

Polite speech emerges from competing social goals

Erica J. Yoon^{a,1,2}, Michael Henry Tessler^{a,1}, Noah D. Goodman^a, and Michael C. Frank^a

^aDepartment of Psychology, Stanford University

This manuscript was compiled on September 27, 2018

Language is a remarkably efficient tool for transmitting information. Yet human speakers make statements that are inefficient, imprecise, or even contrary to their beliefs, all in the service of being polite. What rational machinery underlies polite language use? Here, we show that polite speech emerges from the competition of informational, social, and self-presentational (i.e., *appearing* informational or social) goals. We formalize this tradeoff using a probabilistic model of utterance production, which predicts human choices with high quantitative accuracy. This utility-theoretic approach to speech acts takes a step towards explaining the richness and subtlety of social language use.

politeness | computational modeling | communicative goals | pragmatics

We rarely say exactly what's on our mind. Although *close the window!* could be effective message, we dawdle by adding *can you please...?* or *would you mind...?* Rather than tell an uncomfortable truth, socially-aware speakers lie (*Your dress looks great!*) and prevaricate (*Your poem was so appropriate to the occasion*). Such language use is puzzling for classical views of language as information transfer (1–4). On these views, transfer ought to be efficient and accurate: The speaker is expected to choose a succinct utterance to convey their beliefs (5, 6), and the information transferred is ideally accurate and truthful to the extent of the speaker's knowledge. Polite speech violates these basic expectations about the nature of communication: It is typically inefficient and underinformative, and sometimes even outright false. Yet even young speakers of a language spontaneously produce requests in polite forms (7, 8), and adults use politeness strategies while arguing, preventing unnecessary offense to their interactants (9).

If politeness only gets in the way of effective information transfer, why be polite? Clearly, we have social concerns, and most linguistic theories assume utterance choices are motivated by these concernss, couched as either polite maxims (10), social norms (11), or aspects of a speaker and/or listener's identity referred as *face* (12, 13). This latter theory predicts that when a speaker's intended meaning contains a threat to the listener's face or self-image (and potentially the speaker's face), her messages will be less direct, less efficient, and possibly untruthful. But how does a speaker decide to be indirect (*Your pie could use a bit of salt*) vs. false (*It's delicious!*), and when to tell the blunt truth? Additionally, how does the speaker's own image in the mind of the listener enter into the calculation?

We propose a utility-theoretic solution to the problem of polite language production by quantifying the tradeoff between competing communicative goals. In our model, speakers attempt to maximize utilities that represent their communicative goals: informational utility—derived via classical, effective

information transmission; social utility—derived by being kind and saving the listener's face; and self-presentational utility—derived by appearing in a particular way to save the speaker's own face. Speakers then produce an utterance on the basis of its expected utility (including their cost to utter, approximated by the length of the utterance). The lie that a pie was delicious provides social utility by making the baker feel good, but does not provide information about the true state of the world. Further, if the baker suspects that the pie was in fact terrible, the speaker runs the risk of being seen as uncooperative.

The speaker's utilities are weighed within a probabilistic model of pragmatic reasoning: the Rational Speech Act (RSA) framework (2, 14). Speakers are modeled as agents who choose utterances by reasoning about their effects on a listener relative to their cost, while listeners infer the meaning of an utterance by reasoning about speakers and their goals. This model class has provided a quantitative understanding of a wide variety of complex linguistic behaviors, including vagueness (15), hyperbole (16), and irony (17), among others. In this framework, language use builds on the idea that human social cognition can be approximated via reasoning about others as rational agents who act to maximize their subjective utility (18), a hypothesis which has found empirical support in a wide variety of work with both adults and children (e.g., (19, 20)).

RSA models are defined recursively such that speakers reason about listeners, and vice versa. We use convention in indexing and say a pragmatic listener L_1 reasons about what intended meaning and goals would have led a speaker S_1 to produce a particular utterance. Then S_1 reasons about

Significance Statement

Standard views of language use emphasize its role in information transfer. On these views, the ubiquity of politeness is a puzzle. We present a new quantitative viewpoint on social language use that resolves the puzzle by assuming that speakers balance three potentially conflicting goals: an informational goal (*be informative*), a social goal (*be nice*), and a self-presentational goal (*look nice*). This formal work provides a framework for developing a quantitative understanding of social language use more broadly. Further, by providing a theory of how speakers consider the feelings of others in formulating their utterances, it paves the way for the development of courteous computational agents.

Please provide details of author contributions here.

Please declare any conflict of interest here.

¹E.J.Y. and M.H.T. contributed equally to this work.

²To whom correspondence should be addressed. E-mail: ejyoon@stanford.edu

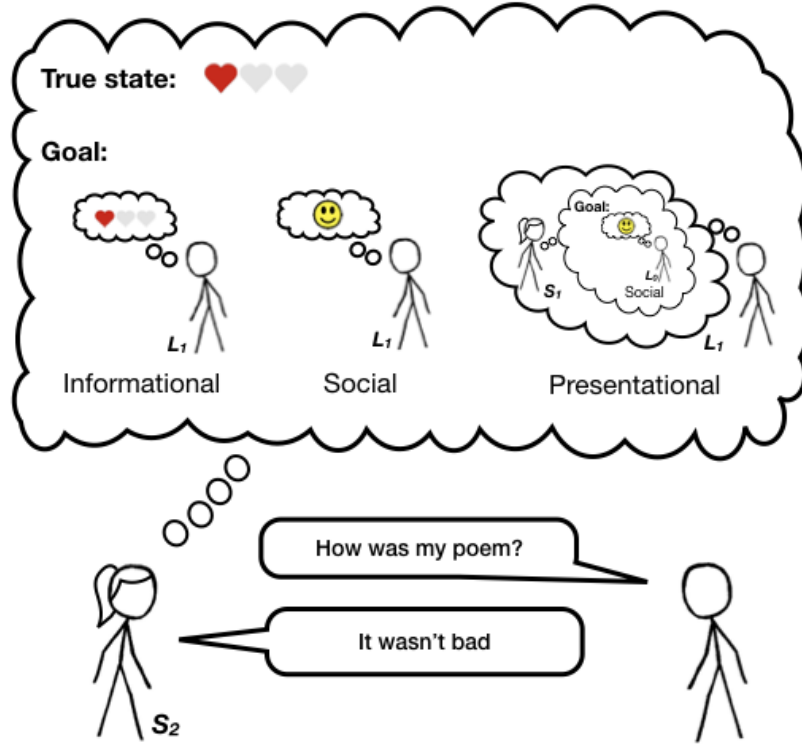


Fig. 1. Diagram of the model: The pragmatic speaker observes the true state and determines her goal between three utilities (informational, social, and presentational), and produces an utterance.

a *literal listener* L_0 , who is modeled as attending only to the literal meanings of words (rather than their pragmatic implications), and hence grounds the recursion. The target of our current work is a model of a polite speaker S_2 : S_2 reasons about what utterance to say to L_1 by considering the set of utilities described above (Figure 1).

We evaluate our model on its ability to predict human utterance choices in situations where polite language use is expected to varying degrees. Imagine Bob baked a pie and asked Ann how good it was. Ann (S_2) produces an utterance w based on the true state of the world s (i.e., the rating, in her mind, truly deserved by Bob's pie) and a set of goal weights $\hat{\phi}$, that determines how much Ann prioritizes each of the three possible goals. Ann's production decision is softmax, which interpolates between maximizing and probability matching (via λ_{S_2} ; (21)):

$$P_{S_2}(w|s, \hat{\phi}) \propto \exp(\lambda_{S_2} \cdot \mathbb{E}[U_{total}(w; s; \hat{\phi})]). \quad [1]$$

We posit that a speaker's utility contains three distinct components: informational, social, and presentational. The total utility U_{total} of an utterance is thus the weighted combination of the three component utilities minus the utterance cost $C(w)$:

$$U_{total}(w; s; \hat{\phi}) = \phi_{inf} \cdot U_{inf}(w; s) + \phi_{soc} \cdot U_{soc}(w; s) + \phi_{pres} \cdot U_{pres}(w; s) - C(w). \quad [2]$$

Foremost, we define *social utility* as the expected subjective utility $V(s)$ of the state implied to the listener by the utterance: $U_{soc}(w) = \mathbb{E}_{P_{L_1}(s|w)}[V(s)]$. The subjective utility function $V(s)$ could vary by culture and context; we test our model

when states are explicit ratings (e.g., on a 4-point scale) and we assume a positive linear value relationship between states and values V to model a listener's preference to be in a highly rated state (e.g., Bob would prefer to have made a pie deserving 4 stars rather than 1 star).

At the same time, a speaker may desire to be epistemically helpful, modeled as standard *informational utility* (U_{inf}). The informational utility indexes the amount of information a literal listener (L_0) would still not know about the state of the world s after hearing the speaker's utterance w (i.e., surprisal; e.g., how likely is Bob to guess Ann's actual opinion of the pie): $U_{inf}(w) = \ln(P_{L_1}(s|w))$. Speakers who optimize for informational utility produce accurate and informative utterances while those who optimize for social utility produce utterances that make the listener feel good.

If a listener is uncertain how their particular speaker is weighing the competing goals to be honest vs. kind, they might try to infer the weighting (e.g., "was she just being nice?"). But then a sophisticated speaker can produce utterances in order to appear *as if* they had certain goals in mind (i.e., a self-presentational goal). The extent to which the speaker *appears* to the listener to have a particular goal in mind (e.g., to be kind) is the utterance's *presentational utility* (U_{pres}) and is the most novel component of our model. The speaker gains presentational utility when her listener believes she has particular goals – that she is trying to be informative or kind. Formally,

$$U_{pres}(w) = \ln(P_{L_1}(\phi_{S_1} | w)) = \ln \int_s P_{L_1}(s, \phi_{S_1} | w). \quad [3]$$

The speaker conveys a particular weighting of informational

vs. social goals (ϕ_{S_1}) by considering the beliefs of listener L1, who hears an utterance and jointly infers the speaker's utilities and the true state of the world:

$$P_{L_1}(s, \hat{\phi}|w) \propto P_{S_1}(w|s, \hat{\phi}) \cdot P(s) \cdot p(\hat{\phi}). \quad [4]$$

The presentational utility is the highest-order term of the model, defined only for a speaker thinking about a listener who evaluates a speaker (i.e., defined for S_2 , but not S_1).

Finally, more complex utterances incur a greater cost, $C(w)$ – capturing the general pressure towards economy in speech. In our work, utterances with negation (e.g., *not terrible*) are assumed to be slightly costlier than their equivalents with no negation (and this cost differential is inferred from data; see Supplementary Information).

Within our experimental domain, we assume there are four possible states of the world corresponding to the value placed on a particular referent (e.g., the presentation the speaker is commenting on): $S = s_1, \dots, s_4$. Since the rating scale is relatively abstract, we assume a uniform prior distribution over possible states of the world. The set of utterances is $\{\text{terrible, bad, good, amazing, not terrible, not bad, not good, and not amazing}\}$. We implemented this model using the probabilistic programming language WebPPL (22).

The model exhibits and explains interpretable behavior. If Bob's pie was good, Ann's utilities align to produce a positive utterance. Saying *[Your pie] was amazing* simultaneously is truthful, kind, and appears as both. If Bob's pie was poor, however, the speaker is in a bind: Ann could be kind and say *It was great*, but at the cost of conveying the wrong information to Bob if he believes her to be truthful. If he does not, he might infer Ann is *just being nice*, but is uninformative. Alternatively, she could say the truth (*It was bad*), but then Bob would think Ann didn't care about him. What is a socially-aware speaker to do? Our quantitative model predicts that indirect speech – like *It wasn't bad* – best navigates Ann's dilemma. Her statement is sufficiently open-ended to include the possibility that the pie was good, but her avoidance of the simpler and less costly *It was good* provides both an inference available to Bob that the pie was mediocre and that Ann cares about his feelings by not saying the blunt truth.

Results

By instantiating our running example in an online experiment, we made a direct, fully pre-registered test of our model and its performance in comparison to a range of alternative models ($N = 202$). Participants read scenarios with information on the speaker's (Ann's, in our example) feelings toward some performance or product (e.g., pie a poem recital; *true state*), on a scale from zero to three hearts (e.g., one out of three hearts). For example, one trial read: *Imagine that Bob gave a poem recital, but he didn't know how good it was. Bob approached Ann, who knows a lot about poems, and asked "How was my poem?"*

We manipulated the speaker's goal across trials: to be *informative* (give accurate and informative feedback); to be *kind* (make the listener feel good); or to be both informative and kind simultaneously. We hypothesized that each of the three experimentally-induced goals would induce a different tradeoff between the three utilities in our model (see Supplementary Information). In a single trial, each scenario was followed by a question asking for the most likely produced utterance by

Table 1. Comparison of variance explained for each model variant and log Bayes Factors quantifying evidence in favor of alternative model in comparison.

Model	Variance explained	log BF
model: informational, social, presentational	0.97	–
model: informational, presentational	0.96	-11.14
model: informational, social	0.92	-25.06
model: social, presentational	0.23	-864
model: presentational only	0.23	-873.83
model: social only	0.22	-885.52
model: informational only	0.83	-274.89

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*.

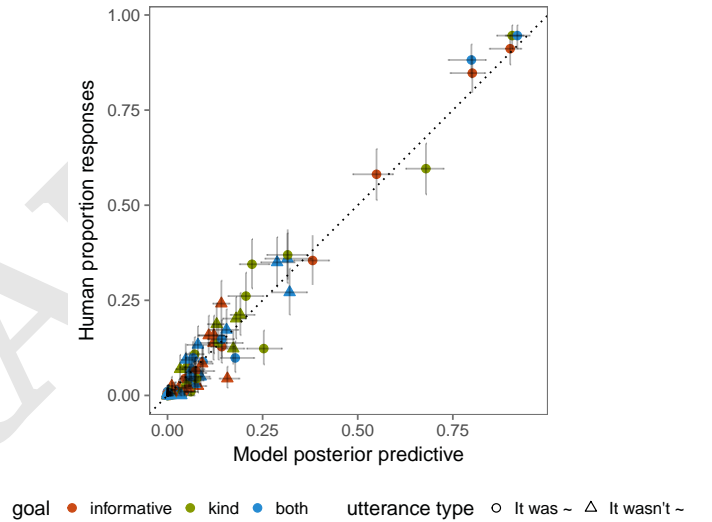


Fig. 2. Full distribution of human responses vs. model predictions. Error bars represent 95% confidence intervals for the data (vertical) and 95% highest density intervals for the model (horizontal).

Our primary behavioral hypothesis was that speakers describing bad states (e.g., Bob's pie deserved 0 hearts) with goals to be both informative and kind would produce more indirect, negative utterances (e.g., *It wasn't terrible*). Such indirect speech acts both save the listener's face and providing a little information about the true state, and thus, are the kind of thing that a socially-conscious speaker would say. This prediction was confirmed: a Bayesian mixed-effects model predicts more negation as a function of true state and goal via an interaction such that a speaker with both goals to be informative and kind produced more negation in worse states compared to a speaker with only the goal to be informative ($M = -1.33, [-1.69, -0.98]$) and goal to be kind ($M = -0.50, [-0.92, -0.07]$). Rather than eschewing one of their goals to increase utility along a single dimension, participants chose utterances that jointly satisfied their conflicting goals by producing indirect, polite speech.

The model parameters (softmax parameters and each goal condition's utility weights) can be inferred from the behavioral data using a Bayesian data analysis model ((23); see Supple-

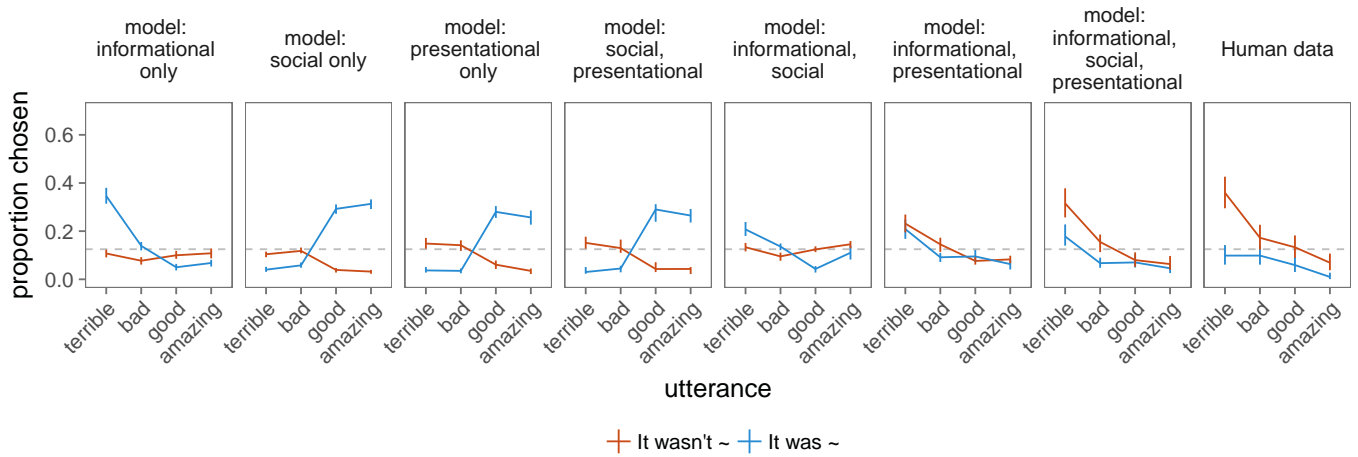


Fig. 3. Comparison of predictions for proportion of utterances chosen by pragmatic speaker from possible model variants (left) and human data (rightmost) for average proportion of negation produced among all utterances, given true state of 0 heart (on a scale of 0 to 3) and speaker with both goals to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

Table 2. Inferred ϕ parameters from all model variants with more than one utility.

Model (utilities)	goal	ϕ_{inf}	ϕ_{soc}	ϕ_{pres}	ϕ_{S1}
inf, soc, pres	both	0.36	0.11	0.54	0.36
inf, soc, pres	informative	0.36	0.02	0.62	0.49
inf, soc, pres	social	0.25	0.31	0.44	0.37
inf, pres	both	0.64	—	0.36	0.17
inf, pres	informative	0.77	—	0.23	0.33
inf, pres	social	0.66	—	0.34	0.04
inf, soc	both	0.54	0.46	—	—
inf, soc	informative	0.82	0.18	—	—
inf, soc	social	0.39	0.61	—	—
soc, pres	both	—	0.38	0.62	0.55
soc, pres	informative	—	0.35	0.65	0.75
soc, pres	social	—	0.48	0.52	0.66

Being informative pushes the weight on social utility close to zero, but the weight on *appearing kind* stays high, suggesting that speakers are expected to manage their own face even when they are not considering others'. *Kind and informative* speakers emphasize informativity slightly more than kindness. In all cases, however, the presentational utilities have greatest weight, which may suggest that appearing honest and kind is more important than actually being so! Overall then, our condition manipulation altered the balance between these weights, but all utilities played a role in all conditions.

Discussion

Politeness is puzzling from an information-theoretic perspective. Incorporating social motivations adds a level of explanation, but so far such intuitions and observations have resisted both formalization and precise testing. We present a utility-theoretic model of language use that captures the interplay between competing informational, social, and presentational goals, and provide preregistered experimental evidence that confirmed its ability to capture human judgments, unlike comparison models that used only a subset of the full utility structure.

To precisely estimate choice behavior in the experiment, it was required to abstract away from natural interactions in a number of ways. Human speakers have access to a potentially infinite set of utterances to select from in order to manage the three-utility tradeoff (*It's hard to write a good poem, That metaphor in the second stanza was so relatable!*). In theory, each utterance will have strengths and weaknesses relative to the speaker's goals, though computation in an unbounded model presents technical challenges (perhaps paralleling the difficulty human speakers feel in finding the right thing to say in a difficult situation; see (14)).

For a socially-conscious speaker, managing listeners' inferences is a fundamental task. Our work extends previous models of language beyond standard informational utilities to address social and self-presentational concerns. Further, our model builds upon the theory of politeness as face management (12) and takes a step towards understanding the complex set of social concerns involved in face management. Our approach can provide insight into a wide range of social

behaviors beyond speech by considering utility-driven inferences in a social context (24, 25) where agents need to take into account concerns about both self and others.

Previous game-theoretic analyses of politeness have either required some social cost to an utterance (e.g., by reducing one's social status or incurring social debt to one's conversational partner; (26)) or a separately-motivated notion of plausible deniability (27). The necessary kind of utterance cost for the first type of account would involve higher-order reasoning about others, and may be able to be defined in terms of the more basic social and self-presentational goals we formalize here. A separate notion of plausible deniability may not be needed to explain politeness, either. Maintaining plausible deniability is in one's own self-interest (e.g., due to controversial viewpoints or covert deception) and goes against the interests of the addressee. Thus, some amount of utility dis-alignment is presumed by these accounts, but intuitively politeness behavior appears present even in the absence of obvious conflict. In our work here, we show that such behaviors can arise from purely cooperative goals (12), though this account is not inconsistent with plausible deniability accounts in cases of conflict.

Utility weights and value functions in our model could provide a framework for understanding systematic cross-cultural differences in what counts as polite. Cross-cultural differences in politeness could be a product of different weightings within the same utility structure. Alternatively, culture could affect the value function V that maps states of the world onto subjective values for the listener (e.g., the mapping from states to utilities may be nonlinear and involve reasoning about the future). Our formal modeling approach with systematic behavior measurements provides an avenue towards understanding the vast range of politeness practices found across languages.

Politeness is but one of the ways that language use deviates from purely informational transmission. We flirt, insult, boast, and empathize by balancing informative transmissions with goals to affect others' feelings or present particular views of ourselves. Our work shows how social and self-presentational motives are integrated with informational concerns more generally, opening up the possibility for a broader theory of social language. Finally, a formal formal account of politeness moves us closer to courteous computation – to machines that can talk with tact.

Materials and Methods

Literal semantic task. We probed judgments of literal meanings of the target words assumed by our model and used in our main experiment. 51 participants with IP addresses in the United States were recruited on Amazon's Mechanical Turk. We used thirteen different context items in which a speaker 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; see Figure 5 for an example of the heart scale). 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 four possible words: *terrible*, *bad*, *good*, and *amazing*, giving rise to eight 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 speaker production 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 (i.e., a standard Beta-Binomial model). Meanings of the words as judged by participants were as one would expect (Figure 4).

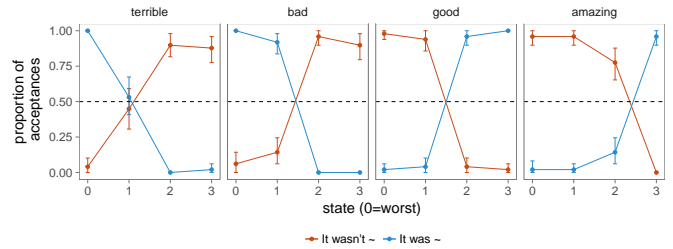


Fig. 4. Semantic measurement results. Proportion of acceptances of utterance types (shown in different colors) combined with target words (shown in different facets) given the true state represented on a scale of hearts. Error bars represent 95% confidence intervals.

Imagine that Fiona filmed a movie, but she didn't know how good it was. Fiona approached Yvonne, who knows a lot about movies, and asked "How was my movie?"

Here's how Yvonne **actually** felt about Fiona's movie, on a scale of 0 to 3 hearts:



If Yvonne wanted to **BOTH** make Fiona feel good **AND** give accurate and informative feedback,

what would Yvonne be most likely to say?

"It wasn't terrible"

Fig. 5. Example of a trial in the speaker production task.

Speaker production task. 202 participants with IP addresses in the United States were recruited on Amazon's Mechanical Turk. As in the literal semantic task 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 (Figure 5). 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; Figure 5). Each participant read twelve scenarios, depicting every possible combination of the three goals and four 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*.

Data availability. Our model, preregistration of hypotheses, procedure, data, and analyses are available at https://github.com/ejyoony/polite_speaker.

ACKNOWLEDGMENTS. This work was supported by NSERC PGS Doctoral scholarship PGSD3-454094-2014 to EJY, NSF Gradu-

ate Research Fellowship DGE-114747 to MHT, ONR grant N00014-13-1-0788 to NDG, and NSF grant BCS 1456077 to MCF.

1. Bühler K (1934) *Sprachtheorie*. (Oxford, England: Fischer).
2. Frank MC, Goodman ND (2012) Predicting pragmatic reasoning in language games. *Science* 336(6084):998–998.
3. Jakobson R (1960) Linguistics and poetics in *Style in language*. (MA: MIT Press), pp. 350–377.
4. Shannon CE (1948) A mathematical theory of communication. *Bell Syst. Tech. J.* 27:623–656.
5. Grice HP (1975) *Logic and conversation*, eds. Cole P, Morgan JL. (Academic Press) Vol. 3, pp. 41–58.
6. Searle J (1975) *Indirect Speech Acts*, eds. Cole P, Morgan JL. (Academic Press) Vol. 3, pp. 59–82.
7. Axia G, Baroni MR (1985) Linguistic politeness at different age levels. *Child Development* pp. 918–927.
8. Clark HH, Schunk DH (1980) Polite responses to polite requests. *Cognition* 8(2):111–143.
9. Holtgraves T (1997) Yes, but... positive politeness in conversation arguments. *Journal of Language and Social Psychology* 16(2):222–239.
10. Leech G (1983) *Principles of pragmatics*. (London, New York: Longman Group Ltd.).
11. Ide S (1989) Formal forms and discernment: Two neglected aspects of universals of linguistic politeness. *Multilingua-journal of cross-cultural and interlanguage communication* 8(2-3):223–248.
12. Brown P, Levinson SC (1987) *Politeness: Some universals in language usage*. (Cambridge university press) Vol. 4.
13. Goffman E (1967) *Interaction ritual: essays on face-to-face interaction*. (Aldine).
14. Goodman ND, Frank MC (2016) Pragmatic language interpretation as probabilistic inference. *Trends in Cognitive Sciences* 20(11):818–829.
15. Lassiter D, Goodman ND (2017) Adjectival vagueness in a bayesian model of interpretation. *Synthese* 194(10):3801–3836.
16. Kao JT, Wu JY, Bergen L, Goodman ND (2014) Nonliteral understanding of number words. *Proceedings of the National Academy of Sciences* 111(33):12002–12007.
17. Kao JT, Goodman ND (2015) Let's talk (ironically) about the weather: Modeling verbal irony in *Proceedings of the 37th Annual Conference of the Cognitive Science Society*.
18. Baker CL, Saxe R, Tenenbaum JB (2009) Action understanding as inverse planning. *Cognition* 113(3):329–349.
19. Jara-Ettinger J, Gweon H, Schulz LE, Tenenbaum JB (2016) The naïve utility calculus: Computational principles underlying commonsense psychology. *Trends in cognitive sciences* 20(8):589–604.
20. Liu S, Ullman TD, Tenenbaum JB, Spelke ES (2017) Ten-month-old infants infer the value of goals from the costs of actions. *Science* 358(6366):1038–1041.
21. Goodman ND, Stuhlmüller A (2013) Knowledge and implicature: Modeling language understanding as social cognition. *Topics in cognitive science* 5(1):173–184.
22. Goodman ND, Stuhlmüller A (2014) The Design and Implementation of Probabilistic Programming Languages (<http://dippl.org>).
23. Lee MD, Wagenmakers EJ (2014) *Bayesian Cognitive Modeling: A Practical Course*. (Cambridge Univ. Press).
24. Baker CL, Jara-Ettinger J, Saxe R, Tenenbaum JB (2017) Rational quantitative attribution of beliefs, desires and percepts in human mentalizing. *Nature Human Behaviour* 1(4):0064.
25. Hamlin KJ, Ullman TD, Tenenbaum JB, Goodman ND, Baker CL (2013) The mentalistic basis of core social cognition: Experiments in preverbal infants and a computational model. *Developmental science* 16(2):209–226.
26. Van Rooy R (2003) Being polite is a handicap: Towards a game theoretical analysis of polite linguistic behavior in *Proceedings of the 9th conference on Theoretical aspects of rationality and knowledge*. (ACM), pp. 45–58.
27. Pinker S, Nowak MA, Lee JJ (2008) The logic of indirect speech. *Proceedings of the National Academy of sciences* 105(3):833–838.