

# Learning about others: Modeling social inference through ambiguity resolution

Asya Achimova                      Gregory Scontras  
asya.achimova@uni-tuebingen.de      g.scontras@uci.edu

Christian Stegemann-Philipps  
christian.stegemann@uni-tuebingen.de

Johannes Lohmann  
johannes.lohmann@uni-tuebingen.de

Martin V. Butz  
martin.butz@uni-tuebingen.de

August 29, 2020

## Abstract

Bayesian accounts of social cognition successfully model the human ability to infer goals and intentions of others on the basis of their behavior. In this paper, we extend this paradigm to the analysis of ambiguity resolutions during brief communicative exchanges. In particular, we model how observed ambiguity resolutions allow speakers to infer aspects of the inner states of their conversation partners, such as their knowledge, beliefs, intentions, or preferences. Moreover, we model how speakers may choose ambiguous utterances over unambiguous ones in an epistemic manner, anticipating social knowledge gain from expected ambiguity resolutions. In a reference game experimental setup, we observed that participants were able to infer listeners' preferences when analyzing their choice of objects given referential ambiguity. Moreover, a subset of speakers was able to strategically choose ambiguous over unambiguous utterances in an epistemic manner, although a different group preferred unambiguous utterances. Both types of inference are well-approximated by a modified rational speech-act model, which focuses on inferring or expecting to infer listeners' preference priors. We conclude that observations of ambiguity resolutions can reveal socially-relevant information about conversation partners, where the involved inference processes are well-approximated by Bayesian inference.

**Keywords:** ambiguity; pragmatics; information gain; event-predictive cognition; Rational Speech Act models; social intelligence

# 1 Ambiguity in natural language

Ambiguity is ubiquitous during conversations: speakers rely on aspects of context and extra-linguistic reasoning to enrich the linguistic signal and deliver their intended meanings. Given that ambiguity can hinder the efficient transfer of information between conversation partners, it is not surprising that linguists have treated the possibility for ambiguity as a bug in the communication system (Chomsky, 2002) and suggested that ambiguity should generally be avoided (Grice, 1975). When we look back at the study of ambiguity, we actually find that the acknowledgment of ambiguity avoidance as a communicative strategy is much older than its pronouncement by modern linguists. Ambiguity in the language system has been on the radars of philosophers starting with Aristotle (Sennet, 2016). Greek and Latin rhetoricians believed that a skillfully-written text allows for a perfectly accurate and lossless transmission of meaning to the reader or listener (Ossa-Richardson, 2019); such a text avoids ambiguities.

The attitude toward ambiguity has at times been quite different between disciplines, in part because the term itself can refer to multiple phenomena. In organizational communication—communication that aids production—ambiguity aligns closely with underspecification: an utterance is ambiguous when it does not provide every detail about the intended meaning, leaving room for the listener to interpret it. This freedom is believed to be important in communication between managers and their employees when managers set future goals that should stimulate rather than limit creativity (Mohr, 1983). Ambiguity allows for the expression of ideas that are broadly true of a large group, as in company slogans or vision statements (Carmon, 2013). In this case, the language needs to be general enough to allow every member of the team to relate the stated general goals to themselves. Similarly, ambiguous descriptions allow speakers to avoid conflict (Pascale & Athos, 1981): interlocutors choose utterances that allow a range of interpretations and do not enforce a particular viewpoint.

For linguistic research, a word is ambiguous if it can have two separate meanings even in the absence of context, simply as a linguistic sign. In that sense, the word “bat” is ambiguous between a winged mammal and a sporting implement. In the case of referential ambiguity, an ambiguous utterance may apply to several possible referents in a scene. Here we use the term ‘referential’ in the sense of Frege (1892), distinguishing the reference of a word—an object/property in the world—and its meaning. For example, a pronoun can be referentially ambiguous if there are multiple potential antecedents in the context. It is this latter type of ambiguity that we focus on in this paper, although the lessons we learn are likely to apply to the broader range of ambiguity phenomena.

In spite of the early advice in linguistics to avoid ambiguity, more recent research has begun to take notice of the efficiency ambiguity affords us: by relying on context to fill in missing information, we can reuse lightweight bits of language rather than fully specifying the intended message (Levinson, 2000; Piantadosi, Tily, & Gibson, 2012; Wasow, 2015). Viewed in this way, ambiguity serves

as a feature—not a bug—of an efficient communication system. This reasoning accords with years of psycholinguistic research documenting that speakers readily produce ambiguous utterances (see Ferreira, 2008, for an overview). Along related lines, Wasow (2015) reviews a large body of evidence and concludes that ambiguity is rarely avoided, even in situations where its avoidance would seemingly be communicatively appropriate. This observation stands at odds with the Gricean maxim to avoid ambiguity (Grice, 1975). However, even Grice recognized a case of strategic ambiguity where it could be the intention of the speaker to communicate more than one possible interpretation afforded by an ambiguous utterance. In such cases, recognition of the ambiguity serves as the communicative purpose of the utterance. Wasow (2015), on the other hand, reviews several cases where ambiguous production serves no obvious communicative purpose.

In the current work, we focus on the effects of resolving – and anticipating the resolution of – ambiguous utterances, identifying an additional benefit to ambiguous language: the social information we gain from observing how listeners resolve ambiguity. We show in which way language users may learn about each other’s inner mental state (in Bayesian terms, their priors) when observing how ambiguity is resolved. When utterances leave room for interpretation, listeners must draw on their opinions, knowledge, beliefs, and preferences to fill in the gaps. On the other hand, when observing how listeners fill in those gaps, speakers thus learn about the opinions, beliefs, and preferences of their conversation partner. Over the course of two studies, we first demonstrate that people are indeed able to infer hidden beliefs of their conversation partners on the basis of observed ambiguity phenomena in an approximately Bayesian manner; second, we show that some speakers can actively create situations of uncertainty, anticipating the epistemic value when observing the resolution of ambiguity. We thus show that humans can interpret the behavior of other people as driven by their motives, intentions, or personal characteristics—an idea that goes back conceptually to the attribution theory (Jones & Davis, 1965; Kelley, 1967; Kelley & Stahelski, 1970).

To explain the behavior we observe in our experiments, we advance a computational cognitive model of the involved probabilistic inference processes. While our model pursues Bayesian inference, or “psychological reasoning” (Frank & Goodman, 2012), we do not focus on the inference of the actor’s knowledge, that is, on *learning from others* (Shafto, Goodman, & Frank, 2012). Rather, we focus on *learning about others*, that is, learning about listeners’ preferences when observing how they resolve ambiguity. We explore interpretive choices and the potential strategic, socially epistemic usage of ambiguous utterances in anticipation of actors’ responses. Interestingly, the modeling results reveal good interpretive abilities but also strong individual differences when the task is to choose (ambiguous) utterances strategically for gaining social knowledge.

In what follows, we first provide computational background on referential ambiguity resolution (Section 2). In Section 3, we develop computational models that are able to infer the preferences of an agent that led her to a particular choice of objects, as well as a model that predicts which utterances are most useful to create the

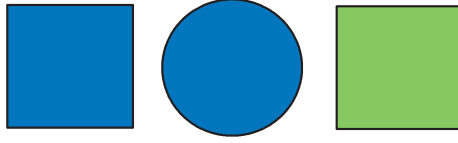


Figure 1: A simple reference game scenario from Frank and Goodman (2012). In the game, speakers are confronted with a collection of objects, which determine the current scenario  $S$ , where  $S = \{\text{solid blue square}, \text{solid blue circle}, \text{solid green square}\}$  in the depicted example. A speaker may choose a single-word utterance  $u$  to signal one of the objects  $s \in S$  to a listener. In the shown scenario, the following set of utterances is available:  $U = \{\text{“solid”}, \text{“blue”}, \text{“green”}, \text{“square”}, \text{“circle”}\}$ .

possibility of learning about the preferences of the conversation partner. Sections 4 and 5 give the results of the behavioral experiments, as well as an evaluation of modeling performance. Section 6 concludes that participants were indeed able to use observable behavior of others to infer their prior beliefs, and hypothesizes why the ability to intentionally create epistemic situations can be found only in part of the population.

## 2 Probabilistic modeling of ambiguity resolution

To see the potential epistemic benefit of ambiguous language, take the scenario in Figure 1. Suppose a speaker produces the single-word utterance to signal one of the objects to a listener. Upon hearing “blue,” the listener faces referential ambiguity: the speaker could mean the blue square or the blue circle. Suppose further that, upon hearing “blue” in this scenario, the listener selects the blue circle. In observing this choice, the speaker learns something about the private thoughts of the listener: what made her select the blue circle instead of the blue square? Perhaps the circle is more salient to the listener, or the listener has a preference for circles, or the listener may believe that the speaker has a preference for circles; there may even be mutual agreement that circles are to be preferred when possible. Importantly, by observing how the listener resolves the ambiguity in reference, the speaker can learn something about the private thoughts of the listener.

However, accessing this added information requires the speaker to reason pragmatically about the pragmatic reasoning of the listener—a higher-order pragmatic reasoning. In order to select a referent, the listener must first interpret the utterance. We follow Frank and Goodman (2012) in treating this interpretation process as active pragmatic, probabilistic reasoning: the listener interprets an utterance by reasoning about the process that generated it, namely the speaker, who selects an utterance by reasoning about how a listener would interpret it. Frank and Goodman model this recursive social reasoning between speakers and listeners, introducing the Rational Speech Act (RSA) modeling framework (Frank & Goodman, 2012;

Franke & Jäger, 2016; Goodman & Frank, 2016).

In this section, we first review Frank and Goodman’s original formulation of RSA. We then explore recent work on social reasoning from a Bayesian perspective.

## 2.1 Original RSA Formalization

Frank and Goodman’s RSA model of the reference game in Figure 1 formalizes a state space, or scenario,  $S$ , as a particular set of objects. The model unfolds computations over the corresponding utterance space  $U$ , which consists of the set of possible utterances. At the base of the reasoning process, there is a hypothetical, naïve literal listener  $L_0$ , who hears an utterance  $u \in U$  and attempts to infer the object  $s \in S$  that  $u$  is meant to reference.  $L_0$  performs this inference by conditioning on the literal semantics of  $u$ ,  $\llbracket u \rrbracket(s)$ , which returns *true* (i.e., 1) for those objects that possess the uttered feature and *false* (i.e., 0), otherwise. As a result, object choice probabilities for the literal listener can be computed by:

$$P_{L_0}(s \mid u) \propto \llbracket u \rrbracket(s), \quad (1)$$

essentially returning a uniform distribution over those objects in  $S$  that contain the uttered feature  $u$ .<sup>1</sup>

One layer up in the reasoning chain, the speaker  $S_1$  observes the scenario  $S$  and is assumed to have the intention to refer to a particular object  $s \in S$ .  $S_1$  chooses an utterance  $u$  on the basis of its expected utility for signaling  $s$  in the scenario  $S$  to  $L_0$ ,  $U_{S_1}(u; s)$ :<sup>2</sup>

$$U_{S_1}(u; s) = \log(P_{L_0}(s \mid u)). \quad (2)$$

Depending on a “greediness” factor  $\alpha$ , the speaker chooses a particular utterance  $u$  with a probability that is exponentially proportional to the utility estimate:

$$P_{S_1}(u \mid s) \propto \exp(\alpha \cdot U_{S_1}(u; s)). \quad (3)$$

At the top layer of the vanilla RSA model, the *pragmatic* listener  $L_1$  infers posteriors over  $s$  on the basis of some observed utterance  $u$ . However, unlike  $L_0$ ,  $L_1$  updates beliefs about the world by reasoning about the process that *generated*  $u$ , namely the utterance choice of speaker  $S_1$ . In other words,  $L_1$  reasons about which object  $s$  would have been most likely to lead  $S_1$  to utter  $u$  given the scenario  $S$ :

$$P_{L_1}(s \mid u) \propto P_{S_1}(u \mid s) \cdot P(s). \quad (4)$$

Frank and Goodman (2012) tested the predictions of their RSA model against behavioral data from reference games as in Figure 1. To model production behavior (that is, which utterance should be chosen to communicate a given object), the

<sup>1</sup>Note that the context  $S$  is typically not made explicit, but rather treated implicitly in the specification of the model.

<sup>2</sup>The original model in Frank and Goodman (2012) also includes a term for the utterance cost,  $C(u)$ . We ignore the term here since we assume uniform cost over all utterances.

authors used the probability distributions from  $S_1$ . To model interpretation behavior (i.e., which object the speaker is trying to communicate on the basis of their utterance), the authors generated predictions from  $L_1$ . Frank and Goodman found strong correlations between model predictions and behavioral data in both cases, confirming the validity of their model of pragmatic reasoning in reference games (see also Qing & Franke, 2015, for a fuller exploration of the modeling choices).

## 2.2 Bayesian Theory of Mind

While the RSA approach is aimed at inferring the meaning of an utterance by reasoning about the process that generated it, the inference process can also apply to aspects of the speaker’s knowledge, or, in other words, her priors. This type of inference is crucial for social interactions, since it helps the conversation partners build more accurate anticipatory models of each others’ behavior. The anticipatory nature of human cognition has been registered in a number of cognitive domains (Butz, 2016; Belardinelli, Stepper, & Butz, 2016; Belardinelli, Lohmann, Farnè, & Butz, 2018; Friston et al., 2015; Hayhoe, Shrivastava, Mruczek, & Pelz, 2003; Lohmann, Belardinelli, & Butz, 2019) and is viewed by some researchers as a core property of the human mind (Butz, 2008; Clark, 2016).

Learning about the goals, beliefs, and preferences (i.e., the priors) of other agents depends on the ability of humans to infer hidden states by observing behavioral choices, thereby engaging in Theory-of-Mind reasoning. Our idea regarding the potential added utility of ambiguous language—that conversation partners can learn about each other as they observe ambiguity-resolution behavior—requires just this sort of reasoning. The large and growing literature on Bayesian Theory of Mind (Baker, Saxe, & Tenenbaum, 2009; Baker, Jara-Ettinger, Saxe, & Tenenbaum, 2017) thus serves as inspiration for our own extended RSA model.

The ability to infer others’ preferences upon observing their behavioral choices develops early in childhood; Lucas et al. (2014), drawing from work in psychology and economics, formalize this process with a Mixed Multinomial Logit model, which is driven by the assumption that, in making choices, agents maximize subjective utility. Jara-Ettinger, Gweon, Schulz, and Tenenbaum (2016) review a large body of experimental work with children and adults, and propose a naïve-utility-calculus model of so-called ‘commonsense psychology’—the ability of people to infer the hidden causes of others’ behavior by treating them as utility-maximizing agents.

The process of mentalizing has been modeled within the Bayesian framework in a number of papers that look at the interpretation of the rational behavior of agents. Shafto et al. (2012) developed a Bayesian model of learning that formalizes the process of inferring others’ knowledge about the world based on their actions and goals. Evans, Stuhlmüller, and Goodman (2016) model the inference of prior preferences when agents are not fully consistent or have restricted knowledge about the choice options. The authors propose a model that can maintain uncertainty over the inferred beliefs. Baker et al. (2017) develop a computational model of Bayesian

Theory-of-Mind reasoning; the authors demonstrate that inference crucially relies on joint reasoning about beliefs, desires, and percepts, since simpler models that consider only a subset of these components are less accurate at predicting human judgments. In a lexicon-learning paradigm, Woensdregt, Kirby, Cummins, and Smith (2016) model how Bayesian inference about the beliefs of speakers can co-develop with the process of determining the likely referents of lexical items. The authors treat word learning as a process of inferring the intended meaning from hearing a word and observing the context in which that word was used. In a series of simulations, the model learns a lexicon and a perspective jointly; where perspective is defined as the distance-based salience of an object for the speaker. Woensdregt et al. further discuss what implications the ability to mentalize carries for reducing referential uncertainty and successful vocabulary learning.

Within the RSA framework, Degen, Tessler, and Goodman (2015) show how inferring the hidden mental states of conversation partners can take the form of inferring priors. The authors explored how listeners infer the speaker’s prior over world states when confronted with an utterance that conflicts with their own prior. If a listener hears that marbles were thrown into a pool and ‘some marbles sank’, she infers that the likely prior over the world states (i.e., that marbles usually sink in water) is not shared by the speaker, thereby reasoning that the world must be ‘wonky’. In their model, listeners actively reason about a wonkiness parameter; in a wonky world, the listener switches from their informative prior to a uniform prior over the world states. As another instance of Theory-of-Mind reasoning within the RSA framework, Yoon, Frank, Tessler, and Goodman (2018) model the comprehension of polite language that features so-called ‘white lies’ (e.g., telling someone that their terrible cookies are ‘okay’). In the model, listeners jointly infer the state of the world (e.g., how good the cookies are) and the goals of the speaker (i.e., how much they prioritize politeness vs. informativity). In work on metaphorical language use, Kao, Bergen, and Goodman (2014) also consider affective goals that justify non-literal interpretation of their utterances. The authors further model the comprehension of non-literal language as a joint inference of the meaning of an utterance and an affective state of the speaker.

### 3 Our model of social inference

Our work contributes to the broad literature on Bayesian Theory of Mind and more specifically to the modeling of prior inference in communicative settings within the RSA framework. Unlike previous work that has focused on modeling the listener’s inference over priors of the pragmatic speaker  $S_1$ , here we focus on actively inferring priors of the listener upon observing – or expecting to observe – their ambiguity resolution behavior. Accordingly, we model the inference process of a higher-order pragmatic speaker upon having observed a listener’s ambiguity-resolution behavior as well as an active utterance inference process, which aims at maximizing expected information gain.

Our model builds on the vanilla version of RSA, modifying the listener’s state prior  $P(s)$  and enhancing the reasoning process towards a social component, yielding a pragmatic social-inference RSA model. By changing  $P(s)$  to a non-uniform distribution, we model prior beliefs about which object the speaker is more likely to refer to, or, when viewed from a more self-centered perspective, which feature preferences  $f$  the listener may have. For example, the listener may like blue things, such that she may be more likely to choose the blue square instead of the green one when hearing the utterance “square” in the scenario shown in Figure 1. As a result, when a pragmatic speaker produces utterance  $u$  and observes the listener’s referent choice  $s$ , the speaker may infer posteriors over possible feature preferences, attempting to explain the observed object choice in this way.

We use  $L_0$  and  $S_1$  from the vanilla model, but we now parameterize  $L_1$ ’s state prior such that it operates given a feature preference  $f$ :

$$P_{L_1}(s | u, f) \propto P_{S_1}(u | s) \cdot P(s | f). \quad (5)$$

We then model a pragmatic speaker  $S_2$ , who updates beliefs about  $L_1$ ’s preferences,  $P(f)$ .  $S_2$  observes  $L_1$ ’s choice of  $s$  given the produced utterance  $u$  and then reasons about the likely feature preference  $f$  that  $L_1$  used to make the observed choice:

$$P_{S_2}(f | u, s) \propto P_{L_1}(s | u, f) \cdot P(f). \quad (6)$$

We also model the reasoning process by which a speaker may select the best utterance to learn about the preferences of the listener, essentially striving to maximize expected information gain concerning the listener’s feature preferences. Starting with no knowledge of the listener’s preferences,  $S_2$  can be assumed to expect a uniform (i.e., flat) feature preference prior  $P(f)$ . The more the speaker’s posterior beliefs about the preferences,  $P_{S_2}(f | u, s)$ , deviate from the uniform prior, the more the speaker will have learned about the listener’s preferences. We can thus model this reasoning in light of expected information gain, which can be equated with the attempt to maximize the KL (Kullback-Leibler) divergence between the speaker’s flat prior and the expected posterior over the listener’s feature preferences  $f$ , integrating over all hypothetically possible object choices  $s \in S$ :

$$P_{S_2}(u) \propto \sum_{s: \llbracket u \rrbracket(s)=1} P_{L_1}(s | u, f) \exp(\lambda \cdot \text{KL}(P(f) || P_{S_2}(f | u, s))), \quad (7)$$

where the factor  $\lambda$  scales the importance of the KL divergence term.

We evaluate two versions of the model. The full model assumes the recursive reasoning process integrating the full RSA formalism. The full model thus assumes that feature-preference inference not only considers the current object choices possible, but also differentiates the choice options further with respect to their pragmatic plausibility. For example, the full model captures the fact that when a speaker utters “blue” in the object situation depicted in the example shown in Figure 1 and has the intention to refer to one particular object, she is more likely to refer to the



blue square than to the blue circle, because in the latter case the utterance choice “circle” would have been unambiguous and thus a better utterance choice.

Recently, it has been shown that even in the original, simpler reference games, fewer layers of reasoning often perform equally well or better than more complex RSA-based models (Sikos, Venhuizen, Drenhaus, & Crocker, 2019). Accordingly, our simple model removes the reasoning about alternative utterances and allows the pragmatic speaker to directly tap into the (expected) interpretation of  $L_0$ , augmenting the literal listener’s choice likelihoods with the feature-preference-dependent object prior  $P(s | f)$ :

$$P_{L_0\text{-simp}}(s | u, f) \propto \llbracket u \rrbracket(s) \cdot P(s | f). \quad (8)$$

The pragmatic speaker  $S_{1\text{-simp}}$  then reasons directly about the modified literal listener  $L_{0\text{-simp}}$ :

$$P_{S_{1\text{-simp}}}(f | u, s) \propto P_{L_{0\text{-simp}}}(s | u, f) \cdot P(f). \quad (9)$$

As a result, the simple model ignores any indirect pragmatic reasoning considerations about which object the speaker may refer to given an utterance and a particular object constellation. The model simply assumes that all objects may be chosen that match the utterance semantics, modifying these choice options depending on the state prior that incorporates feature preferences. The corresponding utterance-selection model simplifies the reasoning process accordingly:

$$P_{S_{1\text{-simp}}}(u) \propto \sum_{s: \llbracket u \rrbracket(s)=1} P_{L_0}(s|u, f) \exp(\lambda \cdot \text{KL}(P(f) || P_{S_{1\text{-simp}}}(f | u, s))). \quad (10)$$

In evaluating our model predictions below, we compare the modeling performance of the full model with that of the simple model.

## 4 Experiment 1: Inferring preferences


Our first task is to check the inferences of the pragmatic speaker having observed that a listener selects some object  $s$  in response to an utterance  $u$ . Is it possible to draw inferences about the most likely preferences the listener had when making her choice? Can this inference process be modeled by our RSA model—that is, by recursive Bayesian inference?

### 4.1 Participants


We recruited 90 participants with US IP addresses through Amazon.com’s Mechanical Turk crowdsourcing service. Participants were compensated for their participation. On the basis of a post-test demographics questionnaire, we identified 82 participants as native speakers of English; their data were included in the analyses reported below. We obtained a confirmation from all the participants that they agreed to participate in the study.

## 4.2 Design and methods







We presented participants with a series of reference game scenarios modeled after Figure 1 from Frank and Goodman (2012). Each scenario featured two people and three objects. One of the people served as the speaker, and the other served as the listener. The speaker asks the listener to choose one of the objects, but in doing so she is allowed to mention only one of the features of the target object. Participants were told that the listener might have a preference for certain object features, and participants were tasked with inferring those preferences after observing the speaker’s utterance and listener’s object choice. A sample trial is shown in Figure 2.

Progress: 

Suppose Maria wants to signal an object in the following scene to Samantha.  
Maria says "**red**" and Samantha chooses the outlined object:



Based on this choice, do you think Samantha has a preference for certain types of objects?

	very unlikely	very likely		very unlikely	very likely
solid things			clouds		
striped things			circles		
polka-dotted things			squares		




Figure 2: A sample trial from *Experiment 1: Inferring preferences*. Each trial portrays a speaker and a listener. The speaker produces an utterance to refer to one of the objects. The listener picks the object with the orange-dotted outline. Participants were tasked with evaluating what preferences of the listener may have led her to the particular object choice, specifying their inference by adjusting the sliders for each of the features.

We followed Frank and Goodman (2012) in our stimuli creation. Objects were allowed to vary along three dimensions: color (blue, red, green), shape (cloud, circle, square), and pattern (solid, striped, polka-dotted). The speaker’s utterance was chosen at random from the properties of the three objects present, and the listener’s choice was chosen at random from the subset of the three objects that possessed the uttered feature. By varying the object properties, the targeted object, and the utterance, we generated a total of 2400 scenes. Speaker and listener names were chosen randomly in each trial. Participants saw the speaker’s utterance in bold (e.g., “red” in Figure 2) and the listener’s choice appeared with a dotted orange outline (e.g., the center object in Figure 2). Based on the observed choice, participants were instructed to adjust a series of six sliders to indicate how likely it is that the listener

had a preference for a given feature. The sliders specified the six feature values of the two feature dimensions that were not mentioned in the speaker’s utterance (e.g., pattern and shape in Figure 2).

To compare our model’s predictions to the human data, we calculated an average value for each slider. We excluded the sliders if their corresponding feature value was not present in a scene. For example, for the trial depicted in Figure 2, we excluded the sliders for solid things and squares since none of these are present, and therefore no learning about them is possible.

To determine model correlations with the gathered data, we partitioned the data into ambiguity classes, similar to Frank and Goodman (2012). Depending on how many features competitor objects share with the chosen object, we were able to identify 48 ambiguity classes, which group the constellations that have the exact same ambiguity pattern. The ambiguity classes identified in Experiment 1 distinguish how many objects are referenced by the utterance, how the referenced objects differ in their two non-uttered features, and how the non-referenced objects differ from the referenced objects and from each other. As a result, each ambiguity class yields unique model predictions for the individual features present (with respect to their “ambiguity role” in the particular ambiguity class) in corresponding scenarios  $S$ , effectively distinguishing all model-relevant cases. Please see the Appendix for examples of different classes.

Participants completed a series of fifteen trials. Objects and utterances were chosen as detailed above, with the constraint that ten trials were potentially informative (i.e., when ambiguity is resolved selectively) and five trials were uninformative (e.g. when the utterance was unambiguous or when choosing amongst identical objects) with respect to the listener’s preferences.

### 4.3 Free parameters and optimization procedure

We fit the model parameters either at the individual level or at the group level by optimizing the KL divergence between the data and the model predictions:

$$\text{KL}(P_{data}(f \mid u, s) \parallel P_{model}(f \mid u, s)), \quad (11)$$

where  $P_{data}(f \mid u, s)$  specifies a participant’s normalized slider-value setting, which offers empirical estimates of the feature-preference posterior given object scene  $S$ , a particular utterance choice  $u$ , and the consequent object choice  $s$ ;  $P_{model}(f \mid u, s)$  specifies the corresponding model posterior, either  $P_{S_2}(f \mid u, s)$  in the case of the full model or  $P_{S_{1-simp}}(f \mid u, s)$  in the case of the simple model. By minimizing the summed KL divergence between the empirical and model-predicted preference posteriors over all considered trials, we maximize the model fit to the participants’ data. Moreover, we can use the minimized KL divergence values to calculate the  $G^2$ -statistic and perform the likelihood-ratio test for nested models, since  $G^2$  values are approximately chi-square distributed (Lewandowsky & Farrell, 2011). Individual vs. global parameter fitting allows us to explore potential differences between

participants. In the case of individual model parameter optimization, parameters were optimized for each individual participant separately, determining the KL divergence with respect to the participant-specific set of trials. In the case of global optimization, all trials of all participants were used to determine the summed KL divergence.

We fit three parameters for our full model and two for the simple model. The soft-max scaling factor  $\alpha$  is only relevant for the full model; it controls how likely speaker  $S_1$  is to maximize utility when choosing utterances. The default value is typically set to  $\alpha = 1$  (i.e., no scaling).

The softness parameter  $\gamma$  regulates the strength of individual feature preferences  $f$ :

$$P(s | f) \propto \begin{cases} 1 + \gamma, & \text{if } s \text{ contains } f \\ \gamma, & \text{otherwise} \end{cases}, \quad (12)$$

controlling the choice probability of those objects  $s$  that contain feature  $f$  compared to those that do not. A value of  $\gamma = 0$  models a hard preference choice; in this case, the speaker always chooses one of the preferred objects. On the other hand, when  $\gamma \rightarrow \infty$ , the choice prior becomes uniform over all objects, thus ignoring feature preferences. For example, in the trial shown in Figure 2, there are two objects that match the utterance  $u = \text{“red”}$ : a red striped cloud and a red dotted circle. When  $\gamma = 1$ ,  $P(s_{\text{red striped cloud}} | f_{\text{“cloud”}}) = 2/3$ , while  $P(s_{\text{red dotted circle}} | f_{\text{“cloud”}}) = 1/3$ , yielding a soft preference for clouds. We use  $\gamma = 0$ —that is, hard preferences—as the default model value.

Finally, we allow for the possibility of noise in our human data introduced by participants not following instructions. Parameter  $\beta$  models the possibility that listeners choose objects that do not pass the semantic filter of the literal listener, allowing for non-literal interpretations that result in choosing objects whose features do not match the received utterance  $u$ . The computation is equivalent to the softness parameter above, in this case softening the object choices of the literal listener  $L_0$  towards a uniform choice over all objects present. Again,  $\beta = 0$  models a hard object choice—that is, full obedience to the uttered instruction  $u$ —while  $\beta \rightarrow \infty$  models a uniform object choice—that is, full ignorance of  $u$ .

The effect of this parameter being higher than the default value of 0 is visible when the properties of objects that do not qualify as a possible referent given an utterance (for example, blue objects when the speaker said “red”) affect the slider values. Consider a scene containing three objects: a red square, a red circle, and a blue cloud. If the listener picks the red square following the utterance “red” and rejects the red circle, the subject is expected to indicate that the listener prefers squares over circles. Since there are no red clouds in that scene, the slider for clouds cannot be adjusted to an informative value. However, if the obedience parameter  $\beta$  is above 0, the modeled participant would lower the slider for clouds, since a cloud was not chosen (despite the fact that the blue cloud did not qualify as a potential referent). This parameter thus allows the model to capture noise associated with adjusting sliders when no explicit evidence for updating a particular type

Model	$r^2$	$df$	$F$	Soft. $\gamma$	Obed. $\beta$	Util. $\alpha$
<i>Simple model</i>						
Softness $\gamma$	0.8607	1,190	1181	0.2204		
Softness $\gamma$ & Obedience $\beta$	0.9788	1,190	8823	0.2210	0.2693	
<i>Full model</i>						
Softness $\gamma$	0.8568	1,190	1144	0.2231		
Softness $\gamma$ & Utility $\alpha$	0.8607	1,190	1144	0.2205		0.1797
Softness $\gamma$ & Utility $\alpha$ & Obedience $\beta$	0.9772	1,190	8170	0.2214	0.0030	0.2657

Table 1: Optimization summary (global) for Experiment 1 model predictions.

of preference was actually present in the scene.

#### 4.4 Results

Linear regression analysis was used to test whether the model values predicted the human data. As Table 1 shows, the full and simple models with softness ( $\gamma$ ) optimized globally provide nearly identically-good fits to the data (full:  $r^2 = 0.86$ ; simple:  $r^2 = 0.86$ ).<sup>3</sup> The simple model with both parameters optimized globally captures nearly all of the variance in the human data ( $r^2 = 0.98$ ); in this model, the obedience parameter  $\beta$  is estimated at 0.27 after optimization, indicating that participants indeed considered (possibly subconsciously) the option to interpret utterances non-literally. The full model with both  $\alpha$  and  $\gamma$  optimized performs slightly better than the version with only  $\gamma$  optimized ( $r^2 = 0.86$ ). When optimizing all three parameters, the full model performs markedly better ( $r^2 = 0.98$ ), on a par with the simple model with both parameters optimized.

Overall, the models accurately capture the participants’ inferences over the feature preferences that lead to the choice of an object. Moreover, the greater complexity and added flexibility of the full model (controlled by the added parameter  $\alpha$ ) does not yield any modeling improvement, implying that performance in our task can be modeled as a more shallow reasoning process by means of the simple model.

We now compare our two model variants further when fitting the parameters to the individual data of each participant separately. In situations when the population is potentially heterogeneous, individual-level modeling in reference games improves the fit of the model despite its increased complexity (Franke & Degen, 2016). We optimized  $\alpha$  and  $\gamma$  in light of the KL divergence between the individual participants’ slider-value choices and the corresponding model predictions. We then again averaged the individualized model prediction values and participants’

<sup>3</sup>Here and throughout the paper we report adjusted  $r^2$  values. All results were significant at the  $p < 0.001$  level if not stated differently in the text.

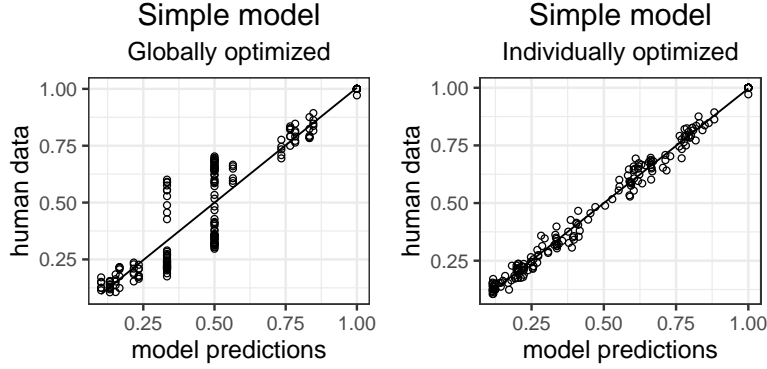


Figure 3: Human data from Experiment 1 plotted against the predictions of our simple model. Each data point indicates the slider values and model-predicted feature-preference posteriors for a particular ambiguity class. Left panel:  $\gamma$  *optimized globally* ( $r^2 = 0.8614$ ); right panel:  $\gamma$  and  $\beta$  *optimized individually* with leave-one-out cross-validation ( $r^2 = 0.9901$ ).

Model	$r^2$	$df$	$F$
<u>Simple model</u>			
Softness $\gamma$	0.8631	1,190	1205
$\gamma$ & Obedience $\beta$	0.9919	1,190	23480
<u>Full model</u>			
$\gamma$ & Utility $\alpha$	0.8627	1,190	1201
$\gamma$ , $\beta$ , & $\alpha$	0.9908	1,190	20620

Table 2: Optimization summary (individual) for Experiment 1 model predictions.

slider values with respect to the particular ambiguity classes and calculated correlations between the data and the model.

As Table 2 shows, the full model optimized at the individual level for the additional parameter  $\alpha$  does not improve the fit compared to the simple model when both models are optimized for  $\gamma$  (simple:  $r^2 = 0.8631$ ; full:  $r^2 = 0.8627$ ) or for  $\gamma$  and  $\beta$  (simple:  $r^2 = 0.9919$ ; full:  $r^2 = 0.9908$ ). Seeing that both models again fit the data nearly equally well (if anything, the simple model performs slightly better), we only consider the predictions of our simple model henceforth. The model fit of the simple model improves considerably when we additionally fit the obedience parameter  $\beta$  at the individual level ( $r^2 = 0.9919$ ). The likelihood ratio test (two-tailed) revealed that a  $\gamma$ - and  $\beta$ -optimized simple model provides a better fit compared to a model optimized only for  $\gamma$  ( $G^2 = 237.36, df = 82, p < 0.01$ ). The more complex model contains one additional parameter  $\beta$  fitted for each subject, giving us 82 degrees of freedom. We additionally checked the generalizability of the model by performing leave-one-out cross-validation on the individual level. Figure 3 shows that the resulting cross-validated model predictions retain the strong

fit ( $r^2 = 0.99, F(1, 190) = 18910$ ).

It bears noting that the individually-fitted parameters do not improve the correlation values much, if at all, when compared to the globally-fitted model. To appreciate the gains obtained by fitting model parameters, Figure 4 shows the average responses of the human participants, of the individually-, two-parameter-optimized simple model, and of the non-optimized simple model for the scene type of the sample trial from Figure 2. In that trial, participants saw that the middle object was chosen following the utterance “red”. There are two potential referents for this description: the red striped cloud and the red dotted circle. Since the cloud was chosen, we infer that the person who chose this object has a preference for clouds over circles, and for striped objects over dotted ones. Note that we cannot learn anything about the preference for solid things or squares in this trial because these features are not present, thus we ignore the respective slider values. Moreover, we cannot really learn anything about color preferences because the color was uttered; thus, sliders for those features were not present. As Figure 4 shows, both humans and the models assign high slider values to clouds and striped things, and low values to circles and dotted things. Indeed, even the non-optimized model fits the qualitative pattern of the results; optimizing  $\beta$  and  $\gamma$  improves the quantitative fit.

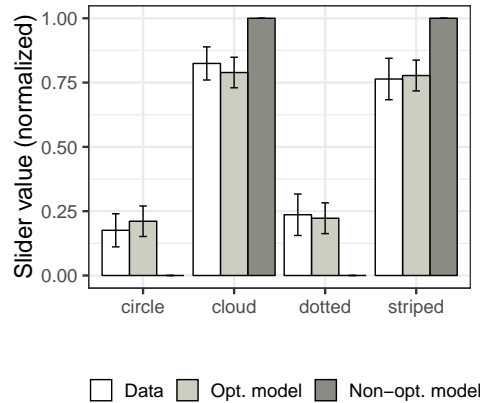


Figure 4: Human data and predictions from the individually-, two-parameter-optimized simple model and the non-optimized simple model for the scenario  $S$  shown in Figure 2. Error bars represent 95% confidence intervals.

We thus find strong empirical support for our simple model, implying that speakers are indeed able to use listener behavior to acquire information about their preferences. We fail to find that the full model predicts the data better. This result suggests that the task in our experiments does not require full-blown pragmatic inference about alternative utterances. The question now turns to whether speakers are able to capitalize on this reasoning when it comes to selecting utterances. In other words, are speakers able to use ambiguous language in a socially epistemic, strategic manner?

## 5 Experiment 2: Epistemic utterance choice

Our next task is to check the predictions of our strategic utterance-selection model: given a set of potential referents  $S$ , will participants reason pragmatically about the anticipated potential epistemic utility of utterances  $u \in U$  for inferring the listener’s preferences?

### 5.1 Participants


We recruited 90 participants with US IP addresses through Amazon.com’s Mechanical Turk crowdsourcing service; participants in Experiment 1 were not eligible to participate in Experiment 2. Participants were compensated for their participation. On the basis of a post-test demographics questionnaire, we again identified 82 participants as native speakers of English; their data were included in the analyses. We obtained a confirmation from all participants that they agree to participate in the study.

### 5.2 Design and methods


Participants encountered a reference game scenario similar to Experiment 1 in which a speaker signals an object to a listener who might have a preference for certain types of objects. However, rather than observing the utterance and referent choice, participants were now tasked with helping the speaker choose an utterance that was “most likely to reveal the listener’s color, shape, or pattern preferences.” Figure 5 shows a sample trial, in which the speaker (“Katie” in the example) is to choose an utterance in order to learn about the listener’s preferences (“Elizabeth” in the example). While the ambiguous utterances “cloud”, “green”, and “striped” may allow inferences about color & texture, shape & texture, and color & shape, respectively, the utterances “solid”, “blue”, and “circle” leave only one response option to the listener, such that the speaker cannot learn about the listener’s preferences when observing the listener’s response (assuming the listener obeys the speaker’s order).

We used the same sets of objects from Experiment 1, which could vary along three dimensions. Each trial featured a set of three objects, as in Figure 5. After observing the objects, participants adjusted sliders to indicate which single-feature utterance the speaker should choose to learn about the preferences of their listener. Potential utterances corresponded to the features of the objects present; depending on the number of unique features, participants adjusted between three and nine sliders. As with Experiment 1, we averaged the data and the respective model predictions across specific ambiguity classes, which include all scenes that yield identical utterance choice options. In this case, 14 distinct conditions can be identified, with a total of 84 slider values to set. Membership within an ambiguity class is defined by how many objects in a scene share each of the features: shape, pattern, and color. If objects share a feature, we also consider whether these objects also















Progress: 

Suppose Katie wants to learn about Elizabeth's preferences in the following scenario:



Katie can choose a single utterance and then watch Elizabeth select an object.

What should Katie say?

	definitely not		definitely
"cloud"			
"solid"			
"green"			
"striped"			
"blue"			
"circle"			

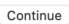


Figure 5: A sample trial from *Experiment 2: Choosing utterances*.

share other features. For example, in Figure 5, two green objects differ in shape, making the utterance “green” informative. If, on the other hand, both green objects were clouds, uttering “green” would not allow the speaker to update their beliefs about the listener’s shape preferences. In the most extreme case, when all objects share all three features, all utterances are ambiguous since multiple objects can always be picked; but no utterance allows the speaker to learn anything about the listener because the object choice is uninformative. Another extreme case is a situation where all objects are unique and do not share any features. In such a case, any utterance will only pick one object, making learning about preferences impossible unless obedience ( $\beta$ ) is not 0—that is, unless listeners have a tendency to disobey the utterance and consider objects that do not satisfy its literal interpretation.

Just like for Experiment 1, each ambiguity class yields unique model predictions for the individual features present in the respective scenarios  $S$ , taking into account their ambiguity role. This grouping strategy effectively distinguishes all model-relevant cases. Please see the Appendix for examples of different classes.

Participants completed a series of fifteen trials. As with Experiment 1, objects were chosen at random, with the constraint that ten trials were potentially informative with respect to the listener’s preferences (as in Figure 5) and five trials were uninformative with respect to the listener’s preferences (e.g., observing a set of three identical objects).

### 5.3 Results

We use our simple model to compute the expected most-informative utterance for inferring preferences. In other words,  $P_{S_1\text{-simp}}(u)$  calculates the probability that a speaker would choose  $u$  for the purpose of inferring preferences.

To generate predictions from  $P_{S_1\text{-simp}}(u)$ , three free parameters can be identified: the preference softness  $\gamma$ , the obedience  $\beta$ , and the  $\lambda$  parameter, which factors the importance of choosing the expected most-informative utterance with respect to the expected KL divergence between preference priors and expected preference posteriors (cf. equation 10). While a positive  $\lambda$  value yields the intention to maximize information gain, a negative value results in a tendency to minimize information gain, that is, a preference for no change in the posterior feature preference estimate  $P_{S_1\text{-simp}}(f | u, s)$ , in comparison to the prior estimate  $P(f)$ . A value of  $\lambda = 0$  effectively ignores information gain and a resulting tendency to choose the object that was most likely referenced given the *utterance*.

We compare the performance of our simple model with the performance of a uniform baseline model, which merely chooses one of the available utterances at random. Seeing that in particular ambiguity cases with particular constellations  $S$  up to nine utterances are possible, the baseline model yields different model predictions for the available utterances in the respective ambiguity classes. As a result, the baseline model is much better in capturing variance in the data than one would expect without this insight ( $r^2 = 0.75$ ; Table 3). The non-optimized simple model, surprisingly, captures very little variance in the human data ( $r^2 = 0.06$ ); for the non-optimized model we set the parameters to hard preference and full obedience ( $\gamma = 0, \beta = 0$ ) and the information gain factor  $\lambda$  to 1, thus preferring to choose those utterances that are expected to yield high information gain. Figure 6 compares the performance of this non-optimized simple RSA model with the baseline model that relies on a simple heuristic of distributing the probability mass equally among all utterances that are available in a scene.

Model	$r^2$	$df$	$F$	Soft. $\gamma$	Obed. $\beta$	Inf. gain $\lambda$
Baseline	0.7466	1,82	245.6			
Non-optimized	0.0595	1,82	6.253			
Softness & Obedience & Inf. gain	0.7991	1,82	331.2	0.0006	0.2758	0.3663

Table 3: Experiment 2. Optimization summary (global).

To examine the reasons for this failure of the non-optimized simple model, we first performed additional global parameter optimization runs. Based on the AIC scores<sup>4</sup>, all the globally-optimized models outperform the uniform model (AIC = 888.048): softness (AIC = 877.47), obedience (AIC = 877.07), KL-value factor (AIC = 879.004). Models that take combinations of several free parameters into

<sup>4</sup>Burnham and Anderson (2002) suggest that AIC can be used for non-nested model comparison.

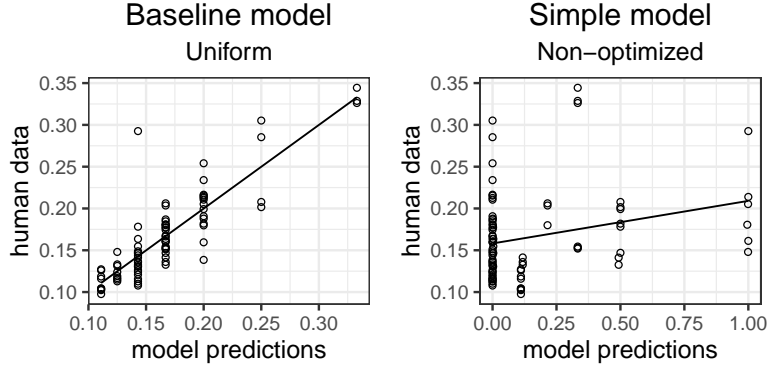


Figure 6: Average human data from Experiment 2 plotted against the predictions of the uniform baseline model (left;  $r^2 = 0.75$ ) and the simple model (right;  $r^2 = 0.06$ ).

account did not yield additional improvements. Thus, when optimizing all of the simple RSA model parameters, there is good evidence that the model accounts for more variance than the uniform base model (Lewandowsky & Farrell, 2011). Figure 7 (left) shows the correlation plot ( $r^2 = 0.7991$ ,  $F(1, 82) = 331.2$ ,  $p < 0.001$ ). The plot reveals that qualitatively, even a model with three-way global optimization is quite similar to the uniform model. Moreover, the analysis of the obtained parameter values (Table 3) suggests that the parameter of primary interest  $\lambda$  has a rather small value of 0.3663 indicating only moderate information gain intention of the subjects. This result could either reflect the overall reasoning strategy of all subjects or come as a consequence of averaging over a heterogeneous group of them, suggesting an analysis of between subject differences.

Turning to individual parameter optimization, we compared three single-parameter individually-optimized simple RSA models to determine which model provides the best fit to the data. All models have similar levels of complexity, with either softness  $\gamma$ , obedience  $\beta$ , or the KL-factor  $\lambda$  being optimized. The results indicate that we get the best fit by optimizing the KL-factor  $\lambda$  ( $r^2 = 0.9059$ ,  $F(1, 82) = 800.2$ ; leave-one-out cross-validated optimization  $r^2 = 0.8902$ ,  $F(1, 82) = 664.8$ ), with other models capturing less variance in the data ( $\beta$ -optimized  $r^2 = 0.8015$ ,  $F(1, 82) = 336.1$ ;  $\gamma$ -optimized  $r^2 = 0.8077$ ,  $F(1, 82) = 349.6$ ).

AIC scores indicate that optimizing for  $\lambda$  individually indeed significantly improves the model fit (783.167) compared to the uniform model (888.047), while optimizing for softness or obedience results in scores actually higher than the uniform model due to the parameter penalty (softness: 1020.668, obedience 995.754). Two- and three-parameter individual optimizations did not yield any significant model improvements when compared to the nested, individually- $\lambda$ -optimized model (best improvement when optimizing  $\gamma$  in addition to  $\lambda$ :  $G^2 = 24.72$ ,  $df = 82$ , n.s.). Figure 7 (right) shows the resulting correlation plot for the individually- $\lambda$ -optimized model.

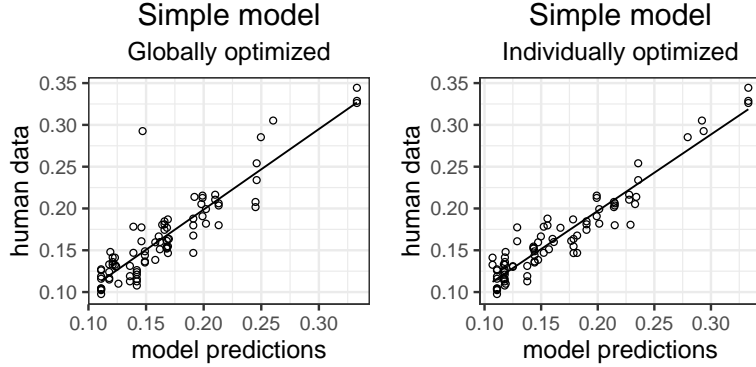


Figure 7: Average human data from Experiment 2 plotted against the predictions of the optimized simple models. Left panel: *globally-optimized three-parameter model* ( $r^2 = 0.7466$ ); right panel: *individually- $\lambda$ -optimized model* ( $r^2 = 0.9059$ ).

Unlike for Experiment 1, where even the non-optimized models provided a good linear fit to the data, optimization produces a large effect on the model predictions in Experiment 2. Figure 8 compares globally-optimized, individually-optimized vs. non-optimized model predictions against the human behavior for the sample trial in Figure 5. We see that the non-optimized model strongly favors ambiguous utterances: only the ambiguous utterances *cloud*, *striped*, and *green* (i.e., the utterances that point to more than one object in the scene) promise information gain about the listener’s preferences when assuming full obedience to the utterance. However, Figure 8 shows that human behavior deviates quite strongly from the non-optimized, ambiguity-selecting baseline. Once we optimize  $\lambda$ , we are able to capture human behavior in the task.

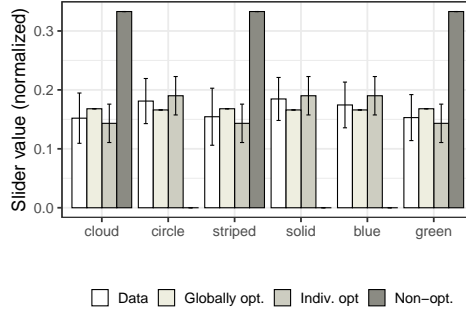


Figure 8: Model predictions and human data for the ambiguity class in in Figure 5. The optimized version of the model is optimized for the KL-factor  $\lambda$ . Error bars represent 95% confidence intervals.

When examining the individually-optimized model values in further detail, we noticed three groups of participants. The first one may be termed a “lazy worker”

or “unpredictable behavior” group: for 28 participants, the KL divergence values of the  $\lambda$ -optimized simple RSA model failed to reach the performance of the baseline model, essentially failing to identify any model-corresponding regularity in the data that goes beyond random utterance choice behavior. The second group of 33 participants yielded more negative values (i.e.,  $-7.11 < \lambda < -0.014$ ,  $\bar{\lambda} = -0.823$ ), indicating that a significant number of participants preferred to systematically choose unambiguous utterances (Uniform model AIC=460;  $\lambda$ -individually optimized model AIC=346). The third group of 21 participants yielded positive values (i.e.,  $.0187 < \lambda < .537$ ,  $\bar{\lambda} = -0.124$ ), indicating that these participants indeed preferred the more ambiguous utterances in a strategic manner (Uniform model AIC=344;  $\lambda$ -individually optimized model AIC=284).

Further experiments with highly similar setups confirmed this trend. In particular, we ran two additional, complementary studies with a blocked design where participants first completed preference-inference trials as in Experiment 1 and then utterance-selection trials as in Experiment 2. In the first complementary study with 10 trials (135 participants, data from 123 native speakers of English included in the analysis, 12 non-native speakers excluded), an identical analysis yielded 42% of participants that preferred ambiguous over unambiguous utterances (37% unpredictable participants; 21% preferred unambiguous utterances). In the second complementary study with 54 participants (two participants excluded as non-native speakers), which contained 30 trials in total and had slightly more general instructions, as many as 64% of the participants systematically preferred ambiguous over unambiguous utterances (21% unpredictable workers; 15% preferred unambiguous utterances).

## 6 Discussion

The ability to infer the intentions and goals of others upon observing their behavior develops early in ontogeny (Liu & Spelke, 2017) and remains a valuable source of information for building predictive models of others. In this paper, we asked whether communicative behavior—specifically the resolution of ambiguous reference—provides interpretable data for speakers to learn more about their listeners. Our primary goal was to develop a computational model of this inference process.

Our model offers an articulated hypothesis about *how* this reasoning proceeds: when speakers are aware of the ambiguity in their utterances, observing how listeners resolve that ambiguity provides clues about the preferences listeners use when doing so. We have found strong support for the idea that speakers can indeed learn about others when observing their interpretation of ambiguous utterances. Experiment 1 demonstrates that naïve speakers are able to reason pragmatically about *why* listeners may take the actions they do in cases of referential ambiguity. This result connects our findings to the Bayesian Theory of Mind literature that describes how children and adults use observable behavior to infer the mental states of others

(Baker et al., 2009, 2017).

Moreover, we have hypothesized that speakers may choose ambiguous over unambiguous utterances in anticipation of social information gain. We formalized this epistemic, information seeking behavior as an utterance choice behavior that attempts to maximize the anticipated averaged Kullback-Leibler divergence between a uniform prior and the possible posterior preference distributions. The results of Experiment 2 are mixed. Some speakers are able to capitalize on this reasoning to strategically select ambiguous utterances that are expected to improve their understanding of the preferences of their listeners. Others did not exhibit any model-conform behavior. Yet others showed a significant preference for unambiguous utterances.

The results of our utterance-choice experiment (Experiment 2) suggest that not all speakers managed to strategically select ambiguous utterances to yield situations of ambiguity resolution. In the case of the participants who significantly preferred unambiguous utterances, two explanations may be warranted. First, it may be the case that these participants think overly egocentrically, thus having the intention to signal their own preferences rather than to give options to the listener. Second, the subjects may strategically choose unambiguous utterances if they anticipate a possibility that the listener will disobey the instructions and choose her preferred object even if it does not strictly qualify, revealing her preferences. For example, if the listener has a strong preference for red objects, she might pick a red square following the utterance “circle” if none of the circles are red. Disobeying instructions may appear an outrageous strategy upon first glance, but in fact the situations when people do not follow instructions or even polite requests are not that uncommon. Particularly in the case where the model did not yield fits better than the uniform model, however, it may simply be the case that these participants did not gain access to the required deeper reasoning process, seeing that it clearly requires additional computational resources (Lieder & Griffiths, 2020). As a result, these participants may have chosen erratically and sometimes preferred to give instructions with predictable outcomes. In fact, the deeper anticipatory reasoning skill might only be present in speakers with high meta-linguistic competence, such as, for example, professional writers, who have been shown to strategically use ambiguity to create an artistic effect of uncertainty about the interpretation of characters and events (Bade, Bauer, Beck, Dörge, & Zirker, 2015; Bauer & Zirker, 2014; Quigley, 2015). Clearly, seeing the mixed results, the floor is open for alternative explanations and models and corresponding future research.

Besides hypothesizing individual differences in the reasoning resources and social focus of our participants, one might question the ecological validity of such epistemically-motivated ambiguity behavior in general. Intuitively, such a strategy appears to be a sub-optimal way to elicit information when compared to a direct question—in the case of our experiments, a question about the preferences of the listener. Yet, considerations of politeness (Yoon et al., 2018) or a general rhetorical strategy of indirectness could make the costs of asking a question prohibitively high. For example, asking a stranger directly about her political views

might appear impolite or even aggressive, while making a vague statement about the course of an election by stating that their results are ‘interesting’ could act as a probe to get the opinion of the listener without forcing her to react. Evidence from a corpus of speed-dating dialogues suggests that direct questions in situations of courtship—supposedly a setup that should promote information-seeking behavior—are a sign that the interaction between the partners is not going particularly smoothly (McFarland, Jurafsky, & Rawlings, 2013); to keep a stalled conversation going, speakers have to ask a question rather than continue the flow of the previous discussion.

However, ambiguity does not need to be an explicit strategy of the speaker to make the inference process possible. Rather, ambiguity is naturally present in language as a consequence of economy considerations and indirectness strategies, providing ample interpretative freedom to the listener. In that sense, observing ambiguity resolution is a consequence not of active but rather retrospective inference that emerges once the speaker realizes that an utterance was potentially ambiguous. Taken together, the results of our experiments and modeling indicate that humans are aware of the fact that by observing responses to ambiguous utterances, information about the listener’s prior preferences can be inferred. They are able to learn about the hidden, internal cognitive states of others, including preferences, but probably also beliefs, knowledge, or intentions.

When viewed from a broader perspective, our analysis can be closely related to a part of the event-predictive mind of the listener and the speaker (Butz, 2016; Butz & Kutter, 2017; Butz, Achimova, Bilkey, & Knott, 2020; Zacks, Speer, Swallow, Braver, & Reynolds, 2007). This paradigm emphasizes the role of predictions in determining the success of our behavior. In order to anticipate the changes in the environment, including the behavior of other human agents, the predictive mind continuously updates its models of others by incorporating novel behavioral evidence. When interpreting an utterance—in our experimental setup, opening up a set of referent choices—the listener’s mind infers the current choices and integrates them with her preference priors, implicitly anticipating possible choice consequences. Moreover, the expected information gain term—computing the utterance choice of the speaker—can be equated with the computation of socially-motivated active inference (Butz, 2017; Clark, 2016; Friston et al., 2015). This inference causes the model to strive for an anticipated epistemic value that quantifies the expected information gain about the preferences of the listener—that is, expecting a form of social information gain.

In conclusion, our results have shown that while epistemic-oriented behavior is hard to elicit, the observation of ambiguity resolution behavior clearly opens possibilities for conversation partners to make approximate Bayesian inferences about each others hidden, internal cognitive states and processes. Here we have focused on preference inferences, where preference co-determined object choice behavior. Future research should enhance on this aspect and further relate to the broader Bayesian theory of mind literature as well as to inverse planning (Baker et al., 2009; Russell, n.d.) hypothesizing more complex cognitive models of the

listener. As a result, we expect that listeners, speakers, and observers will be able to infer aspects of the underlying world knowledge, of the current beliefs about the world and about the common ground of the unfolding conversation, and of the intentions of the involved conversation partners. Intriguingly, such inference processes can theoretically be also implemented in machines.

## Funding

This project has been funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation)–Project number 198647426.

## Data availability

Data supporting the findings of this study are available from the corresponding author upon request.

## References

- Bade, N., Bauer, M., Beck, S., Dörge, C., & Zirker, A. (2015). Ambiguity in Shakespeare's sonnet 138. In S. Winkler (Ed.), *Ambiguity: Language and communication* (pp. 89–109). De Gruyter.
- 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), 1–10.
- Baker, C. L., Saxe, R., & Tenenbaum, J. B. (2009). Action understanding as inverse planning. *Cognition*, 113(3), 329–349. doi: 10.1016/j.cognition.2009.07.005
- Bauer, M., & Zirker, A. (2014). Dickens and ambiguity: the case of a tale of two cities. In C. Huguët & N. Vanfasse (Eds.), *Dickens, modernism and modernity* (pp. 209–229). Paris: Editions du Sagittaire.
- Belardinelli, A., Lohmann, J., Farnè, A., & Butz, M. V. (2018). Mental space maps into the future. *Cognition*, 176, 65–73.
- Belardinelli, A., Stepper, M. Y., & Butz, M. V. (2016). It's in the eyes: Planning precise manual actions before execution. *Journal of vision*, 16(1), 1–18.
- Burnham, K. P., & Anderson, D. R. (2002). *Model selection and multimodel inference. a practical information-theoretic approach* (2nd ed.). Springer.
- Butz, M. V. (2008). How and why the brain lays the foundations for a conscious self. *Constructivist Foundations*, 4(1), 1–42.
- Butz, M. V. (2016). Towards a unified sub-symbolic computational theory of cognition. *Frontiers in Psychology*, 7(925). doi: 10.3389/fpsyg.2016.00925
- Butz, M. V. (2017). Which structures are out there? Learning predictive compositional concepts based on social sensorimotor explorations. In T. K. Metzinger



- & W. Wiese (Eds.), *Philosophy and predictive processing*. Frankfurt am Main: MIND Group. doi: 10.15502/9783958573093
- Butz, M. V., Achimova, A., Bilkey, D., & Knott, A. (2020). Event-predictive cognition: A root for conceptual human thought. *Topics in Cognitive Science*.
- Butz, M. V., & Kutter, E. F. (2017). *How the mind comes into being: Introducing cognitive science from a functional and computational perspective*. Oxford, UK: Oxford University Press.
- Carmon, A. F. (2013). Is it necessary to be clear? An examination of strategic ambiguity in family business mission statements. *Qualitative Research Reports in Communication*, 14(1), 87–96. doi: 10.1080/17459435.2013.835346
- Chomsky, N. (2002). An interview on minimalism. In A. Belletti & L. Rizzi (Eds.), *On nature and language* (p. 92-161). Cambridge: Cambridge University Press.
- Clark, A. (2016). *Surfing uncertainty: Prediction, action and the embodied mind*. Oxford, UK: Oxford University Press.
- Degen, J., Tessler, M. H., & Goodman, N. D. (2015). Wonky worlds: Listeners revise world knowledge when utterances are odd. In D. Noelle et al. (Eds.), *Proceedings of 37th Annual Meeting of the Cognitive Science Society*. Austin, TX.
- Evans, O., Stuhlmüller, A., & Goodman, N. (2016). Learning the preferences of ignorant, inconsistent agents. In V. Rus & Z. Markov (Eds.), *Proceedings of the thirtieth AAAI conference on artificial intelligence*. Palo Alto, California: AAAI Press.
- Ferreira, V. S. (2008). Ambiguity, accessibility, and a division of labor for communicative success. *Psychology of Learning and Motivation: Advances in Research and Theory*, 49, 209-246.
- Frank, M. C., & Goodman, N. D. (2012). Predicting pragmatic reasoning in language games. *Science*, 336, 998-998.
- Franke, M., & Degen, J. (2016). Reasoning in reference games: Individual-vs. population-level probabilistic modeling. *PloS One*, 11(5), e0154854.
- Franke, M., & Jäger, G. (2016). Probabilistic pragmatics, or why Bayes' rule is probably important for pragmatics. *Zeitschrift für Sprachwissenschaft*, 35(1), 3–44.
- Frege, G. (1892). Über Sinn und Bedeutung. *Zeitschrift für Philosophie und philosophische Kritik*, 100, 25–50.
- Friston, K., Rigoli, F., Ognibene, D., Mathys, C., Fitzgerald, T., & Pezzulo, G. (2015). Active inference and epistemic value. *Cognitive Neuroscience*, 6, 187-214. doi: 10.1080/17588928.2015.1020053
- Goodman, N. D., & Frank, M. C. (2016). Pragmatic language interpretation as probabilistic inference. *Trends in Cognitive Sciences*, 20(11), 818-829.
- Grice, H. P. (1975). Logic and conversation. In P. Cole & J. L. Morgan (Eds.), *Syntax and semantics 3: Speech acts* (p. 26-40). New York: Academic Press.
- Hayhoe, M. M., Shrivastava, A., Mruczek, R., & Pelz, J. B. (2003). Visual memory and motor planning in a natural task. *Journal of Vision*, 3(1), 49–63. doi: 10:

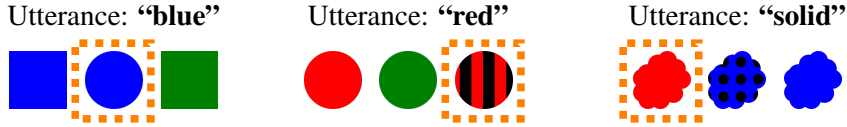
- Jara-Ettinger, J., Gweon, H., Schulz, L. E., & Tenenbaum, J. B. (2016). The naïve utility calculus: Computational principles underlying commonsense psychology. *Trends in Cognitive Sciences*, 20(8), 589–604.
- Jones, E. E., & Davis, K. E. (1965). From acts to dispositions the attribution process in person perception. In *Advances in experimental social psychology* (Vol. 2, pp. 219–266). Elsevier.
- Kao, J., Bergen, L., & Goodman, N. (2014). Formalizing the pragmatics of metaphor understanding. In P. Bello, M. Guarini, M. McShane, & B. Scassellati (Eds.), *Proceedings of the 36th Annual Meeting of the Cognitive Science Society*. Austin, TX: Cognitive Science Society.
- Kelley, H. H. (1967). Attribution theory in social psychology. In *Nebraska Symposium on Motivation* (Vol. 15, pp. 192–238).
- Kelley, H. H., & Stahelski, A. J. (1970). Social interaction basis of cooperators' and competitors' beliefs about others. *Journal of Personality and Social Psychology*, 16(1), 66–91.
- Levinson, S. C. (2000). *Presumptive meanings: The theory of generalized conversational implicature*. Cambridge, MA: MIT Press.
- Lewandowsky, S., & Farrell, S. (2011). *Computational modeling in cognition: Principles and practice*. Thousand Oaks: Sage Publications.
- Lieder, F., & Griffiths, T. L. (2020). Resource-rational analysis: Understanding human cognition as the optimal use of limited computational resources. *Behavioral and Brain Sciences*, 43, e1. doi: 10.1017/S0140525X1900061X
- Liu, S., & Spelke, E. S. (2017). Six-month-old infants expect agents to minimize the cost of their actions. *Cognition*, 160, 35 - 42. Retrieved from <http://www.sciencedirect.com/science/article/pii/S001002771630302X> doi: <https://doi.org/10.1016/j.cognition.2016.12.007>
- Lohmann, J., Belardinelli, A., & Butz, M. V. (2019). Hands ahead in mind and motion: Active inference in peripersonal hand space. *Vision*, 3(2), 15. doi: [doi.org/10.3390/vision3020015](https://doi.org/10.3390/vision3020015)
- Lucas, C. G., Griffiths, T. L., Xu, F., Fawcett, C., Gopnik, A., Kushnir, T., ... Hu, J. (2014). The child as econometrician: A rational model of preference understanding in children. *PloS One*, 9(3), e92160.
- McFarland, D. A., Jurafsky, D., & Rawlings, C. (2013). Making the connection: Social bonding in courtship situations. *American Journal of Sociology*, 118(6), 1596–1649.
- Mohr, L. B. (1983). The implications of effectiveness theory for managerial practice in the public sector. In K. S. Cameron & D. A. Whetten (Eds.), *Organizational effectiveness* (pp. 225–239). Elsevier.
- Ossa-Richardson, A. (2019). *A history of ambiguity*. Princeton University Press.
- Pascale, R. T., & Athos, A. G. (1981). *The art of Japanese management*. New York: Simon & Schuster.
- Piantadosi, S. T., Tily, H., & Gibson, E. (2012). The communicative function of ambiguity in language. *Cognition*, 122, 280–291.

- Qing, C., & Franke, M. (2015). Variations on a Bayesian theme: Comparing Bayesian models of referential reasoning. In H. Zeevat & H.-C. Schmitz (Eds.), *Bayesian natural language semantics and pragmatics* (p. 201-220). Springer.
- Quigley, M. (2015). *Modernist fiction and vagueness: Philosophy, form, and language*. Cambridge University Press.
- Russell, S. (n.d.). The purpose put into the machine. In J. Brockman (Ed.), *Possible minds: 25 ways of looking at ai* (p. 20-32). New York: Penguin Press.
- Sennet, A. (2016). Ambiguity. In E. N. Zalta (Ed.), *The Stanford Encyclopedia of Philosophy* (Spring 2016 ed.). Metaphysics Research Lab, Stanford University. <https://plato.stanford.edu/archives/spr2016/entries/ambiguity/>.
- Shafto, P., Goodman, N. D., & Frank, M. C. (2012). Learning from others: The consequences of psychological reasoning for human learning. *Perspectives on Psychological Science*, 7(4), 341–351.
- Sikos, L., Venhuizen, N., Drenhaus, H., & Crocker, M. (2019, 04). *Reevaluating pragmatic reasoning in web-based language games*. doi: 10.13140/RG.2.2.30535.14249
- Wasow, T. (2015). Ambiguity avoidance is overrated. In S. Winkler (Ed.), *Ambiguity: Language and communication* (p. 29-47). de Gruyter.
- Woensdregt, M., Kirby, S., Cummins, C., & Smith, K. (2016). Modelling the co-development of word learning and perspective-taking. In A. Papafragou, D. Grodner, D. Mirman, & J. Trueswell (Eds.), *Proceedings of 38th Annual Meeting of the Cognitive Science Society* (pp. 1241–1246). Austin, TX: Cognitive Science Society.
- Yoon, E. J., Frank, M. C., Tessler, M. H., & Goodman, N. D. (2018, Dec). *Polite speech emerges from competing social goals*. PsyArXiv. Retrieved from [psyarxiv.com/67ne8](https://psyarxiv.com/67ne8) doi: 10.31234/osf.io/67ne8
- Zacks, J. M., Speer, N. K., Swallow, K. M., Braver, T. S., & Reynolds, J. R. (2007). Event perception: A mind-brain perspective. *Psychological Bulletin*, 133(2), 273–293. doi: 10.1037/0033-2909.133.2.273

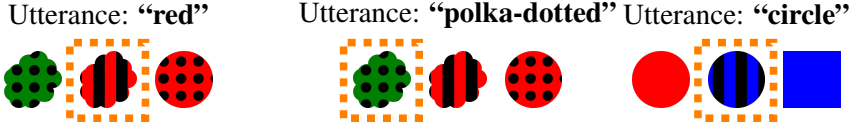
## A Ambiguity classes

### Experiment 1

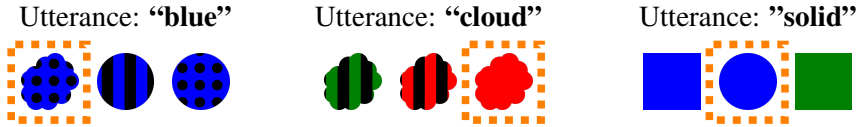
Figure 9 shows three exemplar scenarios for three representative ambiguity classes. Let us consider the first class in more detail. In the scenario  $S$  on the left side of Figure 9a, the utterance “blue” refers either to the blue square or the blue circle. The picked object, that is, the blue circle, is unique in its shape (circle) and shares the other non-referenced property with both other objects (that is, its solid pattern). The referenced but not picked object (that is, the blue square), shares its shape with the non-referenced object. In the scenario  $S$  in the center, the referenced two red objects differ in texture but share shape with the non-referenced object. In the scenario  $S$  on the right, the referenced two solid objects can be contrasted in their color but share their shape with the third object.



(a) The utterance potentially references two objects, the picked object has one non-referenced unique feature, while the other, non-referenced feature is shared among all three objects. The other referenced—but not chosen—object shares its other feature with the non-referenced object.



(b) The utterance  $u$  references two objects whereby both objects only share the uttered feature. The third object shares one feature with each of the two referenced objects.



(c) In this third exemplar ambiguity class, the utterance refers to all three objects. The picked object shares one feature with one other object and has one feature just for itself while the other two objects share it.

Figure 9: Three exemplar scenarios  $S$ , constraining utterance  $u$ , and chosen object  $s$  are shown for three exemplar ambiguity classes for Experiment 1.

### Experiment 2

Figure 10 shows three exemplar scenarios for three representative ambiguity classes. Let us again consider the first class in more detail. In the scenario  $S$  on the left side of Figure 10a, all three objects share the feature pattern (solid), while two share

color (blue), and the other two share shape (square). As a result, uttering *green* or *circle* will unambiguously signal a single object to the listener because the utterance identifies one unique object. On the other hand, uttering *solid* will let the listener choose freely, while uttering *blue* or *square* will give a specific choice between two objects, that is, between the blue circle and the blue square or between the blue square or the green square, respectively. In the scenario *S* in the center, the objects share shape (circle), two share pattern (solid), and the other two share color (red). Here, *circle* references all three objects, *red* or *solid* reference pairs of objects, and *striped* or *green* reference one unique object each. In the scenario *S* on the right, the objects again share shape (cloud), two share pattern (solid), while the other two share color (blue).



(a) In this exemplar ambiguity class, one feature is shared by all three objects, while the two other features allow the distinction between two different pairs of objects and the reference of one of two uniquely-identifiable objects.



(b) In this second exemplar ambiguity class, all three feature types allow the identification of pairs of objects or unique objects, where all three features contain one unique feature type, each. As a result, there are three utterances that each pick out a different pair of objects and three other utterances that each reference one single object—effectively allowing for the unique identification of each object as well as the identification of all three possible pairs.



(c) In this third exemplar ambiguity class, two features have three unique values, while one feature allows the identification of a pair of objects.

Figure 10: Three exemplar scenarios *S* are shown for three exemplar ambiguity classes for Experiment 2.