

# Lexical Clustering and Definite Description Interpretation

**Massimo Poesio**  
University of Edinburgh  
Cognitive Science / HCRC  
2 Buccleuch Place  
Edinburgh EH8 9LW  
Scotland, UK  
poesio@cogsci.ed.ac.uk

**Sabine Schulte im Walde**  
University of Stuttgart  
IMS  
Azenbergstrasse 12  
70174 Stuttgart  
Germany  
schulte@ims.uni-stuttgart.de

**Chris Brew**  
University of Edinburgh  
LTG  
2 Buccleuch Place  
Edinburgh EH8 9LW  
Scotland, UK  
chrisbr@cogsci.ed.ac.uk

## Abstract

We present preliminary results concerning the use of lexical clustering algorithms to acquire the kind of lexical knowledge needed to resolve definite descriptions, and in particular what we call ‘inferential’ descriptions. We tested the hypothesis that the antecedent of an inferential description is primarily identified on the basis of its semantic distance from the description; we also tested several variants of the clustering algorithm. We found that the choice of parameters has a clear effect, and that the best results are obtained by measuring the distance between lexical vectors using the cosine measure. We also found, however, that factors other than semantic distance play the main role in the majority of cases; but in those cases in which the sort of lexical knowledge we acquired is the main factor, the algorithms we used performed reasonably well; several standing problems are discussed.

## Introduction

In order to develop systems for anaphoric resolution whose generality and performance can be evaluated in a quantitative fashion—i.e., by testing them over a corpus including texts from different domains—it is necessary to address the issue of commonsense knowledge. The question we are currently studying is what kind of commonsense knowledge is involved in the resolution of definite descriptions; more specifically, we are trying to identify the commonsense knowledge that is brought to bear when resolving definite descriptions whose head noun is not identical to the antecedent, such as *the vehicle* in:

- (1) *John saw a truck stopped at an intersection. THE VEHICLE’s engine was smoking.*

and in so-called BRIDGING DESCRIPTIONS (Clark 1977), i.e., definite descriptions that refer to an object only indirectly introduced into the common ground as the result of the mention of a related object—such as *the door* in *John walked towards the house. THE DOOR was open.*

Arguably, the minimal hypothesis to pursue in this connection is that resolving these descriptions is purely a matter of lexical knowledge—i.e., that the identification of the antecedent depends solely on the degree of association among lexical items. The assumption that the lexicon is organized like a ‘semantic’ network where some concepts are more closely related than others, originally motivated by semantic priming effects (Meyer & Schvaneveldt 1971; Neely 1991), underlies most current psychological models of the lexicon, including WordNet (Miller *et al.* 1990) and has been adopted in much research on reference resolution: such models assume that the antecedent for *the vehicle* in (1) is found by looking for an antecedent whose concept is semantically close (in some sense), and that *the truck* is chosen because the concept associated with this antecedent subsumes the concept associated with the definite description in the semantic network. We will call this the Main Hypothesis:

**Main Hypothesis** Resolving a definite description is a matter of finding the antecedent in the text that primes the head predicate of the definite description most strongly.

Our self-imposed limitation to processing models that can be quantitatively evaluated as discussed above precludes the possibility of coding the information needed to test this hypothesis ourselves, as this is only possible for a narrow domain. One solution is to use an existing source of lexical knowledge, such as WordNet; however, the results we obtained with this method—reported in (Poesio, Vieira, & Teufel 1997)—were not too satisfactory, owing to the incompleteness of the information hand-coded in WordNet, as well as to several inconsistencies we found in it. As a result, we have been exploring techniques for acquiring this information automatically, which also have the advantage that they could be used to acquire domain-specific knowledge when necessary. We report our preliminary results in this paper.

In this initial phase we have mainly been experimenting with clustering algorithms (Charniak 1993). In particular, the work discussed here was inspired by (Lund, Burgess, & Atchley 1995), who reported that the clusters of words obtained with their HAL model of lexical clustering reflect a notion of distance that correlates well with subjects' results on semantic priming tasks. This work offers therefore the opportunity to test the hypothesis discussed above that resolving bridging descriptions is a matter of semantic priming. We tested therefore the performance of several variants of the HAL method on the task of resolving definite descriptions.

## Background

### Inferential Descriptions

Our studies of definite description use (Poesio & Vieira 1998; Vieira & Teufel 1997; Poesio, Vieira, & Teufel 1997) led to the development of a taxonomy of definite descriptions reflecting—if in a rather crude fashion—the types of commonsense knowledge that appear to be involved in their resolution, rather than assumptions about the way the antecedent is identified by the listener and/or its relation to the definite description, as in the taxonomies habitually presented in the literature (Clark & Marshall 1981; Prince 1981). For the purposes of this paper, we will consider definite descriptions as falling in one of the following three categories:<sup>1</sup>

**Anaphoric same head:** these are the definite descriptions whose resolution involves simply matching the head of the antecedent with the head of the definite description, as in *a car ... the car*;<sup>2</sup>

**Inferential:** this is a semantically eclectic class, including both definite descriptions whose head is not identical to that of the antecedent, and those whose relation with the antecedent is not one of coreference (i.e., 'inferrables' in Prince's taxonomy or 'bridging descriptions' in Clark's). This class also includes references to events (as in *John killed Bill. THE MURDER took place at 5pm*), and to entities introduced by proper names, as in *We are celebrating*

<sup>1</sup>As discussed in (Poesio & Vieira 1998), these categories are not completely mutually exclusive. In that paper we also discuss reliability results for this classification scheme.

<sup>2</sup>In fact, even resolving these cases may involve some form of commonsense inference—e.g., to take into account the effects of pre-modifiers and post-modifiers in recognizing that *a blue car* cannot serve as the antecedent of *the red car* in *I saw a blue car and a red car. The red car was a Ferrari*. At the moment we are using heuristic methods in these cases (Vieira & Poesio 1997), but of course we will eventually have to study how to acquire this sort of information, as well.

*this year 200 years since Franz Schubert's birth. THE FAMOUS COMPOSER was born in 1797..*

**Discourse new:** This class consists of those definite descriptions that do not have an antecedent in the text, and includes both references to 'larger situation' knowledge such as *the sun* and possible first-mention definite descriptions such as *the first man to sail to America* (Hawkins 1978; Prince 1981; Poesio & Vieira 1998).

Our treatment of anaphoric same-head descriptions and discourse-new descriptions is discussed elsewhere (Vieira & Poesio 1997; Vieira 1998); in this paper we are exclusively concerned with the class of inferential descriptions. This class was further analyzed in (Vieira & Teufel 1997; Poesio, Vieira, & Teufel 1997) in order to categorize the types of commonsense knowledge involved in their resolutions and to gain a feeling for how many of the required inferences would be supported by a semantic network. The following classes of inferential descriptions were identified:

- **Synonymy:**  
The antecedent and the bridging descriptions are synonymous, as in *a new album – the record*.
- **Hypernymy/Hyponymy:**  
The antecedent and the bridging description are in a *is-a*-relation, as in *rice – the plant* (superordination/hypernymy) or *a plant – the rice* (subordination/hyponymy).
- **Meronymy:**  
The antecedent and the bridging description stand in a *part-of* relation, as in *a tree – the leaves*.
- **Names:**  
The bridging description refers back to a proper name, as in *Bach – the composer*.
- **Compound Nouns:**  
The 'antecedent' occurs as part of a compound noun, as in *the stock market crash – the markets*.
- **Events:**  
The antecedent is not introduced by a noun phrase, but by either a verb phrase or a sentence, e.g. *they planned – the strategy*.
- **Discourse Topic:**  
The antecedent is the implicit 'discourse topic' of a text, as in *the industry* appearing in a text about oil companies.

• **(General) Inference:**

The bridging description is based on more complex inferential relations such as causal inferences, as in *last week's earthquake – the suffering people*.

The first three classes include the inferential descriptions whose resolution we might expect to be supported by the sort of information stored in a typical semantic network such as WordNet; these networks also include information about individuals of the kind needed to resolve definite descriptions in the 'Names' class, and some information of this type is indeed included in WordNet.

Poesio, Vieira and Teufel ran a test on a corpus of 20 parsed Wall Street Journal articles from the Penn Treebank, including 1040 definite descriptions, of which 204 were classified as inferential. Of these 204 descriptions, 38 fell in a class for which one could expect to find the needed information in WordNet. When trying to resolve an inferential description, the discourse entities in the previous five sentences were considered as potential antecedents, and WordNet was queried to find a relation between the inferential description and each antecedent. WordNet found a relation between an inferential description and an antecedent for 107 of these descriptions, but in only 34 cases the right antecedent was suggested; 15 of these cases fell in the Synonymy / Hyponymy / Meronymy category. Separate heuristic-based techniques were also proposed, so that in total 77 descriptions were identified correctly. The overall results with WordNet are presented in (2), whereas the specific results for the 38 Syn/Hyp/Mer cases are in (3).

(2)

| Relationship           | Resolution |
|------------------------|------------|
| Compound Nouns / Names | 19         |
| Syn/Hyp/Mer            | 4/8/3      |
| Events                 | 0          |
| Discourse Topic        | 0          |
| Inference              | 0          |
| Total                  | 34         |

(3)

| Class | Total | Found | Not Found |
|-------|-------|-------|-----------|
| Syn   | 12    | 4     | 8         |
| Hyp   | 14    | 8     | 6         |
| Mer   | 12    | 3     | 9         |
| Total | 38    | 15    | 23        |

**Acquiring Semantic Networks by Clustering**

'Clustering' is a popular approach to lexical acquisition based on the idea that semantically related words are close to each other in some higher-dimensional space representation where they form 'clusters' of similar words—i.e., the very same intuition behind research on semantic networks. Clustering algorithms

view each word as a point in an  $n$ -dimensional space, i.e., as a vector of size  $n$ , and the similarity between words is measured in terms of the distance between the points that represent them. The goal of clustering algorithms is to construct such a representation automatically, exploiting a corpus. These methods differ depending on the dimensions used and their number, on the metric used to measure the distance among the points, and the algorithm used to construct the vectors (Charniak 1993).

A common approach to clustering is to just use words as dimensions, i.e., to let the vector associated with word  $w$ ,  $C(w)$ , be a record of how frequently  $w$  occurred close to word  $w_i$ ; the underlying idea is that a word is defined by the 'company that it keeps', i.e., by the words with which it is most frequently encountered (e.g., (Brown *et al.* 1992)). Algorithms assigning to words vector representations of this type scan a text and whenever they encounter a word  $w$  they increment all cells of  $C(w)$  corresponding to the words  $w_i$  that occur in the vicinity of  $w$ , typically within a window of fixed size. The words chosen as dimensions are often called CONTEXT WORDS.

Once the vectors associated with each word have been constructed in this way, we can estimate the semantic similarity between words by measuring the distance between the associated vectors. A great number of distance measures have been suggested, but the following three are the best known:

• **Manhattan Metric:**

The Manhattan Metric measures the distance of two points in  $n$ -dimensional space by summing the absolute differences of the vectors' elements:

$$d = \sum_{i=1}^n |x_i - y_i|$$

• **Euclidean Distance:**

The Euclidean Distance is calculated by summing the squared differences of the vectors' elements and then determining the square root:

$$d = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

• **Cosine of the Vectors' Angle:**

This measure does not calculate the distance between points, but the angle  $\alpha$  between the  $n$ -dimensional vectors which determine the points in  $n$ -dimensional space:

$$\cos(\alpha) = \frac{\sum_{i=1}^n x_i y_i}{\sqrt{\sum_{i=1}^n x_i^2} \sqrt{\sum_{i=1}^n y_i^2}}$$

The closer the  $\cos(\alpha)$  is to 1, the smaller the angle  $\alpha$  is and therefore the shorter the distance is.

Other measures proposed in the literature include Spearman Rank correlation coefficient, Hellinger distance, and Kullback-Leibler divergence. Weighted

combinations of different measures have also been used. (See, (Levy, Bullinaria, & Patel 1997) for some discussion.)

### Lund et al's HAL Model

Lund et al. (1995) used a 160 million word corpus of articles extracted from all newsgroups containing English dialogue. They chose as context words the 70,000 most frequently occurring symbols within the corpus.

The co-occurrence counts were calculated as follows. They defined a window size of 10 words to the left and to the right of the target words, and within this window, the co-occurrence values were inversely proportional to the number of words separating a specific pair. So, whenever a target word  $w$  was encountered, the context vector  $C(w)$  was incremented as follows: the count  $C(w)(w_1)$  for the word  $w_1$  next to the target word was incremented by 10, the count  $C(w)(w_2)$  for the next word was incremented by 9, and so forth, thus weighting the closeness of the co-occurring words.

To reduce the amount of data, the column variances of the particular vectors used in each experiment were computed, and the columns with the smallest variances were discarded. This left a 200-element vector for each target word.

### Our Methods

In our experiments we adopted the fundamental aspects from the clustering technique of Lund et al, parameterizing several of its aspects in order to evaluate not only the Main Hypothesis, but also the influence of certain parameters on the results. We briefly discuss our methods here; for more discussion and details, see (Schulte im Walde 1997).

As in the case of Lund et al, our basic clustering algorithm involves associating with each word a vector whose dimensions are other words; and again as in their case, the vectors are constructed by scanning a text, considering for each word  $w$  that is encountered all neighbors  $w_i$  in a window of size  $n$ , and increasing by a factor possibly weighted by distance the cells of  $w$ 's vectors associated with each  $w_i$ . This algorithm was made parametric on window size (we considered sizes 1,2,3,5,10,20 and 30), on whether inflected words or their lemmas were considered, and on whether just words or word / tag pairs were used.

We ran some preliminary experiments to determine two additional parameters: corpus size and number of dimensions of the vectors. We set on a 30 million words corpus; as for the dimension of the vectors, we

followed (Huckle 1996) and used the 2,000 most common content words in our corpus as dimensions.

Our algorithm for resolving inferential definite descriptions is as follows. For each definite, all head nouns and head verbs in the previous five sentences are considered as possible antecedents, as in (Poesio, Vieira, & Teufel 1997). For each antecedent, the distance between the vector associated with the head noun of the definite description and the vector associated with the possible antecedent is measured; the antecedent whose vector is closest to that of the definite description is chosen. Three different measures of distance were tried: Manhattan, Euclidean, and Cosine.

We used the British National Corpus<sup>3</sup> for training and the 20 articles from (Poesio, Vieira, & Teufel 1997) to evaluate the results.

## Experiments and Results

### Experiment 1

In order to get a baseline with respect to which to evaluate the actual performance of the method, we ran an experiment in which the antecedent for each inferential description was chosen randomly. Appropriate<sup>4</sup> antecedents for 12 out of 204 inferential descriptions—5.9% of the total—were found with this method.

### Experiment 2

In this second experiment, we trained and resolved over untagged and lemmatized words. We tried window sizes of 1,2,3,5 and 10 words. The results for the three distance measures were as follows, with the three best results in bold (only 195 inferential descriptions were tested in this experiment):

| Metric     | Window size |            |            |
|------------|-------------|------------|------------|
|            | 1           | 2          | 3          |
| <i>Man</i> | 37 (19.0%)  | 36 (18.5%) | 39 (20.0%) |
| <i>Euc</i> | 37 (19.0%)  | 36 (18.5%) | 39 (20.0%) |
| <i>Cos</i> | 39 (20.0%)  | 36 (18.5%) | 39 (20.0%) |

| Metric     | Window size       |                   |
|------------|-------------------|-------------------|
|            | 5                 | 10                |
| <i>Man</i> | <b>41 (21.0%)</b> | 37 (19.0%)        |
| <i>Euc</i> | 39 (20.0%)        | 40 (20.5%)        |
| <i>Cos</i> | <b>42 (21.5%)</b> | <b>46 (23.6%)</b> |

<sup>3</sup>This is a 100-million words collection of both written and spoken language, see <http://info.ox.ac.uk/bnc/>.

<sup>4</sup>An issue to be kept in mind in what follows is that bridging descriptions, unlike other cases of referential expressions, may be related to more than one 'antecedent' in a text, and therefore evaluating the results of a system is more difficult in this case (Poesio & Vieira 1998).

*Cosine* worked best as a distance measure, and the results were better with bigger windows. The best results for *Manhattan Metric* were achieved at window sizes of three and five; for *Euclidean Distance*, the results seemed to get (slightly) better with larger windows.

The following table summarizes the results for each class of inferential descriptions for the best parameter configuration, measure *Cosine*, and window size 10:

| Relationship           | Resolution |
|------------------------|------------|
| <i>Same Head</i>       | 20         |
| <i>Compound Nouns</i>  | 7          |
| <i>Syn/Hyp/Mer</i>     | 2/2/4      |
| <i>Names</i>           | 1          |
| <i>Events</i>          | 5          |
| <i>Discourse Topic</i> | 2          |
| <i>Inference</i>       | 3          |
| <i>Total</i>           | 46         |

We discuss the results in more detail below; we will just note now that the algorithm identified appropriate same-head antecedents for 20 cases we ourselves had classified as inferential.

### Experiment 3

One problem we observed in the second experiment was that lemmatizing might create two identical word-forms out of two different lexemes, usually noun and verb, as in *to plan* and *the plan*, and since we did not distinguish between different parts of speech, the algorithm could not tell the difference. In our third experiment we ran the clustering algorithm and the resolution algorithms on texts in which each word had been tagged, so as to avoid the problem encountered in the previous experiment; and we tried larger window sizes, since it appeared from the previous experiment that larger windows performed better.<sup>5</sup> The results are summarized by the following two tables:

| Metric     | Window size |            |            |            |
|------------|-------------|------------|------------|------------|
|            | 1           | 2          | 3          | 5          |
| <i>Man</i> | 34 (16.8%)  | 35 (17.2%) | 41 (20.2%) | 41 (20.2%) |
| <i>Euc</i> | 35 (17.2%)  | 37 (18.2%) | 37 (18.2%) | 36 (17.7%) |
| <i>Cos</i> | 41 (20.2%)  | 45 (22.1%) | 46 (22.7%) | 41 (20.2%) |

<sup>5</sup>In this second experiment we also tried varying two additional parameters:

- we ran the clustering algorithm giving equal weight to all words in the window, no matter its distance from the word whose vector was being updated;
- we constructed vectors of twice the size, distinguishing between left and right context.

but neither of these changes affected the results (Schulte im Walde 1997).

| Metric     | Window size |            |            |
|------------|-------------|------------|------------|
|            | 10          | 15         | 20         |
| <i>Man</i> | 42 (20.7%)  | 44 (21.7%) | 44 (21.7%) |
| <i>Euc</i> | 37 (18.2%)  | 38 (18.7%) | 39 (19.2%) |
| <i>Cos</i> | 41 (20.2%)  | 38 (18.7%) | 38 (18.7%) |

The interesting fact about these results is that although *Cosine* was again the most successful measure when a window size of 3 was used, increasing the window size made things worse, not better; unlike for *Manhattan Metric*, whose performance improved with larger windows. Anyway, the total number of correctly resolved inferential descriptions did not change.

The per-class results for the two best-performing combinations were as follows:

| Class                  | Measure          |                   |
|------------------------|------------------|-------------------|
|                        | <i>Cos</i> /WS:3 | <i>Man</i> /WS:15 |
| <i>Same Head</i>       | 9 (100.0%)       | 9 (100.0%)        |
| <i>Syn/Hyp/Mer</i>     | 4/2/2 (22.2%)    | 3/1/2 (16.7%)     |
| <i>Names</i>           | 1 (2.3%)         | 4 (9.1%)          |
| <i>Events</i>          | 5 (16.7%)        | 4 (13.3%)         |
| <i>Compound Nouns</i>  | 16 (66.7%)       | 16 (66.7%)        |
| <i>Discourse Topic</i> | 1 (7.1%)         | 2 (14.3%)         |
| <i>Inference</i>       | 6 (13.0%)        | 3 (6.5%)          |
| <i>Total</i>           | 46 (22.7%)       | 44 (21.7%)        |

## Discussion

### Analysis of the Results

Even the best parameter configuration (measure *Cosine*, window size of 10) only resulted in appropriate antecedents for 23.6% of the inferential descriptions. Why was that?

The cases in which an inferential description was not resolved to its correct antecedent fell in the following categories:

- In some cases, the desired antecedent could not be found since it was not on the list of possible antecedents for the bridging description. This happened if the right word was either before the preceding five sentences, or after the description. In this case, another (incorrect) antecedent was still suggested by the algorithm. There were 34 (17.4%) such cases in Experiment 2, 40 (19.6%) in the third one.
- In several cases, the antecedent found by the algorithm is semantically very close to the definite description—in some cases, even closer—but still not the right antecedent: for example, in one case *market* resolved to *customer* instead of *phone service*. About 40% of the problems in Experiment 2 fell in this category.

An extreme form of this situation are cases in which there is a word-form among the antecedents which is identical to the bridging description, and therefore is always chosen as antecedent, yet is not the

desired antecedent. We already mentioned one reason for that— lemmatization occasionally creates two identical word-forms, e.g., *plan* from *planned*. Another, more interesting reason is that sometimes the desired antecedent is described using a different word-form. This happens, for example, with inferential descriptions referring to names: e.g., one text about companies mentioned the word *company* quite often, and then it mentioned a specific company called *Pinkerton*. The following inferential description *the company* referred to this specific company, but the algorithm picked instead an antecedent explicitly introduced with the word *company* that had appeared in the preceding five sentences.

- Finally, there were cases in which the antecedent suggested by the algorithm did not wear any obvious semantic relation to the definite description. There were 40 such cases (20 % of the total) in the second experiment, 28 (13.72%) in the third one.

### Semantic Priming and Inferential Descriptions

Even though in both Experiments 2 and 3 we got much better results than chance, and even though the results could still be improved by about 14-15% with better clustering algorithms, the fact that in about 40% of the cases the correct antecedent is not the semantically closest one clearly indicates that what we called the Main Hypothesis is false: i.e., that semantic priming is not the only factor involved in the resolution of inferential descriptions.

The most obvious next hypothesis, especially at the light of previous work on definite descriptions, is that attentional mechanisms play a role—i.e., that a focusing mechanism such as those suggested by Grosz (1977) and Sidner (1979) restricts the range of potential antecedents. If this were the case, the ‘long-distance’ cases of inferential descriptions could then be taken as references to previous discourse foci put on the stack (we are thinking here of a model that puts back some ideas from Sidner’s dissertation in Grosz and Sidner’s (1986) model, such as the one discussed in (Poesio, Stevenson, & Hitzeman 1997)).

Identifying the ‘focus’ (or ‘foci’) and tracking focus shifts in a real text in a completely uncontroversial fashion is notoriously difficult, and it is certainly not a task that can be implemented at the moment; we did nevertheless attempt a preliminary verification of this new hypothesis by analyzing 4 of the 20 texts previously studied, identifying the available foci according to the proposal in (Poesio, Stevenson, & Hitzeman 1997), and trying to decide for each inferential description whether its resolution only depended on

lexical knowledge (i.e., the antecedent was clearly not a focus) or whether instead its antecedent was one of the current foci; we didn’t count unclear cases. Surprisingly enough, given all the possible complications just mentioned, the results were fairly clear: of the 44 inferential descriptions in these four texts that we could classify unambiguously, only 15 (about 33%) depended exclusively on lexical knowledge for their resolution; in 29 cases, keeping track of the focus was necessary.

This admittedly very preliminary study suggests that our algorithm in fact performed better than the 22.7% figure mentioned above would suggest. If only about 33% of inferential descriptions can be resolved solely on the ground of lexical knowledge and without keeping track of the current focus, then a fairer evaluation of the performance of our clustering algorithm is that it achieved about 66% of what we could expect it to achieve.

It should also be noted that this analysis indicates that completely ignoring commonsense knowledge during resolution, and just assigning the current focus as antecedent for an inferential description, would not work either: for one thing, about 33% of inferential descriptions do not relate to the current focus, but to some other discourse entity; and anyway when more than one focus is available, the choice among them goes down to lexical knowledge again.<sup>6</sup> In other words, both lexical information and information about the current focus really seem necessary.

### Comparison with WordNet

The results of the two main experiments are summarised in the following table, in which we distinguish between the total number of inferential descriptions being resolved and the performance over the specific relationships of synonymy, hypernymy/hyponymy and meronymy:

| Experiment | IDs        | <i>Syn/Hyp/Mer</i> |
|------------|------------|--------------------|
| 2          | 46 (22.7%) | 2/2/4 (22.2%)      |
| 3          | 46 (22.7%) | 4/2/2 (22.2%)      |

The techniques discussed in this paper resolved correctly a greater number of inferential descriptions (46, 22.7%) than were resolved just using WordNet in (Poesio, Vieira, & Teufel 1997) (34 cases, 16.7%). Worse results were obtained for those inferential descriptions whose relationship to the antecedent is based on synonymy, hypernymy/hyponymy or meronymy than had been obtained with WordNet (22.2% instead

<sup>6</sup>Not to mention that lexical knowledge plays a crucial in some of the best-known algorithms for determining the structure of a text, such as (Morris & Hirst 1991; Hearst 1997).

of 39.5%), but better results were obtained for all other cases (22.8% in the second experiment, as opposed to 9.6% with WordNet).

## Evaluation of the Methods Considered

Since the commonsense knowledge acquired by the methods discussed in this paper does seem to be crucial for resolving inferential descriptions, and the choice of parameters does seem to have an impact on the results,<sup>7</sup> we intend to continue our investigation of different ways of choosing the dimensions, other measures, and at combinations of measures, to see if we can improve the method's performance in this respect.

We expect that performance will be improved by using the same corpus for training and evaluation (already, we had to correct for differences between British and American lexicon). We are also considering whether more than one type of clustering algorithm may be needed. The particular way of computing similarity we have adopted here looks like a good method for acquiring synonymy relations and subtyping relations, i.e., the information used for resolving descriptions that co-refer with their antecedent without being same-head, such as those definites that are expressed via a synonymous or hyponymous predicate (as in *the home / the house*) or that refer to an event (as in *John killed Mary. THE MURDER took place ....*). However, words that are merely associated such as *door / house* do not necessarily always occur in the same contexts; in order to learn this sort of information it may be better to simply look at how often the words themselves occur in the same context, instead of looking at which other words they occur with. E.g., one could tell that 'door' and 'house' are related because they occur often together, especially if they occur together in certain constructions. Vector distance does not expose the desired similarity between door and house; we are investigating the possibility of adding further factors, such as a direct measure of association between the target words, in the decision process. Information from parsing could be useful in the same way.

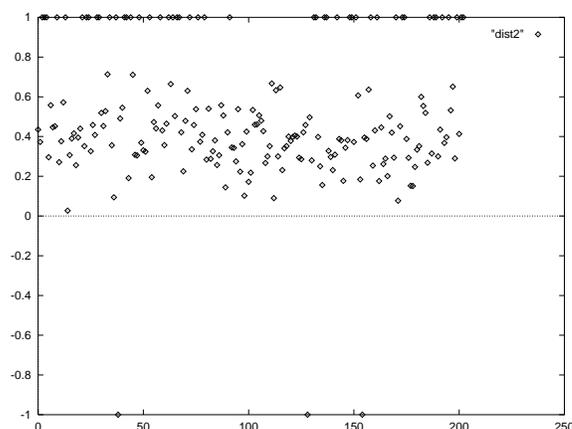
### A problem: improving precision

A rather significant problem of the technique adopted in this paper, as opposed to the methods used in (Poesio, Vieira, & Teufel 1997), is that the method just discussed always identifies an antecedent, even when it

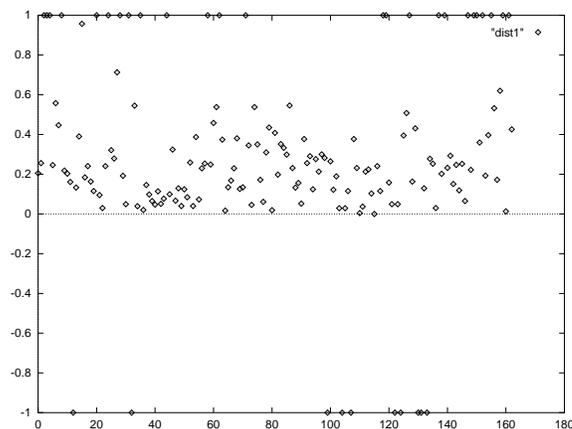
<sup>7</sup>See (Levy, Bullinaria, & Patel 1997) for a more thorough discussion of the impact of various parameter configurations on different tasks.

is not particularly close to the definite description, so the precision is very low.

As a first way around the problem, we considered whether it would be possible to reject candidates that did not make it above a predefined threshold. The answer, unfortunately, seems to be negative. The following figure shows the distances of the 203 chosen antecedents to the bridging descriptions. The distances vary from -1 to +1, but are concentrated in the area between 0.2 and 0.6:



The next figure shows the distances between the inferential descriptions and the 163 desired antecedents. The distances also vary from -1 to +1, but are concentrated in the area between 0 and 0.4:



As the figures indicate, there was no clear separation between the cases in which the resolution was right and wrong. The fact that the desired antecedents lie between 0 and 0.4 indicates again that in some cases the required antecedent was not the closest word. Once again, the addition of attentional factors and the explicit casting of the problem in terms of statistical pattern matching may turn it into one for which a suitable threshold may be identified.

We also considered whether the desired antecedent was generally closer to the inferential description

than the wrongly chosen one, in the sense of being mentioned more recently; again, the result was only partially satisfactory. 73 (61.9%) of the desired antecedents would have been mentioned more recently to the description than the actually chosen antecedent, but 45 (38.1%) would not.

## Acknowledgments

We wish to thank Will Lowe, Scott McDonald and Renata Vieira for much help with the algorithms and with testing the system. Massimo Poesio is supported by an EPSRC Advanced Fellowship.

## References

- Brown, P.; Della Pietra, V. J. D.; DeSouza, P. V.; Lai, J. C.; and Mercer, R. L. 1992. Class-based n-grams models of natural language. *Computational Linguistics* 18(4):467–479.
- Charniak, E. 1993. *Statistical Language Learning*. Cambridge, MA: The MIT Press.
- Clark, H. H., and Marshall, C. R. 1981. Definite reference and mutual knowledge. In Joshi, A.; Webber, B.; and Sag, I., eds., *Elements of Discourse Understanding*. New York: Cambridge University Press.
- Clark, H. H. 1977. Bridging. In Johnson-Laird, P. N., and Wason, P., eds., *Thinking: Readings in Cognitive Science*. London and New York: Cambridge University Press.
- Grosz, B. J., and Sidner, C. L. 1986. Attention, intention, and the structure of discourse. *Computational Linguistics* 12(3):175–204.
- Grosz, B. J. 1977. *The Representation and Use of Focus in Dialogue Understanding*. Ph.D. Dissertation, Stanford University.
- Hawkins, J. A. 1978. *Definiteness and Indefiniteness*. London: Croom Helm.
- Hearst, M. A. 1997. Texttiling: Segmenting text into multi-paragraph subtopic passages. *Computational Linguistics* 23(1):33–64.
- Huckle, C. 1996. *Unsupervised categorization of word meanings using statistics and neural network models*. Ph.D. Dissertation, University of Edinburgh.
- Levy, J.; Bullinaria, J.; and Patel, M. 1997. Using co-occurrence statistics for modelling language processes. In Crocker, M., ed., *Proc. of the AMLAP Conference*.
- Lund, K.; Burgess, C.; and Atchley, R. A. 1995. Semantic and associative priming in high-dimensional semantic space. In *Proc. of the 17th Annual Conference of the Cognitive Science Society*, 660–665.
- Meyer, D. E., and Schvaneveldt, R. 1971. Facilitation in recognizing pairs of words: evidence of a dependence between retrieval operations. *Journal of Experimental Psychology* 90:227–235.
- Miller, G. A.; Beckwith, R.; Fellbaum, C.; Gross, D.; and Miller, K. J. 1990. Introduction to WordNet: An on-line lexical database. *Journal of Lexicography* 3(4).
- Morris, J., and Hirst, G. 1991. Lexical cohesion computed by thesaural relations as an indicator of the structure of text. *Computational Linguistics* 17(1):21–48.
- Neely, J. H. 1991. Semantic priming effects in visual word recognition: A selective review of current findings and theories. In Besner, D., and Humphreys, G. W., eds., *Basic Processes in Reading: Visual Word Recognition*. Lawrence Erlbaum. 264–336.
- Poesio, M., and Vieira, R. 1998. A corpus-based investigation of definite description use. *Computational Linguistics*. To appear. Also available as Research Paper CCS-RP-71, Centre for Cognitive Science.
- Poesio, M.; Stevenson, R.; and Hitzeman, J. 1997. Global focus and pronominalization. In Jurafsky, D., and Regier, T., eds., *Proc. of First Conference on Computational Psycholinguistics*. Berkeley, CA: University of California.
- Poesio, M.; Vieira, R.; and Teufel, S. 1997. Resolving bridging references in unrestricted text. In Mitkov, R., ed., *Proc. of the ACL Workshop on Operational Factors in Robust Anaphora Resolution*, 1–6.
- Prince, E. F. 1981. Toward a taxonomy of given-new information. In Cole, P., ed., *Radical Pragmatics*. New York: Academic Press. 223–256.
- Schulte im Walde, S. 1997. Resolving bridging descriptions in high-dimensional space. Studienarbeit, Universities of Stuttgart and Edinburgh.
- Sidner, C. L. 1979. *Towards a computational theory of definite anaphora comprehension in English discourse*. Ph.D. Dissertation, MIT.
- Vieira, R., and Poesio, M. 1997. Processing definite descriptions in corpora. In Botley, S., and McEnery, T., eds., *Corpus-based and Computational Approaches to Discourse Anaphora*. UCL Press. To appear. Also available as HCRC Research Paper HCRC/RP-86, University of Edinburgh.
- Vieira, R., and Teufel, S. 1997. Towards resolution of bridging descriptions. In *Proc. of the 35th Joint Meeting of the Association for Computational Linguistics*.
- Vieira, R. 1998. *Definite Description Resolution in Unrestricted Texts*. Ph.D. Dissertation, University of Edinburgh, Centre for Cognitive Science.