

Let's talk (ironically) about the weather: A computational model of verbal irony

Justine T. Kao (justinek@stanford.edu), Noah D. Goodman (ngoodman@stanford.edu)

Department of Psychology, Stanford University

Abstract

Verbal irony plays an important role in how we communicate and express our opinions about the world. While there exist many interesting theories and empirical findings about how people use and understand verbal irony, there is to our knowledge no formal model of how people incorporate shared background knowledge and linguistic information to communicate ironically. Here we present two behavioral experiments that examine people's interpretations of utterances given different contexts. We then describe a computational model that reasons about background knowledge, affect, and the speaker's communicative goals to interpret ironic utterances and their rich affective subtexts. We show that by accounting for two types of affect goals—valence and arousal—our model produces interpretations that closely match humans'. Finally, we discuss the implications of our model on irony and its relationship to other types of nonliteral language understanding.

Keywords: irony; computational modeling; pragmatics; nonliteral language understanding

Introduction

For better or for worse, verbal irony—defined as utterances whose apparent meanings are opposite in polarity to the speaker's intended meaning (Roberts & Kreuz, 1994; Colston & O'Brien, 2000)—is a major figurative trope of our time. From popular sitcoms to political satire to *#sarcasm* on Twitter and casual conversations among friends, verbal irony plays an important role in how we communicate and express opinions about the world. The prevalence of verbal irony poses a puzzle for theories of language understanding: Why would speakers ever use an utterance to communicate its opposite meaning, and how can listeners appropriately interpret such an utterance? Previous work has shown that verbal irony serves several important communicative goals, such as to heighten or soften criticism (Colston, 1997b), elicit emotional reactions (Leggett & Gibbs, 2000), highlight group membership (Gibbs, 2000), and express affective attitudes (Colston & Keller, 1998; Colston, 1997a). These findings suggest that while ironic statements are false under their literal meanings, they are often highly informative with respect to social and affective meanings. In this paper, we present a computational model and behavioral experiments to show that people may use inferences about these alternative dimensions of meaning and the speaker's communicative goals to understand ironic utterances.

Linguists and psychologists have proposed several informal theories of how people understand verbal irony. According to a classic Gricean analysis, listeners first need to recognize that an ironic utterance blatantly violates the maxim of quality (to be truthful); they then arrive at a conversational implicature that the intended meaning is contrary to the utterance's literal meaning (Grice, 1967; Wilson, 2006). While Grice's account is appealing in its treatment of verbal irony

as arising naturally from conversational maxims, it does not provide a detailed or satisfactory explanation for how the appropriate implicature is derived from these maxims, or why it is ever rational to deliver an utterance that is opposite from the truth (Wilson, 2006). On the other hand, previous work suggests that it is possible and indeed desirable to consider nonliteral language understanding as a product of general principals of communication, such as reasoning about informativeness with respect to the speaker's communicative goals (Kao, Wu, Bergen, & Goodman, 2014; Kao, Bergen, & Goodman, 2014). In particular, a model that reasons about the speaker's affect is able to interpret hyperbolic utterances and infer the appropriate affective subtext (Kao, Wu, et al., 2014). Given the fact that many researchers believe hyperbole and irony to be closely related phenomena (cite), we suggest that similar principles may account for verbal irony understanding as well. Our goal in this paper is to identify these communicative principles and provide a precise formal account of how they interact to produce ironic interpretations.

Rational Speech Act (RSA) models are a family of computational models that formalize language understanding as recursive reasoning between speaker and listener, and have been shown to account for many phenomena in pragmatics (Frank & Goodman, 2012; Goodman & Stuhlmüller, 2013). Kao, Wu, et al. (2014) introduces a critical extension to basic RSA models by considering the idea that speakers may aim to address different questions under discussion (QUDs) when formulating an utterance. An important task for the listener is to then jointly infer the QUD as well as the speaker's intended meaning. For example, a speaker may want to communicate negative affect about a situation (e.g. unhappiness about the temperature outside) instead of the precise situation (e.g. the temperature outside), in which case choosing an exaggerated utterance (e.g. "It's freezing outside!") effectively addresses the QUD. A listener who reasons about the speaker and QUD is then able to use his background knowledge about temperatures to correctly infer that the speaker is upset about the temperature, but that it is unlikely to be literally freezing outside (especially if she is in California). Extending this model to consider QUDs opens up the possibility for a speaker to produce an utterance that is literally false but satisfies her goal to convey affect. While this model—which we will refer to as qRSA—produces nonliteral interpretations of hyperbolic utterances that closely match humans', Kao, Wu, et al. (2014) considered only a unidimensional space of affects, namely the presence or absence of negative feeling. This overlooks the range of attitudes and emotions that speakers could express with nonliteral utterances. In particular, since verbal irony involves expressing negative meanings with positive utterances and vice versa, a richer space of affect that includes both posi-

tive and negative emotions is necessary. Here we examine the consequence of considering a range of emotions in an empirically derived affect space within the qRSA model, and show that this minimal change is able to capture many of the rich inferences resulting from verbal irony.

In what follows, we will examine interpretations of potentially ironic utterances in an innocuous domain—the weather. We chose the weather as the victim of irony for several reasons. First, people are quite familiar with talking (and complaining) about the weather. Second, we can visually represent the weather to participants with minimal linguistic description in order to obtain measures of nonlinguistic contextual knowledge. Finally, we can vary the weather states to observe how the same utterance is interpreted differently given different contextual knowledge. We first describe an extension to the qRSA model and show that an enriched space of affect enables the model to produce ironic interpretations. We then present two behavioral experiments that examine people’s interpretations of utterances given different weather contexts. We show that by accounting for two types of affect goals, valence and arousal, our model produces interpretations that closely match humans’. Finally, we discuss the implications of our model on irony and its relationship to other types of nonliteral language understanding.

Computational Model

In this section, we describe the qRSA model and compare different spaces of affect to test the conditions for producing ironic interpretations. Following the qRSA model described in Kao, Wu, et al. (2014), a speaker chooses an utterance that most effectively communicates information regarding the question under discussion (QUD) to a literal listener. We consider a meaning space consisting of the variables s, A , where s is the state of the world, and A represents the speaker’s (potentially multidimensional) affect towards the state. We formalize a QUD as a projection from the full meaning space to the subset of interest to the speaker, which could be s or any of the dimensions of A . We define the speaker’s utility as the negative surprisal of the true state under the listener’s distribution given an utterance along the QUD dimension, leading to the following utility function¹:

$$U(u|s, A, q) = \log \sum_{s', A'} \delta_{q(s, A) = q(s', A')} L_{\text{literal}}(s, A|u) \quad (1)$$

where q is the QUD and L_{literal} is the literal listener. The speaker S chooses an utterance according to a softmax decision rule (?, ?):

$$S(u|s, A, q) \propto e^{\lambda U(u|s, A, q)}, \quad (2)$$

where λ is an optimality parameter. A pragmatic listener $L_{\text{pragmatic}}$ then takes into account prior knowledge and his internal model of the speaker to determine the state of the world as well as the speaker’s affect. Because $L_{\text{pragmatic}}$ is uncertain

¹See Kao, Wu, et al. (2014) for details.

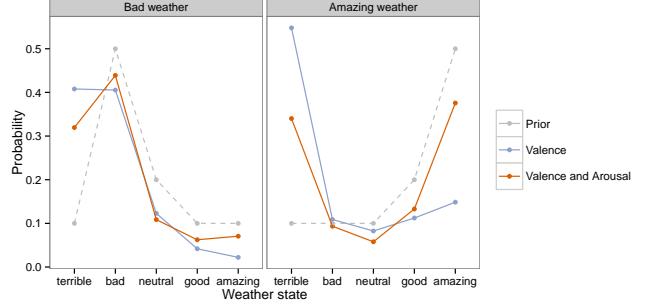


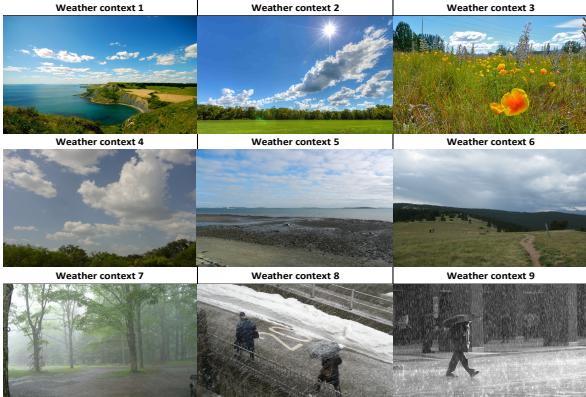
Figure 1: Model interpretations of “The weather is terrible” given different prior beliefs about the weather state and affect dimensions. Gray dotted lines indicate prior beliefs about weather states given a weather context; blue lines indicate interpretations when reasoning only about the speaker’s valence; orange lines indicate interpretations when reasoning about both valence and arousal.

about the QUD, he marginalizes over the possible QUDs under consideration:

$$L_{\text{pragmatic}}(s, A|u) \propto P(s)P(A|s) \sum_q P(q)S(u|s, A, q)$$

The resulting distribution over world states and speaker affects is an *interpretation* of the utterance.

We examine the model’s behavior using affect spaces with different dimensions. We first consider a one-dimensional affect space, where the dimension is emotional valence—whether the speaker feels negative or positive about the world state. Suppose there is strong prior belief that the weather state is bad; also suppose that the QUD is the speaker S ’s emotional valence towards the weather. Based on S ’s understanding of the literal listener’s prior knowledge, she knows that if she produces the utterance “The weather is terrible,” L_{literal} will believe that the weather is terrible and that the speaker is likely to feel negative about it. Since the QUD is successfully addressed if L_{literal} believes that S feels negative towards the weather state, S is motivated to produce the utterance “The weather is terrible.” However, suppose that the pragmatic listener $L_{\text{pragmatic}}$ has strong prior belief that the weather state is bad. Since $L_{\text{pragmatic}}$ reasons about S and her goals, he realizes that S chose the utterance “The weather is terrible” to communicate her negative affect and not the true state of the weather. He will then infer that the weather is likely bad, and that S is extremely likely to feel negative towards it, which successfully produces a hyperbolic interpretation. However, suppose there is strong prior belief that the weather is amazing. Given the utterance “The weather is terrible,” $L_{\text{pragmatic}}$ interprets it literally to mean that the weather is indeed terrible. This is because if it were *not* terrible, then it would be unlikely for S to choose the utterance, as it is unlikely to communicate either the true state of the world or her valence. The blue lines in Figure 1 show simulations of the model with this one-dimensional affect space. While the



(a) The nine weather images shown to participants in Experiments 1 and 2.

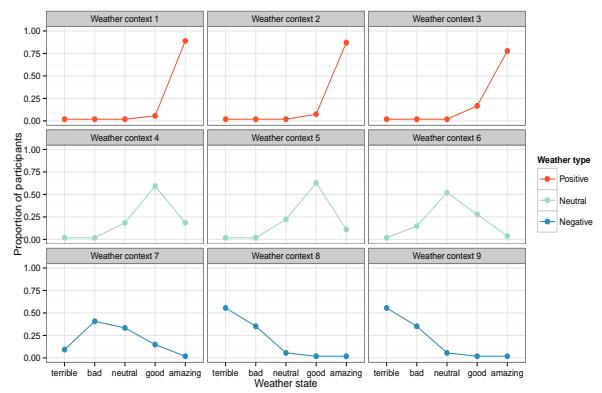


Figure 2: Nine weather contexts and their empirically measured priors over weather states.

model produces a hyperbolic interpretation given strong prior belief about *bad* weather (left panel), it produces a literal interpretation given strong prior belief about *amazing* weather (right panel). In other words, a model that only considers a single affect dimension (valence) is unlikely to infer a positive world state from a highly negative utterance (and vice versa), thus failing on many cases of verbal irony.

We now consider a more complex affect space with two dimensions—valence and arousal—to observe its consequence on interpretation. Affective science identifies valence and arousal as two main dimensions underlying the slew of emotions that people experience (cite). For example, *anger* is a negative valence and high arousal emotion, while *contentment* is a positive valence and low arousal emotion (cite). We suggest that perhaps speakers leverage the arousal dimension to convey high arousal and positive affect (e.g. excitement) using utterances whose literal meanings are associated with high arousal and negative affect (e.g. “The weather is terrible”). The orange lines in Figure 1 show simulations of the qRSA model with a two-dimensional affect space, valence and arousal. Given strong prior belief that the weather state is *bad*, the model interprets “The weather is terrible” to mean that the weather is likely to be *bad*, producing a hyperbolic interpretation. However, given strong prior belief that the weather is *amazing*, the model now interprets “The weather is terrible” ironically to mean that the weather is likely *amazing*. This is because with the enriched two-dimensional affect space, the pragmatic listener realizes that the speaker may be using “terrible” to communicate high emotional arousal. These model simulations suggest that reasoning about a two-dimensional affect space motivated by emotion theory enables the model to appropriately interpret ironic utterances.

To produce an interpretation of an utterance in context, the model requires the following input values: (1) $P(s)$: the prior probability of a weather state s given a weather context. (2) $P(A|s)$: the probability of affect A (positive/negative valence

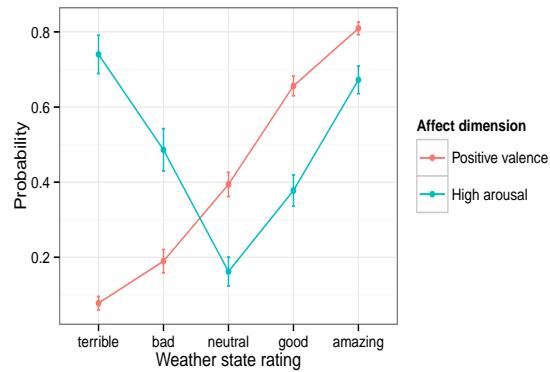


Figure 3: The average probabilities of positive valence and high arousal associated with each weather state. Error bars are 95% confidence intervals.

and high/low arousal) given a weather state. (3) $P(q)$: the prior probability of a particular QUD (4) The speaker optimality parameter λ . We derived the values for (1) and (2) from Experiment 1 and fit (3) and (4) to the data from Experiment 2, which we describe below.

Behavioral Experiments

To quantitatively test our model with the enriched affect space described above, we conducted the following two experiments. In Experiment 1, we measured the prior beliefs over weather states for various weather contexts. We also measured various emotions associated with different weather contexts in order to empirically extract the affective dimensions relevant to this domain. In Experiment 2, we collected people’s ratings of how a speaker perceives and feels about the weather given what she says (e.g. “The weather is terrible!” when the context clearly depicts sunny weather).

Experiment 1: Prior elicitation

Materials and methods We selected nine images from Google Images that depict the weather. To cover a range of weather states, three of the images were of sunny weather, three of cloudy weather, and three of rainy or snowy weather. We refer to these images as weather contexts. Figure 2a shows these nine images. 49 native English speakers with IP addresses in the United States were recruited on Amazon’s Mechanical Turk. Each participant saw all nine images in random order. In each trial, participants were told that a person (e.g. Ann) looks out the window and sees the view depicted by the image. They then indicated how Ann would rate the weather using a labeled 5-point Likert scale, ranging from *terrible*, *bad*, *neutral*, *good*, to *amazing*. Finally, participants used slider bars to rate how likely Ann is to feel each of the following seven emotions about the weather: *excited*, *happy*, *content*, *neutral*, *sad*, *disgusted*, and *angry*. The order of the emotions was randomized for each participant but remained consistent across trials for the same participant. The end points of the slider bars were labeled as “Impossible” and “Absolutely certain.” A link to the experiment is here: http://stanford.edu/~justinek/irony_exp/priors/priors.html

Results For each of the nine weather contexts, we obtained the number of participants who gave each of the weather state ratings and performed add-one Laplace smoothing on the counts. This allowed us to compute a smoothed prior distribution over weather states given each context. From Figure 2b, we see that the sunny and positive weather contexts were more likely to be rated as *amazing*, while the negative weather contexts were more likely to be rated as *bad* or *terrible*.

To examine participants’ ratings of the affects associated with each context, we first performed Principal Component Analysis (PCA) on the seven emotion category ratings. This allowed us to compress the ratings onto a lower-dimensional space and reveal the main affective dimensions that are important in this domain. We found that the first two principal components accounted for 69.14% and 13.86% of the variance in the data. In addition, these components corresponded roughly to the dimensions of emotional valence (positive or negative) and emotional arousal (high or low). To approximate the probabilities of Ann feeling positive or negative affect and high or low arousal given different weather states, we converted the PCA scores into probabilities as follows. We first normalized the scores in each dimension to have zero mean and unit variance. Treating these normalized scores as quantiles of a standard normal distribution, we used the cumulative distribution function to convert the normalized scores into values between 0 and 1. Figure 3 shows the probability of positive valence and high arousal given each weather state.

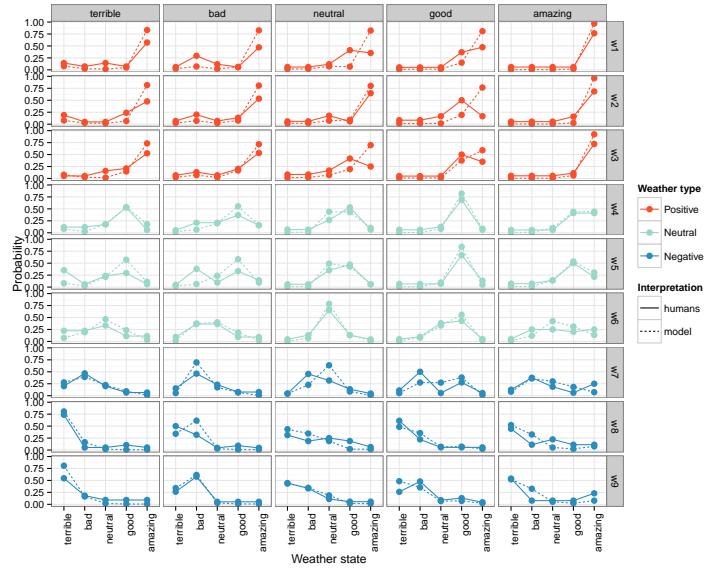
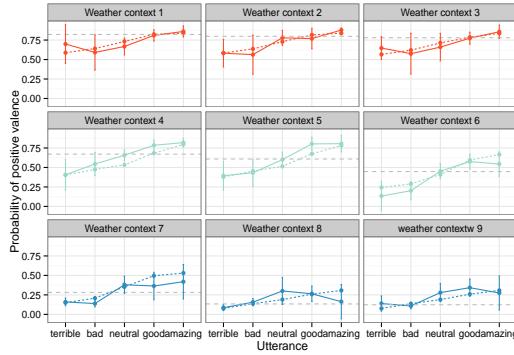


Figure 4: Model’s and participants’ inferences about the weather state (x-axis) given a weather context (row) and an utterance (column). Each panel represents interpretations of an utterance in a weather context. The solid lines are participants’ ratings; the dotted lines are model’s posterior distributions over weather states.

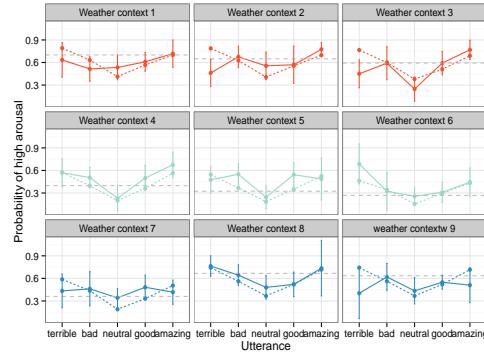
Experiment 2: Irony understanding

In Experiment 1, we obtained the prior distribution over weather states for each weather context as well as the prior probabilities of positive valence and high arousal given each weather state. This gave us the necessary components to generate interpretations of utterances from our model. Here we describe an experiment that elicits people’s interpretations of utterances, which we then use to evaluate model predictions.

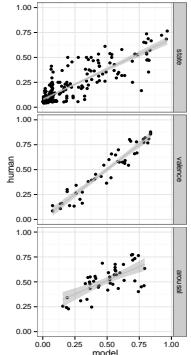
Materials and methods 59 native English speakers with IP addresses in the United States were recruited on Amazon’s Mechanical Turk. Each participant saw all nine images from Figure 2a in random order. In each trial, participants were told that a person (e.g. Ann) and her friend are in a room looking out the window together and see the view depicted by the image. Ann says, “The weather is _____!” where the adjective is randomly selected at each trial from the following set: “*terrible*,” “*bad*,” “*ok*,” “*good*,” and “*amazing*.” Participants first rated how likely it is that Ann’s statement is ironic using a slider with end points labeled “Definitely NOT ironic” and “Definitely ironic.” They then indicated how Ann would actually rate the weather using a labeled 5-point Likert scale, ranging from *terrible*, *bad*, *neutral*, *good*, to *amazing*. Finally, participants used sliders to rate how likely it is that Ann feels each of seven emotions about the weather. A link to the experiment is here: http://stanford.edu/~justinek/irony_exp/interpretation/interpretation_askIrony.html



(a) Average probabilities of speaker feeling positive valence given her utterance in a weather context.



(b) Average probabilities of speaker feeling high arousal given her utterance in a weather context.



(c) Scatter plot of human versus model interpretations.

Results We first examined participants’ irony ratings for each of the weather context and utterance pairs. We found a basic irony effect, where utterances whose polarities are inconsistent with the polarity of the weather context are rated as significantly more ironic than utterances whose polarities are consistent with the weather context (STATS) (Figure XX?). For example, “The weather is terrible” (a negative utterance) is rated as more ironic in weather context 1 (positive context) than in weather context 7 (negative context). A linear regression model with the polarity of the utterance, the polarity of the weather context, and their interaction as predictors of irony ratings produced an adjusted R-squared of 0.91, capturing most of the variance in the data. This suggests that participants’ lay judgments of irony align with its basic definition: utterances whose apparent meanings are opposite in polarity to the speaker’s intended meaning.

Given the fact that participants can identify verbal irony based on its inconsistency with context, how do they then use context to determine the speaker’s intended meaning? We examined participants’ interpretations of utterances given contexts. For each of the 45 weather context (9) \times utterance (5) pairs, we obtained the number of participants who gave each of the five weather state ratings (terrible, bad, neutral, good, amazing). We performed add-one Laplace smoothing on the counts to obtain a smoothed distribution over weather states given each context and utterance. The solid lines in Figure 4 show these distributions of ratings. We see that participants produce ironic interpretations of utterances, such that the weather is most likely to be amazing given that the speaker said “The weather is terrible” in weather context 1. Participants also produce hyperbolic interpretations, such that the weather is most likely to be bad given that the speaker said “The weather is terrible” in weather context 7. This suggests that people are highly sensitive to context when interpreting utterances, and use it both to determine when an utterance is not meant literally and to appropriately recover the intended meaning.

Finally, we examine participants’ inferences about the speaker’s affect given utterances in context. We used the load-

ings from the PCA on emotion ratings from Experiment 1 to project the emotion ratings from Experiment 2 onto the same dimensions. We then normalized and used the cumulative distribution function to convert the scores into values between 0 and 1. Figure 5a shows the average probability of *positive valence* given an utterance in a weather context. The dotted gray lines are the average probabilities of positive valence associated with each weather context without any linguistic input, taken from Experiment 1. Figure 5b shows the average probability of *high arousal* given an utterance in a weather context. The dotted gray lines are the average probabilities of high arousal associated with each weather context without any linguistic input, taken from Experiment 1.

JTK: figure out data take-home point here

Model Evaluation

We now evaluate the model’s performance against these behavioral results. From Experiment 1, we obtained the prior probability of a weather state given a context as well as the probability of affect given a weather state. After fitting four free parameters to maximize correlation with data from Experiment 2 ($\lambda = 1$, $P(q_{state}) = 0.2$, $P(q_{valence}) = 0.3$, $P(q_{arousal}) = 0.4$), the model produced an interpretation for each of the 45 utterance and weather context pairs. Each interpretation is a joint posterior distribution $P(s, A|u)$, where the affect A can be further broken down into valence and arousal dimensions. We will examine the model’s performance on each of these state and affect dimensions by marginalizing over the other dimensions.

Figure 4 shows participants’ and the model’s inferences about the actual weather state given an utterance and a weather context. The model predictions closely match humans’ interpretations, with a correlation of 0.86 (Figure 5c). Figure 5a shows participants’ and the models’ inference about the speaker’s valence given an utterance and a weather context. From the tight correspondence between model and human interpretations of valence ($r = 0.96$), we see that the model is able to incorporate the valence associated with the utterance’s literal meaning (e.g. “The weather is terrible”)

is the model using the prior too strongly here?

and the valence associated with the weather context (e.g. weather context 1) to interpret the probability of the speaker feeling positive valence. In other words, the model infers the appropriate valence even when it is inconsistent with the valence of the utterance's literal meaning, thus capturing the essence of irony. Finally, Figure 5b shows participants' and the model's inferences about the speaker's arousal. The model's prediction for emotional arousal match humans' with a correlation of 0.66. These results suggest that the model is able to incorporate background knowledge and reasoning about multiple affective goals to produce the appropriate ironic interpretations as well as the associated affects.

Discussion

In this paper, we explored the consequences of expanding the space of affect considered in Rational Speech Act models to account for verbal irony. We showed that by making a minimal extension to Kao, Wu, et al. (2014)'s hyperbole model, we can capture people's fine-grained interpretations of ironic utterances. The similarities between these models suggest that understanding hyperbole and irony requires the same underlying principles of communication, which aligns with other informal accounts of the pragmatics of nonliteral language understanding (cite).

While we present evidence of shared communicative principles that unify irony and other forms of nonliteral language, there remain important qualities of verbal irony that may be unique. For example, the echoic mention theory of irony claims that speakers often use verbal irony to remind the listener of previous utterances that turned out to be false or irrelevant, or of positive norms that were explicitly violated (Sperber & Wilson, 1981; Jorgensen, Miller, & Sperber, 1984). On the other hand, pretense theory argues that when a speaker produces an ironic utterance, she is not genuinely making the utterance, but only pretending to be someone who would make such an utterance (Clark & Gerrig, 1984). While our model so far does not explicitly incorporate elements of echoic mention or pretense, it is able to capture many of the main characteristics of verbal irony. By addressing these additional components in future research, we hope to further improve our model's performance and enrich its understanding of the social aspects of irony.

In addition to shedding light on the communicative principles underlying irony understanding, our work also has interesting connections to natural language processing. Given the prevalence of irony in natural language, many researchers aim to automatically detect sarcasm in large bodies of texts in order to recover the correct sentiment from an ostensibly positive or negative utterance (e.g. "I was overjoyed to pay \$30 for an overcooked steak") (Davidov, Tsur, & Rappoport, 2010; Filatova, 2012). A critical insight that emerged from these efforts is that irony detection requires information far beyond surface linguistic cues, often calling upon a deep understanding of context and common knowledge

between speaker and listener that computers currently lack (González-Ibáñez, Muresan, & Wacholder, 2011; Wallace, Do Kook Choe, & Charniak, 2014). By integrating background knowledge and linguistic meaning in a principled manner, we present a formal model that predicts ironic interpretations in a way that is highly sensitive to context and common ground.

We believe that our experimental paradigm and modeling framework lend itself well to a more detailed and precise account of irony understanding. Given the prevalence of irony in everyday language and the social functions it serves, it would be *amazing* (literally) to understand people's ability to use and interpret utterances that mean the opposite of what they say.

References

- Clark, H. H., & Gerrig, R. J. (1984). On the pretense theory of irony. *Metaphor and Symbol*, 12(1), 43–58.
- Colston, H. L. (1997a). "i've never seen anything like it": Overstatement, understatement, and irony. *Metaphor and Symbol*, 12(1), 43–58.
- Colston, H. L. (1997b). Salting a wound or sugarizing a pill: The pragmatic functions of ironic criticism. *Discourse Processes*, 23(1), 25–45.
- Colston, H. L., & Keller, S. B. (1998). You'll never believe this: Irony and hyperbole in expressing surprise. *Journal of psycholinguistic research*, 27(4), 499–513.
- Colston, H. L., & O'Brien, J. (2000). Contrast of kind versus contrast of magnitude: The pragmatic accomplishments of irony and hyperbole. *Discourse Processes*, 30(2), 179–199.
- Davidov, D., Tsur, O., & Rappoport, A. (2010). Semi-supervised recognition of sarcastic sentences in twitter and amazon. In *Proceedings of the fourteenth conference on computational natural language learning* (pp. 107–116).
- Filatova, E. (2012). Irony and sarcasm: Corpus generation and analysis using crowdsourcing. In *Lrec* (pp. 392–398).
- Frank, M. C., & Goodman, N. D. (2012). Predicting pragmatic reasoning in language games. *Science*, 336(6084), 998–998.
- Gibbs, R. W. (2000). Irony in talk among friends. *Metaphor and symbol*, 15(1-2), 5–27.
- González-Ibáñez, R., Muresan, S., & Wacholder, N. (2011). Identifying sarcasm in twitter: a closer look. In *Proceedings of the 49th annual meeting of the association for computational linguistics: Human language technologies: short papers-volume 2* (pp. 581–586).
- Goodman, N. D., & Stuhlmüller, A. (2013). Knowledge and implicature: Modeling language understanding as social cognition. *Topics in cognitive science*, 5(1), 173–184.
- Grice, H. P. (1967). Logic and conversation. *The Semantics-Pragmatics Boundary in Philosophy*, 47.
- Jorgensen, J., Miller, G. A., & Sperber, D. (1984). Test of the mention theory of irony. *Journal of Experimental Psychology: General*, 113(1), 112.

- Kao, J. T., Bergen, L., & Goodman, N. D. (2014). Formalizing the pragmatics of metaphor understanding. In *Proceedings of the 36th annual meeting of the cognitive science society*.
- Kao, J. T., Wu, J. Y., Bergen, L., & Goodman, N. D. (2014). Nonliteral understanding of number words. *Proceedings of the National Academy of Sciences*, 111(33), 12002–12007.
- Leggitt, J. S., & Gibbs, R. W. (2000). Emotional reactions to verbal irony. *Discourse processes*, 29(1), 1–24.
- Roberts, R. M., & Kreuz, R. J. (1994). Why do people use figurative language? *Psychological Science*, 5(3), 159–163.
- Sperber, D., & Wilson, D. (1981). Irony and the use-mention distinction. *Radical pragmatics*, 49.
- Wallace, B. C., Do Kook Choe, L. K., & Charniak, E. (2014). Humans require context to infer ironic intent (so computers probably do, too). *ACL*.
- Wilson, D. (2006). The pragmatics of verbal irony: Echo or pretence? *Lingua*, 116(10), 1722–1743.