

Action Anticipation Based on an Agent's Epistemic State in Toddlers and Adults

The ManyBabies2 Consortium¹

¹ See author note

The ManyBabies2 Consortium consists of Tobias Schuwerk (Ludwig-Maximilians-Universität München), Dora Kampis* (University of Copenhagen), Renée Baillargeon (University of Illinois at Urbana-Champaign), Szilvia Biro (Leiden University), Manuel Bohn (Max Planck Institute for Evolutionary Anthropology), Krista Byers-Heinlein (Concordia University), Sebastian Dörrenberg (University of Bremen), Cynthia Fisher (University of Illinois at Urbana-Champaign), Laura Franchin (University of Trento), Tess Fulcher (University of Chicago), Isa Garbisch (University of Göttingen), Alessandra Geraci (University of Trento), Charlotte Grosse Wiesmann (Max Planck Institute for Human Cognitive and Brain Sciences), J. Kiley Hamlin (University of British Columbia), Daniel Haun (Max Planck Institute for Evolutionary Anthropology), Robert Hepach (University of Oxford), Sabine Hunnius (Radboud University Nijmegen), Daniel C. Hyde (University of Illinois at Urbana-Champaign), Petra Kármán (Central European University), Heather L Kosakowski (MIT), Ágnes M. Kovács (Central European University), Anna Krämer (University of Salzburg), Louisa Kulke (Friedrich-Alexander-University Erlangen-Nürnberg), Crystal Lee (Princeton University), Casey Lew-Williams (Princeton University), Ulf Liszkowski (Universität Hamburg), Kyle Mahowald (University of California, Santa Barbara), Olivier Mascaró (Integrative Neuroscience and Cognition Center, CNRS UMR8002/University of Paris), Marlene Meyer (Radboud University Nijmegen), David Moreau (University of Auckland), Josef Perner (University of Salzburg), Diane Poulin-Dubois (Concordia University), Lindsey J. Powell (University of California, San Diego), Julia Prein (Max Planck Institute for Evolutionary Anthropology), Beate Priewasser (University of Salzburg), Marina Proft (Universität Göttingen), Gal Raz (MIT), Peter Reschke (Brigham Young University), Josephine Ross (University of Dundee), Katrin Rothmaler (Max Planck Institute for Human Cognitive and Brain Sciences), Rebecca Saxe (MIT), Dana Schneider (Friedrich-Schiller-University Jena, Germany), Victoria Southgate (University of Copenhagen), Luca Surian (University of

31 Trento), Anna-Lena Tebbe (Max Planck Institute for Human Cognitive and Brain
32 Sciences), Birgit Träuble (Universität zu Köln), Angeline Sin Mei Tsui (Stanford
33 University), Annie E. Wertz (Max Planck Institute for Human Development), Amanda
34 Woodward (University of Chicago), Francis Yuen (University of British Columbia),
35 Amanda Rose Yuile (University of Illinois at Urbana-Champaign), Luise Zellner
36 (University of Salzburg), Lucie Zimmer (Ludwig-Maximilians-Universität München),
37 Michael C. Frank (Stanford University), and Hannes Rakoczy (University of Göttingen).

38 Correspondence concerning this article should be addressed to The ManyBabies2
39 Consortium, Leopoldstr. 13, 80802 München, Germany. E-mail:
40 tobias.schuwerk@psy.lmu.de

Abstract

Do toddlers and adults engage in spontaneous Theory of Mind (ToM)? Evidence from anticipatory looking (AL) studies suggests they do. But a growing body of failed replication studies raised questions about the paradigm's suitability, urging the need to test the robustness of AL as a spontaneous measure of ToM. In a multi-lab collaboration we examine whether 18- to 27-month-olds' and adults' anticipatory looks distinguish between two basic forms of epistemic states: knowledge and ignorance. In toddlers [ANTICIPATED $n = 520$ 50% FEMALE] and adults [ANTICIPATED $n = 408$, 50% FEMALE], we found [SUPPORT/NO SUPPORT] for epistemic state-based action anticipation. Future research can probe whether this conclusion extends to more complex kinds of epistemic states, such as true and false beliefs.

Keywords: anticipatory looking; spontaneous Theory of Mind; replication

Word count: 10243

Action Anticipation Based on an Agent’s Epistemic State in Toddlers and Adults

The capacity to represent epistemic states, known as Theory of Mind (ToM) or mentalizing, plays a central role in human cognition (Dennett, 1989; Frith & Frith, 2006; Premack & Woodruff, 1978). Although ToM has been under intense scrutiny in the past decades, its nature and ontogeny are still the subjects of much controversy. At the heart of these debates are questions about the reliability of the tools used to measure ToM (Baillargeon, Buttelmann, & Southgate, 2018; e.g., Poulin-Dubois et al., 2018), among others, anticipatory looking (AL) paradigms. To address this issue, in a collaborative long-term project we assess the robustness of infants’ and adults’ tendency to spontaneously take into account different kinds of epistemic states — what they perceive, know, think, or believe — when predicting others’ behaviors. This paper reports the first foundational step of this project, which focuses on the most basic epistemic state ascription: the capacity to distinguish between knowledgeable and ignorant individuals. Simple forms of knowledge attribution (such as tracking what other individuals have seen or experienced) are typically assumed to develop early and to operate spontaneously throughout the lifespan (Liszkowski, Carpenter, & Tomasello, 2007; e.g., Luo & Baillargeon, 2007; O’Neill, 1996; Phillips et al., 2021). Thus, evaluating whether ToM measures are sensitive to the knowledge-ignorance distinction is a crucial test case to assess their robustness. The present paper investigates this question in an AL paradigm including 18-27-month-old infants and adults.

In the following sections we first establish the background and scientific context of this study, namely the reliability and replicability of spontaneous ToM measures. We then introduce a novel way to approach these issues: a large-scale collaborative project targeting the replicability of ToM findings. Finally, we outline the rationale of the present study which uses an AL paradigm to test whether infants and adults distinguish between two basic forms of an agent’s epistemic state: knowledge and ignorance.

Spontaneous Theory of Mind tasks

Humans are proficient at interpreting and predicting others' intentional actions. Adults as well as infants expect agents to act persistently towards the goal they pursue Woodward & Sommerville (2000), and anticipate others' actions based on their goals even before goals are achieved - that is, humans engage in goal-based action anticipation (for review, see Elsner & Adam, 2021; but see Ganglmayer, Attig, Daum, & Paulus, 2019). To predict others' actions, however, it is essential to consider their epistemic state: what they perceive, know, or believe. A number of seminal studies using non-verbal spontaneous measures have suggested that infants, toddlers, older children, and adults show action anticipation and action understanding not only based on other agents' goals (what they want) but also on the basis of their epistemic status (what they perceive, know, or believe). These studies suggest that from infancy onwards, humans spontaneously engage in ToM or mentalizing. For example, studies using violation of expectation methods have demonstrated that infants look longer in response to events in which an agent acts in ways that are incompatible with their (true or false) beliefs, compared to events in which they act in belief-congruent ways (Onishi & Baillargeon, 2005; Surian, Caldi, & Sperber, 2007; Träuble, Marinović, & Pauen, 2010). Other studies have employed more interactive tasks requiring the child to play, communicate, or cooperate with experimenters and, for example, give an experimenter one of several objects as a function of their epistemic status. Such studies have shown that toddlers spontaneously adjust their behavior to the experimenter's beliefs (D. Buttelmann, Carpenter, & Tomasello, 2009; Király, Oláh, Csibra, & Kovács, 2018; Knudsen & Liszkowski, 2012; Southgate, Johnson, Karoui, & Csibra, 2010).

The largest body of evidence for spontaneous ToM comes from studies using AL tasks. In such tasks, participants see an agent who acts in pursuit of some goal (typically, to collect a certain object) and has either a true or a false belief (for example, regarding the location of the target object). A number of studies have shown that infants, toddlers,

older children, neurotypical adults, and even non-human primates anticipate (indicated by looks to the location in question) that an agent will go where it (truly or falsely) believes the object to be rather than, irrespective of the actual location of the object (Gliga, Jones, Bedford, Charman, & Johnson, 2014; Grosse Wiesmann, Friederici, Singer, & Steinbeis, 2017; Hayashi et al., 2020; Kano, Krupenye, Hirata, Tomonaga, & Call, 2019; Krupenye, Kano, Hirata, Call, & Tomasello, 2016; Meristo et al., 2012; Schneider, Bayliss, Becker, & Dux, 2012; Schneider, Slaughter, Bayliss, & Dux, 2013; Senju et al., 2010; Senju, Southgate, Snape, Leonard, & Csibra, 2011; Senju, Southgate, White, & Frith, 2009; Surian & Franchin, 2020; Thoermer, Sodian, Vuori, Perst, & Kristen, 2012). These studies have revealed converging evidence for spontaneous ToM across the human lifespan and even in other primate species.

Across the different measures, the majority of early works on spontaneous ToM in infants and toddlers have reported positive results in the second year of life, and a few studies even within the first year (Kovács, Téglás, & Endress, 2010; Luo & Baillargeon, 2010; Southgate & Verneti, 2014), yielding a rich body of coherent and convergent evidence (for reviews see e.g., Barone, Corradi, & Gomila, 2019; Kamps, Buttelmann, & Kovács, 2020; Scott & Baillargeon, 2017). This growing body of literature has led to a theoretical transformation of the field. In particular, findings with young infants have paved the way for novel accounts of the development and cognitive foundations of ToM. The previous consensus was that full-fledged ToM emerges only at around age 4, potentially as the result of developing executive functions, complex language skills and other factors (e.g., Perner, 1991; Wellman & Cross, 2001). In contrast, the newer accounts proposed that some basic forms of ToM may be phylogenetically more ancient and may develop much earlier in ontogeny (e.g., Baillargeon, Scott, & He, 2010; Carruthers, 2013; Kovács, 2016; Leslie, 2005).

Recently, however, a number of studies have raised uncertainty regarding the empirical foundations of the early-emergence theories, as we review below. In the following

sections, we present an overview of the current empirical picture of early understanding of epistemic states and then introduce ManyBabies2 (MB2), a large-scale collaborative project exploring the replicability of ToM in infancy, of which the current study constitutes the first step.

Replicability of Spontaneous Theory of Mind Tasks

A number of failures to replicate findings from spontaneous ToM tasks have recently been published with infants, toddlers, and adults Kulke & Rakoczy (2019). Besides conceptual replications, many of these studies involve more direct replication attempts with the original stimuli and procedures. One of these was a two-lab replication attempt of one of the most influential AL studies (Southgate, Senju, & Csibra, 2007). This failure to replicate is especially notable not only because of the influence of the original finding of the field, but also because of the large sample size and the involvement of some of the original authors (Kampis et al., 2021). Additional unpublished replication failures have also been reported. Kulke and Rakoczy (2018) examined 65 published and non-published studies including 36 AL studies (replications of Schneider et al., 2012; Southgate et al., 2007; Surian & Geraci, 2012; and Low & Watts, 2013), as well as studies using other paradigms, and classified them as a successful, partial, or non-replication, depending on whether all, some, or none of the original main effects were found. Although no formal analysis of effect size was carried out, overall, non-replications and partial replications outnumbered successful replications, regardless of the method used. In addition to the failure to replicate spontaneous anticipation of agents' behaviors based on their beliefs, many of the replication studies revealed an even more fundamental problem of spontaneous AL procedures: a failure to adequately anticipate an agent's action in the absence of a belief. That is, researchers did not find evidence for spontaneous anticipation of agents' behaviors based on their goals, even in the initial familiarization trials of the experiments, where the agent's beliefs do not play any role yet (e.g., Kampis et al., 2020; Kulke, Reiß, et al., 2018;

Schuwerk et al., 2018). The familiarization trials are designed to convey the goal of the agent, as well as the general timing and structure of events, to set up participants' expectations in the test trials where the agent's epistemic state is then manipulated. Typically, the last familiarization trial can also be used to probe participants' spontaneous action anticipation; and test trials can only be meaningfully interpreted if there is evidence of above-chance anticipation in the familiarization trials. In several AL studies many participants had to be excluded from the main analyses for failing to demonstrate robust action anticipation during the familiarization trials (e.g., Kamps et al., 2020; Kulke, Reiß, et al., 2018; Schuwerk et al., 2018; Southgate et al., 2007). This raises the possibility that these paradigms may not be suitable for reliably eliciting spontaneous action prediction in the first place (for discussion see Baillargeon et al., 2018). In sum, in light of the complex and mixed state of the evidence, it currently remains unclear whether infants, toddlers, and adults engage in spontaneous ToM. This calls for systematic, large-scale, a priori designed multi-lab study that stringently tests for the robustness, reliability, and replicability of spontaneous measures of ToM.

General Rationale of MB2

To this end, ManyBabies 2 (MB2) was established as an international consortium dedicated to investigating infants' and toddlers' ToM skills. The main aim is to test the replicability and thus reliability of findings from spontaneous ToM tasks. In the long-term, MB2 will build on the initial findings and the aim will be extended to include testing the validity of these experimental designs and addressing theoretical accounts of spontaneous ToM. MB2 operates under the general umbrella of ManyBabies (MB), a large-scale international research consortium founded with the aim of probing the reliability of central findings from infancy research. In particular, MB projects bring together large and theoretically diverse groups of researchers to tackle pressing questions of infant cognitive development, by collaboratively designing and implementing methodologies and

pre-registered analysis plans (Frank et al., 2017). The MB2 consortium involves authors of original studies as well as authors of both successful and failed replication studies, and researchers from very different theoretical backgrounds. It thus presents a case of true “adversarial collaboration” (Mellers, Hertwig, & Kahneman, 2001).

Rationale of the Present Study

Based on both theoretical and practical considerations, the current paper presents the first foundational step in MB2, focusing on AL measures. It investigates whether toddlers and adults anticipate (in their looking behavior) how other agents will act based on their goals (i.e., what they want) and epistemic status (i.e., what they know or do not know). From a practical perspective, we focus on AL since it is a child-friendly and widely used method that is also suitable for humans across the lifespan and even other species. Additionally, as AL is screen-based and standardizable, identical stimuli can be presented in different labs. From a theoretical perspective, given the mixed findings with AL tasks reviewed in the previous section, we take a systematic and bottom-up approach. First, we probe whether AL measures are suitable for measuring spontaneous goal-directed action anticipation. With the aim to improve the low overall rates of anticipatory looks in recent studies, we designed new, engaging stimuli to test whether these are successful in eliciting spontaneous action anticipation. Second, in case reliably elicited action anticipation can be found: we probe whether toddlers and adults take into account the agent’s epistemic status in their spontaneous goal-based action anticipation. That is, do they track whether the agent saw or did not see a crucial event, and therefore whether this agent does or does not know something? In the current study we focus on the most basic form of tracking the epistemic status of agents: considering whether they had access to relevant information, and whether they are thus *knowledgeable* or *ignorant*. We reasoned that only after establishing whether a context can elicit spontaneous tracking of an agent’s epistemic status in a more basic sense (i.e., the agent’s knowledge vs. ignorance) is it eventually

meaningful to ask whether this context also elicits more complex epistemic state tracking (i.e., the agent’s beliefs). Answering these first two questions in the present study will allow us, in the long run, to address a third set of questions in subsequent studies, probing the nature of the representations and cognitive mechanisms involved in infant ToM. Do toddlers and adults engage in full-fledged belief-ascription in their spontaneous goal-based action anticipation? What *kind* of epistemic states do toddlers and adults spontaneously attribute to others in their action anticipation (e.g., Horschler, MacLean, & Santos, 2020; Phillips et al., 2021)? Do the results that prove replicable really assess ToM, or can they be interpreted in alternative ways such as behavioral rules, associations, or simple perceptual preferences (see, e.g., Heyes, 2014; Perner & Ruffman, 2005)? The present study lays the foundation for investigating these questions. Regarding the knowledge-ignorance distinction, many accounts in developmental and comparative ToM research have argued for the ontogenetic and evolutionary primacy of representing *what* agents witness and represent, relative to more sophisticated ways of representing *how* agents represent (and potentially mis-represent) objects and situations (e.g., Apperly & Butterfill, 2009; Flavell, 1988; Kaminski, Call, & Tomasello, 2008; Martin & Santos, 2016; Perner, 1991; Phillips et al., 2021). For example, it is often assumed that young children and non-human primates may be capable of so-called “Level I perspective-taking” (understanding *who* sees *what*) but only human children from around age 4 may finally develop capacities for “Level II perspective-taking” [understanding *how* a given situation may appear to different agents; Flavell, Everett, Croft, and Flavell (1981)]. Empirically, many studies using verbal and/or interactive measures have indicated that children may engage in knowledge-ignorance and related distinctions before they engage in more complex forms of meta-representation (e.g., Flavell et al., 1981; Hogrefe, Wimmer, & Perner, 1986; Moll & Tomasello, 2006; O’Neill, 1996; F. Buttelmann & Kovács, 2019; F. Buttelmann, Suhrke, & Buttelmann, 2015; Kampis et al., 2020; though for some findings indicating Level II perspective-taking at an early age see Scott & Baillargeon, 2009; Scott, Richman, & Baillargeon, 2015), and that

non-human primates seem to master knowledge-ignorance tasks while not demonstrating any more complex, meta-representational form of ToM (e.g., Hare, Call, & Tomasello, 2001; Kaminski et al., 2008; Karg, Schmelz, Call, & Tomasello, 2015). The knowledge-ignorance distinction thus appears to be an ideal candidate for assessing epistemic status-based action anticipation in a wide range of populations. To date, however, no study has probed whether or how children’s (and adults’) spontaneous action anticipation, as indicated by AL, is sensitive to ascriptions of knowledge vs. ignorance. Most studies that have addressed ToM with AL measures have targeted the more sophisticated true/false belief contrast. As reviewed above, the results of those studies yield a mixed picture regarding replicability of the findings. It has been argued that tasks that reliably replicate are ones which can be solved with the more basic knowledge-ignorance distinction, whereas tasks that do not replicate require more sophisticated belief-ascription (Powell et al., 2018)¹, suggesting that only some but not all findings might not be replicable. Based on these considerations, the present study tests whether toddlers and adults engage in knowledge- and ignorance-based AL to probe the most basic form of spontaneous, epistemic state-based action anticipation.

Design and Predictions of the Present Study

The current study presents 18- to 27-month-old toddlers and adults with animated scenarios while measuring their gaze behavior. Testing adults (and not just toddlers) is crucial to address debates about the validity and interpretation of AL measures of ToM throughout the lifespan (e.g., Schneider, Slaughter, & Dux, 2017). Following the structure of previous AL paradigms, participants are first familiarized to an agent repeatedly

¹ For example, some studies have found partial replication results, with patterns of the following kind: participants showed systematic anticipation (or appropriate interactive responses) in true belief trials but showed looking (or interactive responses) at chance level in the false belief trials (e.g., Dörrenberg, Wenzel, Proft, Rakoczy, & Liszkowski, 2019; Kulke, Reiß, et al., 2018; Powell et al., 2018). Such a pattern remains ambiguous since it may merely reflect a knowledge-ignorance distinction.

approaching a target (familiarization trials). AL is measured during familiarization trials to probe whether participants understood the agent’s goal and spontaneously anticipate their actions. Subsequently, during test trials the agent’s visual access is manipulated, leading them to be either *knowledgeable* or *ignorant* about the location of the target. Participants’ AL will be measured during test trials to determine whether or not they take into account the agent’s epistemic access and adjust their action anticipation accordingly. Participants’ looking patterns will be recorded using either lab-based corneal reflection eye-tracking or online recording of gaze patterns. We chose to provide the online testing option to increase the flexibility for data collection given the disruption caused by the Covid-19 pandemic. This option will also provide the opportunity to potentially compare in-lab and online testing procedures (Sheskin et al., 2020). Novel animated stimuli were collectively developed within the MB2 consortium on the basis of previous work (e.g., Clements & Perner, 1994) and based on input from collaborators with experience with both successful and failed replication studies (e.g., Grosse Wiesmann et al., 2017; Surian & Geraci, 2012). These animated 3D scenes feature a dynamic interaction aimed to optimally engage participants’ attention: a chasing scenario involving two agents, a *chaser* and a *chasee* (see Figures 1 and 2). As part of the chase, the chasee enters from the top of an upside-down Y-shaped tunnel with two boxes at its exits. The tunnel is opaque so participants cannot see the chasee after it enters the tunnel, but can hear noises that indicate movement. The chasee eventually exits from one of the arms of the Y, and goes into the box on that side. The chaser observes the chasee exit the tunnel and go into a box, and then follows it through the tunnel. During familiarization trials, the chaser always exits the tunnel on the same side as the chasee, and approaches the box where the chasee is currently located. Thus, if participants engage in spontaneous action anticipation during familiarization trials, they should reliably anticipate during the period when the chaser is in the tunnel that it will emerge at the exit that leads to the box containing the chasee. During test trials, the chasee always first hides in one of the boxes but shortly thereafter

leaves its initial hiding place and hides in the box at the other tunnel exit. Critically, the chaser either does (*knowledge* condition) or does not (*ignorance* condition) have epistemic access to the chasee's location. During *knowledge* trials, the chaser observes all movements of the chasee. During *ignorance* trials, the chaser observes the chasee enter the tunnel, but then leaves and only returns once the chasee is already hidden inside the second box. The event sequences in the two conditions are thus identical with the only difference between conditions pertaining to what the chaser has or has not seen. They were designed in this way with the long-term aim to implement, in a minimal contrast design, more complex conditions of false/true belief contrasts with the very same event sequences (true belief conditions will then be identical to the knowledge conditions here, but in false belief conditions the chaser witnesses the chasee's placement in the first box, but then fails to witness the re-location)². Participants' AL (their gaze pattern indicating where they expect the chaser to appear) will be assessed during the anticipatory period - that is, the period during which the chaser is going through the tunnel and is not visible. There will be two main dependent measures: first looks, and a differential looking score (DLS). The first look measure will be binary, indicating which of the two tunnel exits participants fixate first: the exit where the chasee is actually hiding, or the other exit. DLS is a measure of the proportion of time spent looking at the correct tunnel exit during the entire anticipatory

² There is thus a certain asymmetry with regard to the interpretation and the consequences of potentially positive and negative results of the present knowledge-ignorance contrast: in the case of positive results, we can conclude that subjects spontaneously engage in basic epistemic state ascription and can move on to test, with the minimal contrast comparison of knowledge-ignorance vs. false belief-true belief, whether this extends to more complex forms of epistemic state attribution. In the case of negative results, though, we cannot draw firm conclusions to the effect that subjects do not engage in spontaneous epistemic state ascription. More caution is in order since the present knowledge-ignorance contrast has been designed in order to be comparable to future belief contrasts rather than to be the simplest implementation possible. Simpler implementations would then need to be devised that involve fewer steps (i.e. the chasee just goes to one location and this is or is not witnessed by the chasee).

period. In two pilot studies (see Methods section), we addressed the foundational question of the current study: whether these stimuli reveal spontaneous goal-directed action anticipation as measured by AL in the above-described familiarization trials (i.e., without a change of location by the chasee or manipulation of the chaser’s epistemic state). We found that our paradigm indeed elicited action anticipation and exclusion rates due to lack of anticipation were significantly lower relative to previous (original and replication) AL studies. Both toddlers and adults showed reliable anticipation of the chaser’s exit at the chasee’s location, indicating that in contrast with many previous AL studies the current paradigm successfully elicits spontaneous goal-based action anticipation. Based on these pilot data we concluded that the paradigm is suitable for examining the second and critical question: whether toddlers and adults, in their spontaneous goal-based action anticipation, take into account the agent’s epistemic state. We predict that if participants track the chaser’s perceptual access and resulting epistemic state (knowledge/ignorance) and anticipate their actions accordingly, they should look more in anticipation to the exit at the chasee’s location than the other exit in the *knowledge* condition, but should not do so (or to a lesser degree; see below) in the *ignorance* condition. We anticipate three potential factors that could influence participant’s gaze patterns: Keeping track of the chaser’s epistemic status in the *ignorance* condition might either lead to no expectations as to where the chaser will look (resulting in chance level looking between the two exits) or (if participants follow an “ignorance leads to mistakes”-rule, see e.g., Ruffman, 1996) to an expectation that the chaser will go to the wrong location [longer looking to the exit with the empty box; e.g., Fabricius, Boyer, Weimer, and Carroll (2010)]. Either way, participants may still show a ‘pull of the real’ even in the *ignorance* condition, i.e., reveal a default tendency to look to the side where the chasee is located. But if they truly keep track of the epistemic status of the chaser (*knowledge* vs. *ignorance*), they should show this tendency to look to the side where the chasee really is in the *ignorance* condition to a lesser degree than in the *knowledge* condition. In sum, the research questions of the present study

are the following: First, can we observe in a large sample that toddlers and adults robustly anticipate agents' actions based on their goals in this paradigm, as they did in our pilot study? Second, can we find evidence that they take into account the agent's epistemic access (knowledge vs. ignorance) and adjust their action anticipation accordingly? In addressing these questions, the present study will significantly contribute to our knowledge on spontaneous ToM. It will inform us whether the present paradigm and stimuli can elicit spontaneous goal-based and mental-state-based action anticipation in adults and toddlers, based on a large sample of about 800 participants in total from over 20 labs. In the long run, the present study will lay the foundation for future work to address broader questions of what *kind* of epistemic states toddlers and adults spontaneously attribute to others in their action anticipation and what cognitive mechanisms allow them to do so.

Methods

All materials, and later the collected de-identified data, will be provided on the Open Science Framework (OSF; <https://osf.io/jmuvd/>). All analysis scripts, including the pilot data analysis and simulations for the design analysis, can be found on GitHub (<https://github.com/manybabies/mb2-analysis>). We report how we determined our sample size and we will report all data exclusions, all manipulations, and all measures in the study. Additional methodological details can be found in the Supplemental Material.

Stimuli

Figures 1 and 2 provide an overview of the paradigm. For the stimuli, 3D animations were created depicting a chasing scenario between two agents (chaser and chasee) who start in the upper part of the scene. At the very top of the scene a door leads to outside the visible scene. Below this area, a horizontal fence separates the space, and thus the lower part of the space can be reached by the Y-shaped tunnel only. Additional information on

the general scene setup, events, and timings in the familiarization and the test trials, as well as trial randomization can be found in the Supplemental Material.

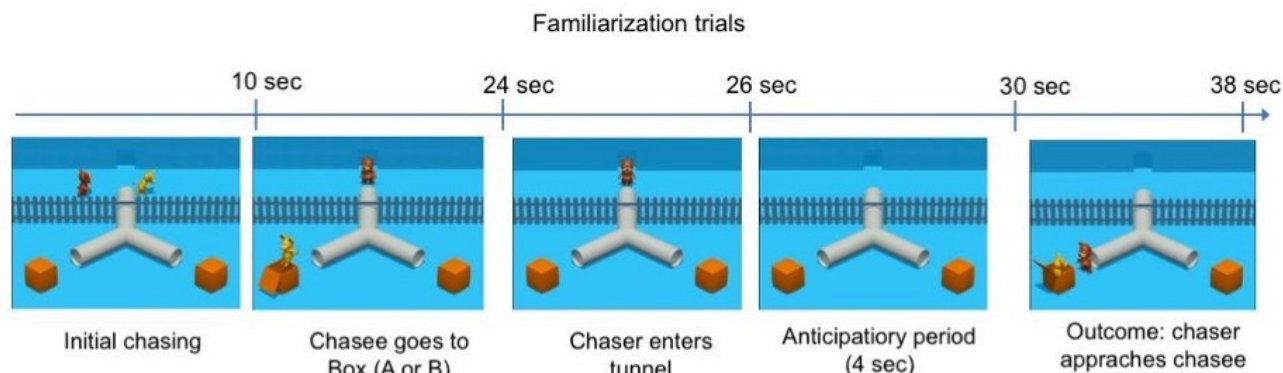


Figure 1. Timeline of the familiarization trials.

Familiarization Trials. All participants will view four familiarization trials (for an overview of key events see Figure 1). During familiarization trials, after a brief chasing introduction, the chasee enters an upside-down Y-shaped tunnel with a box at both of its exits. The chasee then leaves the tunnel through one of the exits and hides in the box on the corresponding side. Subsequently, the chaser enters the tunnel (to follow the chasee), and participants' AL to the tunnel exits is measured before the chaser exits on the side the chasee is hiding, as an index of their goal-based action anticipation. In these familiarization trials, if participants engage in spontaneous action anticipation, they should reliably anticipate that the chaser should emerge at the tunnel exit that leads to the box where the chasee is. After leaving the tunnel, the chaser approaches the box in which the chasee is hiding and knocks on it. Then, the chasee jumps out of the box and the two briefly interact.

Familiarization Phase Pilot Studies. In a pilot study with 18- to 27-month-olds ($n = 65$) and adults ($n = 42$), seven labs used in-lab corneal reflection eye-tracking to collect data on gaze behavior in the familiarization phase. A key desideratum of our paradigm is that it should produce sufficient AL, as a low rate of AL in previous studies has led to high exclusion rates. The goals of the pilot study were to 1) estimate the level of correct goal-based action predictions in the familiarization phase, 2)

determine the optimal number of familiarization trials, 3) check for issues with perceptual properties of stimuli (e.g., distracting visual saliencies), and 4) test the general procedure including preprocessing and analyzing raw gaze data from different eye-tracking systems. We found that the familiarization stimuli elicited a relatively high proportion of goal-directed action anticipations, but we were concerned about the effects of some minor properties of the stimulus (in particular, a small rectangular window in the tunnel tube that allowed participants to see the agents at one point on their path to the tunnel exits). In a second pilot study with 18- to 27-month-olds ($n = 12$, three participating labs), slight changes of stimulus features (the removal of the window in the tube; temporal changes of auditory anticipation cue) did not cause major changes in the AL rates. Sixty-eight percent of toddlers' first looks in the first pilot, 69% of toddlers' first looks in the second pilot, and 69% of adults' first looks were toward the correct area of interest (AOI) during the anticipatory period. The average proportion of looking towards the correct AOI during the anticipatory period was 70.7% ($CI_{95\%} = 67.6\% - 73.8\%$) in toddlers in the first pilot, 70.5% ($CI_{95\%} = 62.8\% - 78.2\%$) in the second pilot for toddlers, and 75.3% ($CI_{95\%} = 71.0\% - 79.5\%$) in adults. In Bayesian analyses, we found strong evidence that toddlers and adults looked more towards the target than towards the distractor during the anticipation period. Based on conceptual and practical methodological considerations while also considering previous studies, we decided to include four trials in the final experiment. The pilot data results of the toddlers supported this decision insofar as we observed a looking bias towards the correct location already in trials 1-4, without additional benefit of trials 5-8. Further, prototypical analysis pipelines were established for combining raw gaze data from different eye-trackers. In short, we developed a way to resample gaze data from different eye-trackers to be at a common Hz rate and to define proportionally correct AOIs for different screen dimensions with the goal to merge all raw data into one data set for inferential statistics. The established analysis procedure is described further in the Data Preprocessing section below. In sum, we concluded that this paradigm sufficiently elicits

goal-directed action predictions, an important prerequisite for drawing any conclusion on AL behavior in the test trials of this study. A detailed description of the two pilot studies can be found in the Supplemental Material.

Test Trials. All participants will see two test trials, one *knowledge* and one *ignorance* trial. However, in line with common practice in ToM studies, the main comparison concerns the first test trial between-participants to avoid potential carryover effects. In addition, in exploratory analyses, we plan to assess whether results remain the same if both trials are taken into account and whether gaze patterns differ between the two trials (see Exploratory Analyses). If the results remain largely unchanged across the two trials, it may suggest that future studies could increase power by including multiple test trials. In test trials, the chasee first hides in one of the boxes, but shortly thereafter the chasee leaves this box and hides in the second box, at the other tunnel exit. Critically, the chaser either witnesses (*knowledge* condition) or does not witness (*ignorance* condition) from which tunnel exit the chasee exited and thus where the chasee is currently hiding (for an overview, see Figure 2). In the *knowledge* trials, the chaser observes all movements of the chasee. The chaser leaves for a brief period of time after the chasee entered the tunnel, but it returns before the chasee exits the tunnel. Therefore, no events take place in the chaser’s absence. In the *ignorance* trials, the chaser sees the chasee enter the tunnel, but then leaves. Therefore, the chaser does not see the chasee entering either box and only returns once the chasee is already hidden in the final location. Finally, the chaser enters the tunnel but does not appear in either exit. Rather, the scene “freezes” for four seconds and participants’ AL is measured. Thus, the *knowledge* and *ignorance* conditions are matched for the chaser leaving for a period of time, but they differ in whether they warrant the chaser’s epistemic access to the location of the chasee. No outcome is shown in either test trials. When designing the *knowledge* and *ignorance* condition, we aimed at keeping all events and their timings parallel, except the crucial manipulation. We show the same events in both conditions. Where possible, all events also have the same duration. In the

case of the chaser’s absence in the *knowledge* condition, there were two main options, both with inevitable trade-offs. First, we could have increased the duration of the chaser’s absence in the *knowledge* condition to match the duration of the chaser’s absence in both conditions. Yet, this would potentially disrupt the flow of events, such as keeping track of the chasee’s actions and the general scene dynamics, since nothing would happen for a substantial amount of time. Second, the chaser can be absent for a shorter time in the *knowledge* than in the *ignorance* condition, in which case the flow of events – the chasee’s actions and the general scene dynamics – remains natural. We chose the second option because we reasoned that the artificial break in the *knowledge* condition could disrupt the participant’s tracking of the chaser’s epistemic state, thus being a confound that would be more detrimental than the difference in the duration of absence. Further, the current contrast has the advantage that the chasee’s sequence and timing of actions are identical in both conditions, thus minimizing the difference between conditions. Finally, with the current design, the duration of the chaser’s absence will be closely matched in the later planned false belief - true belief contrast, because in the future false belief condition, the chaser has to be absent for fewer events (because the chaser witnesses the first hiding events after the chasee reappeared at the other side of the tunnel).

Trial Randomization. We will vary the starting location of the chasee (left or right half of the upper part of the scene) and the box the chasee ended up (left or right box) in both familiarization and test trials. The presentation of the familiarization trials will be counterbalanced in two pseudo-randomized orders. Each lab signs up for one or two sets of 16-trial-combinations, for each of their tested age groups.

Lab Participation Details

Time-Frame. The contributing labs will start data collection as soon as they are able to once our Registered Report receives an in-principle acceptance. The study will be submitted for Stage 2 review within one year after in-principle acceptance (i.e., post-Stage



Figure 2. Schematic overview of stimuli and conditions of the test trials.

Note. After the familiarization phase, participants know about the agent's goal (chaser wants to find chasee), perceptual access (chaser can see what happens on the other side of the fence), and situational constraints (boxes can be reached by walking through the forking tunnel). In the *knowledge* condition, the chaser witnesses the chasee walking through the tunnel and jumping in and out of the first box. While the chasee is in the box, the chaser briefly leaves the scene through the door in the back and returns shortly after. Subsequently, the chaser watches the chasee jumping out of the box again and hiding in the second box. In the *ignorance* condition, the chaser turns around and stands on the other side of the door in the back of the scene, thus unable to witness any of the chasee's actions. The chaser then returns and enters the tunnel to look for the chasee. During the test phase (4 seconds still frame), AL towards the end of the tunnels is measured.

1 review). We anticipate that this time window gives the individual labs enough flexibility to contribute the committed sample sizes; however, if this timeline needs adjusting due to the Covid-19 pandemic this decision will be made prior to any data analysis.

Participation Criterion. The participating labs were recruited from the MB2 consortium. In July 2020, we asked via the MB2 listserv which labs plan to contribute how many participants for the respective age group (toddlers and/or adults). The Supplemental Material provides an overview of participating labs. Each lab made a commitment to collecting data from at least 16 participants (toddlers or adults), but we will not exclude any contributed data on the basis of the total sample size contributed by that lab. Labs will be allowed to test using either in-lab eye-tracking or online methods.

Ethics. All labs will be responsible for obtaining ethics approval from their appropriate institutional review board. The labs will contribute de-identified data for central data analysis (i.e., eye-tracking raw data/coded gaze behavior, demographic information). Video recordings of the participants will be stored at each lab according to the approved local data handling protocol. If allowed by the local institutional review board, video recordings will be made available to other researchers via the video library DataBrary (<https://nyu.databrary.org/>).

Participants. In a preliminary expression of interest, 26 labs signed up to contribute a minimal sample size of 16 toddlers and/or adults. Based on this information, we expect to recruit a total sample of 520 toddlers (ages 18-27 months) and 408 adults (ages 18-55 years). To avoid an unbalanced age distribution in the toddlers sample, labs will sign up for testing at least one of two age bins (bin 1: 18-22 months, bin 2: 23-27 months), and will be asked to ensure approximately equal distribution of participants' age in their collected sample if possible. They will be asked to try to ensure that the mean age of their sample lies in the middle of the range of the chosen bin and that participant ages are distributed across their whole bin. Both for adults and toddlers, basic demographic data will be collected on a voluntary basis with a brief questionnaire (see Supplemental

Material for details). The requested demographic information that is not used in the registered confirmatory and/or exploratory analyses of this study will be collected for further potential follow-up analyses in spin-off projects within the MB framework. After completing the task, adult participants will be asked to fill a funneled debriefing questionnaire. This questionnaire asks what the participant thinks the purpose of the experiment was, whether the participant had any particular goal or strategy while watching the videos, and whether the participant consciously tracked the chaser’s epistemic state. Additionally, we collect details regarding each testing session (see Supplemental Material).

Our final dataset consisted of 1224 participants, with an overall exclusion rate of 24.16% (toddlers: 35.60%, adults: 12.67%). Tables 1 A. and B. show the distribution of included participants across labs, eye-tracking methods, and ages. A final sample of 521 toddlers (49.14% female) that were tested in 37 labs (mean lab sample size = 14.08, SD = 5.56, range: 2 - 32) was analyzed. The average age of toddlers in the final sample was 22.49 months (SD : 2.53, range: 18 - 27.01). The final sample size of included adults was N = 703 (68.85% female), tested in 34 labs (mean lab sample size = 20.68, SD = 12.14, range: 8 - 65). Their mean age was 24.61 years (SD : 7.36, range: 18 - 55).

Apparatus and Procedure

Eye-tracking Methods. We expect that participating labs will use one of three types of eye-tracker brands to track the participant’s gaze patterns: Tobii, EyeLink, or SMI. Thus, apparatus setup will slightly vary in individual labs (e.g., different sampling rates and distances at which the participants are seated in front of the monitor). Participating labs will report their eye-tracker specifications and study procedure alongside

the collected data. To minimize variation between labs, all labs using the same type of eye-tracker will use the same presentation study file specific to that eye-tracker type. The Supplemental Material will provide an overview of employed eye-trackers, stimulus presentation softwares, sampling rates and screen dimensions.

Online Gaze Recording. To allow for the participation of labs that do not have access to an eye-tracker, or are not able to invite participants to their facilities due to current restrictions regarding the COVID-19 pandemic, labs can choose to collect data via online testing. Specifically, labs may choose to manually code gaze direction during stimulus presentation on a frame-by-frame basis from video recordings of a camera facing the participant (e.g., a webcam). Labs that choose to collect data virtually will utilize the platform of their choice (e.g., LookIt, YouTube, Zoom, Labvanced, etc.). Further, labs may also choose to use webcam eye-tracking with tools like WebGazer.js (Papoutsaki et al., 2016). In our analyses, we control for and quantify potential sources of variability due to these different methods.

Testing Procedure. Toddlers will be seated either on their caregiver’s lap or in a highchair. The distance from the monitor will depend on the data collection method. Caregivers will be asked to refrain from interacting with their child and close their eyes during stimulus presentation or wear a set of opaque sunglasses. Adult participants will be seated on a chair within the respective appropriate distance from the monitor. Once the participant is seated, the experimenter will initiate the eye-tracker-specific calibration procedure. Additionally, we will present another calibration stimulus before and after the presentation of the task. This allows for evaluating the accuracy of the calibration procedure across labs (cf., Frank, Vul, & Saxe, 2012).

General Lab Practices

To ensure standardization of procedure, materials for testing practices and instructions will be prepared and distributed to the participating labs. Each lab will be

responsible for maintaining these practices and report all relevant details on testing sessions (for details see the Supplemental Material).

Videos of Participants. As with all MB projects, we strongly encourage labs to record video data of their own lab procedures and each testing session, provided that this is in line with regulations of the respective institutional ethics review board and the given informed consent. Participating labs that cannot contribute participant videos will be asked to provide a video walk-through of their experimental set-up and procedure instead. If no institutional ethics review board restrictions occur, labs are encouraged to share video recordings of the test sessions via DataBrary.

Design Analysis

Here we provide a simulation of the predicted findings because a traditional frequentist power analysis is not applicable for our project for two reasons. First, we use Bayesian methods to quantify the strength of our evidence for or against our hypotheses, rather than assessing the probability of rejecting the null hypothesis. In particular, we compute a Bayes factor (BF; a likelihood ratio comparing two competing hypotheses), which allows us to compare models. Second, because of the many-labs nature of the study, the sample size will not be determined by power analysis, but by the amount of data that participating labs are able to contribute within the pre-established timeframe. Even if the effect size is much smaller than what we anticipate (e.g., less than Cohen's $d = 0.20$), the results would be informative as our study is expected to be dramatically larger than any previous study in this area. If, due to unforeseen reasons, the participating labs will not be able to collect a minimum number of 300 participants per age group within the proposed time period, we plan to extend the time for data collection until this minimum number is reached. Or in contrast, if the effect size is large (e.g., more than Cohen's $d = 0.80$), the resulting increased precision of our model will allow us to test a number of other theoretically and methodologically important hypotheses (see Results section). Although

we did not determine our sample size based on power analysis, here we provide a simulation-based design analysis to demonstrate the range of BFs we might expect to see, given a plausible range of effect sizes and parameters. We focus this analysis on our key analysis of the test trials (as specified below), namely the difference in AL on the first test trial that participants saw. We describe below the simulation for the child sample, but based on our specifications, we expect that a design analysis for adult data would produce similar results. We first ran a simulation for the first look analysis. In each iteration of our simulation, we used a set of parameters to simulate an experiment, using a first look (described below) as the key measure. For the key effect size parameter for condition (*knowledge* vs. *ignorance*), we sampled a range of effect sizes in logit space spanning from small to large effects (Cohen's $d = 0.20 - 0.80$; log odds from 0.36 - 1.45). For each experiment, the betas for age and the age x condition interaction were sampled uniformly between -0.20 and 0.20. The age of each participant was sampled uniformly between 18 and 27 months and then centered. The intercept was sampled from a normal distribution (1, 0.25), corresponding to an average looking proportion of 0.73. Lab intercepts and the lab slope by condition were set to 0.1, and other lab random effects were set to 0 as we do not expect them to be meaningfully non-zero. These values were chosen based on pilot data (average looking proportion), but also to have a large range of possible outcomes (lab intercept, age and age x condition interaction). We are confident that the results would be robust to different choices. We then used these simulated data to simulate an experiment with 22 labs and 440 toddlers and computed the resulting BFs, as specified in the analysis plan below. We adopted all of the priors specified in the results section below³. We ran 349 simulations and, in 72% of them, the BF showed strong evidence in favor of the full model

³ After the design analysis, additional labs expressed their interest in contributing data, which is why the anticipated sample sizes and the numbers this design analysis is based on differ. Given the uncertainty in determining the final sample size in this project, we kept the design analysis as is to have a more conservative estimate of the study's power.

(BF > 10); in 6% the BF showed substantial evidence ($10 > \text{BF} > 3$); it was inconclusive 14% of the time ($1/10 > \text{BF} > 3$), and in 8% of cases the null model was substantially favored (see Supplement). In none of the simulations the BF was $< 1/10$. Thus, under the parameters chosen here for our simulations, it is likely that the planned experiment is of sufficient size to detect the expected effect. We also ran a design analysis for the proportional looking analysis. We used the same experimental parameters (number of labs, participants, ages, etc.). For generating simulated data, we drew the condition effect from a uniform distribution between .05 and .20 (in proportion space). The age and age:condition effects were drawn from uniform distributions between -.05 and .05. Sigma, the overall noise in the experiment, was drawn from a uniform distribution between .05 and .1. The intercept was drawn from a normal distribution with mean .65 and a standard deviation of .05. The by-lab standard deviation for the intercept and condition slope was set to .01. Priors were as described in the main text. We ran 119 simulations, and in all 119 we obtained a BF greater than 10, suggesting that, under our assumptions, the study is well-powered.

Data Preprocessing

Eye-tracking. Raw gaze position data (x- and y-coordinates) will be extracted in the time window starting from the first frame at which the chaser enters the tunnel until the last frame before it exits the tunnel in the last familiarisation trial and in the test trial. For data collected from labs using a binocular eye-tracker, gaze positions of the left and the right eye will be averaged. We will use the peekds R package (<http://github.com/langcog/peekds>) to convert eye-tracking data from disparate trackers into a common format. Because not all eye-trackers record data with the same frequency or regularity, we will resample all data to be at a common rate of 40 Hz (samples per second). We will exclude individual trials if more than 50% of the gaze data is missing (defined as off-screen or unavailable point of gaze during the whole trial, not just the anticipatory

period). Applying this criterion would have caused us to exclude 4% of the trials in our pilot data, which inspection of our pilot data suggested was an appropriate trade-off between not excluding too much usable data and not analyzing trials which were uninformative. For each monitor size, we will determine the specific AOIs and compute whether the specific x- and y-position for each participant, trial, and time point fall within their screen resolution-specific AOIs. Our goal is to determine whether participants are anticipating the emergence of the chaser from one of the two tunnel exits. Thus, we defined AOIs on the stimulus by creating a rectangular region around the tunnel exit that is D units from the top, bottom, left, and right of the boundary of the tunnel exit, where D is the diameter of the tunnel exits. We then expanded the sides of the AOI rectangles by 25% in all directions to account for tracker calibration error. Our rationale was that, if we made the AOI too small, we might fail to capture anticipations by participants with poor calibrations. In contrast, if we made the regions too large, we might capture some fixations by participants looking at the box where the chasee actually is. On the other hand, these chasee looks would not be expected to vary between conditions and so would only affect our baseline level of looking. Thus, the chosen AOIs aim at maximizing our ability to capture between-condition differences. For an illustration of the tunnel exit AOIs see Figure 4. We are not analyzing looks to the boxes, since they can less unambiguously be interpreted as epistemic state-based action predictions and because we observed few anticipatory looks to the boxes in the pilot studies. For more detailed information about the AOI definition process see the description of the pilot study results in the Supplemental Material.

Manual Coding. For data gathered without an eye-tracker (e.g., videos of participants gathered from online administration), precise estimation of looks to specific AOIs will not be possible. Instead, videos will be coded for whether participants are looking to the left or the right side of the screen (or “other/off screen”). In our main analysis, during the critical anticipatory window, we will treat these looks identically to looks to the corresponding AOI. See exploratory analyses for analysis of data collected online.

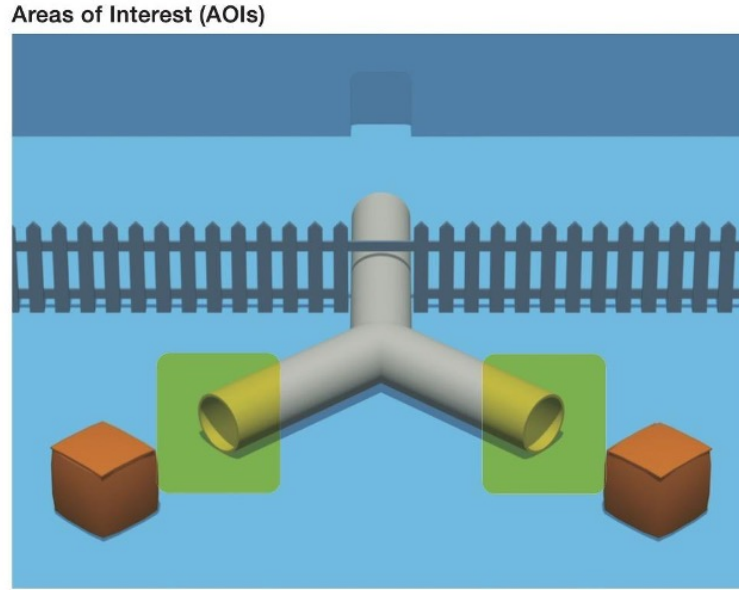


Figure 3. Illustration of Areas of Interest (AOIs) for gaze data analysis during the anticipatory period.

Note. The light green rectangles show the dimensions of the AOIs used for the analysis of AL during the test period.

Temporal Region of Interest. For familiarization trials, we define the start of the anticipatory period (total length = 4000 ms) as starting 120 ms after the first frame after which the chaser has completely entered the tunnel and lasting until 120 ms after the first frame at which the chaser is visible again [we chose 120 ms as a conservative value for cutting off reactive saccades; cf., Yang, Bucci, and Kapoula (2002)]. For test trials, we define the start of the anticipatory period in the same way, with a total duration of 4000 ms.

Dependent Variables. We define two primary dependent variables: 1. First look. First saccades will be determined as the first change in gaze occurring within the anticipatory time window that is directed towards one of the AOIs. The first look is then the binary variable denoting the target of this first saccade (i.e., either the correct or incorrect AOI) and is defined as the first AOI where participants fixated at for at least 150

ms, as in rayner2009eye. The rationale for this definition was that, if participants are looking at a location within the tunnel exit AOIs before the anticipation period, they might have been looking there for other reasons than action prediction. We therefore count only looks that start within the anticipation period because they more unambiguously reflect action predictions. This further prevents us from running into a situation where we would include a lot of fixations on regions other than the tunnel exit AOIs because participants are looking somewhere else before the anticipation period begins. 2. Proportion DLS [also referred to as total relative looking time; Senju et al. (2009)]. We compute the proportion looking (p) to the correct AOI during the full 4000 ms anticipatory window (correct looking time / (correct looking time + incorrect looking time)), excluding looks outside of either AOI.

Results

Confirmatory Analyses

Approach. As discussed in the Methods section, we adopted a Bayesian analysis strategy so as to maximize our ability to make inferences about the presence or absence of a condition effect (i.e., our key effect of interest). In particular, we fit Bayesian mixed effects regressions using the package brms in R (Bürkner, 2017). This framework allows us to estimate key effects of interest while controlling for variability across grouping units (in our case, labs). To facilitate interpretation of individual coefficients, we report means and credible intervals. For key inferences in our confirmatory analysis, we use the bridge sampling approach (Gronau et al., 2017) to compute BFs comparing different models. As the ratio of the likelihood of the observed data under two different models, BFs allow us to quantify the evidence that our data provide with respect to key comparisons. For example, by comparing models with and without condition effects, we can quantify the strength of the evidence for or against such effects. Bayesian model comparisons require the

specification of proper priors on the coefficients of individual models. Here, for our first look analysis, we use a set of weakly informative priors that capture the expectation that the effects that we observe (of condition and, in some cases, trial order) are modest. For coefficients, we choose a normal distribution with mean of 0 and *SD* of 2. Based on our pilot testing and the results of MB1, we assume that lab and participant-level variation will be relatively small, and so for the standard deviation of random effects (i.e., variation in effects across labs and, in the case of the familiarization trials, participants) we set a Normal prior with mean of 0 and *SD* of 0.1. We set an LKJ(2) prior on the correlation matrix in the random effect structure, a prior that is commonly used in Bayesian analyses of this type (Bürkner, 2017). Because the BF is sensitive to the choice of prior, we also ran a secondary analysis with a less informative prior: fixed effect coefficients chosen from a normal distribution with mean 0 and *SD* of 3, and random effect standard deviations drawn from a normal prior with a mean of 0 and *SD* of 0.5 (see Supplement S3). With respect to the specification of random effects, we followed the approach advocated by Barr (2013), that is, specifying the maximal random effect structure justified by our design. Since we are interested in lab-level variation, we will fit random effect coefficients for fixed effects of interest within labs (e.g., condition within lab). Further, where there were participant-level repeated measure data (e.g., familiarization trials), we fitted random effects of participants. For the proportional looking score analysis, we used a uniform prior on the intercept between -0.5 and 0.5 (corresponding to proportional looking scores between 0 and 1: the full possible range). For the priors on the fixed effect coefficients, we used a normal prior with a mean of 0 and an *SD* of 0.1. Because these regressions are in proportion space, 0.10 corresponds to a change in proportion of 10%. For the random effect priors, we used a normal distribution with mean 0 and standard deviation .05. The LKJ prior was specified as above.

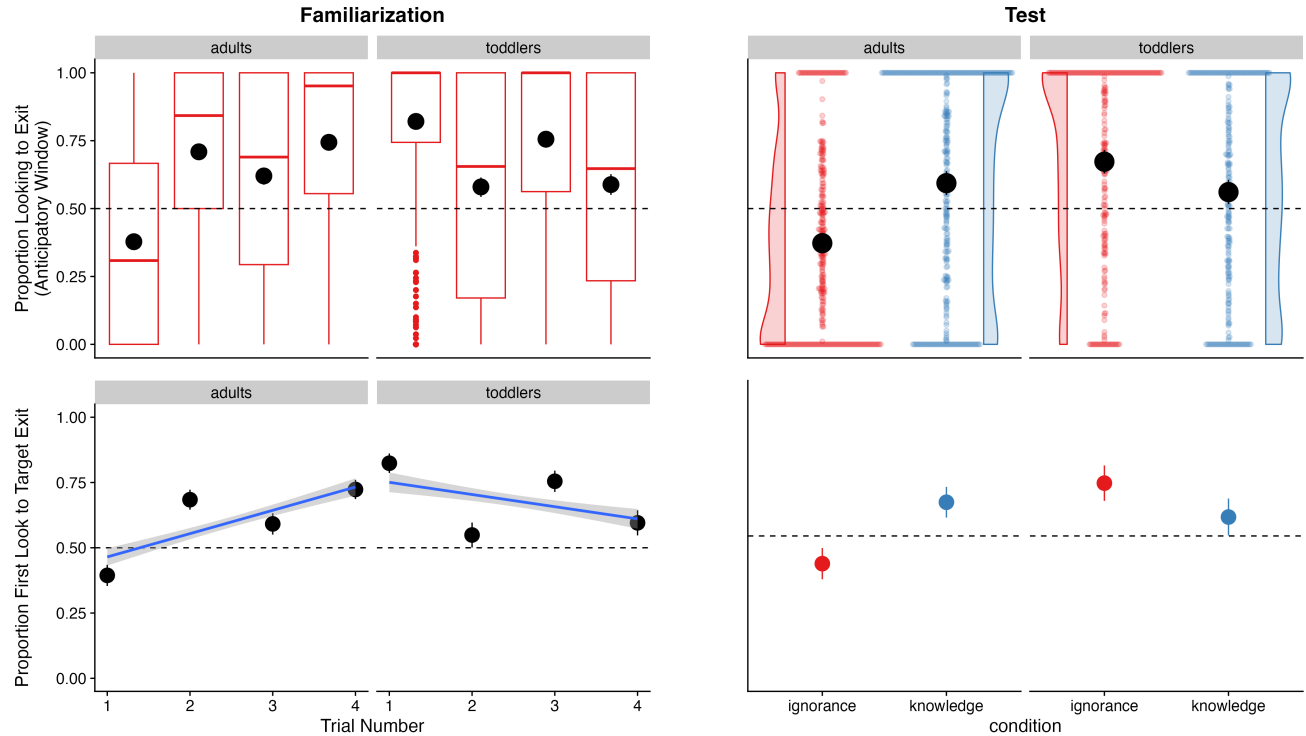


Figure 4. Proportional target looking and proportion of first looks for toddlers and adults during familiarization and test.

Familiarization Trials. Figure 5 shows the proportion of total relative looking time (non-logit transformed) and proportion of first looks for toddlers and adults plotted across familiarization trials and test trials. Our first set of analyses examined data from the four familiarization trials and asked whether participants anticipated the chaser’s reappearance at one of the tunnel exits. In our first analysis, we were interested in whether participants engage in AL during the familiarization trials. To quantify the level of familiarization, we fitted Bayesian mixed effect models predicting target looks based on trial number (1-4) with random effects for lab and participants and random slopes for trial number for each. In R formula notation (which we adopt here because of its relative concision compared with standard mathematical notation), our base model was as follows:

$$measure \sim 1 + trial_number + (trial_number|lab) + (trial_number|participant).$$

We

fitted a total of four instances of this model, one for each age group (toddlers vs. adults) and dependent measure (proportion looking score vs. first look). First look models were fitted using a logistic link function. The proportion looking score models were Gaussian. Our key question of interest was whether overall anticipation is higher than chance levels on the familiarization trial immediately before the test trials, in service of evaluating the evidence that participants are attentive and making predictive looks immediately prior to test. To evaluate this question across the four models, we coded trial number so that the last trial before the test trials (trial 4) was set to the intercept, allowing the model intercept to encode an estimate of the proportion of correct anticipation immediately before test. We then fitted a simpler model for comparison

$$measure \sim 0 + trial_number + (trial_number|lab) + (trial_number|participant),$$

which included no intercept term. We then computed the BF comparing this model to the full model. This BF quantified the evidence for an anticipation effect for each group and measure.

Proportion of total relative looking time.

Toddlers. We used a Bayesian mixed effects models to predict PTL based on trial number (1-4) for toddlers, with random effects for lab and participants and random slopes for trial number for each. The Bayes factor comparing this model to the simpler null model without the intercept was estimated to be $BF > 1000$, strongly favoring the full model over the null model. See also Table 3 for regression coefficients for the full model. These results suggest a significant effect of trial number on PTL, with the negative coefficient indicating a decrease in PTL across the familiarization trials.

Adults. Next, we used a Bayesian mixed effects model to predict PTL based on trial number (1-4) for adults, again with random effects for lab and participants and random slopes for trial number for each. The Bayes factor for the full model against the null model was $BF > 1000$, suggesting strong evidence for the full model. These results suggest a significant effect of trial number on PTL, with the positive coefficient indicating an

increase in target looks across the familiarization trials.

Proportion of first looks.

Toddlers. Investigating proportion of first looks to the target location for toddlers, we again used a Bayesian mixed effects model to predict whether toddlers first look was to the target exit based on trial number (1-4), with random effects for lab and participants and random slopes for trial number for each. The Bayes factor comparing the full model to the simpler model was estimated to be $BF = 15.9$, favoring the full model over the null model. The model also provided support for an effect of trial number on proportion of first looks, with the negative coefficient indicating a decrease in target looks across the familiarization trials.

Adults. Comparing the Bayesian mixed effects model of adults predicting proportion of first looks based on trial number (1-4), with random effects for lab and participants and random slopes for trial number for each with the simpler model without an intercept, we computed a Bayes factor of $BF > 1000$, strongly favoring the full model over the null model. There was again support for an effect of trial number on proportion of first looks, with the positive coefficient indicating an increase in proportion of first target looks across the familiarization trials.

Test Trials. We focused our confirmatory analysis on the first test trial (see Exploratory Analysis section for an analysis of both trials). Our primary question of interest was whether AL differs between conditions (knowledge vs. ignorance, coded as $-.5/.5$) and by age (in months, centered). For child participants, we fitted models with the specification:

$$measure\ 1 + condition + age + condition : age + (1 + condition + age + condition : age | lab).$$

For adult participants, we fitted models with the specification

$$measure\ 1 + condition + (1 + condition | lab).$$

Again, we fitted models with a logistic link for first look analyses and with a standard linear link for DLS. In each case, our key BF

was a comparison of this model with a simpler “null” model that did not include the fixed effect of condition but still included other terms. We take a $BF > 3$ in favor of a particular model as substantial evidence and a $BF > 10$ in favor of strong evidence. A $BF < 1/3$ is taken as substantial evidence in favor of the simpler model, and a $BF < 1/10$ as strong evidence in favor of the simpler model. For the model of data from toddlers, we additionally were interested in whether the model shows changes in AL with age. We assessed evidence for this by computing BFs related to the comparison with a model that did not include an interaction between age and condition as fixed effects

$$measure\ 1 + condition + age + (1 + condition + age + condition : age|lab).$$

These BFs captured the evidence for age-related changes in the difference in action anticipation between the two conditions. It is important to note that in the case of a null effect, there are two main explanations: (1) toddlers and adults in our study do not distinguish between knowledgeable and ignorant agents when predicting their actions. (2) The method used is not appropriate to reveal knowledge/ignorance understanding. By using Bayesian analyses, we are able to better evaluate the first of these two possibilities: The BF provides a measure of our statistical confidence in the null hypothesis, i.e., no difference between experimental conditions, given the data in ways that standard null hypothesis significance testing does not. In other words, instead of merely concluding that we did not find a difference between conditions, we would be able to find no/anecdotal/moderate/strong/very strong/extreme evidence for the null hypothesis that our participants did not distinguish between knowledgeable and ignorant agents when predicting their actions (Schönbrodt & Wagenmakers, 2018). We therefore consider this analysis an important addition to our overall analysis strategy. Yet, even our Bayesian analyses are not able to rule out the second possibility that participants may well show such knowledge/ignorance understanding with different methods, or that this ability may not be measurable with any methods available at the current time. Addressing this alternative explanation warrants follow up experiments.

Proportion of total relative looking time.

Toddlers. As first model, we used a Bayesian mixed effects models to predict toddlers' PTL based on condition, age, and the interaction of condition and age, while accounting for variability across labs. The Bayes factor comparing this model to the simpler null model without the interaction of condition was estimated to be $BF = 23.1$, favoring the full model over the null model. Table 4 shows the statistics for regression coefficients of the full model. These results suggest a significant effect of condition on PTL, with the positive coefficient indicating higher PTL for ignorance trials compared to knowledge trials.

Additionally, we assessed whether toddlers' AL changed with age. Comparing our full model, which included an interaction between age and condition, with a simpler model without this interaction yielded a Bayes factor, $BF = 0.4$, providing modest support for the simpler model. This result suggests that the interaction between age and condition might not be a necessary predictor, as it doesn't provide substantial additional explanatory power. This implies that age-related changes in AL are likely consistent across conditions, rather than differing between them.

Adults. Next, we used a Bayesian mixed effects model to predict PTL based on condition for adults, again with random effects for lab. The Bayes factor comparing this model to the simpler null model without the main effect of condition was estimated to be $BF > 1000$, strongly favoring the full model over the null model. These results suggest a significant main effect of condition on PTL, with the negative coefficient indicating a higher number of target looks for knowledge than for ignorance trials.

Proportion of first looks.

Toddlers. Investigating proportion of first looks for toddlers, we again used a Bayesian mixed effects model to predict target looks based on condition, with random effects for lab. The Bayes factor comparing the full model to the simpler model was estimated to be $BF = 2.5$, providing no substantial evidence in favor of the full model over

the null model.

Again, we examined whether age influenced the difference in action anticipation between knowledge and ignorance trials. To do this, we compared the full model, which included an interaction between age and condition, with a simpler model without this interaction. The computed Bayes factor, $BF = 0.0$, strongly supports the simpler model, suggesting that the interaction term does not substantially improve the model's fit. This implies that age does not appear to significantly affect the difference in action anticipation between the two trial types.

Adults. We compared a Bayesian mixed-effects model predicting the proportion of first looks based on condition, including random effects for lab to a simpler model without the main effect of condition. The analysis yielded a Bayes factor of $BF > 1000$, providing strong evidence in favor of the full model over the null model. Results indicated that first looks to the target were significantly more frequent in the knowledge condition compared to the ignorance condition.

Exploratory Analyses

Spill-over. We will analyze within-participants data from the second test trial that participants saw, using exploratory models to assess whether (1) findings are consistent when both trials are included (overall condition effect), (2) whether effects are magnified or diminished on the second trial (order main effect), and (3) whether there is evidence of “spillover” - dependency in anticipation on the second trial depending on what the first trial is (condition x order interaction effect).

Analyzing condition-effects of within-participants data for both test trials, we fitted a Bayesian mixed-effects model with the dependent variable of PTL and main effects of condition and age and their interaction for toddlers. Comparing this full model to a null model that did not include the fixed effect of condition, we obtained a Bayes Factor of BF

838 = 0.0, providing strong evidence in favor of the null model.

839 For adults, we also fitted a Bayesian mixed-effects model to predict their PTL for
840 both test trials with the main effect of condition and random effects for participant and
841 lab. Again, the data provided very strong evidence for the inclusion of the main effect of
842 condition with a Bayes Factor of $BF > 1000$. The effect of condition was negative and
843 credible, suggesting that PTL was significantly lower in the ignorance condition compared
844 to the knowledge condition.

845 In order to investigate whether there's an interaction of condition and test trial
846 number, we fitted Bayesian mixed-effects model to predict PTL with fixed effects for
847 condition, test trial number, and their interaction, along with random intercepts and slopes
848 for these variables across labs, for toddlers and adults separately. For toddlers, the Bayes
849 factor, $BF = 0.4$, modestly favored the simpler null model without the interaction term,
850 indicating that the interaction between condition and test trial number does not add
851 substantial explanatory power to the model. These results suggest that neither condition
852 nor its interaction with test trial number significantly impacts PTL in this sample.

853 For adults, the Bayes Factor, $BF = 19.7$, provided strong evidence for including the
854 interaction of condition and test trial number as fixed effect. These results indicate that
855 while PTL increased over trials, this effect was moderated by condition, with the ignorance
856 condition showing a slower rate of increase compared to the knowledge condition.

857 To examine whether anticipatory looking during the second test trial influenced by
858 condition and anticipatory looking during the first test trial, we fitted a Bayesian
859 mixed-effects model for each age cohort separately. This model included fixed effects for
860 condition, proportion of target looking during the first test trial, and their interaction.
861 Random intercepts and slopes for these predictors were modeled at the lab level. The
862 Bayes factor, $BF = 1.1$, suggests negligible evidence in favor of including these predictors
863 compared to the null model, indicating that these factors may not strongly influence second

trial anticipatory looking in toddlers. There was a small and non-significant positive main effect of condition, indicating minimal differences in anticipatory looking between conditions and a negligible and non-significant effect of first test trial anticipatory looking. The interaction between condition and first test trial anticipatory looking was also minimal and non-significant. These findings indicate that condition, first trial anticipatory looking, and their interaction do not strongly predict anticipatory looking during the second test trial in toddlers. The Bayes factor close to 1 reflects weak or inconclusive evidence for the inclusion of these predictors, suggesting that the variability in second trial behavior may arise from other unmodeled factors. For adults, the Bayes factor, $BF > 1000$, strongly supports the inclusion of these predictors, suggesting that condition and first trial behavior substantially explain second trial anticipatory looking. The regression results showed a significant negative main effect of condition, indicating reduced anticipatory looking in the ignorance condition compared to the knowledge condition. There was a small negative effect of first trial anticipatory looking, suggesting that higher anticipatory looking in the first test trial is slightly associated with reduced looking in the second test trial. The interaction between condition and first test trial anticipatory looking was negligible and non-significant. These findings indicate that condition strongly impacts anticipatory looking during the second test trial, while anticipatory looking during the first trial has a smaller and more nuanced influence. The interaction between condition and first trial looking appears minimal. The extremely large Bayes factor underscores the importance of considering these predictors in explaining second test trial anticipatory looking behavior.

Relationship between familiarization and test. We will explore whether condition differences vary for participants who show higher rates of anticipation during the four familiarization trials. For example, we might group participants according to whether they did or did not show correct AL at the end of the familiarization phase, defined as overall longer looking at the correct AOI than the incorrect AOI on average in trials 3 and 4 of the familiarization phase.

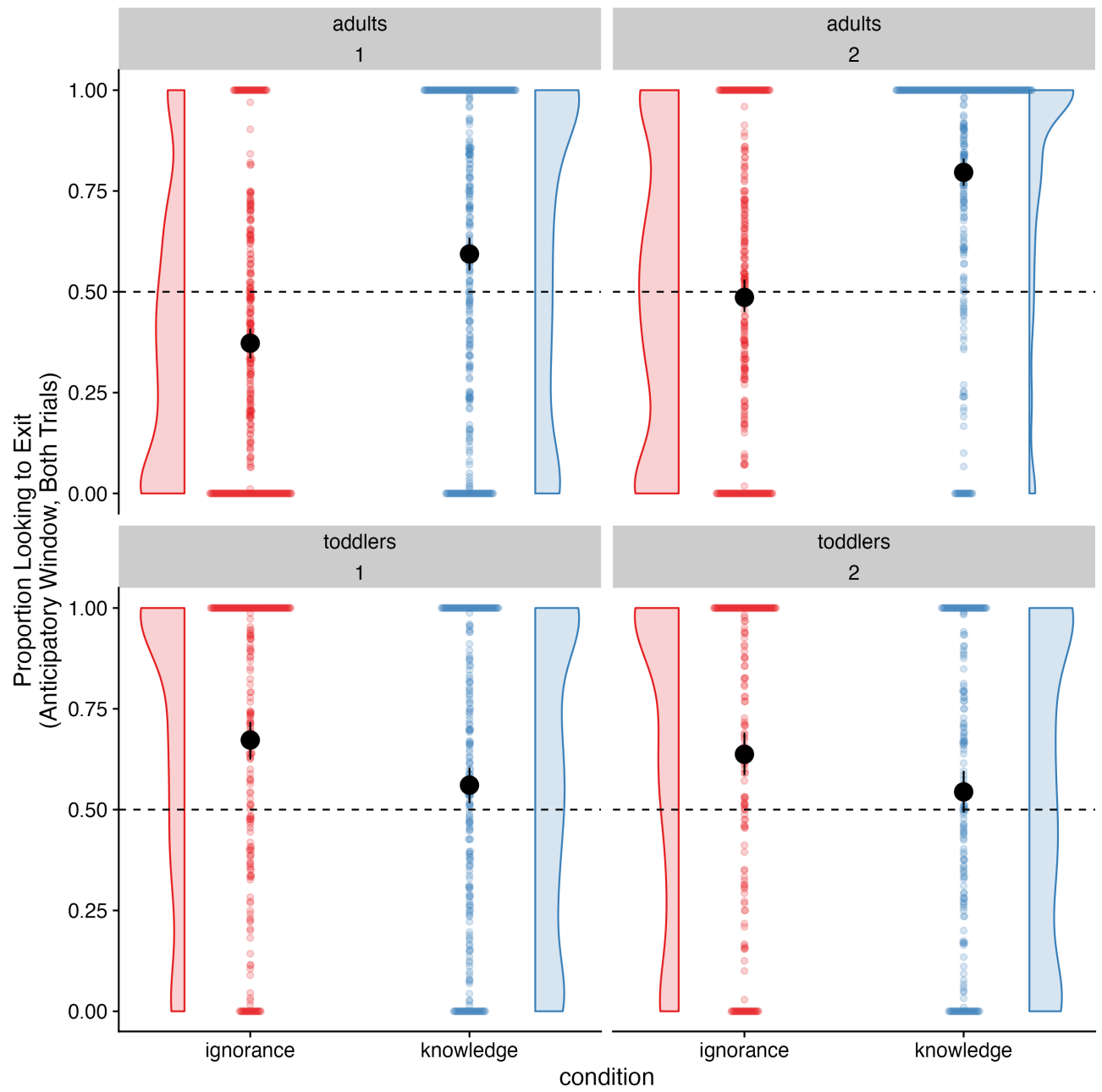


Figure 5. Proportional exit looking for the first and second test trial for toddlers and adults in the ignorance and knowledge condition.

To investigate whether participants who show anticipatory looking during the familiarization phase also exhibit anticipatory looking during the test phase, we explored three different measures. First, we assessed anticipatory looking in participants who successfully anticipated during the final familiarization trial, defined as those whose first fixation was on the target. Second, we examined anticipatory looking in participants who consistently demonstrated anticipatory behavior across all familiarization trials, operationalized as having a PTL greater than 0.5 in each trial. Finally, we computed correlations to explore whether performance in the familiarization phase was related to performance in the test trials.

Relationship between anticipatory looking during the first test trial and first look during final familiarization trial. We fitted a main Bayesian hierarchical model testing the fixed effects of condition (ignorance vs. knowledge), first look during the final familiarization trial (target vs. distractor), and their interaction on first-trial proportion target looking during the anticipatory window for toddlers and adults separately. Random intercepts and slopes for all fixed effects and their interaction were included at the lab level, accounting for variability across different experimental settings. For toddlers, the Bayes factor comparing this model to the simpler null model without the interaction of condition and first look during the final familiarization trial indicated that the data slightly favored the simpler null model over the full model, $BF = 0.7$. The effect of condition was positive, but its confidence interval narrowly included zero, suggesting weak evidence for a condition effect (see Table X). The effect of performance during the final familiarization trial was close to zero, indicating no substantial main effect of prior performance. Similarly, the interaction between condition and performance in the final familiarization trial was small and non-significant. These results suggest that while there is some weak evidence for a main effect of condition on anticipatory looking, neither performance during the final familiarization trial nor its interaction with condition substantially predicted anticipatory looking during the test trial. This result indicates that

the relationship between anticipatory looking and prior familiarization performance does not depend significantly on condition. In other words, toddlers' anticipatory looking during the test trial is likely independent of any conditional effects related to their performance in the familiarization phase. For adults, the Bayes factor comparing this model to the simpler null model without the main effect of condition was estimated to be $BF > 1000$, strongly favoring the base model over the null model. The regression coefficients (see Table X) showed a significant negative effect of condition, indicating that anticipatory looking was lower in ignorance trials compared to knowledge trials. The decisive Bayes factor strongly favors the inclusion of condition and familiarization trial performance in the model, suggesting that these predictors are relevant for understanding anticipatory looking in adults. However, the small and non-significant estimates for the effects of familiarization trial performance and its interaction with condition imply that condition is the primary driver of anticipatory looking differences, with performance in familiarization trials contributing minimally.

Only >50% looking to target during familiarization trials. To examine the effect of condition and successful anticipatory looking during familiarization (above 50% target looking during all familiarization trials) on anticipatory looking during the first test trial, we fitted Bayesian mixed-effects models for each age group separately. The models included fixed effects for condition, anticipatory looking during familiarization trials, and their interaction. Random intercepts and slopes for these predictors were included at the lab level. Comparing the full model to the null model of toddlers revealed a Bayes Factor of $BF = 10.9$, providing moderate evidence favoring the full model over a null model that excludes these predictors, suggesting that these factors contribute meaningfully to explaining the variance in test trial anticipatory looking. The regression analysis showed a positive main effect of condition, indicating higher anticipatory looking in one condition compared to the other. There was a small positive, but non-significant, effect of successful anticipatory looking during familiarization. The interaction between condition and

successful anticipatory looking during familiarization was also small and non-significant. These results indicate that condition is a meaningful predictor of anticipatory looking during test trials in toddlers, with participants showing different levels of anticipatory looking based on condition. However, the successful anticipatory looking during familiarization trials and its interaction with condition appear to have minimal additional impact. The moderate Bayes factor further supports the importance of including these predictors in the model but highlights that condition remains the primary driver of test trial differences. The estimated Bayes factor in favor of the full model of adults over the null model was $BF > 1000$, indicating that the predictors substantially contribute to explaining test trial anticipatory looking. The regression coefficients revealed a significant main effect of condition, with participants showing lower anticipatory looking in the ignorance condition compared to the knowledge condition. There was a small, positive, and non-significant effect of successful anticipatory looking during familiarization and the interaction between condition and successful successful anticipatory looking during familiarization was negligible. These results indicate that condition has a substantial and meaningful impact on anticipatory looking during the first test trial in adults, while successful anticipatory looking in familiarization trials and its interaction with condition have limited additional influence. The extremely large Bayes factor highlights the strong explanatory power of including these predictors in the model, although condition remains the primary driver of the observed differences.

Correlation between familiarization and test. We also examined the correlation between familiarization and test performance across the two age cohorts and conditions (see Figure 7). While no significant correlations were found for adults in either condition, toddlers in the knowledge condition exhibited a significant positive correlation of anticipatory looking in familiarization and test, $r=0.15$, $t(254)=2.35$, $p=0.02$.

Looking patterns during mouse’s change of location. To examine whether participants monitor both the bear and the mouse during the mouse’s location change, and

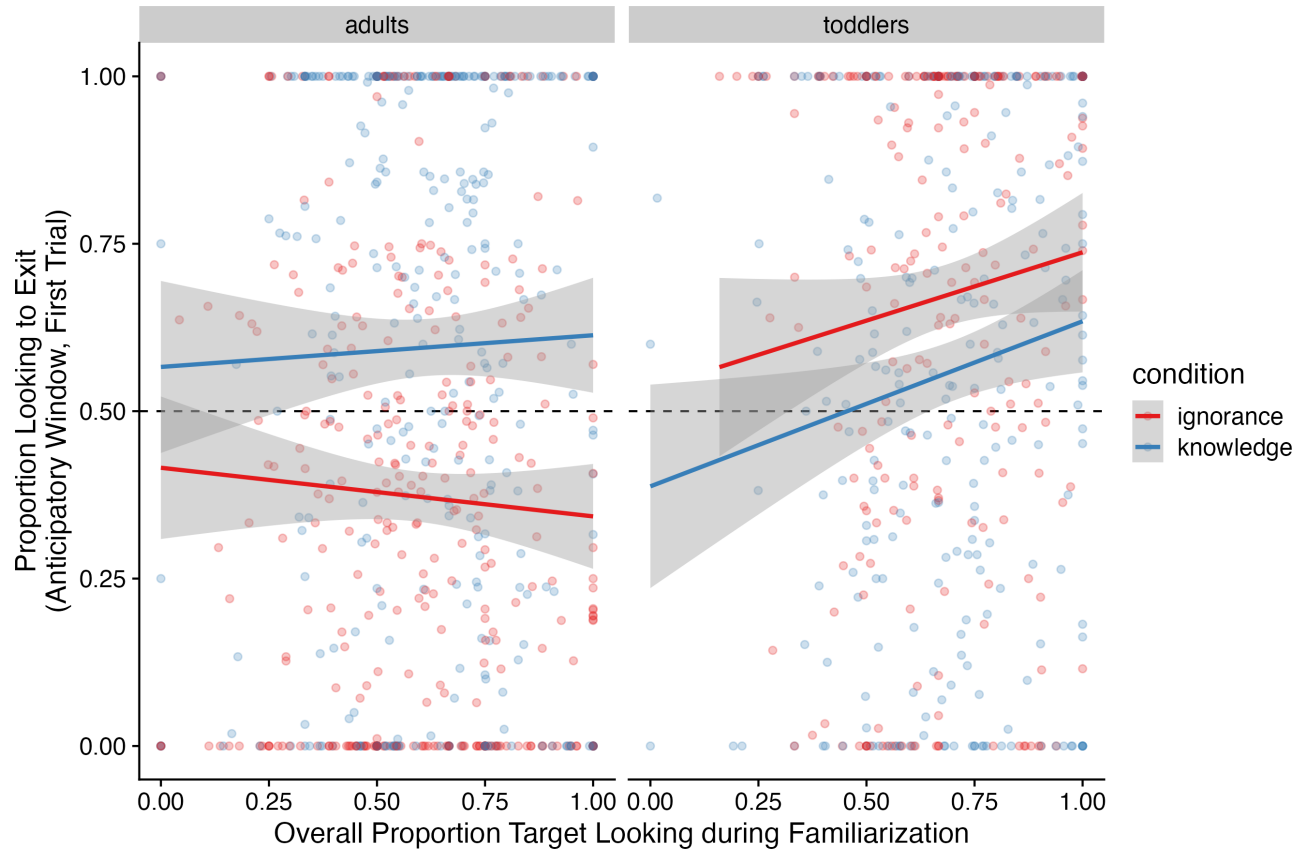


Figure 6. Relationship of anticipatory looking during familiarization and test for both age cohorts and conditions.

how this may influence AL in the test phase, we defined new time windows of interest (TOIs) corresponding to the mouse's location change in each condition and areas of interest (AOIs) for both the mouse and bear. We hypothesized that participants who attend to both AOIs will exhibit greater AL compared to those who predominantly track the mouse during its location change. Specifically, we analyzed the frequency of gaze shifts between the mouse and bear mouse's location change. An additional exploratory analysis of differential gaze duration directed toward mouse and bear during the mouse's location change is provided in the Supplement S3.

Comparing the number of shifts of toddlers and adults during the location change of the mouse. We fitted a Bayesian mixed-effects model to examine

the relationship between the number of shifts between mouse and bear and age cohort during location change of the mouse, while accounting for random effects by lab. The effect of condition was negative and approached significance, suggesting a potential reduction in the number of shifts for the ignorance condition compared to the knowledge condition. The main effect of age cohort was positive and credible, Estimate=0.56, indicating that the number of shifts was higher for adults than for toddlers. Importantly, the interaction between condition and age cohort was negative and credible, indicating that the negative effect of condition was more pronounced in the adult cohort (see Figure 8). Comparing this model to a simpler model without the interaction of condition and age cohort, a Bayes Factor of $BF > 1000$ was computed. This provides strong evidence in favor of including the interaction of condition and age cohort in the model.

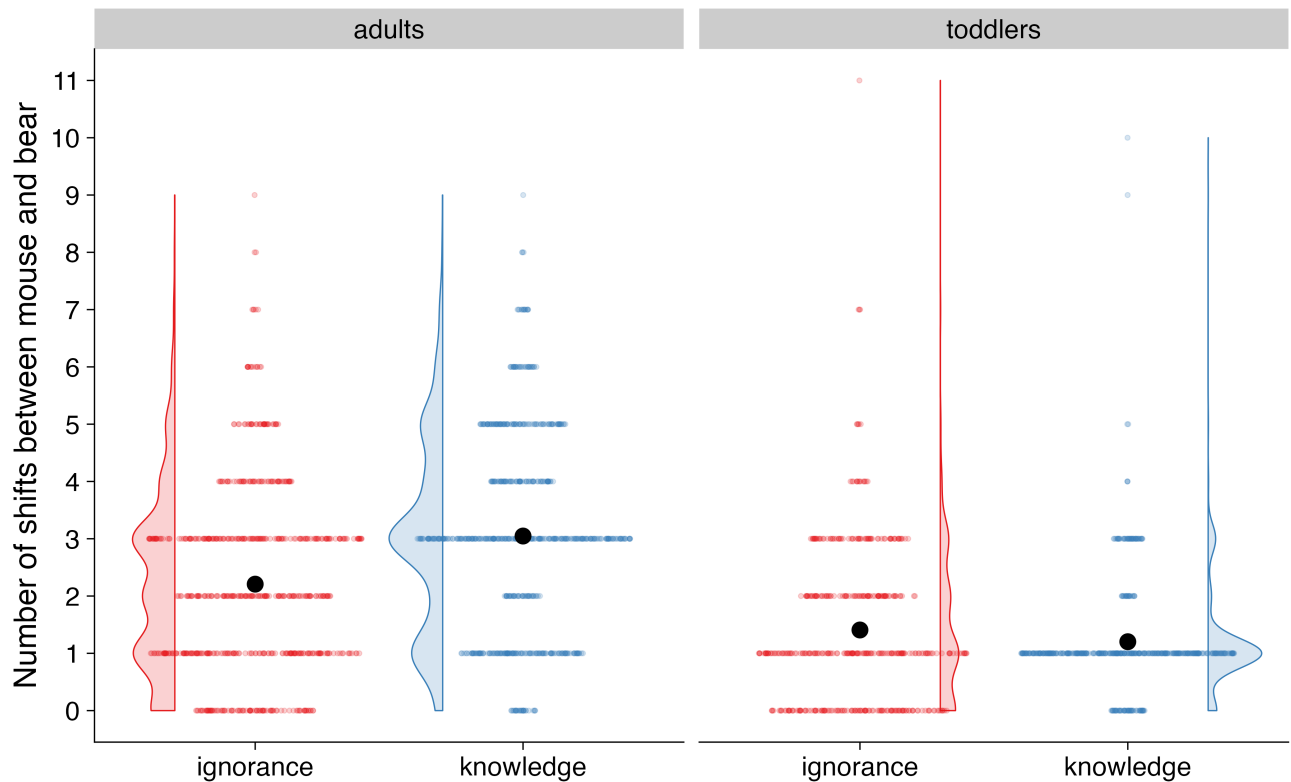


Figure 7. Number of shifts between mouse and bear during location change of mouse in the test phase for toddlers and adults in the ignorance and knowledge condition.

AL as a function of number of gaze shifts between mouse and bear during location change. In order to examine the effect of condition and the number of shifts between mouse and bear during location change of the mouse on anticipatory looking, we fitted Bayesian mixed-effects models for each age cohort separately. The dependent variable was PTL in the anticipation period. The fixed effects included the main effects of condition, the number of shifts, and their interaction. We also included random intercepts and slopes for number of shifts within each participant and within each lab, allowing us to account for the hierarchical structure of the data and potential variability between labs and participants.

For toddlers, comparing this model to a simpler model without the interaction of condition and number of shifts, a Bayes Factor of $BF = 0.0$ was computed, indicating that the data strongly favors the null model over the full model. This suggests that the predictors number of shifts and the interaction with condition included in the full model do not improve the explanation of the observed data compared to the null model.

For adults, the number of shifts showed a small but credible positive effect, suggesting that more shifts were associated with an increase in PTL. The interaction between condition and the number of shifts was negative and credible, indicating that the effect of condition became more negative as the number of shifts increased. The estimated Bayes factor comparing the full model to the null model was approximately $BF > 1000$, providing strong evidence in favor of the full model over the null model.

General Discussion

The current large-scale, multi-lab study set out to examine whether toddlers and adults engage in spontaneous ToM. In particular, we used an anticipatory looking paradigm to explore whether 18- to 27-month-old toddlers and adults distinguish between two basic forms of epistemic states: knowledge and ignorance. Our call for participation

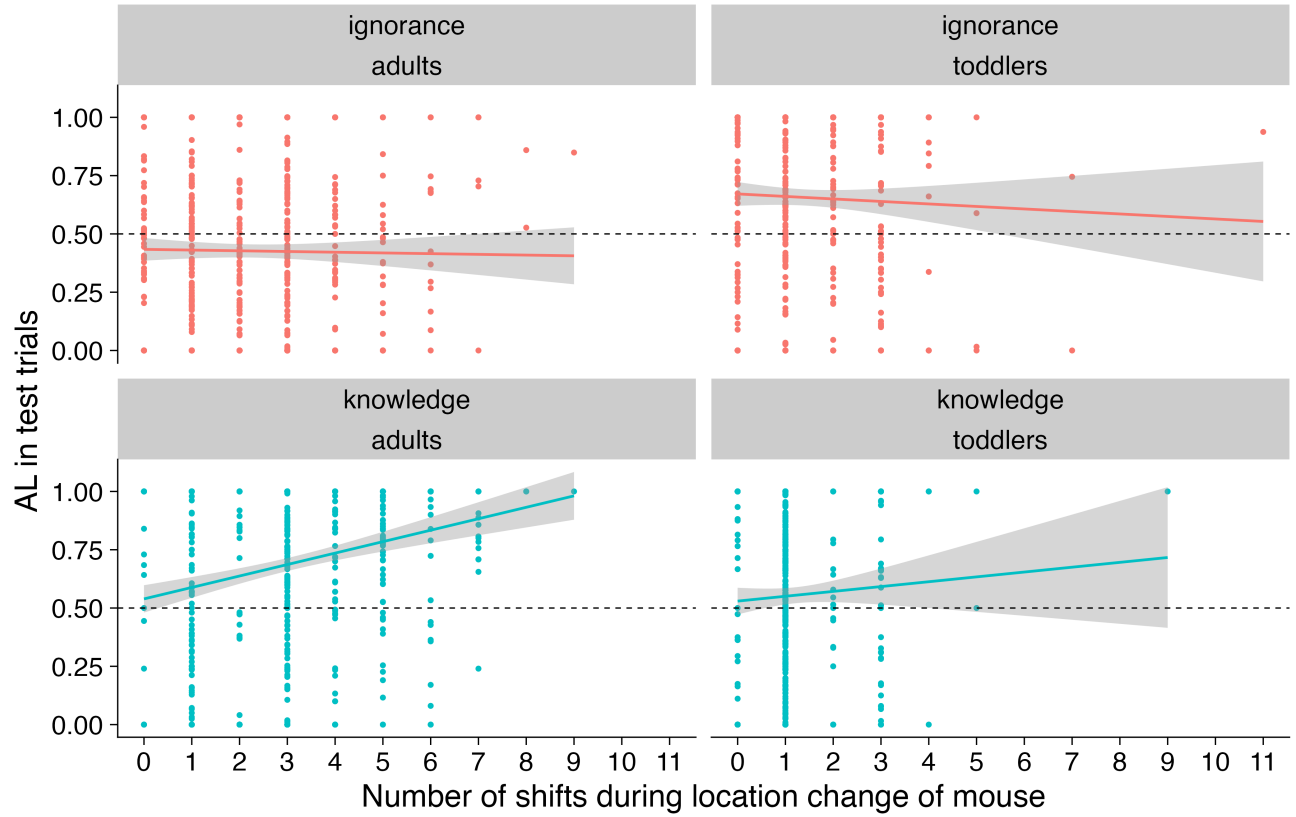


Figure 8. AL as a function of number of shifts between mouse and bear during location change of mouse in the test phase for toddlers and adults.

resulted in contributions from 47 labs, representing a total of 809 toddlers from xyz countries and 805 adults from xyz countries, of which 1224 were included in the final sample used for analysis (see Table 1). We begin our discussion by summarizing the principal results of the study with respect to confirmatory analysis and then discuss limitations of the study as well as future directions.

Conclusion

References

- Apperly, I. A., & Butterfill, S. A. (2009). Do humans have two systems to track beliefs and belief-like states? *Psychological Review*, 116(4), 953.
- Baillargeon, R., Buttelmann, D., & Southgate, V. (2018). Invited commentary: Interpreting failed replications of early false-belief findings: Methodological and theoretical considerations. *Cognitive Development*, 46, 112–124.
- Baillargeon, R., Scott, R. M., & He, Z. (2010). False-belief understanding in infants. *Trends in Cognitive Sciences*, 14(3), 110–118.
- Barone, P., Corradi, G., & Gomila, A. (2019). Infants’ performance in spontaneous-response false belief tasks: A review and meta-analysis. *Infant Behavior and Development*, 57, 101350.
- Barr, D. J. (2013). Random effects structure for testing interactions in linear mixed-effects models. *Frontiers Media SA*.
- Bürkner, P.-C. (2017). Brms: An r package for bayesian multilevel models using stan. *Journal of Statistical Software*, 80, 1–28.
- Burnside, K., Ruel, A., Azar, N., & Poulin-Dubois, D. (2018). Implicit false belief across the lifespan: Non-replication of an anticipatory looking task. *Cognitive Development*, 46, 4–11.
- Buttelmann, D., Carpenter, M., & Tomasello, M. (2009). Eighteen-month-old infants show false belief understanding in an active helping paradigm. *Cognition*, 112(2), 337–342.
- Buttelmann, F., & Kovács, Á. M. (2019). 14-month-olds anticipate others’ actions based on their belief about an object’s identity. *Infancy*, 24(5), 738–751.
- Buttelmann, F., Suhrke, J., & Buttelmann, D. (2015). What you get is what you believe: Eighteen-month-olds demonstrate belief understanding in an unexpected-identity task. *Journal of Experimental Child Psychology*, 131, 94–103.
- Carruthers, P. (2013). Mindreading in infancy. *Mind & Language*, 28(2), 141–172.
- Clements, W. A., & Perner, J. (1994). Implicit understanding of belief. *Cognitive*

Development, 9(4), 377–395.

Csibra, G., & Gergely, G. (2007). “Obsessed with goals”: Functions and mechanisms of teleological interpretation of actions in humans. *Acta Psychologica*, 124(1), 60–78.

Dennett, D. C. (1989). *The intentional stance*. MIT press.

Dörrenberg, S., Rakoczy, H., & Liszkowski, U. (2018). How (not) to measure infant theory of mind: Testing the replicability and validity of four non-verbal measures. *Cognitive Development*, 46, 12–30.

Dörrenberg, S., Wenzel, L., Proft, M., Rakoczy, H., & Liszkowski, U. (2019). Reliability and generalizability of an acted-out false belief task in 3-year-olds. *Infant Behavior and Development*, 54, 13–21.

Elsner, B., & Adam, M. (2021). Infants’ goal prediction for simple action events: The role of experience and agency cues. *Topics in Cognitive Science*, 13(1), 45–62.

Fabricius, W. V., Boyer, T. W., Weimer, A. A., & Carroll, K. (2010). True or false: Do 5-year-olds understand belief? *Developmental Psychology*, 46(6), 1402.

Flavell, J. H. (1988). *The development of children’s knowledge about the mind: From cognitive connections to mental representations*.

Flavell, J. H., Everett, B. A., Croft, K., & Flavell, E. R. (1981). Young children’s knowledge about visual perception: Further evidence for the level 1–level 2 distinction. *Developmental Psychology*, 17(1), 99.

Frank, M. C., Bergelson, E., Bergmann, C., Cristia, A., Floccia, C., Gervain, J., ... Yurovsky, D. (2017). A collaborative approach to infant research: Promoting reproducibility, best practices, and theory-building. *Infancy*, 22(4), 421–435.
<https://doi.org/10.1111/inf.12182>

Frank, M. C., Vul, E., & Saxe, R. (2012). Measuring the development of social attention using free-viewing. *Infancy*, 17(4), 355–375.

Frith, C. D., & Frith, U. (2006). The neural basis of mentalizing. *Neuron*, 50(4), 531–534.

Ganglmayer, K., Attig, M., Daum, M. M., & Paulus, M. (2019). Infants’ perception of

goal-directed actions: A multi-lab replication reveals that infants anticipate paths and not goals. *Infant Behavior and Development*, 57, 101340.

Gergely, G., & Csibra, G. (2003). Teleological reasoning in infancy: The naive theory of rational action. *Trends in Cognitive Sciences*, 7(7), 287–292.

Gergely, G., Nádasdy, Z., Csibra, G., & Biró, S. (1995). Taking the intentional stance at 12 months of age. *Cognition*, 56(2), 165–193.

Gliga, T., Jones, E. J., Bedford, R., Charman, T., & Johnson, M. H. (2014). From early markers to neuro-developmental mechanisms of autism. *Developmental Review*, 34(3), 189–207.

Gronau, Q. F., Sarafoglou, A., Matzke, D., Ly, A., Boehm, U., Marsman, M., ...

Steingroever, H. (2017). A tutorial on bridge sampling. *Journal of Mathematical Psychology*, 81, 80–97.

Grosse Wiesmann, C., Friederici, A. D., Singer, T., & Steinbeis, N. (2017). Implicit and explicit false belief development in preschool children. *Developmental Science*, 20(5), e12445.

Hare, B., Call, J., & Tomasello, M. (2001). Do chimpanzees know what conspecifics know? *Animal Behaviour*, 61(1), 139–151.

Hayashi, T., Akikawa, R., Kawasaki, K., Egawa, J., Minamimoto, T., Kobayashi, K., et al.others. (2020). Macaques exhibit implicit gaze bias anticipating others' false-belief-driven actions via medial prefrontal cortex. *Cell Reports*, 30(13), 4433–4444.

Heyes, C. (2014). Submentalizing: I am not really reading your mind. *Perspectives on Psychological Science*, 9(2), 131–143.

Hogrefe, G.-J., Wimmer, H., & Perner, J. (1986). Ignorance versus false belief: A developmental lag in attribution of epistemic states. *Child Development*, 567–582.

Horschler, D. J., MacLean, E. L., & Santos, L. R. (2020). Do non-human primates really represent others' beliefs? *Trends in Cognitive Sciences*, 24(8), 594–605.

Kaminski, J., Call, J., & Tomasello, M. (2008). Chimpanzees know what others know, but

- not what they believe. *Cognition*, 109(2), 224–234.
- Kampis, D., Buttelmann, F., & Kovács, Á. M. (2020). *Developing a theory of mind: Are infants sensitive to how other people represent the world?*
- Kampis, D., Karman, P., Csibra, G., Southgate, V., & Hernik, M. (2021). A two-lab direct replication attempt of southgate, senju and csibra (2007). *Royal Society Open Science*, 8(8), 210190.
- Kano, F., Krupenye, C., Hirata, S., Tomonaga, M., & Call, J. (2019). Great apes use self-experience to anticipate an agent’s action in a false-belief test. *Proceedings of the National Academy of Sciences*, 116(42), 20904–20909.
- Karg, K., Schmelz, M., Call, J., & Tomasello, M. (2015). The goggles experiment: Can chimpanzees use self-experience to infer what a competitor can see? *Animal Behaviour*, 105, 211–221.
- Király, I., Oláh, K., Csibra, G., & Kovács, Á. M. (2018). Retrospective attribution of false beliefs in 3-year-old children. *Proceedings of the National Academy of Sciences*, 115(45), 11477–11482.
- Knudsen, B., & Liszkowski, U. (2012). 18-month-olds predict specific action mistakes through attribution of false belief, not ignorance, and intervene accordingly. *Infancy*, 17(6), 672–691.
- Kovács, Á. M. (2016). Belief files in theory of mind reasoning. *Review of Philosophy and Psychology*, 7, 509–527.
- Kovács, Á. M., Téglás, E., & Endress, A. D. (2010). The social sense: Susceptibility to others’ beliefs in human infants and adults. *Science*, 330(6012), 1830–1834.
- Krupenye, C., Kano, F., Hirata, S., Call, J., & Tomasello, M. (2016). Great apes anticipate that other individuals will act according to false beliefs. *Science*, 354(6308), 110–114.
- Kulke, L., Duhn, B. von, Schneider, D., & Rakoczy, H. (2018). Is implicit theory of mind a real and robust phenomenon? Results from a systematic replication study. *Psychological Science*, 29(6), 888–900.

- 1132 Kulke, L., & Hinrichs, M. A. B. (2021). Implicit theory of mind under realistic social
1133 circumstances measured with mobile eye-tracking. *Scientific Reports*, *11*(1), 1215.
- 1134 Kulke, L., Johannsen, J., & Rakoczy, H. (2019). Why can some implicit theory of mind
1135 tasks be replicated and others cannot? A test of mentalizing versus submentalizing
1136 accounts. *PloS One*, *14*(3), e0213772.
- 1137 Kulke, L., & Rakoczy, H. (2018). Implicit theory of mind—an overview of current
1138 replications and non-replications. *Data in Brief*, *16*, 101–104.
- 1139 Kulke, L., & Rakoczy, H. (2019). Testing the role of verbal narration in implicit theory of
1140 mind tasks. *Journal of Cognition and Development*, *20*(1), 1–14.
- 1141 Kulke, L., Reiß, M., Krist, H., & Rakoczy, H. (2018). How robust are anticipatory looking
1142 measures of theory of mind? Replication attempts across the life span. *Cognitive*
1143 *Development*, *46*, 97–111.
- 1144 Kulke, L., Wübker, M., & Rakoczy, H. (2019). Is implicit theory of mind real but hard to
1145 detect? Testing adults with different stimulus materials. *Royal Society Open Science*,
1146 *6*(7), 190068.
- 1147 Leslie, A. M. (2005). Developmental parallels in understanding minds and bodies. *Trends*
1148 *in Cognitive Sciences*, *9*(10), 459–462.
- 1149 Liszkowski, U., Carpenter, M., & Tomasello, M. (2007). Pointing out new news, old news,
1150 and absent referents at 12 months of age. *Developmental Science*, *10*(2), F1–F7.
- 1151 Low, J., & Watts, J. (2013). Attributing false beliefs about object identity reveals a
1152 signature blind spot in humans’ efficient mind-reading system. *Psychological Science*,
1153 *24*(3), 305–311.
- 1154 Luo, Y., & Baillargeon, R. (2007). Do 12.5-month-old infants consider what objects others
1155 can see when interpreting their actions? *Cognition*, *105*(3), 489–512.
- 1156 Luo, Y., & Baillargeon, R. (2010). Toward a mentalistic account of early psychological
1157 reasoning. *Current Directions in Psychological Science*, *19*(5), 301–307.
- 1158 Martin, A., & Santos, L. R. (2016). What cognitive representations support primate theory

of mind? *Trends in Cognitive Sciences*, 20(5), 375–382.

Mellers, B., Hertwig, R., & Kahneman, D. (2001). Do frequency representations eliminate conjunction effects? An exercise in adversarial collaboration. *Psychological Science*, 12(4), 269–275.

Meristo, M., Morgan, G., Geraci, A., Iozzi, L., Hjelmquist, E., Surian, L., & Siegal, M. (2012). Belief attribution in deaf and hearing infants. *Developmental Science*, 15(5), 633–640.

Moll, H., & Tomasello, M. (2006). Level 1 perspective-taking at 24 months of age. *British Journal of Developmental Psychology*, 24(3), 603–613.

O'Neill, D. K. (1996). Two-year-old children's sensitivity to a parent's knowledge state when making requests. *Child Development*, 67(2), 659–677.

Onishi, K. H., & Baillargeon, R. (2005). Do 15-month-old infants understand false beliefs? *Science*, 308(5719), 255–258.

Papoutsaki, A., Sangkloy, P., Laskey, J., Daskalova, N., Huang, J., & Hays, J. (2016). *Proceedings of the 25th international joint conference on artificial intelligence (IJCAI)*.

Perner, J. (1991). *Understanding the representational mind*. The MIT Press.

Perner, J., & Ruffman, T. (2005). Infants' insight into the mind: How deep? *Science*, 308(5719), 214–216.

Phillips, J., Buckwalter, W., Cushman, F., Friedman, O., Martin, A., Turri, J., ... Knobe, J. (2021). Knowledge before belief. *Behavioral and Brain Sciences*, 44, e140.

Poulin-Dubois, D., Rakoczy, H., Burnside, K., Crivello, C., Dörrenberg, S., Edwards, K., et al.others. (2018). Do infants understand false beliefs? We don't know yet—a commentary on baillargeon, buttelmann and southgate's commentary. *Cognitive Development*, 48, 302–315.

Powell, L. J., Hobbs, K., Bardis, A., Carey, S., & Saxe, R. (2018). Replications of implicit theory of mind tasks with varying representational demands. *Cognitive Development*, 46, 40–50.

- 1186 Premack, D., & Woodruff, G. (1978). Does the chimpanzee have a theory of mind?
1187 *Behavioral and Brain Sciences*, 1(4), 515–526.
- 1188 Priewasser, B., Fowles, F., Schweller, K., & Perner, J. (2020). Mistaken max befriends
1189 duplo girl: No difference between a standard and an acted-out false belief task. *Journal*
1190 *of Experimental Child Psychology*, 191, 104756.
- 1191 Priewasser, B., Rafetseder, E., Gargitter, C., & Perner, J. (2018). Helping as an early
1192 indicator of a theory of mind: Mentalism or teleology? *Cognitive Development*, 46,
1193 69–78.
- 1194 Ruffman, T. (1996). Do children understand the mind by means of simulation or a theory?
1195 Evidence from their understanding of inference. *Mind & Language*, 11(4), 388–414.
- 1196 Schneider, D., Bayliss, A. P., Becker, S. I., & Dux, P. E. (2012). Eye movements reveal
1197 sustained implicit processing of others' mental states. *Journal of Experimental*
1198 *Psychology: General*, 141(3), 433.
- 1199 Schneider, D., Slaughter, V. P., Bayliss, A. P., & Dux, P. E. (2013). A temporally
1200 sustained implicit theory of mind deficit in autism spectrum disorders. *Cognition*,
1201 129(2), 410–417.
- 1202 Schneider, D., Slaughter, V. P., & Dux, P. E. (2017). Current evidence for automatic
1203 theory of mind processing in adults. *Cognition*, 162, 27–31.
- 1204 Schuwerk, T., Priewasser, B., Sodian, B., & Perner, J. (2018). The robustness and
1205 generalizability of findings on spontaneous false belief sensitivity: A replication
1206 attempt. *Royal Society Open Science*, 5(5), 172273.
- 1207 Scott, R. M., & Baillargeon, R. (2009). Which penguin is this? Attributing false beliefs
1208 about object identity at 18 months. *Child Development*, 80(4), 1172–1196.
- 1209 Scott, R. M., & Baillargeon, R. (2017). Early false-belief understanding. *Trends in*
1210 *Cognitive Sciences*, 21(4), 237–249.
- 1211 Scott, R. M., Richman, J. C., & Baillargeon, R. (2015). Infants understand deceptive
1212 intentions to implant false beliefs about identity: New evidence for early mentalistic

- reasoning. *Cognitive Psychology*, 82, 32–56.
- Senju, A., Southgate, V., Miura, Y., Matsui, T., Hasegawa, T., Tojo, Y., ... Csibra, G. (2010). Absence of spontaneous action anticipation by false belief attribution in children with autism spectrum disorder. *Development and Psychopathology*, 22(2), 353–360.
- Senju, A., Southgate, V., Snape, C., Leonard, M., & Csibra, G. (2011). Do 18-month-olds really attribute mental states to others? A critical test. *Psychological Science*, 22(7), 878–880.
- Senju, A., Southgate, V., White, S., & Frith, U. (2009). Mindblind eyes: An absence of spontaneous theory of mind in asperger syndrome. *Science*, 325(5942), 883–885.
- Sheskin, M., Scott, K., Mills, C. M., Bergelson, E., Bonawitz, E., Spelke, E. S., et al.others. (2020). Online developmental science to foster innovation, access, and impact. *Trends in Cognitive Sciences*, 24(9), 675–678.
- Southgate, V., Johnson, M. H., Karoui, I. E., & Csibra, G. (2010). Motor system activation reveals infants’ on-line prediction of others’ goals. *Psychological Science*, 21(3), 355–359.
- Southgate, V., Senju, A., & Csibra, G. (2007). Action anticipation through attribution of false belief by 2-year-olds. *Psychological Science*, 18(7), 587–592.
- Southgate, V., & Verneti, A. (2014). Belief-based action prediction in preverbal infants. *Cognition*, 130(1), 1–10.
- Surian, L., Caldi, S., & Sperber, D. (2007). Attribution of beliefs by 13-month-old infants. *Psychological Science*, 18(7), 580–586.
- Surian, L., & Franchin, L. (2020). On the domain specificity of the mechanisms underpinning spontaneous anticipatory looks in false-belief tasks. *Developmental Science*, 23(6), e12955.
- Surian, L., & Geraci, A. (2012). Where will the triangle look for it? Attributing false beliefs to a geometric shape at 17 months. *British Journal of Developmental Psychology*, 30(1), 30–44.
- Thoermer, C., Sodian, B., Vuori, M., Perst, H., & Kristen, S. (2012). Continuity from an

1240 implicit to an explicit understanding of false belief from infancy to preschool age.

1241 *British Journal of Developmental Psychology*, 30(1), 172–187.

1242 Träuble, B., Marinović, V., & Pauen, S. (2010). Early theory of mind competencies: Do
1243 infants understand others' beliefs? *Infancy*, 15(4), 434–444.

1244 Wellman, H. M., & Cross, D. (2001). Theory of mind and conceptual change. *Child*
1245 *Development*, 72(3), 702–707.

1246 Wiesmann, C. G., Friederici, A. D., Disla, D., Steinbeis, N., & Singer, T. (2018).

1247 Longitudinal evidence for 4-year-olds' but not 2-and 3-year-olds' false belief-related
1248 action anticipation. *Cognitive Development*, 46, 58–68.

1249 Woodward, A. L., & Sommerville, J. A. (2000). Twelve-month-old infants interpret action
1250 in context. *Psychological Science*, 11(1), 73–77.

1251 Yang, Q., Bucci, M. P., & Kapoula, Z. (2002). The latency of saccades, vergence, and
1252 combined eye movements in children and in adults. *Investigative Ophthalmology &*
1253 *Visual Science*, 43(9), 2939–2949.

Table 1

Lab and Participant information.

Lab	N collected	N included	Sex (N Female)	Mean Age (years)	Method
CogConcordia	21	16	11	22.12	In-lab
CorbitLab	16	15	14	19.87	In-lab
DevlabAU	20	20	15	25.15	In-lab
MEyeLab	53	53	39	24.47	In-lab
MiniDundee	15	13	10	30.23	In-lab
PKUSu	39	32	19	22.66	In-lab
SkidLSDLab	11	8	3	21.62	In-lab
ToMcdlSalzburg	33	31	22	27.23	In-lab
UIUCinfantlab	36	32	25	19.06	In-lab
WSUMARCS	18	13	8	29.85	In-lab
affcogUTSC	23	8	5	20.88	web-based
babyLeidenEdu	20	16	12	23.31	In-lab
babylabAmsterdam	17	16	13	24.00	In-lab
babylabBrookes	67	65	49	21.78	In-lab
babylabINCC	18	18	12	31.00	In-lab
babylabMPIB	16	16	11	27.44	In-lab
babylabNijmegen	19	15	13	22.13	In-lab
babylabTrento	16	16	9	21.69	In-lab
babylabUmassb	33	11	10	19.00	In-lab
babyuniHeidelberg	16	16	14	22.06	In-lab
beinghumanWroclaw	19	16	9	32.75	web-based
careylabHarvard	18	15	12	19.80	In-lab
cclUNIRI	32	32	17	30.53	In-lab
childdevlabAshoka	16	16	8	30.88	In-lab
collabUIOWA	16	16	10	19.19	In-lab
gaugGöttingen	30	28	18	31.71	In-lab
jmuCDL	32	32	22	18.81	In-lab
kidsdevUniofNewcastle	15	14	7	33.57	In-lab
labUNAM	20	11	8	22.45	In-lab

Table 2 continued

Lab	N collected	N included	Sex (N Female)	Mean Age (years)	Method
lmuMunich	31	30	23	22.53	In-lab
mecdmpihcbs	19	19	10	27.79	In-lab
socialcogUmiami	16	15	9	19.27	In-lab
sociocognitivelab	17	17	11	32.12	In-lab
tauccd	15	12	6	24.50	In-lab
Total	803	703	484	24.75	

Table 2

Lab and Participant information.

Lab	N collected	N included	Sex (N Female)	Mean Age (months)	Method
CogConcordia	21	8	4	22.92	web-based
CorbitLab	11	10	5	22.77	In-lab
DevlabAU	18	17	8	19.00	In-lab
PKUSu	50	32	13	20.84	In-lab
SkidLSDLab	8	2	0	20.11	In-lab
ToMedlSalzburg	17	12	6	22.20	In-lab
UIUCinfantlab	18	15	9	21.96	In-lab
babyLeidenEdu	18	12	8	22.59	In-lab
babylabAmsterdam	28	12	6	23.19	In-lab
babylabBrookes	17	12	7	22.15	In-lab
babylabChicago	17	13	4	20.10	In-lab
babylabINCC	16	9	6	23.40	In-lab
babylabNijmegen	19	10	3	23.52	In-lab
babylabOxford	25	19	8	23.42	In-lab
babylabPrinceton	17	11	7	22.15	In-lab
babylabTrento	18	17	10	22.72	In-lab
babylabUmassb	7	6	2	20.35	In-lab
babylingOslo	17	14	7	21.99	In-lab
babyuniHeidelberg	16	12	4	22.69	In-lab
beinghumanWroclaw	24	14	7	23.77	web-based
careylabHarvard	17	12	5	21.99	In-lab
cecBYU	16	14	4	22.39	In-lab
childdevlabAshoka	16	10	6	22.44	In-lab
gaugGöttingen	28	15	9	23.06	In-lab
gertlabLancaster	21	17	8	23.03	In-lab
infantcogUBC	26	19	8	24.39	In-lab
irlConcordia	19	12	5	22.47	In-lab
kidsdevUniofNewcastle	16	14	9	22.36	In-lab
kokuHamburg	19	14	7	25.99	In-lab

Table 2 continued

Lab	N collected	N included	Sex (N Female)	Mean Age (months)	Method
labUNAM	18	12	7	22.68	In-lab
lmuMunich	48	24	16	22.68	In-lab
mecdmpihcbs	25	12	8	23.58	In-lab
mpievaCCP	22	18	10	23.33	In-lab
saxelab	31	15	2	23.13	web-based
socallabUCSD	47	15	4	22.09	web-based
tauccd	15	12	8	22.99	In-lab
unicph	43	29	16	21.50	In-lab
Total	809	521	256	22.48	

Table 3

Results of the Bayesian mixed effects models for the familiarization trials.

model	term	estimate	conf.low	conf.high
PTL toddlers	Intercept Effect	0.12	0.09	0.15
PTL adults	Intercept Effect	0.26	0.23	0.29
First Look toddlers	Intercept Effect	0.44	0.27	0.61
First Look adults	Intercept Effect	1.03	0.86	1.20

Table 4

Results of the Bayesian mixed effects models for the test trials.

model	term	estimate	conf.low	conf.high
PTL toddlers	Condition Effect	0.10	0.03	0.17
PTL adults	Condition Effect	-0.20	-0.26	-0.15
First Look toddlers	Condition Effect	0.53	0.13	0.93
First Look adults	Condition Effect	-0.89	-1.21	-0.56