

#### Profile after Scenario 4

• Openness: 59 (confidence: 0.9)

• Conscientiousness: 72 (confidence: 1.1)

• Extraversion: 34 (confidence: 0.9)

• Agreeableness: 70 (confidence: 0.8)

• **Neuroticism**: 39 (confidence: 0.9)

#### Scenario 5: The Creative Challenge

**System**: You've been invited to participate in a creative challenge where participants must develop an innovative solution to a local environmental issue. What approach do you take? **Options**:

- 1. Research extensively first, then develop a practical, proven solution
- 2. Brainstorm wildly creative ideas, even if they seem impractical at first
- 3. Seek out collaborators with different perspectives to develop ideas together
- 4. Focus on small, incremental improvements to existing approaches

**User selects**: 1 (Research extensively)

#### Profile after Scenario 5

• Openness: 55 (confidence: 1.1)

• Conscientiousness: 78 (confidence: 1.3)

• Extraversion: 35 (confidence: 0.9)

• Agreeableness: 68 (confidence: 0.8)

• **Neuroticism**: 38 (confidence: 1.0)

2591

#### Scenario 6: The Conflict Resolution

**System**: Two friends have had a serious disagreement and haven't spoken in weeks. Both have separately asked for your advice on repairing the relationship. What do you do? **Options**:

- 1. Give each person advice separately based on their perspective
- 2. Listen to both sides and facilitate a conversation between them
- 3. Suggest they give each other space until emotions cool down
- 4. Tell them both to compromise and move on

User selects: 2 (Listen to both sides and facilitate a conversation)

#### Profile after Scenario 6

- Openness: 54 (confidence: 1.1)
- Conscientiousness: 77 (confidence: 1.3)
- Extraversion: 36 (confidence: 1.0)
- Agreeableness: 75 (confidence: 1.1)
- **Neuroticism**: 37 (confidence: 1.0)

#### Scenario 7: The Social Gathering

**System**: You arrive at a party where you only know the host. How do you typically approach this situation?

#### **Options:**

- 1. Actively introduce yourself to new people and initiate conversations
- 2. Find a small group engaged in an interesting discussion and join in
- 3. Find a quiet corner and wait for the host to introduce you
- 4. Help the host with preparations or serving to feel more comfortable

User selects: 3 (Find a quiet corner and wait for the host to introduce you)

#### Profile after Scenario 7

- Openness: 53 (confidence: 1.1)
- Conscientiousness: 76 (confidence: 1.3)
- Extraversion: 30 (confidence: 1.3)
- Agreeableness: 74 (confidence: 1.1)
- Neuroticism: 42 (confidence: 1.2)

#### Scenario 8: The Ethical Dilemma

**System**: You discover that a colleague has taken credit for work you completed. The project was successful and got positive attention from management. What do you do? **Options**:

- 1. Confront the colleague publicly to ensure everyone knows the truth
- 2. Speak privately with the colleague about the situation
- 3. Report the issue to management or HR
- 4. Say nothing but ensure you get proper credit for future work

**User selects**: 2 (Speak privately with the colleague about the situation)

#### Profile after Scenario 8

• Openness: 52 (confidence: 1.1)

• Conscientiousness: 77 (confidence: 1.4)

• Extraversion: 32 (confidence: 1.3)

• Agreeableness: 70 (confidence: 1.3)

• Neuroticism: 40 (confidence: 1.2)

#### Scenario 9: The Life Change

**System**: You have the opportunity to move to a new city for a job that offers better pay but would require leaving your established social network. How do you approach this decision? **Options**:

- 1. Focus primarily on the career advancement opportunity
- 2. Prioritize your existing relationships and social connections
- 3. Take time to extensively research the new city and visit first
- 4. Create a detailed pros/cons list and analyze all factors methodically

User selects: 4 (Create a detailed pros/cons list and analyze all factors methodically)

#### Profile after Scenario 9

• Openness: 50 (confidence: 1.2)

• Conscientiousness: 80 (confidence: 1.5)

• Extraversion: 32 (confidence: 1.4)

• Agreeableness: 68 (confidence: 1.3)

• **Neuroticism**: 41 (confidence: 1.3)

#### Scenario 10: The Unexpected Crisis

**System**: While traveling, you encounter an unexpected emergency situation affecting several people. Resources are limited and tensions are high. How do you respond? **Options**:

- 1. Take charge and direct others to ensure efficiency
- 2. Find experts or authorities who can better handle the situation
- 3. Take a supportive role, helping organize resources and comfort others
- 4. Focus on solving one specific aspect of the problem thoroughly

User selects: 3 (Take a supportive role, helping organize resources and comfort others)

#### Profile after Scenario 10

• Openness: 51 (confidence: 1.3)

• Conscientiousness: 79 (confidence: 1.6)

• Extraversion: 35 (confidence: 1.5)

• Agreeableness: 72 (confidence: 1.5)

• **Neuroticism**: 38 (confidence: 1.4)

#### Scenario 11: The Creative Project

**System**: You have free time to pursue a personal project. What kind of activity are you most likely to choose?

#### **Options:**

- 1. A structured project with clear goals and measurable outcomes
- 2. A creative, open-ended project with room for exploration
- 3. A social activity involving coordination with others
- 4. Learning a specific new skill following an established method

**User selects**: 1 (A structured project with clear goals and measurable outcomes)

#### Profile after Scenario 11

- Openness: 53 (confidence: 1.4)
- Conscientiousness: 82 (confidence: 1.7)
- Extraversion: 34 (confidence: 1.5)
- Agreeableness: 71 (confidence: 1.5)
- Neuroticism: 37 (confidence: 1.4)

#### Scenario 12: The Public Speaking Opportunity

**System**: You've been asked to give a presentation to a large group on a topic you're knowledgeable about. How do you feel and prepare? **Options**:

- 1. Feel excited and prepare a dynamic, engaging presentation
- 2. Feel nervous but prepare extensively to manage anxiety
- 3. Feel neutral and focus on delivering clear, accurate information
- 4. Feel reluctant but try to find ways to make the presentation more interactive

**User selects**: 2 (Feel nervous but prepare extensively to manage anxiety)

#### Profile after Scenario 12

- Openness: 52 (confidence: 1.4)
- Conscientiousness: 83 (confidence: 1.8)
- Extraversion: 32 (confidence: 1.6)
- Agreeableness: 70 (confidence: 1.5)
- **Neuroticism**: 42 (confidence: 1.6)

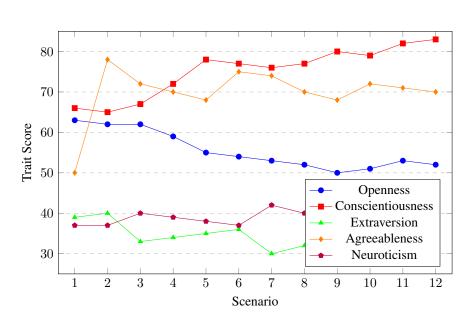


Figure 8: Progression of personality trait scores across scenarios

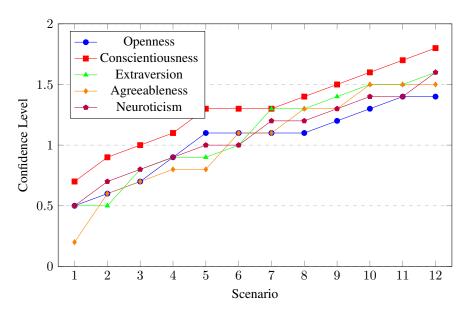


Figure 9: Progression of confidence levels for each trait across scenarios

2779

2780

2781

2782

27832784

2785

2786

2787

2788

2790

2791

2792

2793

2794

2795 2796

2797

2798

2799

2801

2804 2805 2806

#### 2754 2755 2756 52 (Moderconservative) -**Openness:** 2757 Balances practicality with some openness to new experiences 2758 **Conscientiousness:** 83 (Very High) -2759 Highly organized, disciplined, and detail-oriented 2760 **Extraversion:** 32 (Introverted) -Prefers quieter environments and one-on-one interactions 2761 Agreeableness: 70 (High) -2762 Cooperative, empathetic, values harmony 2763 **Neuroticism:** 42 (Low-Moderate) -2764 Generally emotionally stable with occasional anxiety 2765 2766 **Key Traits and Tendencies:** 2767 • Methodical approach to problem-solving (High C + Moderate O) 2768 • Prefers researching before acting (High C) 2769 • Values social harmony but not at the expense of principles (High A + High C) 2770 2771 • Reserved in social situations but empathetic (Low E + High A) 2772 • Prefers structured environments with clear expectations (High C) 2773 • Uncomfortable with sudden changes but adapts through planning (High C + Moder-2774 ate N) 2775 Selective about social engagements but loyal in relationships (Low E + High A) 2776 2777

Figure 10: Complete adventure-based personality assessment session showing progressive trait refinement across 12 scenarios. Each scenario contributes to increasing confidence in trait measurements and ultimately produces a stable, high-confidence personality profile.

This comprehensive example demonstrates how our adventure-based personality assessment system builds a detailed psychological profile through multiple scenarios. The assessment process illustrates several key aspects:

- Progressive refinement: Trait estimates become increasingly stable as evidence accumulates across scenarios
- 2. **Confidence building**: Confidence values steadily increase, reaching robust levels (1.4-1.8) by the end
- Domain coverage: Scenarios span diverse life domains including social situations, work environments, ethical dilemmas, and creative challenges
- 4. **Trait interrelationships**: The system identifies how patterns across scenarios reveal characteristic trait combinations
- 5. Adaptivity: The system selects scenarios to target traits with lower confidence values
- 6. **Format variety**: Scenarios present different types of choices to elicit a comprehensive range of behaviors

The final personality profile provides a nuanced psychological portrait that becomes the foundation for initializing the participant's GNWT-Agent digital twin. By capturing this level of psychological detail, our system ensures that the agent's behavior authentically represents the individual across different social contexts.

#### A.33 DETAILED RESULTS

Table 9: Match decision accuracy using different preference models

Preference Model	Match Prediction Accuracy	Human-Agent Match Correlation
Static (Time 1 only)	$58.9\% \pm 3.0\%$	$0.53 \pm 0.04$
Partial Evolution (Time 1+2)	$69.4\% \pm 2.5\%$	$0.65 \pm 0.03$
Full Evolution (All time points)	$77.8\% \pm 2.0\%$	$0.73 \pm 0.03$
Human (Ground Truth)	100%	1.00

Table 10: Human verification experiment results across two social contexts

Metric	Speed Dating Study (n=20)	Job Interview Study (n=10)
Behavioral fidelity rating	$5.6/7.0 \pm 0.8$	$5.8/7.0 \pm 0.6$
Decision concordance	$74\% \pm 4.2\%$	$81\% \pm 5.3\%$
Personality trait correlation	$0.83 \pm 0.04$	$0.81 \pm 0.05$
Conversational authenticity	$5.4/7.0 \pm 0.9$	$5.6/7.0 \pm 0.7$
Psychological state tracking	$5.7/7.0 \pm 0.6$	$5.5/7.0 \pm 0.8$
Overall agent realism	$5.9/7.0 \pm 0.5$	$5.6/7.0 \pm 0.7$

Table 11: Human-Agent correlation in social dynamics evolution

Evolution	Human (T1→T2)		GNWT-Agent (T1→T2)			
Dimension	Change (%)	Correlation	Change (%)			
Partner Preference	Partner Preference Evolution					
Attractiveness	$+39.0\% \pm 4.8\%$		$+25.0\% \pm 3.2\%$			
Sincerity	$-16.6\% \pm 3.4\%$		$-10.5\% \pm 2.8\%$			
Intelligence	$-24.8\% \pm 3.3\%$	$0.73 \pm 0.04$	$-15.2\% \pm 2.7\%$	$0.86 \pm 0.03$		
Fun	$+1.3\% \pm 2.0\%$	$0.73 \pm 0.04$	$+5.8\% \pm 2.2\%$	$0.00 \pm 0.03$		
Ambition	$-7.0\% \pm 2.7\%$		$-4.5\% \pm 1.9\%$			
Shared Interests	$+9.8\% \pm 3.3\%$		$+9.7\% \pm 1.3\%$			
Self-Perception E	volution					
Attractiveness	$+0.3\% \pm 9.2\%$		$-0.5\% \pm 3.0\%$			
Sincerity	$-3.5\% \pm 12.3\%$		$-2.5\% \pm 2.7\%$			
Intelligence	$-1.9\% \pm 8.8\%$	$0.81 \pm 0.03$	$-1.2\% \pm 2.4\%$	$0.82 \pm 0.04$		
Fun	$-1.3\% \pm 10.8\%$	$0.61 \pm 0.03$	$-0.8\% \pm 2.5\%$	$0.62 \pm 0.04$		
Ambition	$-0.8\% \pm 11.8\%$		$-0.5\% \pm 2.1\%$			
Self-Other Gap	$0.8 \to 0.7$		$0.9 \rightarrow 0.7$			
External Evaluation Correlations (r-value)						
	Time 1	Time 2	Time 1	Time 2		
Attractiveness	$0.67 \pm 0.01$	$-0.01 \pm 0.02$	$0.65 \pm 0.04$	$0.15 \pm 0.03$		
Sincerity	$0.52 \pm 0.01$	$-0.02 \pm 0.02$	$0.50 \pm 0.05$	$0.12 \pm 0.03$		
Intelligence	$0.51 \pm 0.01$	$0.02 \pm 0.02$	$0.48 \pm 0.05$	$0.18 \pm 0.04$		
Fun	$0.69 \pm 0.01$	$0.07 \pm 0.02$	$0.64 \pm 0.05$	$0.20 \pm 0.04$		
Ambition	$0.44 \pm 0.01$	$-0.01 \pm 0.02$	$0.40 \pm 0.06$	$0.10 \pm 0.05$		
Shared Interests	$0.66 \pm 0.01$	$-0.06 \pm 0.02$	$0.62\pm0.05$	$0.08 \pm 0.04$		

Overall Human-Agent Correlation:  $0.72 \pm 0.04$ 

2863

2864

2865 2866

2867

2868

2869

2870

2871

2872 2873

2874

2875

2876

2877

2878

2879

2880

2881 2882

2901

2902

2903

2904

2905

2906

2907

2908

2909

2910

2911

2912

2913

2914

2915

#### A.34 PRIVACY PROTECTION - IMPLEMENTATION DETAILS

# B COMPREHENSIVE PRIVACY PROTECTION FRAMEWORK FOR COGNIPAIR DIGITAL TWINS

The deployment of psychologically authentic digital twins through our GNWT-Agent architecture requires privacy safeguards that go beyond conventional data security. Our enhanced framework addresses the unique challenges of protecting intimate psychological profiles while maintaining the behavioral fidelity necessary for meaningful social pairing in both dating and professional contexts. This comprehensive approach integrates differential privacy, cryptographic watermarking, federated learning, and novel consent mechanisms to create a multi-layered defense against privacy violations while preserving the utility of our digital twin technology.

#### B.1 FOUNDATIONAL PRIVACY ARCHITECTURE

Our privacy protection begins at the moment of personality trait extraction, where we implement a sophisticated differential privacy mechanism that protects individual psychological profiles while maintaining statistical validity for matching purposes. The core principle involves injecting calibrated noise into the Big Five personality traits before GNWT-Agent initialization, ensuring that no adversary can reconstruct exact psychological profiles even with access to multiple agent interactions. This process employs a dynamic privacy budget allocation that adapts based on trait sensitivity and user preferences, allowing users to specify higher protection for traits they consider particularly sensitive.

#### Algorithm 5 Adaptive Differential Privacy for Personality Protection

```
2883
                1: Input: Raw traits T = \{O, C, E, A, N\}, Sensitivity vector S, Global budget \varepsilon_{total}
2884
                2: Initialize: Protected traits T \leftarrow \emptyset, Remaining budget \varepsilon_{rem} \leftarrow \varepsilon_{total}
2885
                3: for each trait t_i \in T do
                          \begin{array}{l} s_i \leftarrow S[i] \cdot \text{UserSensitivityPreference}(t_i) \\ \varepsilon_i \leftarrow \frac{s_i}{\sum_j s_j} \cdot \varepsilon_{total} \end{array}
2887
                5:
                                                                                                                                  \Delta_i \leftarrow \text{ComputeSensitivity}(t_i, \text{HistoricalData})
                6:
2889
                           \sigma_i \leftarrow \sqrt{2\ln(1.25/\delta)} \cdot \frac{\Delta_i}{\varepsilon_i}
                                                                                                                                           2890
                           noise_i \leftarrow \mathcal{N}(0, \sigma_i^2)
2891
                           \tilde{t}_i \leftarrow \text{clip}(t_i + noise_i, 0, 1)
2892
                9:
2893
              10:
                           if |\tilde{t}_i - t_i| > \text{MaxDeviation then}
                           \tilde{t}_i \leftarrow t_i + \text{sign}(noise_i) \cdot \text{MaxDeviation}
              11:
                           end if
              12:
                           \tilde{T} \leftarrow \tilde{T} \cup \{\tilde{t}_i\}
              13:
2896
              14:
                           \varepsilon_{rem} \leftarrow \varepsilon_{rem} - \varepsilon_i
2897
              15:
                           LogPrivacyExpenditure(t_i, \varepsilon_i, timestamp)
              16: end for
2899
              17: Return T, \varepsilon_{rem}
2900
```

Beyond personality traits, our framework implements comprehensive attribute protection through multi-resolution coarsening that adapts based on the attribute's discriminative power and privacy sensitivity. Dating attributes such as attractiveness ratings, sincerity scores, and ambition levels undergo a sophisticated transformation process that preserves matching utility while preventing fine-grained tracking. This coarsening operates at multiple granularities simultaneously, creating a hierarchical representation that allows for progressive disclosure based on interaction depth and mutual consent levels.

The conversation anonymization pipeline represents one of our most innovative privacy mechanisms, transforming raw interactions into privacy-preserving feature vectors through a multi-stage process that preserves behavioral patterns while eliminating identifying information. This pipeline begins with entity recognition and replacement, where proper nouns, specific locations, and temporal references are identified using a fine-tuned NER model and replaced with generalized placeholders. Subsequently, semantic embedding ensures that conversation meaning is preserved while specific phrasings that could identify individuals are transformed into canonical representations. The pipeline concludes with k-anonymization clustering that ensures each conversation pattern is shared by at least 150 other users, preventing unique behavioral fingerprinting.

```
2916
        Algorithm 6 Hierarchical Conversation Anonymization Pipeline
2917
          1: Input: Raw conversation C, Anonymity threshold k, Context window w
2918
          2: Phase 1: Entity Recognition and Masking
2919
          3: entities \leftarrow NER(C)
                                                                                    2920
          4: for each e \in entities do
2921
                 type_e \leftarrow ClassifyEntity(e)
                                                                           ⊳ Person, location, organization, etc.
          5:
2922
          6:
                 gen_e \leftarrow GeneralizeEntity(e, type_e, HierarchyLevel)
2923
          7:
                 C \leftarrow \text{Replace}(C, e, gen_e)
2924
          8: end for
          9: Phase 2: Temporal Quantization
2925
         10: timestamps \leftarrow \text{ExtractTimestamps}(C)
2926
         11: for each t \in timestamps do
2927
         12:
                 \tilde{t} \leftarrow |t/(6 \text{ hours})| \cdot (6 \text{ hours})
2928
                 C \leftarrow \text{Replace}(C, t, \tilde{t})
         13:
2929
         14: end for
2930
         15: Phase 3: Semantic Embedding and Clustering
2931
         16: embeddings \leftarrow []
2932
         17: for each window w_i in SlidingWindow(C, w) do
2933
                 emb_i \leftarrow SBERT(w_i)
         18:

    Sentence transformer embedding

2934
         19:
                 emb_i \leftarrow PCA(emb_i, dims = 50)
                                                                                    2935
         20:
                 embeddings.append(emb_i)
         21: end for
2936
         22: clusters \leftarrow DBSCAN(embeddings, min samples = k)
2937
         23: Phase 4: Pattern Extraction and Validation
2938
         24: patterns \leftarrow \text{ExtractBehavioralPatterns}(clusters)
2939
         25: for each p \in patterns do
2940
                 support \leftarrow CountSupport(p, GlobalDatabase)
         26:
2941
         27:
                 if support < k then
2942
         28:
                     p \leftarrow \text{Generalize}(p)
                                                                             > Further generalization if unique
2943
         29:
                 end if
2944
         30: end for
2945
         31: Return AnonymizedFeatures(patterns)
2946
2947
```

2971

2972

2973

2974

2975

2976

2977

2978

2979

2980

2981

2982

2983

2984

29852986

3013 3014

3015

3016

3017

3018

3019

#### B.2 ADVANCED CRYPTOGRAPHIC WATERMARKING SYSTEM

The potential for malicious actors to impersonate digital twins or falsely attribute generated content represents a critical threat to both individual privacy and system integrity. Our watermarking system embeds undetectable yet verifiable signatures directly into the token generation process of GNWT responses, creating a cryptographic chain of authenticity that survives paraphrasing, translation, and even sophisticated adversarial attacks. This watermarking operates through biased token selection during the language model's decoding process, where specific tokens are subtly favored based on a cryptographic hash of the agent's unique key and the generation context.

The watermark generation process employs a sophisticated multi-layer approach that combines deterministic biasing with adaptive strength calibration. At each token position, we compute a cryptographic hash that determines which tokens receive probability boosts, with the boost magnitude adapting based on the semantic importance of the position and the surrounding context entropy. This adaptive calibration ensures that watermarks remain strong in content-rich portions of text while avoiding detection in formulaic or constrained passages. Furthermore, our system implements redundant watermarking at multiple linguistic levels—lexical, syntactic, and semantic—ensuring robustness against various forms of attack.

#### Algorithm 7 Multi-Layer Cryptographic Watermarking with Adaptive Strength

```
2987
           1: Input: Agent key k, Context C, Base probabilities P
2988
           2: Initialize: Watermarked response R \leftarrow [], Strength factor \alpha_0 \leftarrow 0.1
2989
           3: for position i in generation sequence do
2990
                    Layer 1: Lexical Watermarking
           4:
2991
           5:
                    h_{lex} \leftarrow \text{HMAC-SHA256}(k, i||C)
2992
                    bias_tokens \leftarrow SelectTokens(P, h_{lex} \mod 5)
           6:

    Select 5 tokens

2993
           7:
                    \alpha_i \leftarrow \alpha_0 \cdot \text{ContextEntropy}(C) \cdot \text{ImportanceScore}(i)
2994
           8:
                    for token t in bias_tokens do
2995
           9:
                         P[t] \leftarrow P[t] \cdot (1 + \alpha_i)
2996
          10:
                    end for
                    Layer 2: Syntactic Watermarking
2997
          11:
          12:
                    pos\_tag \leftarrow PredictPOS(C, i)
2998
          13:
                    if pos_{tag} \in \{NOUN, VERB\} then
2999
                    h_{syn} \leftarrow \text{HMAC-SHA256}(k, \text{pos\_tag}||i)
          14:
3000
                    syntax\_bias \leftarrow h_{syn} \mod |ValidTokens(pos\_tag)|
          15:
3001
                    P[\text{syntax\_bias}] \leftarrow P[\text{syntax\_bias}] \cdot 1.05
          16:
3002
          17:
                    end if
3003
          18:
                    Laver 3: Semantic Watermarking
          19:
                    semantic\_vec \leftarrow GetSemanticEmbedding(C)
          20:
                    h_{sem} \leftarrow \text{HMAC}(k, \text{QuantizeVector(semantic\_vec}))
3006
          21:
                    semantic_cluster \leftarrow h_{sem} \mod \text{NumSemanticClusters}
3007
          22:
                    P \leftarrow \text{BiasTowardCluster}(P, \text{semantic\_cluster}, 0.03)
                    t_{selected} \leftarrow \text{Sample}(P)
          23:
3008
          24:
                    R.append(t_{selected})
                    C \leftarrow C || t_{selected}
          25:
                                                                                                              ▶ Update context
3010
          26:
                    LogWatermarkStrength(i, \alpha_i, t_{selected})
3011
          27: end for
3012
          28: Return R
```

Watermark detection employs a multi-hypothesis testing framework that evaluates the statistical likelihood of observing the biased token patterns across all three layers. The detection algorithm maintains robustness against adversarial modifications through redundancy and error correction codes embedded in the watermark pattern. When content undergoes modification, our system can still detect watermarks with high confidence by aggregating evidence across multiple linguistic levels and employing majority voting among independent detection mechanisms.

3025

3026

3027

3028

3030

3031

3032

3033

3034

3035

3036

3037 3038

3062

#### B.3 FEDERATED LEARNING WITH HOMOMORPHIC ENCRYPTION

The evolution of user preferences over time reveals intimate psychological dynamics that require protection from even honest-but-curious servers. Our federated learning framework ensures that preference updates occur locally on user devices, with only encrypted aggregates transmitted to central servers. This system employs partially homomorphic encryption that allows computation on encrypted gradients without decryption, ensuring that individual preference trajectories remain completely private while still enabling system-wide learning and improvement.

The federated learning protocol implements sophisticated privacy amplification through secure multiparty computation and differential privacy injection at multiple stages. Local devices compute preference updates using differentially private stochastic gradient descent, adding calibrated Gaussian noise to gradients before encryption. The aggregation server combines encrypted updates from multiple users using secure aggregation protocols that prevent the server from learning individual contributions even if some participants are malicious. This multi-layer approach ensures privacy against both external adversaries and insider threats while maintaining model convergence properties.

#### Algorithm 8 Privacy-Preserving Federated Preference Learning

```
3039
            1: Client-Side Processing:
3040
            2: Input: Local preferences P_i, Interactions I_i, Model \theta_t
3041
            3: g_i \leftarrow \nabla_{\theta} \mathcal{L}(P_i, I_i, \theta_t)
                                                                                                           3042
            4: \operatorname{clip\_norm} \leftarrow \min(1.0, C/||g_i||_2)
                                                                                                                  3043
            5: g_i \leftarrow g_i \cdot \text{clip\_norm}
3044
            6: \sigma^2 \leftarrow \frac{2C^2 \ln(1.25/\delta)}{2}
                                                                                                                   ▶ Noise calibration
3045
            7: \tilde{g}_i \leftarrow g_i + \mathcal{N}(0, \sigma^2 I)

    Add DP noise

3046
            8: mask_i \leftarrow GenerateSecretMask(seed_i)
3047
            9: \mathsf{masked}_i \leftarrow \tilde{g}_i + \mathsf{mask}_i

    ▶ Apply secret sharing

3048
           10: encrypted_i \leftarrow HomomorphicEncrypt(masked_i, pk)
3049
           11: Send encrypted_i to server
3050
           12:
3051
           13: Server-Side Aggregation:
3052
           14: Wait for minimum n_{min} participants
           15: aggregate \leftarrow \sum_{i=1}^{n} \text{encrypted}_{i}

16: noise_sum \leftarrow \sum_{i=1}^{n} \text{EncryptedNoise}_{i}
3053
                                                                                                           ▶ Homomorphic addition
3054
3055
           17: clean_aggregate ← aggregate − noise_sum
3056
           18: Secure Unmasking Protocol:
           19: for round r in 1 to \log(n) do
3057
           20: Pairs of clients exchange partial unmaskings
3058
           21: Server aggregates partial results homomorphically
3059
           22: end for
3060
           23: \theta_{t+1} \leftarrow \theta_t + \frac{1}{n} \cdot \text{Decrypt}(\text{clean\_aggregate})
3061
           24: Return \theta_{t+1} to all clients
```

3079

3080

3081

3082

3084

3085

3086

3087

3088

3089

3090

3091

3093

3121 3122

3123

3124

3125

3126

3127

3128

3129

3130

3131

#### B.4 APPLICATION-SPECIFIC PRIVACY MECHANISMS

For dating applications, our framework implements progressive disclosure mechanisms that reveal information gradually based on interaction depth and mutual interest signals. Initial interactions operate with maximum privacy protection, revealing only coarse-grained compatibility scores. As users express mutual interest through continued interaction, the system gradually releases more detailed preference alignments while maintaining plausible deniability about specific traits. This progressive disclosure follows a carefully calibrated schedule that balances privacy protection with the need for meaningful connection formation.

The dating platform also employs sophisticated rate limiting and interaction throttling to prevent adversaries from extracting information through repeated queries. Each digital twin maintains an interaction budget that regenerates slowly over time, preventing rapid-fire matching attempts that could be used to infer preferences. Furthermore, all interactions undergo temporal jittering and response caching to prevent timing attacks that could reveal genuine user activity patterns versus simulated twin responses.

#### Algorithm 9 Progressive Disclosure Protocol for Dating Interactions

```
3094
           1: Input: User twins A, B, Interaction history H_{AB}
           2: Initialize: Disclosure level L \leftarrow 0, Info revealed I_R \leftarrow \emptyset
           3: Stage 1: Initial Contact (L=0)
                                                                                                       ⊳ Single score 0-100
3097
           4: compat \leftarrow CoarseCompatibility(A, B)
           5: I_R \leftarrow I_R \cup \{\text{compat} \pm \mathcal{U}(-5,5)\}
                                                                                                                 ▶ Add noise
3098
           6: Stage 2: Mutual Interest (L = 1, after mutual likes)
3099
           7: trait_similarity \leftarrow ComputeTraitSimilarity(A, B)
3100
           8: for each trait dimension d do
3101
                   reveal_d \leftarrow Bernoulli(p = 0.6)
           9:

    Randomly reveal some traits

3102
          10:
                   if reveal<sub>d</sub> then
3103
                   I_R \leftarrow I_R \cup \{\text{GeneralizedTrait}(d, \text{level} = \text{medium})\}
3104
          12:
3105
          13: end for
3106
          14: Stage 3: Extended Interaction (L = 2, after 5+ messages)
3107
          15: interests \leftarrow ExtractSharedInterests(A, B, H_{AB})
          16: I_R \leftarrow I_R \cup SampleInterests(interests, k = 3)
3108
          17: values \leftarrow InferValueAlignment(A, B)
3109
          18: I_B \leftarrow I_B \cup \{\text{CoarsenedValues}(\text{values})\}
3110
          19: Stage 4: Deep Engagement (L=3, after successful video date)
3111
          20: if MutualConsent(A, B) then
3112
          21: I_R \leftarrow I_R \cup \text{DetailedPreferences}(A, B)
3113
          22: but maintain k-anonymity on sensitive attributes
3114
          23: end if
3115
          24: Privacy Budget Check:
3116
          25: \varepsilon_{spent} \leftarrow \text{ComputePrivacyLoss}(I_R)
3117
          26: if \varepsilon_{spent} > \varepsilon_{max} then
3118
          27: HaltDisclosure()
                                                                                               > Prevent further revelation
          28: end if
3119
          29: Return I_R, L, \varepsilon_{spent}
3120
```

For job application scenarios, our privacy framework implements strict audit trails and explainability requirements that ensure fairness while protecting candidate privacy. Every decision made by a digital twin in professional contexts generates a cryptographically signed log entry that records the activated cognitive modules, their weights, and the information considered. These logs employ zero-knowledge proofs to demonstrate compliance with anti-discrimination regulations without revealing specific candidate attributes. The system can prove, for instance, that protected characteristics did not influence a decision without disclosing what those characteristics are.

The professional application framework also implements role-based access control with granular permissions that ensure recruiters and hiring managers only access information necessary for their specific evaluation stage. Initial screening might only reveal skill compatibility scores, while final interview stages could access more detailed cognitive assessments. Throughout this process, all

3133

3134

3135 3136

3137

3138

3139

3140

3141

3142

3143

3144

3145

3146

3147

data access generates immutable audit logs that candidates can review to understand exactly what information was shared and when.

#### B.5 CONSENT MANAGEMENT AND USER CONTROL

Our double consent protocol ensures that both human users and their digital twins must approve significant actions, preventing unauthorized commitments while maintaining user agency. This protocol employs a sophisticated state machine that tracks consent status across multiple interaction types and ensures that neither automated processes nor human impulses can override deliberate decisionmaking. The consent framework integrates with blockchain technology to create an immutable record of consent grants and revocations, providing legal protection for all parties involved.

User control extends beyond simple consent to include fine-grained preference settings for privacyutility trade-offs. Users can specify different privacy levels for different attributes, indicate which types of information can be inferred versus explicitly shared, and set decay functions that automatically increase privacy protection for older interactions. This granular control employs an intuitive interface that visualizes privacy implications in real-time, helping users understand the consequences of their choices without requiring technical expertise.

#### Algorithm 10 Blockchain-Anchored Double Consent Protocol

```
3148
3149
          1: Input: Action request A, User U, Digital Twin T, Context C
3150
          2: Phase 1: Twin Evaluation
3151
          3: T_{eval} \leftarrow T.EvaluateAction(A, C)
3152
          4: T_{confidence} \leftarrow T.ComputeConfidence(A)
3153
          5: T_{risks} \leftarrow T.AssessRisks(A, U.history)
3154
          6: if T_{confidence} < \tau_{min} or T_{risks} > \text{RiskThreshold} then
3155
          7: Return TWIN_REJECT, Explanation(T_{eval}, T_{risks})
3156
          9: T_{consent} \leftarrow \text{SignWithTwinKey}(A, T_{eval}, \text{timestamp})
3157
3158
         11: Phase 2: Human Verification
3159
         12: notification \leftarrow PrepareNotification(A, T_{eval}, T_{confidence})
3160
         13: SendToUser(U, notification, timeout = 24h)
3161
         14: U_{response} \leftarrow \text{WaitForResponse}(U, \text{timeout})
3162
         15: if U_{response} = NULL then
3163
         16: Return TIMEOUT, DefaultAction(A)
         17: end if
3165
3166
         19: Phase 3: Consensus Formation
         20: if U_{response}.decision = APPROVE then
3167
         21: U_{consent} \leftarrow \text{SignWithUserKey}(A, U_{response}, \text{timestamp})
3168
         22: consensus \leftarrow CombineConsents(T_{consent}, U_{consent})
3169
         23: block_hash ← BlockchainRecord(consensus)
3170
         24: ExecuteAction(A, consensus, block hash)
3171
         25: Return SUCCESS, block_hash
3172
         26: else
3173
         27: override_reason \leftarrow U_{response}.reason
3174
         28: LogOverride(U, T, A, override\_reason)
3175
         29: T.UpdatePreferences(override reason)
                                                                                          3176
         30: Return USER_OVERRIDE, override_reason
3177
         31: end if
```

#### B.6 DATA RETENTION AND RIGHT TO ERASURE

3186

3187

3188

3189

3190

3191

3192

3193

3194

3195

3196

3197

3198

3199 3200

3201

3232

3234 3235

3237 3238 Our data retention policies implement cryptographically-enforced expiration where data automatically becomes inaccessible after predetermined periods unless explicitly renewed through user consent. This system employs forward-secure encryption where encryption keys are automatically deleted after expiration, ensuring that even archived data cannot be recovered. The retention framework operates at multiple granularities, with conversation snippets expiring after 30 days, preference patterns after 90 days, and core personality models after one year unless actively renewed.

The right to erasure goes beyond simple deletion to implement cryptographic shredding that renders all derived data permanently inaccessible. When users invoke their right to be forgotten, the system not only deletes their primary data but also retrains affected models to remove their statistical influence. This process employs machine unlearning techniques that provably eliminate the user's contribution to aggregate models while maintaining performance for remaining users. Furthermore, our system maintains deletion certificates that provide cryptographic proof of erasure, giving users verifiable confirmation that their data has been permanently removed.

#### Algorithm 11 Cryptographic Data Shredding with Verified Unlearning

```
1: Input: User U, Deletion request D, Verification key V
3202
          2: Step 1: Immediate Access Revocation
3203
          3: RevokeAllTokens(U)
3204
          4: DisableDigitalTwin(U.twin_id)
3205
          5: NotifyConnectedUsers(U.connections)
3206
3207
          7: Step 2: Cryptographic Key Destruction
          8: keys \leftarrow \text{GetAllUserKeys}(U)
3209
          9: for each k \in keys do
         10:
                OverwriteMemory(k, random, passes = 7)
3211
                DestroyHSMEntry(k.hsm_id)
         11:
3212
         12: end for
         13:
3213
         14: Step 3: Machine Unlearning
3214
         15: affected_models \leftarrow IdentifyAffectedModels(U)
3215
         16: for each model M in affected_models do
3216
                influence \leftarrow ComputeInfluence(U, M)
         17:
3217
         18:
                if influence > \tau then
3218
         19:
                M' \leftarrow \text{RetrainWithout}(M, U.\text{training\_data})
3219
         20:
                else
3220
         21:
                M' \leftarrow M - \text{influence} \cdot \nabla_U M
                                                                               3221
         22:
                end if
3222
         23:
                VerifyUnlearning(M', U)
                                                                                            3223
         24:
                DeployUpdatedModel(M')
         25: end for
3224
         26:
3225
         27: Step 4: Audit and Certification
3226
         28: audit_{log} \leftarrow GenerateAuditTrail(U, D)
3227
         29: deletion_proof \leftarrow ZKProof(DataAbsence(U))
3228
         30: certificate \leftarrow Sign(audit_log, deletion_proof, V)
3229
         31: BlockchainRecord(certificate)
3230
         32: Return certificate, blockchain_hash
3231
```

#### B.7 PRIVACY METRICS AND CONTINUOUS MONITORING

Our privacy framework implements comprehensive metrics that continuously monitor privacy protection effectiveness across multiple dimensions. Re-identification risk is assessed through regular adversarial testing where red teams attempt to link anonymized profiles back to individuals using various attack strategies. These tests employ state-of-the-art linkage attacks, membership inference attempts, and attribute inference techniques to ensure our defenses remain robust against evolving threats. The system maintains a real-time privacy dashboard that tracks key metrics including k-anonymity levels, differential privacy budget consumption, and information leakage rates across all components.

Behavioral uniqueness metrics quantify how distinguishable individual users remain after our anonymization processes. We employ information-theoretic measures to assess the entropy of user behavior patterns and ensure sufficient uncertainty to prevent tracking. The monitoring system automatically triggers additional privacy protections when uniqueness scores exceed thresholds, dynamically adjusting noise levels and generalization parameters to maintain target privacy levels. This adaptive approach ensures consistent privacy protection even as user behavior patterns evolve and new interaction modalities are introduced.

The watermark detection system undergoes continuous calibration through automated testing against adversarial paraphrasing attacks. We maintain a library of evasion techniques including synonym substitution, sentence restructuring, back-translation, and GPT-based paraphrasing, constantly evaluating our watermark robustness against these attacks. Detection thresholds are dynamically adjusted to maintain false positive rates below  $10^{-6}$  while maximizing true positive rates, with regular updates to handle new attack vectors as they emerge.

#### B.8 REGULATORY COMPLIANCE AND LEGAL FRAMEWORK

Our privacy framework ensures comprehensive compliance with international data protection regulations including GDPR, CCPA, and emerging AI governance frameworks. For GDPR compliance, we implement privacy by design principles throughout the system architecture, ensuring that privacy protection is not an afterthought but a fundamental design consideration. Our lawful basis for processing relies on explicit consent with granular opt-in mechanisms for different data uses. The system maintains detailed records of processing activities, implements data protection impact assessments for new features, and ensures cross-border data transfers comply with adequacy decisions and standard contractual clauses.

For employment contexts, our framework adheres to EEOC guidelines and international fair hiring practices through algorithmic auditing and bias mitigation. The system implements disparate impact testing that continuously monitors hiring recommendations for discriminatory patterns across protected categories. When bias is detected, the system automatically adjusts cognitive module weights and decision thresholds to restore fairness while maintaining prediction accuracy. All employment-related decisions generate detailed explanations that can be reviewed by human resources professionals and challenged by candidates through established appeal processes.

The legal framework extends to intellectual property protection for generated content, implementing clear ownership models and usage rights for twin-generated text. Our terms of service establish that users retain ownership of their digital twin's outputs while granting necessary licenses for system operation. The watermarking system provides forensic capabilities for intellectual property disputes, enabling verification of content origin and protecting against false attribution claims. We maintain comprehensive insurance coverage for privacy breaches and implement incident response procedures that ensure rapid notification and remediation in case of any privacy violations.

#### B.9 FUTURE DIRECTIONS AND RESEARCH CHALLENGES

While our current privacy framework provides strong protection against known threats, several research challenges remain for future development. Adversarial machine learning attacks that attempt to extract training data or manipulate model behavior require ongoing defensive research and regular model hardening. We are investigating certified defenses that provide provable robustness guarantees against specific attack classes. The emergence of quantum computing threatens current cryptographic protections, necessitating migration to quantum-resistant algorithms for long-term data protection. Our roadmap includes systematic replacement of classical cryptographic primitives with lattice-based and hash-based alternatives that resist quantum attacks.

The intersection of privacy and explainability presents fundamental tensions that require novel technical solutions. Users deserve explanations for their digital twins' decisions, but detailed explanations

can leak private information about the training process and other users. We are developing privacy-preserving explanation techniques that provide meaningful insights while maintaining differential privacy guarantees. This includes research into synthetic explanations that capture decision logic without revealing specific training examples and federated explanation generation that aggregates insights across multiple users while preserving individual privacy.

The scalability of privacy-preserving technologies remains a significant challenge as user bases grow and interaction complexity increases. Homomorphic encryption and secure multi-party computation introduce substantial computational overhead that can impact system responsiveness. Our research focuses on optimizing cryptographic protocols for specific computation patterns in digital twin interactions and developing hardware acceleration for privacy-preserving operations. We are also investigating selective privacy relaxation where users can choose to trade privacy for performance in non-sensitive contexts while maintaining strong protection for critical interactions.

#### B.10 CONCLUSION

The comprehensive privacy protection framework presented here establishes a new standard for protecting user privacy in psychologically-authentic digital twin systems. Through the integration of differential privacy, cryptographic watermarking, federated learning, and sophisticated consent mechanisms, our framework ensures that CogniPair can deliver meaningful social pairing while respecting user privacy and autonomy. The multi-layered approach provides defense in depth against various attack vectors while maintaining the utility necessary for practical deployment. As digital twins become increasingly sophisticated and prevalent, the privacy protections outlined here will be essential for maintaining user trust and preventing potential harms. Our ongoing research continues to strengthen these protections while exploring new frontiers in privacy-preserving artificial intelligence, ensuring that the benefits of digital twin technology can be realized without compromising fundamental privacy rights.