

A computational model of linguistic humor

Justine T. Kao^{*}, Roger Levy[†] and Noah D. Goodman^{*}

^{*}Stanford University, and [†]University of California, San Diego

Submitted to Proceedings of the National Academy of Sciences of the United States of America

Humor plays an essential role in human interactions. However, its precise nature remains elusive. While research in natural language understanding has made significant advancements in recent years, there has been little direct integration of humor research with computational models of general language understanding. In this paper, we propose two information-theoretic measures of humor—~~ambiguity and distinctiveness~~—~~inspired by~~ humor theories and directly derived from a simple model of sentence processing. We then test these measures on a set of puns and non-pun sentences and show that they correlate significantly with human judgments of funniness. Our model is one of the first to integrate general linguistic knowledge, speaker intent, and humor theory to model humor computationally. We present it as a framework for applying models of language processing to understand higher-level linguistic and cognitive phenomena.

Linguistic humor | Language understanding | Computational modeling

Abbreviations: IR, Incongruity Resolution

Introduction

Love may make the world go round, but humor is the glue that keeps it together. Previous research has shown that humor is ubiquitous across cultures [?, ?], increases interpersonal attraction [?], helps resolve intergroup conflicts [?], and improves psychological wellbeing [?]. Our everyday experiences show that humor plays a critical role in human interactions and composes a significant part of our linguistic, cognitive, and social lives. However, the cognitive basis of humor remains relatively unknown. In this paper, we build upon theories of humor and language processing to computationally model the humor in a set of **phonetic puns**. By providing a formal model of linguistic humor, we aim to shed light on how the mind identifies and processes the inputs that make us laugh.

Of the many theories of humor proposed since Aristotle, **the incongruity theory** comes closest to **a cognitive account**. Incongruity is defined as the incompatible and often schema-violating interpretations of a stimulus. While most humor researchers agree that the experience of incongruity is necessary for generating humor, some argue that it alone is insufficient. These scholars propose a two-stage model of humor termed the Incongruity-Resolution model, **which involves the discovery of a cognitive rule that explains and resolves the incongruity in the stimulus**. However, definitions of both incongruity and resolution are ambiguous and leave much room for disparate interpretations among humor researchers. While informal models of humor have provided important insights into the necessary and sufficient conditions of humor, the lack of computational rigor makes it difficult to operationalize these conditions or to empirically evaluate their contribution to the perception of humor.

More recently, researchers in artificial intelligence and computational linguistics have applied computational tools to examine features of humor. The interest in computational humor ~~was~~ ^{has been} motivated in part by the need for computers to recognize and generate humorous input in order to interact with humans in a more engaging and natural manner [?]. However, most work on computational humor focuses either on

joke-specific templates and schemata [?, ?] or surface features such as innuendo and slang [?, ?]. While these approaches are able to identify and produce humorous stimuli within certain constraints, they fall short of testing and building upon a more general cognitive theory of humor.

In this paper, we directly utilize a computational model of sentence processing to derive theory-driven measures of humor. By basing our measures of humor on existing ideas such as incongruity and resolution, we are able to leverage and test the insights generated by decades of qualitative research. By formalizing these measures, we can quantitatively evaluate how different factors contribute to the experience of humor. Furthermore, by deriving these measures from a model of general sentence processing, **we are able to view linguistic humor as a direct result of language understanding strategies instead of as a separate dedicated process.**

While we aim to develop a model that encompasses a broad range of humorous stimuli, many types of humor rely on rich commonsense knowledge and discourse understanding. To reasonably limit the scope of our task, we focus on testing our model on a subset of linguistic humor for which we are able to obtain reasonable formal representations of meaning. In particular, we focus on phonetic puns, which are puns containing words that sound identical or similar to other words in the English language. Unlike other types of jokes, the space of possible meaning representations of a phonetic pun is relatively constrained and well defined. For example, consider the following sentence:

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

Clarify: do we mean that the phonetic form of this sentence allows for two different interpretations?

This sentence allows for two different interpretations:

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

If the comprehender interprets the word “hare” as *hare*, he will arrive at interpretation (a). If he interprets the word as the homophone *hair*, however, he will arrive at interpretation (b). In general, interpretations of a phonetic pun hinge on a single phonetically ambiguous word, which allows us to use different homophones of the ambiguous word to approximate different interpretations of the entire sentence. **Focusing on phonetic puns thus enables us to formalize measures of humor without first having to represent the meaning of complex sentences and discourse.**

Reserved for Publication Footnotes

I think it would be super helpful to include an example of IR here.

Do we want to say that these two meanings are incongruous?

We observe that although the reader sees the word “hare” explicitly when processing the example given, the “hair” interpretation is still highly accessible. Previous research on sentence comprehension has suggested that people maintain uncertainty about the surface input when processing a sentence. In other words, comprehenders assume that communication happens through a “noisy channel” and that some parts of the input they receive may have been corrupted. In order to successfully infer the true meaning of a sentence, comprehenders consider multiple word-level interpretations during processing and rationally incorporate them to arrive at coherent interpretations at the sentence level. By positing noise in the input and modeling comprehension as rational inference under uncertainty, the noisy-channel model is able to explain a variety of phenomena in language processing.

Here we propose that the humor in phonetic puns may also arise from the assumption of a noisy channel. Because the comprehender maintains uncertainty about the input, it is possible for her to arrive at multiple interpretations of a sentence that are each coherent but incongruous with each other. Previous research has shown that both semantic priming and the sequential structure of language play important roles in sentence processing. We propose a model of sentence comprehension that incorporates the noisy-channel assumption, the semantic relationship between a sentence’s overall meaning and the words that compose it, and the sequential structure of language to infer the meaning of a sentence.

Model. Suppose a sentence is composed of a vector of words $\vec{w} = \{w_1, \dots, w_i, h, w_{i+1}, \dots, w_n\}$, where h is a phonetically ambiguous word. We construct a simple generative model for \vec{w} (see Figure). Given a latent sentence meaning m , each word is generated independently by first deciding if it explicitly reflects the sentence meaning or is corrupted by noise. This process is indicated by an indicator variable vector \vec{f} , where w_i is a meaningful word if $f_i = 1$ and a noise word if $f_i = 0$. When w_i is a meaningful word, it is sampled in proportion to its semantic relevance to m . When it is a corrupted word, we assume that its corruption is affected by the immediately preceding words and sample it from an n-gram language model.

To select the sentence meanings m , we exploit the convenient property of phonetic puns described before. We introduce the simplifying assumption that the sentence meanings m correspond to plausible interpretations of the homophone word h , which are constrained by phonetic similarity. For example, “hare” is a phonetically ambiguous word, and the homophones *hare* and *hair* can approximate the two possi-

ble sentence meanings m_1 and m_2 . While this approximation is admittedly coarse, it makes use of the reasonable assumption that sentences containing the word “hare” will generally be about the topic *hare*, and sentences that have the word “hair” will generally be about the topic *hair*. This assumption reduces the ill-defined space of sentence meanings to the simple proxy of alternate spellings for phonetically ambiguous words.

Using the generative model described, we can infer the joint probability distribution $P(m, \vec{f}|\vec{w})$ of the sentence topic m and the indicator variables \vec{f} given a set of words \vec{w} . This distribution can be factorized as the following:

$$P(m, \vec{f}|\vec{w}) = P(m|\vec{w})P(\vec{f}|m, \vec{w}) \quad [1]$$

Components?

We derive two formal **measure** of humor from the terms on the righthand side, which we call ambiguity and distinctiveness. Ambiguity captures the degree to which sentence meanings are similarly likely. This can be quantified as a summary of the binomial distribution $P(m|\vec{w})$. If the entropy of this distribution is low, then one meaning is much more likely than the other. If the entropy is high, then both meanings are similarly likely.

Distinctiveness measures the degree to which two sentence meanings are supported by distinct parts of the sentence. This can be quantified as a summary of the distribution $P(\vec{f}|m, \vec{w})$. Given \vec{w} and the topic m_1 , which is directly observed in the sentence, we compute the distribution $F_1 = P(\vec{f}|m_1, \vec{w})$. Given \vec{w} and the topic m_2 , we compute the distribution $F_2 = P(\vec{f}|m_2, \vec{w})$. We wish to measure how different the distribution over topic words would be given different sentence meanings. We use a Kullback-Leibler divergence score $D_{KL}(F_1||F_2)$ to measure the distance between F_1 and F_2 . A low KL score indicates that the possible sentence meanings are supported by similar subsets of the sentence. A high KL score indicates the sentence meanings are each strongly supported by distinct subsets of the sentence. Together, our ambiguity and distinctiveness measures constitute a two-dimensional formalization of humor.

While only loosely based on existing definitions of “incongruity” and “resolution,” the two formal measures we derived have natural connections to the incongruity-resolution model of humor. Under the assumption that the two sentence meanings are sufficiently different from each other, a high ambiguity score suggests the coexistence of two incompatible interpretations, which is a characterization of “incongruity.” A high distinctiveness score suggests that given one sentence meaning, the comprehender should regard one set of words as meaningful words, while given another sentence meaning, attention should be directed to a different set of words. This allows the comprehender to pay attention to different sets of words given different candidate sentence meanings, effectively “**resolving**” the incongruity by discovering the cognitive rule that partitions the sentence into distinct parts that are each internally harmonious.

Ok, this is really helpful to me. I would still benefit from more discussion of why distinctiveness actually lines up well with what humor theorists call “resolution”

Results

We evaluated the ambiguity and distinctiveness measures on a set of 435 phonetically ambiguous sentences. Of these sentences, 65 are identical homophone puns and 80 are puns where the two candidate meanings sound similar but are not identical to each other (near homophone puns). The remaining 290 sentences are non-pun control sentences that contain the same phonetically ambiguous words as the puns. Table 1 shows an example of each type of sentence.

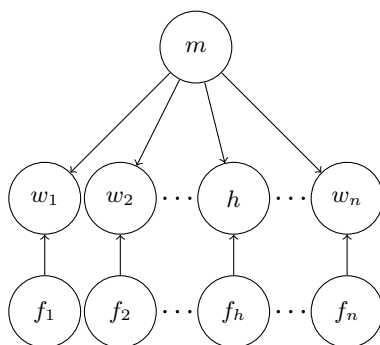


Fig. 1: Generative model of a sentence. Each word w_i is generated based on the sentence topic m if the indicator variable f_i puts it in semantic focus; otherwise it is generated as noise (from a unigram distribution).

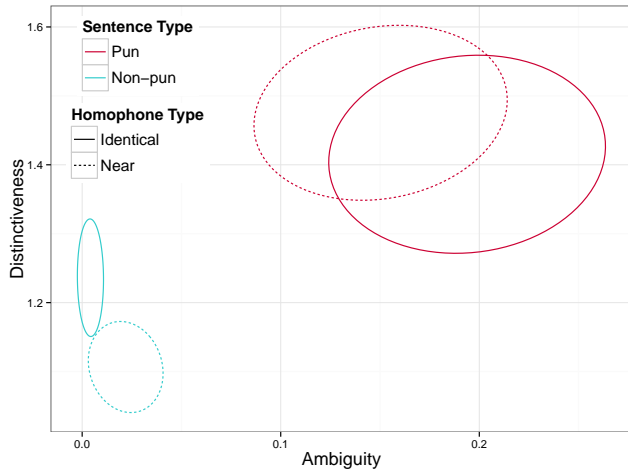


Fig. 3: Standard error ellipses of ambiguity and distinctiveness across sentence types. Puns score higher on ambiguity and distinctiveness; non-puns have low ambiguity and distinctiveness.

Table 1: Example sentences

Type	Hom	Example
Pun	Identical	The magician was so mad he pulled <u>is</u> hare out.
Pun	Near	A dentist has to tell a patient the whole tooth.
Non-pun	Identical	The hare ran rapidly across the field.
Non-pun	Identical	Some people have lots of hair on their heads.
Non-pun	Near	A dentist examines one tooth at a time.
Non-pun	Near	She always speaks the truth.

We obtained human ratings of funniness for each of the 435 sentences (see Materials and Methods). Applying the derivations and relatedness measures described in the Methods section, we computed an ambiguity and distinctiveness score for each sentence. As predicted, ambiguity was significantly higher for pun sentences than non-pun sentences ($F(1, 433) = 108.4, p < 0.0001$). This suggests that our ambiguity measure successfully captures characteristics that distinguish puns from other phonetically ambiguous sentences. Intuitively, while non-pun sentences also contain phonetically ambiguous words, their interpretations are less ambiguous because the words are highly semantically related to only one sentence meaning. Our measure of distinctiveness differs marginally significantly across sentence types ($F(1, 433) = 47.1, p < 0.1$). Using both ambiguity and distinctiveness as dimensions that formalize humor, we can distinguish among pun and non-pun sentences, as shown in Figure . Figure shows the standard error ellipses for the two sentence types in the two-dimensional space of ambiguity and distinctiveness. Although there is a fair amount of noise in the predictors (likely due to simplifying assumptions, the need to use empirical measures of relatedness, and the inherent complexity of humor), we see that pun sentences tend to cluster at a space with higher ambiguity and distinctiveness, while non-pun sentences score significantly lower on both measures. A linear regression model showed that both ambiguity and distinctiveness are significant predictors of funniness ratings across all 435 sentences. Together, the two predictors capture a modest but significant amount of the reliable variance in funniness ratings ($F(2, 432) = 76.79, r = 0.51, p < 0.0001$; see Table 2).

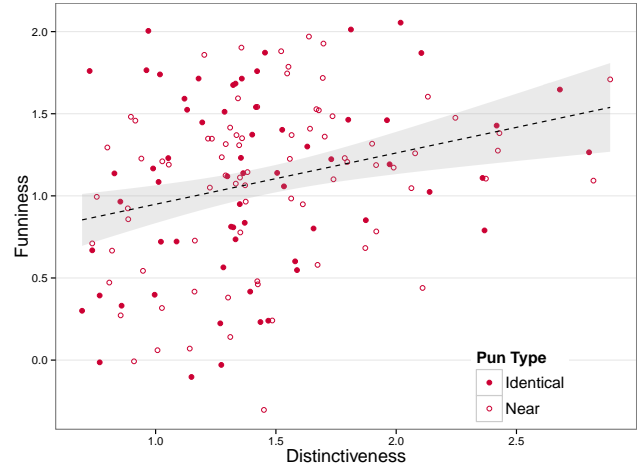


Fig. 2: Figure caption

Table 2: Regression coefficients using ambiguity and distinctiveness to predict funniness ratings

	Estimate	Std. Error	p value
Intercept	-0.830	0.1070	< 0.0001
Ambiguity	1.899	0.212	< 0.0001
Distinctiveness	0.568	0.082	< 0.0001

We now examine the measures' quantitative predictions of funniness within puns. Ambiguity does not correlate with human ratings of funniness within the 145 pun sentences ($r = 0.03, p > 0.05$), suggesting that it alone is unable to distinguish highly funny puns from mediocre puns. On the other hand, distinctiveness ratings correlate significantly with human ratings of funniness within pun sentences ($r = 0.28, p < 0.001$). This suggests that while ambiguity distinguishes puns from non-puns, distinctiveness is needed to separate very funny puns from mediocre ones.

Besides predicting the funniness of a sentence, our model can also tell us which particular features of a pun make it amusing. By finding the most likely sets of meaningful words given each latent sentence meaning m and \vec{w} , we can identify words in a pun that are critical to producing its humor. Table shows the most likely sets of meaningful words given each meaning for two groups of sentences. Sentences in each group contain the same pair of candidate meanings for the target word h . However, they differ in measures of ambiguity, distinctiveness, and funniness. Words in the most likely sets given m_1 are in red; words in the most likely sets given m_2 are in green; and words in the most likely sets of both meanings are in dark blue. We observe that visually, the two pun sentences (which are significantly funnier) have more distinctive and balanced sets of meaningful words for each sentence meaning than other sentences in their groups. Non-pun sentences tend to have no words in support of the meaning that was not observed. Furthermore, imagine if you were asked to explain why the two pun sentences are funny. The colorful words in each pun sentence—for example, the fact that magicians tend to perform magic tricks with hares, and people tend to be described as pulling out their hair when angry—are what one might intuitively use to explain why the sentence is a pun. Our model thus provides a natural way of not only us-

ing ambiguity and distinctiveness to predict when a sentence is a pun, but also to explain what aspects of a pun make it funny.

Discussion

In this paper we presented a simple model of sentence processing and derived formal measures to predict human judgments of humor for a set of sentences. We showed that a noisy-channel model allows for flexible context selection that may give rise to sophisticated linguistic meaning such as humor during normal sentence processing. In particular, the comprehender’s consideration of whether a word is semantically meaningful allows her to switch between sentence meanings and entertain the possibility that the speaker intends both meanings at once. As writer and philosopher Henri Bergson wrote, “A pun is a sentence or utterance in which two ideas are expressed, and we are confronted with only one series of words.” The noisy-channel model allows us to formally capture this duality of meaning and to show how it may be manifested by focusing on different sets of words.

With regards to humor theory, our analysis suggests that the formalizations of ambiguity and distinctiveness we proposed may account for separate aspects of humor appreciation. Ambiguity serves to distinguish humorous input from normal sentences, while distinctiveness predicts the fine-grained degree of funniness within humorous input. We speculate that the degree of funniness captured by distinctiveness may be related to the recognition of a cognitive rule that meaningfully separates one group of words from another. While the connections between our measures and the incongruity-resolution model of humor are tenuous, these results seem to support the idea that incongruity is a necessary condition for humor, while the intensity of humor experienced may depend upon incongruity resolution.

Although our task in this paper is limited in scope, it is a step towards developing models that can explain rich linguistic phenomena such as humor. Future work may incorporate more sophisticated models of language understanding—ones that consider deeper semantic representations and multi-sentence discourse—to derive corresponding measures for different types of jokes. Besides seeking to understand linguistic humor for its own sake, these higher-order phenomena can serve as probes for developing models of language processing that account for a wider range of linguistic behavior, and may even shed light on how language understanding works more generally. We believe that our work contributes to research in humor theory, computational humor, and language understanding, such that some day we can build robots that make us laugh and understand the appreciation for humor that makes us uniquely human.

Materials and Methods

Model details. For simplification, our model disregards function words in the sentences and treats them as only indirectly contributing to the overall meaning of a sentence. Each of the w_i in \vec{w} in this paper is a content word. Here we will describe the derivations for ambiguity and distinctiveness in more detail.

Bernoulli

Ambiguity: $P(m|\vec{w})$ is a *binomial* distribution over the two meaning values m_1 and m_2 given the observed words. If the entropy of this distribution is low, this means that the probability mass is concentrated on only one meaning, and the other meaning is unlikely given the observed words. If entropy is high, 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 sentence. The entropy of $P(m|\vec{w})$ is thus a natural measure of the degree of ambiguity present in a sentence. We compute

$P(m|\vec{w})$ as follows:

$$P(m|\vec{w}) = \sum_{\vec{f}} P(m, \vec{f}|\vec{w}) \quad [2]$$

$$\propto \sum_{\vec{f}} P(\vec{w}|m, \vec{f}) P(m) P(\vec{f}) \quad [3]$$

$$= \sum_{\vec{f}} \left(P(m) P(\vec{f}) \prod_i P(w_i|m, f_i) \right) \quad [4]$$

We approximate $P(m)$ as the unigram frequency of the words that represent m . For example, $P(m = \text{hare})$ is approximated as $P(m = \text{“hare”})$. We also assume a uniform prior probability over all subsets of the words being semantically meaningful. In other words, $P(f_i = 1) = 0.5$, and so $P(\vec{f})$ is a constant. As for $P(w_i|m, f_i)$, the probability depends upon the value of the indicator variable f_i . If $f_i = 1$, w_i is semantically meaningful and was sampled in proportion to its relatedness with the sentence meaning m . If $f_i = 0$, then w_i was generated from a noise process and sampled in proportion to its probability given the previous two words (including function words). In other words, semantically meaningful words are generated based on the meaning, while noise words are generated based on the sequential structure of language. From the generative model,

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

where $P(w_i|m)$ is estimated using measures described in Experiment 2. $P(w_i|\text{bigram}_i)$ was estimated from the Google Ngrams corpus, where bigram_i is the bigram that immediately precedes w_i . Once we derive $M = P(m|\vec{w})$, we compute its information-theoretic entropy as a measure of ambiguity:

$$\text{Amb}(M) = - \sum_i P(m_i|\vec{w}) \log P(m_i|\vec{w})$$

Distinctiveness: We next turn to the distribution over indicator variables \vec{f} given a sentence meaning. Given Bayes’ Rule, $P(\vec{f}|m, \vec{w})$ is computed as follows:

$$P(\vec{f}|m, \vec{w}) \propto P(\vec{w}|m, \vec{f}) P(\vec{f}|m) \quad [5]$$

Since \vec{f} and m are independent, $P(\vec{f}|m) = P(\vec{f})$. Let F_1 denote the distribution $P(\vec{f}|m_1, \vec{w})$ and F_2 denote the distribution $P(\vec{f}|m_2, \vec{w})$. F_1 and F_2 represent the distributions over semantically meaningful sets assuming the sentence topic m_1 and m_2 , respectively. We use the Kullback-Leibler divergence score $D_{KL}(F_1||F_2)$ to measure the distance between F_1 and F_2 . This score measures how “distinct” the semantically meaningful sets are given m_1 and m_2 . A low KL score would indicate that meanings m_1 and m_2 are supported by similar subsets of the sentence; a high KL score would indicate that m_1 and m_2 are supported by distinct subsets of the sentence. Based on the definition of KL,

$$\text{Dist}(F_1||F_2) = \sum_i \ln \left(\frac{F_1(i)}{F_2(i)} \right) F_1(i).$$

Experiment 1. We collected pun sentences from a website called “Pun of the Day” (<http://www.punoftheday.com/>), which contains over a thousand puns submitted by users. Puns were selected such that the ambiguous item in each pun is a single phonetically ambiguous word. We obtained 40 puns where the phonetically ambiguous word has identical homophones (identical-homophone puns). To obtain more identical homophone pun items, a research assistant generated an additional 25 pun sentences based on a separate list of homophone words. We then selected 80 puns where the phonetically ambiguous word has a similar-sounding word that is not an identical homophone (near-homophone puns). This resulted in a total of 145 pun sentences. We selected 290 corresponding non-pun sentences from an online version of Heinle’s Newbury House Dictionary of American English (<http://www.newburyhouse.com/>). We chose sample sentences included in the definition of the homophone word. 145 of the sentences contain the ambiguous words from the pun sentences, and 145 of them contain the alternative homophones. This design ensured that puns and non-pun sentences contain the same set of phonetically ambiguous words. Table 1 shows example sentences from each category.

We obtained funniness ratings for each of the sentences. The 195 identical homophone sentences were rated by 93 subjects on Amazon’s Mechanical Turk. Each subject read roughly 60 sentences in random order, counterbalanced for the sentence types, and rated each sentence on funniness and correctness. The average split-half correlation of funniness ratings was 0.83. Pun sentences were rated as significantly funnier than non-pun sentences ($F(2, 232) = 415.3, p < 0.0001$). The 240

Can we state the precise criterion used to determine when to color a word?



Table 3: Sets of semantically meaningful words, ambiguity/distinctiveness scores, and funniness ratings for two groups of sentences. Words in red are semantically meaningful given m_1 ; green given m_2 ; blue given both.

m_1	m_2	Type	Sentence and Semantic Focus Sets	Amb.	Dist.	Funniness
		Pun	The magician got so mad he pulled his hare out.	0.15	1.36	1.714
hare	hair	Non-pun	The hare ran rapidly through the fields .	$1.43E^{-5}$	1.07	-0.400
		Non-pun	Most people have lots of hair on their heads .	$9.47E^{-11}$	1.55	-0.343
		Pun	A dentist has to tell the patient the whole tooth .	0.10	1.64	1.41
tooth	truth	Non-pun	A dentist examines one tooth at a time.	$8.92E^{-5}$	1.23	-0.45
		Non-pun	She always speaks the truth .	$3.85E^{-10}$	0.72	-0.46

near homophone sentences were rated by 158 subjects on Mechanical Turk. Each subject read 40 sentences in random order, counterbalanced for the sentence types, and rated each sentence on funniness. The average split-half correlation of funniness ratings was X . We z-scored the ratings and used the average z-scored ratings across subjects as human judgments of funniness.

Experiment 2. As described in the model details, computing ambiguity and distinctiveness requires the prior probabilities of meanings $P(m)$ (approximated as the unigram probabilities of the words that denote the meanings), the probabilities of words given the preceding bigram $P(w|bigram)$, and the conditional probabilities of each word in the sentence given a meaning $P(w|m)$. While we computed $P(m)$ and $P(w|bigram)$ directly from the Google Web unigram corpus, $P(w|m)$ is difficult to obtain through traditional topic models trained on corpora due to data sparsity. We thus measure it empirically with an experiment.

We approximate $P(w|m)$ using an empirical measure of the semantic relatedness between w and m , denoted $R(w, m)$. We use $R(w, m)$ as a proxy for point wise mutual information between w and m , defined as follows:

$$R(w, m) = \log \frac{P(w, m)}{P(w)P(m)} = \log P(w|m) - \log P(w) \quad [6]$$

We assume that human ratings of relatedness between two words $R'(w, m)$ approximate true relatedness up to an additive constant z and assume $z = 0$ for our

purposes. With the proper substitutions and transformations,

$$P(w|m) = e^{R'(w, m) + z} P(w) \quad [7]$$

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 word pairs on their semantic relatedness. Since it is difficult to obtain the relatedness rating of a word with itself, we used a free parameter r and fit it to data ($r = 13$). Function words were removed from each of the sentences in our dataset, and the remaining words were paired with each of the interpretations of the homophone sequence (e.g., for the pun in Table ??, "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 using a scale from 1 to 10. The average split-half correlation of the relatedness ratings was 0.916, indicating that semantic relatedness was a reliable measure. We used the average z-scored relatedness measure for each word pair to obtain $R(w, m)'$ and Google Web unigrams to obtain $P(w)$, thus computing $P(w|m)$ for all word and meaning pairs. This completed the values needed to compute ambiguity and distinctiveness for all sentences.

ACKNOWLEDGMENTS. This work was partially supported by a grant from the Spanish Ministry of Science and Technology.

Ah, is this different than what we did before?

1.
2.
3.
4.

5.
6.