

Learning about Others: Pragmatic Social Inference through Ambiguity Resolution

December 20, 2019

Abstract

Grice’s maxim of manner postulates that ambiguity should be avoided to maximize clarity. Nonetheless, ambiguity is ubiquitous during conversations. It has been suggested that some ambiguity may be useful for efficiency reasons in cases when clarity is not affected. Here we investigated whether disambiguations of ambiguous utterances yield another, socially-relevant benefit. In particular, we asked if responses to ambiguous utterances can reveal parts of the internal model of the interpreter. More concretely, we asked if speakers (i) can use responses as a source of information to infer unknown preferences of their conversation partner, and (ii) are able to strategically choose ambiguous utterances over unambiguous utterances for learning about their conversation partner’s preferences. We ran two web-based experiments using a modified version of the original reference game framework (Frank & Goodman, 2012) and modeled the recorded data by modifying the vanilla Rational Speech Act model. The data and modeling results confirm both points. Participants were able to infer Bayesian posteriors of listeners’ preferences when analyzing their choice of objects in situations of referential ambiguity. Moreover, nearly 40% of the speakers were able to strategically choose ambiguous over unambiguous utterances in an epistemic, goal-directed manner, maximizing expected information gain about the listeners’ preferences. Surprisingly, an equally-large number of participants seemed to minimize expected information gain by systematically choosing unambiguous utterances. Our results thus show that ambiguity resolution can reveal aspects of the knowledge, preferences, and beliefs of conversation partners, and that 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

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 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). In the language domain, predictive active inference processes seem to continuously unfold (Christiansen & Chater, 2016), compressing information into event-like units of thought (Gärdenfors, 2014). For example, listeners predict the semantic category of upcoming words (Federmeier

& Kutas, 2002), as evidenced by neurophysiological data. Comprehension of sentences relies not only on the ability of listeners to anticipate subsequent words based on their transitional probabilities, but also takes into account the structural properties of sentences, revealing an even more abstract level of predictions (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 processes seem to be at play, compressing sensorimotor experiences, including language, into event-predictive encodings (Butz, 2016, 2017). Various disciplines associated with cognitive science suggest that our predictive minds develop event-compressive, predictive encodings, which interact with action, including language production and comprehension, essentially determining thought itself in a highly active, epistemic, goal-directed manner (Baldwin & Kosie, 2019; Shin & DuBrow, 2019; Elsner & Adam, 2019; Storck, Ehrlis, & Fallgatter, 2019; Knott & Takac, 2019; Ünal, Ji, & Papafragou, 2019; Stawarczyk, Bezdek, & Zacks, 2019). Here, we reveal socially epistemic comprehension and utterance productions, while observing and generating 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 giving a communicative goal of an agent allows the observer to draw a stronger inference concerning the underlying hypothesis. The 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—that is, 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*, learning about actors’ 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 choice. Interestingly, the modeling work reveals good interpretive abilities but also strong individual differences when the task is to choose (ambiguous) utterances strategically for gaining social knowledge.

We use ambiguity resolution as a paradigm in which learning about others is possible. Intuitively, ambiguity should make understanding each other difficult. 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.

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).

In search of the communicative purpose of ambiguous language, the current work identifies an additional benefit in using such language: the *extra* information we gain from observing how our listeners resolve ambiguity. We propose that language users learn about each other’s private knowledge by observing how they resolve ambiguity. If language does not do the job of specifying the information necessary for full interpretation, then listeners are left to draw on their opinions, beliefs, and preferences to fill in the gaps; by observing how listeners fill those gaps in, speakers learn about the opinions, beliefs, and preferences of their listeners. In a dynamic, naturalistic conversation, speakers can take turns choosing ambiguous statements in order to leave room for their partner to fill the missing information in, thereby revealing opinions, beliefs, and preferences.

By way of illustration, take the scenario in Figure 1. Suppose a speaker produces the single-word utterance “blue” in an attempt to signal one of the objects to a listener. The utterance is referentially ambiguous; the listener can choose either the blue square or the 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.

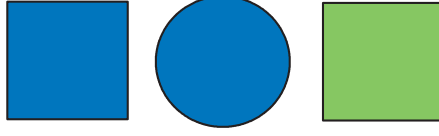


Figure 1: A simple reference game scenario from Frank and Goodman (2012). In the game, speakers are confronted with a collection of objects. The speaker chooses a single-word utterance to signal one of the objects to a listener.

However, accessing this added information requires the speaker to reason pragmatically about the pragmatic reasoning of the listener—a higher-order pragmatic reasoning, as it were. 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 within the Rational Speech Act (RSA) modeling framework (cf. Franke & Jäger, 2016; Goodman & Frank, 2016).

At the base of the reasoning process, there is a hypothetical, naïve literal listener L_0 , who hears an utterance u and attempts to infer the object 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:

$$P_{L_0}(s | u) \propto \llbracket u \rrbracket(s),$$

essentially returning a uniform distribution over those objects in S that contain the uttered feature u .¹

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 situation S , which is determined by the log-likelihood of this particular object choice $U_{S_1}(u; s)$ ²:

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

Depending on a “greediness” factor α , the speaker chooses a particular utterance u with a probability that is exponentially proportional to the utility estimates:

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

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 S_1 . In other words, L_1 reasons about the s that would have been most likely to lead S_1 to utter u :

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

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

²The 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.

Frank and Goodman (2012) tested the predictions of their model against behavioral data from reference games, as in Figure 1. To model production behavior (i.e., 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 . Finding extremely strong correlations between model predictions and behavioral data in both cases, Frank and Goodman have strong support for their model of pragmatic reasoning in reference games (see also Qing & Franke, 2015, for a fuller exploration of the modeling choices).

The current paper builds on the foundational, vanilla RSA model of reference games by introducing uncertainty about the prior beliefs of the listener and modeling a speaker who reasons about these beliefs on the basis of and in anticipation of the observed referent choice.

Results

The 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, and some speakers can actively create situations of uncertainty that can lead to learning. Second, we formalize the human communicative behavior in a probabilistic Bayesian reasoning model.

Our model is a modified version of the vanilla RSA model (Frank & Goodman, 2012). It formalizes a state space S in the form of a particular set of objects (cf. the example in Figure 1) and an utterance space U , which consists of the set of possible utterances; when conceived of as a set of single-word utterances, U amounts to the set of all features present in S . Moreover, the model specifies priors or posteriors over referenced objects, object choices, utterance preferences, and utterance choices. For notational convenience, we denote a particular object choice of the listener by s in S . RSA then models a recursive social reasoning processes, incorporating several levels of probabilistic inference.

Modeling: Pragmatic social inference RSA

Our model builds on the vanilla version of RSA presented above, 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 evaluate two versions of the model. fullPSIRSA assumes a rather deep reasoning process. fullPSIRSA essentially 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 in Figure 1, 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 choice for the speaker.

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 models (Sikos, Venhuizen, Drenhaus, & Crocker, 2019). simplePSIRSA removes this reasoning about alternative utterances and allows the pragmatic speaker to directly tap into the (expected) interpretation of L_0 , directly augmenting the literal listener’s choice likelihoods with the feature-preference-dependent object prior $P(s | f)$:

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

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

$$P_{S_1\text{-simp}}(f | u, s) \propto P_{L_0\text{-simp}}(s | u, f) \cdot P(f).$$

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 also simplifies the reasoning process:

$$P_{S_1\text{-simp}}(u) \propto \sum_{s: \llbracket u \rrbracket(s)=1} P_{L_0}(s | u, f) \exp(\lambda \cdot \text{KL}(P(f) || P_{S_1\text{-simp}}(f | u, s))),$$

In the evaluation section below, we compare the modeling performance of fullPSIRSA with simplePSIRSA.

To test our model, we need to evaluate two main predictions of PSIRSA: first, the pragmatic speaker’s inference about the listener’s feature preferences on the basis of observed object choices in particular situations $P_{S_2}(f | u, s)$; second, the pragmatic speaker’s strategic utterance selection $P_{S_2}(u)$ in light of the anticipated information gain about the listener’s preferences considering the possible object choices.

Experiment 1

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 generative modeling? A sample trial is shown in Figure 2.

Models with global optimization

We fit the following free parameters to optimized the predictions of the models. First, the full model includes a “greediness” parameter α that controls how likely participants are to choose an optimal utterance to signal an object to a speaker. This parameter is absent in the simple model since it relies on fewer layers of reasoning. The second parameter γ controls how soft the preferences are: if preferences are hard, participants will use values at the ends of the scale, while soft preferences will push the slider values towards the middle. Finally, the obedience parameter β allows subjects to consider

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	<input type="range"/>		clouds	<input type="range"/>	
striped things	<input type="range"/>		circles	<input type="range"/>	
polka-dotted things	<input type="range"/>		squares	<input type="range"/>	

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 led her to the particular choice of object. They specify their inference by adjusting the sliders for each of the features.

objects that do not qualify for the utterance, for example, subjects will consider objects that are not blue upon hearing “blue”.

We first present results of the globally-optimized versions of PSIRSA (Figure 3). Here, both simplePSIRSA and fullPSIRSA with softness (γ) optimized globally provide nearly-identically good fits to the data. simplePSIRSA yields a correlation of $r^2 = 0.8614$ when only softness parameter γ is optimized ($\gamma = 0.2204$ after optimization).³ When both parameters are optimized globally, a correlation of $r^2 = 0.9789$ 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. Optimizing only the softness parameter γ , a correlation of $r^2 = 0.8576$ is reached ($\gamma = 0.2231$). Optimizing both, α and γ , a correlation of $r^2 = 0.8614$ is reached ($\alpha = 0.1797$, $\gamma = 0.2205$). When optimizing all three parameters, fullPSIRSA yields a correlation of $r^2 = 0.9773$ ($\alpha = 0.2657$, $\gamma = 0.2214$, $\beta = 0.0030$), not quite reaching the correlation of simplePSIRSA, which may be due to some subtle interactions between parameters α and β . Overall, the results show that participants are indeed able to infer the feature preferences that lead to the choice of an object, and, notably, simplePSIRSA models this inference process very well. The higher model flexibility of fullPSIRSA—controlled via parameter α —does not yield any modeling improvement.

Individually-fitted models

We now compare our two model variants further when fitting the parameters to the individual data of each participant separately. We optimized α and γ in light of the

³All correlations were highly significant, that is, $p < 0.001$, if not stated differently in the text.

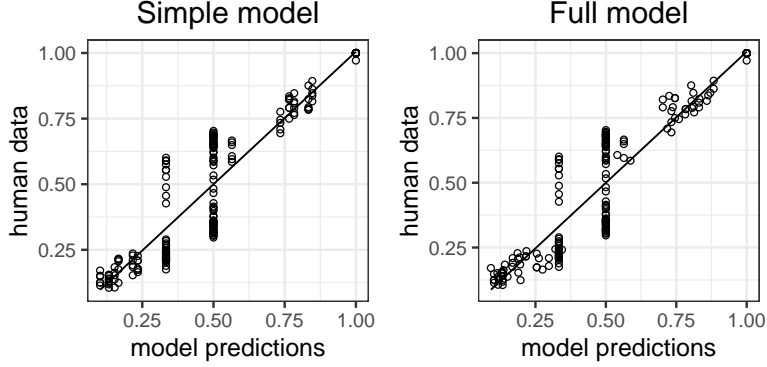


Figure 3: Human data from Experiment 1 plotted against the predictions of simplePSIRSA and fullPSIRSA with γ *optimized globally*. Each data point indicates the slider values and model predicted feature preference posteriors for a particular ambiguity class. Left panel: *simple model* ($r^2 = 0.8614$); right panel *full model* ($r^2 = 0.8576$).

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 α does not improve the fit compared to the simplified model (simplePSIRSA: $r^2 = 0.8631$; fullPSIRSA: $r^2 = 0.8614$). Seeing that both models fit the data nearly equally well (if anything, simplePSIRSA performs slightly better), we will henceforth only consider the predictions of simplePSIRSA. 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 if we additionally fit the obedience parameter β at the individual level. Here the strongest positive correlation between the human judgments and model predictions ($r^2 = 0.992$) can be observed. The likelihood ratio test revealed that a γ - and β -optimized simplePSIRSA model provides a better fit compared to a model optimized only for γ ($G^2 = 237.36, df = 82, p < 0.01$). The more complex model contains one additional parameter β 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. Figure 4 shows that the resulting cross-validated model predictions retain the strong correlation ($r^2 = 0.9901$).

To appreciate the gains obtained by fitting model parameters, Figure 5 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 the uttered, thus sliders

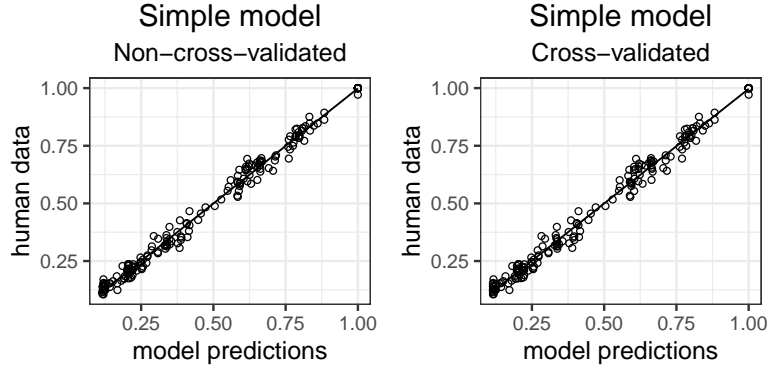


Figure 4: Human data from Experiment 1 plotted against the predictions of the *individually* β - and γ -optimized simplePSIRSA model. Left panel: *non-cross-validated* ($r^2 = 0.992$); right panel: *cross-validated* ($r^2 = 0.9901$).

for those features were not present. As Figure 5 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 results; optimizing β and γ improves the quantitative fit.

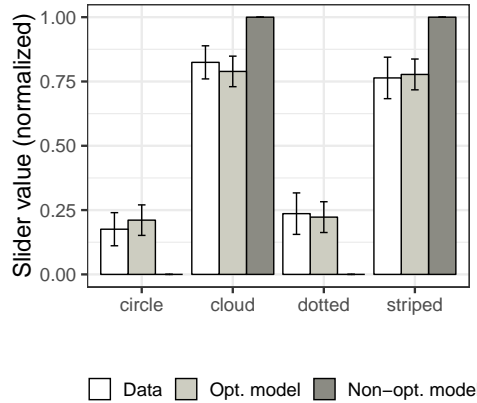



Figure 5: Human data and simplePSIRSA’s (individually-, two-parameter-optimized and non-optimized) feature preference posterior estimates for the scenario S shown in Figure 2.

We thus find strong empirical support for simplePSIRSA: 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 pragmatic inference about alternative utterances, in contrast with Frank and Goodman (2012) but in line with Sikos et al. (2019).


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?

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, are participants able to reason pragmatically about the utility of ambiguous utterances in informing listener preferences? Figure 6 shows a sample trial.

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"		

[Continue](#)

Figure 6: A sample trial from *Experiment 2: Choosing utterances*.

By reasoning about the predictions of S_2 , we are able to use simplePSIRSA to compute the expected most informative utterance with respect to 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 in our reference game scenario.

To generate predictions from $P_{S_1\text{-simp}}(u)$, a total of three free parameters can be identified. We consider different values for γ (i.e., preference softness) and obedience β , as well as for the λ parameter, which factors the importance of choosing the expected most informative utterance with respect to the determined KL divergence values. Note that when allowing negative values for λ , negated information gain essentially minimizes expected information gain. Thus, negative values for λ yield a model that favors unambiguous utterances. Moreover, when $\lambda = 0$, the model collapses to a uniform distribution over the available utterances.

Figure 7 shows the model fits of the non-optimized simplePSIRSA and the simplePSIRSA model with all three parameters optimized globally. Again, we optimized values with respect to KL divergence estimates between the participant data and the model predictions—in this case, for utterance preference distributions. Both models fail to predict the human data ($r^2 = 0.0709$, $p = 0.01439$ [gcs: we need another r^2 value (I think Martin is working on this)]).

Seeing that global optimization does not yield good fits, we moved on to individual optimization. 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 γ , obedience β , or KL-factor λ being optimized. The results indicate that we get the best fit by optimizing the KL-

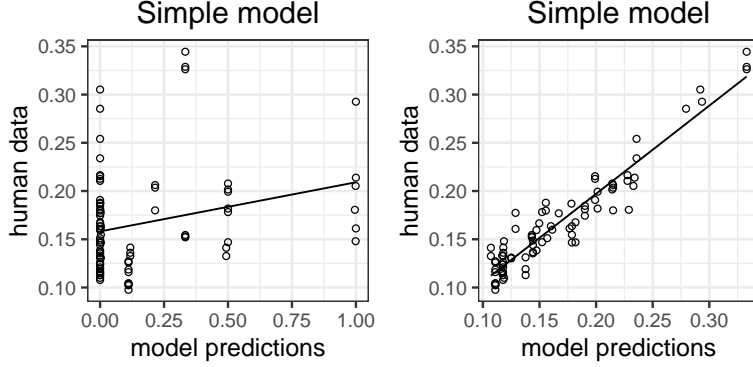


Figure 7: Average human data from Experiment 2 plotted against the predictions of the non-optimized and three-parameter-optimized simplePSIRSA models. [gcs: can we get more informative figure titles that specify non-optimized vs. optimized?] [gcs: it looks like we have the wrong figure on the right side for the globally-optimized model]

factor λ ($r^2 = 0.9071$; leave-one-out cross-validated optimization $r^2 = 0.8902$), with other models capturing less variance in the data (β -optimized $r^2 = 0.8039$; γ -optimized $r^2 = 0.8100$). Two- and three-parameter individual optimizations were unstable due to parameter interactions; therefore, we do not report the results for those models.

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 individually-optimized vs. non-optimized model predictions against the human behavior for the sample trial in Figure 6. 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) and 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 λ , we are able to capture human behavior in the task.

In Figure 9, we compare λ -optimized simplePSIRSA (right panel) with a uniform base model (left panel), which assigns equal probability to each utterance available for a particular context. This baseline model essentially has no utterance preferences whatsoever, but reflects the fact that the ambiguity classes distinguish between cases with a different number of utterance choice options (i.e., three to nine). Since simplePSIRSA with $\lambda = 0$ collapses to the baseline model, optimizing the λ parameter in simplePSIRSA is a nested model. When statistically comparing the individually-optimized λ simplePSIRSA model with the baseline model, the likelihood ratio test confirms that individual λ optimizations yield better fits than the baseline model ($G^2 = 268.87$, $df = 82$, $p < 0.01$).

We were able to distinguish three groups of participants on the basis of the individually-optimized parameter values of λ . The first was a “lazy worker” group of 18 participants whose fitted λ values were close to zero (i.e., $-.02 < \lambda < .02$), indicating that they were randomly selecting utterances. The second group of 32 participants yielded more negative values (i.e., $-7.13 < \lambda < -.02$), indicating that a significant number of participants preferred to systematically choose unambiguous utterances. The third group of again

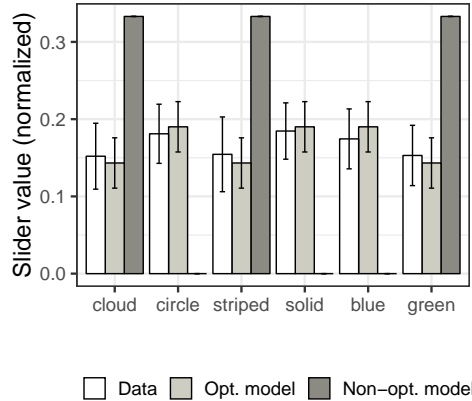


Figure 8: Simple Social Inference model predictions and human data for one of the classes of stimuli *Experiment 2: Picking utterances*. The optimized version of the model is optimized for the KL-factor λ .

32 participants yielded more positive values (i.e., $.02 < \lambda < .54$), indicating that these participants indeed chose the most ambiguous utterance in a strategic manner.

Discussion

We have found strong support for PSIRSA, which infers preference posteriors on the basis of ambiguous language. 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 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 to 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 most likely to inform their understanding of the preferences of their listeners. However, this group of ambiguity-selecting participants included only about 40% of the participants. Further experiments with highly similar setups (not reported in detail here) confirmed this trend. In particular, we ran a complementary study with a blocked design where participants first completed preference-inferences trials as in Experiment 1 and then completed utterance-selection trials as in Experiment 2. Even in such an experimental setup, the trend stayed the same. 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 our model in predicting those results indicate that humans are aware of the fact that by

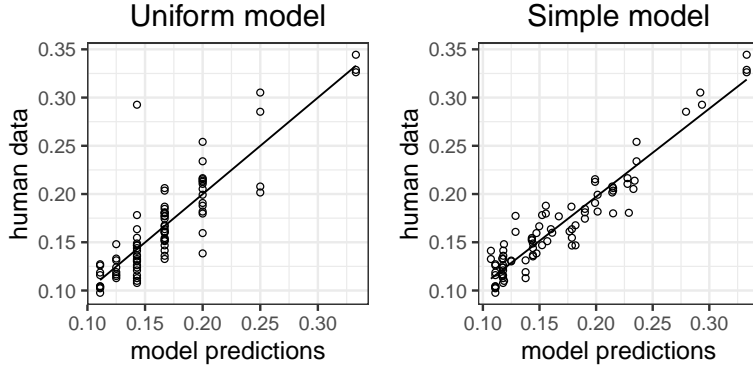


Figure 9: Average human data from Experiment 2 plotted against the predictions of the uniform baseline model and the simplePSIRSA model. Left panel *uniform model* ($r^2 = 0.7497$); right panel *KL-factor λ -optimized model* ($r^2 = 0.9071$).

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 to learn about others 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 to 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 the preferences of others, but may also concern all imaginable knowledge, opinions, beliefs, current trains of thought, and preferences 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 are one device for projecting parts of each other’s privileged grounds into the common ground.

Methods

Experiment 1: Learning about others' preferences

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.

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).

Depending on how many features competitor objects share with the target object, we were able to identify 48 ambiguity classes. Ambiguity classes group trials where a model considers a similar number of alternatives that could qualify for the uttered feature. For example, in Figure 2, the utterance "red" picks out two possible objects. If, however, the utterance were "green", only one object would qualify, and no learning about preferences would be possible. In that case, the model would assign equal probability to the listener's preferring dotted objects, striped objects, clouds, and squares. Once the model establishes that more than one object can be picked, it also needs to consider whether alternative objects share their features with the target object. For example, if both red objects were also striped, the model would not be able to infer any preferences about the pattern. Finally, we also code whether the objects that were not picked are similar in any of their feature values.

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).

Ambiguity classes

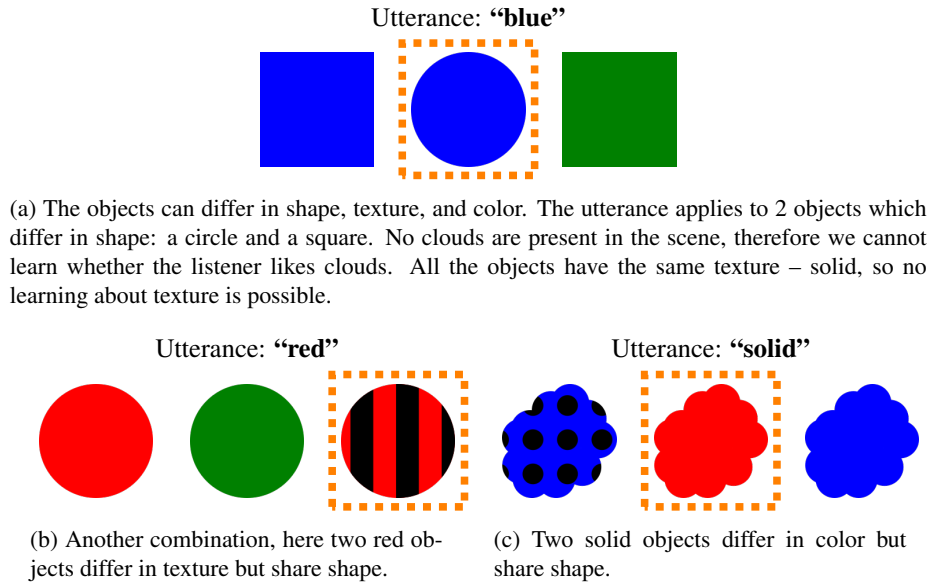


Figure 10: In all these examples, the utterance picks out two objects. The picked object has one feature for itself. The other two objects share that feature between them. The last feature is shared by all objects.

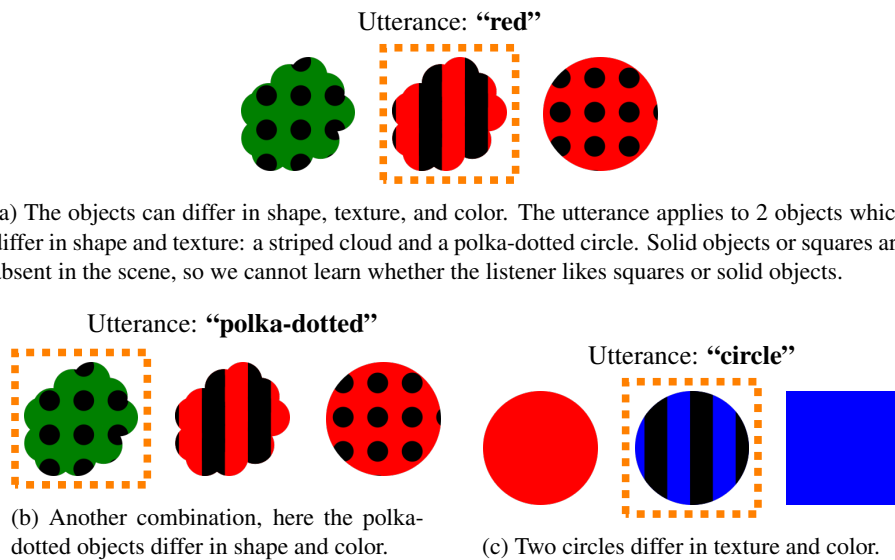


Figure 11: The utterance always picks out two objects and both objects only share that uttered feature. The third object shares one feature with each of those objects.

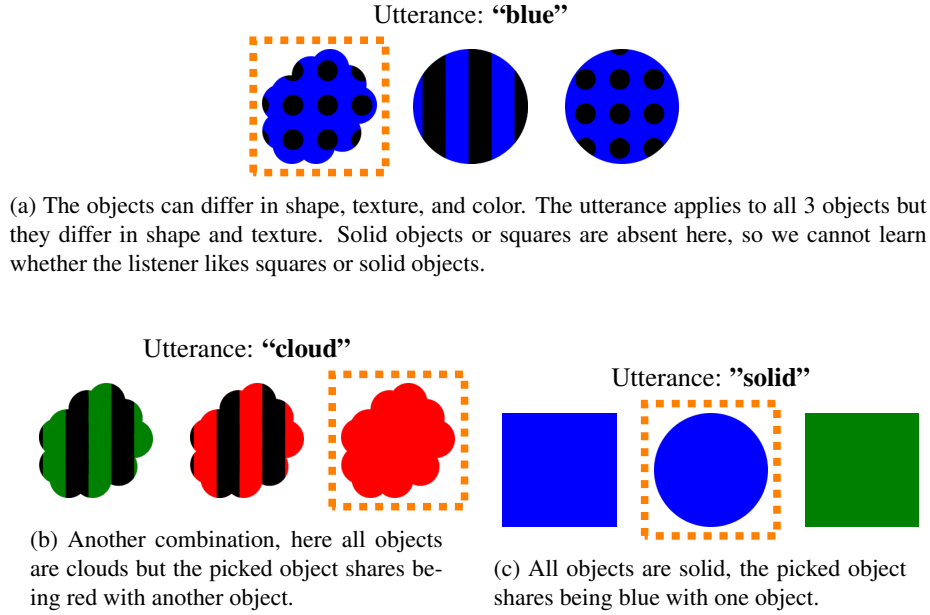


Figure 12: Here the utterance picks out all objects. The picked object shares one feature with one other object and has one feature just for itself while the other two objects share it.

Experiment 2

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.

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 6. 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 6, 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 (β) is not 0—that is, unless listeners have a tendency to disobey the utterance and consider objects that do not satisfy its literal interpretation.

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 6) and five trials were uninformative with respect to the listener’s preferences (e.g., observing a set of three identical objects).

Modeling: fullPSIRSA

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 :

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

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:

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

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 state observations $s \in S$:

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

where the factor λ scales the importance of the KL divergence term.

Optimization procedure

To compare PSIRSA’s predictions to the human data, we calculated an average value for each slider, binning data into 48 ambiguity classes for Experiment 1 and 14 classes

for Experiment 2. 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.

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

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

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 $P_{S_1\text{-simp}}(f | u, s)$ in the case of simplePSIRSA. By minimizing the KL divergence between the empirical and model-predicted preference posteriors for each participant, we maximize the model fit to the participants’ data. Moreover, we can use the minimized KL divergence values to perform the likelihood ratio test for nested models relying on the G^2 -statistic, because the summed KL divergence values are approximately chi-square distributed (Lewandowsky & Farrell, 2011). Individual vs. global-level parameter fitting allows us to explore potential differences between participants, and, more importantly, to evaluate whether the Gricean reasoning strategies apply at the level of individual speakers or only to the population as a whole (Franke & Degen, 2016).

We fit three parameters for fullPSIRSA and two for simplePSIRSA. The soft-max scaling factor α 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 γ regulates the strength of individual feature preferences f :

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

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 assume $\gamma = 0$ —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 β 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 .

References

- Baldwin, D. A., & Kosie, J. E. (2019). How does the mind render streaming experience as events? *Topics in Cognitive Science*. (this volume)
- 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.
- 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.
- 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.
- Eisenberg, E. M. (1984). Ambiguity as strategy in organizational communication. *Communication monographs*, 51(3), 227–242.
- Elman, J. L., & McRae, K. (2019). A model of event knowledge. *Psychological Review*, 126, 252291. 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. (2019). Roles for event representations in sensorimotor experience, memory formation and language processing. *Topics in Cognitive Science*. (this volume)
- 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.
- 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.
- 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 Psy-*

- chological Science*, 7(4), 341–351.
- Shin, Y. S., & DuBrow, S. (2019). Structuring memory through inference-based event segmentation. *Topics in Cognitive Science*. (this volume)
- 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*. (this volume)
- Storchak, H., Ehlis, A.-C., & Fallgatter, A. J. (2019). Action monitoring alterations as indicators of predictive deficits in schizophrenia. *Topics in Cognitive Science*. (this volume)
- Ünal, E., Ji, Y., & Papafragou, A. (2019). From event representation to linguistic meaning. *Topics in Cognitive Science*. (this volume)
- Wasow, T. (2015). Ambiguity avoidance is overrated. In S. Winkler (Ed.), *Ambiguity: Language and communication* (p. 29-47). de Gruyter.