

# Cognitive Science

## Learning about Others: Pragmatic Social Inference through Ambiguity Resolution

--Manuscript Draft--

Manuscript Number:	
Full Title:	Learning about Others: Pragmatic Social Inference through Ambiguity Resolution
Article Type:	Regular Article
Keywords:	ambiguity; pragmatics; information gain; event-predictive cognition; Rational Speech Act models; social intelligence
Corresponding Author:	Asya Achimova, Ph.D. University of Tuebingen Tübingen, GERMANY
Corresponding Author Secondary Information:	
Corresponding Author's Institution:	University of Tuebingen
Corresponding Author's Secondary Institution:	
First Author:	Asya Achimova
First Author Secondary Information:	
Order of Authors:	Asya Achimova
	Gregory Scontras
	Christian Stegemann-Philipps
	Johannes Lohmann
	Martin V. Butz
Order of Authors Secondary Information:	
Abstract:	<p>We investigated whether ambiguity resolution may yield socially-relevant benefits, revealing parts of the privileged ground of the interpreter.</p> <p>In particular, we asked if speakers can (i) use response observations to infer unknown preferences of a listener, and (ii) strategically choose ambiguous utterances for learning about those preferences.</p> <p>We ran experiments in a reference game framework and modeled the data with a pragmatic social inference Rational Speech Act model.</p> <p>Participants were able to infer listeners' preferences when analyzing their choice of objects given referential ambiguity.</p> <p>Moreover, a significant group of speakers were able to strategically choose ambiguous over unambiguous utterances in an epistemic, event-predictive, goal-directed manner, although a different group significantly preferred unambiguous utterances.</p> <p>We conclude that ambiguity resolution indeed reveals aspects of the knowledge, preferences, and beliefs of conversation partners and some of us are able to strategically use ambiguous utterances to gain knowledge about these aspects.</p>

Asya Achimova

Neuro-cognitive Modeling Group,  
Department of Computer Science;  
Research Training Group 1808:  
Ambiguity – Production and Perception

Eberhard Karls Universität Tübingen

asya.achimova@uni-tuebingen.de

January 14, 2020

Dear editors,

We would like to submit a manuscript entitled:  
Learning about others: Pragmatic social inference through ambiguity resolution.

In this paper, we address a fundamental property of human language and communication – ambiguity. Ambiguity seems to be a negative side-effect of an efficient communication system, and the fact that it is so pervasive in language has puzzled researchers for decades. We propose that ambiguity serves an important purpose: it allows the speaker and the listener to reason about hidden beliefs of each other, which lead to individual interpretations of ambiguous phrases.

Here we develop an account of how humans use ambiguity resolution to build more accurate predictive models of each other. This information seeking behavior is aimed at ensuring efficient communication and ultimately efficient interaction between people. Our work brings together several lines of research in linguistics, communication sciences, and mathematical modeling. We develop a formal account of social Bayesian reasoning inspired by Rational Speech Act models. Our analysis offers a foundation for developing precise models of social interaction.

We hope the computational approach we develop will be of interest to the wide readership of the Cognitive Science journal.

Sincerely yours,  
Asya Achimova,

# Learning about Others: Pragmatic Social Inference through Ambiguity Resolution

Asya Achimova  
asya.achimova@uni-tuebingen.de

Gregory Scontras  
gscontra@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

January 14, 2020

## Abstract

We investigated whether ambiguity resolution may yield socially-relevant benefits, revealing parts of the privileged ground of the interpreter. In particular, we asked if speakers can (i) use response observations to infer unknown preferences of a listener, and (ii) strategically choose ambiguous utterances for learning about those preferences. We ran experiments in a reference game framework and modeled the data with a pragmatic social inference Rational Speech Act model. Participants were able to infer listeners' preferences when analyzing their choice of objects given referential ambiguity. Moreover, a significant group of speakers were able to strategically choose ambiguous over unambiguous utterances in an epistemic, event-predictive, goal-directed manner, although a different group significantly preferred unambiguous utterances. We conclude that ambiguity resolution indeed reveals aspects of the knowledge, preferences, and beliefs of conversation partners and some of us are able to strategically use ambiguous utterances to gain knowledge about these aspects.

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

## 1 Introduction

Active inference—that is, the anticipatory, goal-directed, and epistemic invocation of behavior—is closely linked to the predictive mind perspective (Friston

et al., 2015; Hohwy, 2013; Clark, 2016). The anticipatory nature of the human mind reveals itself in many domains. With respect to planning and executing manual sensorimotor interactions, it has been shown that we anticipate future events and event boundaries, revealing anticipatory, event-predictive active inference processes (Belardinelli, Stepper, & Butz, 2016; Belardinelli, Lohmann, Farnè, & Butz, 2018; Friston et al., 2015; Hayhoe, Shrivastava, Mruczek, & Pelz, 2003; Lohmann, Belardinelli, & Butz, 2019). Also in the language domain, active inference processes seem to continuously unfold (Christiansen & Chater, 2016), compressing information into event-like units of thought (Baldwin & Kosie, to appear; Gärdenfors, 2014). For example, neurophysiological data has shown that listeners predict the semantic category of upcoming words (Federmeier & Kutas, 2002). Moreover, the inference process takes the structural properties of sentences into account (Levy, 2008). Dynamic language models show that complex, event-predictive structures guide ambiguity resolution during comprehension and likely also constrain ambiguity generation during language production (Elman & McRae, 2019).

When systematic abstractions become relevant, event-predictive biases seem to be at play, invoking the tendency to compress sensorimotor experiences, including language, into event-predictive encodings (Baldwin & Kosie, to appear; Butz, 2016, 2017; Shin & DuBrow, to appear). Various disciplines associated with cognitive science suggest that our minds develop event-compressed predictive encodings, which are recruited during decision making and action generation, including language production and comprehension, essentially determining thought itself in a highly active, epistemic, goal-directed manner (Baldwin & Kosie, to appear; Shin & DuBrow, to appear; Elsner & Adam, 2019; Knott & Takac, to appear; Ünal, Ji, & Papafragou, to appear; Stawarczyk, Bezdek, & Zacks, 2019). Here, we reveal socially epistemic inferences and utterance productions in scenarios where we observe and actively generate social event-predictive interactions.

In two main studies, we show how speakers update predictive models of the listener’s preferences and beliefs when watching social event interactions, such as when offering a few objects to choose from and observing the object choice of the conversation partner. We thus show that humans can interpret behavior of other people as driven by their motives, intentions, or personal characteristics. Conceptually, this idea goes back to the attribution theory (Jones & Davis, 1965; Kelley, 1967; Kelley & Stahelski, 1970). More recently, Shafto, Goodman, and Frank (2012) developed a Bayesian model of learning that formalizes the process of inferring others’ knowledge about the world based on their actions and goals. They argue that efficient learning is possible if we assume that agents’ actions are driven either by physical (non-social) or communicative goals, but are crucially not random. The authors show that an observer can draw stronger inferences concerning an underlying hypothesis when the acting agent has a communicative goal. The developed model predicts that learners use knowledge of agents’ goals to evaluate how knowledgeable they are, and, as a consequence, how much a learner can trust their actions to be informative about a hypothesis.

While our model also pursues Bayesian inference, or “psychological reasoning”, we do not focus on the inference of the actor’s knowledge, that is, on *learning from others* (Shafto et al., 2012). Rather, we focus on *learning about others*, that is, learning about listeners’ preferences when observing their disambiguating behavioral responses. We explore interpretive choices and the potential strategic, socially epistemic usage of ambiguous utterances in anticipation of actors’ responses. To formalize our hypothesis, we adapt the Rational Speech Act model framework, reliably modeling the involved, probabilistic interpretation processes and socially epistemic action choices. 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 the following, we first review how different disciplines approach ambiguity in natural language and communication, and provide a computational background on referential ambiguity resolution. 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 possibility of learning about the preferences of the conversation partner. Sections 4 and 5 give the results of behavioral experiments and 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 Ambiguity in natural language

### 2.1 Theoretical approaches

If a speaker and a listener understand an ambiguous utterance differently, communication between them might fail. On rare occasions, such communication failure can even be deadly: Pinker (2015) alludes to the Charge of the Light Brigade during the Crimean War as an example of a military disaster that was caused by vague orders. He also mentions how poor wording on a warning light was responsible for the nuclear meltdown at Three Mile Island. Finally, citing Cushing (1994), Pinker describes how the deadliest plane crash in history resulted from pilots and air traffic controllers arriving at different interpretations of the phrase “at takeoff”.

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 (Grice, 1975; Chomsky, 2002). The attitude towards ambiguity has been quite different in other disciplines, in part because the term itself can refer to multiple phenomena. 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 sporting implement. 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. In the case of referential ambiguity, an ambiguous utterance may apply to several possible referents in a scene. For example, a pronoun can be referentially ambiguous if there are multiple potential antecedents in the context. It is the latter type of ambiguity that we are concerned with in this paper.

If we look back at the study of ambiguity, we notice that the strategy of ambiguity avoidance is much older than the pronouncements by modern linguists. Greek and Latin rhetoricians believed that a skillfully-written text allows for a perfectly accurate and lossless transmission of meaning to the listener or reader (Ossa-Richardson, 2019); such a text avoids ambiguities.

Still, despite the teachings of classical philologists, authors continued to create ambiguous texts and readers were faced with the challenge of interpreting them. The Bible is one of the most significant of such texts. In the sixteenth century, the Catholic church responded to the Reformation by proposing that the Bible can contain multiple meanings—Ossa-Richardson (2019) equates these meanings with multiple paths that lead readers to God. In a sense, this proposal contained one of the first acknowledgments of the virtue of ambiguity, though with an important caveat: only God could introduce ambiguity, humans should not.

The search for efficient transmission of meaning has rested on an important assumption: we communicate to transfer knowledge to our conversation partner. It is the efficiency of this transfer that many recent experiments were designed to evaluate. To be more precise, communication was considered efficient if an experimental participant could follow instructions precisely. Yet, ordering actions and following instruction are probably not the most common types of communicative acts (Foppa, 1995), and information-seeking might not be the only communicative task in which we engage (Markova & Graumann, 1995).

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 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 this work, we thus focus on the effects of resolving, or anticipating the res-

olution of, ambiguous utterances, modeling the involved probabilistic inference processes. The main contributions of this paper are two-fold: first, we demonstrate that participants are indeed able to infer hidden beliefs of their conversation partners observing their choices; second, we show that some speakers can actively create situations of uncertainty anticipating the epistemic value when observing the consequent referent choice. We formalize the human communicative behavior in a probabilistic Bayesian model, which approximates the dynamically unfolding reasoning processes, including limits thereof.

## 2.2 Computational modeling

In search of the communicative purpose of ambiguous language, the current work identifies an additional benefit: the *extra* information we gain from observing how listeners resolve ambiguity. We show that language users learn about each other’s private knowledge when observing how ambiguity is resolved. When utterances leave room for interpretation, listeners must draw on their opinions, beliefs, and preferences to fill in the gaps; by observing the concrete interpretation, speakers thus learn about the opinions, beliefs, and preferences of their conversation partner. As a result, in a naturalistic conversation, where speakers take turns, ambiguous utterances open interpretation spaces and the resulting interpretation choices dynamically and mutually reveal individual opinions, beliefs, and preferences.

By way of illustration, take the scenario in Figure 1. Suppose a speaker produces the single-word utterance “blue” – meaning: choose a blue object – creating referential ambiguity for the listener, that is, offering a choice between a blue square and a blue circle. Suppose further that, upon hearing “blue”, 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 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 a Rational Speech Act (RSA) modeling framework.

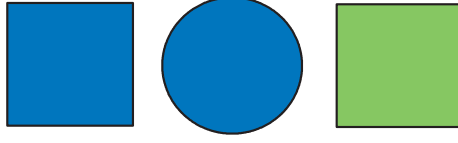


Figure 1: A simple reference game scenario from Frank and Goodman (2012). In the game, speakers are confronted with a collection of objects, which determine the current scenario  $S$ , where  $S = \{\text{solid blue square}, \text{solid blue circle}, \text{solid green square}\}$  in the depicted example. 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:  $U = \{\text{“solid”}, \text{“blue”}, \text{“green”}, \text{“square”}, \text{“circle”}\}$ .

### 2.3 Original RSA Formalization

RSA (cf. Frank & Goodman, 2012; Franke & Jäger, 2016; Goodman & Frank, 2016) formalizes a state space, or scenario,  $S$  in the form of a particular set of objects (cf. the example in Figure 1). Moreover, RSA unfolds computations over the corresponding utterance space  $U$ , which consists of the set of possible utterances, which in turn contains all object features that are present in a particular scenario  $S$ . 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 contain 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, the speaker  $S_1$  observes the state  $S$  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  $S$ , which is determined by the log-likelihood of this particular object choice  $U_{S_1}(u; s)$ :<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:

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

<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.



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 led  $S_1$  to utter  $u$  given the scenario  $S$ :

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

Frank and Goodman (2012) tested the predictions of RSA 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 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).

### 3 Pragmatic social inference RSA model

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 (PSIRSA). By changing  $P(s)$  to a non-uniform distribution, we essentially model prior beliefs of which object the speaker is more likely to refer to, or—when viewed from a more self-centered perspective—which prior object 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 maxi-

mize 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. fullPSIRSA assumes the deep reasoning process integrating the full RSA formalism. It 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, fullPSIRSA includes modeling 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, simplePSIRSA 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, simplePSIRSA ignores any indirect pragmatic reasoning considerations about which object the speaker may refer to given an utterance and a particular object constellation. It simply assumes that all objects may be chosen that match the utterance, modifying these choice options dependent on the feature-preference-dependent object choice priors. 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 the evaluation section below, we compare the modeling performance of fullPSIRSA with simplePSIRSA.

## 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 PSIRSA—that is, by recursive, Bayesian inference? 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		

[Continue](#)

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.

### 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 subjects that they agree 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.

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 PSIRSA’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 prediction values 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 with respect to listener preferences and five trials were uninformative with respect to listener preferences (e.g., observing that the listener chose one of three identical objects).

### 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:

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

where  $P_{data}(f | u, s)$  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  $s$ ;  $P_{model}(f | u, s)$  specifies the corresponding model posterior, either  $P_{S_2}(f | u, s)$  in the case of fullPSIRSA or  $P_{S_{1-simp}}(f | u, s)$  in the case of simplePSIRSA. By minimizing the summed KL divergence between the empirical and model-predicted preference posteriors over all considered trials, we essentially 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 fullPSIRSA and two for simplePSIRSA. The softmax scaling factor  $\alpha$  is only relevant for fullPSIRSA; 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 fit 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  $u$ .

## 4.4 Results

### 4.4.1 Models with global optimization

We fit the following free parameters to optimize the predictions of the models. First, the full model includes a “greediness” parameter  $\alpha$  that controls how likely it is that speakers choose the best-suited utterance to signal a particular object to a listener. This parameter is absent in the simple model since it relies on fewer layers of reasoning. The second parameter  $\gamma$  controls how soft the preferences are. Hard preferences enforce the choice of the preferred object type, while increasing softness converges towards no object preference. Similarly, the obedience parameter  $\beta$  allows subjects to choose objects that do not qualify for the utterance. As for the preference parameter  $\gamma$ , the  $\beta$  range includes hard obedience on the one side of the spectrum – for example, definitely choosing a blue object when hearing “blue” – and full ignorance of the utterance at the other extreme, choosing uniformly from all available objects.

simplePSIRSA and fullPSIRSA with softness ( $\gamma$ ) optimized globally provide nearly identically good fits to the data (Figure 3). Simple linear regression analysis was used to test whether the model values predicted the human data. simplePSIRSA yields a value of  $r^2 = 0.8607^3$  ( $F(1, 190) = 1181$ )<sup>4</sup> when only softness parameter  $\gamma$  is optimized ( $\gamma = 0.2204$  after optimization). When both parameters are optimized globally, a variance estimate of  $r^2 = 0.9788$  ( $F(1, 190) = 8823$ ) is reached ( $\gamma = 0.2210$  and  $\beta = 0.2693$  after optimization), indicating that participants indeed considered (possibly subconsciously) the option to interpret utterances non-literally. fullPSIRSA yields nearly identical values. When optimizing only the softness parameter  $\gamma$ , a value of  $r^2 = 0.8568$  ( $F(1, 190) = 1144$ ) is reached ( $\gamma = 0.2231$ ). Optimizing both,  $\alpha$  and  $\gamma$ , a value of  $r^2 = 0.8607$  ( $F(1, 190) = 1144$ ) is reached ( $\alpha = 0.1797$ ,  $\gamma = 0.2205$ ). When optimizing all three parameters, fullPSIRSA yields a value of  $r^2 = 0.9772$  ( $F(1, 190) = 8170$ ) ( $\alpha = 0.2657$ ,  $\gamma = 0.2214$ ,  $\beta = 0.0030$ ).

Overall, the results show that participants are indeed able to infer the feature preferences that lead to the choice of an object. Moreover, the higher model flexibility of fullPSIRSA—controlled via parameter  $\alpha$ —does not yield any modeling

<sup>3</sup>Here and throughout the paper we report adjusted  $r^2$  values.

<sup>4</sup>All results were significant at  $p < 0.001$  level if not stated differently in the text.

improvement, implying that an approximation of the more shallow reasoning process modeled by simplePSIRSA typically unfolded in the minds of the participants.

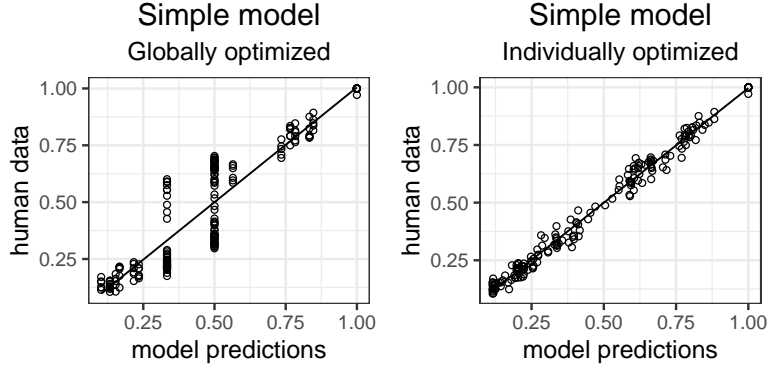


Figure 3: Human data from Experiment 1 plotted against the predictions of simplePSIRSA. Each data point indicates the slider values and model predicted feature preference posteriors for a particular ambiguity class. Left panel:  $\gamma$  *optimized globally* ( $r^2 = 0.8614$ ); right panel:  $\gamma$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

#### 4.4.2 Individually-fitted models

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 the light of the KL divergence between the individual participants' slider value choices and the corresponding model predictions for PSIRSA. We then again averaged the individualized model prediction values and participants' slider values with respect to the particular ambiguity classes and calculated correlations between the data and the model.

The full model optimized at the individual level for the additional parameter  $\alpha$  does not improve the fit compared to the simplified model (simplePSIRSA:  $r^2 = 0.8631$ ,  $F(1, 190) = 1205$ ; fullPSIRSA:  $r^2 = 0.8627$ ,  $F(1, 190) = 1201$ ). Seeing that both models again fit the data nearly equally well (if anything, simplePSIRSA performs slightly better), we only consider the predictions of simplePSIRSA henceforth. Note further that the individually-fitted parameters do not improve the correlation values much, if at all, when compared to the globally-fitted model.

The model fit improves considerably when we additionally fit the obedience parameter  $\beta$  at the individual level. Here the model explains a large proportion of variance in the human judgments ( $r^2 = 0.9919$ ,  $F(1, 190) = 23480$ ). The likelihood ratio test (two-tailed) revealed that a  $\gamma$ - and  $\beta$ -optimized simplePSIRSA model pro-

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

To appreciate the gains obtained by fitting model parameters, Figure 4 shows the average responses of the human participants and of the individually-, two-parameter-optimized simplePSIRSA model and the non-optimized simplePSIRSA 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 can definitely not 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.

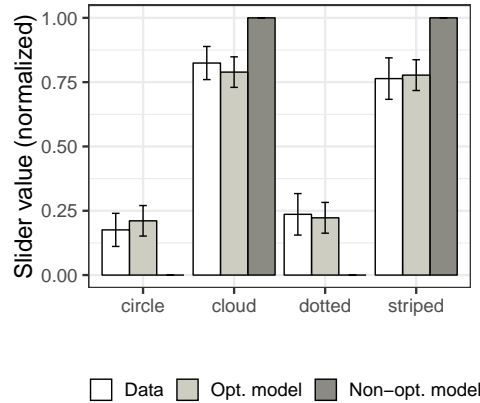


Figure 4: Human data and simplePSIRSA’s (individually-, two-parameter-optimized and non-optimized) feature preference posterior estimates for the scenario  $S$  shown in Figure 2. Error bars represent 95% confidence intervals.


We thus find strong empirical support for simplePSIRSA, implying that speakers are indeed able to use listener behavior to acquire information about their preferences. We fail to find that the fullPSIRSA model predicts the data better. This result suggests that the task in our experiments does not require full-blown prag-




matic 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 aware that ambiguous language is potentially more informative and can thus use ambiguous language in a socially epistemic, strategic manner?

## 5 Experiment 2: Choosing utterances to learn about others

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







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*.

## 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 the subjects 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.”

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 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 prediction values 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 simplePSIRSA 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. equations 7 and 10). While a positive 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 simplePSIRSA with non-optimized parameters and with several parameter optimizations with the performance of a uniform baseline model, which simply chooses one of the available utterances at random. Seeing that in particular ambiguity cases with particular constellations  $S$  three 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 model is much better in capturing variance in the data than one would expect without this insight ( $r^2 = 0.7466$ ,  $F(1, 82) = 245.6$ ,  $p < 0.001$ ). Figure 6 compares this performance to the non-optimized simplePSIRSA, where we set the parameters to hard preference and obedience ( $\gamma = 0$ ,  $\beta = 0$ ) and the information gain factor to  $\lambda = 1$ , thus preferring to choose those utterances that are expected to yield high information gain. Surprisingly, this model captures very little variance in the human data ( $r^2 = 0.0595$ ,  $F(1, 82) = 6.253$ ,  $p < 0.05$ ).

To examine the reasons for this failure, we first performed additional global parameter optimization runs. When optimizing all simplePSIRSA parameters, the model accounts for more variance than the uniform base model ( $r^2 = 0.7991$ ,  $F(1, 82) = 331.2$ ,  $p < 0.001$ ; optimized model parameters:  $\gamma = 0.0006$ ,  $\beta = 0.2758$ ,  $\lambda = 0.3663$ ). Moreover, the nested model comparison test with three free parameters yields a  $G^2$  value of 13.6912, which indicates a more accurate model with  $p < 0.01$ . Figure 7 shows the correlation plot. The parameters indicate that the preference strength is rather high, obedience is not as strong, while the information gain intention is present. We now turn to individual parameter optimization,

suspecting that there may be fundamental differences between the individual participants.

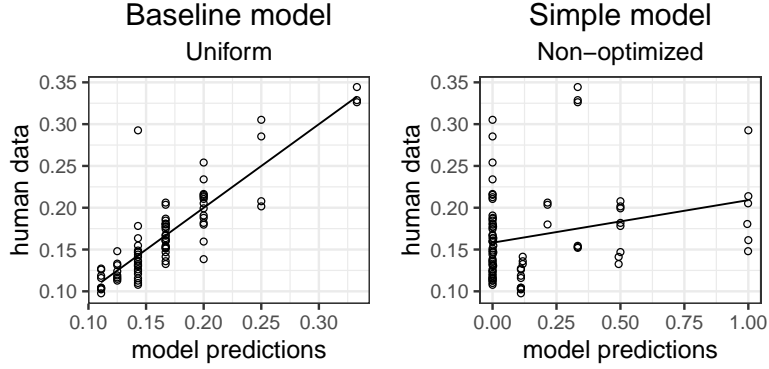


Figure 6: Average human data from Experiment 2 plotted against the predictions of the uniform baseline model and the simplePSIRSA model. Left panel: *uniform model* ( $r^2 = 0.7466$ ); right panel: *non-optimized simplePSIRSA* ( $r^2 = 0.0595$ ).

We compared three single-parameter-individually-optimized simplePSIRSA 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 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  $r^2 = 0.8077$ ,  $F(1, 82) = 349.6$ ). The comparison with the baseline model in terms of nested model statistics confirms that only the individual optimization of  $\lambda$  improves model performance ( $\lambda$ :  $G^2 = 268.88$ ,  $df = 82$ ,  $p < 0.001$ ;  $\gamma$ :  $G^2 = 31.38$ , n.s.;  $\beta$ :  $G^2 = 56.29$ , n.s.). Two- and three-parameter individual optimizations did not yield any significant model improvements when compared to the individually  $\lambda$ -optimized model (best improvement when optimizing  $\gamma$  in addition to  $\lambda$ :  $G^2 = 24.72$ ,  $df = 82$ , n.s.). Figure 7 shows the resulting correlation plot for  $\lambda$ -individually optimized model.

Unlike for Experiment 1, where even the non-optimized models provided a good linear fit to the data, individual optimization produces a large effect on the model predictions in Experiment 2. Figure 8 compares 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: in a situation with a striped green circle, a blue striped cloud, and a solid green cloud, uttering things like *cloud*, *striped*, or *green* (i.e., the utterances that point to more than one object in the scene) could let the speaker learn something about the listener’s preferences. 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.

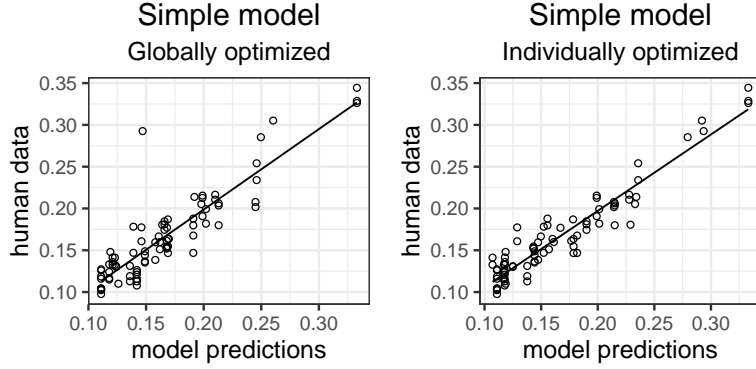


Figure 7: Average human data from Experiment 2 plotted against the predictions of optimized simplePSIRSA models. Left panel: *globally optimized 3 parameter model* ( $r^2 = 0.7466$ ; right panel: *individual KL-factor  $\lambda$ -optimized model* ( $r^2 = 0.9059$ ).

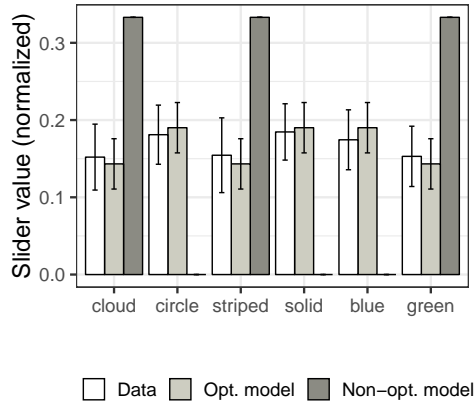


Figure 8: Simple Social Inference model predictions and human data for one of the classes of stimuli *Experiment 2: Picking utterances*. The optimized version of the model is optimized for the KL-factor  $\lambda$ . Error bars represent 95% confidence intervals.

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” group or “unpredictable” behaving group: for 28 participants, the KL divergence values of the  $\lambda$ -optimized simplePSIRSA 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 ( $G^2 = 180, 17$ ,  $df = 33$ ,  $p < 0.001$ ).

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 ( $G^2 = 102.16$ ,  $df = 21$ ,  $p < 0.001$ ).

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-inferences 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), the identical analysis yielded 42% participants that preferred ambiguous over unambiguous utterances (37% unpredictable participants; 21% preferred unambiguous utterances). In the second complementary study with 54 participants (2 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

We have found strong support that we can indeed learn about others when observing their interpretation of ambiguous utterances. The results of Experiment 1 demonstrate that naïve speakers are able to reason pragmatically about *why* listeners may take the actions they do. The success of our computational model PSIRSA in predicting the observed behavior 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. The results of Experiment 2 demonstrate that at least 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.

Currently, we are transferring the experimental setup to more naturalistic interaction scenarios. Even in these cases, though, it appears that we still find participants who consistently prefer to choose unambiguous utterances. Two explanations may be warranted and need to be investigated further. 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, it may simply be the case that these participants do not have access to the required deeper reasoning process, and thus prefer to give instructions with predictable outcomes.

Nonetheless, taken together, the results of our experiments and the success of PSIRSA in modeling these results 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—that is, they are able to learn about the hidden model states of others, including preferences but probably also other aspects of

beliefs.

It should also be noted that ambiguous utterances used in this way are closely related to questions, which may ask directly about considered preferences. Ambiguous utterances provide a ready but more subtle, indirect alternative to asking directly. In normal conversations, a speaker might favor the indirect route, given considerations of politeness and possibly also in an effort to keep the conversation open. With ambiguous language, the conversation partner can choose to disambiguate the ambiguous utterance or, alternatively, choose to continue in a different direction or even change topic.

We note that the analyzed preference prior, viewed from a broader perspective, can be closely related to a part of the event-predictive mind of the listener and the speaker (Butz, 2016; Butz & Kutter, 2017). When interpreting an utterance—in our case, 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; Friston et al., 2015). It 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.

More generally, predictive states of mind about others do not only include considerations of preferences, but may also concern all imaginable knowledge, opinions, beliefs, and current trains of thought of the listener. Moreover, during a conversation, the involved “social” priors will dynamically develop depending on the internal predictive models and the generated utterances, actions, and responses of the speaker and listener. The priors dynamically depend on the privileged grounds of the conversational partners, and also on the common ground in which the conversation unfolds. In that sense, ambiguous utterances and resolutions thereof are one device for projecting parts of each other’s privileged grounds into the common ground.

## **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

- Baldwin, D. A., & Kosie, J. E. (to appear). How does the mind render streaming experience as events? *Topics in Cognitive Science*.
- 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.
- 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., & Kutter, E. F. (2017). *How the mind comes into being: Introducing cognitive science from a functional and computational perspective*. Oxford, UK: Oxford University Press.
- Chomsky, N. (2002). An interview on minimalism. In A. Belletti & L. Rizzi (Eds.), *On nature and language* (p. 92-161). Cambridge: Cambridge University Press.
- Christiansen, M. H., & Chater, N. (2016). The now-or-never bottleneck: A fundamental constraint on language. *Behavioral and Brain Sciences*, 39, 1-18. doi: 10.1017/S0140525X1500031X
- Clark, A. (2016). *Surfing uncertainty: Prediction, action and the embodied mind*. Oxford, UK: Oxford University Press.
- Cushing, S. (1994). *Fatal words: Communication clashes and aircraft crashes*. Chicago: University of Chicago Press.
- Elman, J. L., & McRae, K. (2019). A model of event knowledge. *Psychological Review*, 126, 252-291. doi: 10.1037/rev0000133
- Elsner, B., & Adam, M. (2019). Infants' prediction of action-events for human and non-human agents. *Topics in Cognitive Science*. (this volume)
- Federmeier, K. D., & Kutas, M. (2002). Picture the difference: Electrophysiological investigations of picture processing in the two cerebral hemispheres. *Neuropsychologia*, 40(7), 730–747.
- 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.
- Foppa, K. (1995). On mutual understanding and agreement in dialogues. In I. Markova & F. K. Graumann Carl F. (Eds.), *Mutualities in dialogue*. Cambridge, UK: Cambridge University Press.
- 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.



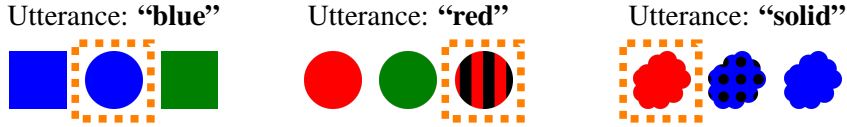
- Franke, M., & Jäger, G. (2016). Probabilistic pragmatics, or why Bayes' rule is probably important for pragmatics. *Zeitschrift für Sprachwissenschaft*, 35(1), 3–44.
- Friston, K., Rigoli, F., Ognibene, D., Mathys, C., Fitzgerald, T., & Pezzulo, G. (2015). Active inference and epistemic value. *Cognitive Neuroscience*, 6, 187–214. doi: 10.1080/17588928.2015.1020053
- Gärdenfors, P. (2014). *The geometry of meaning: Semantics based on conceptual spaces*. Cambridge, MA: MIT Press.
- 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: 10.1167/3.1.6
- Hohwy, J. (2013). *The predictive mind*. Oxford, UK: Oxford University Press.
- 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.
- Kelley, H. H. (1967). Attribution theory in social psychology. In *Nebraska symposium on motivation*.
- 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.
- Knott, A., & Takac, M. (to appear). Roles for event representations in sensorimotor experience, memory formation and language processing. *Topics in Cognitive Science*.
- Levinson, S. C. (2000). *Presumptive meanings: The theory of generalized conversational implicature*. Cambridge, MA: MIT Press.
- Levy, R. (2008). Expectation-based syntactic comprehension. *Cognition*, 106(3), 1126–1177.
- Lewandowsky, S., & Farrell, S. (2011). *Computational modeling in cognition: Principles and practice*. Thousand Oaks: Sage Publications.
- 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
- Markova, I., & Graumann, F. K., Carl F. (1995). Preface. In I. Markova & F. K. Graumann Carl F. (Eds.), *Mutualities in dialogue*. Cambridge, UK: Cambridge University Press.
- Ossa-Richardson, A. (2019). *A history of ambiguity*. Princeton University Press.
- Piantadosi, S. T., Tily, H., & Gibson, E. (2012). The communicative function of ambiguity in language. *Cognition*, 122, 280–291.
- Pinker, S. (2015). *The sense of style: The thinking person's guide to writing in the 21st century*. Penguin Books.

- 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.
- 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.
- Shin, Y. S., & DuBrow, S. (to appear). Structuring memory through inference-based event segmentation. *Topics in Cognitive Science*.
- 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
- Stawarczyk, D., Bezdek, M. A., & Zacks, J. M. (2019). Constructing event representations: The role of the midline default network core. *Topics in Cognitive Science*. doi: doi.org/10.1111/tops.12450
- Ünal, E., Ji, Y., & Papafragou, A. (to appear). From event representation to linguistic meaning. *Topics in Cognitive Science*.
- Wasow, T. (2015). Ambiguity avoidance is overrated. In S. Winkler (Ed.), *Ambiguity: Language and communication* (p. 29-47). de Gruyter.

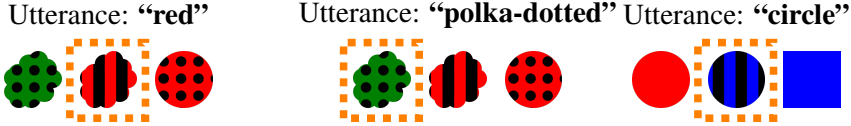
## A Ambiguity classes

### Experiment 1

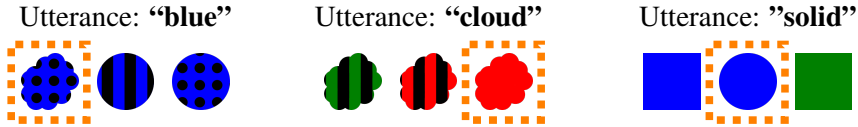
Figure 9 shows three exemplar scenarios for three representative ambiguity classes. Let us consider the first class in more detail. In the scenario  $S$  on the left side of Figure 9a, the utterance “blue” refers either to the blue square or the blue circle. The picked object, that is, the blue circle, is unique in its shape (circle) and shares the other non-referenced property with both other objects (that is, its plain 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 references two objects, the picked object has one non-referenced unique feature, while the other, non-referenced feature is shared amongst 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 the

color (blue), and the other two share the shape (square). As a result, uttering *green* or *circle* will give no choice 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 the shape (circle), two share the pattern (solid), and the other two share the 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 object again share the shape (cloud), two share the pattern (solid), while the other two share the 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 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.