# Towards a Design Guideline for RPA Evaluation:
# A Survey of Large Language Model-Based Role-Playing Agents

**Anonymous ACL submission**

## Abstract

Role-Playing Agent (RPA) is an increasingly popular type of LLM Agent that simulates human-like behaviors in a variety of tasks. But how should we evaluate an RPA? It is hard because of the wide variety of task requirements and the different designs of RPA. This paper aims to propose an evidence-based, actionable, and generalizable evaluation design guideline for LLM-based RPA by systematically reviewing $1,676$ papers published between Jan. 2021 and Dec. 2024. Our analysis synthesizes in total six agent attributes, seven task attributes, and seven evaluation metrics from existing literature. From this finding, we propose an RPA evaluation design guideline to support future researchers in designing their own evaluations in a more systematic and consistent manner.

## 1 Introduction

LLMs have yielded human-like performance in various cognitive tasks (e.g., memorization (Schwarzschild et al., 2025), reasoning (Wang et al., 2023a; Plaat et al., 2024), and planning (Song et al., 2023; Huang et al., 2024)). These emergent capabilities enable an increasingly popular research topic on **Role-Playing Agent** (RPA) (Chen et al., 2024d; Tseng et al., 2024): RPAs are digital intelligent agent systems powered by LLMs, where people provide human-like **agent attributes** (e.g., personas) and **task attributes** (e.g., task descriptions) as input, and prompt the LLM to generate human-like behaviors and the reasoning process. The potential of RPAs is promising and far-reaching, as illustrated by the early results of the massive interdisciplinary studies in social science (Park et al., 2022, 2023), network science (Chen et al., 2024b), psychology(Jiang et al., 2024) and juridical science (He et al., 2024b).

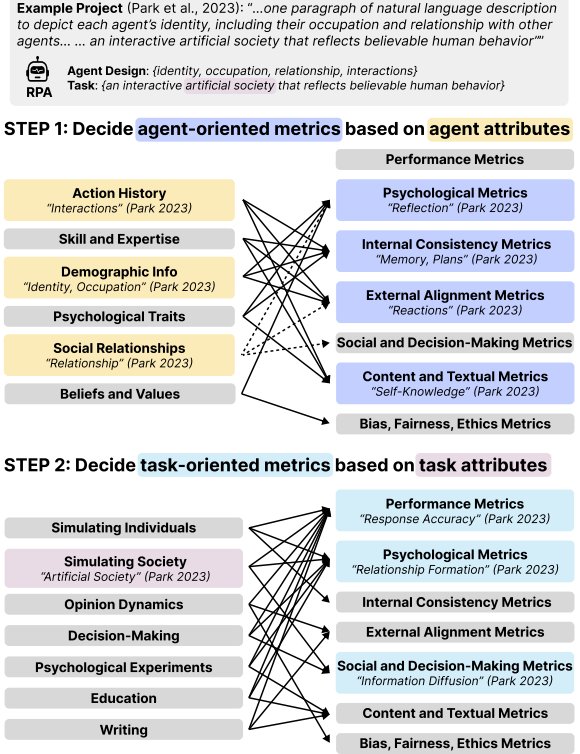Despite the soaring interest, **how can we systematically and consistently evaluate an RPA?** How



Figure 1: RPA evaluation design guideline. To illustrate how to use it in practice, we pretended we were selecting the evaluation metrics for the "Stanford Agent Village" (Park et al., 2023) given agent attributes (yellow) and task attributes (pink). The original authors' selection of evaluation metrics (purple and blue) perfectly aligns with our RPA design guideline, which echoes their work's robustness. More details in Sec 5.1 and a bad example in Sec 5.2.

should we select the evaluation metrics, so that the evaluation results can be comparable or generalizable from one task to another task? It is challenging to find answers to these questions (Dai et al., 2024; Tu et al., 2024; Wang et al., 2024c). One reason is that the variety of the tasks is so broad (e.g., simulating an individual's online browser behavior (Chen et al., 2024b) or simulating a hospital (Li et al., 2024c)) , and the flexibility of RPA design is so high (e.g., an agent persona can be one sentence

or 2-hours of interview log (Park et al., 2024)). Another reason is that researchers often employ arbitrary methods and metrics for the evaluation of their proposed RPAs, which could lead to validity and consistency concerns regarding the evaluation results (Wang et al., 2025; Zhang et al., 2025). As a result, the research community finds it difficult in comparing the performance across multiple RPAs in similar tasks reliably and systematically.

To address this gap, we aim to propose an evidence-based, actionable, and generalizable design guideline for evaluating LLM-based RPAs. To do so, we conducted **a systematic literature review** of 1,676 papers on the LLM Agent topic and identified 122 papers describing its evaluation details. From the expert coding of these papers, we reported that the design of agent attributes interplays with the nature attributes of the downstream tasks (e.g., simulating an individual or simulating a society requires a diverse set of agent attributes). Furthermore, we synthesized common patterns in how prior research successfully (or unsuccessfully) designed their evaluation metrics to correspond to the RPA's agent attributes and task attributes. Based on these common practices, we propose an RPA evaluation design guideline (Fig. 1) and illustrate its generalizability with two case studies.

## 2   Related Work

Existing surveys on the evaluation of RPAs (Gao et al., 2024; Chen et al., 2024d; Tseng et al., 2024; Chen et al., 2024e; Mou et al., 2024a) provide a unified classification of RPA evaluation metrics from the perspective of evaluation approaches. However, they lack a comprehensive and consistent taxonomy for versatile evaluation metrics, leading to arbitrary metrics selection in evaluation practices.

Existing surveys (Gao et al., 2024; Mou et al., 2024a) categorize RPA evaluations into three types: automatic evaluations, human-based evaluations, and LLM-based assessments. Automatic evaluations are efficient and objective, but lack context sensitivity, failing to capture nuances like persona consistency. Human-based evaluations provide deep insight into character alignment and engagement, but they are costly, less scalable, and prone to subjectivity. LLM-based evaluations are automatic and offer scalability and speed, but may not always align with human judgments.

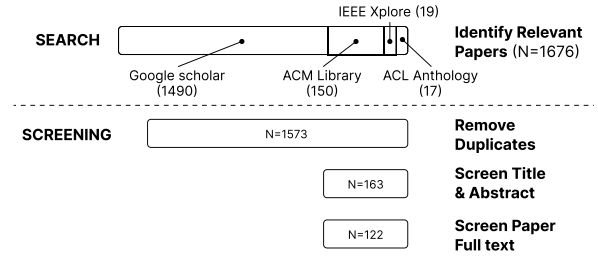The classification of evaluation metrics in prior works varies significantly, leading to inconsistency



Figure 2: Screening process of our literature review. We initially retrieved 1,676 papers published between 2021 and 2024, and systematically narrowed them down to 122 final selections.

and ambiguity. For instance, Gao et al. (2024) focuses on realness validation and ethics evaluation, whereas Chen et al. (2024d) differentiates between character persona and individualized persona. Furthermore, Chen et al. (2024e) classifies evaluation into conversation ability, role-persona consistency, role-behavior consistency, and role-playing attractiveness, which partially overlap with Mou et al. (2024a)'s individual simulation and scenario evaluation. These discrepancies indicate a lack of standardized taxonomy, making it difficult to compare results across studies and select appropriate evaluation metrics for specific applications.

While existing surveys offer different taxonomies of RPA evaluation, they do not provide concrete evaluation design guidelines. Our work aims to bridge this gap by proposing a structured framework that systematically links evaluation metrics to RPA attributes and real-world applications.

## 3   Method

We conduct a systematic literature review to address our research question. Following prior method (Nightingale, 2009), we aim to identify relevant research papers on RPAs and provide a comprehensive summary of the literature. We selected four widely used academic databases: Google Scholar, ACM Digital Library, IEEE Xplore, and ACL Anthology. These databases encompass a broad spectrum of research across AI, human-computer interaction, and computational linguistics. Given the fast-paced nature of LLM research, we did not restrict our selection to peer-reviewed venues, as many impactful studies appear in preprint repositories (e.g., arXiv). Below, we detail our paper selection and annotation process.

2

Table 1: Definition and examples of six agent attributes.

| Agent attributes | Definition | Examples |
|---|---|---|
| Activity History | A record of past actions, behaviors, and engagements, including schedules, browsing history, and lifestyle choices. | Backstory, plot, weekly schedule, browsing history, social media posts, lifestyle |
| Belief and Value | The principles, attitudes, and ideological stances that shape an individual's perspectives and decisions. | Stances, beliefs, attitudes, values, political leaning, religion |
| Demographic Information | Personal identifying details such as name, age, education, career, and location. | Name, appearance, gender, age, date of birth, education, location, career, household income |
| Psychological Traits | Characteristics related to personality, emotions, interests, and cognitive tendencies. | Personality, hobby and interest, emotional |
| Skill and Expertise | The knowledge level, proficiency, and capability in specific domains or technologies. | Knowledge level, technology proficiency, skills |
| Social Relationships | The nature and dynamics of interactions with others, including roles, connections, and communication styles. | Parenting styles, interactions with players |

### 3.1 Literature Search and Screening Method

Our literature review focuses on LLM agents that role-play human behaviors, such as decision-making, reasoning, and deliberating actions. We specifically focus on studies where LLM agents demonstrate the ability to simulate human-like cognitive processes in their objectives, methodologies, or evaluation techniques. To ensure methodological rigor, we defined explicit inclusion and exclusion criteria (Tab. 4 in Appendix A). The inclusion criteria require that an LLM agent in the study exhibits human-like behavior, engages in cognitive activities such as decision-making or reasoning, and operates in an open-ended task environment. We excluded studies where LLM agents primarily serve as chatbots, task-specific assistants, evaluators, or agents operating within predefined and finite action spaces. Additionally, studies focusing solely on perception-based tasks (e.g., computer vision or sensor-based autonomous driving) without cognitive simulation were also excluded.

Following the above survey scope, we searched four databases using the query string provided in Appendix B and initially retrieved $1,676$ papers between January 2021 to December 2024. After removing duplicates, $1,573$ unique papers remained. Two authors independently screened the paper titles and abstracts based on the inclusion criteria. If at least one author deemed a paper relevant, it proceeded to full-text screening, where two authors reviewed the paper in detail and resolved any disagreements through discussion (Fig. 2). The final set of selected studies comprised 122 publications.

### 3.2 Paper Annotation Method

Our team followed established open coding procedures (Brod et al., 2009) and conducted an inductive coding process to identify key themes. Three authors with extensive experience in LLM agents collaboratively annotated the papers, focusing on three dimensions: (1) agent attributes, (2) task attributes, and (3) evaluation metrics.

Two authors independently annotated the same 20% of the sample and then held a meeting to discuss and refine an initial set of categories for the three dimensions. After reaching a consensus, each researcher was responsible for annotating half of the remaining papers. Once the annotations were completed, a third author reviewed the coded data and identified potential discrepancies. Discrepancies were then discussed with the original annotators to ensure consistency until disagreements were resolved. This iterative process helped maintain the reliability and validity of our analysis.

## 4 Survey Findings

Building on the annotated data, we first systematically categorized agent attributes, task attributes, and evaluation metrics. Subsequently, we outline a clear RPA evaluation design guideline for selecting appropriate evaluation metrics based on agent attributes and task attributes.

### 4.1 Agent Attributes

We identified six categories of agent attributes, as shown in Tab. 1. *Activity history* refers to an agent's longitudinal behaviors, such as browsing history (Chen et al., 2024b) or social media activity (Navarro et al., 2024). *Belief and value* encompass the principles, attitudes, and ideological stances that shape an agent's perspectives, including political leanings (Mou et al., 2024c) or religious affiliations (Lv et al., 2024). *Demographic information* includes personal details such as name, age, education, location, career status, and household income. *Psychological traits* include

Table 2: Definition of seven task attributes.

| Task attributes | Definition |
|---|---|
| Simulated Individuals | Simulating specific individuals or groups, such as users and participants. |
| Simulated Society | Simulating social interactions, such as cooperation, competition, and communication. |
| Opinion Dynamics | Simulating political views, legal perspectives, and social media content. |
| Decision Making | Simulating decision-making of stakeholders in investment, public policies, or games. |
| Psychological Experiments | Simulating human traits, including personality, ethics, emotions, and mental health. |
| Educational Training | Simulating teachers and learners to enable personalized teaching and accommodate learner needs. |
| Writing | Simulating readers or characters to support character development and audience understanding. |

Table 3: Definitions and examples of seven evaluation metric categories.

| Evaluation Metrics | Definitions | Examples |
|---|---|---|
| Performance | Assess RPAs' effectiveness in task execution and outcomes. | Prediction accuracy |
| Psychological | Measure human psychological responses to RPAs and the agents' self-awareness and emotional state. | Big Five Invertory |
| External Alignment | Evaluate how closely RPAs align with external ground truth or human behavior and judgments. | Alignment between model and human |
| Internal Consistency | Assess coherence between an RPA's predefined traits (e.g., personality), contextual expectations, and behavior. | Personality-behavior alignment |
| Social and Decision-Making | Analyze RPAs' social interactions and decision-making, including their effects on negotiation, societal welfare, markets, and social dynamics. | Social Conflict Count |
| Content and Textual | Evaluate the quality, coherence, and diversity of RPAs' text, including semantic understanding, linguistic style, and engagement. | Content similarity |
| Bias, Fairness, and Ethics | Assess biases, extreme or unbalanced content, or stereotyping behavior. | Factual error rate |

an agent's personality (Jiang et al., 2023a), emotions, and cognitive tendencies (Castricato et al., 2024). *Skill and expertise* describe an agent's knowledge and proficiency in specific domains, such as technology proficiency or specialized professional skills. Lastly, *social relationships* define the social interactions, roles, and communication styles between agents, including aspects like parenting styles (Ye and Gao, 2024) or relationships between players (Ge et al., 2024).

## 4.2 Task Attributes

For attributes of RPA downstream tasks, we identified seven different types (Tab. 2). Among them, simulated individuals and simulated society primarily use simulation as the ultimate research goal. *Simulated individuals* involve modeling specific individuals or groups, such as end-users (Chen et al., 2024a), to study their behaviors and interactions in a controlled setting. *Simulated Society* focuses on social interactions, including cooperation (Bouzekri et al., 2024), competition (Wu et al., 2024), and communication (Mishra et al., 2023), aiming to explore emergent social dynamics.

In contrast, the other task attributes employ simulation as a means to serve specific research domains. *Opinion dynamics* entails simulating political views (Neuberger et al., 2024), legal perspectives (Chen et al., 2024c), and social media

discourse (Liu et al., 2024c) to analyze the formation and evolution of opinions. *Decision making* addresses the decision-making processes of stakeholders in investment (Sreedhar and Chilton, 2024) and public policy (Ji et al., 2024), providing insights into strategic behaviors. *Psychological experiments* explore human traits such as personality (Bose et al., 2024), ethics (Lei et al., 2024), emotions (), and mental health (De Duro et al., 2025), using simulated scenarios to study cognitive and behavioral responses. *Educational training* supports personalized learning by simulating teachers and learners, enhancing pedagogical approaches and adaptive education systems (Liu et al., 2024d). Finally, *writing* involves modeling readers or characters to facilitate character development (Benharrak et al., 2024) and audience engagement (Choi et al., 2024), contributing to storytelling and content generation research.

## 4.3 Agent- and Task-Oriented Metrics

We derived seven categories of evaluation metrics (Tab. 3) that are shared by both agent- and task-oriented metrics despite the differences in the specific metrics. Agent-oriented metrics focus on intrinsic, task-agnostic properties that define an RPA's essential ability, such as underlying reasoning, consistency, and adaptability. These include *performance* metrics like memorization, *psycho-*
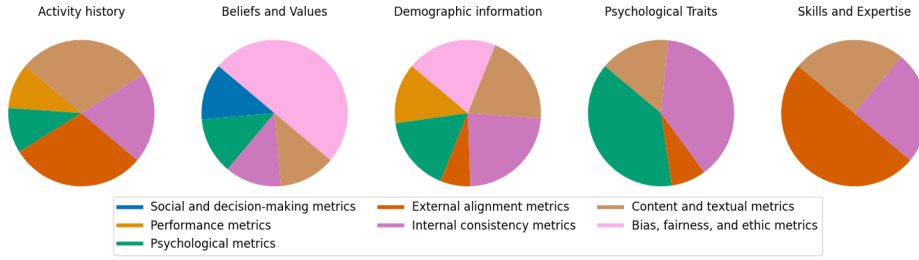
4

Figure 3: Proportional distribution of agent-oriented metrics across different agent attributes.
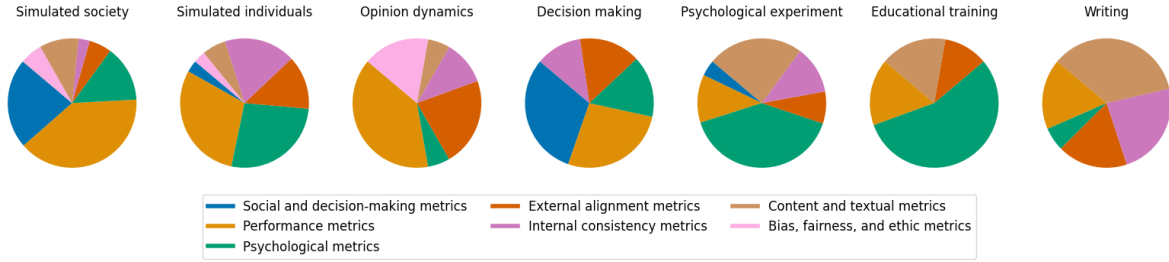


Figure 4: Proportional distribution of task-oriented metrics across different task attributes.

*logical* metrics such as emotional responses measured via entropy of valence and arousal, and *social and decision-making* metrics like social value orientation. Additionally, agent-oriented evaluations emphasize *internal consistency* metrics (e.g., consistency of information across interactions), *external alignment* metrics (e.g., hallucination detection), and *content and textual* metrics such as clarity. These evaluations ensure that RPAs exhibit logical coherence, avoid factual inconsistencies, and align their internal structures with the expected behavioral and cognitive frameworks, independent of any specific downstream task.

In contrast, task-oriented metrics evaluate an RPA's effectiveness in performing specific downstream tasks by assessing various task-related aspects such as accuracy, consistency, social impact, and ethical considerations. *Performance* measures how well RPAs execute designated tasks, such as prediction accuracy. *Psychological* metrics assess human psychological responses to RPAs, including self-awareness and emotional states; for example, the Big Five Inventory. *External alignment* evaluates how closely RPAs align with external ground truth or human behavior; for instance, alignment between model and human. *Internal consistency* ensures coherence between an RPA's predefined traits, contextual expectations, and behavior; for example, personality-behavior alignment. *Social and decision-making* metrics analyze RPAs' influence on negotiation, societal welfare, and social dynamics; for instance, the social conflict count. *Content and textual quality* focuses on the coher-

ence, linguistic style, and engagement of RPAs' generated text, such as content similarity. Lastly, *bias, fairness, and ethics* metrics examine biases, extreme content, or stereotypes; for instance, the factual error rate. Together, these seven metrics provide a comprehensive framework for evaluating RPAs' task performance and broader impact.

## 4.4 RPA Evaluation Design Guideline

Building on our previous classification of agent attributes, task attributes, and evaluation metrics, we observed that both agent design and evaluation can be broadly divided into two categories: agent-related and task-related. This leads us to explore whether there are underlying statistical patterns between agent design and evaluation that could inform the systematic development of design guidelines for evaluation metrics in future research.

**Step 1. Selecting Agent-oriented Metrics Based on Agent Attributes** We analyzed the distribution of agent attributes and agent-oriented metrics, as illustrated in Fig. 3. Our analysis reveals that, for each agent attribute, the top three categories of agent-oriented metrics account for the majority of all metric types. Based on this observation, our first guideline recommends selecting agent-oriented metrics according to agent attributes. Specifically, we suggest referring to Tab. 5 in Appendix D to identify the top three corresponding metrics. For instance, for Activity History, the recommended metrics are external alignment, internal consistency, and content and textual metrics.

5

Likewise, for Beliefs and Values, the most relevant choices are psychological metrics and bias, fairness, and ethics metrics. Notably, no established agent-oriented evaluation metrics exist for social relationships. Based on Social Exchange Theory (Cropanzano and Mitchell, 2005), which explains relationship formation through reciprocal interactions and resource exchanges, we propose assessing social relationships with psychological metrics, external alignment metrics, and social and decision-making metrics.

**Step 2: Selecting Task-Oriented Metrics Based on Task Attributes** Additionally, we analyzed the distribution of task attributes and task-oriented metrics, as shown in Fig. 4. Consistent with our previous findings, we observed that for each category of task attributes, the top three task-oriented metrics account for the vast majority of all metric types. Based on this, our second guideline recommends selecting task-oriented metrics according to task attributes. Specifically, we suggest referring to Tab. 6 in Appendix D to identify the top three corresponding metrics. For instance, for the Simulated Society task, the recommended metrics are social and decision-making, performance, and psychological metrics. Similarly, for the Opinion Dynamics task, the most relevant choices are performance, external alignment, and bias, fairness, and ethics metrics.

However, these two steps are not one-time decisions. As the agent design process evolves, evaluation results may prompt adjustments to agent and task attributes, thereby influencing the selection of evaluation metrics. Therefore, this two-step evaluation guideline should be iteratively used to ensure that the evaluation remains adaptive to changing agent capabilities and task requirements. This iterative process enhances the reliability and relevance of evaluations, ultimately leading to more robust and meaningful assessments.

## 5 Case Study: How to Use RPA Design Guideline to Select Evaluation Metrics

We present **two case studies** to illustrate how following the recommendations of our survey leads to the selection of a comprehensive set of evaluation metrics, while significant deviations may result in incomplete evaluation. By placing ourselves in the role of the authors of these articles, we compare the evaluation outcomes resulting from adhering to or deviating from the RPA evaluation guidelines.

### 5.1 A Good Example: *Generative Agents: Interactive Simulacra of Human Behavior*

As shown in Fig. 1, Park et al. (2023) designed agents with demographic information, action history, and social relationships to create an interactive artificial society. Their evaluation methods are in line with the structured selection process proposed in our survey. Since no established agent-oriented evaluation metrics exist for social relationships, they focused on demographic information and action history. Referring to Fig. 3, they identified four relevant metric categories: Content and textual metrics, Internal consistency metrics, External alignment metrics, and Psychological metrics. Based on Tab. 7 in Appendix F, they selected five specific evaluation metrics: Self-knowledge (Content and textual, Internal consistency), Memory and Plans (Internal consistency), Reactions (External alignment), and Reflections (Psychological).

For task-oriented metrics, they determined that the agents' downstream tasks aligned with *simulated society* and designed the evaluation metrics that are aligned with the top three most relevant metric types reported in Fig. 4. As shown in Tab. 8 in Appendix F, they selected four evaluation metrics: Response accuracy (Performance), Relationship formation (Psychological), Information diffusion and Coordination (Social and decision-making). By systematically aligning evaluation metrics with agent attributes and task objectives, this approach ensured a comprehensive and meaningful assessment.

### 5.2 A Flawed Example: *A Generative Social World for Embodied AI*

A flawed example is presented in Appendix E Fig. 8, which is an ICLR submission and the reviews are publicly available on OpenReview. The authors developed agents with demographic attributes, action history, psychological traits, and social relations for route planning and election campaigns. However, their evaluation deviated significantly from our RPA evaluation design guidelines.

Despite designing agents with clear attributes, they did not include any agent-oriented evaluation metrics. For task-oriented metrics, they identified tasks related to Opinion Dynamics and Decision-Making, which should have been evaluated using five key categories: Performance metrics, Psychological metrics, External alignment metrics, Social and decision-making metrics, and Bias, fairness,
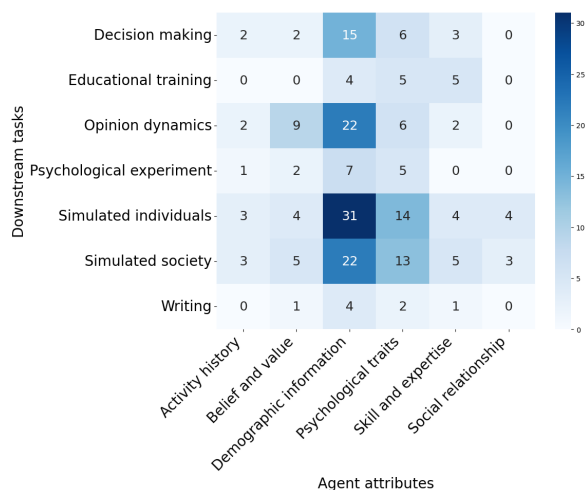
Figure 5: Relationships between agent attributes and downstream tasks. The numbers in the heatmap represent the paper counts.

and ethics metrics. Instead, their evaluation relied only on Arrival rate, Time, and Alignment between campaign strategies, leading to an incomplete assessment. This omission resulted in criticism from reviewers, as one noted: *"The paper performs almost no quantitative experiments... This actually shows that the benchmark cannot cover too many current research methods, which is the biggest weakness of the paper."*

## 6 Relationships Between Agent Attributes and Downstream Tasks

Both agent attributes and downstream tasks play a crucial role in RPA metric selection. Researchers predefine these factors when designing and evaluating RPAs, yet their interrelation remains an open question. In this section, we demonstrate how agent attributes correspond to different downstream tasks by revealing several recurring patterns (Fig. 5).

Demographic information and psychological traits appear fundamental across all downstream tasks. Whether in decision-making, opinion dynamics, or simulated environments, these attributes consistently shape RPA design. As shown in Fig. 5, they are the most frequently incorporated factors, underscoring their central role in modeling agent behavior across diverse applications.

For tasks where simulation itself is the primary objective, such as Simulated Individuals and Simulated Society, the selection of agent attributes becomes broader. In addition to demographic and psychological factors, these tasks frequently incorporate skills, expertise, and social relationships, reflecting the need for richer agent representations to

capture complex social and individual interactions. By contrast, tasks that use simulation as a means to study specific research fields tend to prioritize certain agent attributes. For instance, in Opinion Dynamics, beliefs and values play a distinctive role, as they directly influence how agents interact and form opinions. Similarly, tasks related to Educational Training and Writing exhibit a different pattern, emphasizing skills and expertise over broad demographic or psychological considerations.

In contrast, attributes such as activity history and social relationships receive significantly less emphasis across tasks, raising questions about whether their impact is inherently limited or simply underexplored in current RPA applications.

Overall, these findings highlight the nuanced interplay between agent attributes and downstream tasks. While demographic information and psychological traits are universally relevant, attributes like beliefs and values gain importance in specific contexts. At the same time, the relative absence of activity history and social relationships in current evaluations presents an open research question, particularly in scenarios requiring long-term modeling and complex social interactions.

## 7 Discussion

### 7.1 RPA: an Algorithm v.s. a System

Unlike traditional algorithmic innovations in NLP, the design of RPAs can not only support technical innovations to improve LLMs' humanoid capabilities but also enable RPA-based simulation systems for practical benefits. From the psychology perspective, for instance, RPAs support the exploration of human cognitive and behavioral activities in controlled yet highly scalable experiments, even in hypothetical scenarios. In social science, RPAs can deployed as proxies or pilot experiments to analyze and audit social systems, power dynamics, and human societal behaviors at scale. For the machine learning community, RPAs shed light on dynamic and human-centered model evaluations that are aligned with real-world scenarios by incorporating human and societal factors into consideration. Last but not least, HCI researchers are particularly intrigued by the implications of RPA systems that can provide personalized assistance with human-centered applications in various sectors, such as medicine, healthcare, and education.

Nevertheless, RPAs' capability and flexibility are a double-edged sword; they not only have the

potential to bring benefits to stakeholders but also expose potential risks and even harm if not responsibly designed. To what extent do RPAs' responses align with genuine human cognitive activities, whether the cultural, linguistic, and contextual biases learned from the training data of LLMs impact predicted behaviors, and how to ensure RPAs' robustness and consistency under different scenarios, are critical but underexplored challenges for both technical developers and system designers.

As a result, the design of RPAs should incorporate system design considerations while advancing technical explorations. For instance, RPA design should focus on target users from the very beginning of system design, emphasize the diversity of user backgrounds and perspectives, and iteratively refine the system, as suggested by Gould and Lewis (1985) and Shneiderman and Plaisant (2010) in established design guidelines for system usability. Nevertheless, differences in cultural norms, linguistic subtleties, and domain-specific knowledge can introduce variability in how RPAs are designed and perceived. Designers and developers must focus on a balance between generalization and specificity to ensure RPAs are both adaptable and effective across a wide range of scenarios.

## 7.2 The Design of RPA Persona

A key strength of RPAs lies in their ability to adapt to diverse personas, tasks, and scenarios. How can RPAs' persona be designed to allow LLMs to faithfully and believably reflect the agents' cognitive behaviors with respect to the target task? The persona descriptions of RPAs require careful consideration of both the agents' intrinsic characteristics and the contextual information of the specific environments for which the agents are designed.

The *intrinsic characteristics* of RPAs, such as their personal characteristics, education experience, domain expertise, emotional expressiveness, and decision-making processes, have to be *aligned with the purpose* of the applications of RPAs. For example, an RPA designed for psychological experiments should prioritize cognitive characteristics like personality and empathy ability, whereas an RPA developed for economic simulations might emphasize negotiation tactics, competitive reasoning, and adaptability to changing conditions.

On the other hand, *contextual information*, such as task- and scenario-specific details, factors, and specifications, is equally critical in shaping the behaviors of RPAs. In healthcare applications, for instance, RPAs may simulate caregivers' emotional responses to patients' changing health status but still operate under clinical protocols, such as the ICU visitor rules. The granularity and fidelity of contextual information heavily influence the believability and effectiveness of the agents' behaviors.

## 7.3 The Challenges of RPA Evaluation

The versatility of these agents, which allows them to function in diverse roles and contexts, makes a "one-solution-fits-all" evaluation metric impossible to systematically evaluate RPAs both within and across tasks and user scenarios. One major difficulty lies in designing and determining task-oriented and agent-oriented evaluation metrics. Despite our work recommending an RPA evaluation design guideline based on a comprehensive review of the literature, existing evaluation metrics may not be sufficient to measure the performance of RPAs for different domain-specific applications.

The diversity of user scenarios further exacerbates the evaluation challenge. Different tasks may prioritize different aspects of RPAs, making it difficult to develop a one-size-fits-all evaluation framework. For instance, RPAs designed for psychological research focus on believable emotional responses, whereas RPAs for policymaking simulations underscore robustness to policy changes.

Moreover, cross-task evaluations are particularly challenging due to inconsistencies in how metrics are designed and applied across studies. The lack of standardized evaluation criteria creates barriers to the systematic benchmarking of RPA development and hinders interdisciplinary collaborations across fields. Addressing these challenges will require the development of systematic, multi-faceted evaluation frameworks that can accommodate the diverse applications and capabilities of RPAs while providing consistency and comparability across studies.

## 8 Conclusion

RPA evaluation lacks consistency due to varying tasks, domains, and agent attributes. Our systematic review of $1,676$ papers reveals that task-specific requirements shape agent attributes, while both task characteristics and agent design influence evaluation metrics. By identifying these interdependencies, we propose guidelines to enhance RPA assessment reliability, contributing to a more structured and systematic evaluation framework.

## Limitations

RPAs have widespread applications across various domains and are evolving rapidly. While we strive to examine existing literature on RPAs as comprehensively as possible, we are aware our search is still limited. On the one hand, we may not be able to cover all the varieties of evaluation approaches for RPA in different application domains. On the other hand, new papers about RPAs that were made publicly available after December 2024 are not covered in our work. As a result, our work does not claim to cover all potential evaluation metrics exhaustively, but instead, we aim to offer future researchers a more structured approach and guidelines for designing RPA evaluations.

## Ethics Statement

Our work focuses on summarizing and analyzing the evaluation of RPAs, which we believe will be valuable to researchers in AI, HCI, and related fields such as psychological simulation, educational simulation, and economic simulation. We have made every effort to ensure that this survey is as objective as possible, avoiding both overestimating and underestimating certain trends. We do not anticipate any ethical concerns arising from the research presented in this paper.

## References

Ana Antunes, Joana Campos, Manuel Guimarães, João Dias, and Pedro A. Santos. 2023. Prompting for socially intelligent agents with chatgpt. In *Proceedings of the 23rd ACM International Conference on Intelligent Virtual Agents*, IVA '23, New York, NY, USA. Association for Computing Machinery.

Joshua Ashkinaze, Emily Fry, Narendra Edara, Eric Gilbert, and Ceren Budak. 2024. Plurals: A system for guiding llms via simulated social ensembles. *Preprint*, arXiv:2409.17213.

Sarah Assaf and Timothy Lynar. 2024. Human testing using large-language models: Experimental research and the development of a security awareness controls framework.

Karim Benharrak, Tim Zindulka, Florian Lehmann, Hendrik Heuer, and Daniel Buschek. 2024. Writer-defined ai personas for on-demand feedback generation. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, CHI '24, New York, NY, USA. Association for Computing Machinery.

Ritwik Bose, Mattson Ogg, Michael Wolmetz, and Christopher Ratto. 2024. Assessing behavioral alignment of personality-driven generative agents in social dilemma games. In *NeurIPS 2024 Workshop on Behavioral Machine Learning*.

Elodie Bouzekri, Pascal E Fortin, and Jeremy R Cooperstock. 2024. Chatgpt, tell me more about pilots' opinion on automation. In *2024 IEEE Conference on Cognitive and Computational Aspects of Situation Management (CogSIMA)*, pages 99–106. IEEE.

Meryl Brod, Laura E Tesler, and Torsten L Christensen. 2009. Qualitative research and content validity: developing best practices based on science and experience. *Quality of life research*, 18:1263–1278.

Zhongang Cai, Jianping Jiang, Zhongfei Qing, Xinying Guo, Mingyuan Zhang, Zhengyu Lin, Haiyi Mei, Chen Wei, Ruisi Wang, Wanqi Yin, et al. 2024. Digital life project: Autonomous 3d characters with social intelligence. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 582–592.

Gian Maria Campedelli, Nicolò Penzo, Massimo Stefan, Roberto Dessì, Marco Guerini, Bruno Lepri, and Jacopo Staiano. 2024. I want to break free! persuasion and anti-social behavior of llms in multi-agent settings with social hierarchy. *Preprint*, arXiv:2410.07109.

Louis Castricato, Nathan Lile, Rafael Rafailov, Jan-Philipp Fränken, and Chelsea Finn. 2024. Persona: A reproducible testbed for pluralistic alignment. *Preprint*, arXiv:2407.17387.

Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu, Wei Xue, Shan Zhang, Jie Fu, and Zhiyuan Liu. 2023. Chateval: Towards better llm-based evaluators through multi-agent debate. *ArXiv*, abs/2308.07201.

Chaoran Chen, Leyang Li, Luke Cao, Yanfang Ye, Tianshi Li, Yaxing Yao, and Toby Jia-jun Li. 2024a. Why am i seeing this: Democratizing end user auditing for online content recommendations. *arXiv preprint arXiv:2410.04917*.

Chaoran Chen, Weijun Li, Wenxin Song, Yanfang Ye, Yaxing Yao, and Toby Jia-Jun Li. 2024b. An empathy-based sandbox approach to bridge the privacy gap among attitudes, goals, knowledge, and behaviors. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, CHI '24, New York, NY, USA. Association for Computing Machinery.

Guhong Chen, Liyang Fan, Zihan Gong, Nan Xie, Zixuan Li, Ziqiang Liu, Chengming Li, Qiang Qu, Shiwen Ni, and Min Yang. 2024c. Agentcourt: Simulating court with adversarial evolvable lawyer agents. *arXiv preprint arXiv:2408.08089*.

Jiangjie Chen, Xintao Wang, Rui Xu, Siyu Yuan, Yikai Zhang, Wei Shi, Jian Xie, Shuang Li, Ruihan Yang, Tinghui Zhu, Aili Chen, Nianqi Li, Lida Chen, Caiyu

Hu, Siye Wu, Scott Ren, Ziquan Fu, and Yanghua Xiao. 2024d. From persona to personalization: A survey on role-playing language agents. *Transactions on Machine Learning Research*. Survey Certification.

Nuo Chen, Yan Wang, Yang Deng, and Jia Li. 2024e. The oscars of ai theater: A survey on role-playing with language models. *arXiv preprint arXiv:2407.11484*.

Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang, Chenfei Yuan, Chi-Min Chan, Heyang Yu, Yaxi Lu, Yi-Hsin Hung, Chen Qian, Yujia Qin, Xin Cong, Ruobing Xie, Zhiyuan Liu, Maosong Sun, and Jie Zhou. 2023. Agentverse: Facilitating multi-agent collaboration and exploring emergent behaviors. *Preprint*, arXiv:2308.10848.

Xuzheng Chen, Zhangshiyin, and Guojie Song. 2024f. Towards humanoid: Value-driven agent modeling based on large language models. In *NeurIPS 2024 Workshop on Open-World Agents*.

Haocong Cheng, Si Chen, Christopher Perdriau, and Yun Huang. 2024. Llm-powered ai tutors with personas for d/deaf and hard-of-hearing online learners. *ArXiv*, abs/2411.09873.

Myra Cheng, Tiziano Piccardi, and Diyi Yang. 2023. CoMPosT: Characterizing and evaluating caricature in LLM simulations. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 10853–10875, Singapore. Association for Computational Linguistics.

Yizhou Chi, Lingjun Mao, and Zineng Tang. 2024. Amongagents: Evaluating large language models in the interactive text-based social deduction game. *Preprint*, arXiv:2407.16521.

Yoonseo Choi, Eun Jeong Kang, Seulgi Choi, Min Kyung Lee, and Juho Kim. 2024. Proxona: Leveraging llm-driven personas to enhance creators' understanding of their audience. *arXiv preprint arXiv:2408.10937*.

Yun-Shiuan Chuang, Agam Goyal, Nikunj Harlalka, Siddharth Suresh, Robert Hawkins, Sijia Yang, Dhavan Shah, Junjie Hu, and Timothy T Rogers. 2023a. Simulating opinion dynamics with networks of llm-based agents. *arXiv preprint arXiv:2311.09618*.

Yun-Shiuan Chuang, Siddharth Suresh, Nikunj Harlalka, Agam Goyal, Robert Hawkins, Sijia Yang, Dhavan Shah, Junjie Hu, and Timothy T. Rogers. 2023b. The wisdom of partisan crowds: Comparing collective intelligence in humans and llm-based agents.

Russell Cropanzano and Marie S Mitchell. 2005. Social exchange theory: An interdisciplinary review. *Journal of management*, 31(6):874–900.

Yanqi Dai, Huanran Hu, Lei Wang, Shengjie Jin, Xu Chen, and Zhiwu Lu. 2024. Mmrole: A comprehensive framework for developing and evaluating multimodal role-playing agents. *arXiv preprint arXiv:2408.04203*.

Edoardo Sebastiano De Duro, Riccardo Improta, and Massimo Stella. 2025. Introducing counsellme: A dataset of simulated mental health dialogues for comparing llms like haiku, llamantino and chatgpt against humans. *Emerging Trends in Drugs, Addictions, and Health*, page 100170.

Joost C. F. de Winter, Tom Driessen, and Dimitra Dodou. 2024. The use of chatgpt for personality research: Administering questionnaires using generated personas. *Personality and Individual Differences*.

Jingchao Fang, Nikos Arechiga, Keiichi Namikoshi, Nayeli Bravo, Candice Hogan, and David A Shamma. 2024. On llm wizards: Identifying large language models' behaviors for wizard of oz experiments. In *Proceedings of the 24th ACM International Conference on Intelligent Virtual Agents*, pages 1–11.

Ivar Frisch and Mario Giulianelli. 2024. Llm agents in interaction: Measuring personality consistency and linguistic alignment in interacting populations of large language models. *Preprint*, arXiv:2402.02896.

Chen Gao, Xiaochong Lan, Zhi jie Lu, Jinzhu Mao, Jing Piao, Huandong Wang, Depeng Jin, and Yong Li. 2023. S3: Social-network simulation system with large language model-empowered agents. *ArXiv*, abs/2307.14984.

Chen Gao, Xiaochong Lan, Nian Li, Yuan Yuan, Jingtao Ding, Zhilun Zhou, Fengli Xu, and Yong Li. 2024. Large language models empowered agent-based modeling and simulation: A survey and perspectives. *Humanities and Social Sciences Communications*, 11(1):1–24.

Tao Ge, Xin Chan, Xiaoyang Wang, Dian Yu, Haitao Mi, and Dong Yu. 2024. Scaling synthetic data creation with 1,000,000,000 personas. *arXiv preprint arXiv:2406.20094*.

Tommaso Giorgi, Lorenzo Cima, Tiziano Fagni, Marco Avvenuti, and Stefano Cresci. 2024. Human and llm biases in hate speech annotations: A socio-demographic analysis of annotators and targets. *Preprint*, arXiv:2410.07991.

John D Gould and Clayton Lewis. 1985. Designing for usability: key principles and what designers think. *Communications of the ACM*, 28(3):300–311.

Zhouhong Gu, Xiaoxuan Zhu, Haoran Guo, Lin Zhang, Yin Cai, Hao Shen, Jiangjie Chen, Zheyu Ye, Yifei Dai, Yan Gao, Yao Hu, Hongwei Feng, and Yanghua Xiao. 2024. Agentgroupchat: An interactive group chat simulacra for better eliciting emergent behavior. *Preprint*, arXiv:2403.13433.

George Gui and Olivier Toubia. 2023. The challenge of using llms to simulate human behavior: A causal inference perspective. *ArXiv*, abs/2312.15524.

Shashank Gupta, Vaishnavi Shrivastava, Ameet Deshpande, Ashwin Kalyan, Peter Clark, Ashish Sabharwal, and Tushar Khot. 2024. Bias runs deep: Implicit

reasoning biases in persona-assigned LLMs. In *The Twelfth International Conference on Learning Representations*.

Juhye Ha, Hyeon Jeon, DaEun Han, Jinwook Seo, and Changhoon Oh. 2024. Clochat: Understanding how people customize, interact, and experience personas in large language models. *Proceedings of the CHI Conference on Human Factors in Computing Systems*.

Zhitao He, Pengfei Cao, Chenhao Wang, Zhuoran Jin, Yubo Chen, Jiexin Xu, Huaijun Li, Xiaojian Jiang, Kang Liu, and Jun Zhao. 2024a. Agentscourt: Building judicial decision-making agents with court debate simulation and legal knowledge augmentation. In *Conference on Empirical Methods in Natural Language Processing*.

Zhitao He, Pengfei Cao, Chenhao Wang, Zhuoran Jin, Yubo Chen, Jiexin Xu, Huaijun Li, Kang Liu, and Jun Zhao. 2024b. AgentsCourt: Building judicial decision-making agents with court debate simulation and legal knowledge augmentation. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 9399–9416, Miami, Florida, USA. Association for Computational Linguistics.

Zihong He and Changwang Zhang. 2024. Afspp: Agent framework for shaping preference and personality with large language models. *ArXiv*, abs/2401.02870.

Xu Huang, Weiwen Liu, Xiaolong Chen, Xingmei Wang, Hao Wang, Defu Lian, Yasheng Wang, Ruiming Tang, and Enhong Chen. 2024. Understanding the planning of llm agents: A survey. *arXiv preprint arXiv:2402.02716*.

Yin Jou Huang and Rafik Hadfi. 2024. How personality traits influence negotiation outcomes? a simulation based on large language models. *arXiv preprint arXiv:2407.11549*.

Jiarui Ji, Yang Li, Hongtao Liu, Zhicheng Du, Zhewei Wei, Weiran Shen, Qi Qi, and Yankai Lin. 2024. Srapagent: Simulating and optimizing scarce resource allocation policy with llm-based agent. *arXiv preprint arXiv:2410.14152*.

Jingru Jia, Zehua Yuan, Junhao Pan, Paul E. McNamara, and Deming Chen. 2024. Decision-making behavior evaluation framework for llms under uncertain context. *ArXiv*, abs/2406.05972.

Guangyuan Jiang, Manjie Xu, Song-Chun Zhu, Wenjuan Han, Chi Zhang, and Yixin Zhu. 2023a. Evaluating and inducing personality in pre-trained language models. In *Advances in Neural Information Processing Systems*, volume 36, pages 10622–10643. Curran Associates, Inc.

Hang Jiang, Xiajie Zhang, Xubo Cao, Cynthia Breazeal, Deb Roy, and Jad Kabbara. 2023b. Personallm: Investigating the ability of large language models to express personality traits. In *NAACL-HLT*.

Hang Jiang, Xiajie Zhang, Xubo Cao, Cynthia Breazeal, Deb Roy, and Jad Kabbara. 2024. PersonaLLM: Investigating the ability of large language models to express personality traits. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 3605–3627, Mexico City, Mexico. Association for Computational Linguistics.

Hyoungwook Jin, Seonghee Lee, Hyun Joon Shin, and Juho Kim. 2023. Teach ai how to code: Using large language models as teachable agents for programming education. *Proceedings of the CHI Conference on Human Factors in Computing Systems*.

Yiqiao Jin, Qinlin Zhao, Yiyang Wang, Hao Chen, Kaijie Zhu, Yijia Xiao, and Jindong Wang. 2024. Agentreview: Exploring peer review dynamics with llm agents. In *Conference on Empirical Methods in Natural Language Processing*.

Tianjie Ju, Yiting Wang, Xinbei Ma, Pengzhou Cheng, Haodong Zhao, Yulong Wang, Lifeng Liu, Jian Xie, Zhuosheng Zhang, and Gongshen Liu. 2024. Flooding spread of manipulated knowledge in llm-based multi-agent communities. *ArXiv*, abs/2407.07791.

Zhao Kaiya, Michelangelo Naim, Jovana Kondic, Manuel Cortes, Jiaxin Ge, Shuying Luo, Guangyu Robert Yang, and Andrew Ahn. 2023. Lyfe agents: Generative agents for low-cost real-time social interactions. *Preprint*, arXiv:2310.02172.

Mahammed Kamruzzaman and Gene Louis Kim. 2024. Exploring changes in nation perception with nationality-assigned personas in llms. *Preprint*, arXiv:2406.13993.

Ping Fan Ke and Ka Chung Ng. 2024. Human-ai synergy in survey development: Implications from large language models in business and research. *ACM Transactions on Management Information Systems*.

Kyusik Kim, Hyeonseok Jeon, Jeongwoo Ryu, and Bongwon Suh. 2024. Will llms sink or swim? exploring decision-making under pressure. In *Conference on Empirical Methods in Natural Language Processing*.

Kunyao Lan, Bingrui Jin, Zichen Zhu, Siyuan Chen, Shu Zhang, Kenny Q. Zhu, and Mengyue Wu. 2024. Depression diagnosis dialogue simulation: Self-improving psychiatrist with tertiary memory. *Preprint*, arXiv:2409.15084.

Unggi Lee, Sanghyeok Lee, Junbo Koh, Yeil Jeong, Haewon Jung, Gyuri Byun, Jewoong Moon, Jieun Lim, and † HyeoncheolKim. Generative agent for teacher training: Designing educational problem-solving simulations with large language model-based agents for pre-service teachers.

Yu Lei, Hao Liu, Chengxing Xie, Songjia Liu, Zhiyu Yin, Canyu Chen, Guohao Li, Philip Torr, and Zhen Wu. 2024. Fairmindsim: Alignment of behavior, emotion, and belief in humans and llm agents amid ethical dilemmas. *arXiv preprint arXiv:2410.10398*.

11

Yan Leng and Yuan Yuan. 2024. Do llm agents exhibit social behavior? *Preprint*, arXiv:2312.15198.

Hao Li, Chenghao Yang, An Zhang, Yang Deng, Xiang Wang, and Tat-Seng Chua. 2024a. Hello again! llm-powered personalized agent for long-term dialogue. *arXiv preprint arXiv:2406.05925*.

Jiale Li, Jiayang Li, Jiahao Chen, Yifan Li, Shijie Wang, Hugo Zhou, Minjun Ye, and Yunsheng Su. 2024b. Evolving agents: Interactive simulation of dynamic and diverse human personalities. *ArXiv*, abs/2404.02718.

Junkai Li, Yunghwei Lai, Weitao Li, Jingyi Ren, Meng Zhang, Xinhui Kang, Siyu Wang, Peng Li, Ya-Qin Zhang, Weizhi Ma, et al. 2024c. Agent hospital: A simulacrum of hospital with evolvable medical agents. *arXiv preprint arXiv:2405.02957*.

Nian Li, Chen Gao, Mingyu Li, Yong Li, and Qingmin Liao. 2024d. Econagent: large language model-empowered agents for simulating macroeconomic activities. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15523–15536.

Sha Li, Revanth Gangi Reddy, Khanh Duy Nguyen, Qingyun Wang, May Fung, Chi Han, Jiawei Han, Kartik Natarajan, Clare R. Voss, and Heng Ji. 2024e. Schema-guided culture-aware complex event simulation with multi-agent role-play. *ArXiv*, abs/2410.18935.

Yuan Li, Yixuan Zhang, and Lichao Sun. 2023a. Metaagents: Simulating interactions of human behaviors for llm-based task-oriented coordination via collaborative generative agents. *ArXiv*, abs/2310.06500.

Yuan Li, Yixuan Zhang, and Lichao Sun. 2023b. Metaagents: Simulating interactions of human behaviors for llm-based task-oriented coordination via collaborative generative agents. *Preprint*, arXiv:2310.06500.

Xiaoyu Lin, Xinkai Yu, Ankit Aich, Salvatore Giorgi, and Lyle Ungar. 2024. Diversedialogue: A methodology for designing chatbots with human-like diversity. *Preprint*, arXiv:2409.00262.

Jiaheng Liu, Zehao Ni, Haoran Que, Tao Sun, Noah Wang, Jian Yang, JiakaiWang, Hongcheng Guo, Z.Y. Peng, Ge Zhang, Jiayi Tian, Xingyuan Bu, Ke Xu, Wenge Rong, Junran Peng, and Zhaoxiang Zhang. 2024a. Roleagent: Building, interacting, and benchmarking high-quality role-playing agents from scripts. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.

Ryan Liu, Howard Yen, Raja Marjieh, Thomas L. Griffiths, and Ranjay Krishna. 2023. Improving interpersonal communication by simulating audiences with language models. *Preprint*, arXiv:2311.00687.

Tianjian Liu, Hongzheng Zhao, Yuheng Liu, Xingbo Wang, and Zhenhui Peng. 2024b. Compeer: A generative conversational agent for proactive peer support. In *ACM Symposium on User Interface Software and Technology*.

Xuan Liu, Jie Zhang, Song Guo, Haoyang Shang, Chengxu Yang, and Quanyan Zhu. 2025. Exploring prosocial irrationality for llm agents: A social cognition view. *Preprint*, arXiv:2405.14744.

Yuhan Liu, Zirui Song, Xiaoqing Zhang, Xiuying Chen, and Rui Yan. 2024c. From a tiny slip to a giant leap: An llm-based simulation for fake news evolution. *arXiv preprint arXiv:2410.19064*.

Zhengyuan Liu, Stella Xin Yin, Geyu Lin, and Nancy F. Chen. 2024d. Personality-aware student simulation for conversational intelligent tutoring systems. In *Conference on Empirical Methods in Natural Language Processing*.

Yaojia Lv, Haojie Pan, Zekun Wang, Jiafeng Liang, Yuanxing Liu, Ruiji Fu, Ming Liu, Zhongyuan Wang, and Bing Qin. 2024. Coggpt: Unleashing the power of cognitive dynamics on large language models. *arXiv preprint arXiv:2401.08438*.

Jiří Milička, Anna Marklová, Klára VanSlambrouck, Eva Pospíšilová, Jana Šimsová, Samuel Harvan, and Ondřej Drobil. 2024. Large language models are able to downplay their cognitive abilities to fit the persona they simulate. *Plos one*, 19(3):e0298522.

Kshitij Mishra, Priyanshu Priya, Manisha Burja, and Asif Ekbal. 2023. e-THERAPIST: I suggest you to cultivate a mindset of positivity and nurture uplifting thoughts. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 13952–13967, Singapore. Association for Computational Linguistics.

Konstantinos Mitsopoulos, Ritwik Bose, Brodie Mather, Archna Bhatia, Kevin Gluck, Bonnie Dorr, Christian Lebiere, and Peter Pirolli. 2024. Psychologically-valid generative agents: A novel approach to agent-based modeling in social sciences. *Proceedings of the AAAI Symposium Series*.

Suhong Moon, Marwa Abdulhai, Minwoo Kang, Joseph Suh, Widyadewi Soedarmadji, Eran Kohen Behar, and David M. Chan. 2024. Virtual personas for language models via an anthology of backstories. *Preprint*, arXiv:2407.06576.

Xinyi Mou, Xuanwen Ding, Qi He, Liang Wang, Jingcong Liang, Xinnong Zhang, Libo Sun, Jiayu Lin, Jie Zhou, Xuanjing Huang, et al. 2024a. From individual to society: A survey on social simulation driven by large language model-based agents. *arXiv preprint arXiv:2412.03563*.

Xinyi Mou, Jingcong Liang, Jiayu Lin, Xinnong Zhang, Xiawei Liu, Shiyue Yang, Rong Ye, Lei Chen, Haoyu Kuang, Xuanjing Huang, and Zhongyu Wei. 2024b.

12

Agentsense: Benchmarking social intelligence of language agents through interactive scenarios. *Preprint*, arXiv:2410.19346.

Xinyi Mou, Zhongyu Wei, and Xuanjing Huang. 2024c. Unveiling the truth and facilitating change: Towards agent-based large-scale social movement simulation. In *Annual Meeting of the Association for Computational Linguistics*.

Sonia K. Murthy, Tomer Ullman, and Jennifer Hu. 2024. One fish, two fish, but not the whole sea: Alignment reduces language models' conceptual diversity. *Preprint*, arXiv:2411.04427.

Keiichi Namikoshi, Alexandre L. S. Filipowicz, David A. Shamma, Rumen Iliev, Candice Hogan, and Nikos Aréchiga. 2024. Using llms to model the beliefs and preferences of targeted populations. *ArXiv*, abs/2403.20252.

Alejandro Leonardo Garc'ia Navarro, Nataliia Koneva, Alfonso S'anchez-Maci'an, Jos'e Alberto Hern'andez, and Manuel Goyanes. 2024. Designing reliable experiments with generative agent-based modeling: A comprehensive guide using concordia by google deepmind. *ArXiv*, abs/2411.07038.

Shlomo Neuberger, Niv Eckhaus, Uri Berger, Amir Taubenfeld, Gabriel Stanovsky, and Ariel Goldstein. 2024. Sauce: Synchronous and asynchronous user-customizable environment for multi-agent llm interaction. *arXiv preprint arXiv:2411.03397*.

Alison Nightingale. 2009. A guide to systematic literature reviews. *Surgery (Oxford)*, 27(9):381–384.

Joon Sung Park, Joseph O'Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th annual acm symposium on user interface software and technology*, pages 1–22.

Joon Sung Park, Lindsay Popowski, Carrie Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. 2022. Social simulacra: Creating populated prototypes for social computing systems. In *Proceedings of the 35th Annual ACM Symposium on User Interface Software and Technology*, UIST '22, New York, NY, USA. Association for Computing Machinery.

Joon Sung Park, Carolyn Q. Zou, Aaron Shaw, Benjamin Mako Hill, Carrie Cai, Meredith Ringel Morris, Robb Willer, Percy Liang, and Michael S. Bernstein. 2024. Generative agent simulations of 1,000 people. *Preprint*, arXiv:2411.10109.

Pat Pataranutaporn, Kavin Winson, Peggy Yin, Auttasak Lapapirojn, Pichayoot Ouppaphan, Monchai Lertsutthiwong, Pattie Maes, and Hal E. Hershfield. 2024. Future you: A conversation with an ai-generated future self reduces anxiety, negative emotions, and increases future self-continuity. *ArXiv*, abs/2405.12514.

Giorgio Piatti, Zhijing Jin, Max Kleiman-Weiner, Bernhard Schölkopf, Mrinmaya Sachan, and Rada Mihalcea. 2024. Cooperate or collapse: Emergence of sustainable cooperation in a society of llm agents. *Preprint*, arXiv:2404.16698.

Aske Plaat, Annie Wong, Suzan Verberne, Joost Broekens, Niki van Stein, and Thomas Back. 2024. Reasoning with large language models, a survey. *arXiv preprint arXiv:2407.11511*.

Huachuan Qiu and Zhenzhong Lan. 2024. Interactive agents: Simulating counselor-client psychological counseling via role-playing llm-to-llm interactions. *Preprint*, arXiv:2408.15787.

Yao Qu and Jue Wang. 2024. Performance and biases of large language models in public opinion simulation. *Humanities and Social Sciences Communications*, 11(1):1–13.

Ruiyang Ren, Peng Qiu, Yingqi Qu, Jing Liu, Wayne Xin Zhao, Huaqin Wu, Ji-Rong Wen, and Haifeng Wang. 2024a. Bases: Large-scale web search user simulation with large language model based agents. *ArXiv*, abs/2402.17505.

Siyue Ren, Zhiyao Cui, Ruiqi Song, Zhen Wang, and Shuyue Hu. 2024b. Emergence of social norms in generative agent societies: Principles and architecture. *Preprint*, arXiv:2403.08251.

Joni O. Salminen, João M. Santos, Soon gyo Jung, and Bernard J. Jansen. 2024. Picturing the fictitious person: An exploratory study on the effect of images on user perceptions of ai-generated personas. *Computers in Human Behavior: Artificial Humans*.

Andreas Schuller, Doris Janssen, Julian Blumenröther, Theresa Maria Probst, Michael Schmidt, and Chandan Kumar. 2024. Generating personas using llms and assessing their viability. In *Extended Abstracts of the CHI Conference on Human Factors in Computing Systems*, CHI EA '24, New York, NY, USA. Association for Computing Machinery.

Avi Schwarzschild, Zhili Feng, Pratyush Maini, Zachary Lipton, and J Zico Kolter. 2025. Rethinking llm memorization through the lens of adversarial compression. *Advances in Neural Information Processing Systems*, 37:56244–56267.

Yunfan Shao, Linyang Li, Junqi Dai, and Xipeng Qiu. 2023. Character-llm: A trainable agent for role-playing. *arXiv preprint arXiv:2310.10158*.

Jinxin Shi, Jiabao Zhao, Yilei Wang, Xingjiao Wu, Jiawen Li, and Liangbo He. 2023. Cgmi: Configurable general multi-agent interaction framework. *ArXiv*, abs/2308.12503.

Joongi Shin, Michael A. Hedderich, Bartłomiej Jakub Rey, Andrés Lucero, and Antti Oulasvirta. 2024. Understanding human-ai workflows for generating personas. In *Proceedings of the 2024 ACM Designing Interactive Systems Conference*, DIS '24, page

13

757–781, New York, NY, USA. Association for Computing Machinery.

Ben Shneiderman and Catherine Plaisant. 2010. *Designing the user interface: strategies for effective human-computer interaction*. Pearson Education India.

Chan Hee Song, Jiaman Wu, Clayton Washington, Brian M Sadler, Wei-Lun Chao, and Yu Su. 2023. Llm-planner: Few-shot grounded planning for embodied agents with large language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2998–3009.

Sinan Sonlu, Bennie Bendiksen, Funda Durupinar, and Uğur Güdükbay. 2024. The effects of embodiment and personality expression on learning in llm-based educational agents. *ArXiv*, abs/2407.10993.

Karthik Sreedhar and Lydia Chilton. 2024. Simulating human strategic behavior: Comparing single and multi-agent llms. *arXiv preprint arXiv:2402.08189*.

Libo Sun, Siyuan Wang, Xuanjing Huang, and Zhongyu Wei. 2024. Identity-driven hierarchical role-playing agents. *Preprint*, arXiv:2407.19412.

Eduardo Ryô Tamaki and Levente Littvay. 2024. Chrono-sampling: Generative ai enabled time machine for public opinion data collection.

Yihong Tang, Jiao Ou, Che Liu, Fuzheng Zhang, Di Zhang, and Kun Gai. 2024. Erabal: Enhancing role-playing agents through boundary-aware learning. *Preprint*, arXiv:2409.14710.

Amir Taubenfeld, Yaniv Dover, Roi Reichart, and Ariel Goldstein. 2024. Systematic biases in LLM simulations of debates. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 251–267, Miami, Florida, USA. Association for Computational Linguistics.

Jesus-Pablo Toledo-Zucco, Denis Matignon, and Charles Poussot-Vassal. 2024. Scattering-passive structure-preserving finite element method for the boundary controlled transport equation with a moving mesh. *Preprint*, arXiv:2402.01232.

Haley Triem and Ying Ding. 2024. "tipping the balance": Human intervention in large language model multi-agent debate. *Proceedings of the Association for Information Science and Technology*, 61(1):361–373.

Yu-Min Tseng, Yu-Chao Huang, Teng-Yun Hsiao, Wei-Lin Chen, Chao-Wei Huang, Yu Meng, and Yun-Nung Chen. 2024. Two tales of persona in LLMs: A survey of role-playing and personalization. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 16612–16631, Miami, Florida, USA. Association for Computational Linguistics.

Quan Tu, Shilong Fan, Zihang Tian, and Rui Yan. 2024. Charactereval: A chinese benchmark for role-playing conversational agent evaluation. *arXiv preprint arXiv:2401.01275*.

Deepank Verma, Olaf Mumm, and Vanessa Miriam Carlow. 2023. Generative agents in the streets: Exploring the use of large language models (llms) in collecting urban perceptions. *ArXiv*, abs/2312.13126.

Boshi Wang, Xiang Yue, and Huan Sun. 2023a. Can chatgpt defend its belief in truth? evaluating llm reasoning via debate. *arXiv preprint arXiv:2305.13160*.

Lei Wang, Jingsen Zhang, Hao Yang, Zhi-Yuan Chen, Jiakai Tang, Zeyu Zhang, Xu Chen, Yankai Lin, Ruihua Song, Wayne Xin Zhao, Jun Xu, Zhicheng Dou, Jun Wang, and Ji rong Wen. 2023b. User behavior simulation with large language model based agents.

Qian Wang, Tianyu Wang, Qinbin Li, Jingsheng Liang, and Bingsheng He. 2024a. Megaagent: A practical framework for autonomous cooperation in large-scale llm agent systems. *Preprint*, arXiv:2408.09955.

Qian Wang, Jiaying Wu, Zhenheng Tang, Bingqiao Luo, Nuo Chen, Wei Chen, and Bingsheng He. 2025. What limits llm-based human simulation: Llms or our design? *arXiv preprint arXiv:2501.08579*.

Xiaolong Wang, Yile Wang, Sijie Cheng, Peng Li, and Yang Liu. 2024b. Deem: Dynamic experienced expert modeling for stance detection. *ArXiv*, abs/2402.15264.

Xintao Wang, Yunze Xiao, Jen-tse Huang, Siyu Yuan, Rui Xu, Haoran Guo, Quan Tu, Yaying Fei, Ziang Leng, Wei Wang, et al. 2024c. Incharacter: Evaluating personality fidelity in role-playing agents through psychological interviews. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1840–1873.

Yi Wang, Qian Zhou, and David Ledo. 2024d. Storyverse: Towards co-authoring dynamic plot with llm-based character simulation via narrative planning. In *Proceedings of the 19th International Conference on the Foundations of Digital Games*, FDG '24, New York, NY, USA. Association for Computing Machinery.

Zhengbo Wang, Jian Liang, Ran He, Zilei Wang, and Tieniu Tan. 2024e. Connecting the dots: Collaborative fine-tuning for black-box vision-language models. *arXiv preprint arXiv:2402.04050*.

Zhenhailong Wang, Shaoguang Mao, Wenshan Wu, Tao Ge, Furu Wei, and Heng Ji. 2024f. Unleashing the emergent cognitive synergy in large language models: A task-solving agent through multi-persona self-collaboration. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 257–279, Mexico City, Mexico. Association for Computational Linguistics.

14

Zhenyu Wang, Yi Xu, Dequan Wang, Lingfeng Zhou, and Yiqi Zhou. 2024g. Intelligent computing social modeling and methodological innovations in political science in the era of large language models. *ArXiv*, abs/2410.16301.

Zengqing Wu, Shuyuan Zheng, Qianying Liu, Xu Han, Brian Inhyuk Kwon, Makoto Onizuka, Shaojie Tang, Run Peng, and Chuan Xiao. 2024. Shall we talk: Exploring spontaneous collaborations of competing llm agents. *arXiv preprint arXiv:2402.12327*.

Chengxing Xie, Canyu Chen, Feiran Jia, Ziyu Ye, Kai Shu, Adel Bibi, Ziniu Hu, Philip Torr, Bernard Ghanem, and G. Li. 2024a. Can large language model agents simulate human trust behaviors? *ArXiv*, abs/2402.04559.

Qiuejie Xie, Qiming Feng, Tianqi Zhang, Qingqiu Li, Linyi Yang, Yuejie Zhang, Rui Feng, Liang He, Shang Gao, and Yue Zhang. 2024b. Human simulacra: Benchmarking the personification of large language models. *Preprint*, arXiv:2402.18180.

Zihan Yan, Yaohong Xiang, and Yun Huang. 2024. Social life simulation for non-cognitive skills learning. *ArXiv*, abs/2405.00273.

Frank Tian-fang Ye and Xiaozi Gao. 2024. Simulating family conversations using llms: Demonstration of parenting styles. *arXiv preprint arXiv:2403.06144*.

Leo Yeykelis, Kaavya Pichai, James J. Cummings, and Byron Reeves. 2024. Using large language models to create ai personas for replication and prediction of media effects: An empirical test of 133 published experimental research findings. *Preprint*, arXiv:2408.16073.

Chenxiao Yu, Zhaotian Weng, Yuangang Li, Zheng Li, Xiyang Hu, and Yue Zhao. 2024. Towards more accurate us presidential election via multi-step reasoning with large language models. *ArXiv*, abs/2411.03321.

Zheni Zeng, Jiayi Chen, Huimin Chen, Yukun Yan, Yuxuan Chen, Zhenghao Liu, Zhiyuan Liu, and Maosong Sun. 2024. Persllm: A personified training approach for large language models. *arXiv preprint arXiv:2407.12393*.

Dong Zhang, Zhaowei Li, Pengyu Wang, Xin Zhang, Yaqian Zhou, and Xipeng Qiu. 2024a. Speechagents: Human-communication simulation with multi-modal multi-agent systems. *ArXiv*, abs/2401.03945.

Jintian Zhang, Xin Xu, Ruibo Liu, and Shumin Deng. 2023a. Exploring collaboration mechanisms for llm agents: A social psychology view. *ArXiv*, abs/2310.02124.

Long Zhang, Meng Zhang, Wei Lin Wang, and Yu Luo. 2025. Simulation as reality? the effectiveness of llm-generated data in open-ended question assessment. *arXiv preprint arXiv:2502.06371*.

Qiang Zhang, Jason Naradowsky, and Yusuke Miyao. 2024b. Self-emotion blended dialogue generation in social simulation agents. *Preprint*, arXiv:2408.01633.

Yu Zhang, Jingwei Sun, Li Feng, Cen Yao, Mingming Fan, Liuxin Zhang, Qianying Wang, Xin Geng, and Yong Rui. 2024c. See widely, think wisely: Toward designing a generative multi-agent system to burst filter bubbles. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, CHI '24, New York, NY, USA. Association for Computing Machinery.

Zhaowei Zhang, Ceyao Zhang, Nian Liu, Siyuan Qi, Ziqi Rong, Song-Chun Zhu, Shuguang Cui, and Yaodong Yang. 2023b. Heterogeneous value alignment evaluation for large language models. *arXiv preprint arXiv:2305.17147*.

Qinlin Zhao, Jindong Wang, Yixuan Zhang, Yiqiao Jin, Kaijie Zhu, Hao Chen, and Xing Xie. 2023. Competeai: Understanding the competition dynamics of large language model-based agents. In *International Conference on Machine Learning*.

Xuhui Zhou, Zhe Su, Tiwalayo Eisape, Hyunwoo Kim, and Maarten Sap. 2024a. Is this the real life? is this just fantasy? the misleading success of simulating social interactions with llms. *ArXiv*, abs/2403.05020.

Xuhui Zhou, Hao Zhu, Leena Mathur, Ruohong Zhang, Haofei Yu, Zhengyang Qi, Louis-Philippe Morency, Yonatan Bisk, Daniel Fried, Graham Neubig, and Maarten Sap. 2024b. Sotopia: Interactive evaluation for social intelligence in language agents. *Preprint*, arXiv:2310.11667.

15

Table 4: Inclusion and exclusion criteria.

| Inclusion Criteria (IC) | |
| --- | --- |
| IC-1 | The LLM agents in the paper simulate humanoid behavior with implicit personality (e.g., preference and behavior pattern) or explicit personality (e.g., emotion or characteristics). |
| IC-2 | The LLM agents in the paper have cognitive activities such as decision-making, reasoning, and planning. |
| IC-3 | The LLM agents in the paper are capable of completing complicated and general tasks. |
| IC-4 | The LLM agents' action set in the paper is neither predefined nor finite. |

| Exclusion Criteria (EC) | |
| --- | --- |
| EC-1 | The study does not employ LLM agents for simulation purposes but rather uses them as chatbots, task-specific agents, or evaluators. |
| EC-2 | The paper's research objectives, methodologies, and evaluations are not focused on simulating human-like behavior with LLM agents, but rather on optimizing LLM algorithms. |
| EC-3 | The study primarily investigates the perception or action capabilities of LLM agents without simulating the cognitive process. |
| EC-4 | The LLM agents are restricted to handling specific, close-ended tasks. |
| EC-5 | The LLM agents' actions are either predefined or limited. |



Figure 6: Usage ratio of evaluation approaches for each category of agent-oriented metrics.



Figure 7: Usage ratio of evaluation approaches for each category of task-oriented metrics.

## A    Inclusion and Exclusion Criteria

We summarize the inclusion and exclusion criteria in Table 4. Briefly, the **Inclusion Criteria (IC)** ensure that the reviewed studies focus on LLM agents exhibiting human-like behavior—either implicitly (e.g., preference or behavioral patterns) or explicitly (e.g., emotions or personality)—along with key cognitive processes such as reasoning and decision-making. Moreover, an open-ended action space and the capacity to tackle multifaceted tasks are essential attributes for inclusion.

By contrast, the **Exclusion Criteria (EC)** eliminate studies employing LLMs purely as chatbots, single-purpose systems, or evaluation tools, rather than as agents mimicking human cognition. Likewise, if the LLM agents are restricted to fixed, close-ended tasks or limited to algorithmic optimization without simulating cognitive processes, they fall outside the scope of this work.

## B    Query String

We employed the following query to guide our literature retrieval process:

```
("large language model" OR LLM)
AND (agent OR persona OR "human
digital twin" OR simulacra) AND
(simulat* OR generat* OR eval*)
AND "human behavior" AND cognit*
```

This query was designed to capture a broad spectrum of studies on large language models that simulate or replicate human-like behavior. It combines keywords related to LLM agents (*LLM*, *persona*, *simulacra*), their capabilities (*simulat\**, *generat\**, *eval\**), and the focus on cognitively grounded human behavior (*cognit\**). This ensures that the resulting literature is relevant to our exploration of how LLM-based systems can mimic or exhibit human-like cognition and behavior patterns.

## C    Evaluation Approach Usage for Agent- and Task-Oriented Metrics

We present a breakdown of evaluation approach usage by agent-oriented metrics (Fig. 6) and task-oriented metrics (Fig. 7).

## D    Top Three Metrics for Agent and Task Attributes

We present two tables for referencing the top three frequently used metrics for agent attributes (Tab. 5) and task attributes (Tab. 6).

| Agent attributes | Top 3 agent-oriented metrics |
| --- | --- |
| Activity history | External alignment metrics, internal consistency metrics, content and textual metrics |
| Belief and value | Psychological metrics, bias, fairness, and ethics metrics |
| Demographic information | Psychological metrics, internal consistency metrics, external alignment metrics |
| Psychological traits | Psychological metrics, internal consistency metrics, content and textual metrics |
| Skill and expertise | External alignment metrics, internal consistency metrics, content and textual metrics |
| Social relationship | Psychological metrics, external alignment metrics, social and decision-making metrics |

Table 5: Top 3 frequently used agent-oriented metrics for each agent attribute

| Task attributes | Top 3 task-oriented metrics |
| --- | --- |
| Simulated individuals | Psychological, performance, and internal consistency metrics |
| Simulated society | Social and decision-making metrics, performance metrics, and psychological metrics |
| Opinion dynamics | Performance metrics, external alignment metrics, and bias, fairness, and ethics metrics |
| Decision making | Social and decision-making, performance, and psychological metrics |
| Psychological experiment | Psychological, content and textual, and performance metrics |
| Educational training | Psychological, performance, and content and textual metrics |
| Writing | Content and textual, psychological, and performance metrics |

Table 6: Top 3 frequently used task-oriented metrics for each task attribute



**Example Project**: *"...the LLM generates agent profiles along with their social relationships. The profiles consist of basic attributes such as names, ages, occupations, personalities, and hobbies...generate the daily schedule for each agent"*

**RPA**   **Agent Design**: *{name, age, occupation, hobby, personality}*
**Task**: *{route planning and election campaign}*

**STEP 1: Decide agent-oriented metrics based on agent attributes**

**STEP 2: Decide task-oriented metrics based on task attributes**

**Reviewer comments:** *"The paper performs almost no quantitative experiments...This actually shows that the benchmark cannot cover too many current research methods, which is the biggest weakness of the paper."*

Figure 8: Case study of a flawed example in Section 5.2. Given agent attributes (yellow) and task attributes (pink). The original authors' selection of evaluation metrics (purple and blue). The missing metrics that are recommended by our proposed guideline (orange) align with the reviewer's criticism in red text.

# E   Case Study: Flawed Example

Fig. 8 visualized how the authors in the flawed example selected their evaluation metrics how further evaluation metrics could be uncovered through our proposed guideline.

# F   Metrics Glossary

We present two glossary tables for referencing the source of agent-oriented metrics (Tab. 7) and task-oriented metrics (Tab. 8).

Table 7: Agent-oriented evaluation metrics glossary.

| Attribute | Category | Agent-oriented Metrics | Approach Source |
|---|---|---|---|
| Belief & Value | Bias, fairness, ethics metrics | Exaggeration (normalized average cosine similarity) | Automatic (Cheng et al., 2023) |
| Belief & Value | Bias, fairness, ethics metrics | Individuation (classification accuracy) | Automatic (Cheng et al., 2023) |
| Belief & Value | Bias, fairness, ethics metrics | Bias (performance disparity, prevalence, magnitude, variation, attitude shift) | Automatic (Gupta et al., 2024) |
| Belief & Value | Bias, fairness, ethics metrics | Bias (performance disparity, prevalence, magnitude, variation, attitude shift) | Automatic (Taubenfeld et al., 2024) |
| Demographic Information | Bias, fairness, ethics metrics | Exaggeration (normalized average cosine similarity) | Automatic (Cheng et al., 2023) |
| Demographic Information | Bias, fairness, ethics metrics | Individuation (classification accuracy) | Automatic (Cheng et al., 2023) |
| Demographic Information | Bias, fairness, ethics metrics | Bias (performance disparity, prevalence, magnitude, variation, attitude shift) | Automatic (Gupta et al., 2024) |
| Demographic Information | Bias, fairness, ethics metrics | Bias (performance disparity, prevalence, magnitude, variation, attitude shift) | Automatic (Neuberger et al., 2024) |
| Demographic Information | Bias, fairness, ethics metrics | Bias (performance disparity, prevalence, magnitude, variation, attitude shift) | Automatic (Taubenfeld et al., 2024) |
| Demographic Information | Bias, fairness, ethics metrics | Message toxicity | Automatic (Fang et al., 2024) |
| Activity History | Content and textual metrics | Coherence | LLM (Li et al., 2024e) |
| Activity History | Content and textual metrics | Clarity | Human (Chen et al., 2024b) |
| Activity History | Content and textual metrics | Diversity of dialog (Shannon entropy, intra-remote-clique, inter-remote-clique, semantic similarity, longest common subsequence similarity) | Automatic (Ha et al., 2024) |
| Belief & Value | Content and textual metrics | Diversity of dialog (Shannon entropy, intra-remote-clique, inter-remote-clique, semantic similarity, longest common subsequence similarity) | Automatic (Gu et al., 2024) |
| Demographic Information | Content and textual metrics | Coherence | LLM (Li et al., 2024e) |
| Demographic Information | Content and textual metrics | Attitudes (topic term frequency) | Automatic (Fang et al., 2024) |
| Demographic Information | Content and textual metrics | Diversity of dialog (Shannon entropy, intra-remote-clique, inter-remote-clique, semantic similarity, longest common subsequence similarity) | Automatic (Fang et al., 2024) |
| Demographic Information | Content and textual metrics | Clarity | Human (Chen et al., 2024b) |
| Demographic Information | Content and textual metrics | Diversity of dialog (Shannon entropy, intra-remote-clique, inter-remote-clique, semantic similarity, longest common subsequence similarity) | Automatic (Ha et al., 2024) |
| Demographic Information | Content and textual metrics | Linguistic complexity (utterance length, Kolmogorov complexity) | Automatic (Milička et al., 2024) |
| Psychological Traits | Content and textual metrics | Text similarity (BLEU, ROUGE) | Automatic (Zeng et al., 2024) |
| Psychological Traits | Content and textual metrics | Tone Alignment | LLM (Zeng et al., 2024) |
| Skills and Expertise | Content and textual metrics | Coherence | LLM (Li et al., 2024e) |
| Activity History | External alignment metrics | Hallucination | LLM (Shao et al., 2023) |
| Activity History | External alignment metrics | Entailment | LLM (Li et al., 2024e) |
| Activity History | External alignment metrics | Believability/Credibility(self-knowledge, memory, plans, reactions, reflections) | Human (Park et al., 2023) |

<div align="center">Continued on next page</div>

| Attribute | Category | Agent-oriented Metrics | Approach | Source |
|---|---|---|---|---|
| Demographic Information | External alignment metrics | Entailment | LLM | (Li et al., 2024e) |
| Demographic Information | External alignment metrics | Believability/Credibility(self-knowledge, memory, plans, reactions, reflections) | Human | (Park et al., 2023) |
| Psychological Traits | External alignment metrics | Fact Accuracy | LLM | (Zeng et al., 2024) |
| Skills and Expertise | External alignment metrics | Hallucination | LLM | (Shao et al., 2023) |
| Skills and Expertise | External alignment metrics | Entailment | LLM | (Li et al., 2024e) |
| Activity History | Internal consistency metrics | Stability | LLM | (Shao et al., 2023) |
| Activity History | Internal consistency metrics | Consistency of information | Human | (Chen et al., 2024b) |
| Belief & Value | Internal consistency metrics | Attitude shift | LLM | (Wang et al., 2024e) |
| Demographic Information | Internal consistency metrics | Stability | LLM | (Shao et al., 2023) |
| Demographic Information | Internal consistency metrics | Attitude shift | LLM | (Neuberger et al., 2024) |
| Demographic Information | Internal consistency metrics | Attitude shift | LLM | (Taubenfeld et al., 2024) |
| Demographic Information | Internal consistency metrics | Behavior stability (mean, standard deviation) | Automatic | (Wang et al., 2024g) |
| Demographic Information | Internal consistency metrics | Consistency of information | Human | (Chen et al., 2024b) |
| Demographic Information | Internal consistency metrics | Consistency of psychological state / personalities | Human | (Chen et al., 2024b) |
| Demographic Information | Internal consistency metrics | Consistency of information | Human | (Zeng et al., 2024) |
| Psychological Traits | Internal consistency metrics | Stability | LLM | (Shao et al., 2023) |
| Psychological Traits | Internal consistency metrics | Consistency of information | Human | (Zeng et al., 2024) |
| Psychological Traits | Internal consistency metrics | Consistency of psychological state / personalities | Human | (Zeng et al., 2024) |
| Psychological Traits | Internal consistency metrics | Consistency of information | Human | (Cai et al., 2024) |
| Psychological Traits | Internal consistency metrics | Consistency of psychological state / personalities | Human | (Cai et al., 2024) |
| Skills and Expertise | Internal consistency metrics | Stability | LLM | (Shao et al., 2023) |
| Activity History | Performance metrics | Memorization | LLM | (Shao et al., 2023) |
| Demographic Information | Performance metrics | Memorization | LLM | (Chen et al., 2024b) |
| Demographic Information | Performance metrics | Communication ability (win rates) | Automatic | (Liu et al., 2024a) |
| Demographic Information | Performance metrics | Reaction (accuracy) | Automatic | (Liu et al., 2024a) |
| Demographic Information | Performance metrics | Self-knowledge (accuracy) | Automatic | (Liu et al., 2024a) |
| Activity History | Psychological metrics | Empathy | Human | (Chen et al., 2024b) |
| Belief & Value | Psychological metrics | Value | LLM | (Shao et al., 2023) |
| Demographic Information | Psychological metrics | Personality consistency | Automatic | (Wang et al., 2024c) |
| Demographic Information | Psychological metrics | Measured alignment for personality | Human | (Wang et al., 2024c) |
| Demographic Information | Psychological metrics | Sentiment | Automatic | (Fang et al., 2024) |
| Demographic Information | Psychological metrics | Empathy | Human | (Chen et al., 2024b) |
| Demographic Information | Psychological metrics | Belief (stability, evolution, correlation with behavior) | Automatic | (Lei et al., 2024) |

<div align="center">Continued on next page</div>

| Attribute | Category | Agent-oriented Metrics | Approach Source |
|---|---|---|---|
| Psychological Traits | Psychological metrics | Personality | Automatic (Shao et al., 2023) |
| Psychological Traits | Psychological metrics | Belief (stability, evolution, correlation with behavior) | Automatic (Shao et al., 2023) |
| Psychological Traits | Psychological metrics | Emotion responses (entropy of valence and arousal) | Automatic (Shao et al., 2023) |
| Psychological Traits | Psychological metrics | Personality (Machine Personality Inventory, PsychoBench) | Automatic (Jiang et al., 2023a) |
| Psychological Traits | Psychological metrics | Personality (vignette tests) | Human (Jiang et al., 2023a) |
| Belief & Value | Social and decision-making metrics | Social value orientation (SVO-based Value Rationality Measurement) | Automatic (Zhang et al., 2023b) |

Table 8: Task-oriented evaluation metrics glossary.

| Task | Category | Task-oriented Metrics | Approach Source |
|------|----------|----------------------|-----------------|
| Decision Making | Social and economic metrics | Negotiation (Concession Rate, Negotiation Success Rate, Average Negotiation Round) | Automatic (Huang and Hadfi, 2024) |
| Decision Making | Social and economic metrics | Societal Satisfaction (average per-capita living area size, average waiting time, social welfare) | Automatic (Ji et al., 2024) |
| Decision Making | Social and economic metrics | Societal Fairness (variance in per capita living area size, number of inverse order pairs in house allocation, Gini coefficient) | Automatic (Ji et al., 2024) |
| Decision Making | Social and economic metrics | Macroeconomic (Inflation rate, Unemployment rate, Nominal GDP, Nominal GDP growth, Wage inflation, Real GDP growth, Expected monthly income, Consumption) | Automatic (Li et al., 2024d) |
| Decision Making | Social and economic metrics | Market and Consumer (Purchase probability, Expected competing product price, Customer counts, Price consistency between competitors) | Automatic (Gui and Toubia, 2023) |
| Decision Making | Social and economic metrics | Market and Consumer (Purchase probability, Expected competing product price, Customer counts, Price consistency between competitors) | Automatic (Zhao et al., 2023) |
| Decision Making | Social and economic metrics | Probability weighting | Automatic (Jia et al., 2024) |
| Decision Making | Social and economic metrics | Utility (Intrinsic Utility, Joint Utility) | Automatic (Huang and Hadfi, 2024) |
| Decision Making | Psychological metrics | Level of trust (distribution of amounts sent, trust rate) | Automatic (Xie et al., 2024a) |
| Decision Making | Psychological metrics | Risk preference | Automatic (Jia et al., 2024) |
| Decision Making | Psychological metrics | Loss aversion | Automatic (Jia et al., 2024) |
| Decision Making | Psychological metrics | Selfishness (Selfishness Index, Difference Index) | Automatic (Kim et al., 2024) |
| Decision Making | Performance metrics | Frequency (distribution of expert type) | Automatic (Wang et al., 2024b) |
| Decision Making | Performance metrics | Valid response rate | Automatic (Xie et al., 2024a) |
| Decision Making | Performance metrics | Web search quality (Mean reciprocal rank, Mean reciprocal rank) | Automatic (Ren et al., 2024a) |
| Decision Making | Performance metrics | Performance deviations/alignment from the baseline (accuracy, Jaccard Index, Cohen's Kappa Coefficient, Percentage Agreement, overlapping ratio between prediction and targets) | Automatic (Kim et al., 2024) |
| Decision Making | Performance metrics | Performance deviations/alignment from the baseline (accuracy, Jaccard Index, Cohen's Kappa Coefficient, Percentage Agreement, overlapping ratio between prediction and targets) | Automatic (Jin et al., 2024) |
| Decision Making | Performance metrics | Performance deviations/alignment from the baseline (accuracy, Jaccard Index, Cohen's Kappa Coefficient, Percentage Agreement, overlapping ratio between prediction and targets) | Automatic (Wang et al., 2024b) |
| Decision Making | Performance metrics | Performance deviations/alignment from the baseline (accuracy, Jaccard Index, Cohen's Kappa Coefficient, Percentage Agreement, overlapping ratio between prediction and targets) | Automatic (Wang et al., 2024f) |
| Decision Making | Internal consistency metrics | Behavioral alignment (lottery rate, behavior dynamic, Imitation and differentiation behavior, Proportion of similar and different dishes) | Automatic (Xie et al., 2024a) |

Continued on next page

| Task | Category | Task-oriented Metrics | Approach | Source |
|---|---|---|---|---|
| Decision Making | Internal consistency metrics | Behavioral alignment (lottery rate, behavior dynamic, Imitation and differentiation behavior, Proportion of similar and different dishes) | Automatic | (Zhao et al., 2023) |
| Decision Making | Internal consistency metrics | Cultural appropriateness (Alignment between persona information and its assigned nationality) | LLM | (Li et al., 2024e) |
| Decision Making | External alignment metrics | Factual hallucinations (String matching overlap ratio) | Automatic | (Wang et al., 2024f) |
| Decision Making | External alignment metrics | Simulation capability (Turing test) | Human | (Ji et al., 2024) |
| Decision Making | External alignment metrics | Entailment | LLM | (Li et al., 2024e) |
| Decision Making | External alignment metrics | Realism | LLM | (Li et al., 2024e) |
| Educational Training | Psychological metrics | Perceived reflection on the development of essential non-cognitive skills | Human | (Yan et al., 2024) |
| Educational Training | Psychological metrics | Non-cognitive skill scale | Automatic | (Yan et al., 2024) |
| Educational Training | Psychological metrics | Sense of immersion / Perceived immersion | Human | (Lee et al.) |
| Educational Training | Psychological metrics | Perceived intelligence | Human | (Cheng et al., 2024) |
| Educational Training | Psychological metrics | Perceived enjoyment | Human | (Cheng et al., 2024) |
| Educational Training | Psychological metrics | Perceived trust | Human | (Cheng et al., 2024) |
| Educational Training | Psychological metrics | Perceived sense of connection | Human | (Cheng et al., 2024) |
| Educational Training | Psychological metrics | Personality (Big Five Invertory, MBTI score, SD3 score, Linguistic Inquiry and Word Count framework, HEX-ACO) | Automatic | (Sonlu et al., 2024) |
| Educational Training | Psychological metrics | Personality (Big Five Invertory, MBTI score, SD3 score, Linguistic Inquiry and Word Count framework, HEX-ACO) | Automatic | (Liu et al., 2024d) |
| Educational Training | Psychological metrics | Perceived usefulness | Human | (Cheng et al., 2024) |
| Educational Training | Performance metrics | Density of knowledge-building | Automatic | (Jin et al., 2023) |
| Educational Training | Performance metrics | Effectiveness of questioning | Human | (Shi et al., 2023) |
| Educational Training | Performance metrics | Success criterion function outputs before operation and after operation | Human | (Li et al., 2023a) |
| Educational Training | External alignment metrics | Knowledge level (reconfigurability, persistence, and adaptability) | Automatic | (Jin et al., 2023) |
| Educational Training | External alignment metrics | Perceived human-likeness | Human | (Cheng et al., 2024) |
| Educational Training | Content and textual metrics | Story Content Generation (narratives staging score) | Automatic | (Yan et al., 2024) |
| Educational Training | Content and textual metrics | Willingness to speak | Human | (Shi et al., 2023) |
| Educational Training | Content and textual metrics | Authenticity | Human | (Lee et al.) |
| Opinion Dynamics | Psychological metrics | Opinion change | Human | (Triem and Ding, 2024) |
| Opinion Dynamics | Psychological metrics | Emotional density | Automatic | (Gao et al., 2023) |
| Opinion Dynamics | Performance metrics | Prediction accuracy (F1 score, AUC, MSE, MAE, depression risk prediction accuracy, suicide risk prediction accuracy) | Automatic | (Gao et al., 2023) |

| Task | Category | Task-oriented Metrics | Approach Source |
|------|----------|----------------------|-----------------|
| Opinion Dynamics | Performance metrics | Prediction accuracy (F1 score, AUC, MSE, MAE, depression risk prediction accuracy, suicide risk prediction accuracy) | Automatic (Mou et al., 2024c) |
| Opinion Dynamics | Performance metrics | Prediction accuracy (F1 score, AUC, MSE, MAE, depression risk prediction accuracy, suicide risk prediction accuracy) | Automatic (Yu et al., 2024) |
| Opinion Dynamics | Performance metrics | Classification accuracy | Human (Chan et al., 2023) |
| Opinion Dynamics | Performance metrics | Rephrase accuracy | Automatic (Ju et al., 2024) |
| Opinion Dynamics | Performance metrics | Legal articles evaluation (precision, recall, F1) | Automatic (He et al., 2024a) |
| Opinion Dynamics | Performance metrics | Judgment evaluation for civil and administrative cases (precision, recall, F1) | Automatic (He et al., 2024a) |
| Opinion Dynamics | Performance metrics | Judgment evaluation for criminal cases (accuracy) | Automatic (He et al., 2024a) |
| Opinion Dynamics | Performance metrics | Prediction error rate | Automatic (Gao et al., 2023) |
| Opinion Dynamics | Performance metrics | Locality accuracy | Automatic (Ju et al., 2024) |
| Opinion Dynamics | Performance metrics | Decision probability | Human (Triem and Ding, 2024) |
| Opinion Dynamics | Performance metrics | Decision volatility | Human (Triem and Ding, 2024) |
| Opinion Dynamics | Performance metrics | Case complexity | Human (Triem and Ding, 2024) |
| Opinion Dynamics | Performance metrics | Alignment (compare simulation results with actual social outcomes) | Automatic (Wang et al., 2024g) |
| Opinion Dynamics | Internal consistency metrics | Alignment (stance, content, behavior, static attitude distribution, time series of the average attitude) | Automatic (Mou et al., 2024c) |
| Opinion Dynamics | Internal consistency metrics | Personality-behavior alignment | Human (Navarro et al., 2024) |
| Opinion Dynamics | Internal consistency metrics | Similarity between initial and post preference (KL-divergence, RMSE) | Automatic (Namikoshi et al., 2024) |
| Opinion Dynamics | Internal consistency metrics | Role playing | Human (Lv et al., 2024) |
| Opinion Dynamics | External alignment metrics | Correctness | Human (He et al., 2024a) |
| Opinion Dynamics | External alignment metrics | Accuracy (correctness) | Automatic (Ju et al., 2024) |
| Opinion Dynamics | External alignment metrics | Logicality | Human (He et al., 2024a) |
| Opinion Dynamics | External alignment metrics | Concision | Human (He et al., 2024a) |
| Opinion Dynamics | External alignment metrics | Human likeness index | Automatic (Chuang et al., 2023b) |
| Opinion Dynamics | External alignment metrics | Alignment between model and human (Kappa correlation coefficient, MAE), Authenticity (alignment of ratings between the agent and human annotators) | Human (Chan et al., 2023) |
| Opinion Dynamics | External alignment metrics | Alignment between model and human (Kappa correlation coefficient, MAE), Authenticity (alignment of ratings between the agent and human annotators) | Human (Triem and Ding, 2024) |
| Opinion Dynamics | External alignment metrics | Alignment between model and human (Kappa correlation coefficient, MAE), Authenticity (alignment of ratings between the agent and human annotators) | Human (Lv et al., 2024) |
| Opinion Dynamics | Content and textual metrics | Turn-level Kendall-Tau correlation (naturalness, coherence, engagingness and groundedness) | Automatic (Chan et al., 2023) |

| Task | Category | Task-oriented Metrics | Approach Source |
|---|---|---|---|
| Opinion Dynamics | Content and textual metrics | Turn-level Spearman correlation (naturalness, coherence, engagingness and groundedness) | Automatic (Chan et al., 2023) |
| Opinion Dynamics | Bias, fairness, and ethic metrics | Partisan bias | Automatic (Chuang et al., 2023b) |
| Opinion Dynamics | Bias, fairness, and ethic metrics | Bias (cultural, linguistic, economic, demographic, ideological) | Automatic (Qu and Wang, 2024) |
| Opinion Dynamics | Bias, fairness, and ethic metrics | Bias (mean) | Automatic (Chuang et al., 2023a) |
| Opinion Dynamics | Bias, fairness, and ethic metrics | Extreme values | Automatic (Chuang et al., 2023b) |
| Opinion Dynamics | Bias, fairness, and ethic metrics | Wisdom of Partisan Crowds effect | Automatic (Chuang et al., 2023b) |
| Opinion Dynamics | Bias, fairness, and ethic metrics | Opinion diversity | Automatic (Chuang et al., 2023a) |
| Psychological Experiment | Social and economic metrics | Money allocation | Automatic (Lei et al., 2024) |
| Psychological Experiment | Psychological metrics | Attitude change | Automatic (Wang et al., 2023b) |
| Psychological Experiment | Psychological metrics | Average happiness value per time step | Automatic (He and Zhang, 2024) |
| Psychological Experiment | Psychological metrics | Belief value | Automatic (Lei et al., 2024) |
| Psychological Experiment | Psychological metrics | Personality (Big Five Invertory, MBTI score, SD3 score, Linguistic Inquiry and Word Count framework, HEX-ACO) | Automatic (He and Zhang, 2024) |
| Psychological Experiment | Psychological metrics | Personality (Big Five Invertory, MBTI score, SD3 score, Linguistic Inquiry and Word Count framework, HEX-ACO) | Automatic (de Winter et al., 2024) |
| Psychological Experiment | Psychological metrics | Personality (Big Five Invertory, MBTI score, SD3 score, Linguistic Inquiry and Word Count framework, HEX-ACO) | Automatic (Bose et al., 2024) |
| Psychological Experiment | Psychological metrics | Personality (Big Five Invertory, MBTI score, SD3 score, Linguistic Inquiry and Word Count framework, HEX-ACO) | Automatic (Jiang et al., 2023b) |
| Psychological Experiment | Psychological metrics | Longitudinal trajectories of emotions | Automatic (De Duro et al., 2025) |
| Psychological Experiment | Psychological metrics | Valence entropy | Automatic (Lei et al., 2024) |
| Psychological Experiment | Psychological metrics | Arousal entropy | Automatic (Lei et al., 2024) |
| Psychological Experiment | Performance metrics | Precision of item recommendation | Automatic (Wang et al., 2023b) |
| Psychological Experiment | Performance metrics | Missing rate | Automatic (Lei et al., 2024) |
| Psychological Experiment | Performance metrics | Rejection rate | Automatic (Lei et al., 2024) |
| Psychological Experiment | Internal consistency metrics | Correlation between social dilemma game outcome and agent personality | Automatic (Bose et al., 2024) |
| Psychological Experiment | Internal consistency metrics | Behavioral similarity | Automatic (Li et al., 2024b) |
| Psychological Experiment | Internal consistency metrics | Perception consistency (agent perceived safety, agent perceived liveliness) | LLM (Verma et al., 2023) |
| Psychological Experiment | External alignment metrics | Rationality of the agent memory | Automatic (Wang et al., 2023b) |
| Psychological Experiment | External alignment metrics | Believability of behavior | Automatic (Wang et al., 2023b) |
| Psychological Experiment | Content and textual metrics | Salience of individual words | Automatic (De Duro et al., 2025) |
| Psychological Experiment | Content and textual metrics | Absolutist words | Automatic (De Duro et al., 2025) |

| Task | Category | Task-oriented Metrics | Approach Source |
|---|---|---|---|
| Psychological Experiment | Content and textual metrics | Personal pronouns or emotions | Automatic (De Duro et al., 2025) |
| Psychological Experiment | Content and textual metrics | Information entropy | Automatic (Wang et al., 2023b) |
| Psychological Experiment | Content and textual metrics | Story (readability, personalness, redundancy, cohesiveness, likeability, believability) | Human (Jiang et al., 2023b) |
| Psychological Experiment | Content and textual metrics | Story (readability, personalness, redundancy, cohesiveness, likeability, believability) | LLM (Jiang et al., 2023b) |
| Simulated Individual | Social and economic metrics | Numbers of generated peer support strategies | Automatic (Liu et al., 2024b) |
| Simulated Individual | Social and economic metrics | Perceived social support questionnaire | Human (Liu et al., 2024b) |
| Simulated Individual | Psychological metrics | Emotions | Human (Pataranutaporn et al., 2024) |
| Simulated Individual | Psychological metrics | Agency | Human (Pataranutaporn et al., 2024) |
| Simulated Individual | Psychological metrics | Future consideration | Human (Pataranutaporn et al., 2024) |
| Simulated Individual | Psychological metrics | Self-reflection | Human (Pataranutaporn et al., 2024) |
| Simulated Individual | Psychological metrics | Insight | Human (Pataranutaporn et al., 2024) |
| Simulated Individual | Psychological metrics | Persona Perception Scale | Human (Salminen et al., 2024) |
| Simulated Individual | Psychological metrics | Persona Perception Scale | Human (Shin et al., 2024) |
| Simulated Individual | Psychological metrics | Persona Perception Scale | Human (Ha et al., 2024) |
| Simulated Individual | Psychological metrics | Persona Perception Scale | Human (Chen et al., 2024b) |
| Simulated Individual | Psychological metrics | Engagement | Human (Zhang et al., 2024a) |
| Simulated Individual | Psychological metrics | Safety | Human (Zhang et al., 2024a) |
| Simulated Individual | Psychological metrics | Sensitivity to personalization | Automatic (Giorgi et al., 2024) |
| Simulated Individual | Psychological metrics | Agent self-awareness | LLM (Xie et al., 2024b) |
| Simulated Individual | Psychological metrics | Personality (Big Five Invertory rated by LLM) | LLM (Jiang et al., 2023a) |
| Simulated Individual | Psychological metrics | Positively mention rate | Automatic (Kamruzzaman and Kim, 2024) |
| Simulated Individual | Psychological metrics | Optimism | Human (Pataranutaporn et al., 2024) |
| Simulated Individual | Psychological metrics | Self-esteem | Human (Pataranutaporn et al., 2024) |
| Simulated Individual | Psychological metrics | Pressure perceived scale | Human (Liu et al., 2024b) |
| Simulated Individual | Performance metrics | Error rates (error of average, error of dispersion) | Automatic (Lin et al., 2024) |
| Simulated Individual | Performance metrics | Model fit indices (Chi-square to degrees of freedom ratio, Comparative Fit Index, Tucker-Lewis Index, Root Mean Square Error of Approximation) | Automatic (Ke and Ng, 2024) |
| Simulated Individual | Performance metrics | Knowledge accuracy (WikiRoleEval with human evaluators) | Human (Tang et al., 2024) |
| Simulated Individual | Performance metrics | Knowledge accuracy (WikiRoleEval) | LLM (Tang et al., 2024) |
| Simulated Individual | Performance metrics | Win rates | Automatic (Chi et al., 2024) |
| Simulated Individual | Performance metrics | Comprehension | Automatic (Shin et al., 2024) |
| Simulated Individual | Performance metrics | Completeness | Automatic (Shin et al., 2024) |

| Task | Category | Task-oriented Metrics | Approach | Source |
|---|---|---|---|---|
| Simulated Individual | Performance metrics | Validity (average variance extracted, inter-construct correlations) | Automatic | (Ke and Ng, 2024) |
| Simulated Individual | Performance metrics | Composite reliability | Automatic | (Ke and Ng, 2024) |
| Simulated Individual | Performance metrics | Rated statement quality | Human | (Liu et al., 2023) |
| Simulated Individual | Performance metrics | Rated statement quality | LLM | (Liu et al., 2023) |
| Simulated Individual | Performance metrics | Conversational ability (CharacterEval) | LLM | (Tang et al., 2024) |
| Simulated Individual | Performance metrics | Roleplay subset of MT-Bench | LLM | (Tang et al., 2024) |
| Simulated Individual | Performance metrics | Professional scale (accuracy in replicating profession-specific knowledge) | LLM | (Sun et al., 2024) |
| Simulated Individual | Performance metrics | Language quality | LLM | (Zhang et al., 2024a) |
| Simulated Individual | Performance metrics | Prediction accuracy between real data and generated data (Replication success rate, Kullback-Leibler divergence) | Automatic | (Assaf and Lynar, 2024) |
| Simulated Individual | Performance metrics | Prediction accuracy between real data and generated data (Replication success rate, Kullback-Leibler divergence) | Automatic | (Tamaki and Littvay, 2024) |
| Simulated Individual | Performance metrics | Prediction accuracy between real data and generated data (Replication success rate, Kullback-Leibler divergence) | Automatic | (Park et al., 2024) |
| Simulated Individual | Performance metrics | Prediction accuracy between real data and generated data (Replication success rate, Kullback-Leibler divergence) | Automatic | (Yeykelis et al., 2024) |
| Simulated Individual | Performance metrics | Accuracy of distinguishing between AI-generated and human-built solutions | Automatic | (Schuller et al., 2024) |
| Simulated Individual | Internal consistency metrics | Accuracy of reaction based on social relationship | Automatic | (Liu et al., 2024a) |
| Simulated Individual | Internal consistency metrics | Perceived connection between personas and system outcomes | Human | (Chen et al., 2024b) |
| Simulated Individual | Internal consistency metrics | Representativeness (Wasserstein distance, respond with similar answers to individual survey questions), Consistency (Frobenius norm, the correlation across responses to a set of questions in each survey) | Automatic | (Moon et al., 2024) |
| Simulated Individual | Internal consistency metrics | Role consistency (WikiRoleEval with human evaluators) | Human | (Tang et al., 2024) |
| Simulated Individual | Internal consistency metrics | Role consistency/attractiveness (WikiRoleEval, CharacterEval) | LLM | (Tang et al., 2024) |
| Simulated Individual | Internal consistency metrics | Consistency | Human | (Zhang et al., 2024a) |
| Simulated Individual | Internal consistency metrics | Consistency | Human | (Mishra et al., 2023) |
| Simulated Individual | Internal consistency metrics | Future self-continuity | Human | (Pataranutaporn et al., 2024) |
| Simulated Individual | Internal consistency metrics | Agreement between a synthetic annotator both with and without a leave-one-out attribute (Cohen's Kappa) | Automatic | (Castricato et al., 2024) |
| Simulated Individual | Internal consistency metrics | Consistency with the scenario and characters | Automatic | (Zhang et al., 2024a) |
| Simulated Individual | Internal consistency metrics | Quality and logical coherence of the script content | Automatic | (Zhang et al., 2024a) |
| Simulated Individual | Internal consistency metrics | Nation-related response percentage | Automatic | (Kamruzzaman and Kim, 2024) |

<div align="center">Continued on next page</div>

| Task | Category | Task-oriented Metrics | Approach | Source |
|---|---|---|---|---|
| Simulated Individual | External alignment metrics | Unknown question rejection (WikiRoleEval with human evaluators) | Human | (Tang et al., 2024) |
| Simulated Individual | External alignment metrics | Unknown question rejection (WikiRoleEval) | LLM | (Tang et al., 2024) |
| Simulated Individual | External alignment metrics | Accuracy of self-knowledge | Automatic | (Liu et al., 2024a) |
| Simulated Individual | External alignment metrics | Correctness | Human | (Zhang et al., 2024a) |
| Simulated Individual | External alignment metrics | Correctness | Human | (Milička et al., 2024) |
| Simulated Individual | External alignment metrics | Agreement score between human raters and LLM, | Automatic | (Liu et al., 2023) |
| Simulated Individual | External alignment metrics | Agreement score between human raters and LLM, | Automatic | (Jiang et al., 2023a) |
| Simulated Individual | External alignment metrics | Agreement score between human raters and LLM, | Automatic | (Liu et al., 2024a) |
| Simulated Individual | External alignment metrics | Human-likeness | Human | (Zhang et al., 2024a) |
| Simulated Individual | Content and textual metrics | Content similarity (ROUGE-L, BERTScore, GPT-based-similarity, G-eval) | Automatic | (Shin et al., 2024) |
| Simulated Individual | Content and textual metrics | Entity density of summarization | Automatic | (Liu et al., 2024a) |
| Simulated Individual | Content and textual metrics | Entity recall of summarization | Automatic | (Liu et al., 2024a) |
| Simulated Individual | Content and textual metrics | Dialog diversity | Automatic | (Lin et al., 2024) |
| Simulated Individual | Bias, fairness, and ethic metrics | Hate speech detection accuracy | Automatic | (Giorgi et al., 2024) |
| Simulated Individual | Bias, fairness, and ethic metrics | Population heterogeneity | Automatic | (Murthy et al., 2024) |
| Simulated Society | Social and economic metrics | Social Conflict Count | Automatic | (Ren et al., 2024b) |
| Simulated Society | Social and economic metrics | Social Rules | Human | (Zhou et al., 2024b) |
| Simulated Society | Social and economic metrics | Social Rules | LLM | (Zhou et al., 2024b) |
| Simulated Society | Social and economic metrics | Financial and Material Benefits | Human | (Zhou et al., 2024b) |
| Simulated Society | Social and economic metrics | Financial and Material Benefits | LLM | (Zhou et al., 2024b) |
| Simulated Society | Social and economic metrics | Converged price | Automatic | (Toledo-Zucco et al., 2024) |
| Simulated Society | Social and economic metrics | Information diffusion | Automatic | (Park et al., 2023) |
| Simulated Society | Social and economic metrics | Relationship formation | Automatic | (Park et al., 2023) |
| Simulated Society | Social and economic metrics | Relationship | LLM | (Zhou et al., 2024b) |
| Simulated Society | Social and economic metrics | Coordination within other agents | Automatic | (Park et al., 2023) |
| Simulated Society | Social and economic metrics | Probability of social connection formation | Automatic | (Leng and Yuan, 2024) |
| Simulated Society | Social and economic metrics | Percent of social welfare maximization choices | Automatic | (Leng and Yuan, 2024) |
| Simulated Society | Social and economic metrics | Persuasion (distribution of persuasion outcomes, odds ratios) | Automatic | (Campedelli et al., 2024) |
| Simulated Society | Social and economic metrics | Anti-social behavior (effect on toxic messages) | Automatic | (Campedelli et al., 2024) |
| Simulated Society | Social and economic metrics | Norm Internalization Rate | Automatic | (Ren et al., 2024b) |
| Simulated Society | Social and economic metrics | Norm Compliance Rate | Automatic | (Ren et al., 2024b) |
| Simulated Society | Psychological metrics | NASA-TLX Scores | Human | (Zhang et al., 2024c) |

<div align="center">Continued on next page</div>

| Task | Category | Task-oriented Metrics | Approach | Source |
|---|---|---|---|---|
| Simulated Society | Psychological metrics | Helpfulness rating | Human | (Zhang et al., 2024c) |
| Simulated Society | Psychological metrics | Personality (Big Five Invertory, MBTI score, SD3 score, Linguistic Inquiry and Word Count framework, HEX-ACO) | Automatic | (Frisch and Giulianelli, 2024) |
| Simulated Society | Psychological metrics | Personality (Big Five Invertory, MBTI score, SD3 score, Linguistic Inquiry and Word Count framework, HEX-ACO) | Automatic | (Li et al., 2024b) |
| Simulated Society | Psychological metrics | Degree of reciprocity | Automatic | (Leng and Yuan, 2024) |
| Simulated Society | Psychological metrics | Pleasure rating | Human | (Zhang et al., 2024c) |
| Simulated Society | Psychological metrics | Trend of Favorability Decline | Automatic | (Gu et al., 2024) |
| Simulated Society | Psychological metrics | Negative Favorability Achievement | Automatic | (Gu et al., 2024) |
| Simulated Society | Psychological metrics | Trend of Favorability Decline | Automatic | (Gu et al., 2024) |
| Simulated Society | Psychological metrics | Negative Favorability Achievement | Automatic | (Gu et al., 2024) |
| Simulated Society | Performance metrics | Abstention accuracy | Automatic | (Ashkinaze et al., 2024) |
| Simulated Society | Performance metrics | Accuracy of information gathering | Automatic | (Kaiya et al., 2023) |
| Simulated Society | Performance metrics | Implicit reasoning accuracy | Automatic | (Mou et al., 2024b) |
| Simulated Society | Performance metrics | Prediction accuracy (F1 score, AUC, MSE, MAE, depression risk prediction accuracy, suicide risk prediction accuracy) | Automatic | (Lan et al., 2024) |
| Simulated Society | Performance metrics | Guess accuracy | Automatic | (Leng and Yuan, 2024) |
| Simulated Society | Performance metrics | Classification accuracy | Automatic | (Li et al., 2024a) |
| Simulated Society | Performance metrics | Success rate | Automatic | (Kaiya et al., 2023) |
| Simulated Society | Performance metrics | Success rate | Automatic | (Li et al., 2023b) |
| Simulated Society | Performance metrics | Success rate | Automatic | (Li et al., 2023b) |
| Simulated Society | Performance metrics | Success rate for coordination (identification accuracy, workflow correctness, alignment between job and agent's skill) | Automatic | (Li et al., 2023a) |
| Simulated Society | Performance metrics | Success rate for coordination (identification accuracy, workflow correctness, alignment between job and agent's skill) | Automatic | (Li et al., 2023a) |
| Simulated Society | Performance metrics | Task Accuracy | Automatic | (Zhang et al., 2023a) |
| Simulated Society | Performance metrics | Task Accuracy | Automatic | (Lan et al., 2024) |
| Simulated Society | Performance metrics | Errors in the prompting sequence | Human | (Antunes et al., 2023) |
| Simulated Society | Performance metrics | Error-free execution | Automatic | (Wang et al., 2024a) |
| Simulated Society | Performance metrics | Goal completion | Human | (Mou et al., 2024b) |
| Simulated Society | Performance metrics | Goal completion | LLM | (Zhou et al., 2024a) |
| Simulated Society | Performance metrics | Goal completion | LLM | (Mou et al., 2024b) |
| Simulated Society | Performance metrics | Goal completion | LLM | (Zhou et al., 2024b) |

| Task | Category | Task-oriented Metrics | Approach | Source |
|------|----------|----------------------|----------|--------|
| Simulated Society | Performance metrics | Efficacy | Human | (Ashkinaze et al., 2024) |
| Simulated Society | Performance metrics | Knowledge | Human | (Zhou et al., 2024b) |
| Simulated Society | Performance metrics | Knowledge | LLM | (Zhou et al., 2024b) |
| Simulated Society | Performance metrics | Reasoning abilities | Automatic | (Chen et al., 2023) |
| Simulated Society | Performance metrics | Reasoning abilities | Human | (Chen et al., 2023) |
| Simulated Society | Performance metrics | Efficiency | Automatic | (Piatti et al., 2024) |
| Simulated Society | Performance metrics | Text understanding and creative writing abilities (Dialogue response dataset, Commongen Challenge) | LLM | (Chen et al., 2023) |
| Simulated Society | Performance metrics | Probabilities of receiving, storing, and retrieving the key information across the population | Automatic | (Kaiya et al., 2023) |
| Simulated Society | Performance metrics | Correlation between predicted and real results | Automatic | (Mitsopoulos et al., 2024) |
| Simulated Society | Internal consistency metrics | Behavioral similarity | Automatic | (Li et al., 2024b) |
| Simulated Society | Internal consistency metrics | Semantic consistency (cosine similarity) | Automatic | (Qiu and Lan, 2024) |
| Simulated Society | External alignment metrics | Alignment (Environmental understanding and response accuracy, adherence to predefined settings) | Automatic | (Gu et al., 2024) |
| Simulated Society | External alignment metrics | Strategy accuracy (strategies provided by the models vs. by human experts and evaluate the accuracy) | Automatic | (Zhang et al., 2024b) |
| Simulated Society | External alignment metrics | Believability of behavior | Human | (Zhou et al., 2024b) |
| Simulated Society | External alignment metrics | Believability of behavior | Human | (Park et al., 2023) |
| Simulated Society | Content and textual metrics | Content similarity (ROUGE-L, BERTScore, GPT-based-similarity, G-eval, BLEU-4) | Automatic | (Li et al., 2024a) |
| Simulated Society | Content and textual metrics | Content similarity (ROUGE-L, BERTScore, GPT-based-similarity, G-eval) | Automatic | (Chen et al., 2024f) |
| Simulated Society | Content and textual metrics | Content similarity (ROUGE-L, BERTScore, GPT-based-similarity, G-eval) | Automatic | (Mishra et al., 2023) |
| Simulated Society | Content and textual metrics | Semantic understanding | Automatic | (Gu et al., 2024) |
| Simulated Society | Content and textual metrics | Complexity of generated content | Automatic | (Antunes et al., 2023) |
| Simulated Society | Content and textual metrics | Dialogue generation quality | Automatic | (Antunes et al., 2023) |
| Simulated Society | Content and textual metrics | Number of conversation rounds | Automatic | (Zhang et al., 2024c) |
| Simulated Society | Bias, fairness, and ethic metrics | Bias rate (herd effect, authority effect, ban franklin effect, rumor chain effect, gambler's fallacy, confirmation bias, halo effect) | Human | (Liu et al., 2025) |
| Simulated Society | Bias, fairness, and ethic metrics | Bias rate (herd effect, authority effect, ban franklin effect, rumor chain effect, gambler's fallacy, confirmation bias, halo effect) | LLM | (Liu et al., 2025) |
| Simulated Society | Bias, fairness, and ethic metrics | Bias rate (herd effect, authority effect, ban franklin effect, rumor chain effect, gambler's fallacy, confirmation bias, halo effect) | Automatic | (Liu et al., 2025) |
| Simulated Society | Bias, fairness, and ethic metrics | Equality | Automatic | (Piatti et al., 2024) |

Continued on next page

| Task | Category | Task-oriented Metrics | Approach | Source |
|------|----------|-----------------------|----------|--------|
| Writing | Psychological metrics | Qualitative feedback (expertise, social relation, valence, level of involvement) | Human | (Benharrak et al., 2024) |
| Writing | Performance metrics | Prediction accuracy (F1 score, AUC, MSE, MAE, depression risk prediction accuracy, suicide risk prediction accuracy) | Automatic | (Wang et al., 2024f) |
| Writing | Performance metrics | Success rate | Automatic | (Wang et al., 2024d) |
| Writing | Performance metrics | Behavioral patterns | Human | (Zhang et al., 2024c) |
| Writing | Internal consistency metrics | Consistency (user profile, psychotherapeutic approach) | Automatic | (Mishra et al., 2023) |
| Writing | Internal consistency metrics | Motivational consistency | LLM | (Wang et al., 2024d) |
| Writing | Internal consistency metrics | Audience similarity | Human | (Choi et al., 2024) |
| Writing | Internal consistency metrics | Quality of generated dimension & values (relevance, mutual exclusiveness) | Human | (Choi et al., 2024) |
| Writing | External alignment metrics | Factual error rate | Automatic | (Wang et al., 2024f) |
| Writing | External alignment metrics | Correctness (politeness, interpersonal behaviour) | Automatic | (Mishra et al., 2023) |
| Writing | External alignment metrics | Hallucination (groundedness of the chat responses) | Human | (Choi et al., 2024) |
| Writing | Content and textual metrics | Linguistic similarity | Human | (Choi et al., 2024) |
| Writing | Content and textual metrics | Fluency | Human | (Mishra et al., 2023) |
| Writing | Content and textual metrics | Perplexity | Automatic | (Mishra et al., 2023) |
| Writing | Content and textual metrics | Non-Repetitiveness | Human | (Mishra et al., 2023) |
| Writing | Content and textual metrics | response generation quality | Automatic | (Li et al., 2024a) |
| Writing | Content and textual metrics | Coherency | LLM | (Wang et al., 2024d) |