Contents lists available at ScienceDirect

# Cognition

# Growing out of your own mind: Reexamining the development of the self-other difference in the unexpected contents task[☆]

David M. Sobel [*]

*Brown University, United States of America*

## ARTICLE INFO

## ABSTRACT

The unexpected contents task is a well-established measure for studying young children's developing theory of mind. The task measures whether children understand that others have a false belief about a deceptive container and whether children can track the representational change in their own beliefs about the container's contents. Performance on both questions improves between the ages of 3 and 4. A previous meta-analysis (Wellman, Cross, & Watson, 2001) found little evidence for a difference in children's responses on these questions, but did not investigate the weak effect size that was reported for the interaction between age and question type. The two meta-analyses reported here update the literature review, and find a more robust interaction between question type and age. Three-year-olds showed better performance on questions about their own representational change than others' false belief, while older children showed the reverse pattern. A mega-analysis of a sample of over 1200 children between the ages of 36–60 months then showed the same result. This response pattern requires novel theoretical interpretations, which include reframing the development of children's understanding of false belief.

For over 40 years, researchers have been interested in how children develop a *theory of mind* – an understanding of their own and others' mental states (e.g., Flavell, 1999; Perner, 1991; Premack & Woodruff, 1978; Wellman, 1990). Research in theory of mind has emphasized the different ways children develop knowledge of individual mental states, the interconnections among mental states, and how mental states relate to action, both for understanding themselves and others. Perhaps the most common way theory of mind has been studied is through the understanding of false belief. Understanding false belief allows children to appreciate why individuals act or think in ways that are inconsistent with the actual state of the world.

To illustrate how false belief has been studied, consider the *unexpected contents* task (Gopnik & Astington, 1988; Perner, Leekam, & Wimmer, 1987). In this task, children (usually preschoolers) are introduced to a familiar container, which is revealed to be deceptive, such as a crayon box that contains birthday candles. This measure can be used to consider what children will understand about another's belief state, given that it can differ from the actual state of the world. In such questions (referred to here as *other* questions), children are asked what a person who is not present, and thus lacks the knowledge that the container is deceptive, will think is in the container. But children can also be asked about their understanding of their own mental states, and that revealing the unexpected contents licenses representational change. These questions (referred to here as *self* questions), ask children about changes to their own mental states, for instance, "But what did you <the child> think was in here <the deceptive container>?" (Wording taken from Perner et al., 1987).

Perhaps the most dominant view of theory of mind development, which has been called *theory theory*, involves describing the development of a representation theory of mind. This view suggests that children have intuitive theories of mental states that apply broadly – both to themselves and to others (see e.g., Gopnik, 1993; Gopnik & Wellman, 1994; Gopnik & Meltzoff, 1997; Perner, 1988, 1991, 1996; Wellman, 1990, 2014, among many others). Successful performance on the false belief task indicates the development of the domain-specific conceptual knowledge relevant to understanding the representational nature of mental states. Although the ability to demonstrate that representational understanding is affected by the demand characteristics of the

---

procedure, evidenced for example, by individual differences with children's developing inhibitory control (e.g., Carlson & Moses, 2001; see also Benson, Sabbagh, Carlson, & Zeflazo, 2013; Wellman, 2014), development of the domain-specific knowledge required to succeed on a false belief task occurs from observation and interaction with the world (Wellman & Liu, 2004). Crucial to the discussion here, that domain-specific knowledge allows children to reason successfully on both the self and other questions equivalently.

Although the initial studies that asked children both the self and other questions found conflicting results,[1] strong evidence in favor of this view came from a meta-analysis of over 50 studies conducted by Wellman et al. (2001). This analysis showed that performance on the self and other questions did not significantly differ. Wellman et al. (2001) highlighted this finding as particularly contrastive to another theoretical account of theory of mind development – *simulation theory* – following research in philosophy of mind (e.g., Goldman, 2006; Gordon, 1992; see also Harris, 1992). Simulation theory suggests that inferences about others' mental states are made by accessing and reasoning about one's own mental states.

There are numerous philosophical descriptions of this approach to understanding theory of mind, including the possibility that the two theories reduce to one another (see e.g., Perner & Brandl, 2009), but the one specifically taken up by Wellman et al. regarded the primacy of one's own mental access (Harris, 1992). In particular, children's understanding of their own mental states is salient and necessary to understand similar mental states in others. Wellman et al. (2001), however, argue that the null result between performance on the self and other questions was meaningful for distinguishing between these two accounts: "… simulation accounts emphasizing the primacy of self-experience suggest that self understanding should develop first." And, later in on that page, they write, "on the surface, the fact that children ever systematically err in reporting their own false beliefs seems problematic for simulation accounts." (both p. 678).

Recent evidence, however, suggests that young children do have some first-person access to their mental states, particularly their knowledge states. Harris, Yang, and Cui (2017) demonstrated that 2- and 3-year-olds (children who typically do not succeed on the unexpected contents task) were more likely to deny possessing knowledge (i.e., saying "I don't know") than point out that another person does not know something (i.e., saying "You don't know"). These same children were also more likely to ask questions about whether others know something as opposed to whether they do. The former finding suggests that children might be aware of their own ignorance before that of others, while the latter suggests that children recognize their own lack of knowledge, and that others might be sources of knowledge that they themselves do not possess. These data suggest that although young children might not behave differently on self and other questions on a false belief task, they show differences in their understanding of their own and others' knowledge.

These data suggest the importance of reconsidering developmental findings on the difference between the self and other questions, and in particular the findings from the Wellman et al. (2001) analysis. In particular, while the difference between the self and the other questions was not significant, the interaction between the type of question asked and children's age was the level of a statistical trend ($p > .07$) with a small effect size ($f = 0.12$).[2] Typically, small effects are not investigated in developmental science because of limits on data collection or

argumentation dismissing them as unimportant. However, small effects may be meaningful if they have important theoretical consequences (Ellis, 2010). Here, the theoretical implications of recent research suggesting children have first-person access to their mental states, warrants a reinvestigation of this finding. Doing so also allows us to address a methodological limitation of method used by Wellman et al. (2001), which is that their findings did not directly consider the sample size of each condition they analyzed.

The rest of this paper presents three analyses that reconsider and update the relation between performance on self and other questions about false belief through meta-analysis and a quasi-mega-analysis[3] of performance on the unexpected contents task. Analysis 1 presents a meta-analysis that follows a similar strategy used by Wellman et al. (2001). The goal of this analysis is not to replicate all aspects of Wellman et al. (2001), but rather focus specifically on the self-other difference. To that end, the analysis relies on papers that presented data from the unexpected contents task, which specifically asked either the self or other question. The papers collected for this analysis report how children performed on the self or other questions individually or how children performed on both questions (not combined into a single datum). Analysis 1 uses the same analysis strategy used by Wellman et al. (2001), but also uses a second measure of effect size, which takes the size of the samples into account to ensure the reproducibility of results.

A limitation of this first analysis is that some researchers who administer an unexpected contents task only ask the self question, while other researchers only ask the other question, while still other researchers ask both questions. Combining these between-subject data might not be as robust an analysis as one that contrasts results using a within-subject design. Analysis 2 replicates the metanalytic strategy on the subset of studies from Analysis 1 that present data from the self and other questions from the same participant.

Analysis 3 uses a similar within-subjects approach by considering a dataset of over 1200 preschoolers administered the unexpected contents task. These data were collected in a single lab with similar materials and the same procedure. This mega-analytic strategy not only generalizes and extends the results of the two meta-analyses, but also considers why small effects might be robust and meaningful, even if they are difficult to describe in individual empirical papers.

## 1. Analyses 1 and 2

The goal of these analyses is to examine whether there are differences in how children respond to questions about changes in their previously erroneous belief states (self questions) and measures of others' possessing a false belief (other questions). We considered papers that reported the results of studies that use the unexpected contents task on children between the ages of 36 and 60 months, based on the assumption that the majority of children will eventually succeed on both questions at later ages. Analysis 1 uses a meta-analytic technique similar to that used by Wellman et al. (2001), updating the literature considered as well as extending the analysis to include a more robust measure of effect size. Analysis 2 then considers a subset of the conditions/experiments used in Analysis 1, which specifically contrast performance by the same child on both measures.

## 2. Literature review

To determine papers to analyze, the papers used in the Wellman et al. (2001), as well as a more recent analysis that exclusively considered studies reporting the unexpected contents task (Sobel & Austerweil,

---

[1] Perner et al. (1987) found that children performed better on the self than other question, whereas Gopnik and Gopnik and Astington (1988) found the opposite, that children performed better on the other question than the self question.

[2] The implicature of reporting $p > .07$ is that the significance level of this analysis is between $p = .07$ and $p = .08$. The effect size of this analysis was not reported, but was calculated based on the reported *F*-statistic.

---

[3] A mega-analysis is a reanalysis of raw data from multiple sources (Eisenhauer, 2021). Here, the raw data are all taken from a single lab, hence the "quasi."

2016) were examined. All the papers from those analyses that reported data from the unexpected contents task, which met the inclusion criteria (described below) were included. Google Scholar searches were also performed for the terms "unexpected content task," "representational change task" and "theory of mind scales". The first two are commonly used names for the measure. The latter was included because the scales (Wellman & Liu, 2004) present children with the unexpected contents task. The author also put out a call on the COGDEVSOC listserv (March 2021) for unpublished datasets and recently published data. Fig. 1 shows a flowchart of the literature review process. Because of the intention to include data from the author's published work separately in the mega-analysis presented in Study 3, none of the author's own papers were included in these studies (to reduce the possibility of bias in reproduction).

To summarize this process, after duplicate entries were removed from the initial search, there were 917 entries under consideration. One hundred fifty-two of those entries could be eliminated without review because (a) they were written by the author (8), (b) Google Scholar indicated they were books, book chapters, or published in journals that exclusively published review, and not empirical articles (88), and (c) Google Scholar indicated they were unavailable in English (56). Nine additional entries were removed because they were inaccessible.

The remaining 756 full-text articles and datasets were read by the author to determine whether the paper reported the age of an individual sample and performance on the self and other questions of the unexpected contents task individually (or if only one question was asked, performance on that question). Papers that only reported a composite score of performance on the self and other questions, or "false belief" scores that reported the results of an unexpected content task combined with other measures of false belief (such as the unexpected transfer task or the unexpected identity task) were not included unless the separate scores on the self and other questions of the unexpected contents task could be determined. Moreover, only samples of neurotypical children with a mean age between 36 and 60 months were considered. If the paper reported a longitudinal or training study, only pretest data was used; these kinds of studies were only included if they reported all pretest data as opposed to only including data from children who failed the measure. If the paper contrasted conditions in which deception was and was not used, only cases in which the motive of the study was not to deceive another person were considered (because there are well-known main effects of deceptive motive, see Wellman et al., 2001, and our supplemental materials). Fig. 1 shows the distribution of exclusions.

Following Wellman et al. (2001) and Glass, McGaw and Smith (1981), condition (or age group within a condition where possible) was used as the unit of analysis, instead of individual participant or a full study or paper. This resulted in including 91 papers/datasets used in Analysis 1, which represented 244 age groups/conditions within experiments (9201 data points in total). Analysis 2 focused on the subset of these 91 papers/datasets that asked both the self and other question of the same child participant to consider a within-subject analysis. This resulted in an analysis of the results of 37 papers/datasets (70 unique age groups/conditions, 2466 data points). The papers/datasets used in both analyses are reported in Table 1, as we as the number of conditions each entry contributed to the overall analysis.

## 3. Coding

Analyses 1 and 2 focused on extracting performance on the self and other question of the unexpected contents task, as well as the mean age and sample size of the reported sample. Each condition within each experiment was coded for performance on the test question as well as whether the measure was a self or other question and whether both questions were asked of the same child in the same procedure. If both questions were asked, we reported whether the order of these questions was counterbalanced. Finally, we only included whether the paper reported the mean age of the condition under consideration and the

number of participants in that condition.

The goal of this analysis was not to reproduce all of the analyses performed by Wellman et al. (2001). Many of the analyses that were significant there, which might interact or affect the main result presented here, however, are described in the Supplemental Materials. This includes country of origin, whether a deceptive condition was included in the study, whether the deceptive object was really present in the container when the test question was asked, and whether temporal markers were used when asking the self question. One additional factor – the way in which researchers scored the task – was also considered in the Supplemental Materials, as previous work has shown that the developmental trajectory of children's performance is affected by how the measure is scored (Sobel & Austerweil, 2016).

To investigate differences between the self and other question on the unexpected contents measure, we conducted the same logit analysis as presented in Wellman et al. (2001). The proportion correct on each question ($p$) was transformed via a logit transformation in which the logit score $= \ln(p/1-p)$,[4] which became the dependent variable in the analyses. This analysis, however, does not take the sample size of the condition/age group into account. To account for this, we supplemented this analysis with a second transformation on the data, following the meta-analytic strategy used by Tong, Wang, and Danovitch (2020). Where possible (i.e., if the standard deviation of performance was reported or could be calculated), we transformed performance in each condition/age group to a Hedges' $g$ value, based on the proportion correct and standard deviation of that condition/age group. This was done by assuming that the comparison value of chance had a Mean and Standard Deviation of 0.5 (an idealized representation of chance responding), via the formula:

$$ g = \frac{(Prop_{}\_Correct \quad 0.5)}{\frac{StdDev^2 + .25}{2}} $$

Two readers, blind to the hypotheses of the study coded a subset of 10 of the papers in the same manner. Agreement between these two readers was 100%. Agreement between these readers and the author was 90%. Disagreement was resolved through discussion.

## 4. Results

All data sets described in this manuscript (as well as a list of all articles that were considered) are available at https://osf.io/624av. We first consider the overall dataset (Analysis 1) and then a subset of the papers we considered, which presented both the self and other questions to the same child, using within-subject analyses (Analysis 2).

### 4.1. Analysis 1: Papers reporting findings from either self or other questions on the unexpected contents task

Fig. 2 shows the logit values for the self and other questions across age. A set of hierarchical linear regressions were conducted, with logit scores and $g$-values as the dependent variable. Each first considered age (defined by the mean age of the condition as reported in each paper) and question type (self vs. other). Then the interaction value between age and question type was added in a second model.

For the logit scores, the first model explained a significant amount of

---

[4] One condition presented a proportion of 0, while another presented a proportion of 1, both with standard deviations of 0. These resulted in an undefined logit scores and values of $g$. We entered proportions of 0.01 and 0.99 and standard deviations of 0.0001 for each, which set minimal and maximal values for the dependent variables without skewing the data too greatly. The results presented in the main text are similar if these two entries are not included in the analyses. Please note that the formula used is $\ln(p/1-p)$ and that the value is not raised to the fourth power.
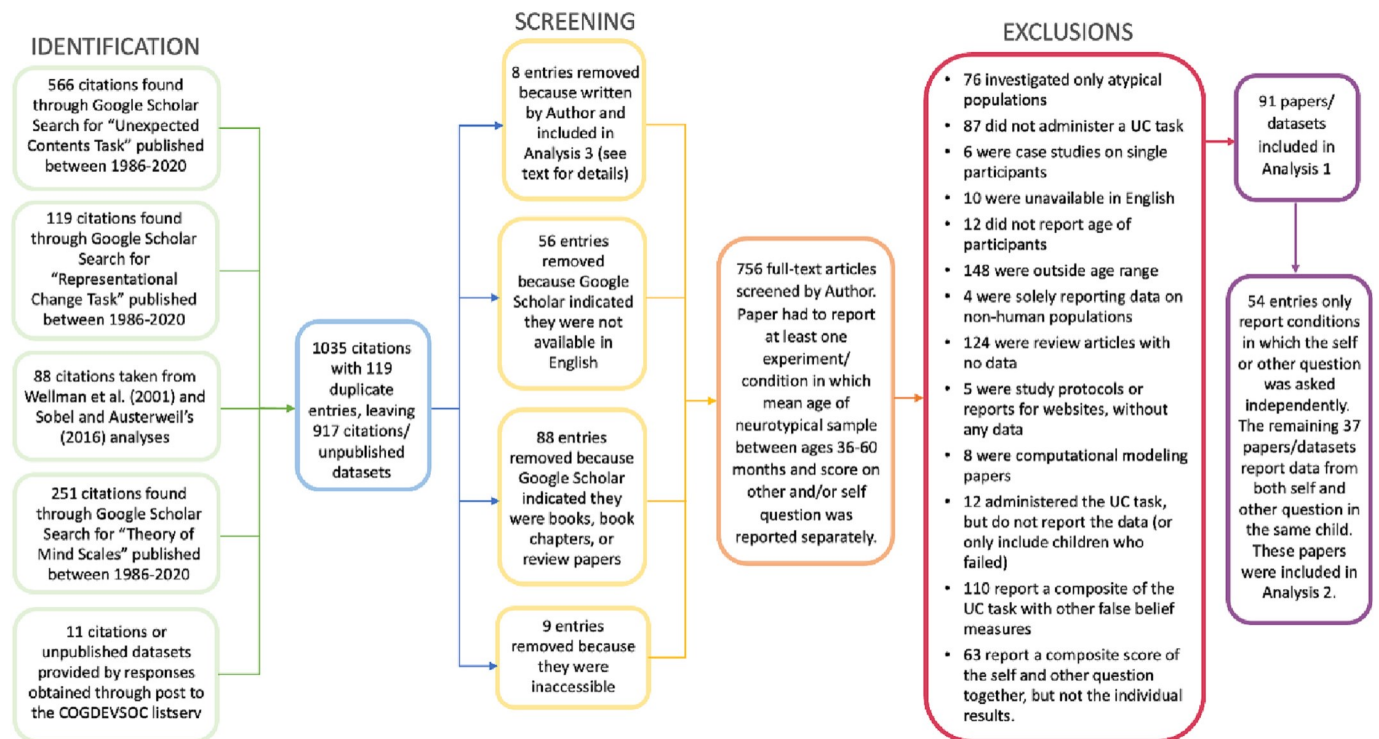
**Fig. 1.** Flowchart representing procedure used for literature identification, screening, and exclusions for Analyses 1–2.

the variance in the logit values over an intercept-only model, $R^2 = 0.28$, $F(2, 241) = 46.21$, $p < .001$. The second model, which added the interaction between age and question type, predicted a significant increase in the amount of variance in the logit values, $\Delta R^2 = 0.01$, $F(1, 240) = 4.07$, $p = .04$. This final model predicted a significant amount of variance over an intercept-only model, $R^2 = 0.29$, $F(3, 240) = 32.56$, $p < .001$. In this final model, age positively predicted performance, $B = 0.07$, $SE = 0.02$, 95% CI [0.04, 0.10], $t = 4.72$, $p < .001$. There was also a significant effect of question type, with performance on the self question higher overall than performance on the other question, $B = -2.42$, $SE = 1.00$, 95% CI [-4.39, -0.45], $t = -2.42$, $p = .02$. There was also a significant interaction between question type and age, $B = 0.04$, $SE = 0.02$, 95% CI [0.00, 0.08], $t = 2.02$, $p = .04$.

A similar analysis for the $g$-values generated by each condition was conducted. The first model (age and question type), explained a significant amount of the variance over an intercept-only model, $R^2 = 0.30$, $F(2, 227) = 49.32$, $p < .001$. The second model, which added the interaction between age and question type, predicted a significant increase in the amount of variance in the logit values, $\Delta R^2 = 0.01$, $F(1, 226) = 4.23$, $p = .04$. This final model predicted a significant amount of variance over an intercept-only model, $R^2 = 0.32$, $F(3, 226) = 34.76$, $p < .001$. In this final model, age positively predicted performance, $B = 0.03$, $SE = 0.01$, 95% CI [0.02, 0.05], $t = 4.84$, $p < .001$, $OR = 1.03$. There was also a significant effect of question type, with performance on the other question lower overall than performance on the self question, $B = -1.12$, $SE = 0.45$, 95% CI [-0.2.00, -0.23], $t = -2.50$, $p = .01$, $OR = 0.33$, and a significant interaction between question type and age, $B = 0.02$, $SE = 0.01$, 95% CI [0.00, 0.04], $t = 2.06$, $p = .04$. $OR = 1.04$.

### 4.2. Analysis 2: Papers reporting findings from both self and other questions on the unexpected contents measure

The second analysis looked at a subset of the papers from Analysis 1 that reported results of the self and other questions for the same children. Table 1 indicates the papers included in this analysis. As in the previous analysis, the logit value and the $g$ value for performance on the

self and other questions in each condition/age groups were considered. These data are shown in Fig. 3. A Generalized Estimating Equation was constructed to control for within-subject effects, with a robust correlation matrix, assuming a linear distribution on the dependent variables with age and question type as independent variables. For both analyses, a factorial model was considered.

Looking at the logit scores, the model revealed a significant effect of age, $B = 0.12$, $SE = 0.02$, 95% CI [0.08, 0.16], Wald $\chi^2(1) = 31.26$, $p < .001$, $OR = 1.13$, and question type, $B = 2.91$, $SE = 0.87$, 95% CI [1.21, 4.61], Wald $\chi^2(1) = 11.28$, $p = .001$, $OR = 18.36$, with children performing better overall on the self questions than the other questions. There was also a significant interaction between age and question type, $B = -0.06$, $SE = 0.02$, 95% CI [-0.09, -0.02], Wald $\chi^2(1) = 11.25$, $p = .001$, $OR = 0.94$. As with the previous analysis, the younger children in this sample initially performed better on the self question than the other question, but as the age of children in the sample increased, performance on the other question started to exceed performance on the self question.

The same pattern of results was obtained when the Hedges' $g$ scores were analyzed. There was a significant effect of age, $B = 0.05$, $SE = 0.01$, 95% CI [0.03, 0.07], Wald $\chi^2(1) = 34.70$, $p < .001$, $OR = 1.05$, and question type, $B = 1.15$, $SE = 0.32$, 95% CI [0.52, 1.78], Wald $\chi^2(1) = 12.91$, $p < .001$, $OR = 3.16$ as well as interaction between age and question type, $B = -0.02$, $SE = 0.01$, 95% CI [-0.03, -0.01], Wald $\chi^2(1) = 10.82$, $p = .001$, $OR = 0.98$.

### 5. Discussion

Both the more inclusive analysis (Analysis 1) and the within-subject analysis (Analysis 2) of the self and other questions showed two findings of note. First, in addition to the expected main effects of age, there was a significant difference between performance on the self and other question, with children performing better on the former over the latter overall. The second was a significant interaction between children's age and question type, which potentially explains the first finding. Younger children performed better on the self question than the other question. However, as children increased in age from 36 to 60 months,

**Table 1**

Lfist of papers/Datasets used fin Anaflysfis 1 (and Anaflysfis 2, shown wfith an *).

| Authors | Date of publflicatfion | Number of age groups/ Condfitfions used fin Anaflysfis 1 |
|---|---|---|
| Afltvater-Mackensen* | Unpubflfished | 4 |
| Amado, Serrat & Sfidera* | 2014 | 2 |
| Appfleton & Reddy* | 1996 | 2 |
| Arreckx | 2007 | 1 |
| Atance | 2001 | 1 |
| Atance, Metcaflf & Thfiessen | 2017 | 2 |
| Atance & ONefiflfl | 2004 | 4 |
| Bafird & Moses | 2001 | 2 |
| Barreto, Osorfio, Baptfista, Fearon, & Martfins* | 2018 | 2 |
| Befißert, Muflvey & Kfiflflen | Unpubflfished | 2 |
| Beflflagamba, Addessfi, Focaroflfi et afl.* | 2015 | 4 |
| Benson, Sabbagh, Carflson, & Zeflazo | 2017 | 1 |
| Bernstefin | 2009 | 1 |
| Bernstefin* | Unpubflfished | 2 |
| Bozbfiyfik | 2016 | 1 |
| CarflsonMoses* | 2001 | 4 |
| Cassfidy | 1998 | 2 |
| Causey & Bjorkflund* | 2014 | 2 |
| Conry-Murray | 2013 | 1 |
| Daflke | 1995 | 2 |
| Davfis* | 2001 | 2 |
| Fabrficfius & Khaflfifl | 2003 | 1 |
| Fflynn* | 2006 | 4 |
| Fflynn, O'Maflfley and Wood | 2004 | 1 |
| Freeman & Lacohee | 1995 | 6 |
| Frye, Zeflazo, and Paflfafi* | 1995 | 8 |
| Gopnfik & Astfington* | 1988 | 8 |
| Gopnfik & Rosatfi* | 2001 | 2 |
| Gopnfik & Sflaughter | 1991 | 2 |
| Guajardo, Parker & Turfley-Ames | 2009 | 1 |
| Guajardo & Turfley-Ames | 2004 | 2 |
| Gut, Haman, Gorbanfiuk, & Chyflfinskfia* | 2020 | 4 |
| Hafla, Hug & Henderson | 2003 | 2 |
| Hansen | 2010 | 2 |
| Hanson & Atance | 2014 | 2 |
| Hasnfi | 2015 | 1 |
| Hfiflfler, Weber & Young | 2014 | 1 |
| Hogrefe, Wfimmer, & Perner | 1986 | 3 |
| Hoflmes, Bflack & Mfiflfler* | 1996 | 2 |
| Hong* | 2016 | 4 |
| Hughes | 1998 | 1 |
| Jackson | 2001 | 1 |
| Jufllien | 2018 | 1 |
| Kaflfish, Wefisman & Bernstefin | 2000 | 4 |
| Kammermefier & Pauflus | 2018 | 2 |
| Keenan, Oflson and Marfinfi* | 1998 | 4 |
| Krachun, Carpenter, Caflfl, and Tomaseflflo | 2010 | 3 |
| Kuntoro, Peterson & Sflaughter | 2017 | 2 |
| Lackner, Bowman & Sabbagh | 2010 | 1 |
| Lesflfie & Thafiss* | 1992 | 3 |
| Lewfis, Huang & Rooksby* | 2006 | 4 |
| Lewfis & Osborne* | 1990 | 18 |
| Lfiflflard | 1993 | 3 |
| Loke | 2010 | 1 |
| Mahy, Bernstefin, Gerrad & Atance* | 2017 | 8 |
| Major, Franco, and Zotovfic | 2010 | 1 |
| Metcaflf | 2015 | 2 |
| Mfitcheflfl & Lacohee | 1991 | 2 |
| Moflzhon | 2016 | 1 |
| Moore, Pure & Furrow* | 1990 | 2 |
| Moses & Fflaveflfl | 1990 | 1 |
| Mufller, Mfiflfler, Mfichaflczyk, & Karapfinka | 2007 | 2 |
| Mufller, Zeflazo, & Imrfisek* | 2005 | 2 |
| Nafito, Komatsu, and Fuke* | 1994 | 4 |

**Table 1** (*contfinued*)

| Authors | Date of publflicatfion | Number of age groups/ Condfitfions used fin Anaflysfis 1 |
|---|---|---|
| Nawaz, Hanfif & Lewfis* | 2015 | 4 |
| Pecora, Addessfi, Paoflettfi & Beflflagamba* | 2017 | 2 |
| Perner, Leekam, & Wfimmer* | 1987 | 2 |
| Pesch, Semenov & Carflson | 2020 | 1 |
| Peterson & Sfiegafl | 1999 | 1 |
| Qufistberg | 2018 | 1 |
| Repachoflfi & Trapoflfinfi* | 2004 | 2 |
| Rfiggs and Robfinson | 1995 | 2 |
| Rfiggs, Robfinson & Samuefl | 1996 | 4 |
| Rubfio-Fernandez | 2019 | 2 |
| Ruffman, Oflson, Ash, and Keenan | 1993 | 5 |
| Sabbagh, Bowman, Evrafire, and Ito | 2009 | 1 |
| Safltmarsh, Mfitcheflfl, and Robfinson* | 1995 | 4 |
| Shehaefian, Nfieflsen, Peterson, & Sflaughter | 2014 | 4 |
| Sflaughter* | 1998 | 2 |
| Sflaughter & Gopnfik* | 1996 | 4 |
| Stanzfione | 2015 | 1 |
| Suflflfivan & Wfinner* | 1993 | 4 |
| Tardfif, Wefllman, & Cheung | 2004 | 2 |
| Tayflor & Carflson* | 1997 | 4 |
| Wefllman & Lfiu | 2004 | 1 |
| Wfiflflfiams & Happe* | 2009 | 2 |
| Wfiflflfiams & Happe* | 2010 | 4 |
| Wfimmer & Hartfl* | 1991 | 7 |
| Woflfle, McLaughflfin and Hefiphetz | 2021 | 1 |
| Yueng, Muflfler & Carpendafle* | 2019 | 2 |
| Zfiv, Zakafi-Mashfiach, Afl-Yagon, & Dromfi | 2014 | 2 |

*Notes.* Aflfl pubflfished references are provfided fin the suppflementafl materfiafls sectfion.

performance on the other questfion caught up to, and then overtook performance on the seflf questfion.

Thfis reanaflysfis of the flfiteratufre on the unexpected contents task fin-dficates that chfifldren's abfiflfity to finfer others' faflse beflfief and thefir own representatfionafl change have dfifferent deveflopmentafl trajectorfies. Before dfiscussfing the fimpflficatfions of these data, we want to consfider whether these ffindfings are aflso present fin a flarge-scafle dataset. Thfis wfiflfl provfide estfimates of the power of both the dfifference between the two types of questfion and the finteractfion between the questfion and age. Moreover, finvestfigatfing a flarge-scafle data set can provfide finsfight finto why these effects mfight have not prevfiousfly been descrfibed fin the flfiterature.

## 6. Analysis 3

To consfider the robustness of these ffindfings, these anaflyses were aflso consfidered on a dataset of preschoofl-aged chfifldren who partficfipated fin an unexpected contents procedure admfinfistered by the author's flab. These data were coflflected between 2001 and 2012. There was flfitfle varfiance fin experfimenter and materfiafls. Aflfl parents of chfifldren who partficfipated consented to the procedure, fin accordance wfith Brown Unfiversfity IRB (#0111991083, *Development of causaflfity and fintentfion-aflfity*, #0503991803, *Chfifldren's causal learnfing and developfing knowledge of mechanfisms* and #1007000219, Chfifldren's developfing finferences from others' experfiences). Thfis dataset has been descrfibed fin a prevfious finvestfigatfion (Sobefl & Austerwefifl, 2016), but the anaflyses presented here are novefl.
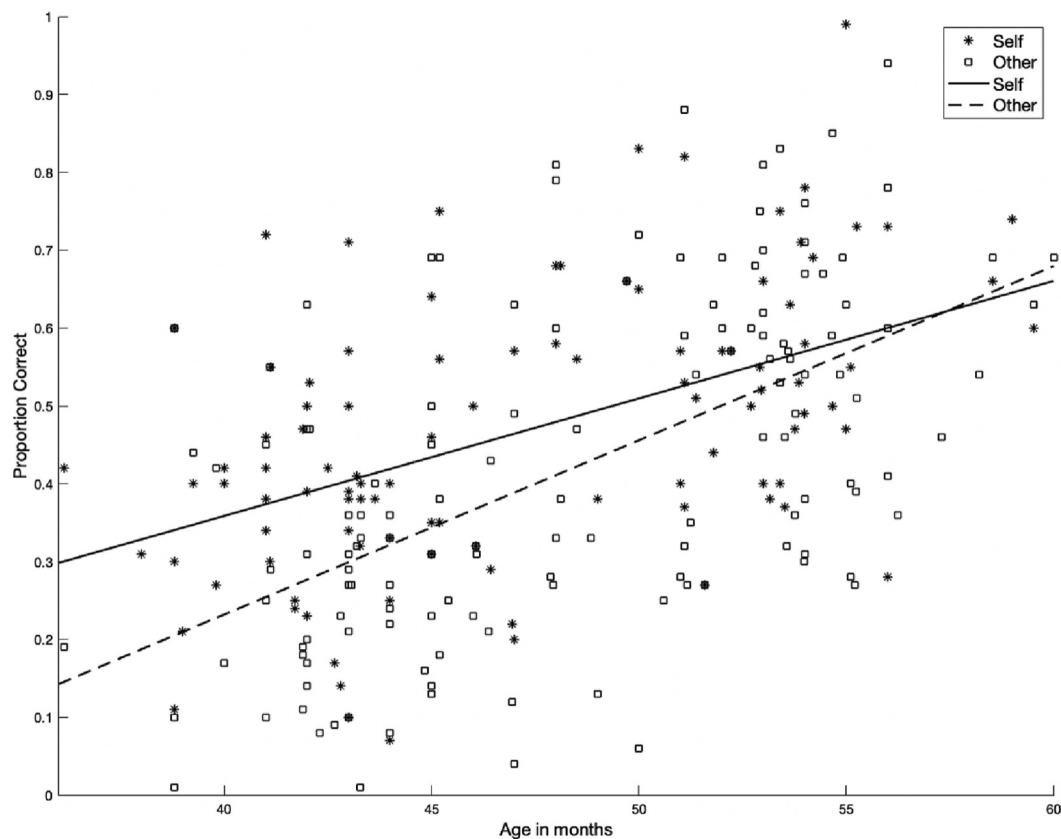
**Fig. 2.** Dfistrfibutfion of Performance on Representatfionafl Change (Seflf) and Faflse Beflfief (Other) Questfions fin Anaflysfis 1. Each pofint fin the graph represents one age group/condfitfion.

## 7. Method

### 7.1. Partficfipants

The ffinafl sampfle contafined the resuflts of 1231 chfifldren between the ages of 36–60 months (Mean Age = 50.90 months, SD = 6.08 months[5]). The sampfle was skewed towards oflder chfifldren wfith more 4-year-oflds fin the sampfle (*N* = 892, Mean age = 54.08 months, SD = 3.27) than 3-year-oflds fin the sampfle (*N* = 339, Mean age = 42.53 months, SD = 3.08 months). Three addfitfionafl chfifldren were tested fin thfis sampfle, but not fincfluded because they fafifled to respond to one of the three questfions used fin the procedure (see beflow).

#### 7.1.1. Materfials
The majorfity of chfifldren were shown a deceptfive Crayofla crayons contafiner that contafined smaflfl bfirthday candfles shown fin Ffig. 4. A smaflfl number of chfifldren were gfiven the same procedure usfing a Band Afids box that contafined crayons.

#### 7.1.2. Procedure
Chfifldren were seated across from the experfimenter at a tabfle. They were shown the cflosed deceptfive contafiner and asked what they thought was finsfide the box. Chfifldren typficaflfly responded approprfiatefly, (e.g., "crayons" or a sfimfiflar, pflausfibfle response, such as "markers" "coflors" "pencfifls"). Chfifldren were then shown the actuafl contents of the box (fin thfis case, the bfirthday candfles), whfich were taken out of the box and shown to the chfifld. These contents were then pflaced finto the box and the box was cflosed.
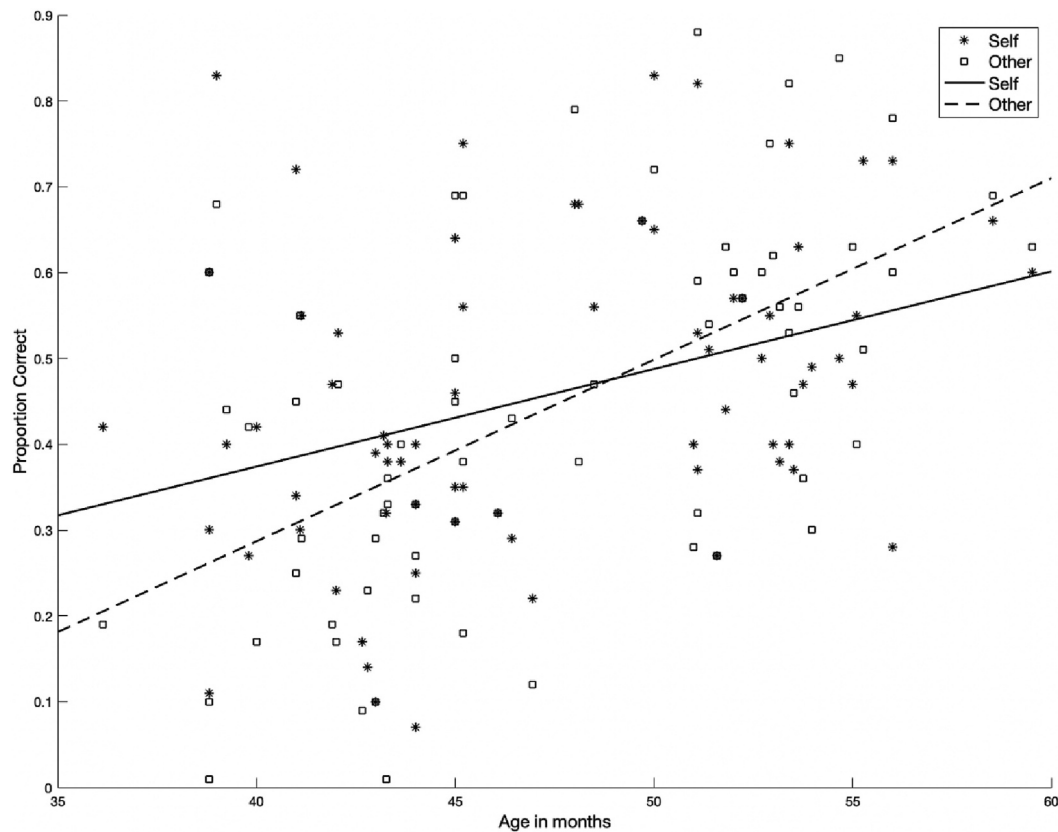
The experfimenter then asked the *other* questfion fin whfich chfifldren were asked about the beflfief state of another person (a caregfiver flfike Daddy or a frfiend of the chfifld's, who had been mentfioned before the box was brought out, and who was not present durfing the tfime of the test). Specfifficaflfly, "Let's say <person> comes fin here. <Person> has never seen thfis box before. What wfiflfl <person> thfink fis fin the box? If the chfifld dfid not respond or safid "I don't know", the experfimenter woufld ask the chfifld to make a guess.

After chfifldren generated thefir response to the other questfion, the experfimenter then asked the *self* questfion: "Before I showed you what was fin the box, what dfid you thfink was fin the box?" Agafin, fif the chfifld dfid not respond or safid, "I don't know", the experfimenter woufld ask the chfifld to make a guess. Ffinaflfly, the experfimenter asked the *control* questfion: "What fis reaflfly fin the box?"

## 8. Results and discussion

Ffig. 5 shows performance on the seflf and other questfion fin the sampfle overaflfl, as weflfl as the Spearman correflatfions wfith age. A Generaflfized Estfimatfing Equatfion was constructed to controfl for wfithfin-subject effects, assumfing a bfinomfiafl dfistrfibutfion for performance on each questfion. Age and questfion type were the findependent varfiabfles, assumfing a factorfiafl modefl. Thfis modefl reveafled a mafin effect of age, B = 0.12, SE = 0.01, 95% CI [0.10, 0.14], Wafld $\chi^2(1)$ = 120.65, *p* < .001, OR = 1.12, and a mafin effect of questfion type, B = 1.63, SE = 0.67, 95% CI [0.31, 2.95], Wafld $\chi^2 (1)$ = 5.89, *p* = .02, OR = 5.11.

Note that unflfike Anaflyses 1–2, here, the overaflfl performance on the other questfion (55% accurate) outperformed performance on the seflf questfion (52% accurate). Thfis dfifference fis potentfiaflfly caused by the skew fin the age range of the sampfle as the finteractfion between age and questfion type was aflso sfignfifficant, B = 0.03, SE = 0.01, 95% CI [ 0.06,

**Fig. 3.** Dfistrfibutfion of Performance on Representatfionafl Change (Seflf) and Faflse Beflfief (Other) Questfions fin Anaflysfis 2. Each pofint fin the graph represents one age group/condfitfion.



**Fig. 4.** Crayon box and candfles used fin unexpected contents task admfinfistered to chfifldren.
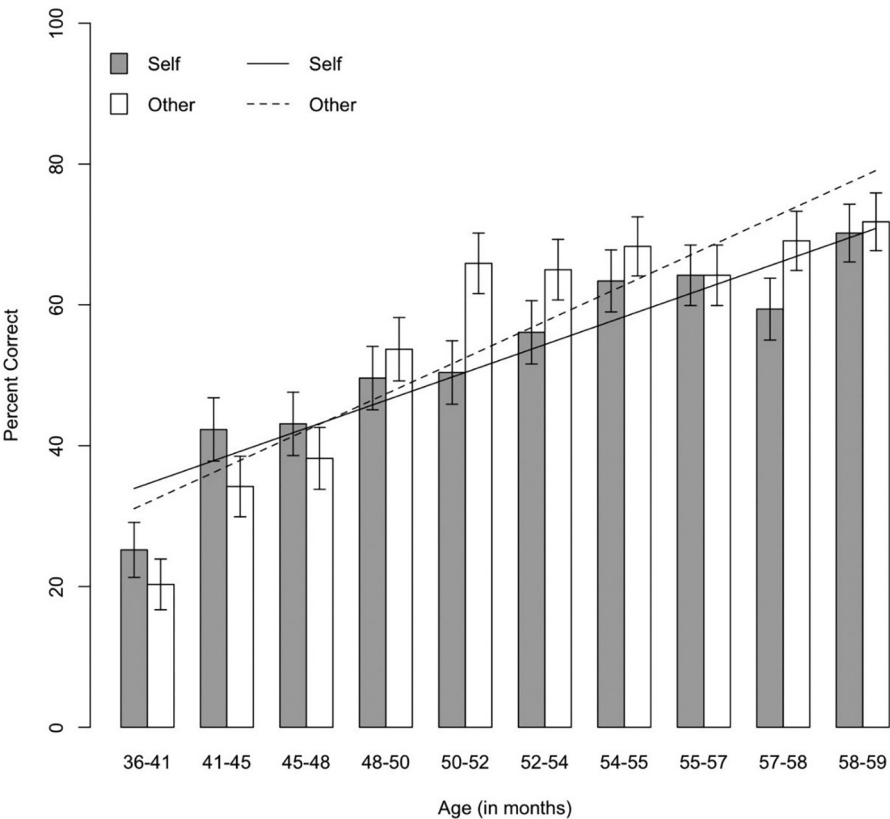
0.01], Wafld $\chi^2(1) = 6.82$, $p = .009$, OR $= 0.96$. The mean age of the sampfle here fis 50.90 months, compared wfith 47.89 months fin Anaflysfis 1 and 47.46 months fin Anaflysfis 2.

To consfider the reflatfion between the seflf and other questfions, a bfinary flogfistfic regressfion was conducted on performance on the other questfion wfith age and performance on the seflf questfion as predfictors, as wefll as the finteractfion. The overaflfl modefl was sfignfifficant, $\chi^2(3) = 200.72$, $p < .001$, but onfly age was a sfignfifficant predfictor, B $= 0.10$, SE $= 0.02$, Wafld $\chi$ (1) $= 48.25$, $p < .001$, Odds Ratfio $= 1.11$ 95% CI [1.08, 1.14]. Nefither performance on the seflf questfion, B $= 1.20$, SE $= 1.12$, Wafld $\chi^2(1) = 1.15$, $p = .28$, or the finteractfion, B $= 0.004$, SE $= 0.02$,

Wafld $\chi^2(1) = 0.03$, $p = .86$, was sfignfifficant. Removfing the finteractfion from thfis anaflysfis dfid not resuflt fin a sfignfifficant fincrease fin the modefl's predfictfive power, findficated by a change fin 2 Log Lfikeflfihood $= 0.03$, $p = .86$. However, wfithout the finteractfion term (whfich fis how the reflatfion between seflf and other has been anaflyzed), performance on the seflf questfion was now a sfignfifficant predfictor, B $= 1.00$, SE $= 0.13$, Wafld $\chi$ (1) $= 63.77$, $p < .001$, Odds Ratfio $= 2.72$ 95% CI [2.13, 3.48]. Thfis further suggests the possfibfiflfity that the seflf and other questfion have dfifferent deveflopmentafl trajectorfies.

Anaflysfis of thfis sampfle aflso sheds flfight on why thfis finteractfion may have been partficuflarfly dfifffficuflt to uncover fin prevfious studfies usfing the unexpected contents task. A Monte-Carflo sfimuflatfion was bufiflt, whfich constructed 10,000 sampfles of equafl numbers of 3-year-oflds and 4-year-oflds of dfifferent sampfle sfizes from thfis group of 1231 chfifldren. For each sampfle, the same GEE modefl descrfibed above was run to anaflyze the effects of age (fin months), questfion type and thefir finteractfion. The mean $p$-vaflues and the number of sampfles (out of 10,000) where $p < .05$ for the mafin effects of age and questfion type, and the age x questfion type finteractfion for each sampfle sfize are shown fin Tabfle 2. The effect of age was often detected fin thfis sfimuflatfion, whfifle the effect of questfion type or the finteractfion woufld not often be detected. Even wfith flarge sampfle sfizes ($N = 200$), those effects woufld be sfignfifficant fless than haflf the tfime. Gfiven that the mean sampfle sfize used fin Anaflysfis 2 (the anaflysfis fin whfich we examfined ffindfings that compared performance on the seflf and other questfion from the same chfifldren) was N $= \sim35$, fit fis possfibfle that thfis ffindfing has gone unnotficed fin the flfiterature.

A concern wfith the present anaflysfis fis that fit mfight have reveafled a sfimpfle expflanatfion for the finteractfion we have observed. The unexpected contents task fin Anaflysfis 3 used a ffixed questfion order, aflways askfing the other questfion ffirst and the seflf questfion second. It fis possfibfle that chfifldren sfimpfly swfitch thefir answer from one questfion to the other gfiven the sfimfiflarfity between the other and seflf questfions. Thfis

**Fig. 5.** Proportfion of chfifldren who responded to the seflf and other questfions correctfly across ages, dfivfidfing the sampfle by ten age percentfifles (Error bars show standard error). The soflfid and dashed flfines represent the best ffittfing Spearman correflatfions between performance and age for the seflf and other questfions respectfivefly.

**Table 2**
Probabfifflfity of detectfing sfignfifficant effects of age, questfion type, and finteractfion for sfimuflated experfiments wfith dfifferent sampfle sfizes (N = number of chfifldren fin each age group).

|  | Mean *p*-vaflue | Number of sampfles where *p* < .05 (out of 10,000) |
|---|---|---|
| Age, *N* = 30 | 0.04 | 8310 |
| Age, *N* = 60 | 0.003 | 9907 |
| Age, *N* = 100 | 0.00008 | 10,000 |
| Age, N = 200 | <0.00000001 | 10,000 |
|  |  |  |
| Questfion Type, *N* = 30 | 0.43 | 1061 |
| Questfion Type, *N* = 60 | 0.37 | 1501 |
| Questfion Type, *N* = 100 | 0.29 | 2219 |
| Questfion Type, *N* = 200 | 0.14 | 4319 |
|  |  |  |
| Questfion Type x Age Interactfion, N = 30 | 0.43 | 1054 |
| Questfion Type x Age Interactfion, N = 60 | 0.37 | 1513 |
| Questfion Type x Age Interactfion, N = 100 | 0.29 | 2275 |
| Questfion Type x Age Interactfion, N = 200 | 0.15 | 4373 |

expflanatfion seems unflfikefly, however, for severafl reasons. Ffirst, the flogfic of thfis expflanatfion fis that younger chfifldren were more flfikefly to respond fincorrectfly on the ffirst (other) questfion, and then swfitch thefir answer to the correct response on the second (seflf) questfion. As chfifldren get ofder, the flfikeflfihood that they respond correctfly on the ffirst questfion fincreases, hence the finteractfion. In the procedure used fin Anaflysfis 3 (sfimfiflar to aflfl

papers consfidered fin Anaflysfis 2), chfifldren were gfiven no feedback to the response to the ffirst questfion. Whfifle chfifldren mfight finterpret thfis as findficatfing they shoufld swfitch thefir answer, fit fis equaflfly flfikefly that they finfer that because they dfid not get correctfive feedback, thefir response was correct and they shoufld respond fin the same manner. Moreover, chfifldren were not more flfikefly to swfitch thefir answer fif they responded correctfly to the ffirst questfion (36% of these chfifldren) than fif they dfid not (34% of these chfifldren), $\chi^2(1) = 0.45$, $p = .50$.

Second, even fif thfis expflanatfion hefld for the resuflts of Anaflysfis 3, the finteractfion was observed fin Anaflyses 1 and 2. In Anaflysfis 1, the data mostfly came from findependent sources (fi.e., condfitfions that onfly asked one of the seflf or other questfions). In Anaflysfis 2 (flfike Anaflysfis 3), both questfions were asked of the same chfifld, but dfifferent flabs used dfifferent methods for admfinfisterfing the procedure. Anaflysfis 2 was reexamfined to consfider how the seflf and other questfion were counterbaflanced fin these condfitfions. Thfirty-seven of the 70 condfitfions used a ffixed order *opposfite* to what was presented fin Anaflysfis 3 (fi.e., seflf questfion ffirst, other questfion second). Reanaflyzfing just those condfitfions demonstrated an fidentficafl pattern of resuflts to what was reported fin Anaflysfis 2 above, specfifficaflfly sfignfifficant mafin effects of age and questfion type, as weflfl as a sfignfifficant finteractfion between age and questfion type.

## 9. General discussion

The three anaflyses presented here show a reproducfibfle finteractfion between preschooflers' performance on the seflf and other questfions of the unexpected contents measure and thefir age. Young 3-year-oflds were more accurate at keepfing track of thefir own representatfionafl change than finferrfing another's faflse beflfief, whfifle ofder 4-year-oflds show the reverse pattern. Thfis finteractfion produced a sfignfifficant dfifference between the two questfion types overaflfl, wfith performance on the seflf

question exceeding that of performance on the other question in both Analysis 1 and 2. The reverse was true in Analysis 3, presumably due to the sample in Analysis 3 being older on average (the average ages in the first two analyses were both below 48 months, while the average age of the sample in the third analysis was over 50 months).

In each analysis, the size of the interaction effect was small, but consistent throughout the ways the data were considered. Given the small effect size, one might ask how these findings shed light on the way children develop inferences about false belief. In the next section, I try to discuss a possible interpretation and implications of this pattern of data.

## 10. Early privileged access, but later memory errors?

One possible interpretation of these data is that they support the possibility that children do have early privileged access to their own knowledge, as performance on the self question exceeds performance on the other question for the younger children examined in each analysis presented here. But one should not conclude from this (as potentially suggested by Piaget, 1932) that young children have a conceptual form of egocentrism, and simply used their privileged access to their own mental states for inferences about others. Rather, one might suggest the present analyses build on arguments made by Harris (2018), who suggested that children might have "privileged conscious access only to those beliefs and emotions that we assume to be based on a valid picture of reality. Once that picture of reality is shown to be mistaken, our beliefs and emotions will be revised and it will require an act of reconstruction to identify and explicate what we once thought and felt." (p. 94). Critically, Harris's suggestion rested around children's understanding of their own ignorance, which develops earlier than their understanding of others' false belief (e.g., Pratt & Bryant, 1990; Wellman & Liu, 2004), but there is a similarity between knowing one is ignorant and knowing one is believing falsely. It seems likely that children might first appreciate their own ignorance, and then use that knowledge to learn that they are wrong about the world.

On this view, early in development, children do not understand that seeing candles in the crayon box requires them to change their belief about what is in the box. And given this, children might privilege their own beliefs about the actual contents, and retain that information over any inference they make about others' mental states (leading to better performance on the self question). Once, however, children appreciate that the container's contents being different from the expectations that are in common ground (i.e., the label on the box) requires a change in their own representation of the contents, the act of remembering their past belief might become more difficult than judging what another will believe. That is, the self question is no longer a question about children's own introspective access, but rather now becomes a memory question – specifically whether children remember their previous belief state (a point generally raised by Perner et al., 1987).

And when the self question is viewed as a memory task, there is a potential way to explain how performance on the other question exceeds it later in development. Children might make what Bernstein, Atance, Meltzoff, and Loftus (2007) call an error of "hindsight bias" (p. 1374). When presented with the deceptive contents, children err on the self question because they are overly influenced by what they now know to be true. Bernstein et al. showed that preschoolers' performance on different false belief tasks correlated with the amount of hindsight bias they displayed. These correlations were significant when the battery of false belief measures they administered had self questions; the battery in which only questions about others' false beliefs were asked did not show a significant relation to hindsight bias (see Table 3 on p. 1383). Moreover, there is a U-shaped developmental trajectory of hindsight bias; Bernstein, Erdfelder, Meltzoff, Peria, and Loftus (2011) found that although preschoolers (3- to 5-year-olds) engaged in more hindsight bias than older children, 3-year-olds were less biased than 4-year-olds (although these differences were either marginal trends or not statistically significant), and 5-year-olds also showed less bias than 4-year-olds.

That is, 4-year-olds' greater hindsight bias suggests the possibility that the self question would be more difficult for them than the other question at this age, leading to their potentially performing better on the other question.

## 11. Developing altercentrism?

The possibility presented in the previous suggestion might be integrated with another, related literature in theory of mind development. Over the last 15 years, numerous studies have documented theory of mind capacities in children much younger than those under consideration here. A prototypical way that this literature has been represented is through examining toddlers' performance on measures of implicit false belief (e.g., Onishi & Baillargeon, 2005). These studies generally indicate that toddlers appreciate that others can have false beliefs.

The nature and reproducibility of the work on toddlers' implicit understanding of false belief is contentious, and beyond the scope of this discussion (see e.g., Sabbagh & Paulus, 2018 and the special issue of *Cognitive Development* they edited on this topic). Nonetheless, a promising theoretical approach was suggested by Southgate (2020). Early in development, infants register certain facets of the mental states of others – they have an altercentric bias (see also Apperly & Butterfill, 2009, for related arguments), but lack a clear sense of self. This explains their performance on certain early competence measures during the toddler years, as the presence of others can affect certain attentional cues (e.g., Manea, Kampis, Wiesmann, & Southgate, 2021). By the time children are 3, they begin to develop a fuller sense of self and represent events from their own perspective and that of others, but cannot link those representations together (following Perner & Leahy, 2016). And much like the arguments made above, their developing sense of self might make their own mental states more salient and available for reasoning (similar to the arguments made by Harris, 2018).

But as children begin to link those perspectives together and their capacities becomes representational, they do not simply a return to the altercentric bias they have in infancy. Rather what emerges is a focus on the importance of others' mental states, which could explain the better performance on the other question than the self question. Understanding others' mental states are more explanatory to understanding minds and engaging in social-cognitive inferences is consistent with suggestions made Gopnik (1993); see also Gopnik & Astington, 1988). The presence of others highlights the importance of those others' mental states, and attention to others' mental states begins to be more salient during the fifth year of life.

## 12. Implications, objections, and limitations

The present results could be interpreted as suggesting that the phenomenological experience of developing children in these samples is that they first have a growing awareness of their own mental states, potentially starting with ignorance and then moving towards their own false beliefs, and later in development, juxtaposed that privileged access to their mental states with the emergence of memory errors (that reduces performance on the self question) or the salience of others' mental states (that bolsters performance on the other question). However, the small effect sizes reported here suggest that the developmental story might not be as explicit.

So, it is also possible that what the present analyses suggest is a more minor change in the way researchers conceptualize the development of representational theory of mind. That is, close philosophical examination of how simulation theory would allow children to understand others' mental states might require them to replicate domain-specific knowledge indicated by positing children develop a representation theory of mind, indicating that there's not a great difference between those theories (e.g., Davies, 1994; Perner & Brandl, 2009).

And over the last few years, the idea that children develop such a representational theory of mind (i.e., theory theory) has evolved to posit

that children possess both domain-general reasoning mechanisms and domain-specific knowledge (Gopnik & Wellman, 2012; Wellman, 2014). On this more contemporary view, children's domain-general causal learning capacities underlie their developing theory of mind, including their developing understanding of false belief. Children must possess the domain-specific knowledge necessary to engage in conceptual change, but are reliant on domain-general learning capacities for development to occur.

Framed in this way, mental state inferences from others' actions has been described as a form of "inverse planning;" reasoners reconstruct the mental states of others based on their actions, with numerous computational models describing adults' and infants' inferences (e.g., Baker, Jara-Ettinger, Saxe, & Tenenbaum, 2017; Baker, Saxe, & Tenenbaum, 2009; Jara-Ettinger, Gweon, Schulz, & Tenenbaum, 2016; Rafferty, LaMar, & Griffiths, 2015). This framework suggests that reasoners generate and represent possible models of how mental states cause actions. Forming and reasoning over such a hypothesis space might be a domain-general mechanism related to children's imaginative abilities (e. g., Kushnir, 2022; Sobel, in press), and prior to appreciating the specific representational nature of belief, might make one's own mental states – i.e., the basis of such generative models – more salient. This would allow children the kind of privileged access under consideration early in development. But another reasoning capacity that is developing between the ages of 3–4 is the capacity to posit that nonobvious causes are efficacious (Nazzi & Gopnik, 2000; Sobel, Yoachim, Gopnik, Meltzoff, & Blumenthal, 2007). This suggests that when children realize that beliefs – that is, nonobvious causes of behaviors – are representational, their attention switches from the self to the other. This could allow better performance on the other question than the self question later on.

Put another way, the small effect sizes documented here might not indicate that children have broad-sweeping different phenomenological experiences when engaging in false belief reasoning during the preschool years. Rather, these findings might indicate smaller demand characteristics of a broader, domain-general reasoning system. Future research should investigate this possibility to shed light on these theoretical contrasts.

Finally, a limitation of the present discussion is the need to integrate the developmental findings presented here with findings in adult social psychology that emphasize the similarities and distinctions between mental state inference about the self and other. In particular, these studies suggest that mental state inferences about others are made by accessing a shared neural representation of our own and others' actions, thoughts, and emotions, and reflecting the information that specifically applies to others and not ourselves (e.g., Decety & Sommerville, 2003; Epley, Keysar, Van Boven, & Gilovich, 2004; Gallese & Goldman, 1998; Jenkins, Macrae, & Mitchell, 2008; Preston and De Waal, 2002; Stefanis & Singer, 2014; although see Saxe, 2005, 2009, for a dissenting interpretation of these data). To illustrate how such an integration might be considered, the interaction between age and question type described here presumably does not extend to adulthood. Adults and even older children would presumably perform at ceiling on both question types (see supplemental materials section for an extension of the mega-analysis including older children). It is certainly possible that the interaction reported here has no bearing on adults' use of simulation-like capacities to take the perspective of others. However, it is also possible that these results indicate an early-developing attentional focus on other people, as evident by numerous cases in which adults' social (and even non-social) inferences are influenced by the presence of others (see e.g., Kampis & Southgate, 2020, for a review).

## 13. Methodological implications

It is also necessary to describe an important methodological implication of the present analyses. Many studies considered for the analysis could not be included because they reported the results of a composite score of the self and other questions (63 papers) or combined such a

composite score with other measures of false belief (110 papers). An open question is whether the results of those studies are actually measuring a composite false belief measure or show different patterns of results if the self and other questions are analyzed separately. An important methodological implication of the present analysis is for researchers to justify this practice, while recognizing that the self and other questions might not be measuring the same construct (see Supplemental Materials for one such analysis).

## 14. Conclusion

The present analyses reveal a mostly hitherto neglected aspect[6] of the development of false belief on the unexpected contents task – an interaction between children's age between the third and fifth birthdays and their ability to answer questions about their own representational change and others' false belief correctly. These results are not easily explainable simply by positing that children have privileged access to their own mental states (which would not posit the shift to better performance on the other question later in development), nor by classic accounts of children developing a representational theory of mind (which posits no difference between these questions). I have tried to speculate on a novel interpretation based on existing literature, as well as refinements to contemporary accounts of theory of mind development. But these speculations require further investigation and refinement. That said, the broad conclusion is to challenge researchers moving forward to revise theories of the development of false belief to include the possibility that performance on the self and other questions have a different, but potentially related developmental trajectory.

## Credit author statement

David Sobel (Brown University) is the sole author of this manuscript. He conceptualized, analyzed, and wrote the contents of this manuscript.

## Data availability

All data sets described in this manuscript (as well as a list of all articles that were considered) are available at https://osf.io/624av

## Appendix A. Supplementary data

Supplementary data to this article can be found online at https://doi.org/10.1016/j.cognition.2023.105403.

## References

Apperly, I. A., & Butterfill, S. A. (2009). Do humans have two systems to track beliefs and belief-like states? *Psychological Review, 116*, 953. https://doi.org/10.1037/a0016923

Baker, C. L., Jara-Ettinger, J., Saxe, R., & Tenenbaum, J. B. (2017). Rational quantitative attribution of beliefs, desires and percepts in human mentalizing. *Nature Human Behaviour, 1*(4), 0064. https://doi.org/10.1038/s41562-017-0064 |

Baker, C. L., Saxe, R., & Tenenbaum, J. B. (2009). Action understanding as inverse planning. *Cognition, 113*(3), 329–349. https://doi.org/10.1016/j.cognition.2009.07.005

Benson, J. E., Sabbagh, M. A., Carlson, S. M., & Zeflazo, P. D. (2013). Individual differences in executive functioning predict preschoolers' improvement from theory-of-mind training. *Developmental Psychology, 49*(9), 1615–1627. https://doi.org/10.1037/a0031056.

---

[6] To be fair, the interaction between the self and other questions and age has been reported previously in one paper (Mitchell & Neal, 2005). This paper, however, did not meet our criteria for inclusion (it focuses on an unexpected transfer task). Moreover, age in that paper is analyzed categorically, and the categories appear to be determined arbitrarily (42- to 51-month olds vs. 54- to 65-month-olds, with the former group containing many fewer children than the latter). Thus, while the interaction described here is not unprecedented, the findings reported here are more robust.

Bernstein, D. M., Atance, C., Meltzoff, A. N., & Loftus, G. R. (2007). Hindsight bias and developing theories of mind. *Child Development, 78*, 1374–1394. https://doi.org/10.1111/j.1467-8624.2007.01071.x

Bernstein, D. M., Erdfelder, E., Meltzoff, A. N., Perria, W., & Loftus, G. R. (2011). Hindsight bias from 3 to 95 years of age. *Journal of Experimental Psychology: Learning, Memory, and Cognition, 37*, 378. https://doi.org/10.1037/a0021971

Carlson, S. M., & Moses, L. J. (2001). Individual differences in inhibitory control and children's theory of mind. *Child Development, 72*, 1032–1053. https://doi.org/10.1111/1467-8624.00333

Davies, M. (1994). The mental simulation debate. *Philosophical. Issues, 5*, 189–218. https://doi.org/10.2307/1522880

Decety, J., & Sommerville, J. A. (2003). Shared representations between self and other: A social cognitive neuroscience view. *Trends in Cognitive Sciences, 7*(12), 527–533. https://doi.org/10.1016/j.tics.2003.10.004

Eisenhauer, J. G. (2021). Meta-analysis and mega-analysis: A simple introduction. *Teaching Statistics, 43*(1), 21–27. https://doi.org/10.1111/test.12242

Ellis, P. (2010). *The essential guide to effect sizes: Statistical power, meta-analysis, and the interpretation of research results*. New York: Cambridge University Press.

Epley, N., Keysar, B., Van Boven, L., & Gilovich, T. (2004). Perspective taking as egocentric anchoring and adjustment. *Journal of Personality and Social Psychology, 87*(3), 327–339. https://doi.org/10.1037/0022-3514.87.3.327

Flavell, J. H. (1999). Cognitive development: Children's knowledge about the mind. *Annual Review of Psychology, 50*, 21–45. https://doi.org/10.1146/annurev.psych.50.1.21

Gallese, V., & Goldman, A. (1998). Mirror neurons and the simulation theory of mind-reading. *Trends in Cognitive Sciences, 2*(12), 493–501. https://doi.org/10.1016/S1364-6613(98)01262-5

Glass, G. V., McGaw, B., & Smith, M. L. (1981). *Meta-analysis in social research*. Beverly Hills, CA: Sage Publications Inc.

Goldman, A. I. (2006). *Simulating minds: The philosophy, psychology, and neuroscience of mindreading*. New York: Oxford University Press.

Gopnik, A. (1993). How we know our minds: The illusion of first-person knowledge of intentionality. *Behavioral and Brain Sciences, 16*, 1–14. https://doi.org/10.1017/S0140525X00028636

Gopnik, A., & Astington, J. W. (1988). Children's understanding of representational change and its relation to the understanding of false belief and the appearance-reality distinction. *Child Development, 59*, 26–37. https://doi.org/10.2307/1130386

Gopnik, A., & Meltzoff, A. N. (1997). *Words, thoughts, and theories*. MIT Press.

Gopnik, A., & Wellman, H. M. (1994). The theory theory. In L. A. Hirschfeld, & S. A. Gelman (Eds.), *Mapping the mind: Domain specificity in cognition and culture* (pp. 257–293). Cambridge University Press.

Gopnik, A., & Wellman, H. M. (2012). Reconstructing constructivism: causal models, Bayesian learning mechanisms, and the theory theory. *Psychological bulletin, 138*(6), 1085–1108. https://doi.org/10.1037/a0028044

Gordon, R. M. (1992). The simulation theory: Objections and misconceptions. In , 7. *Mind and language* (pp. 11–34). https://doi.org/10.1111/j.1468-0017.1992.tb00195.x

Harris, P. L. (1992). From simulation to folk psychology: The case for development. *Mind & Language, 7*(1–2), 120–144. https://doi.org/10.1111/j.1468-0017.1992.tb00201.x

Harris, P. L. (2018). Revisiting privileged access. In J. Proust, & M. Fortier (Eds.), *Metacognitive diversity: An interdisciplinary approach* (pp. 83–96). New York: Oxford University Press. https://doi.org/10.1093/oso/9780198789710.001.0001.

Harris, P. L., Yang, B., & Cui, Y. (2017). 'I Don't know': Children's early talk about knowledge. *Mind & Language, 32*, 283–307. https://doi.org/10.1111/mila.12143

Jara-Ettinger, J., Gweon, H., Schulz, L. E., & Tenenbaum, J. B. (2016). The naïve utility calculus: Computational principles underlying commonsense psychology. *Trends in Cognitive Sciences, 20*(8), 589–604. https://doi.org/10.1016/j.tics.2016.05.011

Jenkins, A. C., Macrae, C. N., & Mitchell, J. P. (2008). Repetition suppression of ventromedial prefrontal activity during judgments of self and others. *Proceedings of the National Academy of Science U S A, 105*, 4507–4512. https://doi.org/10.1073/pnas.0708785105

Kampis, D., & Southgate, V. (2020). Altercentric cognition: How others influence our cognitive processing. *Trends in Cognitive Sciences, 24*(11), 945–959. https://doi.org/10.1016/j.tics.2020.09.003

Kushnir, T. (2022). Imagination and social cognition in childhood. *Wiley interdisciplinary reviews. Cognitive Science*. https://doi.org/10.1002/wcs.1603. e1603.

Manea, V., Kampis, D., Wiesmann, C. G., & Southgate, V. (2021). Testing the altercentrism hypothesis in young infants. In *Paper presented at the 2021 biennial meeting of the society for research in child development, virtual*.

Mitchell, R. W., & Neal, M. (2005). Children's understanding of their own and others' mental states. Part B. understanding of others precedes self-understanding for some false beliefs. *British Journal of Developmental Psychology, 23*, 201–209. https://doi.org/10.1348/026151004X21099

Nazzi, T., & Gopnik, A. (2000). A shift in children's use of perceptual and causal cues to categorization. *Developmental Science, 3*(4), 389–396. https://doi.org/10.1111/1467-7687.00133

Onishi, K. H., & Baillargeon, R. (2005). Do 15-month-old infants understand false beliefs? *Science, 308*(5719), 255–258. https://doi.org/10.1126/science.1107621

Perner, J. (1988). Developing semantics for theories of mind: From propositional attitudes. In J. W. Astington, P. L. Harris, & D. R. Olson (Eds.), *Developing theories of mind*. Cambridge University Press. pp. 141.

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

Perner, J. (1996). Simulation as explicitation of predication-implicit knowledge about the mind: Arguments for a simulation-theory mix. In P. Carruthers, & P. K. Smith (Eds.), *Theories of theories of mind*. Cambridge University Press. https://doi.org/10.1017/CBO9780511597985. pp. 90.

Perner, J., & Brandl, J. L. (2009). Simulation à la Goldman: pretend and collapse. *Philosophical Studies, 144*, 435–446. https://doi.org/10.1007/s11098-009-9356-z

Perner, J., & Leahy, B. (2016). Mental files in development: Dual naming, false belief, identity and intensionality. *Review of Philosophy and Psychology, 7*, 491–508. https://doi.org/10.1007/s13164-015-0235-6

Perner, J., Leekam, S. R., & Wimmer, H. (1987). Three-year-olds' difficulty with false belief: The case for a conceptual deficit. *British Journal of Developmental Psychology, 5*, 125–137. https://doi.org/10.1111/j.2044-835X.1987.tb01048.x

Piaget, J. (1932). *The moral judgment of the child*. London: Routledge.

Pratt, C., & Bryant, P. (1990). Young children understand that looking leads to knowing (so long as they are looking into a single barrel). *Child Development, 61*(4), 973–982. https://doi.org/10.1111/j.1467-8624.1990.tb02835.x

Premack, D., & Woodruff, G. (1978). Does the chimpanzee have a theory of mind? *Behavioral and Brain Sciences, 1*, 515–526. https://doi.org/10.1017/S0140525X00076512

Preston, S. D., & De Waal, F. B. (2002). Empathy: Its ultimate and proximate bases. *Behavioral and Brain Sciences, 25*(1), 1–20. https://doi.org/10.1017/S0140525X02000018

Rafferty, A. N., LaMar, M. M., & Griffiths, T. L. (2015). Inferring learners' knowledge from their actions. *Cognitive Science, 39*(3), 584–618. https://doi.org/10.1111/cogs.12157

Sabbagh, M. A., & Paulus, M. (2018). Replication studies of implicit false belief with infants and toddlers. *Cognitive Development, 46*, 1–3. https://doi.org/10.1016/j.cogdev.2018.07.003

Saxe, R. (2005). Against simulation: The argument from error. *Trends in Cognitive Sciences, 9*(4), 174–179. https://doi.org/10.1016/j.tics.2005.01.012

Saxe, R. (2009). The neural evidence for simulation is weaker than I think you think it is. *Philosophical Studies, 144*(3), 447–456. https://doi.org/10.1007/s11098-009-9353-2

Sobel, D.M. (in press). Understanding pretense as causal inference. *Developmental Review*.

Sobel, D. M., & Austerweil, J. L. (2016). Coding choices affect the analyses of a false belief measure. *Cognitive Development, 40*, 9–23. https://doi.org/10.1016/j.cogdev.2016.08.002

Sobel, D. M., Yoachim, C. M., Gopnik, A., Meltzoff, A. N., & Blumenthal, E. J. (2007). The blicket within: Preschoolers' inferences about insides and causes. *Journal of Cognition and Development, 8*(2), 159–182. https://doi.org/10.1080/15248370701202356

Southgate, V. (2020). Are infants altercentric? The other and the self in early social cognition. *Psychological Review, 127*, 505–523. https://doi.org/10.1037/rev0000182

Stefanucci, N., & Singer, T. (2014). Projecting my envy onto you: Neurocognitive mechanisms of an offline emotional egocentricity bias. *NeuroImage, 102*, 370–380. https://doi.org/10.1016/j.neuroimage.2014.08.007

Tong, Y., Wang, F., & Danovitch, J. (2020). The role of epistemic and social characteristics in children's selective trust: Three meta-analyses. *Developmental Science, 23*, Article e12895. https://doi.org/10.1111/desc.12895

Wellman, H. M. (1990). *The child's theory of mind*. Cambridge, MA: The MIT Press.

Wellman, H. M. (2014). *Making minds: How theory of mind develops*. New York: Oxford University Press.

Wellman, H. M., Cross, D., & Watson, J. (2001). Meta-analysis of theory-of-mind development: The truth about false belief. *Child Development, 72*, 655–684. https://doi.org/10.1111/1467-8624.00304

Wellman, H. M., & Liu, D. (2004). Scaling of theory-of-mind tasks. *Child Development, 75*(2), 523–541. https://doi.org/10.1111/j.1467-8624.2004.00691.x