

Applying the Theory of Make-Believe to Human-Robot Interaction

Matthew RUEBEN^{a,1}, Eitan ROTHBERG^b, and Maja J. MATARIĆ^a

^a*Dept. of Computer Science, University of Southern California, Los Angeles, California, USA*

^b*The Ohio State University, Columbus, Ohio, USA*

Abstract. People often make ascriptions that they believe to be literally false. A robot, for example, may be treated as if it were a dog, or as if it had certain intentions, emotions, or personality traits. How can one do this while also believing that robots cannot really have such traits? In this paper we explore how Kendall Walton’s theory of make-believe might account for this apparent paradox. We propose several extensions to Walton’s theory, some implications for how we make attributions and use mental models, and an informal account of human-robot interaction from the human’s perspective.

Keywords. Human-robot interaction, mental state attribution, theory of mind, intentional stance, anthropomorphism, make-believe, pretence theory, prop theory, fictionalism

1. Introduction

One way that the human mind makes sense of the world is by imposing interpretations upon incoming sensory information [1]. Included in this process is the assignment of causes or explanations to events—scholars call this *attribution* [2]. Surprisingly, however, people often attribute events to causes that they know cannot be correct. Beliefs, intentions, desires, and even personhood are ascribed to inanimate objects to make sense of their behavior [3]. We describe a car as deciding to break down, the weather as unable to make up its mind, or a virus as aggressive. This is not merely a feature of language, as people do actually get angry at cars when they do not start, and this still happens even though most people would admit that they know cars do not really have intentions [4].

People have also been observed to make such attributions about the behavior of robots [5,6]. Robots are often designed to mimic and resemble humans or animals. This can include using human-readable social cues, from a face that appears to make eye contact to gestures and posture. By their appearance and behaviors robots can suggest animacy or sentience to human onlookers—they can appear to think, learn, emote, and attend—and methods for making this more convincing are improving quickly. But today’s robots are not conscious². They lack awareness, knowledge, understanding³, and

¹ Corresponding Author, e-mail: mrueben@usc.edu.

² See Haikonen [7] for a discussion of what machine consciousness would entail, whether it is possible, and how to test for it.

³ In fact, Zlotowski et al. have observed that we are so quick to anthropomorphize robots that we should not be trusted to judge whether they actually become intelligent in the future [5].

feeling. Yet even people who claim to believe this nonetheless respond to robots in social ways [4]. As the role of robots in modern societies continues to grow, it is important that we understand this puzzling behavior.

In this paper we apply Kendall Walton's theory of make-believe to anthropomorphism and related phenomena by which humans make attributions that they know to be false about robots. Walton proposes that much human language and behavior involves thinking about, referring to, or even entering fictional worlds [8]. After summarizing Walton's theory we extend it in several ways and then apply it to the context of human-robot interaction. Our extensions consist of a discussion of communication about fiction during an interaction (Section 4.1) and the addition of our terms "faking" and "imitating" to Walton's idea of participating in a make-believe (Section 4.2). We then explore the role of the imagination in attribution theory, the intentional stance, and mental models (Section **Error! Reference source not found.**) and present an informal account of human-robot interaction (Section 4.4). We discuss some alternative accounts of anthropomorphism in human-robot interaction (Section **Error! Reference source not found.**) before closing with some research implications of our Waltonian account.

2. Background

2.1. Attribution Theory and "Computers Are Social Actors"

Humans sometimes attribute phenomena to causes that they know cannot be the correct ones. For example, people have been observed to explain the movements of mere geometric shapes in terms of agency, intentionality, and relationships [9]. Humans' preference for a predictive but impossible explanation over no explanation at all is one of the psychological puzzles addressed by theories of attribution from early on through recent times [2,10]. Particularly interesting to human-robot interaction researchers is the use of social cognitive processes when working with technology—people often behave as if "Computers Are Social Actors (CASA)" [4,11].

2.2. Attribution of Robot Behaviors to Mental States

Researchers from several fields have studied how people explain the actions of robots as if they had minds (via "mentalizing"). Thellman and Ziemke [12] have defined several of the relevant terms. **Folk psychology** "broadly encompasses all mind-related theories that people have about themselves and others," such as, "people who are tired are generally irritable," and even, "robots do not have minds." **Theory of mind (ToM)** is narrower: "the ability to see other entities as intentional agents, whose behavior is influenced by states, beliefs, desires, etc." [13]. Children develop theories of mind gradually, from following eye gaze to modelling different types of mental states [13]. The **intentional stance** is what Dennett calls the "craft" or practice of folk psychology [12]. Here "intentional" has a technical meaning: "being 'about' or 'directed at' something as an object." This does include the English "intending," but also believing, desiring, hoping, fearing, and having a hunch. Taking the intentional stance is using intentional states like these as an interpretive framework to predict behavior.

The intentional stance is the last of Dennett's three strategies or stances that people use to predict what will happen next [3]. The other two are the **physical stance**, which

applies the laws of physics to the system's physical makeup, and the **design stance**, which reasons that the system will behave as it was designed to behave.

2.3. The Problem of Anthropomorphism

“Anthropomorphism describes the tendency to imbue the real or imagined behavior of nonhuman agents with humanlike characteristics, motivations, intentions, or emotions” [14]. Surprisingly, humans can do this while maintaining the belief that the entity being anthropomorphized actually lacks the things being ascribed to it [4]. Our main focus for this paper is on the question: how do people do these seemingly contradictory mental operations?

Zlotowski et al. [15] address this question via a dual-system approach—i.e., they argue that anthropomorphism can come from two different types of mental processes. “Type 1” processes are fast, intuitive, and automatic, whereas “Type 2” processes are reflective, deliberate, and require more cognitive resources. This approach accounts for, e.g., the observations of Nass & Moon [4] that study participants “interact with technology in a social way” [15] but do not ascribe humanlike characteristics to that technology in an interview afterwards. Zlotowski et al. argue that social responses to robots are due to anthropomorphism produced by automatic Type 1 processes, whereas reflective Type 2 processes are responsible for any lack of anthropomorphization when asked to explicitly state one's beliefs about the robot [15].

Several researchers have suggested that the ability to “hold both humanistic and mechanistic conceptions of the robot simultaneously” [6] seems similar to using metaphors or even make-believe. Zlotowski et al. [5] ask whether the same sort of make-believe or imagination that is used to say, “Sherlock Holmes lives at 221B Baker Street” is also being used to say that a robot, “has a ‘happy’ expression”. In the discussion of their empirical results Fussell et al. [6] wonder if the participants' judgments “reflect beliefs that robots literally possess feelings, attitudes and personality traits or whether they instead are based on metaphoric extension.” Another hint in this direction is the link that Severson & Woodard make between children's role play and anthropomorphism: “We propose that the imaginative process of simulating and projecting internal states is common to both role play and anthropomorphism” [16]. Anthropomorphism of robots is a poorly understood process, but make-believe and the use of metaphors are familiar human experiences that have been studied by scholars from many disciplines and across centuries of history. This suggests the potential fruitfulness of any comparison that can be made between anthropomorphism and other ways that humans use fictions to interact with the world.

3. Walton's Theory of Make-Believe

Scholarly work has investigated how people experience, talk about, and create stories that contain fictional elements. The scholars whose work we read came from diverse fields: psycholinguistics, the philosophy of art, aesthetics, art criticism, and cognitive science [17–21]. For this paper the account of fictional experience is drawn mainly from the work of philosopher Kendall Walton [8,18,22].

Walton says that representational art (e.g., a novel, story, painting, sculpture, stage play, or film) creates a **fictional world** [18]. For example, there is a world of Swan Lake—one can imagine that the events being portrayed in the ballet were true. Of course,

audience members might imagine Swan Lake differently, thereby creating their own versions of the fictional world even though they are all watching the same performance. Nevertheless, Walton argues that it is still reasonable to refer to a single world of Swan Lake comprised of all the details that are fixed by the story of the ballet itself—e.g., that the princess is turned into a swan and not some other bird. Inasmuch as a robot is what Walton calls a **prop**, robots "prompt and prescribe certain imaginings" to which people may choose to add original details.

Walton defines **make-believe** as an activity wherein participants inhabit a fictional world [18]. When people become "caught up in the story" or "emotionally involved" in the fictional world," it is because of their psychological **participation** in a make-believe [8,18]. Scholars have proposed different explanations for this apparent "paradox of emotional response to fiction"—i.e., that the person experiencing the emotions also knows that the fiction is not really true. An audience member may feel something like terror at a horror movie, including an elevated heart rate and cold sweat, despite knowing "perfectly well that the [monster] is not real and that he is in no danger" [22]. Novitz [19] reviews four accounts of how this might happen, but the important point is that participation in a make-believe can affect us deeply—physiologically as well as intellectually.

Participating in a make-believe also enables the participant to act within the fictional world through the fictional character they are playing, changing the state of affairs therein [22]. Walton explains that, "what happens in the fictional world is just what is portrayed as happening," so even an audience member can jump on stage, enter the play as a new or existing character, and act out a new twist in the story [22]. This impromptu participant would thereby make it **fictional** that the new twist happened. Walton uses this attribute "fictional" for statements about the contents of fictional worlds just like "true" is typically used for statements about the real world [22]. He stipulates that fictional states of affairs are imagined, not believed: "Imagining aims at the fictional as belief aims at the true." [18]

Walton extends his theory of make-believe to also describe our pervasive use of metaphors [8]. His examples include our metaphorical references to Italy as a boot, to male and female connectors for pipes, and to the saddle of a mountain. He argues that metaphors such as these are best understood as referring to, implying, introducing, or reminding us about games of make-believe—e.g., in which Italy is a boot or a giant horse's saddle connects the peaks of two mountains⁴. Even so, in most cases we do not actually play the game of make-believe suggested by each of these metaphors—Walton says, "Our participation is minimal at best." [8] Usually we are just trying to locate the city of Rome on a map or indicate a certain part of the mountain. A mere "recognition of the possibility of make-believe" can suffice when using the metaphor to communicate [8]. In some cases participation may be required, but it can be quite minimal—e.g., to fictionally recognize a man or woman in the icons designating the two different restrooms.

⁴ Using an especially robotics-relevant example Walton notes that using one metaphor can unintentionally lead people to use other metaphors that follow or are related: "A comment that a computer remembers a phone number may prepare me to think of computers behaving so as to make it fictional that they forget things, that they calculate, make decisions, and even lose patience or complain about their handlers or give up on a task." [8]

4. Synthesis and Suggestions

In this section we propose some extensions to Walton's theory of make-believe and discuss their implications for attribution, the intentional stance, and mental models. After that, we propose an informal account of human-robot interaction inspired by Walton's theory.

4.1. *Communicating about the Fictional World*

There are several ways that participants in a make-believe can communicate things about the fictional world to other participants. Sometimes these messages are emitted naturally as a result of full participation in the fictional world—i.e., something said or done by a player while *in character* communicates some things that are fictional of that world. For example, a child pretending to be Peter Pan might gasp and draw a plastic knife—this might be sufficient to communicate that fictionally one of Peter's enemies, such as Captain Hook, is approaching. These in-character behaviors can be insufficient, however, to maintain the necessary amount of common ground about what is fictional in the world of the make-believe. At these times participants often come at least partly out of character—break the fourth wall, so to speak—to give other participants the information they need. For example, it would be important to avoid confusion about who is playing which character: “no, *you're* Captain Hook, and *I'm* Peter Pan!”

People can also signal to others that a make-believe is occurring⁵ and invite them to play along. As above, these meta-communications can happen in or out of character. Breaking character, when it does happen, can be subtle: a sarcastic tone of voice or meaningful glance could suffice. For plays, films, and television shows, the context alone of the screen or stage could be enough of a signal to viewers that the actors are only behaving that way for the purposes of the production.

4.2. *“Behaving As If” without Make-Believing: Faking and Imitating*

Walton explains how someone can “refer to” a fictional world without participating in it—i.e., without becoming “caught up in the story” or “emotionally involved” [8,18]. It seems like one can also *behave as if* inside a fictional world—go through the motions, so to speak—without actually entering it via make-belief. They would calculate or simulate⁶ the appropriate fictional thoughts and actions from *outside* the fictional world instead of assuming a character and allowing them to happen naturally. This should not be called “acting” because professional actors are often *in character*—i.e., make-believing.

We will call this either **faking** or **imitating**—the distinction is as follows. Fakers avoid emitting signals that their actions are not indicative of their real beliefs. Instead, they try to send the same signals as if they were really participating in the make-believe. This can be used to fool someone maliciously, but also more innocuously, as when going through the motions of participating in a child's make-believe to humor them.

When imitating, on the other hand, it is acceptable or even desired for others to detect that the behavior is not sincere. The imitator might send intentional signals to make this clear. They might act out participating in someone else's game, but with

⁵ Cf. “framing” in Tannen [23].

⁶ See the controversy between theory-theorists and simulation-theorists in theory of mind [13].

obvious sarcasm. They might pantomime a skillful motion to a pupil, exaggerating the important parts. They might choose the context of a stage or already-established game. Imitating and faking are two ways one can “behave as if” without actually participating in a make-believe.

4.3. *Attributions and Mental Models can be Imagined instead of Believed*

Are two contradictory beliefs being held simultaneously when, for example, a robot’s behavior is attributed to mental states by someone who also confesses that robots lack minds? Although it could be possible that this person is holding two contradictory beliefs, Walton’s theory of make-believe offers an alternative to this. Instead of having two contradictory beliefs—one that the robot has, and another that the robot does not have mental states—they instead *imagine* that the robot has mental states while also *believing* that it does not.

In this view, attribution is a process conducted by the imagination. When we attribute the robot’s behavior to mental states, for example, we consider a fictional world in which the robot has mental states that would make sense of its actions. In terms of Dennett’s stances [3], perhaps we are taking the intentional stance, but with respect to a fictional world in which the robot has intentionality. Note that it is not necessary to participate in a make-believe to make an attribution; the fictional world is used to explain the robot’s behavior, which it can do even when considered from the outside.

Using mental models could also involve the imagination. For instance, suppose someone uses their mental model of a dog to help them learn to play with the Aibo robot. This person is considering a fictional world in which the robot *is* a dog, and then using this fiction to predict what the robot might do next, or how a human should respond. As with attribution, the mental model can be useful whether these answers are obtained by participating in the make-believe or by calculating them from outside the fictional world.

4.4. *An Informal Account of Human-Robot Interaction from the Human’s Perspective*

The following informal account describes human-robot interaction in terms of Walton’s theory of make-believe, including our concepts of faking and imitating. Alternative accounts are considered in the next section.

1. Robots are often “appearing or behaving as if”, but lack intentions and awareness. For example, an anthropomorphic or zoomorphic robot appears as if it were a human or animal, and may also exhibit behaviors to improve the resemblance. The engineers and designers that created the robot probably intended for people to ascribe to it characteristics such as intelligence or animacy that it lacks in fact. But the robot does not have intentions or awareness; it is not intending to “appear or behave as if”, nor is it aware of it when it does so.

The robot is also not *make-believing*, or even faking or imitating a make-belief, that it has these characteristics. Currently most robot software lacks any implementation of the concept of make-believe; even if there is a representation of the robot’s true capabilities or of the character the robot is supposed to be playing, there is no concept of it assuming a fictional persona or becoming its true (i.e., merely mechanical) self again. Therefore, most robots are not capable of breaking the fourth wall—i.e., signaling apparent awareness⁷ that their behavior is only an act.

⁷ See Haikonen [7] for a discussion of machine consciousness.

2. A fictional world is generated via the robot's appearance and behaviors. The robot's designers may have had different conceptions of what should be fictional in the world generated via the robot, and each person who interacts with the robot may perceive a yet different world than any of those. Each such world is defined by what is fictional of the robot, the humans, and the surrounding environment therein.

The real robot can be considered to be playing a fictional character that is collocated with it in space and time. It is *true* that humans interact with the real robot that is really there; it is simultaneously made *fictional* that the character is there and that those same humans are interacting with it. This is how a fictional world can be involved while real, physical tasks such as cleaning or physical therapy are also being accomplished.

We hypothesize that people who interact with a robot tend to try to figure out what is fictional of *the world that the designers intended* to generate via the robot. They might say, "I see: it's meant to be a dog," even if this supposed intention does not seem to have led to a very convincingly doglike robot. This suggests that an *anthropomorphic* robot be defined as a robot *seemingly intended to make it fictional that* it has "humanlike properties, characteristics, or mental states" [14].

People might guess incorrectly about what the designers intended to make fictional. Usually this is considered a problem with the robot—its behaviors are not properly suggesting the fiction. Another perspective, though, is that the fiction should have been designed to better match the robot's behaviors—its capabilities, affordances, and idiosyncrasies. A robot that moves slowly and needs physical assistance with certain tasks could be elderly; a robot that asks a lot of questions could be a curious child.

3. Each person may or may not feel invited into the fictional world. There are several possible sources of a perceived invitation.

The robot can seem to extend either a verbal or a non-verbal invitation—e.g., by how much it seems to be engaging with each person, expecting them to participate. Future research should explore how the perception of being invited into the robot's fictional world is influenced by the person's understanding that the robot lacks any awareness that it is only acting, that its persona is not real. If we believed this would we feel bad to ignore or decline a robot's invitation to pretend it is, e.g., a playful dog or a bubbly receptionist?

The type of interaction and other contextual factors can also influence whether someone feels obliged to enter the fiction. These could include the roles assumed by the human and the robot, and the task or activity being done. For example, if a robot is leading a group of people in a physical exercise regimen, its position of authority might make people feel pressured to act as if the fiction were true even if the robot does not explicitly encourage this.

Other people could also make someone feel pressured to participate. For example, a person could feel motivated to fit in with a certain group of people who have all seemed to enter the robot's fiction.

4. Each person responds with beliefs, stances, and actions.

4a. What each person believes. It is possible for one to believe that some of the propositions apparently made fictional by the robot—either about itself or about its environment—are actually true. For those about the robot, each person's ability to correctly separate fact from fiction could depend on how much they know about robots and related technologies.

Other people will only *imagine* that these propositions are *fictional* while simultaneously *believing* that what they know about the robot as mere mechanism is *true*. The two mental processes do not necessarily conflict, and people do not contradict

themselves when they switch between acting according to the truth and acting according to the fiction.

4b. Which fictional world each person chooses. People can consider fictional worlds besides the one being generated by the robot. They could conceive of a modified version of that world and then act as if that slightly different set of fictionalities were true. They could also substitute a completely different world. Finally, it is possible that different people will have different perceptions of what exactly is being made fictional via the robot, so there could be multiple different worlds influencing the interaction even if nobody is trying to introduce something new.

4c. How much each person acts as if the fiction were true. Once someone has a fiction in mind it can influence their behavior to different extents. Some people might remain still and silent, just observing. Others might talk as if the fiction were true, or use non-verbal behaviors such as facial expressions and gestures. People might express emotions according to what is fictional at a certain moment. Task execution, including *how* actions are done, could also be influenced.

The extent to which someone is “acting as if” can also change over time. They could do so only occasionally when it is helpful, expedient, or humorous. On the other hand, someone could mostly act according to a fiction but occasionally break the fourth wall by making a comment from the real world. Socially savvy robots should be able to detect when people switch between acting according to a fiction and acting according to reality.

4d. Whether each person “acts as if” by participating in a make-believe or by faking or imitating such participation. As discussed above, make-believing includes becoming “‘emotionally involved’ in the fictional world” and “caught up in the story” [8,18], whereas faked or imitated make-believing only involves going through the motions (Section 4.2). Each person’s behaviors could result from a mix of these three modes—a person might even be able to make-believe certain elements of the fiction but fake or imitate make-belief for the others.

It will be important for certain kinds of robots to be effective at encouraging people to actually participate in a make-believe and not just fake it. Robots could also encourage participation in more modalities, such as speech, gestures, facial expressions, and task-specific actions. Researchers could even explore how someone would react to a robot that accuses them of faking.

4e. People could also just consider the fiction without acting as if it were true. Using, e.g., a metaphor to talk or think about the robot can be brief and require little or no participation [8]. For example, one could refer to the sounds made by a dog-shaped robot as “barks,” “growls,” and “panting.” Walton would call this *prop oriented make-believe* because its purpose is to understand the prop, which is the robot in this case [8].

5. Alternative Accounts of Anthropomorphism in Human-Robot Interaction

Johanna Seibt has done some of the early work on robo-ontology. In several of her papers [24–26], but especially in “Towards an Ontology of Simulated Social Interaction” [24], she argues against accounts of anthropomorphism like Walton’s that appeal to fictionality. For example, Seibt argues that to speak of fictional social actions and interactions is incoherent because a real social (inter)action is always performed in the course of trying to perform a fictional one [24]. She proposes an alternative account in which robots perform *simulated* social actions instead of fictionally performing real social actions. Seibt and her colleagues also make the related suggestion that humans are

not always anthropomorphizing (i.e., humanizing) robots at all but rather sociomorphing them—i.e., ascribing to the robot social capacities that it actually has [27]. These two possibilities would reframe or even eliminate the problem addressed in this paper.

We think Seibt's account of simulated social actions makes sense. We propose, however, that "acting as if" in the sense of making a social interaction fictional is often *comprised of* "acting as if" in the sense of simulating social actions [24]. To act something out—i.e., to make it fictional—requires going through the motions—i.e., simulation. The relationship between the simulation and the fiction might be well described by Clark as the relationship between his base, proximal, and distal layers [28] as discussed below. Whereas Seibt focuses on what sorts of actions the robot is actually doing, we focus on how these actions suggest to people that the robot has capabilities that it does not really have. We apply Walton's theory to describe this perception via a fictional world in which the robot does have those capabilities.

Clark's idea of "layering" in language could also help to explain anthropomorphism. Clark was inspired by Walton (among others) to describe some actions as having multiple layers, one on top of another—for example, playful *teasing* can be described as some statement that is not meant to be taken literally (the bottom layer) but is instead being used as a joke (the top layer) [28]. Clark applies the concept of layering to mechanized speech (e.g., from a robot), which he understands as helping us imagine a virtual speaker in a second layer. By extension, if a robot could be described as implying it is a real social actor, but only in Clark's "distal layer," then this accords with and builds upon our account. Clark also uses layering to describe the use of depictions in communication with his "staging theory": "depictions are physical scenes that people stage for others to use in imagining the scenes they are depicting" [29]. He uses a theater metaphor (e.g., "staging") because scholars also use theater metaphors to describe children's make-believe play.

In a similar vein, Remmers argues in this volume that the perception that robots perform social actions can neither be viewed as mere illusion nor as a full reality; instead, he analogizes them to depictions of real people [30]. He applies the phenomenological theory of image perception by Husserl and Wiesing to argue that a person simultaneously perceives (1) the robot as mere mechanism, (2) the robot performing the social action being depicted, and (3) the conflict between these first two perceptions. Remmers' focus is on the human's *perception* process, whereas we are talking more about cognition and imagination—perhaps our work and his could be reconciled to produce a more complete account.

Musiał proposes that animism towards robots stems from "magical thinking" in which people disregard the limits imposed by logic or the distinction between fiction and reality [31]. Future experimentation should focus on testing the specific claims that distinguish these different accounts—this one, ours, and each of the others discussed above.

Finally, anthropomorphism might be a type of mental modeling. Mental models are representations used to make sense of something in the real world, and might be approximations or metaphors [21,32]. An attempt to explain anthropomorphism as mental modeling would be useful future work. In a related account, Airenti describes anthropomorphism as a mode of interaction that we learn to use as a tool to interact with things that are not human [33]. These two accounts do not rule out that anthropomorphism involves the imagination or even make-believe, but the purpose of anthropomorphism as a useful tool is emphasized.

6. Conclusion

This paper discussed how humans make attributions that they know to be false about robots. Applying Kendall Walton's theory of make-believe adds something new to the picture: humans are considering fictional states of affairs when making these attributions. Using the additional concepts of faking and imitating, the implications of this new picture were discussed for how we make attributions and use mental models, and then for human-robot interaction from the human's perspective.

It has been mentioned that whereas anthropomorphism is mysterious, make-believe and metaphor have long been studied in several fields. This connection can now be leveraged: inasmuch as the above Waltonian account of human-robot interaction is true, interacting with a robot is analogous to reading a novel or watching a play or film. These potential analogies suggest that well-developed fields such as acting theory, narrative theory, and film theory could be very valuable to human-robot interaction researchers. Since those fields are more mature than human-robot interaction, exploring these comparisons could yield novel research directions or even breakthroughs.

Acknowledgments

This work was supported by an NSF NRI grant for "Socially Aware, Expressive, and Personalized Mobile Remote Presence: Co-Robots as Gateways to Access to K-12 In-School Education", NSF IIS-1528121. The authors also gratefully acknowledge Tim Buchanan for the early discussions that led to this paper, Sam Thellman and Joel Oesch for their advice and encouragement, and Kathryn Nelson for inspiring the idea of faking or imitating participation in a make-believe. Finally, the authors would like to thank the two anonymous reviewers of this manuscript for their thoughtful and detailed comments.

References

- [1] Blake R, Sekuler R. Perception. 5th ed. McGraw-Hill; 2005.
- [2] Heider F. The psychology of interpersonal relations. Psychology Press; 2013.
- [3] Dennett DC. Précis of the intentional stance. Behavioral and brain sciences. 1988;11(3):495–505.
- [4] Nass C, Moon Y. Machines and Mindlessness: Social Responses to Computers. Journal of Social Issues. 2000;56(1):81–103.
- [5] Zlotowski J, Proudfoot D, Yogeewaran K, Bartneck C. Anthropomorphism: opportunities and challenges in human–robot interaction. International Journal of Social Robotics. 2015;7(3):347–360.
- [6] Fussell SR, Kiesler S, Setlock LD, Yew V. How People Anthropomorphize Robots. In: Proceedings of the 3rd ACM/IEEE International Conference on Human Robot Interaction. New York, NY, USA: ACM; 2008. p. 145–152.
- [7] Haikonen PO. Consciousness and robot sentience. World Scientific; 2012. (Series on Machine Consciousness; vol. 2).
- [8] Walton KL. Metaphor and prop oriented make-believe. Fictionalism in metaphysics. 2005;65–87.
- [9] Heider F, Simmel M. An experimental study of apparent behavior. The American journal of psychology. 1944;57(2):243–259.
- [10] Malle BF, Korman J. Attribution theory. Oxford University Press; 2013.
- [11] Nass C, Steuer J, Tauber ER. Computers Are Social Actors. In: Proceedings of the SIGCHI Conference on Human Factors in Computing Systems. New York, NY, USA: ACM; 1994. p. 72–78. (CHI '94).
- [12] Thellman S, Ziemke T. The Intentional Stance Toward Robots: Conceptual and Methodological Considerations. In: The 41st Annual Conference of the Cognitive Science Society, July 24-26, Montreal, Canada. 2019. p. 1097–1103.

- [13] Krämer NC, von der Pütten A, Eimler S. Human-agent and human-robot interaction theory: similarities to and differences from human-human interaction. In: Human-computer interaction: The agency perspective. Springer; 2012. p. 215–240.
- [14] Epley N, Waytz A, Cacioppo JT. On seeing human: a three-factor theory of anthropomorphism. *Psychological review*. 2007;114(4):864.
- [15] Zlotowski J, Sumioka H, Eyssel F, Nishio S, Bartneck C, Ishiguro H. Model of Dual Anthropomorphism: The Relationship Between the Media Equation Effect and Implicit Anthropomorphism. *International Journal of Social Robotics*. 2018;10(5):701–714.
- [16] Severson RL, Woodard SR. Imagining others' minds: The positive relation between children's role play and anthropomorphism. *Frontiers in psychology*. 2018;9:2140.
- [17] Clark HH, Van Der Wege MM. Imagination in Narratives. In: *The handbook of discourse analysis*. 2nd ed. John Wiley & Sons, Inc.; 2015. p. 406–421.
- [18] Walton KL. *Mimesis as make-believe: On the foundations of the representational arts*. Harvard University Press; 1990.
- [19] Novitz D. Fiction, Imagination and Emotion. *The Journal of Aesthetics and Art Criticism*. 1980;38(3):279–88.
- [20] Gerrig R. *Experiencing Narrative Worlds*. Routledge; 2018.
- [21] Byrne RM. Towards a model theory of imaginary thinking. In: *Mental models in cognitive science: Essays in honour of Phil Johnson-Laird*. Psychology Press; 1996. p. 155–174.
- [22] Walton KL. Appreciating fiction: suspending disbelief or pretending belief? *Dispositio*. 1980;5(13/14):1–18.
- [23] Tannen D. What's in a frame? Surface evidence for underlying expectations. *Framing in discourse*. 1993;14:56.
- [24] Seibt J. Towards an Ontology of Simulated Social Interaction. In: *Sociality and Normativity for Robots: Philosophical Inquiries into Human-Robot Interactions*. Springer; 2017. p. 11–39.
- [25] Seibt J, Damholdt MF, Vestergaard C. Five Principles of Integrative Social Robotics. In: *Robophilosophy/TRANSOR*. 2018. p. 28–42.
- [26] Seibt J. "Integrative Social Robotics" - A New Method Paradigm to Solve the Description Problem And the Regulation Problem? In: *Robophilosophy/TRANSOR*. 2016. p. 104–115.
- [27] Damholdt MF, Vestergaard C, Seibt J. Testing for 'Anthropomorphization': A Case for Mixed Methods in Human-Robot Interaction. In: Jost C, Le Pèvédic B, Belpaeme T, Bethel C, Chrysostomou D, Crook N, et al., editors. *Human-Robot Interaction*. Springer, Cham; 2020. p. 203–227. (Springer Series on Bio- and Neurosystems; vol. 12).
- [28] Clark HH. *Using language*. Cambridge university press; 1996.
- [29] Clark HH. Depicting as a method of communication. *Psychological Review*. 2016;123(3):324–47.
- [30] Remmers P. The artificial nature of social robots: A phenomenological interpretation of two conflicting tendencies in human-robot interaction. In: *Robophilosophy/TRANSOR*. 2020.
- [31] Musiał M. Magical thinking and empathy towards robots. In: *What social robots can and should do Proceedings of Robophilosophy*. 2016. p. 347–355.
- [32] Johnson-Laird PN. *Mental models: Towards a cognitive science of language, inference, and consciousness*. Harvard University Press; 1983.
- [33] Airenti G. The development of anthropomorphism in interaction: intersubjectivity, imagination, and theory of mind. *Frontiers in psychology*. 2018;9:2136.