

# The Funny Thing About Incongruity: A Noisy Channel Model of Puns

**Justine Kao (justinek@stanford.edu)**

Department of Psychology  
Stanford, USA

**Noah Goodman (ngoodman@stanford.edu)**

Department of Psychology  
Stanford, USA

**Roger Levy (rlevy@ucsd.edu)**

Department of Linguistics  
UCSD, USA

## Abstract

“Researchers showed the robot ten puns, hoping that one of them would make it laugh. Unfortunately, no pun in ten did.”

**Keywords:** Humor; language understanding; noisy channel; probabilistic models; sentence processing

## Introduction

Humor plays an essential role in human interactions. In a study on gender differences in desired characteristics of relationship partners, both men and women rated sense of humor as more important than physical attractiveness and earning potential (Stewart, Stinnett, & Rosenfeld, 2000). Humor has important positive effects on children’s development (Frank & McGhee, 1989), success in the work place (Duncan, Smeltzer, & Leap, 1990), marital satisfaction (Ziv & Gadish, 1989), and coping with illness and traumatic events (Johnson, 2002; Gelkopf & Kreitler, 1996). In this paper, we are interested in understanding how this fundamental and ubiquitous phenomenon works from the perspective of cognitive science. What makes something funny? How might the defining characteristics of humor shed light on the ways in which the mind processes and evaluates information?

A leading theory of humor posits that incongruity is a necessary condition for sensations of mirth (Veale, 2004; Forabosco, 1992; McGhee, 1979; Martin, 2007; Hurley, Dennett, & Adams, 2011). The concept of incongruity—perceiving a situation from different viewpoints and finding the subsequent interpretations of the same situation to be incompatible—and its role in humor dates back to Kant’s theories about laughter and the sublime (Veatch, 1998; Forabosco, 1992). Although there is disagreement in the literature about whether incongruity alone is sufficient, most theorists accept that incongruity is necessary for producing sensations of humor. As Veale (2004) states, “Of the few sweeping generalizations one can make about humor that are neither controversial or trivially false, one is surely that humor is a phenomenon that relies on incongruity.”

While the importance of incongruity has long been recognized, its definitions are often ambiguous and difficult to operationalize in empirical research. In this paper, we propose a computational model of humor that formalizes the concept of

incongruity and tests its relationship to humor. Since humor in its entirety relies on commonsense knowledge, reasoning, and discourse understanding that are beyond the scope of this paper, we focus on applying formalizations of incongruity to a subset of linguistic humor: puns.

Writer and philosopher Henri Bergson defined a pun as “a sentence or utterance in which two ideas are expressed, and we are confronted with only one series of words.” This definition highlights the fact that one sentence must evoke two different interpretations in order to count as a pun, which aligns with the concept of incongruity as a requisite of humor. We chose to develop our model on homophone puns—puns that contain homophone words—because the space of possible interpretations of a homophone pun is relatively constrained and well-defined. Since distinct interpretations of a homophone pun hinges on one phonetically ambiguous word, interpretations at the sentence level can be approximated by the two lexical forms of the homophone word. An example helps illustrate this idea:

*“The magician got so mad he pulled his hare out.”*

This sentence allows for two interpretations:

- (a) The magician got so mad he performed the trick of pulling a rabbit out of his hat.
- (b) The magician got so mad he (idiomatically) pulled out the hair on his head.

If the comprehender interprets the word “hare” as itself, he will arrive at interpretation (a); if he interprets the word as its homophone “hair,” he will arrive at interpretation (b). In other words, the sentence-level differences between interpretations (a) and (b) can be approximated by the two interpretations of the observed word “hare.”

Critically, even though the example we gave was a written pun and the reader sees the word “hare” explicitly on the page, the “hair” interpretation is still present and even salient in the context of the sentence. We propose that the reason why two ideas can be communicated through one set of words is because comprehenders maintain uncertainty about the input and consider alternative meanings distinct from what is directly observed. This is inspired by noisy channel models of

sentence processing, which posit that language comprehension is a rational process that incorporates uncertainty about the input at the word level to arrive at sentence-level interpretations that are globally coherent (Levy, 2008; Levy, Bicknell, Slattery, & Rayner, 2009). This process allows the comprehender to maintain multiple word-level interpretations and arrive at more than one interpretations of a sentence that may each be coherent but incongruous with each other. The notion of incongruity thus fits naturally into a noisy channel model of sentence comprehension and can be formalized as such to explain its role in humor.

Our purposes for developing a formal model of linguistic humor are two-fold. First, we wish to formalize the concept of incongruity and test assumptions adopted by leading theories in humor research. Secondly, we aim to show that a noisy channel of language processing allows for flexible context selection and sentence comprehension that gives rise to sophisticated linguistic and social meaning such as humor.

## Model

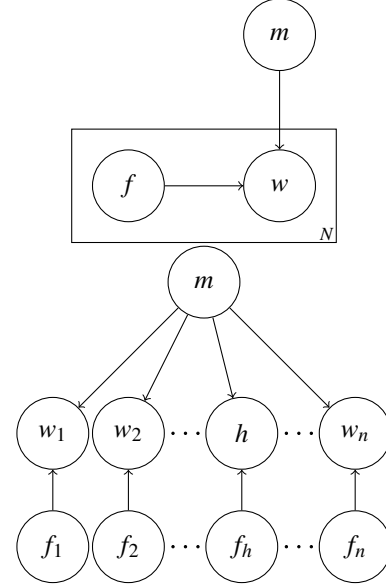
Our model adopts a noisy channel approach in which a comprehender maintains uncertainty about an observed word and considers alternative interpretations. Given two candidate interpretations of the word, the comprehender determines whether one or both interpretations are likely by selectively using contextual information from the sentence.

Previous research suggests that semantic priming plays an important role in lexical disambiguation during sentence processing. In particular, efficient processing of sentences that contain lexical ambiguity involves leveraging semantic association with neighboring words (Seidenberg, Tanenhaus, Leiman, & Bienkowski, 1982; Simpson, 1981; Burke & Yee, 1984; Smith & Levy, 2011). Based on these findings, we propose that comprehenders use semantic information from the surrounding words to compute the probabilities of a candidate interpretation. Motivated by the incongruity theory of humor, we then predict that a sentence in which the target homophone word has two candidate interpretations that are similarly likely and supported by different viewpoints of the sentence is more likely to be funny.

Suppose a sentence is composed of a vector of content words  $\vec{w} = \{w_1, \dots, w_n\}$  and a phonetically ambiguous word  $h$  (remains to be decided if we want to say that  $h$  is part of  $\vec{w}$  or distinct). Since we assume that the candidate interpretations of  $h$  are constrained by phonetic similarity, we make the simplifying assumption that only the two homophone interpretations of  $h$  are salient and probable:  $m_1$  and  $m_2$ . For example, in the magician pun described above,  $h$  is the observed lexical item “hare,” and  $m_1$  and  $m_2$  are the interpretations *hare* and *hair*.

We propose a generative model of a sentence in which the latent variable  $m$  and a latent indicator variable  $\vec{f}$  are responsible for generating the observed words.  $\vec{f}$  is a vector of indicator values that determines whether each word in a sentence is in semantic focus. When  $f_i = 1$ ,  $w_i$  is in semantic focus

and is generated due to semantic relevance to  $m$ . Otherwise if  $f_i = 0$ ,  $w_i$  is simply drawn from a default unigram distribution. Similar approaches have been used in generative models of language to account for words that serve purposes in a sentence beyond providing semantic information, such as topic models that incorporate syntax (Griffiths, Steyvers, Blei, & Tenenbaum, 2005). For our purposes, the indicator variable allows the model to consider subsets of the sentence as viewpoints that differentially support the two candidate interpretations of  $h$ . (Decide which graphical model to use based on role of  $h$ .)



Given our generative model and the observed words  $w$  and  $h$ , we can infer the joint probability distribution of the latent variables  $m$  and  $\vec{f}$  given the observed words. This distribution can be factorized into:

$$P(m, \vec{f} | \vec{w}) = P(m | \vec{w}) P(\vec{f} | m, \vec{w})$$

We derive two measures separately from the factors  $P(m | \vec{w})$  and  $P(\vec{f} | m, \vec{w})$ , which we will call ambiguity and disjointedness. Together, these two measures represent our formalization of incongruity. A sentence should support two interpretations (and is therefore ambiguous) in order to generate incongruous interpretations. In addition, the two interpretations should also be supported by different viewpoints of the sentence, which we approximate as the disjointedness of the subsets of the sentence that are in semantic focus with the two values of  $m$ . We derive these two measures more formally in the following paragraphs.

Let  $M$  denote the distribution  $P(m | \vec{w})$ , a binomial distribution over the two meaning values  $m_1$  and  $m_2$  given the observed words. If the entropy of this distribution  $H(M)$  is low, this means that the probability mass is concentrated on only one meaning, and the alternative meaning is unlikely given the observed words. If  $H(M)$  is high, on the other hand, this means that the probability mass is more evenly distributed among  $m_1$  and  $m_2$ , and the two interpretations are similarly

likely given the contexts.  $H(M)$  is thus a natural measure of the degree of ambiguity present in a sentence.

From conditional probability and Bayes' theorem, we compute  $P(m|\vec{w})$  as follows:

$$\begin{aligned} P(m|\vec{w}) &= \sum_f P(m, f|\vec{w}) \propto P(\vec{w}|m, \vec{f}) P(m) P(\vec{f}) \\ &= P(m) P(\vec{f}) \prod_i P(w_i|m, f_i) \end{aligned}$$

From the generative model,

$$P(w_i|m, f_i) = \begin{cases} P(w_i), & \text{if } f = 0 \\ P(w_i|m), & \text{if } f = 1 \end{cases}$$

We then compute the entropy of  $P(m|w)$  as a measure of ambiguity.

Let  $F_1$  denote the distribution  $P(f|m_1, \vec{w})$  and  $F_2$  denote the distribution  $P(f|m_2, \vec{w})$ .  $F_1$  and  $F_2$  represent the distributions over semantic focus sets given the observed words and  $m_1$  and  $m_2$ , respectively. We use a symmetrised KullbackLeibler divergence score  $D_{KL}(F_1||F_2) + D_{KL}(F_2||F_1)$  to measure the distance between the two distributions. This score approximates the disjointedness of the semantic focus sets given  $m_1$  and  $m_2$ . A low KL score would indicate that  $m_1$  and  $m_2$  are likely to be in semantic focus with similar subsets of the sentence; a high KL score would indicate that  $m_1$  and  $m_2$  are likely to be in semantic focus with different subsets of the sentence.

From conditional probability,

$$P(\vec{f}|m, \vec{w}) \propto P(\vec{f}, \vec{w}|m) = P(\vec{w}|m, \vec{f}) P(\vec{f}|m)$$

From the generative model, since  $\vec{f}$  and  $m$  are independent,  $P(\vec{f}|m) = P(\vec{f})$ . We then compute the symmetrised KL divergence for  $F_1$  and  $F_2$  as a measure of disjointedness.

We have described a generative model of sentences containing phonetically ambiguous words that incorporates the idea of semantic focus sets. We then used the model to derive two measures that formalize incongruity. By combining a noisy channel model of language processing and standard information theoretic measures, our model presents a formalization of incongruity to be tested on humans' judgments of potentially humorous sentences.

## Evaluation

We evaluate our model and measures on a set of 235 sentences, consisting of 65 puns, 130 control non-pun sentences that match the puns in containing the same phonetically ambiguous words, and 40 "de-punned" control sentences that are matched with a subset of the puns with certain manipulations described below. We evaluate our model based on how well it predicts people's funniness ratings of these sentences.

## Materials

We selected 40 pun sentences from a large collection of puns on a website called Pun of the Day, which contains over one

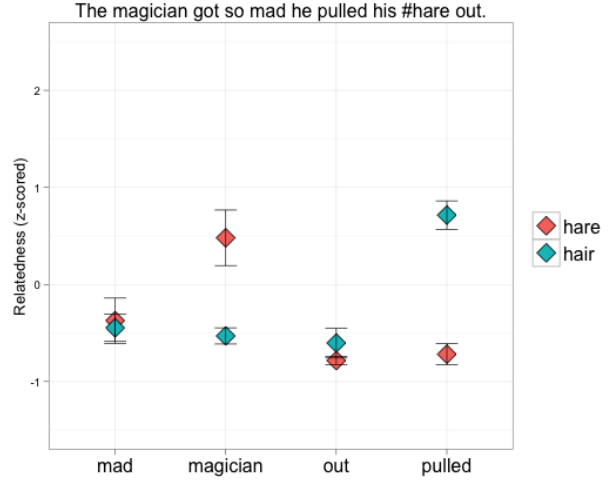


Figure 1: Relatedness of word pairs in example pun.

thousand puns. Punns were selected such that the ambiguous item is a single phonetically ambiguous word, and such that no two puns in the collection have the same ambiguous item. To obtain more homophone pun items, a research assistant generated an addition 25 pun sentences based on a separate list of homophone words.

The 130 non-pun sentences were chosen to match each pun sentence on its ambiguous word as well as the alternative homophone. The sentences were taken from an online version of Heinle's Newbury House Dictionary of American English (<http://nhd.heinle.com/>). We selected sample sentences included in the definition of the homophone word. This design ensured that puns and non-pun sentences contain the same phonetically ambiguous words.

We also constructed 40 sentences to be minimally different from the pun sentences that we collected from "Pun of the Day," which we will call de-punned sentences. A second research assistant who was blind to the hypothesis was asked to replace one word in each of the pun sentences (without changing the homophone word itself) so that the sentence is still grammatical but is no longer a pun. This resulted in sentences that were only one word different from the pun sentences. Below are example sentences from each category.

Type	Example
Pun	"The magician got so mad he pulled his hare out."
De-pun	"The teacher got so mad he pulled his hare out."
Normal	"The hare ran rapidly across the field."
Normal	"Some people have lots of hair on their heads."

## Human ratings of semantic relatedness

As described in the model section, in order to compute our measures, we need the prior probabilities of the interpretations  $P(m)$ , the prior probabilities of the words  $P(w)$ , and the

conditional probabilities of each word in the sentence given an interpretation  $P(w|m)$ . We computed  $P(w)$  and  $P(m)$  directly from the Google unigram corpus. However,  $P(w|m)$  is difficult to obtain through corpora (due to data sparsity) as well as through empirical measures of association strength. Given these constraints, we approximate  $P(w|m)$  using empirical measures of the semantic relatedness between  $w$  and  $m$ , denoted as  $R(c, m)$ . We use  $R(c, m)$  as a proxy for point wise mutual information between  $c$  and  $m$ , defined as follows:

$$R(w, m) = \log \frac{P(w, m)}{P(w)P(m)} = \log \frac{P(w|m)}{P(w)} = \log P(w|m) - \log P(w)$$

With the proper substitutions and transformations,

$$P(w|m) = e^{R(w, m)} P(w)$$

To obtain  $R(w, m)$  for each of the words  $w$  in the stimuli sentences, we recruited 200 subjects on Amazon’s Mechanical Turk to rate distinct word pairs on their semantic relatedness. Function words were removed from the sentences, and the remaining words were paired with each of the interpretations of the homophone sequence (e.g., “magician” and “hare” is a legitimate word pair, as well as “magician” and “hair”). This resulted in 1460 distinct word pairs. Each subject saw 146 pairs of words in random order and were asked to rate how related each word pair is from 1 to 10. The average split-half correlation of the relatedness ratings was 0.916.

Figure 1 and 2 show the relatedness of content words in the sentence with the two homophone interpretations. We see that in the pun sentence, the word “magician” is rated as significantly more related to “hare” than it is to “hair”, while the word “pulled” is rated as significantly more related to “hair” than it is to “hare.” On the other hand, all words in the non-pun example are significantly more related to the word “hare” than to “hair.”

Figure 3 shows the average relatedness ratings of words and the two homophone interpretations across the three types of sentences. In pun sentences, the average relatedness of words to the two homophone interpretations are roughly equivalent. In the non-pun sentences, the average relatedness of words to the observed homophone is significantly higher than to the alternative homophone. This analysis of human ratings of relatedness supports the intuition on which our model is based that funnier sentences are those in which different contexts support incongruous interpretations of the homophone.

### Human Ratings of Funniness

We obtained funniness ratings of the sentences from 100 subjects on Amazon’s Mechanical Turk. Each subject read 58 - 59 sentences in random order, counterbalanced for the sentence types, and rated each sentence on funniness and correctness. The average split-half correlation of the funniness ratings was 0.83. Figure 4 shows the average funniness ratings of puns, non-puns, and de-punned sentences. Pun sentences are rated as significantly funnier than de-punned sentences,

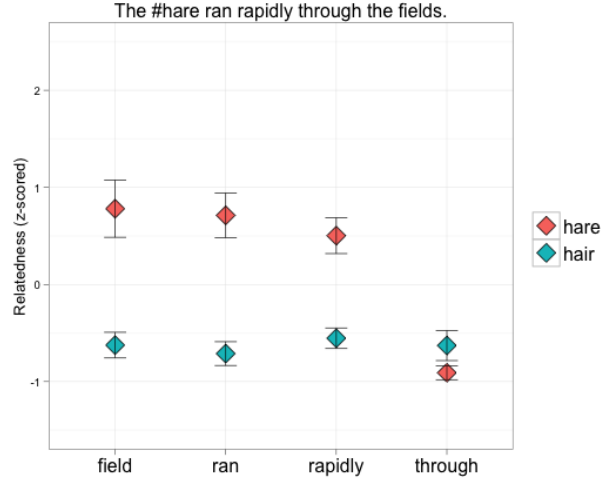


Figure 2: Relatedness of word pairs in example non pun.

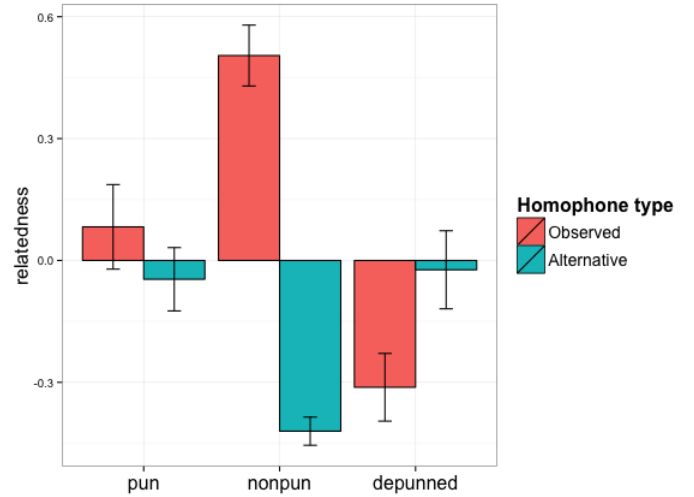


Figure 3: Average relatedness ratings across sentence types

and de-punned sentences are rated as significantly funnier than non-pun sentences.

## Results

We computed an ambiguity and disjointedness value for each of the 235 sentences following the derivations described in the model section and using the relatedness measures described above.

As predicted, ambiguity differs significantly across sentence types ( $F(2, 232) = 11.94, p < 0.0001$ ) and correlates significantly with human ratings of funniness across the 235 sentences ( $r = 0.27, p < 0.0001$ ). Furthermore, disjointedness differs significantly across sentence types as well ( $F(2, 232) = 5.76, p < 0.005$ ) and correlates significantly with human ratings of funniness, although less strongly ( $r = 0.21, p < 0.005$ ).

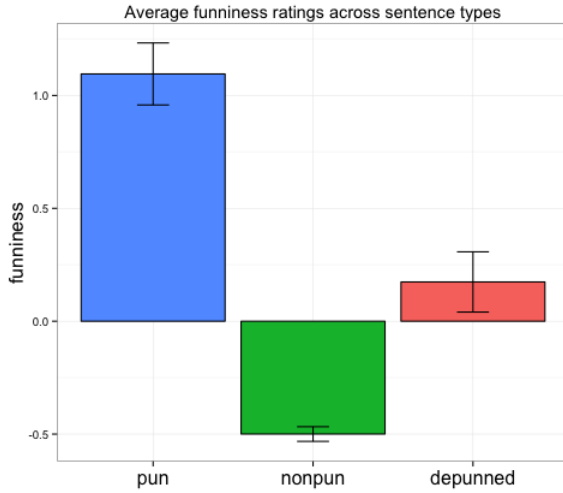


Figure 4: Average funniness ratings across sentence types

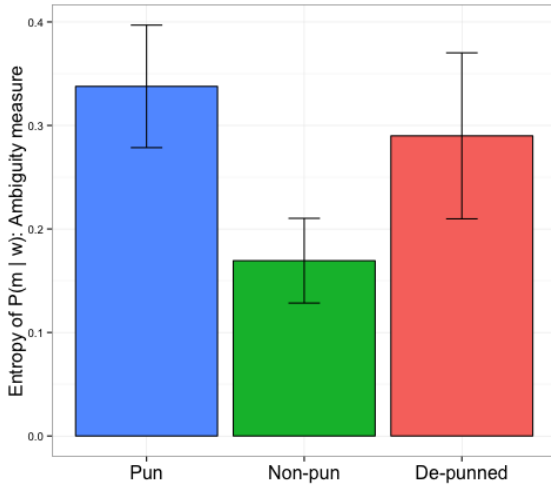


Figure 5: Ambiguity across sentence types

Using both ambiguity and disjointedness as the two dimensions that formalize incongruity, we can distinguish among puns, non-puns, and de-punned sentences. Both ambiguity and disjointedness measures are significant predictors of human ratings of funniness ( $F(2, 232) = 19.6, R^2 = 0.137, p < 0.001$ ).

	Estimate	Std. Error	p value
Intercept	-0.687	0.140	$1.75e^{-06}$ ***
Ambiguity	1.011	0.193	$3.49e^{-07}$ ***
Disjointedness	0.237	0.054	$1.81e^{-05}$ ***

**Incomplete** Show that our model is able to select contexts that give maximal incongruity and maximal justification, which can tell us *why* a pun is funny in addition to when

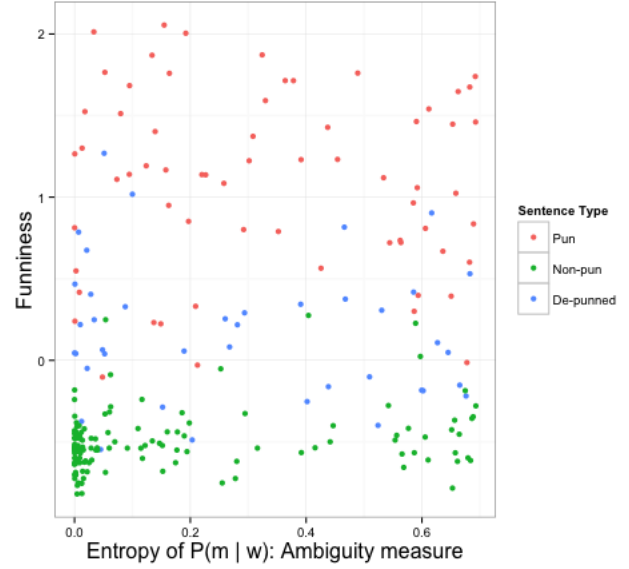


Figure 6: Ambiguity correlated with funniness

it is funny.

## Discussion

### Mostly incomplete

Given the prevalence of humor in human communication, researchers in artificial intelligence have argued that computers should be able to generate and detect humor in order to interact with humans more naturally and effectively (Mihalcea & Strapparava, 2006). As a result, computational humor has made important progress and also received attention from popular press in the last decade (insert NYT citation here). However, most of the work in computational humor has focused either on utilizing joke-specific templates and schemata (Binsted, 1996; Kiddon & Brun, 2011), or identifying linguistic features such as slang and alliteration that strongly predict humorous intent (Mihalcea & Strapparava, 2006; Reyes, Buscaldi, & Rosso, 2010). The former type of studies is restricted to identifying jokes with a very specific format and structure, while the latter type falls short of testing or building upon deeper and more general theories of humor involving the management of incongruity.

Our work moves beyond these two types of approaches and directly utilizes incongruence to identify humorous texts. Given that humor theorists view incongruity as an essential component of jokes, we examined whether human judgments of funniness can be predicted by the presence of incongruous interpretations of the same input. In particular, we developed a formal model of linguistic humor that fits naturally into the framework of normal language processing. We propose that a noisy channel approach to language processing allows the comprehender to consider alternative viewpoints and interpretations of the same linguistic input and can account for the possibility of incongruity.

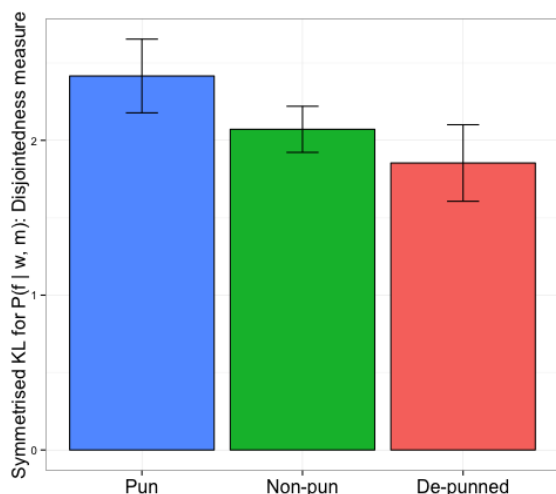


Figure 7: Disjointedness across sentence types

**Incomplete** Describe its advantages over previous work on computational humor. Describe potential applications and generalizations to other forms of humor.

## Acknowledgments

## References

- Binsted, K. (1996). Machine humour: An implemented model of puns.
- Burke, D., & Yee, P. (1984). Semantic priming during sentence processing by young and older adults. *Developmental Psychology*, 20(5), 903.
- Duncan, W., Smeltzer, L., & Leap, T. (1990). Humor and work: Applications of joking behavior to management. *Journal of Management*, 16(2), 255–278.
- Forabosco, G. (1992). Cognitive aspects of the humor process: The concept of incongruity. *Humor: International Journal of Humor Research*; *Humor: International Journal of Humor Research*.
- Frank, M., & McGhee, P. (1989). *Humor and children's development: A guide to practical applications*. Routledge.
- Gelkopf, M., & Kreidler, S. (1996). Is humor only fun, an alternative cure or magic? the cognitive therapeutic potential of humor. *Journal of Cognitive Psychotherapy*, 10(4), 235–254.
- Griffiths, T., Steyvers, M., Blei, D., & Tenenbaum, J. (2005). Integrating topics and syntax. *Advances in neural information processing systems*, 17, 537–544.
- Hurley, M., Dennett, D., & Adams, R. (2011). *Inside jokes: Using humor to reverse-engineer the mind*. Mit Pr.
- Johnson, P. (2002). The use of humor and its influences on spirituality and coping in breast cancer survivors. In *Oncology nursing forum* (Vol. 29, pp. 691–695).
- Kiddon, C., & Brun, Y. (2011). That's what she said: double entendre identification. In *Proceedings of the 49th annual*

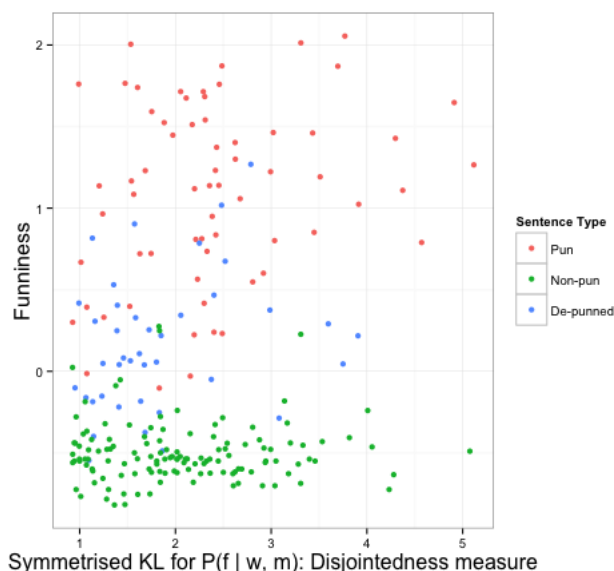


Figure 8: Disjointedness correlated with funniness

*meeting of the association for computational linguistics: Human language technologies* (pp. 89–94).

- Levy, R. (2008). A noisy-channel model of rational human sentence comprehension under uncertain input. In *Proceedings of the conference on empirical methods in natural language processing* (pp. 234–243).
- Levy, R., Bicknell, K., Slattery, T., & Rayner, K. (2009). Eye movement evidence that readers maintain and act on uncertainty about past linguistic input. *Proceedings of the National Academy of Sciences*, 106(50), 21086–21090.
- Martin, R. (2007). *The psychology of humor: An integrative approach*. Academic Press.
- McGhee, P. (1979). *Humor, its origin and development*. WH Freeman San Francisco.
- Mihalcea, R., & Strapparava, C. (2006). Learning to laugh (automatically): Computational models for humor recognition. *Computational Intelligence*, 22(2), 126–142.
- Reyes, A., Buscaldi, D., & Rosso, P. (2010). The impact of semantic and morphosyntactic ambiguity on automatic humour recognition. In H. Horacek, E. Mtais, R. Muoz, & M. Wolska (Eds.), *Natural language processing and information systems* (Vol. 5723, p. 130–141). Springer Berlin / Heidelberg.
- Seidenberg, M., Tanenhaus, M., Leiman, J., & Bienkowski, M. (1982). Automatic access of the meanings of ambiguous words in context: Some limitations of knowledge-based processing. *Cognitive psychology*, 14(4), 489–537.
- Simpson, G. (1981). Meaning dominance and semantic context in the processing of lexical ambiguity. *Journal of Verbal Learning and Verbal Behavior*, 20(1), 120–136.
- Smith, N., & Levy, R. (2011). Cloze but no cigar: The complex relationship between cloze, corpus, and subjective probabilities in language processing. In *Proceedings of the*

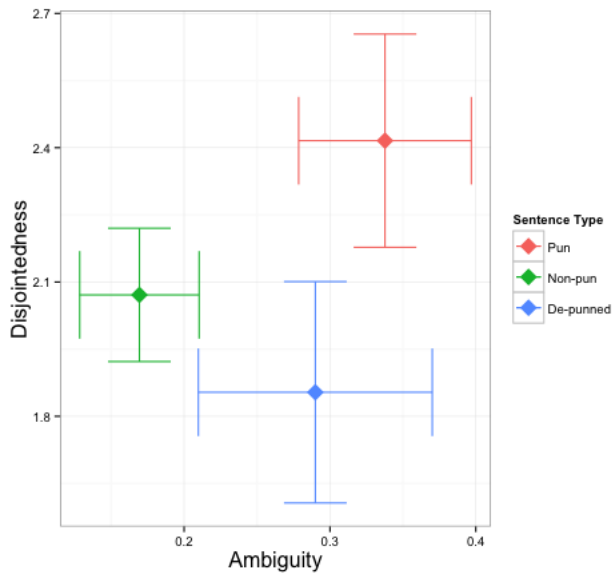


Figure 9: Ambiguity and disjointedness across sentence types

- 33rd annual meeting of the cognitive science conference.
- Stewart, S., Stinnett, H., & Rosenfeld, L. (2000). Sex differences in desired characteristics of short-term and long-term relationship partners. *Journal of Social and Personal Relationships*, 17(6), 843–853.
- Veale, T. (2004). Incongruity in humor: Root cause or epiphenomenon? *Humor-International Journal of Humor Research*, 17(4), 419–428.
- Veatch, T. (1998). A theory of humor. *Humor*, 11, 161–215.
- Ziv, A., & Gadish, O. (1989). Humor and marital satisfaction. *The journal of social psychology*, 129(6), 759–768.