

# Quantifying the role of discourse topicality in speakers' choices of referring expressions

**Naho Orita**

Department of Linguistics  
University of Maryland  
nahoo@umd.edu

**Eliana Vornov**

Departments of Computer Science and Linguistics  
University of Maryland  
evornov@umd.edu

**Naomi H. Feldman**

Department of Linguistics  
University of Maryland  
nhf@umd.edu

**Jordan Boyd-Graber**

College of Information Studies and UMIACS  
University of Maryland  
jbg@umiacs.umd.edu

## Abstract

The salience of an entity in the discourse is correlated with the type of referring expression that speakers use to refer to that entity. Speakers tend to use pronouns to refer to salient entities, whereas they use lexical noun phrases to refer to less salient entities. We propose a novel approach to formalize the interaction between salience and choices of referring expressions using topic modeling, focusing specifically on the notion of topicality. We show that topic models can capture the observation that topical referents are more likely to be pronominalized. This lends support to theories of discourse salience that appeal to latent topic representations and suggests that topic models can capture aspects of speakers' cognitive representations of entities in the discourse.

## 1 Introduction

Speakers' choices of referring expressions (pronouns, demonstratives, full names, and so on) have been used as a tool to understand cognitive representations of entities in a discourse. Many researchers have proposed a correlation between the type of a referring form and *salience* (or *accessibility*, *prominence*, *focus*) of the entity in the discourse (Chafe, 1976; Gundel et al., 1993; Brennan, 1995; Ariel, 1990). Because a pronoun carries less information compared to more specified forms (e.g., *she* vs. *Hillary Clinton*), theories predict that speakers tend to use pronouns when they

think that a referent is sufficiently salient in the discourse. When the referent is less salient, more specified forms are used. In other words, the likelihood of pronominalization increases as referents become more salient.

Topic modeling (Blei et al., 2003; Griffiths et al., 2007) uses a probabilistic model that recovers a latent topic representation from observed words in a document. The model assumes that words appearing in documents have been generated from a mixture of latent topics. These latent topics have been argued to provide a coarse semantic representation of documents and to be in close correspondence with many aspects of human semantic cognition (Griffiths et al., 2007). This previous work has focused on semantic relationships among words and documents. While it is often assumed that the topics extracted by topic models correspond to the gist of a document, and although topic models have been used to capture discourse-level properties in some settings (Nguyen et al., 2013), the ability of topic models to capture cognitive aspects of speakers' discourse representations has not yet been tested.

In this paper we use topic modeling to formalize the idea of salience in the discourse. We focus specifically on the idea of topicality as a predictor of salience (Ariel, 1990; Arnold, 1998) and ask whether the latent topics that are recovered by topic models can predict speakers' choices of referring expressions. Simulations show that the referents of pronouns belong, on average, to higher probability topics than the referents of full noun phrases, indicating that topical referents are more likely to be pronominalized. This suggests that

the information recovered by topic models is relevant to speakers' choices of referring expressions and that topic models can provide a useful tool for quantifying speakers' representations of entities in the discourse.

The structure of this paper is as follows. Section 2 briefly reviews studies that look at the correlation between saliency and choices of referring expression, focusing on topicality, and introduces our approach to this problem. Section 3 describes a model that learns a latent topic distribution and formalizes the notion of topicality within this framework. Section 4 describes the data we used for our simulation. Section 5 shows simulation results. Section 6 discusses implications and future directions.

## 2 Saliency and referring expressions

Various factors have been proposed to influence referent saliency (Arnold, 1998; Arnold, 2010). These factors include **givenness** (Chafe, 1976; Gundel et al., 1993), **grammatical position** (Brennan, 1995; Stevenson et al., 1994), **order of mention** (Järvikivi et al., 2005; Kaiser and Trueswell, 2008), **recency** (Givón, 1983; Arnold, 1998), **syntactic focus and syntactic topic** (Cowles et al., 2007; Foraker and McElree, 2007; Walker et al., 1994), **parallelism** (Chambers and Smyth, 1998; Arnold, 1998), **thematic role** (Stevenson et al., 1994; Arnold, 2001; Rohde et al., 2007), **coherence relation** (Kehler, 2002; Rohde et al., 2007) and **topicality** (Ariel, 1990; Arnold, 1998; Arnold, 1999). Psycholinguistic experiments (Arnold, 1998; Arnold, 2001; Kaiser, 2006) show that determining the salient referent is a complex process which is affected by various sources of information, and that these multiple factors have different strengths of influence.

Among the numerous factors influencing the saliency of a referent, this study focuses on *topicality*. In contrast to surface-level factors such as grammatical position, order of mention, and recency, the representation of topicality is latent and requires inference. Because of this latent representation, it has been challenging to investigate the role of topicality in discourse.

Many researchers have observed that there is a correlation between a linguistic category “topic” and referent saliency and have suggested that topical referents are more likely to be pronominalized (Ariel, 1990; Dahl and Fraurud, 1996). How-

ever, Arnold (2010) points out that examining the relation between topicality and choices of referring expressions is difficult for two reasons. First, identifying the topic is known to be hard. Arnold (2010) shows that it is hard to determine what the topic is even in a simple sentence like *Andy brews beer* (Is the topic *Andy*, *beer*, or *brewing*?). Second, researchers have defined the notion of “topic” differently as follows.

- The topic is often defined as what the sentence is about (Reinhart, 1981).
- The topic can be defined as prominent characters such as the protagonist (Francik, 1985).
- The topic is often associated with old information (Gundel et al., 1993).
- The subject position is considered to be a topical position (Chafe, 1976).
- Repeated mentions are topical (Kameyama, 1994).
- Psycholinguistic experiments define a discourse topic as a referent that has already been mentioned in the preceding discourse as a pronoun/the topic of a cleft (Arnold, 1999) or realized in subject position (Cowles, 2003).
- Centering theory (Grosz et al., 1995; Brennan, 1995) formalizes the topic as a backward-looking center that is a single entity mentioned in the last sentence and in the most salient grammatical position (the grammatical subject is the most salient, and followed by the object and oblique object).
- Givón (1983) suggests that all discourse entities are topical but that topicality is defined by a gradient/continuous property. Givón shows that three measures of topicality – *recency* (the distance between the referent and the referring expression), *persistence* (how long the referent would remain in the subsequent discourse), and *potential interference* (how many other potential referents of the referring expression there are in the preceding discourse) – correlate with the types of reference expressions. Note that these scales measure topicality of the *referring expression*, but not the referent per se.

The variation in the literature seems to derive from three fundamental properties. First, as Arnold (2010) pointed out, there is variation in the

linguistic unit that bears the topic. For example, Reinhart (1981) defines each *sentence* as having a single topic, whereas Givón (1983) defines each *entity* as having a single topic. Second, there is a variation in type of variable. For example, Givón (1983) defines topicality as a continuous property, whereas Centering seems to treat topicality as categorical based on the grammatical position of the referent. Third, many studies define ‘topic’ as a combination of surface linguistic factors such as grammatical position and recency. When topicality is defined in terms of meaning, as in Reinhart (1981), we face difficulty in identifying *what the topic is*, as summarized in Arnold (1998). None of the existing definitions/measures seem to provide a way to capture latent topic representations, and this makes it challenging to investigate their role in discourse representations. It is this idea of latent topic representations that we aim to formalize.

Our study investigates whether topic modeling (Blei et al., 2003; Griffiths et al., 2007) can be used to formalize the relationship between topicality and choices of referring expressions. Because of their structured representations, consisting of a set of topics as well as information about which words belong to those topics, topic models are able to capture topicality by means of semantic associations. For example, observing a word *Clinton* increases the topicality of other words associated with the topic that *Clinton* belongs to, e.g., *president*, *Washington* and so on. In other words, topic models can capture not only the salience of referents within a document, but also the salience of referents via the structured topic representation learned from multiple texts.

We use topic modeling to verify the prevailing hypothesis that topical referents are more likely to be pronominalized than lexical nouns. Examining the relationship between topicality and referring expressions using topic modeling provides an opportunity to test how well the representation recovered by topic models corresponds to the cognitive representation of entities in a discourse. If we can recover the observation that topical referents are more likely to be pronominalized than more specified forms, this could indicate that topic models can capture not only aspects of human semantic cognition (Griffiths et al., 2007), but also aspects of a higher level of linguistic representation, discourse.

### 3 Model

#### 3.1 Recovering latent topics

We formalize topicality of referents using topic modeling. Each document is represented as a probability distribution over topics. Each topic is represented as a probability distribution over possible referents in the corpus. In training our topic model, we assume that all lexical nouns in the discourse are potential referents. The topic model is trained only on lexical nouns, excluding all other words. This ensures that the latent topics capture information about which referents typically occur together in documents.<sup>1</sup>

Rather than pre-specifying a number of latent topics, we use the hierarchical Dirichlet process (Teh et al., 2006), which learns a number of topics to flexibly represent input data. The summary of the generative process is as follows.

1. Draw a global topic distribution  $G_0 \sim \text{DP}(\gamma, H)$  (where  $\gamma$  is a hyperparameter and  $H$  is a base distribution).
2. For each document  $d \in \{1, \dots, D\}$  (where  $D$  denotes the number of documents in the corpus),
  - (a) draw a document-topic distribution  $G_d \sim \text{DP}(\alpha_0, G_0)$  (where  $\alpha_0$  is a hyperparameter).
  - (b) For each referent  $r \in \{1, \dots, N_d\}$  (where  $N_d$  denotes the number of referents in document  $d$ ),
    - i. draw a topic parameter  $\phi_{d,r} \sim G_d$ .
    - ii. draw a word  $x_{d,r} \sim \text{Mult}(\phi_{d,r})$ .

This process generates a distribution over topics for each document, a distribution over referents for each topic, and a topic assignment for each referent. The distribution over topics for each document represents what the topics of the document are. The distribution over referents for each topic represents what the topic is about. An illustration of this representation is in Table 3.1. Topics and words that appear in the second and third columns are ordered from highest to lowest. We can represent topicality of the referents using this

<sup>1</sup>Excluding pronouns from the training set introduces a confound, because it artificially lowers the probability of the topics corresponding to those pronouns. However, in this paper our predicted effect goes in the opposite direction: we predict that topics corresponding to the referents of pronouns will have higher probability than those corresponding to the referents of lexical nouns. Excluding pronouns thus makes us less likely to find support for our hypothesis.

probabilistic latent topic representation, measuring which topics have high probability and assuming that referents associated with high probability topics are likely to be topical in the discourse.

Word	Top 3 topic IDs	Associated words in the 1st topic
Clinton	5, 26, 61	president, meeting, peace, Washington, talks
FBI	148, 73, 67	Leung, charges, Katrina, documents, indictment
oil	91, 145, 140	Burmah, Iraq, SHV, coda, pipeline

Table 1: Illustration of the topic distribution

Given this generative process, we can use Bayesian inference to recover the latent topic distribution. We use the Gibbs sampling algorithm in Teh et al. (2006) to estimate the conditional distribution of the latent structure, the distributions over topics associated with each document, and the distributions over words associated with each topic. The state space consists of latent variables for topic assignments, which we refer to as  $\mathbf{z} = \{z_{d,r}\}$ . In each iteration we compute the conditional distribution  $p(z_{d,r} | \mathbf{x}, \mathbf{z}_{-d,r}, *)$ , where the subscript  $-d, r$  denotes counts without considering  $z_{d,r}$  and  $*$  denotes all hyperparameters. Recovering these latent variables allows us to determine what the topic of the referent is and how likely that topic is in a particular document. We use the latent topic and its probability to represent topicality.

### 3.2 A measure of topicality

Discourse theories predict that topical referents are more likely to be pronominalized than more specified expressions.<sup>2</sup> We can quantify the effect of topicality on choices of referring expressions by comparing the topicality of the referents of two types of referring expressions, pronouns and lexical nouns. If topical words are more likely to be pronominalized, then the topicality of the referents of pronouns should be higher than the topicality of the referents of lexical nouns.

Annotated coreference chains in the corpus, described below, are used to determine the referent of each referring expression. We look at the topic assigned to each referent  $r$  in document  $d$  by the topic model,  $z_{d,r}$ . We take the log probability

<sup>2</sup>Although theories make more fine-grained predictions on the choices of referring expressions with respect to saliency, e.g., a full name is used to refer to less salient entity compared to a definite description (c.f. accessibility marking scale in Ariel 1990), we focus here on the coarse contrast between pronouns and lexical nouns.

of this topic within the document,  $\log p(z_{d,r} | G_d)$ , as a measure of the topicality of the referent. We take the expectation over a uniform distribution of referents, where the uniform distributions are denoted  $u(\text{lex})$  and  $u(\text{pro})$ , to obtain an estimate of the average topicality of the referents of lexical nouns,  $\mathbb{E}_{u(\text{lex})} [\log p(z_{d,r} | G_d)]$ , and the average topicality of the referents of pronouns,  $\mathbb{E}_{u(\text{pro})} [\log p(z_{d,r} | G_d)]$ , within each document. The expectation for the referents of the pronouns in a document is computed as

$$\mathbb{E}_{u(\text{pro})} [\log p(z_{d,r} | G_d)] = \frac{\sum_{r=1}^{N_{d,\text{pro}}} \log p(z_{d,r} | G_d)}{N_{d,\text{pro}}} \quad (1)$$

where  $N_{d,\text{pro}}$  denotes the number of pronouns in a document  $d$ . Replacing  $N_{d,\text{pro}}$  with  $N_{d,\text{lex}}$  (the number of lexical nouns in a document  $d$ ) gives us the expectation for the referents of lexical nouns.

To obtain a single measure for each document of the extent to which our measure of topicality predicts speakers' choices of referring expressions, we subtract the average topicality for the referents of lexical nouns from the average topicality for the referents of pronouns within the document to obtain a log likelihood ratio  $q_d$ ,

$$q_d = \mathbb{E}_{u(\text{pro})} [\log p(z_{d,r} | G_d)] - \mathbb{E}_{u(\text{lex})} [\log p(z_{d,r} | G_d)] \quad (2)$$

A value of  $q_d$  greater than zero indicates that the referents of pronouns are more likely to be topical than the referents of lexical nouns.

## 4 Annotated coreference data

Our simulations use a training set of the Ontonotes corpus (Pradhan et al., 2007), which consists of news texts. We use these data because each entity in the corpus has a coreference annotation. We use the coreference annotations in our evaluation, described above. The training set in the corpus consists of 229 documents, which contain 3,648 sentences and 79,060 word tokens. We extract only lexical nouns (23,084 tokens) and pronouns (2,867 tokens) from the corpus as input to the model.<sup>3</sup>

Some preprocessing is necessary before using these data as input to a topic model. This necessity arises because some entities in the corpus are represented as phrases, such as in (1a) and (1b) below,

<sup>3</sup>In particular, we extracted words that are tagged as NN, NNS, NNP, NNPS, and for pronouns as PRP, PRPS.

where numbers following each expression represent the entity ID that is assigned to this expression in the annotated corpus. However, topic models use bag-of-words representations and therefore assign latent topic structure only to individual words, and not to entire phrases. We preprocessed these entities as in (2). This enabled us to attribute entity IDs to individual words, rather than entire phrases, allowing us to establish a correspondence between these ID numbers and the latent topics recovered by our model for the same words.

1. Before preprocessing
  - (a) a tradition in Betsy’s family: 352
  - (b) Betsy’s family: 348
  - (c) Betsy: 184
2. After preprocessing
  - (a) tradition: 352
  - (b) family: 348
  - (c) Betsy: 184

Annotated coreference chains in the corpus were used to determine the referent of each pronoun and lexical noun. The annotations group all referring expressions in a document that refer to the same entity together into one coreference chain, with the order of expressions in the chain corresponding to the order in which they appear in the document. We assume that the referent for each pronoun and lexical noun appears in its coreference chain. We further assume that the referent needs to be a lexical noun, and thus exclude all pronouns from consideration as referents. If a lexical noun does not have any other words before it in the coreference chain, i.e., that noun is the first or the only word in that coreference chain, we assume that this noun refers to itself (the noun itself is the referent). Otherwise, if a coreference chain has multiple referents, we take its referent to be the lexical noun that is before and closest to the target word.

## 5 Results

To recover the latent topic distribution, we ran 5 independent Gibbs sampling chains for 1000 iterations.<sup>4</sup> Hyperparameters  $\gamma$ ,  $\alpha_0$ , and  $\eta$  were fixed at 1.0, 1.0, and 0.01, respectively.<sup>5</sup> The model re-

<sup>4</sup>We used a Python version of the hierarchical Dirichlet process implemented by Ke Zhai (<http://github.com/kzhai/PyNPB/tree/master/src/hdp>).

<sup>5</sup>Parameter  $\gamma$  controls how likely a new topic is to be created in the corpus. If the value of  $\gamma$  is high, more topics are

covered an average of 161 topics (range: 160–163 topics).

We computed the log likelihood ratio  $q_d$  (Equation 2) for each document and took the average of this value across documents for each chain. The formula to compute this average is as follows.

For each chain  $g$ ,

1. get the final sample  $s$  in  $g$ .
2. For each document  $d$  in the corpus,
  - i. compute  $q_d$  based on  $s$ .
3. Compute the average of all  $q_d$  in the corpus.

The average log likelihood ratio in each chain consistently shows values greater than zero across the 5 chains. The average log likelihood ratio across chains is 1.0625 with standard deviation 0.7329. As an example, in one chain, the average of the expected values for the referents of pronouns across documents is  $-1.1849$  with standard deviation 0.8796. In the same chain, the average of the expected values for the referents of lexical nouns across documents is  $-2.2356$  with standard deviation 0.5009.

We used the median test<sup>6</sup> to evaluate whether the two groups of the referents are different with respect to the expected values of the log probabilities of topics. The test shows a significant difference between two groups ( $p < 0.0001$ ).

We also computed the probability density  $p(q)$  from the log likelihood ratio  $q_d$  for each document using the final samples from each chain. Graph 1 shows the probability density  $p(q)$  from each chain. The peak after zero confirms the observed effect.

Table 2 shows examples of target pronouns and lexical nouns, their referents, and the topic assigned to each referent from a document. Table 3 shows the distribution over topics in the document obtained from one chain. Topics in Table 3 are ordered from highest to lowest. Only four topics were present in this document. The list of referents associated with each topic in Table 3 is recovered from the topic distribution over referents. This list shows what the topic is about.

discovered in the corpus. Parameter  $\alpha_0$  controls the sparseness of the distribution over topics in a document, and parameter  $\eta$  controls the sparseness of the distribution over words in a topic.

<sup>6</sup>The median test compares medians to test group differences (Siegel, 1956).

Topic ID	Associated words	Probability
1	Milosevic, Kostunica, Slobodan, president, Belgrade, Serbia, Vojislav, Yugoslavia, crimes, parliament	0.64
2	president, Clinton, meeting, peace, Washington, talks, visit, negotiators, region, . . . , Albanians	0.16
3	people, years, U.S., president, time, government, today, country, world, way, year	0.16
4	government, minister, party, Barak, today, prime, east, parliament, leader, opposition, peace, leadership	0.04

Table 3: The document-topic distribution

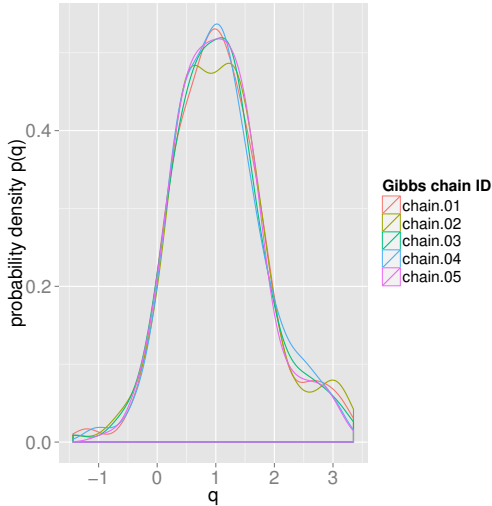


Figure 1: The probability density of  $p(q)$

Target	Referent	Referent's Topic ID
his	Spilanovic	1
he	Spilanovic	1
its	Belgrade	1
Goran	Minister	4
Albanians	Albanians	2
Kosovo	Kosovo	1

Table 2: Target words, their corresponding referents, and the assigned topics of the referents

The topics associated with the pronouns *his*, *he* and *its* have the highest probability in the document-topic distribution, as shown in Table 3. In contrast, although the topic associated with the word *Kosovo* has the highest probability in the document-topic distribution, the topics associated with nouns *Goran* and *Albanians* do not have high probability in the document-topic distribution. This is an example from one document, but this tendency is observed in most of the documents in the corpus.

These results indicate that the referents of pronouns are more topical than the referents of lexical nouns using our measure of topicality derived from the topic model. This suggests that our measure of topicality captures aspects of salience that influence choices of referring expressions.

However, there is a possibility that the effect we observed is simply derived from referent frequencies and that topic modeling structure does not play a role beyond this. Tily and Piantadosi (2009) found that the frequency of referents has a significant effect on predicting the upcoming referent. Although their finding is about comprehender’s ability to predict the upcoming referent (not the type of referring expression), we conducted an additional analysis to rule out the possibility that referent frequencies alone were driving our results.

In order to quantify the effect of referent frequency on choices of referring expressions, we computed the same log likelihood ratio  $q_d$  with referent probabilities. The probability of a referent in a document was computed as follows:

$$p(r_i | doc_d) = \frac{C_{d,r_i}}{C_{d,\cdot}} \quad (3)$$

where  $C_{d,r_i}$  denotes the number of mentions that refer to referent  $r_i$  in document  $d$  and  $C_{d,\cdot}$  denotes the total number of mentions in document  $d$ . We can directly compute this value by using the annotated coreference chains in the corpus.

The log likelihood ratio for this measure is 2.3562. The average of the expected values for the referents of pronouns across documents is  $-1.1993$  with standard deviation 0.6812. The average of the expected values for the referents of lexical nouns across documents is  $-3.5556$  with standard deviation 0.9742. The median test shows a significant difference between two groups. ( $p < 0.0001$ ). These results indicate that the frequency of a referent captures aspects of its salience that influence choices of referring expressions, raising the question of whether our latent topic representations capture something that simple referent frequencies do not.

In order to examine to what extent the relationship between topicality and referring expressions captures information that goes beyond simple referent frequencies, we compare two logistic regres-

sion models.<sup>7</sup> Both models are built to predict whether a referent will be a full noun phrase or a pronoun. The first model incorporates only the log probability of the referent as a predictor, whereas the second includes both the log probability of the referent and our topicality measure as predictors.<sup>8</sup>

The null hypothesis is that removing our topicality measure from the second model makes no difference for predicting the types of referring expressions. Under this null hypothesis, twice the difference in the log likelihoods between the two models should follow a  $\chi^2(1)$  distribution. We find a significant difference in likelihood between these two models ( $\chi^2(1) = 118.38, p < 0.0001$ ), indicating that the latent measure of topicality derived from the topic model predicts aspects of listeners' choices of referring expressions that are not predicted by the probabilities of individual referents.

## 6 Discussion

In this study we formalized the correlation between topicality and choices of referring expressions using a latent topic representation obtained through topic modeling. Both quantitative and qualitative results showed that according to this latent topic representation, the referents of pronouns are more likely to be topical than the referents of lexical nouns. This suggests that topic models can capture aspects of discourse representations that are relevant to the selection of referring expressions. We also showed that this latent topic representation has an independent contribution beyond simple referent frequency.

This study examined only two independent factors: topicality and referent frequency. However, discourse studies suggest that the salience of a referent is determined by various sources of information and multiple discourse factors with different strengths of influence (Arnold, 2010). Our framework could eventually form part of a more complex model that explicitly formalizes the interaction of information source and various discourse factors. Having a formal model would help by allowing us to test different hypotheses and develop a firm theory regarding cognitive representations of entities in the discourse.

<sup>7</sup>Models were fit using `glm` in R. For the log-likelihood ratio test, `lrtest` in R package `epicalc` was used.

<sup>8</sup>We also ran a version of this comparison in which frequency of mention was included as a predictor in both models, and obtained similar results.

One possibility for exploring the role of various discourse factors in our framework is to use recent advances in topic modeling. For example, TagLDA (Zhu et al., 2006) includes part-of-speech as part of the model, and syntactic topic models (Boyd-Graber and Blei, 2008) incorporate syntactic information. Whereas simulations in our study only used nouns as input, it has been observed that the thematic role of the entity influences referent salience (Stevenson et al., 1994; Arnold, 2001; Rohde et al., 2007). Using part-of-speech and syntactic information together with the topic information could help us approximate the influence of the thematic role and allow us to simulate how this factor interacts with latent topic information and other factors.

It has been challenging to quantify the influence of latent factors such as topicality, and the simulations in this paper represent only a first step toward capturing these challenging factors. The simulations nevertheless provide an example of how formal models can help us validate theories of the relationship between speakers' discourse representations and the language they produce.

## Acknowledgments

We thank Ke Zhai, Viet-An Nguyen, and four anonymous reviewers for helpful comments and discussion.

## References

- Mira Ariel. 1990. *Accessing noun-phrase antecedents*. Routledge, London.
- Jennifer Arnold. 1998. *Reference form and discourse patterns*. Ph.D. thesis, Stanford University Stanford, CA.
- Jennifer Arnold. 1999. Marking salience: The similarity of topic and focus. *Unpublished manuscript, University of Pennsylvania*.
- Jennifer Arnold. 2001. The effect of thematic roles on pronoun use and frequency of reference continuation. *Discourse Processes*, 31(2):137–162.
- Jennifer Arnold. 2010. How speakers refer: the role of accessibility. *Language and Linguistics Compass*, 4(4):187–203.
- David M Blei, Andrew Y Ng, and Michael I Jordan. 2003. Latent Dirichlet allocation. *Journal of Machine Learning Research*, 3:993–1022.
- Jordan L Boyd-Graber and David M Blei. 2008. Syntactic topic models. In *Neural Information Processing Systems*, pages 185–192.

- Susan E Brennan. 1995. Centering attention in discourse. *Language and Cognitive Processes*, 10(2):137–167.
- Wallace Chafe. 1976. Givenness, contrastiveness, definiteness, subjects, topics, and point of view. In C. N. Li, editor, *Subject and Topic*. Academic Press, New York.
- Craig G Chambers and Ron Smyth. 1998. Structural parallelism and discourse coherence: A test of Centering theory. *Journal of Memory and Language*, 39(4):593–608.
- H Wind Cowles, Matthew Walenski, and Robert Klender. 2007. Linguistic and cognitive prominence in anaphor resolution: topic, contrastive focus and pronouns. *Topoi*, 26(1):3–18.
- Heidi Wind Cowles. 2003. *Processing information structure: Evidence from comprehension and production*. Ph.D. thesis, University of California, San Diego.
- Osten Dahl and Kari Fraurud. 1996. Animacy in grammar and discourse. *Pragmatics and Beyond New Series*, pages 47–64.
- Stephani Foraker and Brian McElree. 2007. The role of prominence in pronoun resolution: Active versus passive representations. *Journal of Memory and Language*, 56(3):357–383.
- Ellen Palmer Francik. 1985. *Referential choice and focus of attention in narratives (discourse anaphora, topic continuity, language production)*. Ph.D. thesis, Stanford University.
- Talmy Givón. 1983. *Topic continuity in discourse: A quantitative cross-language study*, volume 3. John Benjamins Publishing.
- Thomas L Griffiths, Mark Steyvers, and Joshua B Tenenbaum. 2007. Topics in semantic representation. *Psychological Review*, 114(2):211.
- Barbara J Grosz, Scott Weinstein, and Aravind K Joshi. 1995. Centering: A framework for modeling the local coherence of discourse. *Computational Linguistics*, 21(2):203–225.
- Jeanette K Gundel, Nancy Hedberg, and Ron Zacharski. 1993. Cognitive status and the form of referring expressions in discourse. *Language*, pages 274–307.
- Juhani Järvi­kivi, Roger PG van Gompel, Jukka Hyönä, and Raymond Bertram. 2005. Ambiguous pronoun resolution contrasting the first-mention and subject-preference accounts. *Psychological Science*, 16(4):260–264.
- Elsi Kaiser and John C Trueswell. 2008. Interpreting pronouns and demonstratives in Finnish: Evidence for a form-specific approach to reference resolution. *Language and Cognitive Processes*, 23(5):709–748.
- Elsi Kaiser. 2006. Effects of topic and focus on saliency. In *Proceedings of Sinn und Bedeutung*, volume 10, pages 139–154. Citeseer.
- Megumi Kameyama. 1994. Indefeasible semantics and defeasible pragmatics. In *CWI Report CS-R9441 and SRI Technical Note 544*. Citeseer.
- Andrew Kehler. 2002. *Coherence, reference, and the theory of grammar*. CSLI publications, Stanford, CA.
- Viet-An Nguyen, Jordan Boyd-Graber, Philip Resnik, Deborah A Cai, Jennifer E Midberry, and Yuanxin Wang. 2013. Modeling topic control to detect influence in conversations using nonparametric topic models. *Machine Learning*, pages 1–41.
- Sameer S Pradhan, Eduard Hovy, Mitch Marcus, Martha Palmer, Lance Ramshaw, and Ralph Weischedel. 2007. Ontonotes: A unified relational semantic representation. *International Journal of Semantic Computing*, 1(4):405–419.
- Tanya Reinhart. 1981. Pragmatics and linguistics: An analysis of sentence topics in pragmatics and philosophy I. *Philosophica*, 27(1):53–94.
- Hannah Rohde, Andrew Kehler, and Jeffrey L Elman. 2007. Pronoun interpretation as a side effect of discourse coherence. In *Proceedings of the 29th Annual Conference of the Cognitive Science Society*, pages 617–622.
- Sidney Siegel. 1956. *Nonparametric statistics for the behavioral sciences*. McGraw-Hill.
- Rosemary J Stevenson, Rosalind A Crawley, and David Kleinman. 1994. Thematic roles, focus and the representation of events. *Language and Cognitive Processes*, 9(4):519–548.
- Y. W. Teh, M. I. Jordan, M. J. Beal, and D. M. Blei. 2006. Hierarchical Dirichlet processes. *Journal of the American Statistical Association*, 101:1566–1581.
- Harry Tily and Steven Piantadosi. 2009. Refer efficiently: Use less informative expressions for more predictable meanings. In *Proceedings of the workshop on the production of referring expressions: Bridging the gap between computational and empirical approaches to reference*.
- Marilyn Walker, Sharon Cote, and Masayo Iida. 1994. Japanese discourse and the process of centering. *Computational Linguistics*, 20(2):193–232.
- Xiaojin Zhu, David Blei, and John Lafferty. 2006. TagLDA: Bringing document structure knowledge into topic models. Technical report, Technical Report TR-1553, University of Wisconsin.