

Speakers try to be helpful and look helpful: Modeling epistemic and social goals in speech

Erica J. Yoon,^{1*†} Michael Henry Tessler,^{1*} Noah D. Goodman,¹ Michael C. Frank¹

¹Department of Psychology, Stanford University,
450 Serra Mall, Stanford, CA 94305.

*These authors contributed equally to this work.

†To whom correspondence should be addressed; E-mail: ejyoon@stanford.edu.

Through language, people not only communicate information about the world, but also form, maintain, and improve relationships with others. We propose a computational model of speech in which the speaker’s goals reflect these functions of language: to be *epistemically* helpful, conveying the true state to the listener; to be *socially* helpful, making the listener feel good; and to be *presentational*, having herself appear to be helpful and have these two goals in mind. Tested against a case of simple polite speech, our model is able to predict people’s speech production judgments. Our extension of formal theories of communication to account for speakers’ social goals represents an advance in understanding of human speech.

With language, we can convey to others our knowledge and beliefs about the world, but we don’t always say what is on our mind explicitly or even truthfully. Imagine your friend tries on her new dress and excitedly asks you, “How does this look?” If her dress is truly hideous, the most informative answer will be “It looks terrible.” Yet intuitively we are more likely to lie and say “It looks great!” or carefully remark “Oh, it doesn’t look too bad, but I really loved the other

dress you wore last time.” But these indirect or false utterances can mislead the listener away from the truth, so why do people deviate from producing the most direct, truthful utterance possible? In this report, we present a computational model of speech that is indirect or false to some degree, but thereby offers a balance between informational and social goals of the speaker.

Language is a versatile tool for both information exchange and social rapport. Through language, people can communicate information to improve the listeners’ epistemic knowledge, but also form and improve their social relationships with others. On one hand, language can be seen as a transmission device that transfers information that reflects context or the state of affairs from a sender to a receiver (?, ?, ?). Informativity then is an integral assumption of language use, on which speakers rely to convey more than what their spoken words literally suggest (?, ?). On the other hand, language performs important social roles, whereby people make contact with others and form relationships (?, ?). Language users tend to speak in a manner that abides by expectations of the community (?), promoting a sense of social solidarity. Thus, epistemic and social functions of language translate to speakers’ goals to be informative and social in their language use. [mht: I think the first paragraph can either (1) be tightened up a bit or (2) merged with the second paragraph... probably (1) is easier]

Informational and social goals often conflict with one another. Courteous language to make the listener feel good (“Your dress is gorgeous” about a truly hideous dress; “It’s hard to give a good presentation” after an awful talk) will risk potential loss of their intended message, suffer inefficiencies, or even convey wrong information. Conversely, maximally informative and efficient utterances can seem too blunt or rude (“Be quiet!”; “Read my essay and give feedback now.”) and fail to promote or maintain good relationships. How do people then speak to find balance between goals to be informative versus socially apt?

Polite speech, in which people try to preserve other people’s feelings by softening their potentially harsh message, reflects the compromise between informative and social goals in

conflict. On a daily basis, adults and even young children spontaneously produce polite, indirect requests (“Can you please close the window?”) rather than direct orders (“Close the window.”) (?), and speakers use polite forms of speech even while arguing, preventing unnecessary offense to their interactants (?). Based on informativity-focused accounts of language, politeness clashes with the idea of cooperative communication as it violates the goal to exchange information efficiently and accurately (?); if information transfer were the only currency in communication, politeness would be both infelicitous and undesirable. A polite speaker, however, can be seen to be cooperative and helpful via both an *epistemic* goal to improve the listener’s knowledge and a *social* goal to minimize potential damage to the hearer’s (and the speaker’s own) self-image, called *face* (?). If the speaker’s intended meaning contains no threat to the speaker or listener’s face, then the speaker will choose to convey the meaning in an explicit and efficient manner (putting it “on the record”). As the degree of face-threat becomes more severe, however, a speaker will choose to be polite by producing more indirect utterances. [mht: I feel like paragraphs (3) & (4) can be merged]

In this paper, using a case study of simple polite speech, we formalize the idea that language reflects a principled tradeoff between speakers’ epistemic and social communicative goals. Previous informal theories of politeness tried to explain how speakers’ social goals give rise to polite speech: For example, a prominent theory of politeness (?) proposes that deviation from informativity increases the level of polite face-saving. But there has not yet been a formal account of the social goals in speech, thus no systematic, quantitative predictions of polite speech production have been available. [mht: ← this issue (perhaps all of these issues) should be brought up earlier]

On the other hand, formal theories of language have extensively accounted for speakers’ desires to be informative, but not for their potential social goals. The Rational Speech Act (RSA) framework describes language understanding as recursive probabilistic inference be-

tween a pragmatic listener and an informative speaker (?). This framework has been successful at capturing the quantitative details of a number of language understanding tasks but it neglects the social goals a speaker may pursue. Also, work with RSA models have focused mostly on listener inference tasks in which people infer meanings of utterances or goals behind them based on the assumption that people produce corresponding utterances, and there have been fewer speaker production tasks that verify that people actually produce utterances as predicted by the models. [mht: not sure how much we want to say about “issues within the RSA family of models”... maybe there can be a paragraph about theoretical and computational issues, merging the previous 2 paragraphs]

Here we propose a computational model that unifies formal theories of informative communication and informal theories of polite speech, and accounts for both epistemic and social goals of speakers; We then verify the model predictions against human predictions for speaker production. Our model represents a context in which someone (e.g., Bob the listener) gave a performance of some kind, such as a presentation, and asks another person (e.g., Ann the speaker) how well he did. The speaker considers two factors to choose her utterance: the true state, or the rating truly deserved by the performance, and her goal (Figure ??). The speaker’s goal is represented as a utility function, based on which she chooses utterances approximately optimally (?). The speaker’s utility function in the current model is made up of three components. First, *epistemic utility* refers to the standard, informative utility in RSA: the amount of information a literal listener (L_0) would still not know about world state s after hearing a speaker’s utterance w . In our model, besides the standard epistemic utility, we add utilities to represent the speaker’s social desires that lead her to produce utterances that deviate from informativity but help maintain good relationships and self-reputation: *Social utility* is the expected subjective utility of the state inferred given the utterance w , and is related to the intrinsic value of the state. We use a value function (V) to map states to subjective utility values, which captures

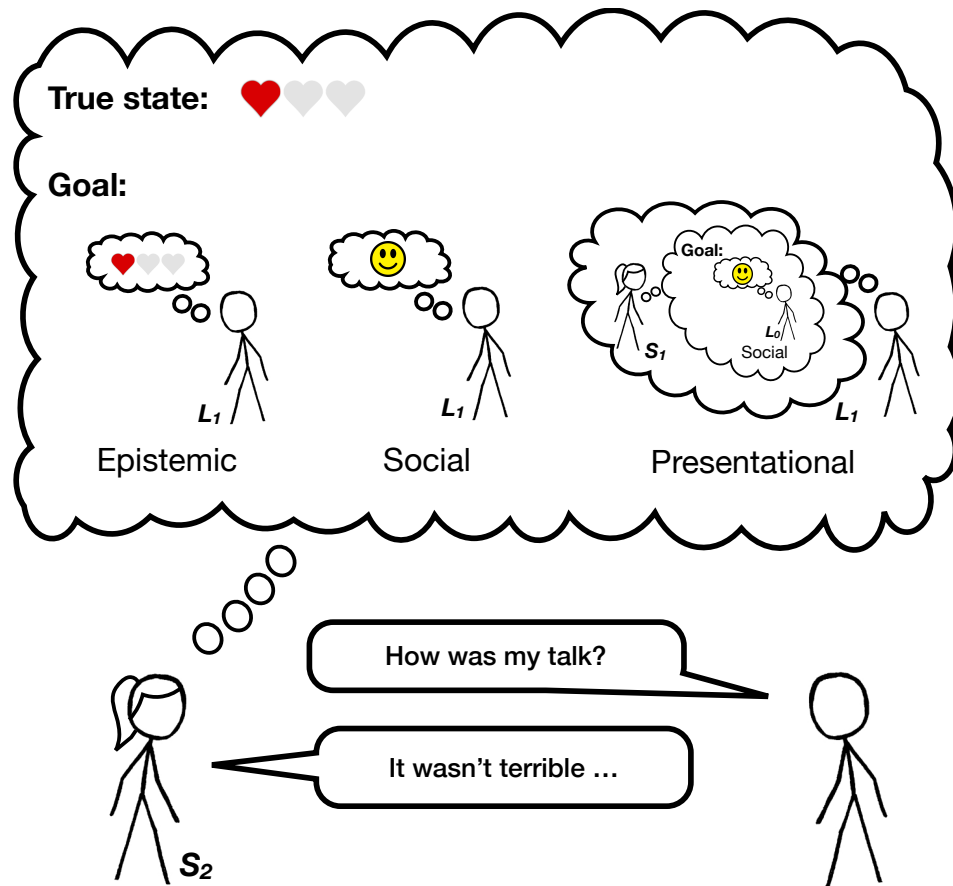


Figure 1: Diagram of the pRSA model: The “pragmatic speaker” observes the true state and determines her goal between three utilities (epistemic, social, and presentational), and produces an utterance.

the affective consequences for the listener of being in state s . *Presentational utility* concerns conveying a particular utility (ϕ_{S_1}) to the listener, such that the speaker appears to have a particular goal (e.g. to be socially helpful) in mind. A weight is assigned to each utility component to determine which goal is to be prioritized ($\phi_{epistemic}$, ϕ_{social} , $\phi_{presentational}$ respectively). Finally, some utterances might be costlier than others. The total utility of an utterance subtracts the cost $C(w)$ from the weighted combination of the epistemic, social, and presentational utilities).

$$U(w; s; \hat{\phi}) = \phi_{epistemic} \cdot \ln(P_{L_1}(s | w)) + \phi_{social} \cdot \mathbb{E}_{P_{L_1}(s|w)}[V(s)] + \phi_{presentational} \cdot \ln(P_{L_1}(s, \phi_{S_1} | w)) - C(w) \quad (1)$$

Based on the utility function above, the speaker (S_2) chooses utterances w approximately optimally (as per speaker optimality parameter λ_{S_2}) given the state s and his goal weights $\hat{\phi}$.

$$P_{S_2}(w | s, \hat{\phi}) \propto \exp(\lambda_{S_2} \cdot \mathbb{E}[U(w; s; \hat{\phi})]) \quad (2)$$

We used a simple procedure to empirically test whether our model is able to predict production of polite utterances. Participants read scenarios in which someone (e.g. Bob) gave a performance of some kind, and another person (Ann) evaluated it (see Supplementary Materials for a detailed description of the task). We provided information on Ann’s feelings toward the presentation (*true state*), which were shown on a scale from zero to three hearts (e.g. one out of three hearts filled in red color). We also presented Ann’s *goal*, which was one of the following: to be *informative* and give accurate feedback; to be *social* and to make Bob feel good; or to be *both* informative and social at the same time. We assumed that the goal descriptions conveyed to the participants a particular set of goal weights ($\phi_{epistemic}$, ϕ_{social} , $\phi_{presentational}$, ϕ_{S_1}) that the speaker was using. We put uninformative priors on these weights ($\beta \sim Uniform(0, 1)$) and inferred their credible values separately for each goal condition (see Supplementary Materials). We hypothesized that speakers with both goals to be informative and social given bad true

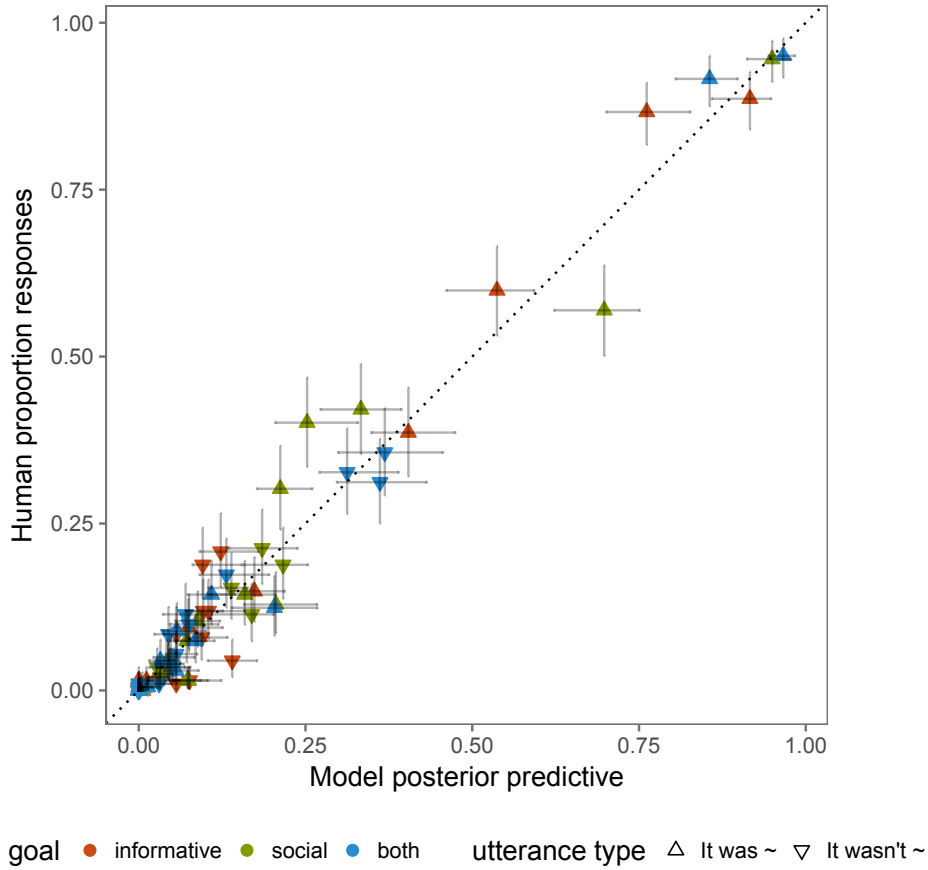


Figure 2: Full distribution of human responses vs. fitted model predictions for pragmatic speaker production. Error bars represent 95% confidence intervals for the data (vertical) and 95% highest density intervals for the model (horizontal).

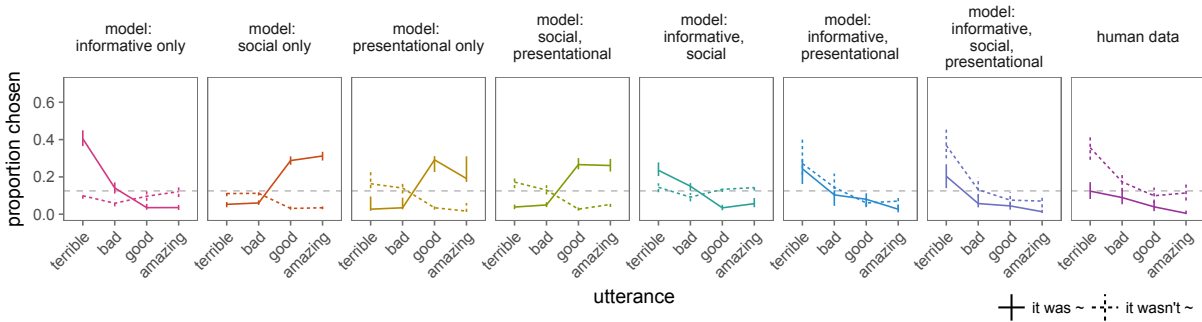


Figure 3: Comparison of predictions for proportion of utterances chosen by pragmatic speaker from possible model variants (left) and human data (rightmost) for average proportion of negation produced among all utterances, given true state of 0 heart (on a scale of 0 to 3) and speaker with both goals. Gray dotted line indicates chance level at 12.5%.

states (i.e. Bob’s performance was poor) would produce more negation (“It wasn’t”) to save the listener’s face while vaguely conveying the bad true state (see our pre-registered model, hypothesis, and procedure at FIXME). Each participant read 12 scenarios total (4 true states \times 3 goals). In a single trial, each scenario was followed by a question that asked for the most likely utterance by Ann. Participants selected one of eight possible utterances, by choosing between *It was* vs. *It wasn’t* and then among *terrible*, *bad*, *good*, and *amazing*. We separately gathered the literal meaning judgments for the eight possible utterances, by measuring how likely each utterance is to be true given each true state, to set expected literal meanings of utterances in our model (see Supplementary Materials for literal semantic results).

Mean proportion of utterances chosen by participants in each true-state \times goal condition were overall highly consistent with the our model predictions. The posterior predictive of the model explained almost all of the variance in the production data ($r^2(96) = 0.97$; Figure ??). We compared the predictions of our model with its variants, which had different combinations of utility components: with one or two of the three utilities to be epistemic, social, and presentational (Figure ??). The current model yielded the closest predictions to the human data, especially for the condition in which the speaker was described to have both goals to be informative and social. The model fitting for the both-goal condition suggested that the speaker had a high weight on epistemic utility, low weight on social utility, and high weight on presentational utility to convey high social utility, meaning that the speaker wanted to be genuinely informative while appearing to care about the listener’s feelings (see Supplementary Materials for inferred weights for all goal conditions). Unlike other model variants, the current model was able to capture the key pattern of the human data that given a truly terrible performance (e.g. with the true state of 0 or 1 heart) a speaker with both informative and social goals would prefer to produce “It wasn’t terrible” the most, which would avoid saying a directly negative utterance but allow the listener to infer that his performance was actually (close to being) terrible. Both

| Model | Variance explained | Marginal likelihood |
|-------------------------------------|--------------------|---------------------|
| informative only | 0.811 | |
| social only | 0.228 | |
| presentational only | 0.237 | |
| social, presentational | 0.243 | |
| informative, social | 0.909 | |
| informative, presentational | 0.95 | |
| informative, social, presentational | 0.969 | |

Table 1: Inferred parameters from all model variants.

the variance explained ($r^2 = 0.969$) and the Bayes factor (FIXME) were the highest for the current model compared to its alternatives (see Table 1 in Supplementary Materials).

Our work unifies previous formal models of communication and informal theories of social uses of language. Our findings suggest that neither epistemic nor social motives alone motivate polite speech; instead, production of polite speech results from the conflict between these two, combined with a self-presentational desire to look epistemically and socially helpful.

Beyond language understanding, our model casts new light on understanding cognitive mechanisms that involve both information processing and interactions among agents. Many processes of human learning and inference have been extensively and fruitfully examined under computational modeling approaches similar to our own (, , ,). However, these processes have been mostly described as purely information-driven, separate or independent of social goals of the agents involved. Our work provides the next step for addressing the issue of how social goals influence cognitive processes that involve more than one agent. [mht: i’m not sure what this paragraph conveys... the citations make it seem about cognition generally, but there might be a more targeted approach about language and then broadening to NLP applications]

Our formal model also has important implications for building an artificial system that is more human-like. Social interactions require more than merely following simple rules (say “please”, “thank you”) at the right moments; they require a balance of more intricate goals to

be informative and social. [mht: maybe we can speculate about learning speaker-level parameters over repeated interactions, analagous to hierarchical learning about speaker’s articulatory patterns, e.g., Kleinschmidt & Jaeger (2015 ?) psych review] We successfully modeled the speaker who considers the tradeoff between communicative informativity, kindness and presentation of herself as helpful. This work takes a concrete step toward quantitative models of the nuances of human speech. And it moves us closer to courteous computation – to computers that communicate with tact.

References

1. K. Bühler, *Sprachtheorie* (Oxford, England: Fischer, 1934).
2. C. E. Shannon, *Bell Syst. Tech. J.* **27**, 623 (1948).
3. R. Jakobson, *Style in language* (MA: MIT Press, 1960), pp. 350–377.
4. H. P. Grice, *Logic and conversation* (Academic Press, 1975), vol. 3, pp. 41–58.
5. J. Searle, *Indirect Speech Acts* (Academic Press, 1975), vol. 3, pp. 59–82.
6. M. A. K. Halliday, *Learning How to Mean: Explorations in the Development of Language* (London: Edward Arnold, 1975).
7. S. M. Ervin-Tripp, *Journal of social issues* **23**, 78 (1967).
8. S. M. Ervin-Tripp, *Advances in experimental social psychology* **4**, 91 (1969).
9. H. H. Clark, D. H. Schunk, *Cognition* **8**, 111 (1980).
10. G. Axia, M. R. Baroni, *Child Development* pp. 918–927 (1985).
11. T. Holtgraves, *Journal of Language and Social Psychology* **16**, 222 (1997).

12. P. Brown, S. C. Levinson, *Politeness: Some universals in language usage*, vol. 4 (Cambridge university press, 1987).
13. N. D. Goodman, M. C. Frank, *Trends in Cognitive Sciences* **20**, 818 (2016).
14. N. D. Goodman, A. Stuhlmüller, *Topics in cognitive science* **5**, 173 (2013).
15. J. B. Tenenbaum, C. Kemp, T. L. Griffiths, N. D. Goodman, *Science* **331**, 1279 (2011).
16. F. Xu, J. B. Tenenbaum, *Psychological review* **114**, 245 (2007).
17. E. Bonawitz, *et al.*, *Cognition* **120**, 322 (2011).
18. C. L. Baker, R. Saxe, J. B. Tenenbaum, *Cognition* **113**, 329 (2009).
19. M. D. Lee, E. J. Wagenmakers, *Bayesian Cognitive Modeling: A Practical Course* (Cambridge Univ. Press, 2014).

Acknowledgments

This work was supported by NSERC PGS Doctoral scholarship PGSD3-454094-2014 to EJY, NSF Graduate Research Fellowship DGE-114747 to MHT, ONR grant N00014-13-1-0788 to NDG, and NSF grant BCS 1456077 to MCF.

Supplementary materials

Materials and Methods

Literal semantic task

We probed judgments of literal meanings of the target words assumed by our model and used in all our experiments. 51 participants with IP addresses in the United States were recruited on Amazon’s Mechanical Turk. We used 13 different context items in which someone evaluated a

performance of some kind. For example, in one of the contexts, Ann saw a presentation, and Ann’s feelings toward the presentation (*true state*) were shown on a scale from zero to three hearts (e.g., two out of three hearts filled in red color). The question of interest was ”Do you think Ann thought the presentation was / wasn’t X?” and participants responded by choosing either “no” or “yes.” The target could be one of five possible words: *terrible*, *bad*, *good*, and *amazing*, giving rise to ten different possible utterances (with negation or no negation). Each participant read 32 scenarios, depicting every possible combination of states and utterances. The order of context items was randomized, and there were a maximum of four repeats of each context item per participant. For this and the subsequent experiment, we analyzed the data by collapsing across context items.

For each utterance-state pair, we computed the posterior distribution over the semantic weight (i.e., how consistent X utterance is with Y state) assuming a uniform prior over the weight. Meanings of the words as judged by participants were as one would expect (see Figure ??). We used the fraction of participants that endorsed utterance w for state s to set informative priors to infer posterior credible values of the literal meanings from data in the speaker production experiment.

Speaker production task

202 participants with IP addresses in the United States were recruited on Amazon’s Mechanical Turk. As in the semantics measurements above, we used scenarios in which a person (e.g., Bob) gave some performance and asked for another person (e.g., Ann)’s opinion on the performance (see Fig. ??). Additionally, we provided information on the speaker Ann’s goal – *to make Bob feel good*, or *to give as accurate and informative feedback as possible*, or *both* – and the true state – how Ann actually felt about Bob’s performance (e.g., two out of three hearts, on a scale from zero to three hearts). Each participant read 12 scenarios, depicting every possible

combination of goals and states. The order of context items was randomized, and there were a maximum of two repeats of each context item per participant.

Each scenario was followed by a question that read, "If Ann wanted *to make Bob feel good* but not necessarily give informative feedback (or *to give accurate and informative feedback* but not necessarily make Bob feel good, or *BOTH make Bob feel good AND give accurate and informative feedback*), what would Ann be most likely to say?" Participants indicated their answer by choosing one of the options on the two dropdown menus, side-by-side, one for choosing between *It was* vs. *It wasn't* and the other for choosing among *terrible*, *bad*, *good*, and *amazing*.

Supplementary Text

Model fitting and inferred parameters

In the speaker production task, participants were told what speakers' intentions were (e.g. wanted to make Bob feel good). We assume that the intention descriptions conveyed the weight mixtures $\phi_{epistemic}$, ϕ_{social} , ϕ_{self} and ϕ_{S_1} that the speaker was using. We put uninformative priors on each of these mixtures ($\phi \sim Uniform(0, 1)$) and inferred their credible values separately for each goal condition ("wanted to X") using Bayesian data analytic techniques (?). We ran 4 MCMC chains for 40,000 iterations, discarding the first 20,000 for burnin. The inferred values of weight mixtures for each model variant (with different phi components) are shown in Table ??.

There were two additional parameters of the model, on which we put uninformative priors: the value scale parameter ($\alpha \sim Unif(0, 10)$) in the utility function; and the cost parameter ($C(u) \sim Unif(1, 10)$). We inferred their posterior credible values from the data. The Maximum A-Posteriori (MAP) estimates and 95% Highest Probability Density Intervals (HDI) are shown in Table ??.

| Model | informative speaker condition | | | | social speaker condition | | | | both-goal speaker condition | | | | speaker optimality | cost |
|-------------------------------------|-------------------------------|--------------|---------------|--------------|--------------------------|--------------|---------------|--------------|-----------------------------|--------------|---------------|--------------|--------------------|------|
| | ϕ_{inf} | ϕ_{soc} | ϕ_{pres} | ϕ_{S_1} | ϕ_{inf} | ϕ_{soc} | ϕ_{pres} | ϕ_{S_1} | ϕ_{inf} | ϕ_{soc} | ϕ_{pres} | ϕ_{S_1} | | |
| informative only | 0.07 | 0.51 | 0.43 | 0.51 | 0.04 | 0.48 | 0.48 | 0.5 | 0.07 | 0.47 | 0.45 | 0.5 | 11.15 | 1.64 |
| social only | 0.46 | 0.09 | 0.44 | 0.53 | 0.23 | 0.5 | 0.27 | 0.49 | 0.3 | 0.32 | 0.38 | 0.51 | 5.46 | 1.76 |
| presentational only | 0.45 | 0.37 | 0.18 | 0.67 | 0.27 | 0.35 | 0.38 | 0.66 | 0.32 | 0.37 | 0.31 | 0.63 | 13.75 | 3.67 |
| social, presentational | NA | 0.43 | 0.57 | 0.05 | NA | 0.63 | 0.37 | 0.31 | NA | 0.42 | 0.58 | 0.51 | 3.43 | 1.62 |
| informative, social | 0.74 | 0.26 | NA | 0.51 | 0.39 | 0.61 | NA | 0.54 | 0.53 | 0.47 | NA | 0.51 | 3.65 | 1.18 |
| informative, presentational | 0.49 | NA | 0.51 | 0.63 | 0.46 | NA | 0.54 | 0.86 | 0.53 | NA | 0.47 | 0.8 | 4.26 | 3.04 |
| informative, social, presentational | 0.32 | 0.02 | 0.66 | 0.56 | 0.27 | 0.26 | 0.46 | 0.71 | 0.41 | 0.03 | 0.56 | 0.73 | 4.98 | 3.14 |

Table 2: Inferred parameters from all model variants.

Data analysis tools

We used R (3.4.2, R Core Team, 2017) and the R-packages *bindrcpp* (0.2, Miller, 2017), *binom* (1.1.1, Dorai-Raj, 2014), *coda* (0.19.1, Plummer, Best, Cowles, & Vines, 2006), *dplyr* (0.7.4, Wickham, Francois, Henry, & Miller, 2017), *forcats* (0.2.0, Wickham, 2017a), *ggplot2* (2.2.1, Wickham, 2009), *ggthemes* (3.4.0, Arnold, 2017), *gridExtra* (2.3, Auguie, 2017), *jsonlite* (1.5, Ooms, 2014), *langcog* (0.1.9001, Braginsky, Yurovsky, & Frank, n.d.), *magrittr* (1.5, Bache & Wickham, 2014), *papaja* (0.1.0.9492, Aust & Barth, 2017), *purrr* (0.2.4, Henry & Wickham, 2017), *readr* (1.1.1, Wickham, Hester, & Francois, 2017), *rwebppl* (0.1.97, Braginsky, Tessler, & Hawkins, n.d.), *stringr* (1.2.0, Wickham, 2017b), *tibble* (1.3.4, Miller & Wickham, 2017), *tidyr* (0.7.2, Wickham & Henry, 2017), and *tidyverse* (1.2.1, Wickham, 2017c) for all our analyses.

Figs. ?? to ??

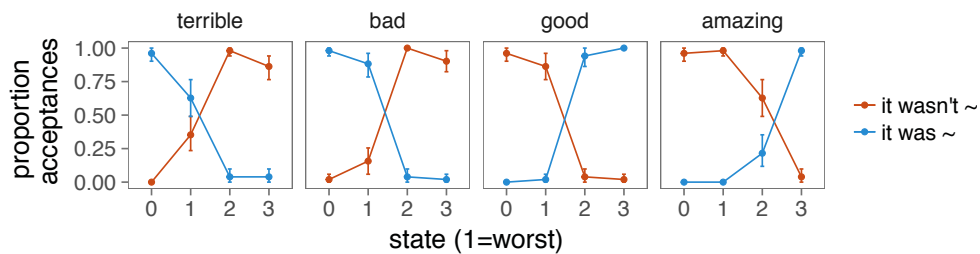



Figure S1: Semantic measurement results. Proportion of acceptances of utterance types (colors) combined with target words (facets) given the true state represented on a scale of hearts. Error bars represent 95% confidence intervals.

Imagine that Kelly baked some cookies, but she didn't know how good it was. Kelly approached Justine, who knows a lot about baking, and asked "How did my cookie taste?"

Here's how Justine **actually** felt about Kelly's cookie, on a scale of 0 to 3 hearts:



If Justine wanted to BOTH make Kelly feel good AND give accurate and informative feedback,

what would Justine be most likely to say?

"It

Figure S2: Example of a trial in the speaker production task.

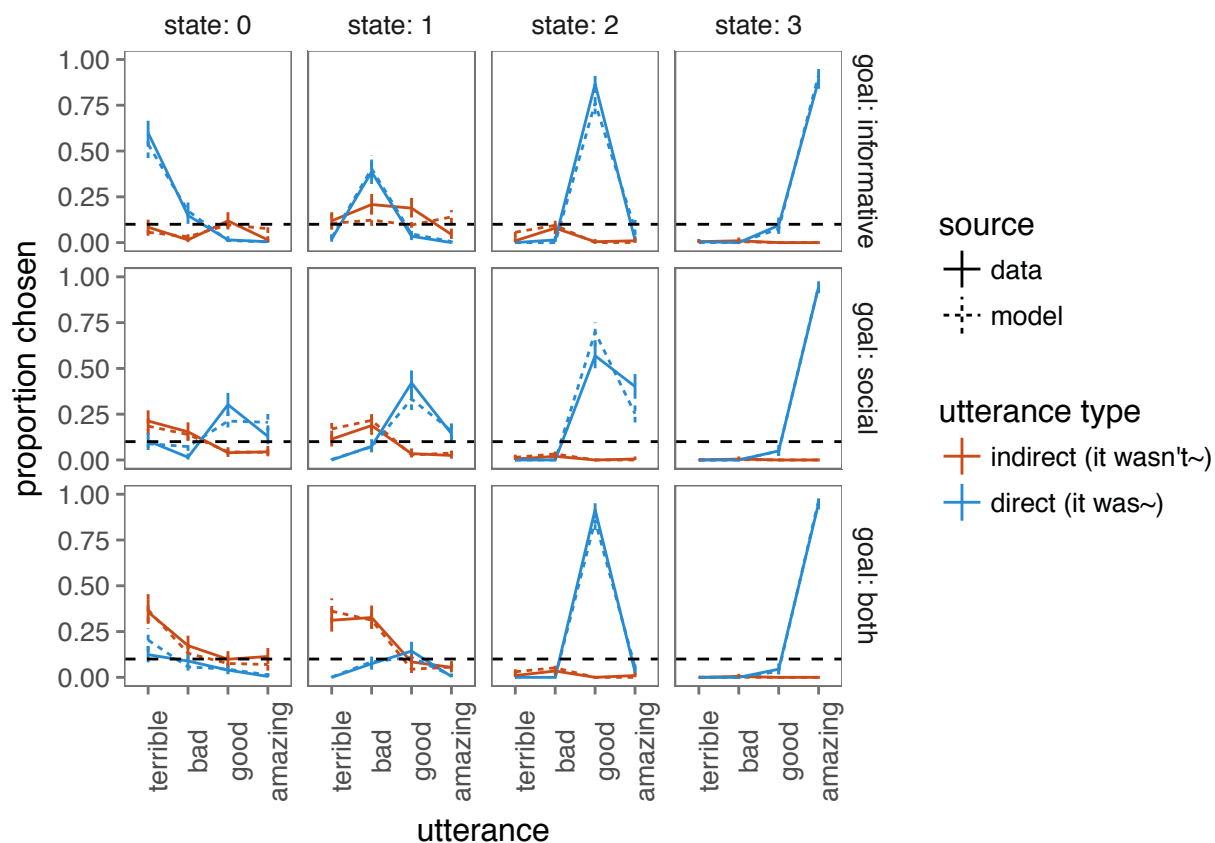


Figure S3: Experimental results (solid lines) and fitted model predictions (dashed lines) for speaker production. Proportion of utterances chosen (utterance type – direct vs. indirect – in different colors and words shown on x-axis) given the true states (columns) and speaker goals (rows). Error bars represent 95% confidence intervals for the data and 95% highest density intervals for the model. Black dotted line represents the chance level.

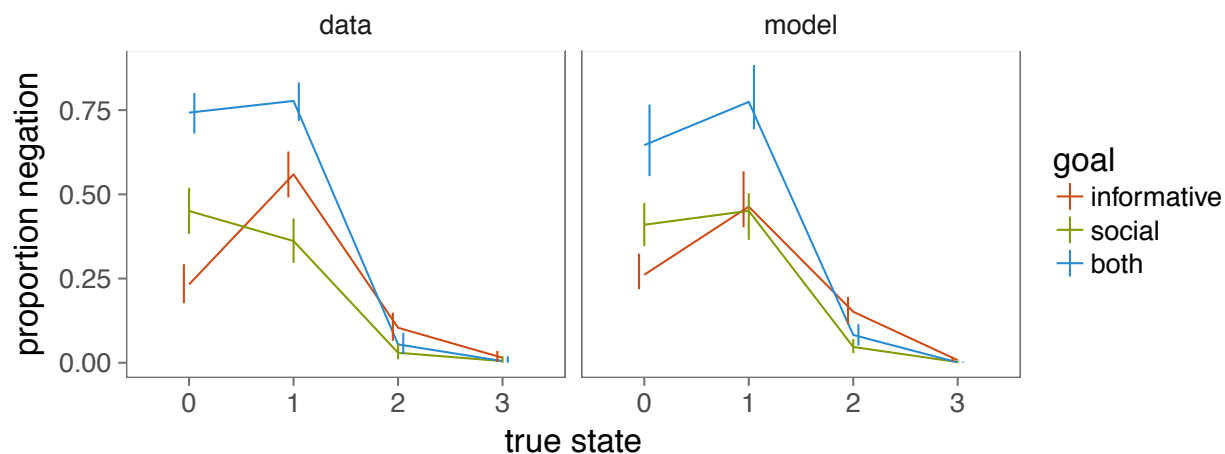


Figure S4: Experimental results (left) and fitted model predictions (right) for average proportion of negation produced among all utterances, given true states (x-axis) and goals (colors).

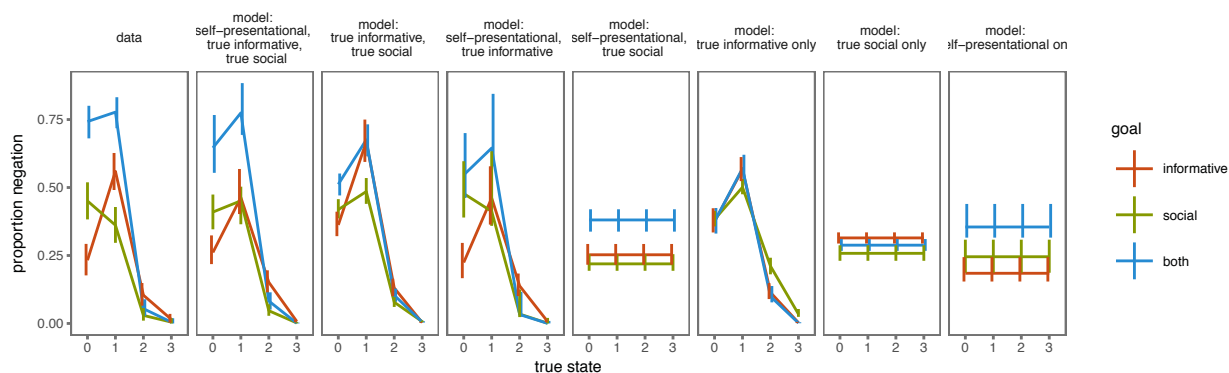


Figure S5: Experimental results (leftmost) and predictions from different model alternatives for average proportion of negation produced among all utterances, given true states (x-axis) and goals (colors).