

1 Modeling individual differences in children's information integration during pragmatic word
2 learning

3 Manuel Bohn¹, Louisa S. Schmidt², Cornelia Schulze², Michael C. Frank³, & Michael Henry
4 Tessler^{4,5}

5 ¹ Department of Comparative Cultural Psychology, Max Planck Institute for Evolutionary
6 Anthropology, Leipzig, Germany

7 ² Leipzig Research Center for Early Child Development, Leipzig University, Leipzig, Germany

8 ³ Department of Psychology, Stanford University, Stanford, USA

9 ⁴ DeepMind, London, UK

10 ⁵ Department of Brain and Cognitive Sciences, Massachusetts Institute of Technology,
11 Cambridge, USA

M. Bohn received funding from the European Union’s Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement no. 749229. M. H. Tessler was funded by the National Science Foundation SBE Postdoctoral Research Fellowship Grant No. 1911790. M. C. Frank was supported by a Jacobs Foundation Advanced Research Fellowship and the Zhou Fund for Language and Cognition. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

The authors made the following contributions. Manuel Bohn: Conceptualization, Methodology, Formal Analysis, Visualization, Writing – original draft, Writing – review & editing; Louisa S. Schmidt: Conceptualization, Methodology, Investigation, Writing – review & editing; Cornelia Schulze: Conceptualization, Methodology, Writing – review & editing; Michael C. Frank: Conceptualization, Writing – review & editing; Michael Henry Tessler: Conceptualization, Methodology, Formal Analysis, Writing – review & editing.

Correspondence concerning this article should be addressed to Manuel Bohn, Max Planck Institute for Evolutionary Anthropology, Deutscher Platz 6, 04103 Leipzig, Germany. E-mail: manuel_bohn@eva.mpg.de

Abstract

Computational cognitive models have made an important contribution to understanding pragmatic language learning. The focus of this approach has been on explaining adult behavior on a group level. We extend this work to predicting word learning in 3- to 5-year-old children ($N = 60$) on an individual level. In Part 1, we use data from four independent tasks to estimate child-specific sensitivity parameters to three information sources: semantic knowledge, expectations about speaker informativeness, and sensitivity to common ground. In Part 2, we use these parameters to generate participant-specific trial-by-trial predictions about how the same children should behave in a new task that jointly manipulated all three information sources. The model accurately predicted children's behavior in the majority of trials and provided a better explanation of the data compared to two alternative models. As such, this work advances a substantive and testable theory of individual differences in pragmatic word learning.

Keywords: Pragmatics, language development, individual differences, cognitive modeling

Word count: X

Modeling individual differences in children’s information integration during pragmatic word learning

Introduction

A defining feature of human communication is its flexibility. Conventional languages – signed and spoken – allow for expressing a near infinite number of messages in thousands of different ways. In the absence of a shared language, humans can produce and understand novel signals which can rapidly be transformed into structured communication systems (Bohn, Kachel, & Tomasello, 2019; Brentari & Goldin-Meadow, 2017; Goldin-Meadow & Feldman, 1977). The flexibility stems from a powerful social-cognitive infrastructure that underlies human communication (Sperber & Wilson, 2001; Tomasello, 2008). Interlocutors can recruit and integrate a range of different information sources, conventional language being one of them, in order to successfully communicate. For example, to infer what a speaker means by a simple utterance like “she would like the blue one”, the listener has to integrate the semantics of the words with social information available in context such as gestures or gaze and the common ground shared between interlocutors. Such inferences about intended messages are often called pragmatic inferences. They play an important role during everyday language use (H. H. Clark, 1996) and, even more so, during language acquisition (Bohn & Frank, 2019; E. V. Clark, 2009).

Theoretical accounts of language use and learning postulate that pragmatic inferences require information integration. However, they often fail to specify how exactly this happens. This special case mirrors a general issue in psychology and – even more so — in developmental science: a paucity of strong, explicit theories that explain and predict behavior (Muthukrishna & Henrich, 2019). Computational cognitive modeling is often invoked as a way to overcome this issue (Rooij & Baggio, 2021; Simmering, Triesch, Deák, & Spencer, 2010). Cognitive models formalize the computational processes that generate the observed behavior (Rooij, 2022; Ullman & Tenenbaum, 2020). The modelling process forces

researchers to explicitly state their assumptions and intuitions which may result in stronger theories (Guest & Martin, 2021). The field of pragmatic language comprehension has been comparatively active from a computational modelling perspective (Anderson, 2021; Cummins & Ruiter, 2014; Degen, Hawkins, Graf, Kreiss, & Goodman, 2020; Franke & Bergen, 2020; see e.g., Heller, Parisien, & Stevenson, 2016; Tessler & Goodman, 2019; Yoon, Tessler, Goodman, & Frank, 2020). A very productive framework is the Rational Speech Act (RSA) framework, which sees pragmatic language comprehension as a special case of Bayesian social reasoning (Frank & Goodman, 2012; Goodman & Frank, 2016). RSA models are characterized by their recursive structure in which a listener reasons about a cooperative speaker – sensu Grice (1991) – who reasons about a literal listener who interprets words according to their literal semantics.

Most of the time, computational cognitive models – including RSA – are used to explain phenomena in a principled and abstract sense. That is, researchers develop algorithms that reproduce well-known effects from the literature or patterns in already existing data. For example, Frank, Goodman, and Tenenbaum (2009) modeled word learning as inferences about speaker’s intentions and were thereby able to reproduce a range of different effects in early child language (e.g. cross-situational word learning, mutual exclusivity). Such work makes an important contribution to explaining these phenomena in computational terms. However, for a comprehensive theory, models should also be able to *predict* new data (Hofman et al., 2021; Shmueli, 2010; Yarkoni & Westfall, 2017). Recent work has therefore explored how computational models of pragmatic reasoning can be used to make quantitative predictions about *new* data. For example, Bohn, Tessler, Merrick, and Frank (2021) studied young children’s information integration during pragmatic word learning (see also Bohn, Tessler, Merrick, & Frank, 2022). They measured children’s developing sensitivity to three information sources and used an RSA model to generate predictions about situations in which these information sources need to be integrated. Newly collected data aligned closely with what the model predicted. These results offer support for

the theoretical assumptions built into the model, namely that children rationally integrate all available information sources in a stable manner across development.

This line of work critically tests the scope and validity of models of pragmatic reasoning. However, they face yet another fundamental problem. Cognitive models often explain and predict behavior on an aggregated level. The model generates predictions for prototypical agents, which are evaluated in comparison to data that is aggregated across individuals. The assumption is that the “average person” behaves like the prototypical agent. This approach leaves open the question of whether these models are able to predict behavior on an individual level (Estes & Todd Maddox, 2005). In other words, it is unclear if any real individual behaves like the prototypical agent whose cognitive processes are – computationally – simulated. Most likely, there are differences between individuals. For example, Franke and Degen (2016) studied quantity implicatures and found that participant data was best captured by a model that assumes a population in which individuals differ in the depth of their Theory of Mind reasoning. A central question is therefore whether models that accurately predict group-level results can also be used to predict individual differences. In the present study, we address this issue and use a computational cognitive model of pragmatic reasoning to predict individual differences between children.

We build on the work by Bohn et al. (2021) and study how children integrate different information sources in a word learning situation. We focus on how children’s semantic knowledge interacts with their expectations about informative communication and sensitivity to common ground. We formalized this integration process in a model derived from the RSA framework. Importantly, the model was designed to capture individual differences, which we conceptualize as differences between children in sensitivity to the different information sources. In Part 1, we collected data in four tasks from which we estimated child-specific sensitivity parameters. In Part 2, we used these parameters to predict – on a trial-by-trial basis – how the same children should behave in a new task that required information integration. We compared the model predictions to the data and found that, in the majority

of trials, the model accurately predicted children’s behavior.

Part 1: Sensitivity

Methods

Methods, sample size and analyses were pre-registered at: <https://osf.io/pa5x2>. All data, analysis scripts, model code and experimental procedures are publicly available in the following online repository: <https://github.com/manuelbohn/spin-within>.

Participants. We collected complete data for 60 children ($m_{age} = 4.11$, $range_{age}$: 3.06 - 4.93, 30 girls). As per our pre-registration, children who provided valid data for fewer than half of the test trials in any of the three experiments were excluded from the analysis. This was the case for five additional children (two 3-year-olds, three 4-year-olds) due to 1) disinterest in the experiments ($n = 2$), 2) parental interference due to fussiness ($n = 2$), 3) withdrawal from the study after the first testing session ($n = 1$). Children came from an ethnically homogeneous, mid-size German city (~550,000 inhabitants, median income €1,974 per month as of 2020); were mostly monolingual and had mixed socioeconomic backgrounds. The study was approved by an internal ethics committee at the Max Planck Institute for Evolutionary Anthropology. Data was collected between March and July 2021.

Procedure. Children were recruited via a database and participated with their parents via an online conferencing tool. The different tasks were programmed as interactive picture books in JavaScript/HTML and presented on a website. During the video call, participants would enter the website with the different tasks and share their screen. The experimenter guided them through the procedure and told caregivers when to advance to the next task. Children responded by pointing to objects on the screen, which their caregivers would then select for them via mouse click. For the production task, the experimenter shared their screen and presented pictures in a slide show. For the mutual exclusivity, discourse novelty, and combination tasks, pre-recorded sound files were used to address the child.

Figure 1 shows screenshots from the different tasks.

In the *discourse novelty* task, children saw a speaker (cartoon animal) standing between two tables. On one table, there was a novel object (drawn for the purpose of this study) while the other was empty. The speaker sequentially turned to both sides (order counterbalanced) and either commented on the presence or absence of an object (without using any labels). Then, the speaker disappeared and – while the speaker was gone – another novel object appeared on the previously empty table. Next, the speaker re-appeared and requested one of the objects using a novel non-word as the label. We assumed that children would take the novel word to refer to the object that was new to the speaker. Children received 16 trials, each with a new pair of novel objects. The location of the empty table was counterbalanced.

In the *mutual exclusivity* task, children again saw a speaker and two tables. On one table, there was a novel object while on the other there was a (potentially) familiar object. The speaker used a novel non-word to request one of the objects. We assumed that children would take the novel word to refer to the novel object. In line with previous work (Bohn et al., 2021; Grassmann, Schulze, & Tomasello, 2015; Lewis, Cristiano, Lake, Kwan, & Frank, 2020) we assumed this inference would be modulated by children’s lexical knowledge of the familiar object. Children received 16 trials, each with a new pair of novel and familiar objects. The location of the familiar object was counterbalanced. Both the discourse novelty as well as the mutual exclusivity showed good re-test reliability in a previous study and seem well-suited for individual-level measurement (Bohn, Tessler, Kordt, Hausmann, & Frank, 2022).

In the *word production* task, the experimenter showed the child each of the 16 familiar objects from the mutual exclusivity task and asked them to name it. We used a pre-defined list of acceptable labels per object to categorize children’s responses as either correct or incorrect.

In the *word comprehension* task, the child saw four slides with six objects. Four objects per slide were taken from the 16 familiar objects that also featured in the mutual exclusivity and word production tasks. Two objects were unrelated distractors. The experimenter labelled one familiar object after the other and asked the child to point to it.

Data collection was split into two sessions which were scheduled for consecutive days or at most within two weeks of each other. For the majority of children, the second session was scheduled within a week of the first one. On day one, children completed the mutual exclusivity and the discourse novelty tasks. On day two, they completed the combination task followed by the word comprehension and production tasks.

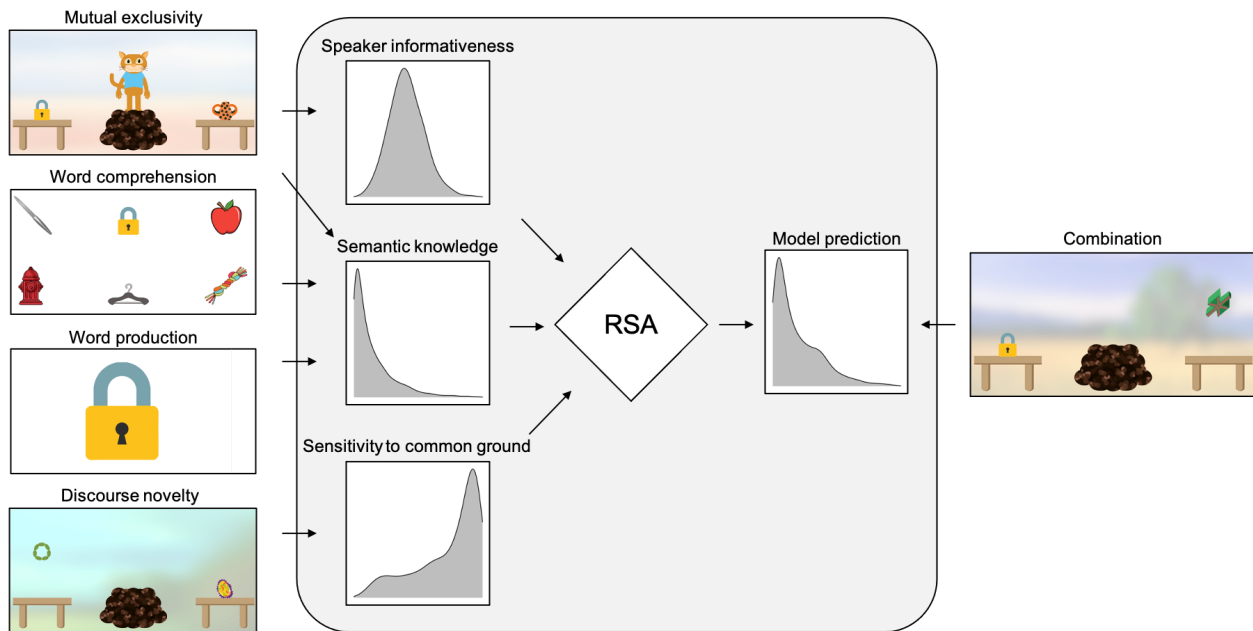


Figure 1. Schematic overview of the study and the model. Pictures on the left show screenshots from the four sensitivity tasks. Arrows indicate which tasks informed which parameter in the model (grey area). Based on the data from the sensitivity tasks, child specific parameter distributions for each information source were estimated. These sources were integrated via an RSA model, which generated predictions for each trial of the combination task. These predictions were then evaluated against new data from the combination task.

Analysis

The focus of the analysis was on estimating person-specific parameters for each information source. Models to estimate parameters were implemented in the probabilistic programming language `webpp1` (Goodman & Stuhlmüller, 2014). The three information sources were: sensitivity to common ground (ρ_i), expectations about speaker informativeness (α_i), and semantic knowledge (θ_{ij}). Figure 1 shows which tasks informed which parameters. All parameters were estimated via hierarchical regression (mixed-effects) models. That is, for each parameter, we estimated an intercept and slope (fixed effects) that best described the developmental trajectory for this parameter based on the available data. Participant-specific parameters values (random effects) were estimated as deviations from the value expected for a participant based on their age. Details about the estimation procedure can be found in the supplementary material. The code to run the models can be found in the associated online repository.

The parameters for semantic knowledge (θ_{ij}) were simultaneously informed by the data from the mutual exclusivity, the comprehension and the production experiments. To leverage the mutual exclusivity data, we adopted the RSA model described in Part 2 to a situation in which both objects (novel and familiar) had equal prior probability (i.e., no common ground information). In the same model, we also estimated the parameter for speaker informativeness (see below). For the comprehension experiment, we simply assumed that the child was able to select the correct word with probability θ_{ij} . If the child did not know the word, we assumed they would select the correct word at a rate expected by chance (1/6). For the production experiment, we assumed that if the child knew the word (a function of θ_{ij}), they produced the word with probability γ . This successful-production-probability γ was the same for all children and was inferred based on the data. This adjustment reflects the finding that children’s receptive vocabulary for nouns tends to be larger than the productive (E. V. Clark & Hecht, 1983; Frank, Braginsky, Yurovsky, & Marchman, 2021). Taken together, for

each child i and familiar object j there were three data points to inform θ : one trial from the mutual exclusivity, one from the comprehension and one from the production experiment.

The parameter representing a child’s expectations about how informative speakers are (α_i), was estimated based on the data from the mutual exclusivity experiment. As mentioned above, this was done jointly with semantic knowledge in a RSA model adopted to a situation with equal prior probability of the two objects (novel and familiar). Thus, for each child, there were 16 data points to inform α .

We estimated children’s sensitivity to common ground (ρ_i) based on the data from the discourse novelty experiment. This was done via simple logistic regression and based on the 12 data points from this task.

Results

Figure 2 visualizes the results for the four sensitivity tasks and the person specific model parameters estimated from the data. In all four tasks, we saw that children performed above chance (not applicable in the case of word production), suggesting that they made the alleged pragmatic inference or knew (some) of the words for the objects involved. With respect to age, performance in raw test scores seemed to increase with age in the three tasks relying on semantic knowledge (mutual exclusivity, word production and word comprehension). Performance in these tasks was also correlated (see supplementary material). For discourse novelty, performance did not increase with age. Most importantly, however, we saw considerable variation between individuals. When focusing on the individual-specific parameter estimates (Figure 2B), we saw that parameters that were estimated based on more data (sensitivity to common ground – 12 trials, and expectations about speaker informativeness – 16 trials) had better defined posterior distributions compared to semantic knowledge (3 trials per object).

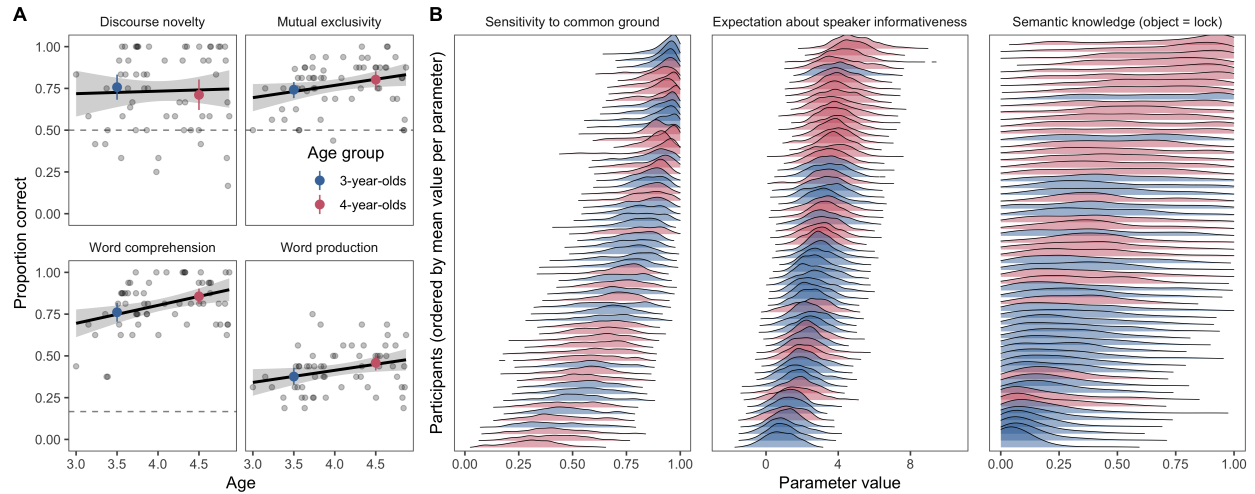


Figure 2. Results for the sensitivity tasks. A: proportion of correct responses in each task by age. Colored dots show the mean proportion of correct responses (with 95% CI) binned by year. Regression lines show fitted generalized linear models with 95% CIs. B: posterior distributions for each parameter (information source) and participant, ordered by mean value, separate for each parameter. Color shows age group.

Discussion

The goal of Part 1 was to estimate person-specific parameters representing each individual's sensitivity to the three information sources. We found that, as a group, children were sensitive to the different information sources. Furthermore, there was substantial variation between individuals in *how* sensitive they were to each information source. These results provided a solid basis for studying information integration in Part 2.

Part 2: Integration

In Part 2, we studied how children integrate the three information sources. We incorporated the parameters estimated in Part 1 in a computational cognitive model of pragmatic reasoning to generate participant-specific predictions about how the three information sources should be integrated. We then compared these predictions to new data

collected with a task in which all three information sources were manipulated. We used Bayesian model comparisons to compare our focal *rational integration model* to alternative models that made different theoretical assumptions about the integration process.

Methods

The study was pre-registered and all data, analysis script and materials are publicly available (see Part 1 for more information).

Participants. Participants were the same as in Part 1.

Procedure. The task was implemented in the same environment as the tasks in Part 1. Each child completed the combination task on the second testing day. The general procedure followed that of the novelty task, however, only one of the objects was unknown while the other was familiar. The combination task had two conditions. In the *congruent condition*, the object that was new to discourse was the novel object. As a consequence, mutual exclusivity and discourse inferences pointed to the same object as the referent of the novel word were aligned. In the *incongruent condition*, the familiar object was new to discourse and thus, the two inferences pointed to different objects. We created matched pairs for the 16 familiar objects and assigned one object of each pair to one of the two conditions. Thus, there were eight trials per condition in the combination task in which each trial was with a different familiar object. We counterbalanced the order of conditions and the side on which the discourse-novel object appeared. Responses were coded from a mutual exclusivity perspective (choosing novel object = 1). All children received the same order of trials. There was the option to terminate the study after 8 trials (two children).

Analysis

We adopted the modelling framework used by Bohn et al. (2021). Our models are situated in the Rational Speech Act (RSA) framework (Frank & Goodman, 2012; Goodman

270 & Frank, 2016). RSA models treat language understanding as a special case of Bayesian
 271 social reasoning. A listener interprets an utterance by assuming it was produced by a
 272 cooperative speaker who has the goal to be informative. Being informative is defined as
 273 producing messages that increase the probability of the listener inferring the speaker’s
 274 intended message. The focal *rational integration* model, including all data-analytic
 275 parameters, is formally defined as:

$$P_{L_1}(r \mid u; \{\rho_i, \alpha_i \theta_{ij}\}) \propto P_{S_1}(u \mid r; \{\alpha_i, \theta_{ij}\}) \cdot P(r \mid \rho_i) \quad (1)$$

276 The model describes a listener (L_1) reasoning about the intended referent of a
 277 speaker’s (S_1) utterance. This reasoning is contextualized by the prior probability of each
 278 referent $P(r \mid \rho_i)$. This prior probability is a function of the common ground ρ shared
 279 between speaker and listener in that interacting around the objects changes the probability
 280 that they will be referred to later.

281 To decide between referents, the listener (L_1) reasons about what a rational speaker
 282 (S_1) would say given an intended referent. This speaker is assumed to compute the
 283 informativity for each available utterance and then choose the most informative one. The
 284 expectation of speaker informativeness may vary and is captured by the parameter α :

$$P_{S_1}(u \mid r; \{\alpha_i \theta_{ij}\}) \propto P_{L_0}(r \mid u; \{\theta_{ij}\})^{\alpha_i} \quad (2)$$

285 The informativity of each utterance is given by imagining which referent a literal
 286 listener (L_0), who interprets words according to their lexicon \mathcal{L} , would infer upon hearing
 287 the utterance. This reasoning depends on what kind of semantic knowledge (word–object
 288 mappings, $\theta_{i,j}$) the speaker thinks the literal listener has. As noted above, for familiar
 289 objects, we take semantic knowledge to be a function of the degree-of-acquisition of the
 290 associated word.

$$P_{L_0}(r \mid u; \{\theta_{ij}\}) \propto \mathcal{L}(u, r \mid \theta_{ij}) \quad (3)$$

This modelling framework allows us to generate predictions for each participant and trial in the combination task based on the participant-specific parameters estimated in Part 1. That is, for each combination of ρ , α , and θ for participant i and familiar object j , the model returns a distribution for the probability with which the child should choose the novel object. We contrasted the predictions made by the *rational integration* model described above to those made by two plausible alternative models which assume that children selectively ignore some of the available information sources (Gagliardi, Feldman, & Lidz, 2017). These models generated predictions based on the same parameters as the *rational integration* model, the only difference lay in how the parameters were used. The *no speaker informativeness* model assumed that the speaker does not communicate in an informative way and therefore focused on the sensitivity to common ground. The *no common ground* model ignores common ground information and focused on the mutual exclusivity inference (speaker informativeness and semantic knowledge instead). A detailed description of all the models along with technical information about parameter estimation can be found in the supplementary material.

We evaluated the model predictions in two steps. First, we replicated the group-level results of Bohn et al. (2021). That is, we compared the three models in how well they predict the data of the combination task when aggregating across individuals. For this, we correlated model predictions and the data (aggregated by trial and age group) and computed pairwise Bayes Factors based on the marginal likelihood of the data given the model.

Second, and most importantly, we evaluated how well the model predicted performance on an *individual* level. For each trial, we converted the (continuous) probability distribution returned by the model into a binary prediction (the structure of the data) by flipping a coin

with the Maximum a posteriori estimate (MAP) of the distribution as its weight¹. For the focal and the two alternative models, we then computed the proportion of trials for which the model predictions matched children’s responses and compared them to a level expected by random guessing using a Bayesian t-test. Finally, for each child, we computed the Bayes Factor in favor of the *rational integration* model and checked for how many children this value was above 1 (log-Bayes Factors > 0). Bayes Factors larger than 1 present evidence in favor of the *rational integration* model. We evaluated the distribution of Bayes Factors following the classification of Lee and Wagenmakers (2014).

Results

On a group-level, the results of the present study replicated those of Bohn et al. (2021). The predictions made by the *rational integration* model were highly correlated with children’s responses in the combination task. The model explained around 74% of the variance in the data and with that more compared to the two alternative models (Figure 3A). Bayes Factors computed via the marginal likelihood of the data (Figure 3B) strongly favored the *rational integration* model in comparison to the *no common ground* ($BF_{10} = 9.1e+53$) as well as the *no speaker informativeness* model ($BF_{10} = 1.2e+44$).

Finally, we turned to the individual-level results. When looking at the proportion of correct predictions (for one run of the coin-flipping procedure), we saw that the *rational integration* model correctly predicted children’s responses in the combination task in 72% of trials, which was well above chance ($BF_{10} = 2.15e+14$) and higher compared to the two alternative models (Figure 4A). Note that the alternative models also predicted children’s responses at a level above chance (*no common ground*: 61%, $BF_{10} = 220251$; *no speaker informativeness*: 60%, $BF_{10} = 55.4$), emphasizing that they constitute plausible alternatives.

¹ Note that this procedure is not deterministic and the results will slightly vary from one execution to the next (see also Figure 4).

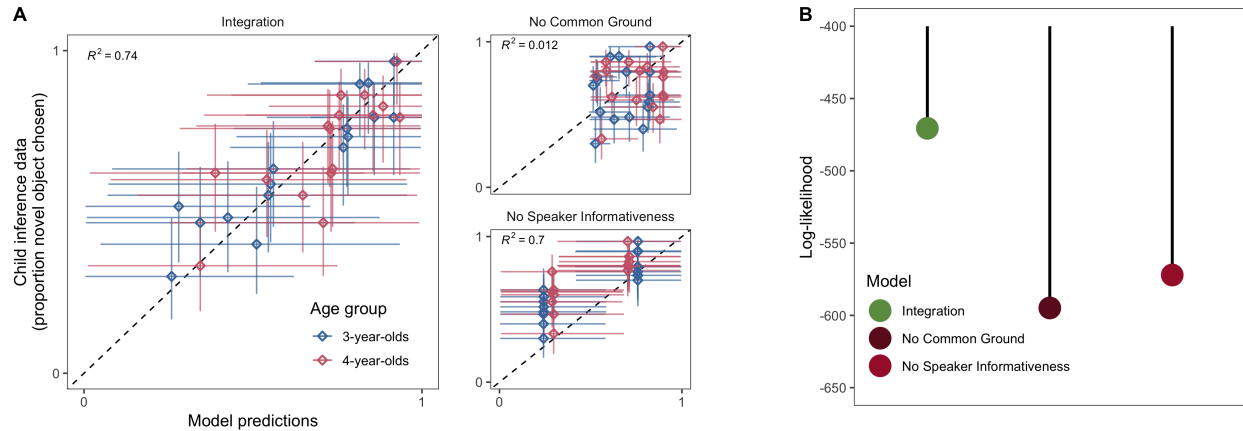


Figure 3. Group-level model comparison. A: Correlation between model predictions and data (aggregated across individuals and binned by year with 95%HDI) for each trial in the combination experiment. B: log-likelihood for each model given the data.

In the supplementary material we also compared models with respect to the situations in which they did or did not correctly predict children's responses.

When directly comparing the models on an individual level, we found that the *rational integration* model provided the best fit for the majority of children. In comparison to the *no common ground* model, 62% of Bayes Factors were larger than 1 and 35% were larger than 10. In comparison to the *no speaker informativeness* model, 68% of Bayes Factors were larger than 1 and 45% were larger than 10 (Figure 4B).

Discussion

The results of Part 2 show that the *rational integration* model accurately predicted children's responses in the combination task. Importantly, this was the case not just on a group level, but also on an individual level. Based on the sensitivity measures obtained for each child in Part 2, the model correctly predicted children's responses in the majority of trials. Furthermore, it was more likely to be correct and provided a better explanation of the data compared to two alternative models that assumed that children selectively ignored some of the information sources.

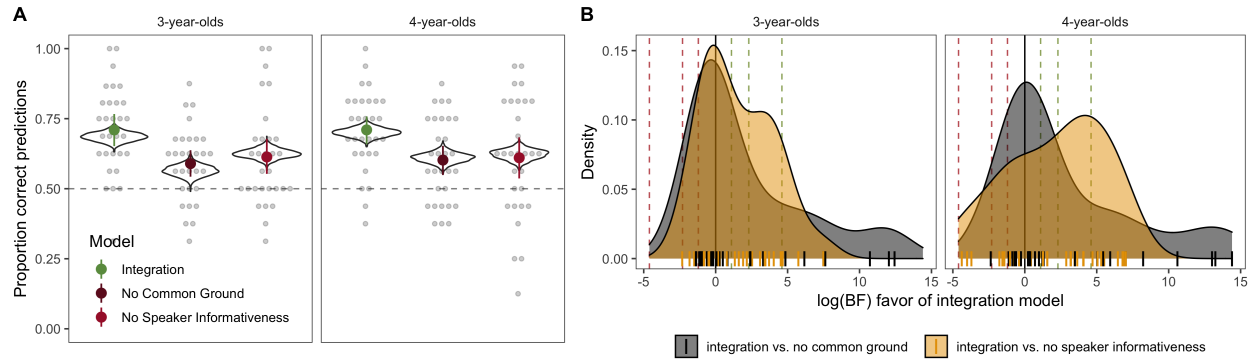


Figure 4. Individual-level model comparison. A: proportion of correct predictions for each model. Solid colored dots show mean with 95%CI for one run of the coin flip procedure. Light dots show aggregated individual data for the same run. Violins show distribution of means for 1000 runs of the procedure. B: distribution of log-Bayes Factors for each individual. Dashed lines show Bayes Factor thresholds of 3, 10 and 100.

General discussion

In this study, we used a computational cognitive model of pragmatic reasoning to make out-of-sample predictions about children’s behavior on a trial-by-trial basis. In Part 1, we used data from four tasks to estimate child-specific sensitivity parameters capturing their semantic knowledge, expectations about speaker informativeness and sensitivity to common ground. In Part 2, we used these parameters to predict how the same children should behave in a new task in which all three information sources were jointly manipulated. We found strong support for our focal *rational integration* model in that this model accurately predicted children’s responses in the majority of trials and provided a better fit compared to two alternative models. Taken together, this work provides a strong test of the theoretical assumptions built into the model.

The *rational integration* model was built around three main theoretical assumptions. First, it assumes that children integrate all available information sources. The model comparison, in which we compared the focal model to two models that selectively ignore

some of the information sources, strongly supported this assumption. For the majority of individuals – as well as on a group level – this model provided the best fit. However, for some individuals, one of the alternative models provided a better fit. Finding out why this is the case and what characterizes these individuals would be an interesting avenue for future research. Second, the model assumes that the integration process does not change with age. We did not probe this assumption in the present study because, in order to do so on an individual level, it would require longitudinal data. We think this would be anotehr interesting extension of our work. Finally, the model assumes that children differ in their sensitivity to the different information sources but *not* in the way they integrate information. Even though the model built around this assumption predicted the data well, it would also be interesting to explore structural differences between individuals (e.g. Franke & Degen, 2016).

Even though the model explains *and* predicts the data well, it is first and foremost a computational model, meaning that we should be careful with granting the processes and parameters in it too much psychological realism. Nevertheless, we think that when studying individual differences, the model parameters can be interpreted in a psychologically *more* plausible way compared to raw performance scores that are otherwise used to describe individuals in the field. They are estimated taking into account the structure of the task and the different processes that are involved in it. This allows for informing a parameter based on data from multiple tasks, as, for example, semantic knowledge was estimated based on the mutual exclusivity, comprehension and production tasks. Support for such an approach comes from a recent study that used an RSA-type model to link performance in different pragmatic reasoning tasks in order to jointly estimate a single parameter capturing children’s pragmatic abilities (Bohn et al., 2022). This parameter was correlated with measures of executive functions which have been repeatedly suggested to play an important role in pragmatic reasoning (e.g., Matthews, Biney, & Abbot-Smith, 2018). Taken together we think that computational modelling can make an important contribution to studying individual differences on a process level.

Our study is limited in terms of generalizability because we tested one sample of children growing up in a western, affluent setting. However, the modelling approach put forward here provides an interesting way of studying and theorizing about cross-cultural differences. Following Bohn and Frank (2019), our *prima facie* assumption is that children from different cultural settings might differ in terms of their sensitivity to different information sources – just like individuals differ within cultural settings – but the way that information is integrated is the same across cultures. This prediction could be tested by comparing alternative models that make different assumptions about the integration process.

Taken together, we have shown that children’s pragmatic word learning can be predicted on a trial-by-trial basis by a computational cognitive model. Taken together with previous work that focused on aggregated developmental trajectories (Bohn et al., 2021), this suggests that the same computational processes can be used to predict group- and individual-level data. Furthermore, we have offered a substantive and testable theory of how individuals might differ from one another in the alleged processes.

References

- Anderson, C. J. (2021). Tell me everything you know: A conversation update system for the rational speech acts framework. *Proceedings of the society for computation in linguistics 2021*, 244–253.
- Bohn, M., & Frank, M. C. (2019). The pervasive role of pragmatics in early language. *Annual Review of Developmental Psychology*, 1(1), 223–249.
- Bohn, M., Kachel, G., & Tomasello, M. (2019). Young children spontaneously recreate core properties of language in a new modality. *Proceedings of the National Academy of Sciences*, 116(51), 26072–26077.
- Bohn, M., Tessler, M. H., Kordt, C., Hausmann, T., & Frank, M. C. (2022). An individual differences perspective on the development of pragmatic abilities in the preschool years.
- Bohn, M., Tessler, M. H., Merrick, M., & Frank, M. C. (2021). How young children integrate information sources to infer the meaning of words. *Nature Human Behaviour*, 5(8), 1046–1054.
- Bohn, M., Tessler, M. H., Merrick, M., & Frank, M. C. (2022). Predicting pragmatic cue integration in adults’ and children’s inferences about novel word meanings. *Journal of Experimental Psychology: General*.
- Brentari, D., & Goldin-Meadow, S. (2017). Language emergence. *Annual Review of Linguistics*, 3, 363–388.
- Clark, E. V. (2009). *First language acquisition*. Cambridge: Cambridge University Press.
- Clark, E. V., & Hecht, B. F. (1983). Comprehension, production, and language acquisition. *Annual Review of Psychology*, 34(1), 325–349.
- Clark, H. H. (1996). *Using language*. Cambridge: Cambridge University Press.
- Cummins, C., & Ruiter, J. P. de. (2014). Computational approaches to the pragmatics problem. *Language and Linguistics Compass*, 8(4), 133–143.

- 434 Degen, J., Hawkins, R. D., Graf, C., Kreiss, E., & Goodman, N. D. (2020). When
435 redundancy is useful: A bayesian approach to “overinformative” referring
436 expressions. *Psychological Review*, 127(4), 591.
- 437 Estes, W. K., & Todd Maddox, W. (2005). Risks of drawing inferences about
438 cognitive processes from model fits to individual versus average performance.
439 *Psychonomic Bulletin & Review*, 12(3), 403–408.
- 440 Frank, M. C., Braginsky, M., Yurovsky, D., & Marchman, V. A. (2021). *Variability
441 and consistency in early language learning: The wordbank project*. MIT Press.
- 442 Frank, M. C., & Goodman, N. D. (2012). Predicting pragmatic reasoning in language
443 games. *Science*, 336(6084), 998–998.
- 444 Frank, M. C., Goodman, N. D., & Tenenbaum, J. B. (2009). Using speakers’
445 referential intentions to model early cross-situational word learning. *Psychological
446 Science*, 20(5), 578–585.
- 447 Franke, M., & Bergen, L. (2020). Theory-driven statistical modeling for semantics
448 and pragmatics: A case study on grammatically generated implicature readings.
449 *Language*, 96(2), e77–e96.
- 450 Franke, M., & Degen, J. (2016). Reasoning in reference games: Individual-vs.
451 Population-level probabilistic modeling. *PloS One*, 11(5), e0154854.
- 452 Gagliardi, A., Feldman, N. H., & Lidz, J. (2017). Modeling statistical insensitivity:
453 Sources of suboptimal behavior. *Cognitive Science*, 41(1), 188–217.
- 454 Goldin-Meadow, S., & Feldman, H. (1977). The development of language-like
455 communication without a language model. *Science*, 197(4301), 401–403.
- 456 Goodman, N. D., & Frank, M. C. (2016). Pragmatic language interpretation as
457 probabilistic inference. *Trends in Cognitive Sciences*, 20(11), 818–829.
- 458 Goodman, N. D., & Stuhlmüller, A. (2014). *The design and implementation of
459 probabilistic programming languages*. <http://dippl.org>.
- 460 Grassmann, S., Schulze, C., & Tomasello, M. (2015). Children’s level of word

knowledge predicts their exclusion of familiar objects as referents of novel words.

Frontiers in Psychology, 6, 1200.

Grice, H. P. (1991). *Studies in the way of words*. Cambridge, MA: Harvard University Press.

Guest, O., & Martin, A. E. (2021). How computational modeling can force theory building in psychological science. *Perspectives on Psychological Science*, 16(4), 789–802.

Heller, D., Parisien, C., & Stevenson, S. (2016). Perspective-taking behavior as the probabilistic weighing of multiple domains. *Cognition*, 149, 104–120.

Hofman, J. M., Watts, D. J., Athey, S., Garip, F., Griffiths, T. L., Kleinberg, J., ... others. (2021). Integrating explanation and prediction in computational social science. *Nature*, 595(7866), 181–188.

Lee, M. D., & Wagenmakers, E.-J. (2014). *Bayesian cognitive modeling: A practical course*. Cambridge University Press.

Lewis, M., Cristiano, V., Lake, B. M., Kwan, T., & Frank, M. C. (2020). The role of developmental change and linguistic experience in the mutual exclusivity effect. *Cognition*, 198, 104191.

Matthews, D., Biney, H., & Abbot-Smith, K. (2018). Individual differences in children’s pragmatic ability: A review of associations with formal language, social cognition, and executive functions. *Language Learning and Development*, 14(3), 186–223.

Muthukrishna, M., & Henrich, J. (2019). A problem in theory. *Nature Human Behaviour*, 3(3), 221–229.

Rooij, I. van. (2022). Psychological models and their distractors. *Nature Reviews Psychology*, 1–2.

Rooij, I. van, & Baggio, G. (2021). Theory before the test: How to build high-verisimilitude explanatory theories in psychological science. *Perspectives on*

488 *Psychological Science*, 16(4), 682–697.

489 Shmueli, G. (2010). To explain or to predict? *Statistical Science*, 25(3), 289–310.

490 Simmering, V. R., Triesch, J., Deák, G. O., & Spencer, J. P. (2010). A dialogue on
491 the role of computational modeling in developmental science. *Child Development*
492 *Perspectives*, 4(2), 152–158.

493 Sperber, D., & Wilson, D. (2001). *Relevance: Communication and cognition* (2nd
494 ed.). Cambridge, MA: Blackwell Publishers.

495 Tessler, M. H., & Goodman, N. D. (2019). The language of generalization.
496 *Psychological Review*, 126(3), 395.

497 Tomasello, M. (2008). *Origins of human communication*. Cambridge, MA: MIT Press.

498 Ullman, T. D., & Tenenbaum, J. B. (2020). Bayesian models of conceptual
499 development: Learning as building models of the world. *Annual Review of*
500 *Developmental Psychology*, 2, 533–558.

501 Yarkoni, T., & Westfall, J. (2017). Choosing prediction over explanation in
502 psychology: Lessons from machine learning. *Perspectives on Psychological Science*,
503 12(6), 1100–1122.

504 Yoon, E. J., Tessler, M. H., Goodman, N. D., & Frank, M. C. (2020). Polite speech
505 emerges from competing social goals. *Open Mind*, 4, 71–87.