
Mindless Dialogue? A Critical Note on Theory of Mind and Communicative Generative Agents

Effat Farhana

E.FARHANA@AUBURN.EDU

Computing Science and Software Engineering, Auburn University, Auburn, AL 36849 USA

Abstract

Theory of Mind (ToM), the ability to attribute beliefs, desires, and intentions to others, plays a foundational role in human communication, language acquisition, and social interaction. The rise of generative models capable of open-ended dialogue raises important questions about whether these systems truly use ToM-like reasoning or simply mimic the surface-level behaviors associated with it, such as predicting others' thoughts or intentions. While recent work suggests that large language models (LLMs) exhibit some competencies associated with ToM, such abilities often emerge from statistical pattern recognition rather than grounded reasoning about mental states. In this commentary, we critically examine the assumption that generative agents "understand" or "reason" about others' minds during communication. We highlight the distinction between genuine ToM capacities and the appearance of ToM arising from training on massive datasets. Drawing on insights from both cognitive science and machine learning, we argue that many current benchmarks and evaluations fail to capture the nuance of ToM as a developmental and socially grounded process. We also examine the risks of attributing human-like understanding to generative agents, particularly in socially sensitive settings where the illusion of comprehension may lead to overtrust or potential manipulation. In alignment with the ACS focus, this commentary urges the community to rethink how ToM is conceptualized, implemented, and evaluated in communicative artificial intelligence (AI) agents. We advocate for interdisciplinary approaches that go beyond behavioral proxies, emphasize developmental insights, and consider the broader social implications of deploying agents that appear to "know what we mean" even when they do not.

1. Introduction

Communication extends beyond the mere exchange of signals; it fundamentally relies on the ability to understand others' beliefs, intentions, and knowledge, capacities collectively known as Theory of Mind (ToM). ToM is defined as the capacity to attribute mental states such as beliefs, desires, and intentions to oneself and others (Wellman, 1990), while recognizing that these mental states can differ between individuals. Humans utilize ToM to interpret implicit meanings, anticipate others' behavior, and coordinate effectively in complex social interactions. This cognitive ability enables individuals to interpret, predict, and explain the actions of others and plays a crucial role across a wide range of communicative contexts.

With the increasing integration of artificial intelligence (AI) systems into diverse domains, these technologies are frequently engaging with humans in interactive settings. Examples include chatbots in e-commerce platforms and virtual assistants in educational and healthcare environments.



This work is licensed under a Creative Commons Attribution International 4.0 License.

The emergence of generative AI (GenAI) agents has further accelerated the prevalence of such interactions, driving a growing interest in evaluating and enhancing AI systems’ social cognitive capabilities and ToM. For instance, Sap et al. (2022) conducted a comprehensive evaluation of GenAI such as GPT-3 and GPT-4 on ToM benchmarks like SocialIQa (Sap et al., 2019) and ToMi (Le et al., 2020), which assess understanding of social interactions and mental states conveyed through natural language. Despite their advanced linguistic abilities, these models attained accuracy rates of only 55% and 60%, respectively, substantially below human performance. While these systems hold considerable promise, their deployment in human-centered contexts introduces unique challenges, including concerns regarding the nature of learned knowledge from user data and the opacity of their decision-making processes.

In this commentary, we begin by examining the role of ToM across various communication scenarios involving humans, AI systems, and non-verbal contexts. We then focus on the critical importance of ToM within GenAI systems designed for human-centric applications. Finally, we present a call to action centered on three key areas to advance ToM capabilities in GenAI communicative agents: (1) evaluation and measuring of ToM in GenAI systems, (2) design and modeling of ToM in GenAI systems, and (3) ethical and social alignment. The first area highlights the need for robust methods to measure and evaluate ToM abilities in GenAI systems, identifying existing gaps. The second addresses design considerations essential for developing GenAI agents with effective ToM tailored to human interaction. The third emphasizes the importance of ensuring that ToM-enabled AI systems operate safely, transparently, and respectfully within social contexts, addressing potential risks such as manipulation, privacy violations, and erosion of user trust.

2. Background

ToM is essential for effective communication and social coordination. In this work, we explore the role of ToM across four distinct communication settings: human-to-human, human-to-AI, AI-to-AI, and non-verbal communication. Figure 1 presents an overview of these settings and illustrates how ToM operates within each context.

2.1 Human–Human Communication and ToM

Human communication fundamentally depends on ToM to interpret implicit meanings, disambiguate messages, and predict others’ behavior. The ability to reason about the mental states of others facilitates cooperative interactions, and enables complex social exchanges. ToM has been extensively investigated within cognitive science and linguistics, with significant insights drawn from developmental psychology. For example, research in child development has elucidated the progressive acquisition of ToM abilities (Beaudoin et al., 2020). Moreover, proficiency in ToM has been associated with positive social outcomes, including stronger interpersonal relationships and enhanced social competence (Liddle & Nettle, 2006; Peterson & Siegal, 2002; Fink et al., 2015). Conversely, deficits in ToM have been documented among individuals on the autism spectrum, contributing to social communication challenges characteristic of this population (Baron-Cohen, 2002, 2000).

2.2 Human–AI Communication and ToM

As AI systems become increasingly integrated into daily life, the importance of incorporating ToM capabilities grows correspondingly. Examples of human-AI communication platforms include voice assistants such as Siri and Alexa, as well as virtual assistants deployed in educational and e-commerce settings. These interactions pose unique challenges related to ToM, as users tend to anthropomorphize AI agents and expect them to understand human goals, beliefs, and contextual subtleties. However, many current AI systems lack explicit models of human mental states, which constrains their adaptability and may result in miscommunication or suboptimal user experiences. Enhancing AI systems with ToM-like functionalities can improve their interpretability and responsiveness, foster greater user trust, and enable more effective and fluid collaboration between humans and machines.

2.3 AI–AI Communication and ToM

As AI agents increasingly operate in multi-agent environments, the ability to model other agents’ beliefs and intentions becomes critical. ToM facilitates coordination, negotiation, and competition among agents. SOTOPIA, developed by researchers at Carnegie Mellon University, represents a significant advancement in the integration of ToM within AI systems (Zhou et al., 2024). SOTOPIA is a system in which AI agents can be simulated to possess distinct personas. Multiple AI agents communicate with one another, mimicking human social interactions. During communication, the AI agents employ ToM to dynamically model and reason about the beliefs, intentions, and knowledge of other agents, thereby enabling more sophisticated and context-aware interactions.

2.4 Non-Verbal Communication and ToM

Research on non-human animals provides compelling evidence that non-verbal ToM exists across a range of species. Behaviors such as gaze following, perspective taking, and deception suggest that ToM can emerge without reliance on language, developing instead through embodied experience and contextual inference. Animals not only learn from their interactions but also demonstrate the ability to generalize learned knowledge to novel situations. One widely used framework for studying this phenomenon is the *competitive feeding paradigm* in psychology (Hare et al., 2000), which examines how animals reason about an opponent’s visual perspective, as well as their true or false beliefs about the environment, as these evolve over a sequence of events. Building on *competitive feeding paradigm*, Michelson et al. (2024) introduced a benchmark to assess non-verbal ToM in AI agents. Non-verbal ToM has also been studied in competitive game setting where each player has to infer the opponent’s intention and beliefs from the opponent’s action (Riemer et al.; Yao et al., 2025).

3. GenAI in Communicating Agents and ToM

Recent progress in GenAI, particularly large language models (LLMs), presents new perspectives on ToM in artificial agents. LLMs generate coherent and context-sensitive responses that can simulate some aspects of mindreading. However, these models typically lack grounded representations of

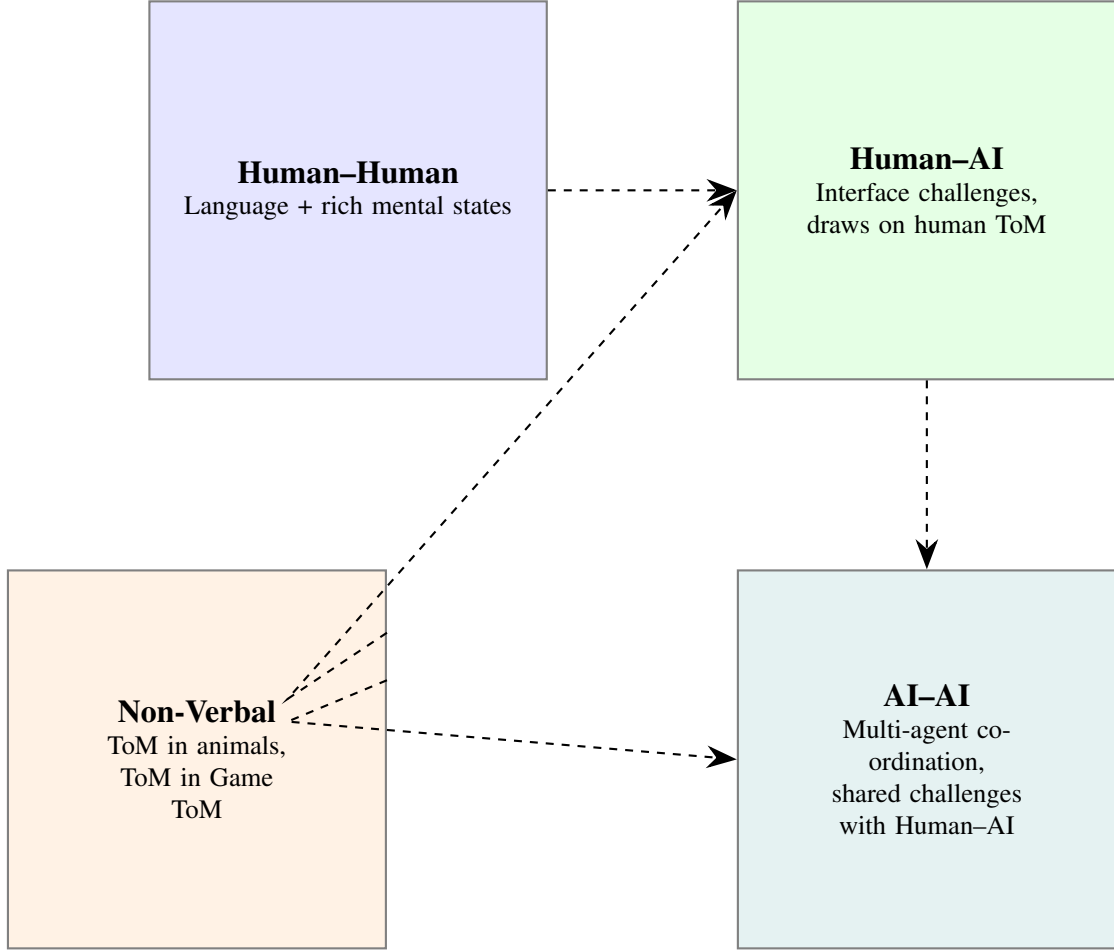


Figure 1: Conceptual overview of ToM across four communication contexts: human-human ToM, characterized by rich language and recursive mental state reasoning. Human-AI and AI-AI ToM reflect shared challenges in modeling beliefs and intentions for effective interaction. Non-verbal animal ToM, while distinct and non-linguistic, provides inspiration for embodied and non-verbal mental state inference in artificial agents. Arrows indicate key influences and relationships among these domains.

beliefs or intentions, which can lead to inconsistent or misleading outputs. Incorporating explicit mental state reasoning into generative models remains an open research challenge essential for advancing trustworthy and socially competent AI.

This section discusses applications of GenAI in human-AI and AI-AI communication, emphasizing the essential role of ToM in improving these interactions. Table 1 summarizes key aspects of these communication scenarios with respect to ToM.

Table 1: Comparison of Human–AI and AI–AI Communication Scenarios with Respect to Theory of Mind.

Aspect	Human–AI Communication	AI–AI Communication
Application Examples	Intelligent tutoring, voice assistants, assistive technologies	Multi-agent coordination, social simulations (e.g., SOTOPIA), autonomous vehicles
Communication Goals	Personalization, user understanding, effective assistance	Cooperation, negotiation, social reasoning, joint decision-making
Role of ToM	Model user beliefs, intentions, knowledge gaps	Model other agents’ mental states to predict behaviors and adapt strategies
Benefits of ToM	Improved interaction quality, trust, contextual appropriateness	Enhanced coordination, anticipation of others’ actions, socially intelligent behavior
Key Challenges	Implicit mental states, ambiguous user signals, risk of misinterpretation	Computational complexity, recursive reasoning, partial observability, ethical concerns
Current Limitations	Lack of explicit mental state models, inconsistency, bias propagation	Scalability issues, simplified models, uncertainty handling, ethical alignment

3.1 Challenges of Integrating ToM: Human–AI Communication Scenarios

GenAI supports a variety of human–AI communication applications, including intelligent tutoring systems, voice assistants, and assistive technologies for people with disabilities. These systems benefit from GenAI’s ability to interpret user inputs and generate contextually appropriate responses. However, to achieve effective and trustworthy communication, agents must model users’ beliefs, intentions, and knowledge gaps. Without ToM, agents may generate linguistically correct yet contextually inappropriate responses, leading to misunderstandings and reduced user trust.

3.2 Challenges of Integrating ToM: AI–AI Communication Scenarios

In multi-agent systems, communication among AI agents relies heavily on ToM to support coordination, negotiation, and social reasoning. For example, the SOTOPIA framework from Carnegie Mellon University simulates social environments where agents model others’ beliefs and intentions to anticipate actions and adapt strategies. This enables agents to engage in more socially intelligent and cooperative behaviors. Despite its promise, this approach faces challenges including computational complexity, simplified assumptions, and ethical concerns regarding mental state manipulation.

Table 1 outlines the applications, goals, benefits, and challenges of ToM integration in human–AI and AI–AI communication, highlighting areas that require further research and development.

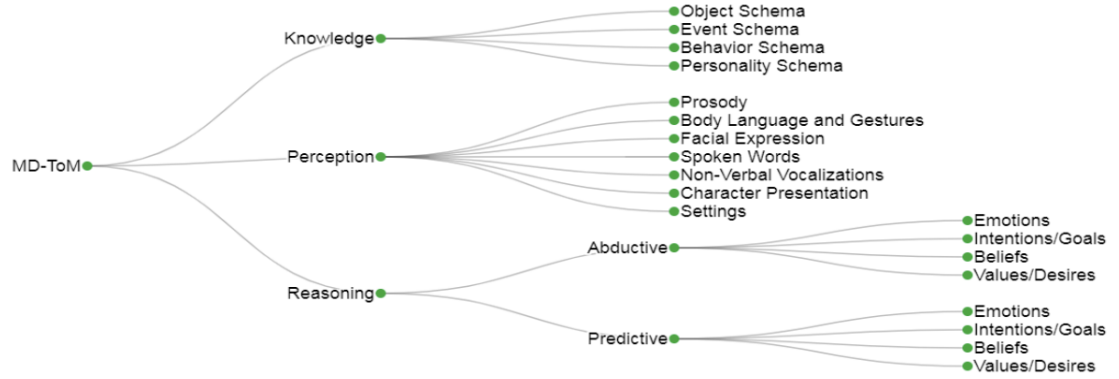


Figure 2: Multi-dimensional ToM, High Level Overview: 3-Component Model

4. A Call for Action: Enhancing GenAI Communication with ToM

This commentary paper focuses on the integration of ToM into generative AI systems designed for communication. We target two core scenarios: human–AI and AI–AI interaction. Both involve agents that must reason about the knowledge, beliefs, and intentions of others to communicate effectively. To advance this area, we propose three key areas to advance ToM capabilities in GenAI communicative agents: (1) evaluation and measuring of ToM in GenAI systems, (2) design and modeling of ToM in GenAI systems, and (3) ethical and social alignment.

4.1 Evaluation and Measuring ToM in GenAI Systems

We discuss these in following subsections.

4.1.1 Designing Diverse Cognitive Tests for ToM.

GenAI systems and LLMs are trained on vast amount of data from internet. As machine learning memorize patterns from training data, it is plausible that LLMs may have encountered and memorized widely available cognitive tests during pretraining. A prominent example is the Sally-Anne false belief test, a standard assessment in developmental psychology that evaluates ToM by testing whether an individual can understand that others may hold beliefs that differ from reality (Frith, 2004). In this task, a character named Sally places an object (e.g., a marble) in a basket and leaves the scene. While she is away, another character, Anne, moves the object to a different location (e.g., a box). The Sally-Anne false belief test assesses whether a person can correctly infer that Sally will look for the object where she originally left it, reflecting an understanding of her false belief. LLMs tend to perform remarkably well on the Sally-Anne and its close variants, such as those in which the characters’ names are altered, which suggests memorization or pattern recognition rather than genuine cognitive reasoning. However, when evaluated on alternative ToM tasks that differ substantially from the false belief paradigm, LLM performance declines considerably, as reported by Ullman (2023). To more comprehensively assess the ToM capabilities of LLMs, Shapira et al.

(2024) introduced a diverse set of six novel ToM tasks. Their findings indicate that LLM performance varies significantly across these tasks, suggesting that current models may lack a robust, generalizable understanding of ToM.

We advocate that evaluating GenAI systems for ToM capabilities should involve a diverse set of task types rather than relying solely on standard or widely known cognitive assessments. In particular, commonly used psychology tests for ToM may not be reliable to test ToM capabilities in LLMs, as LLMs are capable of memorizing these tasks and successfully solving even modified variants through pattern recognition rather than true mental state inference.

4.1.2 *Measuring ToM Across Multiple Dimensions in AI Systems.*

A key limitation of designing ToM abilities in GenAI systems and LLM is to assess only one or two isolated dimensions of ToM. However, ToM is a multi-dimensional in nature. Stack et al. (2022) introduced a comprehensive three-component framework for ToM that is applicable to both *humans* and *AI systems*, as illustrated in Figure 2. This framework identifies perception, knowledge, and belief as the core components of ToM, each encompassing several sub-components. The authors emphasize that effectively modeling ToM requires the integration of diverse perceptual inputs, knowledge structures, and reasoning mechanisms. Similarly, Ma et al. (2023) proposed a multidimensional framework specifically for *AI systems*, drawing inspiration from the psychological model ATOMS developed by Beaudoin et al. (2020). ATOMS categorizes mental states into seven distinct types based on a meta-analysis of ToM research in child development, offering a broader foundation for machine-level ToM evaluation.

When developing GenAI systems with ToM capabilities, we emphasize the importance of adopting a holistic evaluation framework that captures the multifaceted nature of ToM, rather than limiting assessment to one or two isolated dimensions.

4.2 ToM with Human Centric and Cognitive Science-Inspired Model Design

The design of ToM mechanisms in GenAI systems can benefit from cognitive science theories about how humans acquire and apply mental state reasoning. To better facilitate cognitive theory-based approach, we advocate following two approaches.

4.2.1 *Mutual ToM.*

This approach reflects the dynamics of human communication, where individuals are typically aware of each other’s beliefs, intentions, and mental states. The concept of mutual ToM extends this idea to human-AI interaction, advocating for a bidirectional understanding between the human and the AI system. It encompasses two key dimensions: the human’s ToM of the AI (i.e., the human’s understanding of how the AI operates and interprets behavior) and the AI’s ToM of the human (i.e., the AI’s ability to infer and respond to human beliefs, goals, and mental states in context). The goal is to foster effective and adaptive communication where both parties—human and AI—possess some awareness of the other’s cognitive processes.

Initial work in this direction has already emerged. For example, Wang et al. (2021) implemented mutual ToM in the context of a conversational agent: a virtual teaching assistant (TA) named Jill Watson, deployed in an online class discussion forum. Jill Watson was designed to assist students by leveraging ToM principles. The study investigated how students perceived and interacted with the virtual TA, thereby exploring the human-to-AI dimension of mutual ToM. A related concept was proposed by Bara et al. (2021), who emphasized the importance of establishing common ground in communication as a prerequisite for effective collaboration. Their research focused on a 3D game-based collaborative task that required participants to engage in mutual reasoning about each other’s mental states, thereby highlighting the practical relevance of mutual ToM in shared environments.

In the development of GenAI systems with ToM capabilities, we advocate for a unified framework centered on mutual understanding, which includes both the *human’s ToM of the AI* and the *AI’s ToM of the human*. This bidirectional perspective emphasizes that effective interaction depends not only on the AI’s ability to infer human mental states, but also on the human’s understanding of the AI’s behavior, reasoning processes, and limitations.

4.2.2 Recursive ToM.

Human beliefs and perceptions about others are inherently dynamic rather than static. These mental states continuously evolve as individuals observe and interpret the behavior of those around them. This ongoing process of updating beliefs about another person’s thoughts and intentions is referred to as recursive ToM (recursive ToM). Building on this concept, Kleiman-Weiner et al. (2025) introduced the Bayesian Reciprocator, a computational framework designed to capture the dynamic evolution of human cooperation. This approach leverages rational ToM inferences within a Bayesian framework to model how individuals iteratively update their beliefs and expectations based on observed social interactions, thereby offering a principled method for understanding and predicting cooperative behavior. Wang’s recent work on mutual ToM also entails recursive ToM Wang & Goel (2022), which enables continuously refine users’ understanding of each others’ minds through behavioral and verbal feedback.

When designing an GenAI systems enabling ToM, we advocate for a dynamic and recursive ToM to allow the system to continuously update its beliefs and mental state representations.

4.3 Social Impact and Responsible Deployment of ToM-Enabled Agents

4.3.1 Responsible and robust GenAI Deployment.

As GenAI systems increasingly demonstrate ToM-like behaviors, it becomes imperative to carefully assess their broader impact on users and society. Such agents must tailor their communication strategies to accommodate the diverse needs of different user groups, including children, older adults, and individuals with disabilities, to ensure appropriate and ethical interactions. Furthermore, the introduction of ToM capabilities raises important ethical concerns related to manipulation, privacy, and the potential for misaligned behavior. For example, hallucinated beliefs generated by GenAI sys-

tems can lead to harmful or misleading outputs, posing significant risks. Mitigating these risks necessitates the implementation of robust safeguards that constrain model behavior and maintain reliability under distributional shifts.

When designing a GenAI systems enabling ToM, we advocate for considering the ethical and user privacy aspects of ToM to ensure positive social impact by such systems.

4.3.2 Integrating Logical AI.

The deployment of ToM-enabled systems in real-world contexts demands enhanced interpretability of their decision-making processes to foster transparency and trust. Researchers have begun to address these challenges; for instance, (Sclar et al., 2023) proposed an explicit symbolic representation to elucidate the reasoning and ToM capabilities of LLMs. Researchers can also integrate logical AI models as knowledge base to improve GenAI’s reasoning capabilities. An example is the seminal work of Allen and Perrault Allen & Perrault (1980) on model of cooperative behavior applied to natural language understanding system. A relative recent work by Gabaldon and colleagues ? describes how to frame a task task-oriented dialogue in which the interacting agents adopt and pursue a shared goal.

Integrating logical AI models as knowledge base can enhance GenAI systems’ interpretability and reasoning capabilities.

5. Conclusion

In this commentary, we examine the role of ToM in communication agents within the context of the GenAI era. To this end, we analyze the significance of ToM across four distinct communication settings: human-human, human-AI, AI-AI, and non-verbal interactions. Given the prominence of generative AI in communication, our primary focus lies on Human-AI and AI-AI interactions to evaluate ToM capabilities. Building on existing research, we put forward a call to action for the design of AI systems that incorporate robust ToM functionalities to enhance communication in the GenAI era.

References

- Allen, J. F., & Perrault, C. R. (1980). Analyzing intention in utterances. *Artificial intelligence*, 15, 143–178.
- Bara, C.-P., Sky, C.-W., & Chai, J. (2021). Mindcraft: Theory of mind modeling for situated dialogue in collaborative tasks. *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing* (pp. 1112–1125).
- Baron-Cohen, S. (2000). Theory of mind and autism: A fifteen year review. In S. Baron-Cohen, H. Tager-Flusberg, & D. J. Cohen (Eds.), *Understanding other minds: Perspectives from developmental cognitive neuroscience (2nd ed.)*, 3–20. New York, NY, US: Oxford University Press.

- Baron-Cohen, S. (2002). The extreme male brain theory of autism. *Trends in Cognitive Sciences*, 6, 248–254. From <http://www.sciencedirect.com/science/article/pii/S1364661302019046>.
- Beaudoin, C., Leblanc, É., Gagner, C., & Beauchamp, M. H. (2020). Systematic review and inventory of theory of mind measures for young children. *Frontiers in psychology*, 10, 2905.
- Fink, E., Begeer, S., Peterson, C. C., Slaughter, V., & de Rosnay, M. (2015). Friendlessness and theory of mind: A prospective longitudinal study. *British Journal of Developmental Psychology*, 33, 1–17.
- Frith, C. D. (2004). Schizophrenia and theory of mind. *Psychological medicine*, 34, 385–389.
- Hare, B., Call, J., Agnetta, B., & Tomasello, M. (2000). Chimpanzees know what conspecifics do and do not see. *Animal Behaviour*, 59, 771–785.
- Kleiman-Weiner, M., Vientós, A., Rand, D. G., & Tenenbaum, J. B. (2025). Evolving general cooperation with a bayesian theory of mind. *Proceedings of the National Academy of Sciences*, 122, e2400993122.
- Le, T., Hajishirzi, H., & Zettlemoyer, L. (2020). ToMi: Theory of mind reasoning in natural language processing. *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)* (pp. 2475–2487). Association for Computational Linguistics.
- Liddle, B., & Nettle, D. (2006). Higher-order theory of mind and social competence in school-age children. *Journal of Cultural and Evolutionary Psychology*, 4, 231–244.
- Ma, Z., Sansom, J., Peng, R., & Chai, J. (2023). Towards a holistic landscape of situated theory of mind in large language models. *Findings of the Association for Computational Linguistics: EMNLP 2023* (pp. 1011–1031). Singapore: Association for Computational Linguistics. From <https://aclanthology.org/2023.findings-emnlp.72/>.
- Michelson, J., Sanyal, D., Ainooson, J., Farhana, E., & Kunda, M. (2024). Standoff: benchmarking representation learning for nonverbal theory of mind tasks. *2024 IEEE International Conference on Development and Learning (ICDL)* (pp. 1–6). IEEE.
- Peterson, C. C., & Siegal, M. (2002). Mindreading and moral awareness in popular and rejected preschoolers. *British Journal of Developmental Psychology*, 20, 205–224.
- Riemer, M., Ashktorab, Z., Bouneffouf, D., Das, P., Liu, M., Weisz, J. D., & Campbell, M. (????). Position: Theory of mind benchmarks are broken for large language models. *Forty-second International Conference on Machine Learning Position Paper Track*.
- Sap, M., Le Bras, R., Allaway, E., Bhagavatula, C., Lourie, N., Rashkin, H., Roof, B., Smith, N. A., & Choi, Y. (2019). Socialiqa: Commonsense reasoning about social interactions. *Proceedings of the AAAI Conference on Artificial Intelligence* (pp. 4647–4654). AAAI.
- Sap, M., Le Bras, R., Fried, D., & Choi, Y. (2022). Neural theory-of-mind? on the limits of social intelligence in large lms. *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing* (pp. 3762–3780).
- Sclar, M., Kumar, S., West, P., Suhr, A., Choi, Y., & Tsvetkov, Y. (2023). Minding language models’ (lack of) theory of mind: A plug-and-play multi-character belief tracker. *Proceedings*

- of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)* (pp. 13960–13980). Toronto, Canada: Association for Computational Linguistics. From <https://aclanthology.org/2023.acl-long.780/>.
- Shapira, N., Levy, M., Alavi, S. H., Zhou, X., Choi, Y., Goldberg, Y., Sap, M., & Shwartz, V. (2024). Clever hans or neural theory of mind? stress testing social reasoning in large language models. *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)* (pp. 2257–2273).
- Stack, C. H., Farhana, E., Shen, X., Zhao, S., & Maliakal, A. (2022). Framework for a multi-dimensional test of theory of mind for humans and ai systems. *The Tenth Annual Conference on Advances in Cognitive Systems*.
- Ullman, T. (2023). Large language models fail on trivial alterations to theory-of-mind tasks. *arXiv preprint arXiv:2302.08399*.
- Wang, Q., & Goel, A. K. (2022). Mutual theory of mind for human-ai communication. *IJCAI Workshop on Communication in Human-AI Interaction (CHAI)*.
- Wang, Q., Saha, K., Gregori, E., Joyner, D., & Goel, A. (2021). Towards mutual theory of mind in human-ai interaction: How language reflects what students perceive about a virtual teaching assistant. *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. New York, NY, USA: Association for Computing Machinery. From <https://doi.org/10.1145/3411764.3445645>.
- Wellman, H. M. (1990). *The child’s theory of mind*. MIT Press.
- Yao, J., Wang, K., Hsieh, R., Zhou, H., Zou, T., Cheng, Z., Wang, Z., & Viswanath, P. (2025). Spin-bench: How well do llms plan strategically and reason socially? *arXiv preprint arXiv:2503.12349*.
- Zhou, X., et al. (2024). Sotopia: Interactive evaluation for social intelligence in language agents. *The Twelfth International Conference on Learning Representations*.