# Determining Definiteness from Context

**Christopher Brown**
University of Texas at Austin
chrisbrown@utexas.edu

## Abstract

Noun phrases occurring in natural language can be described as either definite or indefinite. Definiteness is a semantic concept that relates to givenness, familiarity, topicality, and focus. It can be denoted by a variety of determiners, the absence of a determiner, or by context. This project focuses on the definite markers, *a* and *the*, and investigates what types of models are best at predicting the most suitable determiner when it is not available.

## 1 Problem and Motivation

Definiteness is an attribute of a noun phrase that describes the nature of its intended referent.

(1) **A dog** was barking all night.

(2) **Dogs** were barking all night.

(3) **This dog** was barking all night.

(4) **The dog** was barking all night.

In each case, the bold noun phrase picks out a different set of dogs or implies something different about which dog the speaker means to refer to. (1) is the least definite; (2) is similar, but allows the dogs to take shifts; (3) suggests the speaker has some way of uniquely identifying the dog, but the hearer would not know which dog without more information; (4) requires both hearer and speaker to know which dog is referred to.

In English, determiners are the primary overt markers of definiteness / indefiniteness. Definiteness is not completely specified by determiners, though; some uses of *a* can be definite ("A man just proposed to me in the orangery."[1]), while *this* and bare plurals are ambiguous. There is a wide variety of other determiners and quantifiers, as well as most pronouns, which mark definiteness in some way.

I could not find any corpora annotated for definiteness; **?** describes a corpus where 'Information Status' is encoded, but not specifically definiteness. In many cases, definiteness is clear from the determiner, as with *a* and *the*. Noun phrases with *a* (and *an*) are predominantly indefinite, and those with *the* are always definite. The analysis described in this paper assumes that these particular determiners cover enough cases to serve as a significant portion of the problem, and that *a/an* always denotes indefiniteness, and *the* always denotes definiteness.

### 1.1 Ubiquity

Determiners are very common in the English language (Table 1); like many other functional words, determiners are often redundant—misuse is very quickly detected, and the intended use can usually be determined from context.

---

[1] Fodor and Sag, 'Referential and Quantificational Indefinites', *Linguistics and Philosophy* 5. 1982.

| Token | Count | Percentage |
|---|---|---|
| **the** | 69968 | 6.83% |
| of | | |
| and | | |
| to | | |
| **a** | 23490 | 2.29% |
| in | | |
| **that** | 10786 | 1.05% |
| … | | |
| **this** | 5145 | 0.50% |
| … | | |
| **an** | 3746 | 0.36% |

Table 1: Brown corpus counts and proportions of selected definiteness markers (in bold).

## 1.2 Simplifying the problem

Definiteness effects overlap with other semantic features, particularly presupposition and the scope of context. As with many linguistic aspects of semantics / pragmatics, a complete solution would require a much fuller understanding of language than any current NLP system provides.

In this paper, I show that it's possible to recover some information about definiteness from the noun phrases and

## 2 Applications

The task developed in this paper and detecting definiteness in general fill a very specific niche of natural language processing and understanding. Determiners serve as useful features in POS (part of speech) tagging, NER (named entity recognition), and likely many other NLP tasks. Such methods regard determiners as a side-effect, though, and don't necessarily distinguish one determiner from another or understand the determiners as independent linguistic features.

This linguistic analysis definiteness processing primarily contributes to two common tasks: language generation and machine translation.

### 2.1 Language generation

When generating language from a logical structure, co-referential noun phrases representing a single entity must surface fluently. It is unnatural to repeat a proper noun or an extended noun phrase when a pronoun would suffice. But context and focus shift as the discourse develops, eventually making any pronoun ambiguous, forcing the speaker to use a more explicit reference.

- Centering

Definiteness is just one part of anaphora generation, but this paper's task does not precisely tackle this question. For present purposes, the aspect that is relevant here requires shallower understanding. If we use a full noun phrase more than once, we need to know when to use an indefinite description and when to use the definite form.

### 2.2 Translation

Some languages, like Russian and Korean, do not have determiners to denote definiteness; others, like Japanese and Hindi, have a few articles that convey some semantic level of *in*definiteness, usually, or are used for other special purposes (such as to denote humans). When translating between these languages and languages with overt definiteness markers, we must insert determiners that often have no parallel marker in the source text.

For example, Google Translate renders each of the four alternations of "A/The man bit a/the dog" as Человек укусил собаку. Translating back into English produces "A person bitten by a dog." I don't know Russian; it may be genuinely ambiguous, and only determinable via context; but I know that translating back into English accurately requires choosing the correct determiners, for which we need a model of definiteness (as well as anaphora / centering).

There are a few papers on the second of these applications, translation, but they are specific to only a few language pairs (**??**).

### 2.3 General contribution

This can explain why so little computational work exists on the subject, but the linguistic and philosophical literature on definiteness is much more developed. However

These applications are useful to two specific niches of NLP, , and apply only in some cases of those areas.

- simplifications used here. - 'a' definite use: say we'll just roll it into noise, here

## 3  Experiment

The basic research question will be to determine what type of language models produce the best prediction of definiteness. Because English has determiners, any parsed document is "labeled" data. I will simply replace *a* and *the*'s with the placeholder «DET», retaining the original token as the label for that token. While this is an artificial evaluation metric, it could be useful in both of the application scenarios. In anaphora generation, each instance has to be filled by either a full noun phrase, either definite or indefinite, or a pronoun. In the noun phrase case, a placeholder would be inserted by a first pass of the language generation model, and then resolved by the anaphora resolution pass. In machine translation, noun phrases translated from a determiner-less language to one with determiners could be padded with a placeholder, which would then be resolved by some post-processing step performed on just the target language.

My project will consist of evaluating different models for predicting the deleted marker; it is not a typical sequence labeling problem, since relatively few of the tokens need to be resolved. But neither is each instance independent of the previous; in fact, a prior mention is presumably a strong indicator that an entity is now in the active context, and subsequent instances should use the definite determiner. The sequence is crucial in a sense of building and maintaining a 'center,' which is a common approach when resolving anaphora (??). But with an intelligent set of features, we might achieve sufficient accuracy with a simple logistic regression; given a set of features, evaluate whether the placeholder should be definite or indefinite?

My investigation will also compare feature sets, to determine which are the most efficient at predicting the deleted definiteness markers. Is the parse structure important, or just the surrounding tokens? Position in sentence / document? One presumably useful feature would be to use a cumulative index for each token type in the document—1 for the first instance, 2 for the second, and so on—which might handle the prior mention issue without requiring a sequence model.

## 4

## 5  Semantics of definiteness

The theoretical semantics of definite noun phrases has an extensive history; we can start with Gottlob Frege and Bertrand Russell and names. Frege said that the meaning of a name (or definite noun phrase) was the object denoted by that description. Russell complicated matters with his "The King of France is bald" example, in which "the King of France" cannot refer to anything (at least not since 1848). Russell would call this sentence false, despite the description's failure to refer to anything; the failure of the noun phrase propagates out to the failure of the sentence.

Keith Donnellan and P. F. Strawson disagree; they say that referencing is something speakers, not words, do, and that "the King of France" is simply unsuccessful. Compare "a King of France is bald," which can refer to any (dead) past King of France, and thus have a truth value. Definite noun phrase usually denote a unique entity, and if this fails, the sentence is incoherent. But compare "After you enter the lobby, take the elevator to the 13th floor," which is acceptable even if there is more than one elevator.

This is just the beginning, but it's clear that the difference between definite and indefinite noun phrases involves a number of factors, such as context, real world knowledge, and other pragmatic phenomena. It is an active field of study in linguistics or philosophy of language, but has received very little attention in computational linguistics.