

A Survey on Large Language Model based Human-Agent Systems

Anonymous ACL submission

Abstract

Recent advances in large language models (LLMs) have sparked growing interest in building fully autonomous agents. However, fully autonomous LLM-based agents still face significant challenges, including limited reliability due to hallucinations, difficulty in handling complex tasks, and substantial safety and ethical risks, all of which limit their feasibility and trustworthiness in real-world applications. To overcome these limitations, LLM-based human-agent systems (LLM-HAS) incorporate human-provided information, feedback, or control into the agent system to enhance system performance, reliability and safety. This paper provides the first comprehensive and structured survey of LLM-HAS. It clarifies fundamental concepts, systematically presents core components shaping these systems, including environment & profiling, human feedback, interaction types, orchestration and communication, explores emerging applications, and discusses unique challenges and opportunities. By consolidating current knowledge and offering a structured overview, we aim to foster further research and innovation in this rapidly evolving interdisciplinary field. Paper lists and resources are available at [GitHub repository](#).

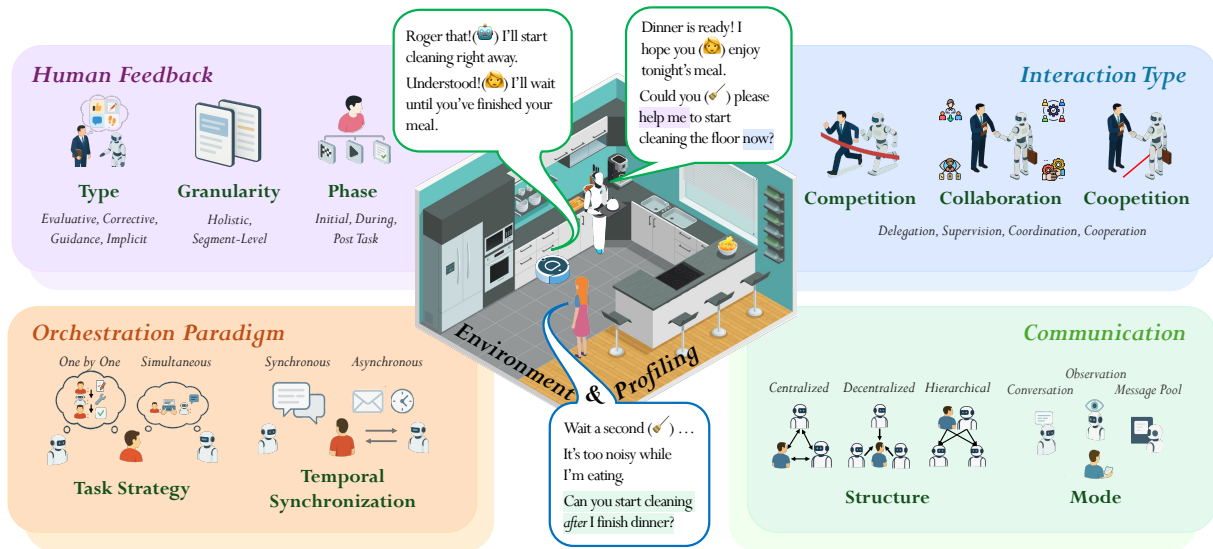
1 Introduction

Recent advances in Large Language Models (LLMs) have led to growing enthusiasm for building fully autonomous agent systems that use LLMs as a central engine to perceive environments, make decisions, and execute actions to achieve goals (Wang et al., 2024a; Li et al., 2024a). These agents are often equipped with modules for memory, planning, and tool use, aiming to automate complex workflows with minimal human involvement (Xie et al., 2024a; Xi et al., 2025). However, the pursuit of *full autonomy* faces critical hurdles. (1) **Reliability** remains a major concern due to LLMs’ propensity for hallucination, gener-

ating plausible but factually incorrect or nonsensical outputs, which undermines trust and can lead to significant errors, especially when actions are chained (Gosmar and Dahl, 2025; Xu et al., 2024; Glickman and Sharot, 2025). (2) **Complexity** often stalls autonomous agents; they struggle with very complicated tasks requiring deep domain expertise, long multi-step execution, nuanced reasoning, dynamic adaptation, or strict long-context consistency dependencies, as seen in scientific research (Feng et al., 2024; Yehudai et al., 2025). (3) **Safety and Ethical Risks** escalate with autonomy; agents can take unintended harmful actions, amplify societal biases present in training data, or create accountability gaps, particularly in critical decision-making scenarios involving finance, healthcare, or security (Mitchell et al., 2025; Deng et al., 2024; Wang et al., 2024c).

The persistence of these challenges suggests that full autonomy may be unsuitable for many real-world applications (Mitchell et al., 2025; Natarajan et al., 2025) and underscores a crucial insight often overlooked in the drive for pure automation: the indispensable role of human involvement. Humans are frequently needed to provide additional information, essential clarification, or domain knowledge, offer vital feedback and corrections, and exercise necessary oversight and control. These motivate a paradigm shift towards systems explicitly designed for human-agent collaboration: **LLM-based Human-Agent Systems (LLM-HAS)**.

While surveys on LLM-based autonomous agents (Wang et al., 2024a; Li et al., 2024a), multi-agent systems (Tran et al., 2025; Wu et al., 2025b), and specific applications exist (Wang et al., 2025b; Peng et al., 2025), a dedicated synthesis focusing specifically on LLM-based human-agent systems is lacking. This survey fills the gap by providing a comprehensive and structured overview of the LLM-HAS. It clarifies the fundamental concepts (Section 2) and systematically presents its



LLM-based Human-Agent Systems (LLM-HAS)

Figure 1: Overview of LLM-based Human-Agent Systems (LLM-HAS). LLM-HAS are interactive frameworks where humans actively provide additional information, feedback, or control during interaction with an LLM-powered agent to enhance system performance, reliability, and safety. The system is composed of five core components: **Environment & Profiling** (including environment settings, and role definitions, goals, and agent capabilities such as planning and memory), **Human Feedback** (with varying types, timing, and granularity), **Interaction Types** (collaborative, competitive, cooperative, or mixed), **Orchestration** (task strategy and temporal synchronization), and **Communication** (information flow structure and mode).

084 core components (Section 3), major implementa- 110
 085 tion strategies (Section 4), emerging applica- 111
 086 tions and resources (Section 5), and unique chal- 112
 087 lenges and opportunities (Section 6) within this specific 113
 088 niche. To the best of our knowledge, this is still the 114
 089 first survey on LLM-based human-agent systems. 115
 090 We aim to consolidate current knowledge and in- 116
 091 spire further research and application in this rapidly 117
 092 evolving field. A [GitHub repository](#) is maintained 118
 093 to provide a sustainable resource complementing 119
 094 our survey paper. 120

095 2 LLM-Based Human-Agent Systems 121

096 We define LLM-based human-agent systems as in- 124
 097 teractive frameworks where humans actively pro- 125
 098 vide additional information, feedback, or control 126
 099 during interaction with an LLM-powered agent 127
 100 to enhance system performance, reliability, and 128
 101 safety (Feng et al., 2024; Shao et al., 2024; Mehta 129
 102 et al., 2024). The core idea is **synergy**: combining 130
 103 unique human strengths—like intuition, creativity, 131
 104 expertise, ethical judgment, and adaptability—with 132
 105 LLM agent capabilities such as vast knowledge 133
 106 recall, computational speed, and sophisticated lan- 134
 107 guage processing. LLM-HAS builds upon core 135
 108 LLM agent components but places critical empha- 136
 109 sis on the human’s interactive role and capabilities:

- (1) **Providing Information / Clarification:** Hu- 110
 mans provide additional information that 111
 agents lack or cannot reliably infer, such as 112
 login credentials, payment details, domain ex- 113
 pertise, constraints, or resolve ambiguities, 114
 helping agents interpret situations more ac- 115
 curately (Naik et al., 2025; Kim et al., 2025b). 116
- (2) **Providing Feedback / Error Correction:** Hu- 117
 mans evaluate agent outputs and provide feed- 118
 back, ranging from simple ratings to complex 119
 critiques, demonstrations or corrections, effec- 120
 tively guiding agents’ adjustment (Gao et al., 121
 2024b; Dutta et al., 2024; Li et al., 2024b). 122
- (3) **Taking Control / Action:** In high-stakes or 123
 sensitive scenarios (e.g., healthcare, privacy, 124
 or ethics), humans retain the authority to over- 125
 ride, redirect, or halt agent actions, ensuring 126
 accountability, safety, and alignment with hu- 127
 man values (Chen et al., 2025b; Natarajan 128
 et al., 2025; Xiao and Wang, 2023). 129

130 Figure 1 provides a generalized overview of 130
 131 LLM-based human-agent systems. These systems 131
 132 operate within a defined **Environment** (e.g., phys- 132
 133 ical world, simulation) that provides context and 133
 134 stimuli. **Human & Agent Profiling** characterizes 134
 135 the participants’ roles and goals, and the agent’s 135
 136 core LLM engine augmented with capabilities like 136

137 planning, memory, and tool use. **Human Feed-**
 138 **back** can occur during different phases in various
 139 types and granularities. Human-Agent **Interaction**
 140 **Types** may be collaborative (most common), com-
 141 petitive, cooperative, or mixed. The **Orchestration**
 142 layer governs high-level coordination, choosing a
 143 task strategy (e.g., sequential one-by-one versus
 144 parallel simultaneous execution) and a temporal
 145 synchronization mode (real-time synchronous ex-
 146 changes versus delayed asynchronous workflows)
 147 so that each actor acts at the right moment. The
 148 **Communication** layer specifies how information
 149 flows, defining message structure (centralized, de-
 150 centralized, hierarchical) and mode (conversation,
 151 observation signals, or shared message pools). The
 152 effective interplay and configuration of these com-
 153 ponents, along with various human feedback, are
 154 critical for tailoring the system to specific tasks and
 155 optimizing the overall system’s performance. The
 156 taxonomy of LLM-based human-agent systems is
 157 outlined in Figure 3. A detailed and structured cat-
 158 egorization of representative works is provided in
 159 the Table 6 and Table 7.

160 3 Core Components

161 In this section, we examine LLM-HAS through
 162 five core aspects: environment & profiling, human
 163 feedback, interaction type, orchestration paradigm,
 164 and communication. These dimensions provide a
 165 unified standard for analyzing existing work and
 166 guiding the design of future systems.

167 3.1 Environment and Profiling

168 **Environment Setting.** The environment in
 169 LLM-HAS defines a shared interaction space that
 170 can exist either in the physical world, such as
 171 offices (Sun et al., 2024b), or in fully simulated
 172 virtual environments where agents and humans
 173 engage under controlled conditions (Sun et al.,
 174 2024b; Zhang et al., 2024a; Guo et al., 2024b).
 175 These systems can be configured in various ways,
 176 including single-human single-agent, single-
 177 human multi-agent, multi-human single-agent, and
 178 multi-human multi-agent setups, each reflecting
 179 different collaboration dynamics and complexities.

180 **Human & Agent Profiling.** Human participants
 181 can be broadly categorized as *lazy* or *informative*
 182 users. Lazy users provide minimal guidance, typ-
 183 ically offering evaluative feedback such as binary
 184 correctness or scalar rating. In contrast, informa-

186 tive users engage deeply by offering demonstra-
 187 tions, detailed guidance, refinements, or even tak-
 188 ing over parts of the task (Wang et al., 2024b; Liu
 189 et al., 2024b; Han et al., 2025). On the other side,
 190 agents are profiled by their roles and capabilities,
 191 which range from versatile general assistants to
 192 specialized experts in mathematics, engineering,
 193 medicine, or robotic cleaning, each adapted to the
 194 particular demands of its operational context (Guo
 195 et al., 2024a; Samuel et al., 2024).

196 3.2 Human Feedback

197 **Human Feedback Type.** We categorize human
 198 feedback as *evaluative*, *corrective*, *guidance*, and
 199 *implicit* feedback. (1) **Evaluative Feedback** pro-
 200 vides an assessment of the agent’s output quality,
 201 typically as preference ranking, scalar rating, or
 202 binary assessment. A prime example is preference
 203 ranking, where users compare agent outputs,
 204 forming the basis of Reinforcement Learning
 205 from Human Feedback (RLHF) (Chaudhari et al.,
 206 2024). Alternatively, platforms like Uni-RLHF
 207 (Yuan et al., 2024) support scalar ratings or binary
 208 assessments. (2) **Corrective Feedback** offers direct
 209 edits or fixes to the agent’s behavior. For instance,
 210 the PRELUDE (Gao et al., 2024a) framework
 211 learns latent preferences from user edits made
 212 to agent-generated text. (3) **Guidance Feedback**
 213 means the human proactively provides instructions,
 214 critiques, or demonstrations to shape the agent’s
 215 behavior. Agents like InteractGen (Sun et al.,
 216 2024b), AutoManual (Chen et al., 2024a) can be
 217 bootstrapped using initial demonstrations, while
 218 methods like Self-Refine (Choudhury and Sodhi,
 219 2025) employ iterative critiques and refinements to
 220 improve outputs. (4) **Implicit Feedback** is inferred
 221 by the agent observing user actions or control
 222 signals, rather than explicitly stated or direct output
 223 modifications. For example, an agent might learn
 224 user priorities by observing how a user adjusts
 225 control sliders in a system like VeriPlan (Lee et al.,
 226 2025a), or infer preferences by analyzing user
 227 behaviors like clicks and purchases in frameworks
 228 such as AgentA/B (Wang et al., 2025a). This
 229 contrasts with corrective feedback, where the user
 230 directly edits the output; here, the agent interprets
 231 the user’s independent actions or control choices.

232 **Human Feedback Granularity.** Human feed-
 233 back also varies in granularity, from coarse-grained,
 234 holistic judgments to fine-grained, segment-level
 235 critiques. (1) **Coarse-grained/Holistic feedback**
 236

Dimension	Category	Definition Summary	Key Characteristics / Trade-offs	Example Work
Type	<i>Evaluative</i>	User provides an assessment of the agent’s output quality, typically as binary assessment, scalar rating, or preference ranking .	① Easy to collect, scalable. ② Less specific signal for improvement.	<i>EmoAgent</i> (Qiu et al., 2025), <i>MINT</i> (Wang et al., 2024b), <i>SOTOPIA</i> (Zhou et al., 2024)
	<i>Corrective</i>	User offers edits or fixes to the agent’s behavior.	① Highly informative, clear signal for improvement. ② Higher user effort, often fine-grained & interactive.	<i>SymbioticRAG</i> (Sun et al., 2025a), <i>SWEET-RL</i> (Zhou et al., 2025), <i>AI Chains</i> (Wu et al., 2022a)
	<i>Guidance</i>	User proactively provides instructions, demonstrations, or critiques to shape the agent’s behavior.	① Bootstraps learning, conveys complex goals, proactive alignment. ② Requires clear specification from user.	<i>Drive As You Speak</i> (Cui et al., 2024), <i>Hierarchical Agent</i> (Liu et al., 2023b), <i>Ask Before Plan</i> (Zhang et al., 2024c)
	<i>Implicit</i>	Inferred by the agent observing user actions or control signals , rather than explicitly stated or direct output modifications.	① Natural, unobtrusive collection. ② Ambiguous, requires careful interpretation.	<i>MTOM</i> (Zhang et al., 2024b), <i>Attentive Supp.</i> (Tanneberg et al., 2024), <i>MineWorld</i> (Guo et al., 2025)
Granularity	<i>Coarse-grained / Holistic</i>	Single assessment/signal for an entire agent output, trajectory, or task outcome .	① Simple for user, good for overall assessment ② Obscures specific errors, less precise learning signal.	<i>AssistantX</i> (Sun et al., 2024a), <i>Help Feedback</i> (Mehta et al., 2024), <i>AXIS</i> (Lu et al., 2024)
	<i>Fine-grained / Segment-Level</i>	Feedback targeting specific parts of agent output, actions, or process .	① Precise learning signal, crucial for debugging complex skills ② Potentially higher user effort/burden.	<i>Collaborative Gym</i> (Shao et al., 2024), <i>Prison Dilemm</i> (Jiang et al., 2025), <i>FineArena</i> (Xu et al., 2025)
Phase	<i>Initial Setup & Goal Definition</i>	Feedback provided before task execution, configuring the agent system and defining the task, goals, constraints, and preference .	① Initial and proactive alignment, prevents costly errors, sets constraints ② Requires upfront user input.	<i>AgentCoord</i> (Pan et al., 2024a), <i>GDfC</i> (Wang et al., 2025c), <i>SMALL</i> (Wang et al., 2024c)
	<i>During Task Execution</i>	Online, interactive feedback while the agent is actively performing the task , enabling real-time adaptation .	① Enables real-time adaptation, crucial for dynamic/collaborative tasks ② Requires timely notification and responsive interfaces.	<i>InteractGen</i> (Sun et al., 2024b), <i>CowPilot</i> (Huq et al., 2025), <i>EasyLAN</i> (Pan et al., 2024b)
	<i>Post-Task Eval. & Refinement</i>	Feedback provided after task completion to assess outcomes and provide suggestions for future use .	① Non-disruptive, good for aggregate data/offline learning ② No impact on completed task.	<i>HRT-ML</i> (Liu et al., 2024b), <i>M3HF</i> (Wang et al., 2025d), <i>MAIH</i> (Wang et al., 2024c)

Table 1: Dimensions of Human Feedback in LLM-based human-agent systems, including feedback type, granularity, and phase. For each dimension, a summary, key characteristics, and example works are provided for comparison. A detailed overview of human feedback types and their subtypes is provided in our appendix (Table 5).

provides a single assessment for the entire agent output. Standard RLHF often relies on holistic preferences between complete responses, which simplifies feedback collection but struggles with credit assignment in complex tasks. (2) **Fine-grained/Segment-Level Feedback** by contrast, targets specific parts (e.g., sentences, paragraphs, code blocks). This is crucial in environments like ConvCodeWorld (Han et al., 2025), where feedback pertains to specific conversational turns or generated code segments, or in annotation tasks like PDFChatAnnotator (Tang et al., 2024), where feedback applies to specific annotations or parts of the document. This finer granularity provides more precise learning signals, crucial for debugging complex behaviors.

Human Feedback Phase. Human feedback can be incorporated at different phases of the LLM-

agent pipeline (Wang et al., 2025d). (1) **Initial Setup & Goal Definition** occurs before task execution, configuring the agent system and defining goals, such as setting coordination strategies (AgentCoord (Pan et al., 2024a)) or critiquing plans before execution (Ask-before-Plan (Zhang et al., 2024c)). (2) **During Task Execution** involves online, interactive feedback while the agent is actively performing the task, enabling real-time adaptation. Examples include interactive instruction editing (InstructEdit (Wang et al., 2023)), mid-task refinements (Mutual Theory of Mind (Zhang et al., 2024b), Collaborative Gym (Shao et al., 2024)), or online interventions (HG-Dagger (Kelly et al., 2019)). (3) **Post-Task Evaluation & Refinement** happens after task completion to assess outcomes and provide feedback for future use. Frameworks like MAIH (Wang et al., 2024c) and EmoAgent

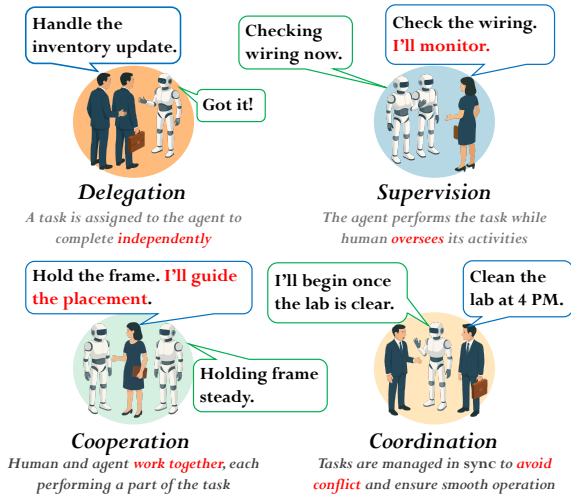


Figure 2: The subtype of the collaboration between humans and LLM-based agents.

(Qiu et al., 2025) apply feedback loops after initial generation for benchmarking or offline learning, while AdaPlanner (Sun et al., 2023) archives successful plans post-task as skills for future use. Table 1 summarizes different dimensions of human feedback, key characteristics, and example work.

3.3 Human-Agent Interaction Types

Interaction types define how individuals communicate, exchange information, and take actions with one another. In LLM-HAS, interactions tend to be more dynamic and complex compared to multi-agent systems. This complexity arises from the various roles and responsibilities assigned to both human agents and those based on LLMs, necessitating a finer-grained framework to describe their collaborative behaviors. The following categorization highlights the three key interaction types: **Collaboration**, **Competition**, and **Cooperation**.

3.3.1 Collaboration

Collaborations are by far the most common interaction and foundational interaction, which involve humans and LLM-based agents working together to achieve a common goal. This partnership combines human creativity and contextual understanding with LLM-based agents to address challenges and improve the efficiency and quality of results (Vats et al., 2024; Du et al., 2024; Sun et al., 2025a). Depending on the type of collaboration considered, it can be categorized into four main fine-grained subtypes (Figure 2): (1) **Delegation & Direct Command** (Kiewiet and McCubbins, 1991), (2) **Supervision** (Loganbill et al., 1982), (3) **Cooperation** (Rand and Nowak,

2013), and (4) **Coordination** (Turvey, 1990).

Delegation & Direct Command. In this interaction modality, a controlling party, usually a human, assigns explicit tasks to the LLM-based agent by providing clear and direct instructions. The agent is expected to execute these directives autonomously or on behalf of humans, ensuring that responsibilities are well-defined and actions align with the system’s overarching objectives. Unlike supervision, where strategies can be dynamically adjusted in response to new situations, delegation involves providing instructions upfront. This means the agent follows a predetermined set of tasks rather than adapting to the situation. For instance, an investor specifies their risk preference to the agent executing the investment strategy like FineArena (Xu et al., 2025), or a driver utters the command to LLM-based agent like Drive as you Speak (Cui et al., 2024).

Supervision. Supervision is the process by which one party, usually a human operator, oversees, monitors, and guides the actions of an LLM-based agent. This involves real time evaluation and intervention to ensure the agent’s output aligns with established goals and quality standards. Supervision also encompasses setting alert thresholds and providing corrective inputs when deviations occur. By maintaining a continuous feedback loop between the human and the agent, supervision helps calibrate agent behavior, catch and mitigate errors before they propagate and build confidence in the system. It also enables agents to handle routine tasks with increasing independence. For instance, agents notify humans to verify alignment (Liu et al., 2023b), and teleoperators monitor the LLM-generated motion plans (Liu et al., 2023a).

Cooperation. Cooperation refers to the voluntary and joint efforts of multiple parties to achieve agreed-upon goals. This collaboration type combines the various efforts and outcomes of different individuals and LLM-based agents toward a common objective. It emphasizes collective commitment, mutual assistance, and the pooling of resources to attain a shared result, thereby fostering a collaborative problem-solving environment. For instance, the human robot coordination in household activities (Chang et al., 2024), the cooperative embodied language agent (CoELA) (Zhang et al., 2024a), human designers collaborat-

ing with an LLM-based agent (Sharma et al., 2024).

Coordination. Coordination is the organized process of aligning and synchronizing the actions of humans and LLM-based agents to achieve a shared objective. Unlike cooperation, the key idea behind coordination is to avoid conflict and bias in both humans and LLM-based agents to reach the final goal. It involves clear communication, strategic planning, and the intentional division of tasks, ensuring that individual efforts are harmonized and effectively integrated to support common goals. For example, humans and agents work in a shared workspace to complete interdependent tasks (Zhang et al., 2024b), human-agent integration supports adaptive decision-making (Sun et al., 2024b), and the collaborative framework facilitates coordination between humans and agents (Pan et al., 2024a).

3.3.2 Competition

Competition is a form of interaction where participants aim to achieve their own goals, which often conflict with the objectives of others. In the LLM-HAS, competition emerges when agents or humans seek to enhance their personal performance or obtain resources, even if it negatively impacts collective results. In addition, competition also necessitates effective balancing mechanisms, like performance regulation or conflict resolution strategies, to prevent unproductive behaviors and ensure that the overall goals of the system remain intact. For instance, the SOTOPIA framework simulates social behaviors between humans and LLM-based agents (Zhou et al., 2024).

3.3.3 Coopetition

Coopetition is an interaction where cooperation and competition coexist at the same time. Within this interaction, participants collaborate on shared tasks or mutual goals while also seeking to outdo each other to improve their own performance or gain extra advantages. In terms of the LLM-HAS, this dual aspect implies that agents and human may join forces to address complex issues while competing in specific domains such as efficiency or precision. This approach not only combines the strengths of both collaboration and competition, but also fosters innovation driven by competitive incentives while also reaping the benefits of cooperative synergy. Successfully managing coopetition typically requires mechanisms for building trust and adaptable strategies that reconcile collec-

Orchestration Paradigm	Description
Task Strategy	What order and grouping of tasks do participants perform?
<i>One-by-One</i>	Actors take turns (e.g., human plans → agent executes → human reviews → agent refines).
<i>Simultaneous</i>	Actors work in parallel (e.g., agent streams partial suggestions while human types).
Temporal Synchronization	When and how tightly do actors’ steps need to align in time?
<i>Synchronous</i>	(1) Real-time interaction: Humans and agents communicate simultaneously or in immediate sequence; (2) Immediate response: Participants expect or require prompt feedback. (e.g., live chat session, real-time voice assistant).
<i>Asynchronous</i>	(1) Delayed interaction: Communication occurs without the expectation of immediate responses; (2) Flexible timing: Participants can respond at their convenience. (e.g., email queues, human leaves comments, agent processes offline).

Table 2: Orchestration paradigms in LLM-based human-agent systems encompass two orthogonal dimensions: task strategy, which can be one-by-one or simultaneous, and temporal synchronization, which can be synchronous or asynchronous.

tive advantages with personal aspirations, which is a challenge for the LLM-HAS. For example, humans and agents play the prisoner’s dilemma in the shared workspace (Jiang et al., 2025).

3.4 Orchestration Paradigm

The orchestration paradigm in LLM-HAS refers to *how* tasks and interactions are managed between humans and agents, covering two dimensions in our survey: **Task Strategy** (*ordering*) and **Temporal Synchronization** (*timing*). Table 2 summarizes the two dimensions of the orchestration paradigm.

3.4.1 Task Strategy

In LLM-HAS, the chosen task strategy, defined by the order and grouping of tasks performed by humans and agents, significantly impacts task execution effectiveness and efficiency. These strategies can typically be categorized into *one-by-one* and *simultaneous* paradigms.

One-by-One. The one-by-one strategy requires participants (humans and LLM-based agents) to perform tasks sequentially, taking clearly defined

turns. For example, a human first outlines a plan, the agent then executes the task, the human subsequently reviews the output, and finally, the agent refines its work based on feedback (Liu et al., 2024a; Zhou et al., 2025). Such sequential interaction helps maintain a clear order of execution and reduces the complexity associated with concurrent task management. However, this rigidity may limit overall efficiency and flexibility, especially in dynamic scenarios requiring parallel processing or rapid interaction cycles (Bansal et al., 2024; Guo et al., 2024b).

Simultaneous. Simultaneous strategy describes an interaction pattern in which LLM-based agents and humans respond concurrently in real time. Compared to the one-by-one strategy, the simultaneous approach more closely mirrors real-world conditions encountered in many simulation tasks (Wang et al., 2025d; Zhang et al., 2025). However, this strategy demands sophisticated mechanisms to handle latency mitigation and seamless coordination between participants.

3.4.2 Temporal Synchronization

Temporal synchronization in LLM-HAS refers to the timing and coordination of interactions between humans and agents. It significantly influences system responsiveness, user experience, and task efficiency. It can be broadly categorized into two modes: *synchronous* and *asynchronous*.

Synchronous. Synchronous interaction involves real-time interactions where humans and agents engage simultaneously. Immediate response is expected, facilitating dynamic exchanges. Examples include live chat sessions, real-time voice assistants (e.g., Siri, Alexa), and collaborative decision-making scenarios (Zhang et al., 2024b; Liu et al., 2023b). This mode is advantageous for tasks requiring rapid responses, immediate clarification, or real-time collaboration (Mehta et al., 2024; Han et al., 2025).

Asynchronous. In contrast, asynchronous interaction occurs without the necessity for immediate responses. Participants can engage at their convenience, allowing for flexibility in communication. Examples include email exchanges, message queues, ticket-based support systems, and task assignments where agents process and report outcomes independently (Shao et al., 2024; Zhang

et al., 2025). Asynchronous communication is beneficial for complex issues that require thoughtful analysis or when participants are in different time zones (Sun et al., 2024b,a).

3.5 Communication

The communication layer in LLM-HAS specifies how information flows, defining **communication structure** (*centralized, decentralized, hierarchical*) and **mode** (*conversation, observation signals, or shared message pools*). Due to space constraints, a detailed introduction is provided in Appendix C.

3.6 Human Agency Scale

The five components discussed above collectively characterize *how humans and agents collaborate or interact*. However, they do not directly address a fundamental question: *to what extent should humans be involved in task completion?* Different tasks and application contexts call for varying degrees of human participation, ranging from full automation to essential human involvement. Drawing on recent work that examines worker preferences and technological capabilities across occupational tasks (Shao et al., 2025), we introduce the Human Agency Scale, a system-level framework that quantifies the desired or required level of human involvement in LLM-HAS (Shao et al., 2025). This scale defines five levels based on the degree of human involvement required for effective task completion: *A1: Full Automation*, *A2: Minimal Human Input*, *A3: Equal Partnership*, *A4: Agent-Assisted* and *A5: Human-Driven*. Levels A1–A2 correspond to **automation**, where agent replaces human effort, while A3–A5 represent **augmentation**, where agent enhances human capabilities. A detailed discussion of each level is provided in Appendix D.

4 Implementation Strategies

This section compares major implementation strategies adopted in LLM-based human-agent systems. Specifically, we include three widely-used categories: 1) Prompting-based methods, (2) Supervised Fine-Tuning (SFT)-based methods, and (3) Reinforcement Learning (RL)-based methods. For each category, we summarize representative methods and analyze their strengths and limitations.

Prompting-based collaboration remains the most widely adopted strategy due to its flexibility, easy implementation and minimal training overhead. Recent work demonstrates that carefully

529	structured prompts can elicit sophisticated collabora-	581
530	tive behaviors, such as proactive clarification,	582
531	shared planning, and theory-of-mind reasoning.	583
532	Systems like MToM (Zhang et al., 2024b) and	584
533	Collaborative Gym (Shao et al., 2024) show that	585
534	explicit role, belief, or goal modeling in prompts	586
535	enables agents to anticipate user intent and adapt	587
536	responses accordingly. Interactive benchmarks and	588
537	interfaces, such as RECODE-H (Miao et al., 2025),	589
538	Magentic-UI (Mozannar et al., 2025), and ARIA	590
539	(He et al., 2025), further illustrate how real-time	591
540	human feedback (e.g., critiques, corrections, or	592
541	preferences) can be injected into the agent loop	
542	at inference time to guide task execution and self-	
543	improvement. Analyses of real-world usage, such	
544	as PATHs (Mysore et al., 2025), reveal recurring	
545	human–AI collaboration patterns that prompting	
546	can exploit without modifying model parameters.	
547	However, despite their agility, prompting-based	
548	methods are often brittle: behaviors are sensitive	
549	to prompt design, have limited controllability, and	
550	struggle to accumulate learning across sessions.	
551	Supervised fine-tuning (SFT) addresses these	
552	limitations by converting human interaction traces,	
553	such as edits, revisions, or clarifications, into per-	
554	sistent behavioral improvements. Works like PRE-	
555	LUDE (Gao et al., 2024a) and XtraGPT (Chen	
556	et al., 2025a) demonstrate how user edits can be	
557	treated as supervision signals, allowing agents to	
558	learn latent user preferences or revision strategies	
559	beyond single-turn prompting. Hybrid systems,	
560	such as Ask-before-Plan (Zhang et al., 2024c) and	
561	CollabLLM (Wu et al., 2025a), combine prompt-	
562	ing with SFT to balance adaptability and stability,	
563	enabling agents to proactively ask questions while	
564	grounding behavior in learned collaboration poli-	
565	cies. Compared to prompting, SFT yields more con-	
566	sistent agent behavior and stronger performance on	
567	specific tasks, but it incurs higher data and engineer-	
568	ing costs and remains constrained by the coverage	
569	and bias of collected interaction data.	
570	Reinforcement learning (RL) formulates hu-	
571	man–agent interaction as a sequential decision-	
572	making problem with explicit reward objectives.	
573	Recent RL-based work, such as UserRL (Qian	
574	et al., 2025b), MUA-RL (Zhao et al., 2025),	
575	SWEET-RL (Zhou et al., 2025), and ReHAC (Feng	
576	et al., 2024), optimizes agents for proactive help-	
577	seeking, tool use, and multi-turn coordination un-	
578	der delayed rewards. Interactive environments like	
579	UserBench (Qian et al., 2025a) provide controlled	
580	testbeds for evaluating user-centric policies, mov-	
	ing beyond static benchmarks toward longitudinal	581
	interaction. Compared to prompting and SFT, RL	582
	enables agents to reason over long horizons and	583
	trade off immediate assistance against future user	584
	satisfaction. However, RL approaches face chal-	585
	lenges in reward specification, sample efficiency,	586
	and training stability. Thus, many recent works	587
	(Zhao et al., 2025; Qian et al., 2025b) adopt hybrid	588
	pipelines that bootstrap RL from prompting or SFT,	589
	suggesting that effective human–agent collabora-	590
	tion arises from complementary learning signals	591
	rather than a single paradigm.	592
	5 Applications and Resources	593
	A diverse range of applications, tools, and re-	594
	sources has emerged for LLM-HAS. We elaborate	595
	on the five most frequent application domains in	596
	Appendix F, summarize the corresponding datasets	597
	and benchmarks in Table 4, and provide a detailed	598
	introduction to representative open-source LLM-	599
	HAS frameworks in Appendix G.	600
	6 Challenges and Opportunities	601
	We provide an in-depth analysis of current chal-	602
	lenges and future opportunities for LLM-HAS in	603
	Appendix I. We highlight the challenges posed by	604
	human flexibility and variability, the limitations	605
	of mostly agent-centered work, inadequate eval-	606
	uation methodologies, unresolved safety vulnera-	607
	bilities, and discuss corresponding opportunities	608
	across these dimensions, applications and beyond.	609
	7 Conclusion	610
	This paper presents a comprehensive review of	611
	LLM-based Human-Agent Systems. We introduce	612
	a structured taxonomy covering five core dimen-	613
	sions: environment and profiling, human feed-	614
	back, interaction types, orchestration paradigms,	615
	and communication, and use it to classify and an-	616
	alyze existing research on LLM-HAS. We also sum-	617
	marize representative implementation frameworks,	618
	benchmark datasets, and evaluation metrics to sup-	619
	port reproducibility and comparative analysis. Fi-	620
	nally, we identify key challenges and unresolved	621
	issues in current LLM-HAS research. These issues	622
	remain major obstacles to the development of effec-	623
	tive, adaptive, safe and trustworthy human-agent	624
	systems. We hope this review offers a comprehen-	625
	sive understanding of the LLM-HAS landscape and	626
	serves as a practical guide for future research.	627

628 Limitations

629 Although we strive to include a wide range of rep-
630 resentative works (e.g., ACL, EMNLP, NAACL,
631 EACL, COLM, NeurIPS, ICLR, ICML, etc.), some
632 relevant research may not be included, especially
633 recent preprints or interdisciplinary research in
634 fields such as cognitive science.

635 References

636 Adem Ait, Javier Luis Cánovas Izquierdo, and Jordi
637 Cabot. 2024. Towards modeling human-agent col-
638 laborative workflows: A bpmn extension. *arXiv*
639 *preprint arXiv:2412.05958*.

640 Hao Bai, Yifei Zhou, Mert Cemri, Jiayi Pan, Alane
641 Suhr, Sergey Levine, and Aviral Kumar. 2024. *Di-*
642 *girl: Training in-the-wild device-control agents with*
643 *autonomous reinforcement learning*. In *Advances in*
644 *Neural Information Processing Systems*, volume 37,
645 pages 12461–12495. Curran Associates, Inc.

646 Gagan Bansal, Jennifer Wortman Vaughan, Saleema
647 Amershi, Eric Horvitz, Adam Fournery, Hussein
648 Mozannar, Victor Dibia, and Daniel S Weld. 2024.
649 Challenges in human-agent communication. *arXiv*
650 *preprint arXiv:2412.10380*.

651 Victor Barres, Honghua Dong, Soham Ray, Xujie Si,
652 and Karthik Narasimhan. 2025. *τ 2-bench: Evaluat-*
653 *ing conversational agents in a dual-control environ-*
654 *ment*. *ArXiv*, abs/2506.07982.

655 Uwe M Borghoff, Paolo Bottoni, and Remo Pareschi.
656 2025. Human-artificial interaction in the age of agen-
657 tic ai: a system-theoretical approach. *Frontiers in*
658 *Human Dynamics*, 7:1579166.

659 Matthew Chang, Gunjan Chhablani, Alexander Clegg,
660 Mikael Dallaire Cote, Ruta Desai, Michal Hlavac,
661 Vladimir Karashchuk, Jacob Krantz, Roozbeh Mot-
662 taghi, Priyam Parashar, and 1 others. 2024. Partnr: A
663 benchmark for planning and reasoning in embodied
664 multi-agent tasks. *arXiv preprint arXiv:2411.00081*.

665 Shreyas Chaudhari, Pranjal Aggarwal, Vishvak Mura-
666 hari, Tanmay Rajpurohit, Ashwin Kalyan, Karthik
667 Narasimhan, Ameet Deshpande, and Bruno Castro
668 da Silva. 2024. Rlhf deciphered: A critical analysis
669 of reinforcement learning from human feedback for
670 llms. *arXiv preprint arXiv:2404.08555*.

671 Minghao Chen, Yihang Li, Yanting Yang, Shiyu Yu,
672 Binbin Lin, and Xiaofei He. 2024a. *Automanual:*
673 *Generating instruction manuals by LLM agents via*
674 *interactive environmental learning*. In *The Thirty-*
675 *eighth Annual Conference on Neural Information*
676 *Processing Systems*.

677 Nuo Chen, Andre Lin Huikai, Jiaying Wu, Junyi Hou,
678 Zining Zhang, Qian Wang, Xidong Wang, and Bing-
679 sheng He. 2025a. *Xtragpt: Context-aware and con-*
680 *trollable academic paper revision*.

Ying-Jung Chen, Chi-Sheng Chen, and Ahmad Albar-
qawi. 2025b. Reinforcing clinical decision support
through multi-agent systems and ethical ai govern-
ance. *arXiv preprint arXiv:2504.03699*.

Yixin Chen, Guoxi Zhang, Yaowei Zhang, Hongming
Xu, Peiyuan Zhi, Qing Li, and Siyuan Huang. 2024b.
Synergai: Perception alignment for human-robot col-
laboration. *arXiv preprint arXiv:2409.15684*.

Sanjiban Choudhury and Paloma Sodhi. 2025. *Better*
than your teacher: LLM agents that learn from privi-
leged AI feedback. In *The Thirteenth International*
Conference on Learning Representations.

Michelle Cohn, Mahima Pushkarna, Gbolahan O
Olanubi, Joseph M Moran, Daniel Padgett, Zion
Mengesha, and Courtney Heldreth. 2024. Believing
anthropomorphism: examining the role of anthropo-
morphic cues on trust in large language models. In
Extended Abstracts of the CHI Conference on Human
Factors in Computing Systems, pages 1–15.

Can Cui, Yunsheng Ma, Xu Cao, Wenqian Ye, and Ziran
Wang. 2024. Drive as you speak: Enabling human-
like interaction with large language models in au-
tonomous vehicles. In *Proceedings of the IEEE/CVF*
Winter Conference on Applications of Computer Vi-
sion, pages 902–909.

Chengyuan Deng, Yiqun Duan, Xin Jin, Heng Chang,
Yijun Tian, Han Liu, Yichen Wang, Kuofeng Gao,
Henry Peng Zou, Yiqiao Jin, Yijia Xiao, Sheng-
hao Wu, Zongxing Xie, Weimin Lyu, Sihong He,
Lu Cheng, Haohan Wang, and Jun Zhuang. 2024. *De-*
constructing The Ethics of Large Language Models
from Long-standing Issues to New-emerging Dilem-
mas: A Survey. *arXiv e-prints*, arXiv:2406.05392.

Danny Driess, Fei Xia, Mehdi S. M. Sajjadi, Corey
Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan
Wahid, Jonathan Tompson, Quan Vuong, Tianhe
Yu, Wenlong Huang, Yevgen Chebotar, Pierre Ser-
manet, Daniel Duckworth, Sergey Levine, Vincent
Vanhoucke, Karol Hausman, Marc Toussaint, Klaus
Greff, and 3 others. 2023. Palm-e: an embodied
multimodal language model. In *Proceedings of the*
40th International Conference on Machine Learning,
ICML’23. JMLR.org.

Jiangshu Du, Yibo Wang, Wenting Zhao, Zhongfen
Deng, Shuaiqi Liu, Renze Lou, Henry Peng Zou,
Pranav Narayanan Venkit, Nan Zhang, Mukund Sri-
nath, Haoran Ranran Zhang, Vipul Gupta, Yinghui Li,
Tao Li, Fei Wang, Qin Liu, Tianlin Liu, Pengzhi Gao,
Congying Xia, and 21 others. 2024. *LLMs assist*
NLP researchers: Critique paper (meta-)reviewing.
In *Proceedings of the 2024 Conference on Empiri-*
cal Methods in Natural Language Processing, pages
5081–5099, Miami, Florida, USA. Association for
Computational Linguistics.

Subhabrata Dutta, Timo Kaufmann, Goran Glavaš, Ivan
Habernal, Kristian Kersting, Frauke Kreuter, Mira
Mezini, Iryna Gurevych, Eyke Hüllermeier, and

738	Hinrich Schuetze. 2024. Problem solving through human-ai preference-based cooperation. <i>arXiv preprint arXiv:2408.07461</i> .	793
739		794
740		795
741	Selin S Everett, Bryan J Bunning, Priyank Jain, Ivan Lopez, Anup Agarwal, Manisha Desai, Robert Gallo, Ethan Goh, Vinay B Kadiyala, Zahir Kanjee, and 1 others. 2025. From tool to teammate: A randomized controlled trial of clinician-ai collaborative workflows for diagnosis. <i>MedRxiv</i> .	796
742		797
743		798
744		799
745		800
746		801
747	Xueyang Feng, Zhi-Yuan Chen, Yujia Qin, Yankai Lin, Xu Chen, Zhiyuan Liu, and Ji-Rong Wen. 2024. Large language model-based human-agent collaboration for complex task solving. In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 1336–1357, Miami, Florida, USA. Association for Computational Linguistics.	802
748		803
749		804
750		805
751		806
752		807
753		808
754	George Fragiadakis, Christos Diou, George Kousiouris, and Mara Nikolaidou. 2024. Evaluating human-ai collaboration: A review and methodological framework. <i>arXiv preprint arXiv:2407.19098</i> .	809
755		810
756		811
757		812
758	Ge Gao, Alexey Taymanov, Eduardo Salinas, Paul Mineiro, and Dipendra Misra. 2024a. Aligning LLM agents by learning latent preference from user edits . In <i>The Thirty-eighth Annual Conference on Neural Information Processing Systems</i> .	813
759		814
760		815
761		816
762		817
763	Jie Gao, Simret Araya Gebreegziabher, Kenny Tsu Wei Choo, Toby Jia-Jun Li, Simon Tangi Perrault, and Thomas W Malone. 2024b. A taxonomy for human-llm interaction modes: An initial exploration. In <i>Extended Abstracts of the CHI Conference on Human Factors in Computing Systems</i> , pages 1–11.	818
764		819
765		820
766		821
767		822
768		823
769	Yiming Gao, Feiyu Liu, Liang Wang, Zhenjie Lian, Dehua Zheng, Weixuan Wang, Wenjin Yang, Siqin Li, Xianliang Wang, Wenhui Chen, and 1 others. 2024c. Enhancing human experience in human-agent collaboration: A human-centered modeling approach based on positive human gain. <i>arXiv preprint arXiv:2401.16444</i> .	824
770		825
771		826
772		827
773		828
774		829
775		830
776	Christos Gkournelos, Christos Konstantinou, and Sotiris Makris. 2024. An llm-based approach for enabling seamless human-robot collaboration in assembly. <i>CIRP Annals</i> , 73(1):9–12.	831
777		832
778		833
779		834
780	Moshe Glickman and Tali Sharot. 2025. How human-ai feedback loops alter human perceptual, emotional and social judgements. <i>Nature Human Behaviour</i> , 9(2):345–359.	835
781		836
782		837
783		838
784	Ran Gong, Qiuyuan Huang, Xiaojian Ma, Hoi Vo, Zane Durante, Yusuke Noda, Zilong Zheng, Song-Chun Zhu, Demetri Terzopoulos, Li Fei-Fei, and 1 others. 2023. Mindagent: Emergent gaming interaction. <i>arXiv preprint arXiv:2309.09971</i> .	839
785		840
786		841
787		842
788		843
789	Diego Gosmar and Deborah A Dahl. 2025. Hallucination mitigation using agentic ai natural language-based frameworks. <i>arXiv preprint arXiv:2501.13946</i> .	844
790		845
791		846
792		847
	Junliang Guo, Yang Ye, Tianyu He, Haoyu Wu, Yushu Jiang, Tim Pearce, and Jiang Bian. 2025. Mineworld: a real-time and open-source interactive world model on minecraft. <i>arXiv preprint arXiv:2504.08388</i> .	848
		849
		850
	Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V Chawla, Olaf Wiest, and Xianliang Zhang. 2024a. Large language model based multi-agents: a survey of progress and challenges. In <i>Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence</i> , pages 8048–8057.	
	Xudong Guo, Kaixuan Huang, Jiale Liu, Wenhui Fan, Natalia Vélez, Qingyun Wu, Huazheng Wang, Thomas L. Griffiths, and Mengdi Wang. 2024b. Embodied LLM agents learn to cooperate in organized teams . In <i>Language Gamification - NeurIPS 2024 Workshop</i> .	
	Hojae Han, Seung-won Hwang, Rajhans Samdani, and Yuxiong He. 2025. Convcodeworld: Benchmarking conversational code generation in reproducible feedback environments. <i>arXiv preprint arXiv:2502.19852</i> .	
	Allyson I Hauptman, Beau G Schelble, Nathan J McNeese, and Kapil Chalil Madathil. 2023. Adapt and overcome: Perceptions of adaptive autonomous agents for human-ai teaming. <i>Computers in Human Behavior</i> , 138:107451.	
	Yufei He, Ruoyu Li, Alex Chen, Yue Liu, Yulin Chen, Yuan Sui, Cheng Chen, Yi Zhu, Luca Luo, Frank Yang, and Bryan Hooi. 2025. Enabling self-improving agents to learn at test time with human-in-the-loop guidance . In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing: Industry Track</i> , pages 1625–1653, Suzhou (China). Association for Computational Linguistics.	
	Sirui Hong, Xiawu Zheng, Jonathan Chen, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, and 1 others. 2023. Metagpt: Meta programming for multi-agent collaborative framework. <i>arXiv preprint arXiv:2308.00352</i> , 3(4):6.	
	Wei-Chieh Huang, Henry Peng Zou, Yaozu Wu, Dongyuan Li, Yankai Chen, Weizhi Zhang, Yangning Li, Angelo Zangari, Jizhou Guo, Chunyu Miao, and 1 others. 2025. Deepresearchguard: Deep research with open-domain evaluation and multi-stage guardrails for safety. <i>arXiv preprint arXiv:2510.10994</i> .	
	Faria Huq, Zora Zhiruo Wang, Frank F. Xu, Tianyue Ou, Shuyan Zhou, Jeffrey P. Bigham, and Graham Neubig. 2025. Cowpilot: A framework for autonomous and human-agent collaborative web navigation . <i>Preprint</i> , arXiv:2501.16609.	
	Guanxuan Jiang, Yuyang Wang, and Pan Hui. 2025. Experimental exploration: Investigating cooperative interaction behavior between humans and large language model agents . <i>Preprint</i> , arXiv:2503.07320.	

851	Yucheng Jiang, Yijia Shao, Dekun Ma, Sina Semnani, and Monica Lam. 2024. Into the unknown unknowns: Engaged human learning through participation in language model agent conversations . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 9917–9955, Miami, Florida, USA. Association for Computational Linguistics.	2020 CHI conference on human factors in computing systems, pages 1–15.	906 907
852			
853			
854			
855			
856			
857			
858			
859	Seth Karten, Mycal Tucker, Huao Li, Siva Kailas, Michael Lewis, and Katia Sycara. 2023. Interpretable learned emergent communication for human-agent teams. <i>IEEE Transactions on Cognitive and Developmental Systems</i> , 15(4):1801–1811.	Jonghan Lim, Sujani Patel, Alex Evans, John Pimley, Yifei Li, and Ilya Kovalenko. 2024. Enhancing human-robot collaborative assembly in manufacturing systems using large language models . In <i>2024 IEEE 20th International Conference on Automation Science and Engineering (CASE)</i> , pages 2581–2587.	908 909 910 911 912 913
860			
861			
862			
863			
864	Michael Kelly, Chelsea Sidrane, Katherine Driggs-Campbell, and Mykel J Kochenderfer. 2019. Hgdagger: Interactive imitation learning with human experts. In <i>2019 International Conference on Robotics and Automation (ICRA)</i> , pages 8077–8083. IEEE.	Haokun Liu, Yaonan Zhu, Kenji Kato, Izumi Kondo, Tadayoshi Aoyama, and Yasuhisa Hasegawa. 2023a. Llm-based human-robot collaboration framework for manipulation tasks. <i>arXiv preprint arXiv:2308.14972</i> .	914 915 916 917 918
865			
866			
867			
868			
869	D Roderick Kiewiet and Mathew D McCubbins. 1991. <i>The logic of delegation</i> . University of Chicago Press.	Haokun Liu, Yaonan Zhu, Kenji Kato, Atsushi Tsukahara, Izumi Kondo, Tadayoshi Aoyama, and Yasuhisa Hasegawa. 2024a. Enhancing the llm-based robot manipulation through human-robot collaboration. <i>IEEE Robotics and Automation Letters</i> .	919 920 921 922 923
870			
871	Been Kim, John Hewitt, Neel Nanda, Noah Fiedel, and Oyvind Tafjord. 2025a. Because we have llms, we can and should pursue agentic interpretability. <i>arXiv preprint arXiv:2506.12152</i> .	Haoyang Liu, Yijiang Li, and Haohan Wang. 2025. Genomas: A multi-agent framework for scientific discovery via code-driven gene expression analysis . <i>ArXiv</i> , abs/2507.21035.	924 925 926 927
872			
873			
874			
875	JiWoo Kim, Minsuk Chang, and JinYeong Bak. 2025b. Beyond turn-taking: Introducing text-based overlap into human-llm interactions. <i>arXiv preprint arXiv:2501.18103</i> .	Jijia Liu, Chao Yu, Jiakuan Gao, Yuqing Xie, Qingmin Liao, Yi Wu, and Yu Wang. 2023b. Llm-powered hierarchical language agent for real-time human-ai coordination . <i>ArXiv</i> , abs/2312.15224.	928 929 930 931
876			
877			
878			
879	Christine Lee, David J. Porfirio, Xinyu Jessica Wang, Kevin Zhao, and Bilge Mutlu. 2025a. Veriplan: Integrating formal verification and llms into end-user planning . <i>ArXiv</i> , abs/2502.17898.	Shipeng Liu, FNU Shrutika, Boshen Zhang, Zhehui Huang, and Feifei Qian. 2024b. Effect of adaptive communication support on human-ai collaboration. <i>arXiv preprint arXiv:2412.06808</i> .	932 933 934 935
880			
881			
882			
883	Dong Won Lee, Yubin Kim, Denison Guvenoz, Sooyeon Jeong, Parker Malachowsky, Louis-Philippe Morency, Cynthia Breazeal, and Hae Won Park. 2025b. The human robot social interaction (hsri) dataset: Benchmarking foundational models’ social reasoning. <i>arXiv preprint arXiv:2504.13898</i> .	Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, and 1 others. 2024c. Agentbench: Evaluating llms as agents. In <i>The Twelfth International Conference on Learning Representations</i> .	936 937 938 939 940
884			
885			
886			
887			
888			
889	Ido Levy, Ben Wiesel, Sami Marreed, Alon Oved, Avi Yaeli, and Segev Shlomov. 2024. Stwebagentbench: A benchmark for evaluating safety and trustworthiness in web agents. <i>arXiv preprint arXiv:2410.06703</i> .	Carol Loganbill, Emily Hardy, and Ursula Delworth. 1982. Supervision: A conceptual model. <i>The counseling psychologist</i> , 10(1):3–42.	941 942 943
890			
891			
892			
893			
894	Xinyi Li, Sai Wang, Siqi Zeng, Yu Wu, and Yi Yang. 2024a. A survey on llm-based multi-agent systems: workflow, infrastructure, and challenges. <i>Vicinity</i> , 1(1):9.	Bowen Lou, Tian Lu, Raghu Santanam, and Yingjie Zhang. 2025. Unraveling human-ai teaming: A review and outlook. <i>arXiv preprint arXiv:2504.05755</i> .	944 945 946
895			
896			
897			
898	Youquan Li, Miao Zheng, Fan Yang, Guosheng Dong, Bin Cui, Weipeng Chen, Zenan Zhou, and Wentao Zhang. 2024b. Fb-bench: A fine-grained multi-task benchmark for evaluating llms’ responsiveness to human feedback. <i>arXiv preprint arXiv:2410.09412</i> .	Ryan Lowe, Yi I Wu, Aviv Tamar, Jean Harb, OpenAI Pieter Abbeel, and Igor Mordatch. 2017. Multi-agent actor-critic for mixed cooperative-competitive environments. <i>Advances in neural information processing systems</i> , 30.	947 948 949 950 951
899			
900			
901			
902			
903	Q Vera Liao, Daniel Gruen, and Sarah Miller. 2020. Questioning the ai: informing design practices for explainable ai user experiences. In <i>Proceedings of the</i>	Junting Lu, Zhiyang Zhang, Fangkai Yang, Jue Zhang, Lu Wang, Chao Du, Qingwei Lin, Saravan Rajmohan, Dongmei Zhang, and Qi Zhang. 2024. Turn every application into an agent: Towards efficient human-agent-computer interaction with api-first llm-based agents. <i>arXiv preprint arXiv:2409.17140</i> .	952 953 954 955 956 957
904			
905			
		Xing Han Lù, Zdeněk Kasner, and Siva Reddy. 2024. Weblinx: Real-world website navigation with multi-turn dialogue. <i>arXiv preprint arXiv:2402.05930</i> .	958 959 960

961	Junyu Luo, Weizhi Zhang, Ye Yuan, Yusheng Zhao,	<i>on Empirical Methods in Natural Language Process-</i>	1017
962	Junwei Yang, Yiyang Gu, Bohan Wu, Binqi Chen,	<i>ing</i> , pages 16830–16857, Suzhou, China. Association	1018
963	Ziyue Qiao, Qingqing Long, and 1 others. 2025.	for Computational Linguistics.	1019
964	Large language model agent: A survey on method-		
965	ology, applications and challenges. <i>arXiv preprint</i>	Riya Naik, Ashwin Srinivasan, Estrid He, and Swati	1020
966	<i>arXiv:2503.21460</i> .	Agarwal. 2025. An empirical study of the role	1021
		of incompleteness and ambiguity in interactions	1022
967	Qianou Ma, Dora Zhao, Xinran Zhao, Chenglei Si,	with large language models. <i>arXiv preprint</i>	1023
968	Chenyang Yang, Ryan Louie, Ehud Reiter, Diyi	<i>arXiv:2503.17936</i> .	1024
969	Yang, and Tongshuang Wu. 2025. Sphere: An eval-		
970	uation card for human-ai systems. <i>arXiv preprint</i>	Sriraam Natarajan, Saurabh Mathur, Sahil Sidheekh,	1025
971	<i>arXiv:2504.07971</i> .	Wolfgang Stammer, and Kristian Kersting. 2025.	1026
		Human-in-the-loop or ai-in-the-loop? automate or	1027
972	Andrea Madotto, Zhaojiang Lin, Chien-Sheng Wu, and	collaborate? In <i>Proceedings of the AAAI Conference</i>	1028
973	Pascale Fung. 2019. Personalizing dialogue agents	<i>on Artificial Intelligence</i> , volume 39, pages 28594–	1029
974	via meta-learning. In <i>Proceedings of the 57th annual</i>	28600.	1030
975	<i>meeting of the association for computational</i>		
976	<i>linguistics</i> , pages 5454–5459.	Bo Pan, Jiaying Lu, Ke Wang, Li Zheng, Zhen Wen,	1031
		Yingchaojie Feng, Minfeng Zhu, and Wei Chen.	1032
977	Nikhil Mehta, Milagro Teruel, Xin Deng, Sergio	2024a. Agentcoord: Visually exploring coordina-	1033
978	Figueroa Sanz, Ahmed Awadallah, and Julia Kiseleva.	tion strategy for llm-based multi-agent collaboration.	1034
979	2024. Improving grounded language understanding	<i>arXiv preprint arXiv:2404.11943</i> .	1035
980	in a collaborative environment by interacting with		
981	agents through help feedback . In <i>Findings of the As-</i>	Lihang Pan, Yuxuan Li, Chun Yu, and Yuanchun Shi.	1036
982	<i>sociation for Computational Linguistics: EACL 2024</i> ,	2024b. A human-computer collaborative tool for	1037
983	pages 1306–1321, St. Julian’s, Malta. Association	training a single large language model agent into	1038
984	for Computational Linguistics.	a network through few examples. <i>arXiv preprint</i>	1039
		<i>arXiv:2404.15974</i> .	1040
985	Yannick Metz, David Lindner, Raphaël Baur, and	Pat Pataranutaporn, Sheer Karny, Chayapatr Archi-	1041
986	Mennatallah El-Assady. 2024. Mapping out the	waranguprok, Constanze Albrecht, Auren R Liu, and	1042
987	space of human feedback for reinforcement learn-	Pattie Maes. 2025. "my boyfriend is ai": A computa-	1043
988	ing: A conceptual framework. <i>arXiv preprint</i>	tional analysis of human-ai companionship in reddit’s	1044
989	<i>arXiv:2411.11761</i> .	ai community. <i>arXiv preprint arXiv:2509.11391</i> .	1045
990	Chunyu Miao, Henry Peng Zou, Yangning Li, Yankai	Qiyao Peng, Hongtao Liu, Hua Huang, Qing Yang,	1046
991	Chen, Yibo Wang, Fangxin Wang, Yifan Li,	and Minglai Shao. 2025. A survey on llm-powered	1047
992	Wooseong Yang, Bowei He, Xinni Zhang, and 1 oth-	agents for recommender systems. <i>arXiv preprint</i>	1048
993	ers. 2025. Recode-h: A benchmark for research code	<i>arXiv:2502.10050</i> .	1049
994	development with interactive human feedback. <i>arXiv</i>		
995	<i>preprint arXiv:2510.06186</i> .	Cheng Qian, Zuxin Liu, Akshara Prabhakar, Zhiwei Liu,	1050
		Jianguo Zhang, Haolin Chen, Heng Ji, Weiran Yao,	1051
996	Margaret Mitchell, Avijit Ghosh, Alexandra Sasha Luc-	Shelby Heinecke, Silvio Savarese, Caiming Xiong,	1052
997	cioni, and Giada Pistilli. 2025. Fully autonomous	and Huan Wang. 2025a. Userbench: An interac-	1053
998	ai agents should not be developed. <i>arXiv preprint</i>	tive gym environment for user-centric agents . <i>ArXiv</i> ,	1054
999	<i>arXiv:2502.02649</i> .	abs/2507.22034.	1055
1000	Hussein Mozannar, Gagan Bansal, Cheng Tan, Adam	Cheng Qian, Zuxin Liu, Akshara Prabhakar, Jieliu Qiu,	1056
1001	Fourney, Victor Dibia, Jingya Chen, Jack Gerrits,	Zhiwei Liu, Haolin Chen, Shirley Kokane, Heng Ji,	1057
1002	Tyler Payne, Matheus Kunzler Maldaner, Madeleine	Weiran Yao, Shelby Heinecke, Silvio Savarese, Caim-	1058
1003	Grunde-McLaughlin, Eric Zhu, Griffin Bassman, Ja-	ming Xiong, and Huan Wang. 2025b. Userrl: Training	1059
1004	cob Alber, Peter Chang, Ricky Loynd, Friederike	interactive user-centric agent via reinforcement learn-	1060
1005	Niedtner, Ece Kamar, Maya Murad, Rafah Hosn,	ing . <i>ArXiv</i> , abs/2509.19736.	1061
1006	and Saleema Amershi. 2025. Magentic-ui: To-		
1007	wards human-in-the-loop agentic systems . <i>ArXiv</i> ,	Jiahao Qiu, Yinghui He, Xinzhe Juan, Yiming Wang,	1062
1008	abs/2507.22358.	Yuhan Liu, Zixin Yao, Yue Wu, Xun Jiang, Ling	1063
		Yang, and Mengdi Wang. 2025. Emoagent: Assess-	1064
1009	Anirban Mukherjee and Hannah Hanwen Chang. 2025.	ing and safeguarding human-ai interaction for mental	1065
1010	Stochastic, dynamic, fluid autonomy in agentic ai:	health safety. <i>arXiv preprint arXiv:2504.09689</i> .	1066
1011	Implications for authorship, inventorship, and liabil-		
1012	ity. <i>arXiv preprint arXiv:2504.04058</i> .	David G Rand and Martin A Nowak. 2013. Human	1067
		cooperation. <i>Trends in cognitive sciences</i> , 17(8):413–	1068
1013	Sheshera Mysore, Debarati Das, Hancheng Cao, and	425.	1069
1014	Bahareh Sarrafzadeh. 2025. Prototypical human-AI		
1015	collaboration behaviors from LLM-assisted writing		
1016	in the wild . In <i>Proceedings of the 2025 Conference</i>		

1070	Vinay Samuel, Henry Peng Zou, Yue Zhou, Shreyas Chaudhari, Ashwin Kalyan, Tanmay Rajpurohit, Ameet Deshpande, Karthik Narasimhan, and Vishvak Murahari. 2024. Personagym: Evaluating persona agents and llms. <i>arXiv preprint arXiv:2407.18416</i> .	1124
1071		1125
1072		1126
1073		1127
1074		1128
1075	SeungWon Seo, SeongRae Noh, Junhyeok Lee, SooBin Lim, Won Hee Lee, and HyeongYeop Kang. 2025. Reveca: Adaptive planning and trajectory-based validation in cooperative language agents using information relevance and relative proximity. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 39, pages 23295–23303.	1129
1076		1130
1077		1131
1078		1132
1079		1133
1080		1134
1081		1135
1082	Kathrin Seßler, Arne Bewersdorff, Claudia Nerdel, and Enkelejda Kasneci. 2025. Towards adaptive feedback with ai: Comparing the feedback quality of llms and teachers on experimentation protocols. <i>arXiv preprint arXiv:2502.12842</i> .	1136
1083		1137
1084		1138
1085		1139
1086		1140
1087	Yijia Shao, Vinay Samuel, Yucheng Jiang, John Yang, and Diyi Yang. 2024. Collaborative gym: A framework for enabling and evaluating human-agent collaboration. <i>arXiv preprint arXiv:2412.15701</i> .	1141
1088		1142
1089		1143
1090		1144
1091	Yijia Shao, Humishka Zope, Yucheng Jiang, Jiaxin Pei, David Nguyen, Erik Brynjolfsson, and Diyi Yang. 2025. Future of work with ai agents: Auditing automation and augmentation potential across the us workforce. <i>arXiv preprint arXiv:2506.06576</i> .	1145
1092		1146
1093		1147
1094		1148
1095		1149
1096	Ashish Sharma, Sudha Rao, Chris Brockett, Akanksha Malhotra, Nebojsa Jojic, and Bill Dolan. 2024. Investigating agency of llms in human-ai collaboration tasks. <i>Preprint</i> , arXiv:2305.12815.	1150
1097		1151
1098		1152
1099		1153
1100	Amanpreet Singh, Tushar Jain, and Sainbayar Sukhbaatar. 2018. Learning when to communicate at scale in multiagent cooperative and competitive tasks. <i>arXiv preprint arXiv:1812.09755</i> .	1154
1101		1155
1102		1156
1103		1157
1104	Jonathan Spencer, Sanjiban Choudhury, Matthew Barnes, Matthew Schmittle, Mung Chiang, Peter Ramadge, and Siddhartha Srinivasa. 2020. Learning from interventions: Human-robot interaction as both explicit and implicit feedback. In <i>16th robotics: science and systems, RSS 2020</i> . MIT Press Journals.	1158
1105		1159
1106		1160
1107		1161
1108		1162
1109		1163
1110	Haotian Sun, Yuchen Zhuang, Ling kai Kong, Bo Dai, and Chao Zhang. 2023. Adaplaner: Adaptive planning from feedback with language models. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	1164
1111		1165
1112		1166
1113		1167
1114		1168
1115	Nan Sun, Bo Mao, Yongchang Li, Lumeng Ma, Di Guo, and Huaping Liu. 2024a. Assistantx: An llm-powered proactive assistant in collaborative human-populated environment. <i>arXiv preprint arXiv:2409.17655</i> .	1169
1116		1170
1117		1171
1118		1172
1119		1173
1120		1174
1121		1175
1122		1176
1123		1177
		1178
		1179
		1180
		1181
		1182
		1183
		1184
		1185
		1186
		1187
		1188
		1189
		1190
		1191
		1192
		1193
		1194
		1195
		1196
		1197
		1198
		1199
		1200
		1201
		1202
		1203
		1204
		1205
		1206
		1207
		1208
		1209
		1210
		1211
		1212
		1213
		1214
		1215
		1216
		1217
		1218
		1219
		1220
		1221
		1222
		1223
		1224
		1225
		1226
		1227
		1228
		1229
		1230
		1231
		1232
		1233
		1234
		1235
		1236
		1237
		1238
		1239
		1240
		1241
		1242
		1243
		1244
		1245
		1246
		1247
		1248
		1249
		1250
		1251
		1252
		1253
		1254
		1255
		1256
		1257
		1258
		1259
		1260
		1261
		1262
		1263
		1264
		1265
		1266
		1267
		1268
		1269
		1270
		1271
		1272
		1273
		1274
		1275
		1276
		1277
		1278
		1279
		1280
		1281
		1282
		1283
		1284
		1285
		1286
		1287
		1288
		1289
		1290
		1291
		1292
		1293
		1294
		1295
		1296
		1297
		1298
		1299
		1300
		1301
		1302
		1303
		1304
		1305
		1306
		1307
		1308
		1309
		1310
		1311
		1312
		1313
		1314
		1315
		1316
		1317
		1318
		1319
		1320
		1321
		1322
		1323
		1324
		1325
		1326
		1327
		1328
		1329
		1330
		1331
		1332
		1333
		1334
		1335
		1336
		1337
		1338
		1339
		1340
		1341
		1342
		1343
		1344
		1345
		1346
		1347
		1348
		1349
		1350
		1351
		1352
		1353
		1354
		1355
		1356
		1357
		1358
		1359
		1360
		1361
		1362
		1363
		1364
		1365
		1366
		1367
		1368
		1369
		1370
		1371
		1372
		1373
		1374
		1375
		1376
		1377
		1378
		1379
		1380
		1381
		1382
		1383
		1384
		1385
		1386
		1387
		1388
		1389
		1390
		1391
		1392
		1393
		1394
		1395
		1396
		1397
		1398
		1399
		1400
		1401
		1402
		1403
		1404
		1405
		1406
		1407
		1408
		1409
		1410
		1411
		1412
		1413
		1414
		1415
		1416
		1417
		1418
		1419
		1420
		1421
		1422
		1423
		1424
		1425
		1426
		1427
		1428
		1429
		1430
		1431
		1432
		1433
		1434
		1435
		1436
		1437
		1438
		1439
		1440
		1441
		1442
		1443
		1444
		1445
		1446
		1447
		1448
		1449
		1450
		1451
		1452
		1453
		1454
		1455
		1456
		1457
		1458
		1459
		1460
		1461
		1462
		1463
		1464
		1465
		1466
		1467
		1468
		1469
		1470
		1471
		1472
		1473
		1474
		1475
		1476
		1477
		1478
		1479
		1480
		1481
		1482
		1483
		1484
		1485
		1486
		1487
		1488
		1489
		1490
		1491
		1492
		1493
		1494
		1495
		1496
		1497
		1498
		1499
		1500

1180	Qian Wang, Biao Zhang, Michael Birsak, and Peter Wonka. 2023. Instructedit: Improving automatic masks for diffusion-based image editing with user instructions. <i>arXiv preprint arXiv:2305.18047</i> .	Zhenyu Wu, Ziwei Wang, Xiuwei Xu, Jiwen Lu, and Haibin Yan. 2023. Embodied task planning with large language models. <i>arXiv preprint arXiv:2307.01848</i> .	1234
1181			1235
1182			1236
1183			1237
1184	Wenxuan Wang, Zizhan Ma, Zheng Wang, Chenghan Wu, Wenting Chen, Xiang Li, and Yixuan Yuan. 2025b. A survey of llm-based agents in medicine: How far are we from baymax? <i>arXiv preprint arXiv:2502.11211</i> .	Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwu Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, and 1 others. 2025. The rise and potential of large language model based agents: A survey. <i>Science China Information Sciences</i> , 68(2):121101.	1238
1185			1239
1186			1240
1187			1241
1188			1242
1189	Xingyao Wang, Zihan Wang, Jiateng Liu, Yangyi Chen, Lifan Yuan, Hao Peng, and Heng Ji. 2024b. MINT: Evaluating LLMs in multi-turn interaction with tools and language feedback. In <i>The Twelfth International Conference on Learning Representations</i> .	Hengjia Xiao and Peng Wang. 2023. Llm a*: Human in the loop large language models enabled a* search for robotics. <i>arXiv preprint arXiv:2312.01797</i> .	1244
1190			1245
1191			1246
1192			
1193			
1194	Xingzhi Wang, Zhoumingju Jiang, Yi Xiong, and Ang Liu. 2025c. Human-llm collaboration in generative design for customization. <i>Journal of Manufacturing Systems</i> , 80:425–435.	Chengxing Xie, Canyu Chen, Feiran Jia, Ziyu Ye, Shiyang Lai, Kai Shu, Jindong Gu, Adel Bibi, Ziniu Hu, David Jurgens, and 1 others. 2024a. Can large language model agents simulate human trust behavior? In <i>The Thirty-eighth Annual Conference on Neural Information Processing Systems</i> .	1247
1195			1248
1196			1249
1197			1250
1198	Ziyan Wang, Meng Fang, Tristan Tomilin, Fei Fang, and Yali Du. 2024c. Safe multi-agent reinforcement learning with natural language constraints. <i>arXiv preprint arXiv:2405.20018</i> .	Jian Xie, Kai Zhang, Jiangjie Chen, Tinghui Zhu, Renze Lou, Yuandong Tian, Yanghua Xiao, and Yu Su. 2024b. Travelplanner: A benchmark for real-world planning with language agents. In <i>ICLR 2024 Workshop on Large Language Model (LLM) Agents</i> .	1251
1199			1252
1200			1253
1201			1254
1202	Ziyan Wang, Zhicheng Zhang, Fei Fang, and Yali Du. 2025d. M3hf: Multi-agent reinforcement learning from multi-phase human feedback of mixed quality. <i>arXiv preprint arXiv:2503.02077</i> .	Congluo Xu, Zhaobin Liu, and Ziyang Li. 2025. Finarena: A human-agent collaboration framework for financial market analysis and forecasting. <i>arXiv preprint arXiv:2503.02692</i> .	1255
1203			1256
1204			1257
1205			
1206	Joseph Weizenbaum. 1966. Eliza—a computer program for the study of natural language communication between man and machine. <i>Communications of the ACM</i> , 9(1):36–45.	Hongshen Xu, Zichen Zhu, Lei Pan, Zihan Wang, Su Zhu, Da Ma, Ruisheng Cao, Lu Chen, and Kai Yu. 2024. Reducing tool hallucination via reliability alignment. <i>arXiv preprint arXiv:2412.04141</i> .	1258
1207			1259
1208			1260
1209			1261
1210	Terry Winograd. 1972. Understanding natural language. <i>Cognitive psychology</i> , 3(1):1–191.	Wei Xu, Marvin J Dainoff, Liezhong Ge, and Zaifeng Gao. 2023. Transitioning to human interaction with ai systems: New challenges and opportunities for hci professionals to enable human-centered ai. <i>International Journal of Human-Computer Interaction</i> , 39(3):494–518.	1262
1211			1263
1212	Shirley Wu, Michel Galley, Baolin Peng, Hao Cheng, Gavin Li, Yao Dou, Weixin Cai, James Zou, Jure Leskovec, and Jianfeng Gao. 2025a. CollabLLM: From passive responders to active collaborators. In <i>Forty-second International Conference on Machine Learning</i> .	Bingyu Yan, Xiaoming Zhang, Litian Zhang, Lian Zhang, Ziyi Zhou, Dezhuang Miao, and Chaozhuo Li. 2025. Beyond self-talk: A communication-centric survey of llm-based multi-agent systems. <i>arXiv preprint arXiv:2502.14321</i> .	1264
1213			1265
1214			1266
1215			1267
1216			1268
1217			1269
1218	Tongshuang Wu, Michael Terry, and Carrie Jun Cai. 2022a. Ai chains: Transparent and controllable human-ai interaction by chaining large language model prompts. In <i>Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems</i> , CHI '22, New York, NY, USA. Association for Computing Machinery.	Lixiang Yan. 2025. From passive tool to socio-cognitive teammate: A conceptual framework for agentic ai in human-ai collaborative learning. <i>arXiv preprint arXiv:2508.14825</i> .	1270
1219			1271
1220			1272
1221			1273
1222			1274
1223			1275
1224			1276
1225	Xingjiao Wu, Luwei Xiao, Yixuan Sun, Junhang Zhang, Tianlong Ma, and Liang He. 2022b. A survey of human-in-the-loop for machine learning. <i>Future Generation Computer Systems</i> , 135:364–381.	Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik R Narasimhan. 2025. τ -bench: A benchmark for <u>T</u> ool- <u>A</u> gent- <u>U</u> ser interaction in real-world domains. In <i>The Thirteenth International Conference on Learning Representations</i> .	1277
1226			1278
1227			1279
1228			1280
1229	Yaozu Wu, Dongyuan Li, Yankai Chen, Renhe Jiang, Henry Peng Zou, Liancheng Fang, Zhen Wang, and Philip S Yu. 2025b. Multi-agent autonomous driving systems with large language models: A survey of recent advances. <i>arXiv preprint arXiv:2502.16804</i> .		1281
1230			1282
1231			1283
1232			1284
1233			1285
			1286

1287	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. 2022. React: Synergizing reasoning and acting in language models. In <i>The eleventh international conference on learning representations</i> .	<i>The Twelfth International Conference on Learning Representations</i> .	1343
1288			1344
1289			
1290			
1291			
1292	Yang Ye, Hengxu You, and Jing Du. 2023. Improved trust in human-robot collaboration with chatgpt. <i>IEEE Access</i> , 11:55748–55754.	Yifei Zhou, Song Jiang, Yuandong Tian, Jason Weston, Sergey Levine, Sainbayar Sukhbaatar, and Xian Li. 2025. Sweet-rl: Training multi-turn llm agents on collaborative reasoning tasks. <i>Preprint</i> , arXiv:2503.15478.	1345
1293			1346
1294			1347
1295	Asaf Yehudai, Lilach Eden, Alan Li, Guy Uziel, Yilun Zhao, Roy Bar-Haim, Arman Cohan, and Michal Shmueli-Scheuer. 2025. Survey on evaluation of llm-based agents. <i>arXiv preprint arXiv:2503.16416</i> .	Henry Peng Zou, Wei-Chieh Huang, Yaozu Wu, Chunyu Miao, Dongyuan Li, Aiwei Liu, Yue Zhou, Yankai Chen, Weizhi Zhang, Yangning Li, and 1 others. 2025. A call for collaborative intelligence: Why human-agent systems should precede ai autonomy. <i>arXiv preprint arXiv:2506.09420</i> .	1348
1296			1349
1297			
1298			
1299	Steve Young, Milica Gašić, Blaise Thomson, and Jason D Williams. 2013. Pomdp-based statistical spoken dialog systems: A review. <i>Proceedings of the IEEE</i> , 101(5):1160–1179.		1350
1300			1351
1301			1352
1302			1353
1303	Yifu Yuan, Jianye HAO, Yi Ma, Zibin Dong, Hebin Liang, Jinyi Liu, Zhixin Feng, Kai Zhao, and YAN ZHENG. 2024. Uni-RLHF: Universal platform and benchmark suite for reinforcement learning with diverse human feedback. In <i>The Twelfth International Conference on Learning Representations</i> .		1354
1304			1355
1305			
1306			
1307			
1308			
1309	Hongxin Zhang, Weihua Du, Jiaming Shan, Qinhong Zhou, Yilun Du, Joshua B. Tenenbaum, Tianmin Shu, and Chuang Gan. 2024a. Building cooperative embodied agents modularly with large language models. In <i>The Twelfth International Conference on Learning Representations</i> .		
1310			
1311			
1312			
1313			
1314			
1315	Shao Zhang, Xihuai Wang, Wenhao Zhang, Yongshan Chen, Landi Gao, Dakuo Wang, Weinan Zhang, Xinbing Wang, and Ying Wen. 2024b. Mutual theory of mind in human-ai collaboration: An empirical study with llm-driven ai agents in a real-time shared workspace task. <i>ArXiv</i> , abs/2409.08811.		
1316			
1317			
1318			
1319			
1320			
1321	Shao Zhang, Xihuai Wang, Wenhao Zhang, Chaoran Li, Junru Song, Tingyu Li, Lin Qiu, Xuezhi Cao, Xunliang Cai, Wen Yao, and 1 others. 2025. Leveraging dual process theory in language agent framework for real-time simultaneous human-ai collaboration. <i>arXiv preprint arXiv:2502.11882</i> .		
1322			
1323			
1324			
1325			
1326			
1327	Xuan Zhang, Yang Deng, Zifeng Ren, See-Kiong Ng, and Tat-Seng Chua. 2024c. Ask-before-plan: Proactive language agents for real-world planning. In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 10836–10863, Miami, Florida, USA. Association for Computational Linguistics.		
1328			
1329			
1330			
1331			
1332			
1333	Weikang Zhao, Xili Wang, Chengdi Ma, Lingbin Kong, Zhaohua Yang, Mingxiang Tuo, Xiaowei Shi, Yitao Zhai, and Xunliang Cai. 2025. Mua-rl: Multi-user-interacting agent reinforcement learning for agentic tool use. <i>ArXiv</i> , abs/2508.18669.		
1334			
1335			
1336			
1337			
1338	Xuhui Zhou, Hao Zhu, Leena Mathur, Ruohong Zhang, Haofei Yu, Zhengyang Qi, Louis-Philippe Morency, Yonatan Bisk, Daniel Fried, Graham Neubig, and Maarten Sap. 2024. SOTOPIA: Interactive evaluation for social intelligence in language agents. In		
1339			
1340			
1341			
1342			

Appendix

Table of Contents

A Motivation: Why LLM-based Human-Agent Systems	16
B Evolution of Human-Agent Collaboration and Interaction Systems	16
C Communication	18
C.1 Communication Structure	18
C.2 Communication Mode	18
D Human Agency Scale	19
D.1 Definition	19
D.2 Human Agency Scale for Design and Analysis	20
E Empirical Distributions of Core Taxonomy Dimensions	20
F Applications	21
G Implementation Tools and Resources	23
G.1 Human-Agent Framework	23
G.2 Datasets and Benchmarks	23
H Evaluation Metrics	23
I Challenges and Opportunities	24
J Ethical and Societal Issues	25
K Human Feedback Type and Subtype	26
L Difference with Traditional Human-in-the-Loop and Human-Computer Interaction Systems	26
M Difference with Multi-Agent Systems	28
N Tables	28

A Motivation: Why LLM-based Human-Agent Systems

Despite the rapid development of fully autonomous agents based on LLMs in recent years, such systems face persistent challenges in terms of reliability, complexity, and safety and ethical risks. Autonomous agents frequently generate incorrect or misleading information and often lack a true

understanding of human goals or contextual nuances. These limitations suggest that full autonomy may not be suitable for many real-world applications (Mitchell et al., 2025; Natarajan et al., 2025; Zou et al., 2025) and highlight a critical but often overlooked insight: the indispensable role of human involvement.

Humans can provide important complementary capabilities such as disambiguation, domain-specific knowledge, feedback, corrections, and high-level supervision that are difficult for automated systems to replicate (Wang, 2024; Dutta et al., 2024). These factors have led to growing interest in a new class of systems explicitly designed for human-agent collaboration: *LLM-based Human-Agent Systems (LLM-HAS)*. Rather than aiming to replace humans, LLM-HAS frameworks are presupposed with active human involvement, leveraging human expertise, supervision, and guidance to compensate for the limitations of the autonomy agent system.

By integrating human collaborators, LLM-HAS systems become more trustworthy, adaptable, and context-aware. In high-stakes fields such as healthcare or finance, collaborative systems are better able to handle complex and sensitive tasks than standalone agent systems. For example, the success of the "FTT" (Everett et al., 2025) (a hybrid team of human clinicians and agents) in the diagnostic reasoning challenge proves that human-agent collaboration can surpass both humans and agents alone. Human-agent collaboration allows building intelligent systems that are more robust, ethical, and value-aligned than either humans or agents could achieve alone.

B Evolution of Human-Agent Collaboration and Interaction Systems

The evolution of human-agent collaboration and interaction systems has been shaped by a series of major paradigm shifts. Early systems were grounded in rule-based and symbolic AI approaches, relying on predefined rules and handcrafted scripts to simulate human interactions. Iconic systems such as ELIZA (Weizenbaum, 1966) and SHRDLU (Winograd, 1972) enabled basic language interactions, but only within highly constrained environments. These agents acted as deterministic tools, offering little room for flexibility, adaptation, or learning from user behaviors. Interactions were largely one-directional, with humans issuing commands and

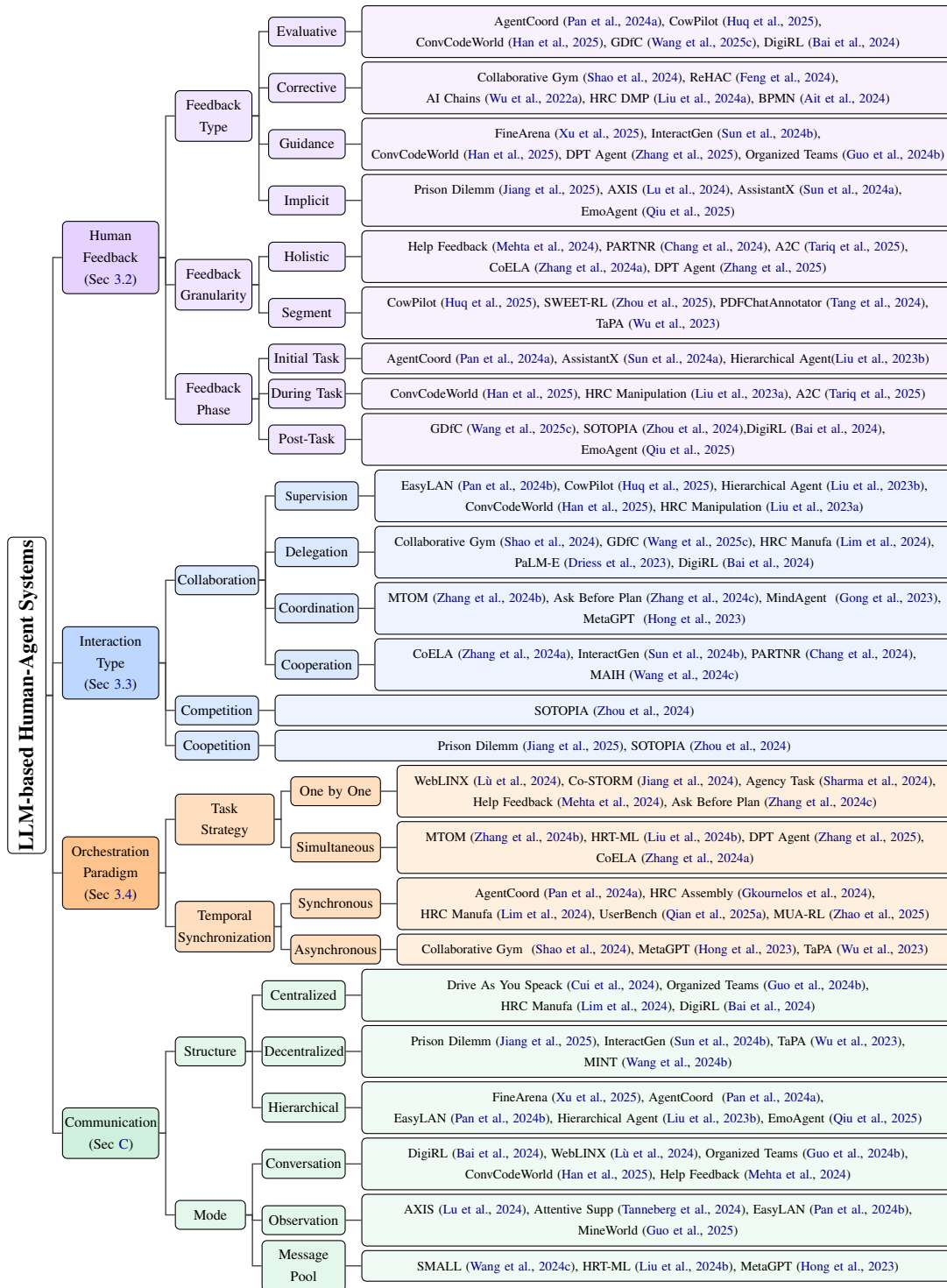


Figure 3: Taxonomy of LLM-based Human-Agent Systems. A more detailed and structured categorization of representative works is provided in the appendix (Table 6 and 7).

agents executing tasks in a fixed, scripted manner.

The second major shift emerged with the rise of machine learning and data-driven NLP in the 2010s. Fueled by large annotated corpora and supervised learning techniques, agents began to exhibit more robust and flexible behavior, particularly in tasks such as speech recognition and dialogue

generation (Young et al., 2013). Commercial systems like Siri and Google Assistant exemplified this transition, allowing broader-domain conversations and more natural interactions. However, these systems remained primarily agent-centric: optimization efforts largely focused on improving task performance, with less attention given to human

1461	adaptivity or user-centered design. Users were often	characterizes the micro-level methods of message	1511
1462	treated as passive input providers, with limited	exchange.	1512
1463	opportunities for active collaboration, personaliza-		
1464	tion, or mutual learning (Madotto et al., 2019; Liao	C.1 Communication Structure	1513
1465	et al., 2020).	Communication structure refers to the organiza-	1514
1466	In recent years, the advent of large language	tional structure of agents, including both humans	1515
1467	models (LLMs) has fundamentally transformed	and agents, in LLM-HAS. It determines how in-	1516
1468	the landscape of human-agent collaboration. Sys-	formation flows at the macro level and shapes the	1517
1469	tems powered by LLMs, such as ChatGPT and	rules of interaction at the micro level. While origi-	1518
1470	Claude, exhibit remarkable abilities in reason-	nally developed for LLM-based multi-agent envi-	1519
1471	ing, co-creation, and open-ended problem solving.	ronments (Guo et al., 2024a), these structures have	1520
1472	These agents move beyond reactive responses to	been effectively adapted to human-agent scenar-	1521
1473	actively engage in collaborative workflows with hu-	ios by treating humans as specialized agents. In	1522
1474	mans, supporting iterative refinement, clarification,	such systems, the communication structure not only	1523
1475	and joint decision-making (Lou et al., 2025; Yan,	governs the efficiency of information exchange but	1524
1476	2025). This new paradigm emphasizes human-AI	also significantly impacts the system’s adaptability,	1525
1477	collaboration as a core design principle, spurring	scalability, and robustness to human variability. We	1526
1478	a growing focus on human-centered design, per-	classify the representative structures into three	1527
1479	sonalized adaptation, and interactive learning. As	types: Centralized , Decentralized , and Hierar-	1528
1480	research shifts from optimizing isolated agent ca-	archical .	1529
1481	pabilities toward co-optimizing human-agent sys-	In Centralized structure, one primary agent or a	1530
1482	tems as integrated, adaptive teams, human adaptiv-	group of core agents acts as a central node to coor-	1531
1483	ity, transparency, and control are increasingly pri-	ordinate all communications within the system. This	1532
1484	oritized as central components of effective collabora-	central agent manages interactions among other	1533
1485	tion (Qian et al., 2025a; Sun et al., 2025b).	agents, simplifying coordination and minimizing	1534
1486		conflicts (Cui et al., 2024). Decentralized structure	1535
1487	C Communication	employs peer-to-peer communication, enabling di-	1536
1488	In LLM-HAS, communication serves as the fun-	rect interactions among agents without centralized	1537
1489	damental mechanism defining the transmission,	control. Agents autonomously manage their com-	1538
1490	reception, and transformation of information be-	munications based on systemic information, en-	1539
1491	tween humans and LLM-based agents. It focuses	hancing system flexibility, adaptability, and robust-	1540
1492	specifically on how <i>information flows</i> across partic-	ness (Shao et al., 2024; Driess et al., 2023). In ad-	1541
1493	ipants to support effective interaction and mutual	dition, Hierarchical structure organizes agents into	1542
1494	understanding. Unlike LLM-based multi-agent sys-	clearly defined levels, assigning distinct roles and	1543
1495	tems (Yan et al., 2025), human-agent systems in-	responsibilities according to their position within	1544
1496	troduce a unique dimension (i.e., flexible, and cog-	the hierarchy (Liu et al., 2023b; Pan et al., 2024b).	1545
1497	nitively diverse human participation). This leads	High-level agents typically fulfill managerial or	1546
1498	to a broader and more complex communication	strategic roles, providing overarching guidance and	1547
1499	landscape, encompassing both human-to-agent and	supervision, while lower-level agents perform spe-	1548
1500	agent-to-agent exchanges, each influenced by hu-	cialized tasks and execute detailed operations.	1549
1501	man interpretability, feedback style, and interaction		
1502	latency.	C.2 Communication Mode	1550
1503	To systematically analyze communication behav-	Communication mode defines the manner through	1551
1504	ior in such systems, we propose a two-dimensional	which humans and agents exchange information	1552
1505	taxonomy that captures the communication behav-	within LLM-HAS. Specifically, communication	1553
1506	ior characteristics of humans and agents from	mode describes the methods employed by partic-	1554
1507	macro-structures to micro-interaction rules. Specif-	ipants to transmit, acquire, and utilize information,	1555
1508	ically, we divide this section into the follow-	critically shaping interaction efficiency and the	1556
1509	ing parts: Communication Structure , which de-	overall performance of the system. Broadly, com-	1557
1510	scribes the macro-level organization of informa-	munication modes can be categorized into three	1558
	tion channels, and Communication Mode , which	primary approaches: Conversation , Observation ,	1559
		and Shared Message Pool .	1560

1561
1562
1563
1564
1565
1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612

Conversation. The conversation-based mode is currently the most prevalent and intuitive approach in LLM-HAS, wherein agents and humans directly engage through natural language dialogues. This interaction format typically utilizes conversational interfaces that allow iterative exchanges, questions, clarifications, and dynamic responses, facilitating efficient collaboration and mutual understanding (Shao et al., 2024). For instance, conversational LLM agents can assist users by answering queries, explaining complex concepts, or collaboratively solving reasoning tasks through iterative dialogues (Wang et al., 2024b). While intuitive and flexible, conversational interactions rely significantly on the communicative clarity and dialogue management capabilities of LLM agents.

Observation. In the observation-based communication mode, agents acquire information implicitly by observing participants behaviors, decisions, or interactions within their environment, rather than through explicit verbal communication. This mode leverages indirect signals, including user actions, feedback cues, or behavioral traces, to infer intentions, preferences, or states (Lu et al., 2024). For example, an LLM-driven tutoring system may adaptively provide targeted instructions by continuously observing student problem-solving behaviors without explicit verbal queries (Pan et al., 2024b). However, relying solely on observational signals can introduce ambiguity, potentially impacting inference accuracy unless complemented by robust inferential mechanisms.

Message Pool. The shared message pool mode involves agents and humans exchanging information through a common information repository. Participants publish messages or data into a message pool, subscribing and retrieving relevant messages based on specific interests or tasks (Sun et al., 2024a). This approach significantly simplifies direct agent-to-agent or human-to-agent interactions, reduces communication complexity, and enhances information management efficiency. A prominent example includes the MetaGPT framework (Hong et al., 2023), where LLM-based agents collaboratively retrieve information dynamically from a shared message pool, streamlining cooperation and information dissemination. Despite these advantages, shared message pools must carefully manage access control to avoid information conflicts or in-

efficient retrieval.

1613
1614
1615
1616
1617
1618
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660

D Human Agency Scale

The development of LLM-based human-agent systems raises a fundamental question: *how much human involvement is appropriate for a given task?* Traditional perspectives on LLM-based Agents often adopt an "Agent-first" view, focusing primarily on the extent to which tasks can be automated. While useful for assessing technological capabilities, such perspectives do not adequately capture the human-centered considerations essential for responsible agent deployment, including human preferences, decision-making authority, and accountability (Zou et al., 2025). Drawing on recent work that examines worker preferences and technological capabilities across occupational tasks (Shao et al., 2025), human agency scale provides a unified framework for quantifying the degree of human involvement required or desired in human-agent collaboration. This framework centers on **human agency** (i.e., the capacity for humans to exercise meaningful control, judgment, and decision-making authority within the system).

Human agency scale serves three key purposes in the context of LLM-HAS:

- **System Design:** Helping designers determine appropriate configurations of interaction types, feedback mechanisms, and orchestration paradigms based on target agency levels.
- **System Analysis:** Providing a unified lens for comparing and categorizing existing LLM-HAS implementations.
- **Responsible Deployment:** Ensuring that human oversight and control are appropriately calibrated to task requirements, particularly in high-stakes domains.

D.1 Definition

The human agency scale defines five levels (A1–A5) based on the degree of human involvement required for effective task completion (as shown in Table 3).

Different human agency scale levels suit different situations depending on various factors. For example, routine, well-structured tasks may be suitable for A1–A2, while open-ended or creative tasks benefit from A3–A5. Tasks requiring specialized domain knowledge or tacit expertise typically demand higher human involvement (A4–A5), as do

Table 3: Human Agency Scale

Level	Name	Description	Agent Role
A1	Full Automation	Agent handles task entirely without human involvement	Automation
A2	Minimal Human Input	Agent needs human input only at key points (e.g., spot-checking or exception handling)	Automation
A3	Equal Partnership	Human and agent collaborate closely, outperforming either alone (e.g., planning/analysis tasks requiring iterative back-and-forth)	Augmentation
A4	Agent-Assisted	Agent requires substantial human input to complete task	Augmentation
A5	Human-Driven	Task fully relies on continuous human involvement	Augmentation

¹ **Automation:** The agent takes primary responsibility for task execution with minimal human oversight.

² **Augmentation:** The human retains meaningful involvement with Agent providing collaborative support.

high-stakes decisions in healthcare, finance, or legal domains where human oversight and accountability are essential.

D.2 Human Agency Scale for Design and Analysis

The Human agency scale framework is not merely a descriptive tool, it also serves as a practical guide for system design. By identifying the appropriate agency level for a given task, designers can make informed decisions about interaction types, feedback mechanisms, and orchestration strategies (Discussed in Section 3). Conversely, observing a system’s component configuration allows researchers to infer its effective human agency level.

Each human agency level implies distinct requirements for system configuration. At lower levels (A1–A2), systems typically employ delegation-based interaction, asynchronous orchestration, and centralized communication with minimal human touchpoints. Additionally, feedback tends to be evaluative and occurs post-task (Yao et al., 2022; Xie et al., 2024b; Liu et al., 2024c). As human agency increases toward A3, cooperation and coordination become the dominant interaction patterns, with feedback shifting to guidance and corrective types during task execution. Systems at this level often balance synchronous and asynchronous modes and adopt decentralized communication to facilitate equal partnership (Shao et al., 2024; Feng et al., 2024; Barres et al., 2025). At higher levels (A4–A5), supervision emerges as the primary interaction type, requiring synchronous orchestration and continuous feedback loops. Communication structures become hierarchical with richer interaction modes such as conversation and observation to support sustained human engagement (Qiu et al., 2025).

E Empirical Distributions of Core Taxonomy Dimensions

In this section, we provide quantitative analysis and empirical distribution of communication structures, orchestration paradigms, and human feedback types across the surveyed literature.

Communication Patterns. The empirical evidence indicates a strong preference for flexible and direct interaction architectures. As shown in Figure 4, Communication Structure is predominantly Decentralized (65.6%, $N = 40$), significantly outpacing Hierarchical (21.3%) and Centralized (14.8%) arrangements. This suggests a trend towards autonomous agent behaviors rather than rigid command-and-control topologies from the conventional setting. This decentralization is mirrored in the Communication Mode, where Conversation is the overwhelming standard, utilized in 88.5% ($N = 54$) of the surveyed papers. This indicates that current LLM-HAS designs focus on natural language dialogue, mimicking human interpersonal interactions, over more systemic approaches like Observation (13.1%) or shared Message Pools (4.9%).

Orchestration Strategies. In terms of Orchestration, the landscape is heavily skewed towards linear, sequential workflows. A substantial 90.2% ($N = 55$) of works employ a One-by-One strategy, whereas only 9.8% attempt Simultaneous execution. This sequential bias is reinforced by the synchronization protocols, with 78.7% ($N = 48$) of systems operating Synchronously. These figures suggest that, despite the potential for parallel processing in LLM agents, current human-agent workflows are designed to align with the linear, turn-taking cognitive limitations of human collaborators, rather than leveraging asynchronous (21.3%)

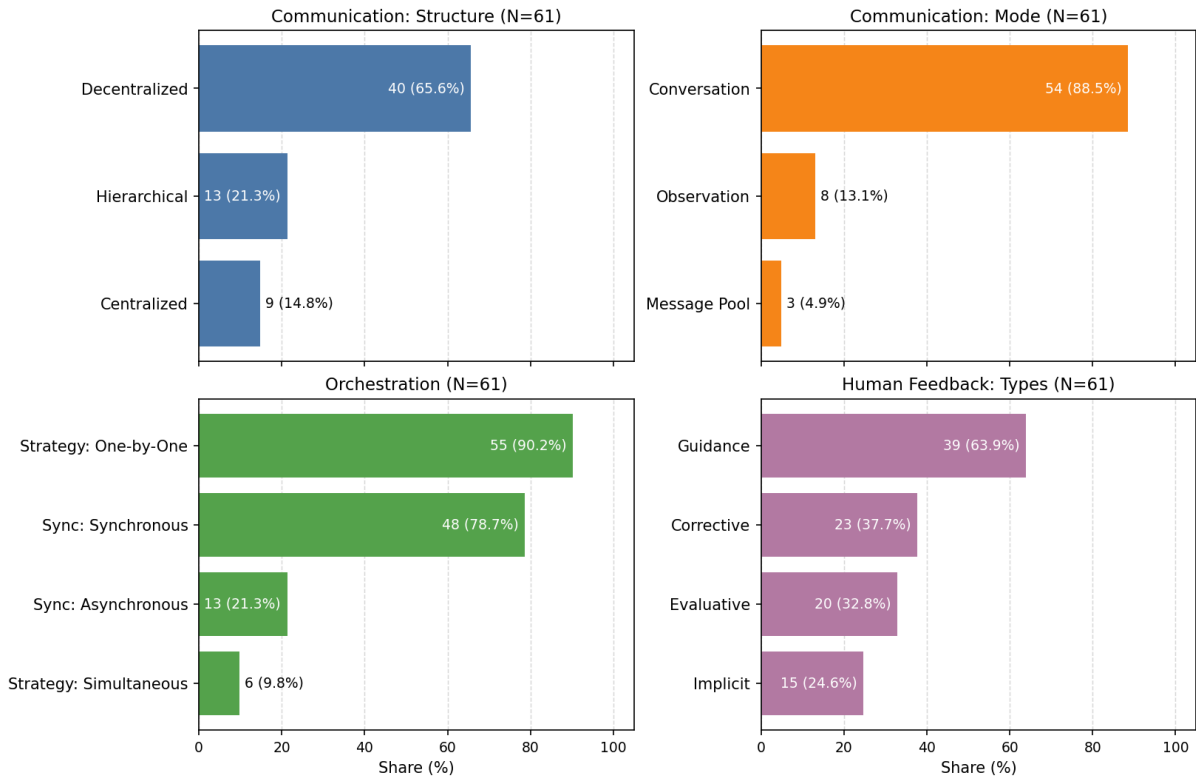


Figure 4: Distributions of core LLM-HAS design dimensions of the collected papers ($N = 61$). Bar labels indicate paper count followed by percentage.

or concurrent operations.

Human Feedback Dynamics. The distribution of Human Feedback highlights that humans are primarily viewed as directors. Guidance is the most prevalent type, appearing in 63.9% ($N = 39$) of works, indicating that humans are frequently involved in proactively steering task execution. In contrast, Corrective (37.7%) and Evaluative (32.8%) feedback are less frequent, suggesting that while error recovery and final grading are essential, the core value proposition of humans in current loops is to provide intermediate direction. Implicit feedback remains the least utilized category at 24.6%, pointing to a significant, under-explored opportunity for systems to learn from passive human signals without demanding explicit user effort.

F Applications

A diverse range of applications has emerged for LLM-HAS. We elaborate on the five most frequent domains below and summarize corresponding datasets and benchmarks in Table 4. With new applications appearing almost weekly in this fast-growing field, we maintain a [GitHub repository](#) to track recent developments.

Embodied AI. Applications in Embodied AI involve various aspects of dynamic and complex real-world tasks, benefiting from valuable human feedback and interactions in LLM-HAS. [Ye et al. \(2023\)](#) explores incorporating LLMs in human-robotic collaboration assembly tasks, allowing seamless communication between robots and humans and increasing trust in human operators. To address the challenges of false planning due to suboptimal environment changes, [Seo et al. \(2025\)](#) proposes REVECA to enable efficient memory management and optimal planning. Additionally, [Tanneberg et al. \(2024\)](#) extends the agents’ collaboration with a group of humans via Attentive Support, enabling agents’ ability to remain silent to not disturb the group if desired.

Software Development. The inherently collaborative nature of software development makes human-agent collaboration vital to improve development efficiency ([Lu et al., 2024](#); [Han et al., 2025](#); [Zhou et al., 2025](#)). [Feng et al. \(2024\)](#) introduces ReHAC framework, wherein agents are trained to determine the optimal stages for human intervention within the problem-solving process,

Domain	Datasets & Benchmarks	Proposed or Used by	Data Link
Embodied AI	TaPA	TaPA (Wu et al., 2023)	Link
	EmboInteract	InteractGen (Sun et al., 2024b)	–
	AssistantX	AssistantX (Sun et al., 2024a)	–
	IGLU Multi-Turn	Help Feedback (Mehta et al., 2024)	Link
	PARTNR	PARTNR (Chang et al., 2024)	Link
	MINT	MINT (Wang et al., 2024b)	Link
	C-WAH	REVECA (Seo et al., 2025)	Link
Conversational Systems	HSRI	HSRI (Lee et al., 2025b)	–
	WEBLINX	WebLINX (Lù et al., 2024)	–
	Ask-before-Plan	Ask Before Plan (Zhang et al., 2024c)	Link
	Agency Dialogue	Agency Task (Sharma et al., 2024)	–
	WildSeek	Co-STORM (Jiang et al., 2024)	Link
	MINT	MINT (Wang et al., 2024b)	Link
	HOTPOTQA	ReHAC (Feng et al., 2024)	Link
Software Development	StrategyQA	ReHAC (Feng et al., 2024)	Link
	MINT	MINT (Wang et al., 2024b)	Link
	InterCode	ReHAC (Feng et al., 2024)	Link
	ColBench	SWEET-RL (Zhou et al., 2025)	Link
	ConvCodeWorld	ConvCodeWorld (Han et al., 2025)	Link
	ConvCodeBench	ConvCodeWorld (Han et al., 2025)	Link
Gaming	RECODE-H	RECODE-H (Miao et al., 2025)	Link
	CuisineWorld	MindAgent (Gong et al., 2023)	Link
Finance	MineWorld	MineWorld (Guo et al., 2025)	Link
	FinArena-Low-Cost	FineArena (Xu et al., 2025)	Link
Healthcare	EmoEval	EmoAgent (Qiu et al., 2025)	Link
	GenoTEX	GenoMAS (Liu et al., 2025)	Link
Retail	τ 2-Bench	τ 2-Bench (Barres et al., 2025)	Link
	τ -Bench	τ -Bench (Yao et al., 2025)	Link
Travel	UserBench	UserBench (Qian et al., 2025a)	Link
	τ 2-Bench	τ 2-Bench (Barres et al., 2025)	Link
	τ -Bench	τ -Bench (Yao et al., 2025)	Link

Table 4: Datasets and Benchmarks across various domains.

offering improved generalizability over the traditional heuristic-based approaches. Building on this direction, Zhou et al. (2025); Han et al. (2025); Wang et al. (2024b) investigate a broader spectrum of human feedback types via multi-turn human-agent interactions. These approaches incorporate carefully designed optimization objectives to effectively capture more diverse and nuanced interactions between humans and agents.

Conversational Systems. In conversational systems, due to the frequent presence of ambiguity and lack of necessary information that agents cannot reliably infer, such as login credentials and payment details, effective human-agent collaboration constitutes a critical component of the system. Zhang et al. (2024c) introduces Proactive Agent Planning, wherein agents are trained to predict classification needs based on the user-agent conversational interactions and current environment, thereby leading to improved reasoning efficacy. Wu et al. (2022a) introduces Chaining the LLM to improve the quality of task outcomes and enhance the transparency and

controllability of the conversational systems.

Gaming. LLM-HAS are naturally well-suited to simulated gaming environments due to their dynamicity and sophistication. Proper human-agent interactions have been shown to enhance humans’ experience, satisfaction and understanding of both the environment and agents (Gong et al., 2023; Gao et al., 2024c). Collaborative interactions also contribute to improved agents’ task performance and decision-making capabilities. For instance, MindAgent framework (Gong et al., 2023) illustrates the efficacy of human-agent collaboration through measurable improvements in task outcomes when humans and agents work together. Mehta et al. (2024) demonstrates agents achieve improved outcomes when interacting with humans via autonomous confusion detection and clarification questions and inquiries. Ait et al. (2024) introduces Meta-Command Communication-based framework to enable effective human-agent collaboration. To address challenges related to execution latency while maintaining strong reasoning capabilities, Liu et al. (2023a) proposes

Hierarchical Language Agent that promotes faster responses, stronger cooperation, and more consistent language communications.

Finance. Given the complexity of stock markets and financial systems, where investors' strategies and risk preferences are critical determinants of outcomes, human-agent collaboration is increasingly recognized as a valuable paradigm. FinArena (Xu et al., 2025) demonstrates the potential of integrating experienced investors with advanced AI agents to support stock prediction tasks. This collaborative framework has been shown to improve investment performance, yielding competitive annualized returns and Sharpe ratios (Xu et al., 2025).

G Implementation Tools and Resources

G.1 Human-Agent Framework

This section provides a detailed introduction to three representative open-source LLM-HAS frameworks: Collaborative Gym (Shao et al., 2024), COWPILOT (Huq et al., 2025), and DPT-Agent (Zhang et al., 2025). They differ in key configuration aspects, including environment settings, interaction types, orchestration paradigms, and communication strategies. Specifically, **Collaborative Gym** (Shao et al., 2024) facilitates asynchronous interactions among humans, agents, and task environments, supporting various simulated and real-world tasks such as travel planning, data analysis, and academic writing. It emphasizes flexible, real-time collaboration and evaluates both outcomes and interaction quality, making it a robust tool for studying human-agent dynamics. **COWPILOT** (Huq et al., 2025) provides a framework for human-agent collaborative web navigation through a Chrome extension. It employs a "Suggest-then-Execute" model under human supervision, allowing dynamic interventions to enhance task completion rates and reduce human workload. It effectively demonstrates how human intervention can significantly improve agent performance. **DPT-Agent** (Zhang et al., 2025) applies Dual Process Theory (DPT) to enable real-time simultaneous human-agent interactions. It features intuitive, fast decision-making and deliberative reasoning components, employing Theory of Mind and asynchronous reflection to manage latency and adapt dynamically to human actions. This approach excels in environments requiring immediate and adaptive responses.

Other frameworks, such as **A2C** (Tariq et al.,

2025), **FinArena** (Xu et al., 2025), and **human-robot collaboration framework** (Liu et al., 2023a), also contribute significantly to specific domains like cybersecurity, financial forecasting, and robotic manipulation, respectively. These frameworks further demonstrate the diverse potential and adaptability of LLM-HAS.

G.2 Datasets and Benchmarks

We summarize the commonly used datasets and benchmarks for Large Language Model-based Human-Agent Systems in Table 4. Diverse domains employ distinct methodologies for evaluating these systems, aligned closely with their unique application contexts. Within the domain of embodied AI, the primary approach involves simulated environments (Sun et al., 2024b,a; Mehta et al., 2024), designed to assess how effectively agents cooperate and execute tasks in dynamic, interactive scenarios. Another significant domain, Conversational Systems, encompasses applications such as question answering (Feng et al., 2024), website navigation (Lù et al., 2024; Levy et al., 2024), design decision assistance (Sharma et al., 2024), and travel planning (Zhang et al., 2024c), adopting benchmarks that evaluate the ability of language models to function as user-aligned conversational assistants, ensuring interactions meet user expectations. Despite the extensive application coverage of current benchmarks, there remains a clear necessity for the development of more comprehensive and standardized benchmarking frameworks.

H Evaluation Metrics

In this section, we introduce evaluation metrics specifically designed for human-agent systems across four key aspects: feedback mechanisms, adaptability, trust and safety, and communication methods. To evaluate feedback mechanisms, (Liu et al., 2024b) assesses a human-robot teaming framework using multi-modal language feedback at varying frequency levels (inactive, passive, active, superactive). (Metz et al., 2024) proposes seven metrics, expressiveness, ease, definiteness, context independence, precision, unbiasedness, and informativeness, to evaluate feedback quality. In the education domain, (Seßler et al., 2025) adopts six dimensions based on educational feedback theory. (Spencer et al., 2020) evaluates the Expert Intervention Learning (EIL) method by comparing robot performance with and without expert intervention.

For adaptability, (Hauptman et al., 2023) examines how human-LLM agents respond to cyber incidents under different levels of autonomy across five NIST-defined phases. For trust and safety, (Levy et al., 2024) introduces a benchmark that evaluates web agents on their ability to comply with policies, avoid unsafe behavior, respect security constraints, and handle errors gracefully, including seeking user input when needed. Finally, (Karten et al., 2023) assess four categories of communication methods in human-agent teaming, focusing on effectiveness and interpretability within simulated environments of Predator-Prey (Lowe et al., 2017) and Traffic Junction (Singh et al., 2018).

In addition to these aspects, AXIS (Lu et al., 2024) and SYNERGAI (Chen et al., 2024b) evaluate the effectiveness and robustness of human-LLM agent systems in the domains of operating systems and embedded AI, respectively. These studies highlight how evaluation criteria can vary significantly depending on the specific task or application context, reflecting differences in system constraints, performance expectations, and interaction complexity.

I Challenges and Opportunities

In this section, we highlight some existing challenges and opportunities for LLM-HAS.

Human Flexibility and Variability. Human feedback varies widely in terms of role, timing, and style across various LLM-HAS. Humans are often subjective, influenced by their personalities, which means different individuals interacting with an LLM-HAS may lead to different outcomes and conclusions. This highlights the need and opportunity for i) thorough investigations or benchmarks on how varied human feedback affects entire systems, and ii) flexible frameworks that can support and adapt to diverse human feedback. In addition, humans, regarded as a "special agent" in the LLM-HAS, are subject to fewer restrictions and evaluations than LLM-based agents. This limits how the LLM-HAS can be improved because the impedance may be on the human side instead of the agent. This concern remains and requires a refined strategy to define the strict, fine interaction rule and evaluation equally for both human and LLM-based agents. Also, many studies today substitute real human participants with LLM simulated human proxies, failing to capture human

input's variety and unpredictability. For example, CollabLLM (Wu et al., 2025a) employs a user simulator to mimic human interaction according to a predefined linguistic style. Nevertheless, the model still relies on fixed prompts to reproduce the requested actions, and its internal knowledge far exceeds that of an average human. As a result, the simulated conversation rarely involves extensive classification and verification steps (Yao et al., 2025), which are often expressed in imprecise language in the human perspective. In contrast, real users frequently produce grammatical errors or struggle to articulate their intentions clearly, behaviors that are rarely observed in LLM agents. The performance gap between humans and the simulated human remains unknown, potentially making the comparison incomparable.

Mostly Agent-Centered Work. In most LLM-HAS studies, guidance flows in a single direction, with humans evaluating agent outputs and providing corrective or evaluative feedback. Namely, the current studies are mostly agent-centered. This agent-centered framework relegates humans to passive evaluators and overlooks the potential for agents to proactively monitor and guide human actions, thereby undercutting bidirectional collaboration (Zhang et al., 2024a). For example, ConvCodeWorld (Han et al., 2025) treats humans as scripted evaluators. Within the framework, LLM-simulated humans provide feedback logs to the agent, rather than empowering agents to observe and coach human coding actions dynamically. However, enabling agents to observe human actions, detect errors or inefficiencies, and offer timely suggestions can transform collaboration and reduce human effort by leveraging agent intelligence. When agents act as instructors by proposing alternative strategies, drawing attention to overlooked risks, and reinforcing effective practices as tasks unfold in real time, both humans and agents benefit. We believe that exploring human-centered LLM-HAS, or shifting toward an equalized LLM-HAS, will unlock the full promise of teamwork between humans and agents.

Inadequate Evaluation Methodologies. In existing evaluation frameworks for LLM-HAS, improvements focus primarily on agent accuracy and static benchmarks, which ignore the real burden placed on human collaborators (Ma et al., 2025). People dedicate varying amounts of time,

2034	attention and cognitive effort depending on the	these systems need clear safeguards around data	2085
2035	type and frequency of feedback they must provide,	sharing, error recovery protocols when agents	2086
2036	yet no standard metric captures this human	behave unpredictably and privacy protections that	2087
2037	workload or its impact on overall efficiency. For	cover every stage of the interaction. Robustness	2088
2038	example, frameworks like CoELA (Jiang et al.,	measures must ensure agents handle ambiguous	2089
2039	2024) evaluate success purely by metrics such as	or adversarial inputs without passing harm on	2090
2040	transport rate improvements, yet ignore entirely	to their human partners (Glickman and Sharot,	2091
2041	the invisible coordination costs and cognitive	2025). Without studies that emphasize human	2092
2042	load on humans (Qiu et al., 2025). Evaluation	experience in safety and privacy design, real-world	2093
2043	methods should measure factors such as time spent	deployments will struggle to gain trust or meet	2094
2044	offering feedback, feedback quality, frequency,	acceptable risk thresholds. Rigorous investigation	2095
2045	and impact (Fragiadakis et al., 2024), perceived	of how safety, robustness and privacy shape human	2096
2046	mental workload and effort required to detect and	agent workflows from design through deployment	2097
2047	correct errors, and they should cover every phase	is essential to build collaborations that are both	2098
2048	of the human agent collaboration from initial	effective and respectful of human needs.	2099
2049	task assignment through post execution review,		2100
2050	to systematically evaluate LLM-HAS. Evaluation	Applications and Beyond. The potential of LLM-	2101
2051	methods should measure factors such as time spent	HAS extends well beyond current applications.	2102
2052	offering feedback, perceived mental workload and	Many opportunities remain to be explored in chal-	2103
2053	effort required to detect and correct errors, and	lenging domains such as healthcare, finance, sci-	2104
2054	they should cover every phase of the human agent	entific research, education, and so on (Luo et al.,	2105
2055	collaboration from initial task assignment through	2025; Guo et al., 2024a). While fully autonomous	2106
2056	post execution review. As human expertise and	LLM-based agent systems encounter difficulties in	2107
2057	LLM-based agent capabilities merge to deliver	handling complex, long-term tasks and earning full	2108
2058	unprecedented performance, both uncertainty and	trust in safety and reliability, the involvement of hu-	2109
2059	variability grow. A new evaluation approach or set	mans to provide additional information, feedback,	2110
2060	of metrics that systematically and comprehensively	and control allows LLM-HAS to greatly improve	2111
2061	quantifies contributions and costs for both humans	overall system performance and safety. This opens	2112
2062	and agents is essential to ensure truly efficient	the door to impactful applications across a broad	2113
2063	collaboration.	range of critical fields.	2114
2064			
2065	Unresolved Safety Vulnerabilities. Most	J Ethical and Societal Issues	2115
2066	LLM-HAS works emphasize improving agent	Although LLM-based human-agent systems have	2116
2067	performance and have left safety, robustness and	demonstrated impressive capabilities in different	2117
2068	privacy underexplored in the context of human	fields, there are still some unresolved social and	2118
2069	interaction (Qiu et al., 2025). As people and LLM-	ethical issues. These problems do not stem purely	2119
2070	based agents collaborate in dynamic workflows, the	from model behavior, but from the process of	2120
2071	risk of misaligned behavior, unexpected failures,	agents interacting with humans and transmitting	2121
2072	or unintended disclosure of sensitive information	information. Over time, agents will subtly influ-	2122
2073	grows. For example, the MetaGPT agent-centered	ence human cognition, emotions, and behavior.	2123
2074	framework (Hong et al., 2023), while integrating	Emotional connection and dependence. One	2124
2075	task decomposition and communication, fails to	point of concern is that LLM-based human-agent	2125
2076	integrate essential safety measures such as input	systems can establish emotional connections with	2126
2077	sanitization, privacy-preserving data handling, and	users, allowing people to have emotional projec-	2127
2078	robust error-containment protocols. The MINT	tions and trust in agents similar to those between	2128
2079	benchmark (Wang et al., 2024b), though it quanti-	people (Cohn et al., 2024). As agents are in-	2129
2080	fies performance gains from multi-turn tool use and	creasingly able to maintain long-term, emotionally	2130
2081	language feedback, omits any analysis of whether	charged interactions, users may begin to anthropo-	2131
2082	these interactive protocols might be exploited for	morphize them or view them as social partners. Re-	2132
2083	code-injection attacks, data exfiltration, or other	cent empirical studies have shown that while users	2133
2084	emergent safety failures. Humans engaging with	report an increase in sense of support and engage-	2134

2135 ment when interacting with artificial intelligence
2136 partners, such relationships may also weaken real-
2137 world social connections and exacerbate loneliness
2138 or emotional dependence, especially among so-
2139 cially isolated people (Pataranutaporn et al., 2025)
2140 . These risks suggest that we need to be wary of
2141 users’ over-reliance and unrealistic expectations on
2142 intelligent agents and balance the boundaries be-
2143 tween human and agent interaction and real social
2144 interaction.

2145 **Responsibility gaps and ambiguous autonomy.**
2146 LLM-driven agents often act with partial autonomy,
2147 planning and executing tasks without full human
2148 oversight and participation. As these systems be-
2149 come more capable, it becomes increasingly diffi-
2150 cult to separate the user’s intent from the agent’s au-
2151 tonomous behavior or to assign blame when things
2152 go wrong (Zou et al., 2025; Mukherjee and Chang,
2153 2025). This problem means that harm may occur
2154 without a clearly identifiable responsible party. If
2155 errors or harmful results occur, it is often difficult
2156 to clearly identify the responsible party. In most
2157 current LLM-HAS architectures, such mechanisms
2158 are still inadequate or missing. Solving this prob-
2159 lem requires systematic efforts in interpretability,
2160 procedural transparency, and governance standards.

2161 **Privacy and data-protection risks.** Because
2162 LLMs’ generative outputs rely on extensive train-
2163 ing corpora and user inputs, they have the poten-
2164 tial to leak private information. Sensitive informa-
2165 tion, including identity numbers or medical records,
2166 may unintentionally be replicated in generated re-
2167 sponses due to the generative nature of these mod-
2168 els, according to a recent survey of LLM-based
2169 agents. When data moves through several mod-
2170 ules, such as the core LLM controller, multi-source
2171 inputs, and long-term memory, privacy risks are
2172 increased in the long and complex agentic work-
2173 flow. Sensitive information may be disseminated to
2174 other users or outside tools as a result of these com-
2175 ponents’ unregulated data flow. Therefore, strong
2176 protections such as strict data usage, safeguards
2177 (Huang et al., 2025), and processing filters are nec-
2178 essary to stop LLM agents from disclosing private
2179 information. These issues highlight the need for
2180 a comprehensive strategy for developing, imple-
2181 menting, and policing LLM-based human-agent
2182 systems. It is possible to ensure that such agents
2183 promote rather than undermine human well-being
2184 by paying close attention to the emotional effects,
2185 establishing explicit accountability structures, and

enforcing strict privacy protections. 2186

2187 **K Human Feedback Type and Subtype**

2188 In this appendix, we present a detailed overview
2189 of human feedback types and their subtypes, as
2190 summarized in Table 5. This table provides con-
2191 cise definitions and illustrates how humans pro-
2192 vide feedback to LLM-based agents in LLM-HAS.
2193 While the main paper introduced the broad cate-
2194 gories of evaluative, corrective, guidance, and im-
2195 plicit feedback, here we expand each category into
2196 more granular subtypes, ranging from scalar ratings
2197 and preference rankings to direct edits, demonstra-
2198 tions, and inferred behavioral signals. Recognizing
2199 these subtypes clarifies the ways in which humans
2200 interact with LLM agents, by offering precise in-
2201 structions and well-defined tasks, to enhance the
2202 accuracy and quality of generated outputs. This
2203 deeper understanding empowers users to optimize
2204 their interactions with LLM-based agents. Addi-
2205 tionally, the systematic breakdown of human feed-
2206 back provides a foundation for cross-study compar-
2207 isons. It underscores the diverse strategies through
2208 which human users can guide, correct, or collab-
2209 orate with LLM-based agents in a more detailed
2210 way.

2211 **L Difference with Traditional** 2212 **Human-in-the-Loop and** 2213 **Human-Computer Interaction Systems**

2214 LLM-based human-agent systems (LLM-HAS) dif-
2215 fer from traditional human-in-the-loop (HITL) sys-
2216 tems and classic human-computer interaction (HCI)
2217 frameworks. They vary in system structure, interac-
2218 tion dynamics, and the way they use feedback. Al-
2219 though all three involve human participation, they
2220 have different ideas about the role of humans, the
2221 independence of intelligent systems, and how col-
2222 laboration works (Wu et al., 2022b; Borghoff et al.,
2223 2025).

2224 **LLM-HAS vs. HITL.** Traditional HITL sys-
2225 tems often include humans at fixed and predictable
2226 stages of the machine learning pipeline, like data
2227 labeling, model selection, or post-correction (Kim
2228 et al., 2025a). Human involvement is usually occa-
2229 sional and specific to tasks, and feedback is mostly
2230 gathered offline in structured formats, such as la-
2231 bels, binary corrections, or rankings. Because of
2232 this, HITL frameworks focus on control, supervi-
2233 sion, and reducing errors but provide little support
2234 for ongoing, interactive, or two-way collaboration

Human Feedback Type	Description	How it Helps Agents
Evaluative Feedback	User provides an assessment of the agent’s output quality.	Signals overall correctness or preference, guiding general alignment.
<i>Preference Ranking</i>	User compares two or more agent outputs and selects the preferred one.	Helps the agent learn relative quality and subjective nuances.
<i>Scalar Rating</i>	User assigns a numerical score (e.g., 1–5) to the agent’s output.	Provides a quantitative measure of satisfaction or quality.
<i>Binary Assessment</i>	User indicates simple correctness (e.g., yes/no, thumbs up/down).	Offers a basic signal of success or failure.
Corrective Feedback	User modifies or directly improves the agent’s output.	Provides explicit examples of desired output, enabling direct learning from errors.
<i>Direct Edits / Refinements</i>	User manually changes the agent’s generated text or code.	Shows the agent the precise correction needed.
Guidance Feedback	User provides instructions or explanations to steer the agent.	Offers deeper context, reasoning, or demonstrations for learning complex behaviors.
<i>Demonstrations</i>	User shows the agent how to perform a task correctly.	Teaches specific procedures or desired interaction patterns.
<i>Instructions / Critiques</i>	User provides natural language explanations, critiques, or step-by-step guidance.	Helps the agent understand why an output is wrong and how to improve.
Implicit Feedback	Agent infers user preference from their behavior.	Reveals preferences and usability issues without explicit feedback requests.
<i>Human Action / Control</i>	Human directly takes actions and control.	Collaborates with humans to effectively finish tasks or learns from human actions.

Table 5: Human Feedback Types and Subtypes. The subtypes of evaluative feedback includes preference ranking, scalar rating, and binary assessment. The subtypes of corrective feedback includes the direct edits or refinement. The subtypes of guidance feedback includes the demonstration and instructions or critiques. The subtypes of implicit feedback include the human action or control.

during task execution. In contrast, LLM-HAS allow continuous, multi-round interaction using natural language. This lets humans guide, critique, refine, or redirect agent behavior as tasks progress. Instead of mainly being supervisors or annotators, humans in LLM-HAS become active collaborators whose input influences both the process and the results of agent actions.

LLM-HAS vs. Traditional HCI Systems. Classic HCI systems are often set up for direct manipulation or command-response interaction, where users specify actions that systems respond to in a fixed way. Even though modern HCI research increasingly focuses on user-centered design and interactive experiences, most HCI systems do not see computational components as independent agents with their own initiative or reasoning abilities (Xu et al., 2023). In contrast, LLM-HAS introduce agentic elements that can create plans, start actions, ask for clarification, and change their behavior based on ongoing interactions. This change moves the interaction from simple tool use to collaboration, allowing for smoother and more flexible

human-agent workflows that go beyond traditional interface-based interaction models.

Feedback and Adaptation. Another major difference is how feedback is represented and used. HITL systems typically depend on infrequent, structured feedback gathered for training or evaluation, while traditional HCI systems often see user feedback as temporary signals that don’t directly affect system behavior during a session. LLM-HAS, on the other hand, can take in rich, natural language feedback in real time. This allows users to express complex intentions, preferences, and judgments as they work on tasks. Supported by large language models, agents in LLM-HAS can learn from minimal input, modify their responses immediately, and change their behavior without needing retraining. This ability for quick adaptation and personalization sets LLM-HAS apart from both HITL and traditional HCI models.

Together, these differences position LLM-HAS as a new type of interactive intelligent system that blends adaptive intelligence with human-centered design. Rather than just enhancing existing HITL

2281 or HCI frameworks, LLM-HAS function under a
2282 different model where humans and agents work
2283 together in reasoning, decision-making, and action
2284 throughout the interaction process.

of representative works, respectively. Both tables
present the representative work. For all the col-
lected work, please refer to the Github page.

2330
2331
2332
2333

2285 **M Difference with Multi-Agent Systems**

2286 While both LLM-HAS and MAS involve collabora-
2287 tion among multiple entities, the key distinction
2288 lies in the nature and role of the collaborating par-
2289 ties (Feng et al., 2024; Shao et al., 2024). Multi-
2290 agent systems are typically composed exclusively
2291 of autonomous agents—each designed to make de-
2292 cisions, communicate, and coordinate tasks with
2293 one another. In these MAS, each agent operates
2294 based on its own set of objectives and algorithms,
2295 and the overall behavior emerges from their inter-
2296 actions (Tran et al., 2025; Guo et al., 2024a).

2297 In contrast, LLM-based human-agent systems
2298 explicitly incorporate humans as active partici-
2299 pants within the decision-making loop (Feng et al.,
2300 2024). Rather than letting the system run purely on
2301 the combined strategies of several LLM-powered
2302 agents, these systems are engineered with mech-
2303 anisms to allow human supervision, intervention,
2304 and feedback (Mehta et al., 2024). This human-
2305 in-the-loop design is critical when balancing the
2306 strengths of LLMs, such as processing vast amounts
2307 of knowledge and performing rapid reasoning,
2308 with the need for contextual, ethical, and domain-
2309 specific judgments that humans uniquely provide
2310 (Vats et al., 2024).

2311 Furthermore, multi-agent systems often assume
2312 that the collaboration among agents can lead to
2313 a form of “collective intelligence” where agents
2314 work toward shared objectives (Sun et al., 2024b).
2315 In many such frameworks, the communication pro-
2316 tocols, coordination strategies, and role dynamics
2317 are all defined among non-human entities. In con-
2318 trast, in human-agent systems, the interaction pro-
2319 tocols are designed to enhance transparency and
2320 provide control for human decision-makers (Shao
2321 et al., 2024). The system can selectively escalate
2322 issues for human review, enable corrective actions
2323 when the automated decision may be off-mark, and
2324 integrate human feedback to iteratively improve the
2325 agent’s performance over time (Mehta et al., 2024).

2326 **N Tables**

2327 Table 6 catalogs the environmental configuration
2328 and human feedback type, and Table 7 categorizes
2329 the interaction, orchestration, and communication

Table 6: ① Environment Configuration and ② Human Feedback to LLM-based agents in human-agent systems. Environment Configuration specifies whether a single or multiple humans collaborate with one or more LLM-based agents, while Human Feedback characterizes the type, subtype, granularity, and interaction phase of the human feedback to the LLM-based agents.

Paper	Venue	Code/ Data	Environment Configuration		Human Feedback to LLM-based Agent			
			Human	LLM Agent	Type	Subtype	Granularity	Phase
Collaborative Gym (Shao et al., 2024)	Arxiv'24	Link	Single	Single	Corrective, Guidance	Refinements, Instructions	Segment	During Task
MTOM (Zhang et al., 2024b)	Arxiv'24	-	Single	Single	Implicit	Human Action	Segment	During Task
FineArena (Xu et al., 2025)	Arxiv'25	-	Single	Multiple	Guidance	Demonstrations	Segment, Holistic	Initial Setup, During Task
Prison Dilemm (Jiang et al., 2025)	Arxiv'25	-	Single	Single	Implicit	Human Action	Segment	During Task
PPP (Sun et al., 2025b)	Arxiv'25	Link	Single	Single	Guidance, Evaluate	Scalar rating, Refinements	Segment, Holistic	During Task
AI Chains (Wu et al., 2022a)	CHI'24	-	Single	Single	Corrective	Refinements	Segment	During Task
Drive As You Speak (Cui et al., 2024)	WACV'24	-	Single	Single	Guidance	Demonstrations	Holistic	Initial Setup
AgentCoord (Pan et al., 2024a)	Arxiv'24	Link	Single	Multiple	Evaluative, Corrective	Preference Ranking, Refinements	Segment, Holistic	Initial Setup, During Task
CowPilot (Huq et al., 2025)	Arxiv'25	Link	Single	Single	Corrective, Evaluative	Binary Assessment, Refinements	Segment	During Task
EasyLAN (Pan et al., 2024b)	Arxiv'24	-	Single	Multiple	Corrective, Guidance	Demonstrations, Refinements	Segment, Holistic	During Task
Hierarchical Agent (Liu et al., 2023b)	AAMAS'24	-	Single	Multiple	Guidance	Demonstrations	Segment	During Task
SWEET-RL (Zhou et al., 2025)	Arxiv'25	Link	Single	Single	Corrective, Implicit	Refinements, Human Action	Segment	Initial Setup, During Task
HRC Assembly (Gkoumelos et al., 2024)	CIRP'24	-	Single	Multiple	Guidance	Demonstrations	Segment	During Task
REVECA (Seo et al., 2025)	Arxiv'24	-	Single	Multiple	Guidance	Demonstrations	Holistic	Initial Setup
AssistantX (Siu et al., 2024a)	Arxiv'24	Link	Multiple	Multiple	Implicit, Guidance	Human Action, Demonstrations	Holistic, Segment	Initial Setup, During Task
MINT (Wang et al., 2024b)	ICLR'24	Link	Multiple	Single	Evaluative, Corrective, Guidance	Binary Assessment, Refinements, Instructions	Holistic	During Task
Help Feedback (Mehta et al., 2024)	EACL'24	-	Single	Single	Evaluative, Guidance	Demonstrations, Instructions, Binary Assessment	Holistic, Segment	During Task
ConvCodeWorld (Han et al., 2025)	ICLR'25	Link	Single	Single	Guidance, Evaluative	Demonstrations, Instructions, Binary Assessment	Segment, Holistic	During Task
ReHAC (Feng et al., 2024)	ACL'24	Link	Single	Single	Corrective	Refinements	Segment	During Task
DPT Agent (Zhang et al., 2025)	Arxiv'25	Link	Single	Single	Guidance	Instructions	Holistic	During Task
HRC Manipulation (Liu et al., 2023a)	IEEE'23	-	Single	Single	Corrective, Guidance	Demonstrations, Refinements	Segment	During Task
HRC DMP (Liu et al., 2024a)	IEEE'24	-	Single	Single	Corrective, Guidance	Refinements, Demonstrations	Segment	During Task
PARTNR (Chang et al., 2024)	ICLR'25	Link	Single	Single	Guidance	Demonstrations	Holistic	Initial Setup
Organized Teams (Guo et al., 2024b)	Arxiv'24	Link	Single	Multiple	Guidance	Demonstrations	Holistic, Segment	Initial Setup, During Task
CoELA (Zhang et al., 2024a)	ICLR'23	-	Single	Multiple	Guidance	Demonstrations	Holistic, Segment	Initial Setup, During Task
Agency Task (Sharma et al., 2024)	EACL'24	Link	Single	Single	Guidance	Demonstrations	Segment	During Task
GDIC (Wang et al., 2025c)	SME'25	-	Single	Multiple	Guidance, Evaluative	Demonstrations, Binary Assessment, Preference Ranking	Holistic, Segment	Initial Setup, During Task, Post Task
PDFChatAnnotator (Tang et al., 2024)	IUI'24	-	Single	Single	Corrective, Guidance	Demonstrations, Refinements	Segment	During Task
Attentive Supp. (Tanneberg et al., 2024)	IEEE'24	Link	Multiple	Single	Implicit, Guidance	Demonstrations, Human Action	Segment	During Task
HRC Trust (Ye et al., 2023)	IEEE'23	-	Single	Single	Guidance	Demonstrations, Instructions	Segment	During Task
BPMN (Ait et al., 2024)	Arxiv'24	Link	Multiple	Multiple	Guidance, Corrective	Instructions, Refinements	Segment	During Task, Post Task
Co-STORM (Jiang et al., 2024)	EMNLP'24	Link	Single	Multiple	Guidance	Demonstrations	Segment	During Task
HRC Manufa. (Lim et al., 2024)	IEEE'24	-	Single	Single	Corrective, Guidance	Demonstrations, Refinements, Instructions	Segment	Initial Setup, During Task
A2C (Tariq et al., 2025)	Arxiv'24	Link	Multiple	Multiple	Guidance, Evaluative	Binary Assessment, Instructions	Holistic, Segment	During Task
MindAgent (Gong et al., 2023)	NAACL'24	Link	Single	Multiple	Guidance	Demonstrations	Segment	During Task
Ask Before Plan (Zhang et al., 2024c)	EMNLP'24	Link	Single	Multiple	Guidance	Demonstrations	Segment	Initial Setup, During Task
SOTOPIA (Zhou et al., 2024)	ICLR'24	-	Multiple	Multiple	Evaluative, Implicit	Scaler Rating, Human Action	Holistic, Segment	During Task, Post Task
PaLM-E (Driess et al., 2023)	ICML'23	Link	Single	Single	Guidance, Implicit	Demonstrations, Human Action	Holistic, Segment	Initial Setup, During Task
TaPA (Wu et al., 2023)	Arxiv'23	Link	Single	Single	Guidance	Demonstrations	Holistic, Segment	Initial Setup
MetaGPT (Hong et al., 2023)	ICLR'24	Link	Single	Multiple	Guidance	Demonstrations	Holistic	Initial Setup
DigiRL (Bai et al., 2024)	NeurIPS'24	Link	Single	Single	Evaluative, Guidance	Binary Assessment, Demonstrations	Holistic	During Task, Post Task
WebLINX (Lù et al., 2024)	Arxiv'24	Link	Single	Multiple	Guidance	Demonstrations	Holistic, Segment	Initial Setup, During Task
MineWorld (Guo et al., 2025)	Arxiv'25	Link	Multiple	Single	Implicit	Human Action	Segment	During Task
M3HF (Wang et al., 2025d)	ICML'25	-	Multiple	Multiple	Evaluative, Guidance	Binary Assessment, Instructions	Segment, Holistic	During Task, Post Task
UserBench (Qian et al., 2025a)	Arxiv'25	Link	Single	Single	Implicit, Guidance	Human Action, Refinement	Segment	Initial Setup, During Task
τ^2 -Bench (Barres et al., 2025)	Arxiv'25	Link	Single	Single	Evaluative, Implicit	Human Action, Binary assessment	Segment, Holistic	Initial Setup, During Task
Magentic-UI (Mozannar et al., 2025)	Arxiv'24	Link	Single	Multiple	Evaluative, Corrective, Guidance, Implicit	Binary Assessment, Refinement, Corrective, Human Action	Segment	During Task, Post Task
RECODE-H (Miao et al., 2025)	Arxiv'25	Link	Single	Single	Guidance, Corrective	Refinements, Corrective, Demonstration	Segment	During Task
EmoAgent (Qiu et al., 2025)	Arxiv'25	-	Single	Multiple	Corrective, Implicit, Guidance	Human Action, Instructions, Binary Assessment	Segment, Holistic	During Task, Post Task
SymbioticRAG (Sun et al., 2025a)	Arxiv'25	-	Single	Single	Corrective, Implicit, Evaluative	Binary Assessment, Refinements, Demonstrations, Instructions, Human Action	Segment	Initial Setup, During Task, Post Task

Table 7: ① Interaction ② Orchestration ③ Communication in LLM-based human-agent systems. Interaction types capture the human and agent collaboration type; Orchestration covers task strategy and temporal synchronization; Communication describes how messages are structured and delivered in the system.

Paper	Venue	Code/ Data	Interaction		Orchestration		Communication	
			Types	Variant	Strategy	Synchronization	Structure	Mode
Collaborative Gym (Shao et al., 2024)	Arxiv'24	Link	Collaboration	Cooperation, Delegation	One-by-One	Asynchronous	Decentralized	Conversation
MTOM (Zhang et al., 2024b)	Arxiv'24	-	Collaboration	Coordination, Cooperation	Simultaneous	Synchronous	Decentralized	Conversation
FineArena (Xu et al., 2025)	Arxiv'25	-	Collaboration	Delegation, Cooperation	One-by-One	Synchronous	Hierarchical	Conversation
Prison Dilemm (Jiang et al., 2025)	Arxiv'25	-	Coopetition	-	One-by-One	Asynchronous	Decentralized	Conversation
PPP (Sun et al., 2025b)	Arxiv'25	Link	Collaboration	Cooperation	One-by-One	Asynchronous	Decentralized	Conversation
AI Chains (Wu et al., 2022a)	CHI'24	-	Collaboration	Cooperation	One-by-One	Synchronous	Decentralized	Conversation
Drive As You Speak (Cui et al., 2024)	WACV'24	-	Collaboration	Delegation	One-by-One	Synchronous	Centralized	Conversation
AgentCoord (Pan et al., 2024a)	Arxiv'24	Link	Collaboration	Coordination	One-by-One	Synchronous	Hierarchical	Conversation
CowPilot (Huq et al., 2025)	Arxiv'25	Link	Collaboration	Supervision, Delegation, Cooperation	One-by-One	Synchronous	Decentralized	Conversation
EasyLAN (Pan et al., 2024b)	Arxiv'24	-	Collaboration	Delegation, Supervision	One-by-One	Synchronous	Hierarchical	Observation
Hierarchical Agent (Liu et al., 2023b)	AAMAS'24	-	Collaboration	Supervision, Delegation, Cooperation	One-by-One	Synchronous	Hierarchical	Conversation
SWEET-RL (Zhou et al., 2025)	Arxiv'25	Link	Collaboration	Delegation	One-by-One	Synchronous	Centralized	Conversation
HRC Assembly (Gkourmelos et al., 2024)	CIRP'24	-	Collaboration	Delegation, Cooperation	One-by-One	Synchronous	Decentralized	Conversation
REVECA (Seo et al., 2025)	Arxiv'24	-	Collaboration	Cooperation	One-by-One	Synchronous	Decentralized	Conversation
AssistantX (Sun et al., 2024a)	Arxiv'24	Link	Collaboration	Delegation, Cooperation	One-by-One	Asynchronous	Decentralized	Message Pool
MINT (Wang et al., 2024b)	ICLR'24	Link	Collaboration	Delegation, Cooperation	One-by-One	Synchronous	Decentralized	Conversation
Help Feedback (Mehta et al., 2024)	EACL'24	-	Collaboration	Supervision, Delegation, Cooperation	One-by-One	Asynchronous	Decentralized	Conversation
ConvCodeWorld (Han et al., 2025)	ICLR'25	Link	Collaboration	Supervision, Delegation	One-by-One	Asynchronous	Decentralized	Conversation
ReHAC (Feng et al., 2024)	ACL'24	Link	Collaboration	Coordination, Supervision	One-by-One	Synchronous	Decentralized	Conversation
DPT Agent (Zhang et al., 2025)	Arxiv'25	Link	Collaboration	Coordination	Simultaneous	Asynchronous	Decentralized	Observation
HRC Manipulation (Liu et al., 2023a)	IEEE'23	-	Collaboration	Supervision, Delegation	One-by-One	Synchronous	Decentralized	Conversation
HRC DMP (Liu et al., 2024a)	IEEE'24	-	Collaboration	Delegation, Supervision	One-by-One	Synchronous	Decentralized	Conversation
PARTNR (Chang et al., 2024)	ICLR'25	Link	Collaboration	Coordination, Cooperation	Simultaneous	Synchronous	Decentralized, Centralized	Observation
Organized Teams (Guo et al., 2024b)	Arxiv'24	Link	Collaboration	Cooperation, Coordination	One-by-One	Synchronous	Decentralized, Centralized, Hierarchical	Conversation
CoELA (Zhang et al., 2024a)	ICLR'23	-	Collaboration	Cooperation, Coordination	Simultaneous	Synchronous	Decentralized	Conversation
Agency Task (Sharma et al., 2024)	EACL'24	Link	Collaboration	Cooperation, Delegation	One-by-One	Synchronous	Decentralized	Conversation
GDIC (Wang et al., 2025c)	SME'25	-	Collaboration	Delegation	One-by-One	Synchronous	Decentralized	Conversation
PDFChatAnnotator (Tang et al., 2024)	IUI'24	-	Collaboration	Delegation	One-by-One	Synchronous	Decentralized	Conversation
Attentive Supp. (Tanneberg et al., 2024)	IEEE'24	Link	Collaboration	Coordination	One-by-One	Synchronous	Decentralized	Observation
HRC Trust (Ye et al., 2023)	IEEE'23	-	Collaboration	Delegation	One-by-One	Synchronous	Decentralized	Conversation
BPMN (Ait et al., 2024)	Arxiv'24	Link	Collaboration	Coordination	Simultaneous	Asynchronous	Decentralized	Message Pool
Co-STORM (Jiang et al., 2024)	EMNLP'24	Link	Collaboration	Coordination	One-by-One	Synchronous	Centralized	Conversation
HRC Manufa. (Lim et al., 2024)	IEEE'24	-	Collaboration	Delegation, Cooperation	One-by-One	Synchronous	Centralized	Conversation
A2C (Tariq et al., 2025)	Arxiv'24	Link	Collaboration	Cooperation	One-by-One	Asynchronous	Hierarchical	Conversation
MindAgent (Gong et al., 2023)	NAACL'24	Link	Collaboration	Coordination	Simultaneous	Synchronous	Centralized	Conversation
Ask Before Plan (Zhang et al., 2024c)	EMNLP'24	Link	Collaboration	Coordination, Delegation	One-by-One	Synchronous	Hierarchical	Conversation
SOTOPIA (Zhou et al., 2024)	ICLR'24	-	Collaboration, Competition	Coordination, Cooperation	One-by-One	Synchronous	Decentralized	Conversation
PaLM-E (Driess et al., 2023)	ICML'23	Link	Collaboration	Delegation	One-by-One	Synchronous	Decentralized	Conversation
TaPA (Wu et al., 2023)	Arxiv'23	Link	Collaboration	Delegation	One-by-One	Asynchronous	Decentralized	Conversation
MetaGPT (Hong et al., 2023)	ICLR'24	Link	Collaboration	Coordination	One-by-One	Asynchronous	Decentralized	Message Pool
DigiRL (Bai et al., 2024)	NeurIPS'24	Link	Collaboration	Delegation	One-by-One	Synchronous	Centralized	Conversation
WebLIXX (Lü et al., 2024)	Arxiv'24	Link	Collaboration	Delegation	One-by-One	Synchronous	Hierarchical	Conversation
MineWorld (Guo et al., 2025)	Arxiv'25	Link	Collaboration	Delegation	One-by-One	Synchronous	Decentralized	Observation
M3HF (Wang et al., 2025d)	ICML'25	-	Collaboration	Cooperation	One-by-One, Simultaneous	Synchronous	Centralized	Message Pool
UserBench (Qian et al., 2025a)	Arxiv'25	Link	Collaboration	Cooperation	One-by-One	Asynchronous	Decentralized	Conversation
τ^2 -Bench (Barres et al., 2025)	Arxiv'25	Link	Collaboration	Cooperation, Coordination	One-by-One, Simultaneous	Synchronous	Decentralized, Hierarchical	Conversation
Magentic-UI (Mozannar et al., 2025)	Arxiv'24	Link	Collaboration	Cooperation, Coordination	Simultaneous	Asynchronous, Synchronous	Hierarchical, Centralized	Conversation, Observation
RECODE-H (Miao et al., 2025)	Arxiv'25	Link	Collaboration	Supervision, Cooperation	One-by-One	Synchronous	Hierarchical	Conversation
EmoAgent (Qiu et al., 2025)	Arxiv'25	-	Collaboration	Supervision, Coordination, Cooperation	One-by-One	Synchronous	Hierarchical, Centralized	Conversation, Observation
SymbioticRAG (Sun et al., 2025a)	Arxiv'25	-	Collaboration	Cooperation, Supervision, Delegation	One-by-One	Synchronous	Centralized	Conversation