

# Underspecification and Anaphora: Theoretical Issues and Preliminary Evidence

Massimo Poesio  
*University of Essex*

Patrick Sturt  
*University of Edinburgh*

Ron Artstein  
*University of Essex*

Ruth Filik  
*University of Glasgow*

Much experimental work in psycholinguistics suggests that fully specified syntactic and semantic interpretations are obtained incrementally. The finding that interpretation takes place incrementally is very robust and underlies our own view of sentence processing as well; however, most of this work tends to test very simple interpretive judgments, and using materials which have very clean-cut interpretations, which makes the view expressed above more questionable when applied to semantic interpretation. This article discusses a class of anaphoric expressions that do not appear to have a clear antecedent, referring to data from both corpus analysis and psycholinguistic experiments. We argue that these cases of anaphora are similar to cases of lexical polysemy, and propose an explicit semantic representation for such cases.

There is very strong evidence that most language interpretation processes take place quickly and incrementally (Marslen-Wilson, 1975; Swinney, 1979; Frazier, 1987; Tanenhaus, Spivey-Knowlton, Eberhard, & Sedivy, 1995). Most such papers also suggest that the interpretation thus derived is fully disambiguated; however, this further conclusion tends to be based on very simplified views about ambiguity and interpretation, which do not take into account the diversity of ambiguous expressions that are found in natural language. We will argue in this paper that with certain types of ambiguity, the ambiguous expressions may be left unresolved in the right context.

The starting point for the work presented in this paper is the study reported by Frazier and Rayner (1990), who demonstrated the need to distinguish between two types of lexical ambiguity: *homonymy* and *polysemy*. As we will discuss in more detail below, Frazier and Rayner showed that readers commit to preferred interpretations of homonymous words, like *bank*, which have multiple meanings, while they do not immediately commit to any specific interpretation of polysemous words, like *newspaper*, which have multiple related senses.

The research project discussed in this paper aims to extend such investigations

to another aspect of language use, namely anaphora. We show evidence suggesting that in the case of anaphoric interpretation, just as in the cases of lexical ambiguity studied by Frazier and Rayner (1990), a distinction can be drawn between two types of ‘ambiguity,’ and that in one of these cases at least, the final interpretation may not be fully specified, but only ‘good enough’ in the sense of Ferreira, Bailey, and Ferraro (2002). A second goal of our research project is to go further than Frazier and Rayner, or Ferreira and colleagues, in specifying what a ‘good enough interpretation’ may be; we make some preliminary theoretical suggestions in this respect. This paper builds on previous work (Poesio, Reyle, & Stevenson, 2001), but goes further in at least two directions: first of all, we undertook a large scale ambiguity annotation experiment, confirming that indeed these anaphoric expressions can be identified as ambiguous, and secondly, online studies were carried out.

The structure of the paper is as follows. We begin by discussing the just mentioned papers suggesting that initial interpretations may not be fully disambiguated (Frazier & Rayner, 1990; Ferreira et al., 2002). We also discuss our own preliminary work (Poesio et al., 2001) in which we used corpus analysis and corpus annotation to identify ‘interesting’ expressions from the point of view of semantic interpretation. In the following section we describe new evidence from corpora concerning the class of anaphoric expressions that we call *mereological references*, which appear to be ambiguous in a different way from other types of anaphoric expressions. We go on to show that these anaphoric expressions behave differently from other anaphoric expressions in both on-line and off-line measures of processing. Finally, we consider what the initial ‘good enough’ interpretation of these expressions may be.

## BACKGROUND

As we have mentioned above, much work on language processing is based on the implicit assumption that ambiguity is always fully resolved in natural language. However, the evidence underlying this assumption derives from work which involves very simple interpretive judgments, using materials which have very clean cut interpretations (Marslen-Wilson, 1975; Swinney, 1979; Frazier, 1987; Tanenhaus et al., 1995). The work reviewed in this section suggests, in contrast, that when more complicated judgments are tested, this view of interpretation may need to be revised.

### Frazier and Rayner: Polysemy vs Homonymy

Work on lexical access (Swinney, 1979; Simpson, 1994) tends to focus on *homonymy*—cases of lexical ambiguity like *pitcher* or *records*, in which the two interpretations of a word are clearly distinct (often also from an etymological perspective), and can be argued to be associated with distinct lexical entries. However, Frazier and Rayner (1990) found that the interpretation of homonymous words could be experimentally distinguished from that of *polysemous* words such

as *newspaper*, whose senses are closely interrelated (“a business firm that publishes newspapers” and “the physical object that is the product of a newspaper publisher”, according to WordNet 2.1). Whereas garden paths could be observed when the initially preferred interpretation of an homonym such as *record* was disconfirmed by subsequent context, as in (1d), no garden paths were observed for polysemous words like *newspaper* in (2d).

- (1)
  - a. After they were scratched, the records were carefully guarded.
  - b. After the political takeover, the records were carefully guarded.
  - c. The records were carefully guarded after they were scratched.
  - d. The records were carefully guarded after the political takeover.
- (2)
  - a. Lying in the rain, the newspaper was destroyed.
  - b. Managing advertising so poorly, the newspaper was destroyed.
  - c. Unfortunately the newspaper was destroyed, lying in the rain.
  - d. Unfortunately the newspaper was destroyed, managing advertising so poorly.

Although Frazier and Rayner do not explain what the initial interpretation of *newspaper* may be, their study clearly suggests that in cases where such a ‘preliminary’ interpretation may be available, the requirement that all phrases are immediately and fully interpreted may be weakened; Frazier and Rayner call this the *Immediate Partial Interpretation Hypothesis*. They also suggest that this initial partial interpretation in some sense ‘covers’ the disambiguated interpretations.

(As an aside it is worth noting that, as any lexicographer will point out, matters are much more complicated, and the distinction between homonymy and polysemy is sometimes very difficult to draw (Pinkal, 1995; Lyons, 1995). We take this point to provide further support for our main point, that the findings about incremental interpretation may be based on an overly simplified view of interpretation.)

## Good enough syntactic representations

Ferreira et al. (2002) discuss several experiments challenging “. . . the assumption that utterance interpretations are compositionally built up from words clustered into hierarchically organized constituents”. In one such experiment, Duffy, Henderson, and Morris (1989) measured naming times for the final word in sentences such as those in (3), finding that the naming times for the word *cocktails* were significantly slower in (3b) than in (3a) (in which *bartender* serves as a prime for *cocktails*). However, they also found no difference between the naming times for (3a) and (3c).

- (3)
  - a. The boy watched the *bartender* serve the cocktails.
  - b. The boy saw that the *person* liked the cocktails.
  - c. The boy who watched the *bartender* served the cocktails.

According to Ferreira et al., priming is usually taken to operate at the ‘conceptual’ level—a level at which argument structure has already been identified. However,

these results suggest that priming in fact either operates at a level in which not even syntactic structure is realized, or at the very least is not affected by syntactic structure (see Morris (1994), however, for contrasting results using eye fixation times as a measure of priming). Further evidence questioning the extent to which semantic interpretation is affected by syntactic structure was reported in a second series of experiments cited by Ferreira et al. In one study, Christianson, Hollingworth, and Ferreira (2001) were concerned with the question of whether after reanalysis people delete the initial incorrect interpretation. They tested this by presenting their subjects with garden path sentences such as (4):

(4) While Anna dressed the baby played in the crib.

and afterwards asking them questions such as *Did the baby play in the crib?* (which should be answered positively if *the baby* is successfully re-interpreted as the subject of *played*) and *Did Anna dress the baby?* (which, should be answered negatively if *dressed* is successfully re-interpreted as an intransitive verb). What Christianson et al. found is that both questions tended to be answered positively, compared with unambiguous controls. This shows that part of the final interpretation persisted from the initial misanalysis, and was not consistent with the globally correct analysis of the sentence. Therefore, either the participants did not always complete the syntactic reanalysis of the sentence, or they did not always complete the change of interpretation following syntactic reanalysis.

On the basis of these and other experiments, one could argue that semantic interpretation may not always be entirely parasitic on syntactic structure, or perhaps that subjects sometimes only derive incomplete syntactic representations, which are ‘good enough’ for their purposes (Ferreira et al.’s conclusion). No specific hypotheses concerning the form of these representations are made by Ferreira et al.

## Underspecification and anaphoric reference

Anaphoric reference is generally viewed as a case of *h-ambiguity* (Pinkal, 1995; Poesio, 1999): the type of ambiguity in which the alternative interpretations of an expression are semantically distinct, as in the case of homonymy. Accordingly, most evidence on anaphoric interpretation suggests that it is quick and incremental, both in the case of definite descriptions (Tanenhaus et al., 1995) and in the case of pronouns and reflexives (Arnold, Eisenband, Brown-Schmidt, & Trueswell, 2000; Sturt, 2003) as would be predicted by Frazier and Rayner’s Immediate Partial Interpretation Hypothesis, although there is some evidence for a delay in processing of pronouns (Corbett & Chang, 1983; Garrod, Freudenthal, & Boyle, 1994). However, ambiguous anaphoric references investigated in this project do not appear to suffer from some of the processing difficulty caused by ambiguity, and in this way, these expressions resemble the cases of lexical access studied by Frazier and Rayner. In addition, it is intuitively plausible that these expressions may only be assigned a ‘good enough’ interpretation, like the cases reported by Ferreira et al.

Poesio et al. (2001) argued that in dialogue corpora, a fairly powerful test can be used to find candidates for underspecification: look for anaphoric expressions that (i) don't seem to have a preferred antecedent, yet (ii) do not result in the listener signalling a misunderstanding. Poesio et al. carried out an analysis of the anaphoric expressions in the TRAINS corpus collected at the University of Rochester (Gross, Allen, & Traum, 1993), identifying several classes of expressions passing this test. We concentrate here on what Poesio et al. called the *mereological* cases. An example of this class is the pronoun *it* in utterance 5.1 in the fragment in (5).

- (5) 3.1 M: can we .. kindly hook up  
 3.2 : uh  
 3.3 : engine E2 to the boxcar at .. Elmira  
 4.1 S: ok  
 5.1 M: +and+ send it to Corning  
 5.2 : as soon as possible please  
 6.1 S: okay  
 [2sec]  
 7.1 M: do let me know when it gets there  
 8.1 S: okay it'll /  
 8.2 : it should get there at 2 AM  
 9.1 M: great  
 9.2 : uh can you give the  
 9.3 : manager at Corning instructions that  
 9.4 : as soon as it arrives  
 9.5 : it should be filled with oranges  
 10.1 S: okay  
 10.2 : then we can get that filled

In this example, it is not clear whether the pronoun *it* in 5.1 refers to *the engine E2* which has been hooked up to *the boxcar at Elmira*, or to the boxcar itself, or indeed whether that matters. It's only at utterance 9.5 that we get evidence that *it* probably referred to *the boxcar at Elmira*, since it is only boxcars that can be filled with oranges; yet, if anything, focusing theories would predict engine E2 to be the antecedent, since engine E2 is the direct object, the THEME, and comes first (Sidner, 1979; Stevenson, Crawley, & Kleinman, 1994; Grosz, Joshi, & Weinstein, 1995; Pearson, Stevenson, & Poesio, 2001). This didn't only happen with the verb *hook up*, but more generally whenever two objects were put together, e.g., by loading a cargo into some sort of transport vehicle, as in the following example:

- (6) 26.1 S: okay  
 27.1 M: so then we'll  
 27.2 : ... we'll be in a position to  
 27.3 : load the orange juice into the tanker car  
 27.4 : ... and send that off

There are important similarities and differences between these cases of anaphoric ambiguity and the cases studied by Frazier and Rayner. Just like the cases of lexical polysemy studied by Frazier and Rayner (and unlike the cases of lexical homonymy), mereological pronouns can receive an interpretation that ‘covers’ both possible antecedents: the ‘train’ formed by the engine and the boxcar in (5), the ‘tanker car full of orange juice’ in (6). As said above, this is not normally the case with anaphoric expressions, so we would expect to find a contrast between mereological pronouns and cases in which no such ‘merged’ interpretation has been created. However, there is also an important difference: whereas in the cases studied by Frazier and Rayner the existence of a superinterpretation was part of lexical knowledge, in the case of mereological interpretation the merged interpretation is the result of actions discussed in the text.

Mereological pronouns are an interesting case of expression for which only a ‘good enough’ interpretation may be reached, for in this case, we can make a fairly precise hypothesis concerning what this interpretation may be. From the point of view of the plan, once the engine and the boxcar are hooked up, or once the orange juice is loaded in the tanker car, the two objects will move together. So strictly speaking, all that the listener needs to do in this case is to somehow restrict the interpretation to the object obtained by joining together the two possible antecedents and its parts. The resulting reference will be ‘good enough’ for the listener’s purposes, as long as it specifies the engine, the boxcar, the object formed by combining the two, or any part of that object. This idea is made more precise in the formulation of the Justified Sloppines Hypothesis below.

## CORPUS EVIDENCE

Poesio et al. (2001) only carried out informal corpus studies of the mereological cases. In this section we describe a systematic analysis of the data in the corpus, which we achieved through an annotation experiment with many naïve subjects (Poesio & Artstein, 2005). In this study, our annotators could mark ambiguity by indicating multiple antecedents for an anaphoric expression, to reveal whether indeed subjects perceived the mereological cases as ambiguous. In addition to this *explicit* marking of ambiguity, we also assessed *implicit* ambiguity, where the pattern of annotation by multiple subjects revealed that an item was ambiguous even when the annotators were not aware of this themselves. As the corpus consists of samples of real dialogue, we feel that our method can be informative about realistic language use. Thus, such corpus studies can complement more controlled experimental studies, such as those of (Frazier & Rayner, 1990) or our own experimental studies described below, where experimental control inevitably involves a sacrifice of realistic contexts.

## Methods

*Coding Scheme.* The coding manual used in this experiment is based on the approach to anaphoric annotation developed in the projects MATE, which developed standards for Multilingual AnnoTation in dialoguE (Poesio, Bruneseaux, &

Romary, 1999), and GNOME, which carried out annotation of discourse and semantic information for the purposes of Generating Nominal Expressions (Poesio, 2004). The task and instructions were simplified, the primary simplification being that we did not annotate bridging references. On the other hand, we added instructions for marking multiple antecedents for ambiguous anaphoric expressions, and a simple way for marking discourse deixis.

*Materials.* The data used for this coding experiment come from the first edition of the TRAINS corpus collected at the University of Rochester (Gross et al., 1993). This corpus consists of transcripts of dialogues between two humans. One of the humans plays the ‘manager’ of a railway company, and has the aim of developing a plan to achieve a transportation goal (delivering a certain amount of goods at a given town by a given deadline). The other participant in the dialogue plays a ‘system’, whose role is to help managers develop this plan by providing them with required information such as timetables and equipment availability (the role of the system in all of the dialogues in the TRAINS 91 corpus was played by Derek Gross). Most dialogues in the corpus contain examples of mereological references; the text chosen for the experiment was dialogue 3.2 from the corpus. This dialogue contains 16 instances of *it* (counting all tokens, whether referential, expletives, or of unclear use) and 16 instances of *that*. Subjects were trained on dialogue 3.1.

*Annotation tools.* The subjects entered their annotations directly on a computer, using the MMAX2 annotation tool (Müller & Strube, 2003) (available from <http://mmax.eml-research.de/>). This tool has a graphical interface that makes it relatively easy both to do the basic annotation task of assigning labels to portions of text—the so-called *markables*—and to mark anaphoric antecedents. MMAX 2’s interface includes two main windows. One window shows the text, with the current markable highlighted in yellow. Using this window it is possible to identify the antecedent(s) of an anaphoric expression by creating one or more *pointers*; this is done by clicking on them with the mouse. A second window, the *markable browser*, makes it possible for the annotators to go to the next markable in the text by hitting the space bar.

*Subjects.* Eighteen paid subjects participated in the experiment, all students at the University of Essex (mostly undergraduates from the Departments of Psychology and Language and Linguistics). They were paid £30 for their participation.

*Instructions.* Subjects were given the map of the ‘TRAINS world’ used in the original dialogues and were instructed to go through the markables in order using MMAX 2’s markable browser, and to label each markable as belonging to one of the following four classes.

*phrase*: a markable referring to an object mentioned earlier in the dialogue;

*segment*: a markable referring to a plan, event, action, or fact discussed earlier in the dialogue;

place: one of the five railway stations in the ‘TRAINS world’ (Avon, Bath, Corning, Dansville, and Elmira) explicitly mentioned by name (this was done to save time, since references to places are never ambiguous in this domain);

none: a markable that does not fit any of the above criteria, for instance one referring to a novel object, or a non-referential noun phrase.

For markables designated as *phrase* or *segment*, subjects were instructed to create a pointer (in the sense discussed above) to the antecedent, which could be either a phrase markable or a dialogue turn. In case an expression was considered ambiguous, subjects were instructed to create more than one pointer. (Due to the limitations of the coding scheme they could only indicate ambiguity between two phrase antecedents or two discourse antecedents, not an ambiguity between a phrase antecedent and a discourse antecedent.)

## Results

The dialogue used in this study contains 8 instances of *it* or *that* (about 25% of the total number of tokens of these types) which, according to our own judgment, fit the ‘mereological reference’ pattern of Poesio et al. (2001). It also exhibits two other types of anaphoric ambiguity: ambiguity in discourse deictic references to actions / plans, and an ambiguity that coders were not able to mark explicitly in our scheme—the ambiguity between *discourse new* and *discourse old* interpretations for anaphoric expressions (Prince, 1992). (An expression is considered discourse new if it introduces a new object in the discourse; discourse old if it is a reference to an already introduced object.)

We found perfect agreement among the 18 annotators for 65 of the 148 markables in the dialogue (43.9%), and near perfect agreement (no more than two disagreeing coders) on another 14 markables (9.5%)—i.e., there was no real disagreement on 53.4% of markables. Of the remaining markables, we considered 67 (45.3%) implicitly ambiguous in that at least two distinct antecedents were marked by at least two coders each. (The remaining two markables were assigned a single label by 14 or 15 coders, and distinct labels by each of the remaining coders.) Of the 67 implicitly ambiguous markables, 24 (16.2% of the total number of markables) were marked as explicitly ambiguous by at least one annotator; overall, there were 131 explicit ambiguity markings by all coders. (In addition, 27 of these implicitly ambiguous markables (18.2% of all the markables) were marked as implicitly ambiguous between a discourse-old and a discourse-new interpretation—a type of ambiguity which, as we said above, our coders were not able to mark explicitly in the scheme used in this particular experiment.) These results are summarized in Table 1.

The results of this experiment clearly show that annotators are able to identify ambiguity and mark it: of the 18 annotators, only 3 did not mark any item as explicitly ambiguous. However, ambiguity is seldom noticed by many coders. Of the 24 items marked as ambiguous by at least two coders (38 by at least one coder), only 6 were marked as ambiguous by seven coders or more. (5 of those

Table 1  
*Coder judgments concerning ambiguity in TRAINS 91 Dialogue 3.2*

	First Half		Second half		Total	
	<i>n</i>	%	<i>n</i>	%	<i>n</i>	%
Total	72	100.0	76	100.0	148	100.0
Perfect agreement	27	37.5	38	50.0	65	43.9
Almost perfect	10	13.9	4	5.3	14	9.5
Ambiguous (total)	35	48.6	32	42.1	67	45.3
Explicit ambiguity	18	25.0	6	7.9	24	16.2
Old/new ambiguity	8	11.1	19	25.0	27	18.2

were cases of mereological reference.) The distribution of items according to the number of coders that marked them as ambiguous is shown in Figure 1.

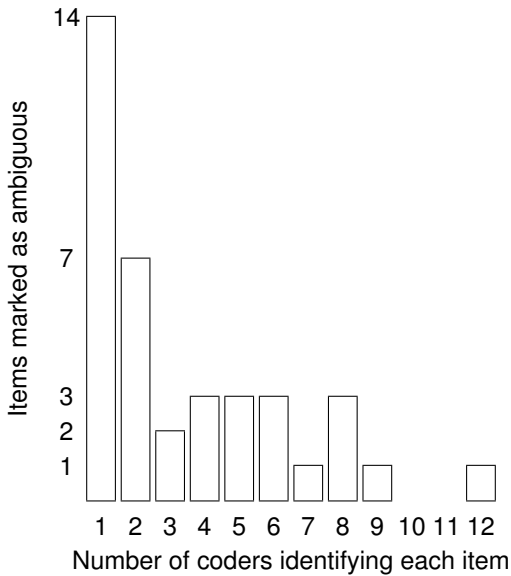


Figure 1. Distribution of ambiguity markings

The 8 instances of *it* or *that* that we ourselves classified as ‘mereological references’ were all classified as explicitly ambiguous by at least 2 coders. We can see an example of an explicitly marked ambiguous ‘mereological’ pronoun in utterance unit 3.1 of the following fragment, where the two *it* pronouns display the same type of ambiguity already seen in (5): they could refer to engine E2, the boxcar, or both.

- (7) 1.4 M: first thing I'd like you to do  
 1.5 is send engine E2 off with a boxcar to Corning  
 to pick up oranges  
 1.6 uh as soon as possible  
 2.1 S: okay  
 [6 sec]  
 3.1 M: and while it's there it should pick up the tanker

All of our subjects marked the two *it* pronouns in 3.1 as 'phrase' references in the sense introduced above (i.e., as referring to objects previously introduced using a nominal expression). The first pronoun was marked by 9 coders as explicitly ambiguous between engine E2 and the boxcar, by 6 coders as unambiguous and referring to engine E2, and by 3 coders as unambiguous and referring to the boxcar. The second pronoun was marked by 8 coders as ambiguous between engine E2 and the boxcar—2 did so directly, and 6 marked the previous, ambiguous pronoun as an antecedent; the remaining coders marked it as unambiguous—8 as referring to engine E2, and 2 as referring to the boxcar. In other words, even without explicit marking of ambiguity, we would still be able to infer, implicitly, that these pronouns are ambiguous, by the fact that different interpretations are chosen by at least two coders.

The dialogue used in this study also contained two examples of anaphoric references to a mereological structure which was not the result of an action performed in the plan, but was pre-existing. These two examples are shown in the following fragment.

- (8) 18.6 S: it turns out that the boxcar at Elmira  
 18.7 has a bad wheel  
 18.8 and they're .. gonna start fixing that at midnight  
 18.9 but it won't be ready until 8  
 19.1 M: oh what a pain in the butt

The fact that the boxcar under discussion has a wheel is not the result of an action performed during the plan, yet just as in the mereological cases, it is not clear to us whether the pronouns *that* in utterance unit 18.8 and *it* in utterance unit 18.9 refer to the boxcar or to its bad wheel. However, these cases were interpreted by our annotators in a different way than the cases discussed earlier. None of our annotators identified the pronoun *that* as ambiguous: 15 marked it as referring to the wheel, and one as referring to the boxcar (two others marked it as discourse deictic). The pronoun *it* was marked as ambiguous by a single annotator; 11 marked it as referring to the boxcar alone, and 6 to the wheel alone.

In summary, this study shows that mereological references are fairly common, at least in this type of dialogue, and that in each such case, the alternative interpretations identified by Poesio et al. are indeed available: in fact, our coders did not display any particular preference among the interpretations annotated. We also found that even with limited training it is possible to get coders to explicitly identify many of these ambiguous cases, which suggests that a systematic study of this phenomenon is possible. What the study does *not* tell us, however, is whether

these mereological cases are somehow problematic during interpretation, and how they are interpreted; this is the task of off-line and on-line behavioral studies, to which we now turn.

## PSYCHOLOGICAL EVIDENCE

### Offline evidence

If mereological references genuinely behaved from an interpretation perspective like the cases of lexical polysemy studied by Frazier and Rayner, we would expect to observe processing differences between such cases and other cases of anaphoric reference. Poesio et al. (2001) tested whether sentences that contain a potentially ambiguous anaphoric reference are easier to process when the two potential antecedents are part of a single mereological structure which makes the possible interpretations equivalent for the purposes of the plan, compared to a condition where the two potential antecedents are not joined. Poesio et al. tested this hypothesis using the offline Magnitude Estimation technique proposed in Bard, Robertson, and Sorace (1996). They asked participants to judge whether ‘MERE-**OLOGY**’ sentences such as (9a), which depicts *the engine* and *the boxcar* being attached together, are ‘more acceptable’ (in that less ambiguous) than ‘NON-MERE**OLOGY**’ sentences like (9b), in which *the engine* and *the boxcar* are not attached together:

- (9) a. The engineer hooked up the engine to the boxcar and sent it to London.
- b. The engineer separated the engine from the boxcar and sent it to London.

The experiment was run using WebExp, a software package for running experiments on the Web developed at the Universities of Edinburgh and Saarbruecken (<http://www.webexp.info/>). A significant effect of MERE**OLOGY** was found, such that mereological sentences like (9a) were judged to be reliably more acceptable than non-mereology sentences like (9b). The participants were also asked to judge the acceptability of the initial parts of the sentences, to ensure that any differences in acceptability were not due to uncontrolled features of the sentences before the anaphoric section:

- (10) a. The engineer hooked up the engine to the boxcar.
- b. The engineer separated the engine from the boxcar.

It was found that the shorter control sentences like (10a) and (10b) did not differ in acceptability, leading to an interaction between mereology and sentence length. Therefore Poesio et al. (2001) concluded that the differences between (9a) and (9b) were indeed due to the anaphoric content of the sentences, and they interpreted the difference in terms of a perceived difficulty associated with the referential ambiguity in the non-mereological case (9b). The lack of such perceived difficulty in (9a) was interpreted as evidence that the initial interpretation was underspecified.

This off-line acceptability experiment left at least two questions open, however: how do the patterns of acceptability map on to on-line processing, and what was the interpretation assigned to the anaphoric expression?

## Online evidence

Filik, Sanford, and Sturt (2005) conducted an eye-tracking study to examine the on-line processing of anaphoric references to mereological structures. Experimental items were as in (11):

- (11) a. **Mereology/singular**  
 There were many delays. / The railwayman hooked up / the engine and the boxcar, / and sent it / quickly/ to the central station.
- b. **Mereology/plural**  
 There were many delays. / The railwayman hooked up / the engine and the boxcar, / and sent them / quickly/ to the central station.
- c. **Neutral/singular**  
 There were many delays. / The railwayman saw / the engine and the boxcar, / and sent it / quickly/ to the central station.
- d. **Neutral/plural**  
 There were many delays. / The railwayman saw / the engine and the boxcar, / and sent them / quickly/ to the central station.

The experimental design involved a manipulation of whether the main predicate was mereology-constructing (*hooked up*) or neutral (*saw*). The experiment also manipulated whether the subsequent pronoun was singular *it* or plural *them*. Previous work has found that a singular pronoun causes difficulty when it is forced to refer to a single conjunct in a coordinated noun phrase (the so-called “conjunction cost”) while a plural pronoun is relatively easy in such cases (Albrecht & Clifton, 1998; Moxey, Sanford, Sturt, & Morrow, 2004). Therefore more processing difficulty is expected following the pronoun in (11c) than in (11d), as these are the two neutral conditions, involving non-mereological predicates. From the point of view of underspecification, the interesting contrast is between (11a) and (11b), which involve mereological predicates. Here, the prediction is that the mereological structure makes alternative interpretations for a singular pronoun available, for the reasons discussed above. This should eliminate the conjunction cost, making (11a) relatively easy to process compared with (11c).

Eye-movement data from the region containing the pronoun (*and sent it*) showed that singular reference was indeed easier for mereology-constructing than neutral sentences: there were more initial regressions from this region in the singular neutral condition than the singular mereology condition. A similar pattern of effects was found in the final region of the sentence (*to the central station*) in regression path time (the time taken to “go past” the region). It was also the case that plural reference was equally easy to process for mereology constructing and neutral sentences. One possible reason for this is that the use of the conjoined noun phrase *the engine and the boxcar* may automatically make a plural reference object available (Kamp & Reyle, 1993; Moxey et al., 2004), facilitating

subsequent reference through a plural pronoun, whether a mereology constructing predicate is used or not.

Filik et al. further investigated this issue using the *text change detection* paradigm (Sanford & Sturt, 2002). In this procedure, participants read a piece of text at their own pace. This piece of text is then shown for a second time, with a possible change to one of the words. The task for the participant is to signal whether they noticed this change, and if so, which word it is that has changed. The likelihood of noticing a change may reflect the level of specification applied to the readers' mental representation of this word. If fewer changes are detected this may indicate that readers have underspecified their representations in some way. In the current change detection study, participants read sentences based on those used in the eye-tracking study described above. In these sentences, the plural pronoun 'them' appeared in the first display of the text, and was changed to the singular pronoun 'it' for the second presentation. Results showed that readers noticed fewer changes from 'them' to 'it' in mereology-constructing than neutral sentences, suggesting that readers may have underspecified the referent of the pronoun in mereology-constructing cases, supporting the results of the eye-tracking study.

Thus, the off-line and on-line experimental results provide strong support for the effect of mereological structures on anaphoric interpretation.

## THE JUSTIFIED SLOPPINESS HYPOTHESIS

Up until now, we have discussed evidence suggesting that cases of anaphoric reference to structured entities such as those discussed above do indeed behave differently from other cases of anaphoric ambiguity. We have hinted that this may be because, as in the cases of lexical polysemy, an interpretation that supersedes all alternatives appears to be available, and that assigning this interpretation to the pronoun may be 'good enough' for the purposes of the plan. In this section we will attempt to characterize in a more systematic fashion the possible interpretations of these mereological pronouns, borrowing some notation from the theory of plurals and parts proposed by Link (1983).

We will write, e.g.,  $oj \oplus tc$  to indicate the object that has *oj* (orange juice) and *tc* (tanker car) as subparts, and  $a \triangleleft b$  to say that *a* is a mereological part of *b*. With this notation, we can formalize the first and most obvious property of examples (5) and (6): namely, that actions like *hooking up* and *loading* are performed that create a new object  $a \oplus b$  out of the potential antecedents *a* and *b* (e.g.,  $oj \oplus tc$  in (6)).

The second property of these examples is that four interpretations for the pronominal expression are possible. The complete list of the possible interpretations of *that* in (6), 27.4 is as follows:

1. the orange juice, *oj*;
2. the tanker car, *tc*;
3. the composite object formed by loading the orange juice in the tanker car,  $oj \oplus tc$ ;

4. in addition to the three fully specified interpretations above, the interpretation might also have an underspecified interpretation—an indeterminate  $x \triangleleft (oj \oplus tc)$

This latter interpretation ( $x \triangleleft (oj \oplus tc)$ ) is what has been called a *p-underspecified interpretation* in Poesio (1999)—i.e., a ‘disjunctive’ interpretation that ‘covers’ all of the alternative interpretations, similar to those proposed for certain cases of lexical polysemy in Copestake and Briscoe (1995). We will hypothesize below that the existence of such an underspecified interpretation may be a further important property of these contexts. Either this interpretation or interpretation 3, the composite object, may be what Frazier and Rayner had in mind when they talked of an ‘interpretation that covers all alternatives’.

The third property that these examples have in common is that both in situations involving attaching two objects together and in situations involving loading objects into other objects, all of the alternative interpretations of the anaphoric expression are equivalent as far as the plan of moving these objects to a new location is concerned: after the two explicitly mentioned potential antecedents are joined, if one of them gets moved, the other one must be moved as well. (Note that these two interpretations might not be equivalent for the purposes of a different plan—e.g., lifting the two objects with a crane.) E.g., in (6), 27.4, all interpretations of the instruction *send that off* will achieve the same result irrespective of how the pronoun is interpreted. We will write  $X \sim Y$  to indicate that interpretation X is equivalent to interpretation Y for the purpose of the plan: i.e., we write

$$oj \sim tc$$

to indicate that from the point of view of the plan, interpreting the pronoun *that* as referring to the orange juice or the tanker car are equivalent. Similarly, in the case of (5), we will write  $e \sim b$  to say that from the point of view of the plan, the interpretation of *it* in which it refers to engine E2 and that in which it refers to the boxcar are equivalent. In fact, in this case, *all* four situations discussed above are equivalent in the sense just discussed: they are ‘good enough’ in the sense of Ferreira et al.

To summarize, the mereological references in the TRAINS dialogues have at least three aspects in common:

1. Both explicitly mentioned potential antecedents  $x$  and  $y$  are elements of an underlying mereological structure with summum  $\sigma = x \oplus y$  which has been explicitly constructed (and made salient) in the dialogue ( $\sigma = oj \oplus tc$  in (6));
2. The existence of this structure makes it possible to construct a *p-underspecified interpretation* in which the anaphoric expression is interpreted as denoting an element  $z$  included in the mereological structure—i.e., part-of its summum  $\sigma$ . In the notation used in Discourse Representation Theory (Kamp & Reyle, 1993), the presence in the discourse model of this *p-underspecified interpretation*

is indicated as follows.

$x$	$y$	$\sigma$	$z$
...			
		$\sigma = x \oplus y$	
		$z \triangleleft^* \sigma$	
		...	

3. All possible interpretations  $(x, y, z, x \oplus y)$  are equivalent for the purposes of the plan.

The evidence discussed in the previous sections may now be viewed as providing support for the following generalization:

**Justified Sloppiness Hypothesis** Ambiguous anaphoric expressions are not perceived as infelicitous provided that Conditions 1–3 hold.

This may be because if these three conditions hold, the speaker’s sloppiness in using an anaphoric expression in an ambiguous context is not problematic; we will therefore use the term *justified sloppiness* to indicate cases such as those discussed in the previous section, and refer to the hypothesis above as the *Justified Sloppiness Hypothesis*.

Of course, the fact that a p-underspecified interpretation exists does not mean that the listener will adopt it as its final interpretation; however, this possibility is what makes these examples interesting from an underspecification perspective. A questionnaire was run by Filik et al., which showed some evidence of a preference for either the SUM interpretation ( $engine \oplus boxcar$ ) or the underspecified representation  $z \triangleleft^* engine \oplus boxcar$  (where, in this context,  $z$  can be instantiated either as the engine or the boxcar, or the two combined).

The questionnaire showed that people tend to prefer the SUM interpretation of *it* in the context of (11a), and that they only rarely chose *the engine* or *the boxcar* as the referent on its own. However, this asymmetry may be due to the use of the conjunction; as mentioned above, we already know that single conjuncts of coordinated noun phrases are relatively inaccessible to subsequent singular reference. If this preference holds also in the context of mereological predicates, it is possible that the SUM reference is the only element of the underspecified interpretation that remains salient in the context of conjunction.

## DISCUSSION

The evidence presented in this paper suggests that linguists (computational and theoretical) and psychologists should pay attention to cases in which the interpretation of an ambiguous expression is not clear-cut. Such cases are fairly numerous, and subjects do find them ambiguous, at least implicitly. This suggests the need for relaxing the assumption that natural language expressions have a single, well-specified interpretation; some story also has to be said about how humans can nevertheless understand each other. The corpus methodology we are proposing, although not yet fully developed, appears nonetheless to be usable to identify at

least some of these cases, thus preparing the way for more controlled psychological experiments.

Secondly, our findings, particularly from the off-line and on-line behavioral studies, extend Frazier and Rayner's argument about lexical ambiguity to anaphoric ambiguity. We take the main point to be that not all types of ambiguity are alike, and the existence of a 'super-interpretation' appears to make certain types of pronominal ambiguity more acceptable, just as it did for certain types of lexical ambiguity in the Frazier and Rayner study. An interesting difference in the cases of anaphoric reference we are studying is that this 'super-interpretation' is made available through context instead of lexically, as in the cases studied by Frazier and Rayner.

One of the goals of this project is to make more specific hypotheses about what these 'super interpretations' or 'good enough interpretations' may be. We sketched out in some detail what we think the alternatives might be, and discussed several bits of evidence that we feel will shed some light on this point, even if they are not conclusive. Some evidence that either the SUM interpretation (*engine*⊕*boxcar*) or the underspecified representation  $z \leftarrow^* \textit{engine} \oplus \textit{boxcar}$  (where, in this context,  $z$  can be instantiated either as the engine or the boxcar, or the two combined) is provided by the offline questionnaire run by Filik et al. (2005).

Further research will examine in more detail the effect of mereological predicates in sentences with and without coordination. However, finding ways of discriminating between the 'underspecified' and the 'sum' reading may well prove difficult.

Given the available evidence, the Justified Sloppiness Hypothesis is the strongest generalization that we may draw at the moment. The Justified Sloppiness Hypothesis is weaker than Frazier and Rayner's Immediate Partial Interpretation Hypothesis (which is specifically an hypothesis about interpretation), but may still be applicable to cases other than mereological references; for example, to cases of reference to events and other abstract objects (see Poesio et al. (2001)).

## ACKNOWLEDGMENTS

This work was in part supported by EPSRC grant GR/S76434/01, ARRAU; the pilot work was carried out using the Royal Society grant *Cases of Unresolved Underspecification*, joint with Uwe Reyle. We wish to thank our earlier collaborators Uwe Reyle and Rosemary Stevenson, as well as Michael Strube and Christoph Müller for the use of MMAX and much help with it.

## REFERENCES

- Albrecht, J. E., & Clifton, C., Jr. (1998). Accessing singular antecedents in conjoined phrases. *Memory and Cognition*, 26, 599–610.
- Arnold, J. E., Eisenband, J. G., Brown-Schmidt, S., & Trueswell, J. C. (2000). The rapid use of gender information: Evidence of the time course of pronoun resolution from eyetracking. *Cognition*, 76, B13–B26.
- Bard, E. G., Robertson, D., & Sorace, A. (1996). Magnitude estimation of linguistic acceptability. *Language*, 72, 32–68.

- Christianson, K., Hollingworth, A., & Ferreira, F. (2001). Thematic roles assigned along the garden path linger. *Cognitive Psychology*, *42*, 368–407.
- Copetake, A., & Briscoe, T. (1995). Semi-productive polysemy and sense extension. *Journal of Semantics*, *12*, 15–68. (Special Issue on Lexical Semantics)
- Corbett, A., & Chang, F. (1983). Pronoun disambiguating: Accessing potential antecedents. *Memory and Cognition*, *11*, 283–294.
- Duffy, S. A., Henderson, J. M., & Morris, R. K. (1989). Semantic facilitation of lexical access during sentence processing. *Journal of Experimental Psychology: Learning, Memory and Cognition*, *15*, 791–801.
- Ferreira, F., Bailey, K. G. D., & Ferraro, V. (2002). Good-enough representations in language comprehension. *Current Directions in Psychological Science*, *11*, 11–15.
- Filik, R., Sanford, A., & Sturt, P. (2005). *Anaphoric reference to structured entities*. Manuscript in Preparation.
- Frazier, L. (1987). Sentence processing: a tutorial review. In M. Coltheart (Ed.), *Attention and performance XII: The psychology of reading* (pp. 559–586). Hove: Erlbaum.
- Frazier, L., & Rayner, K. (1990). Taking on semantic commitments: Processing multiple meanings vs. multiple senses. *Journal of Memory and Language*, *29*, 181–200.
- Garrod, S. C., Freudenthal, D., & Boyle, E. (1994). The role of different types of anaphor in the on-line resolution of sentences in a discourse. *Journal of Memory and Language*, *33*, 39–68.
- Gross, D., Allen, J., & Traum, D. (1993, June). *The TRAINS 91 dialogues* (TRAINS Technical Note No. 92-1). Computer Science Dept. University of Rochester.
- Grosz, B. J., Joshi, A. K., & Weinstein, S. (1995). Centering: A framework for modeling the local coherence of discourse. *Computational Linguistics*, *21*, 202–225. ((The paper originally appeared as an unpublished manuscript in 1986.))
- Kamp, H., & Reyle, U. (1993). *From discourse to logic: Introduction to modeltheoretic semantics of natural language, formal logic and discourse representation theory*. Dordrecht: Kluwer.
- Link, G. (1983). The logical analysis of plurals and mass terms: A lattice-theoretical approach. In R. Bäuerle, C. Schwarze, & A. von Stechow (Eds.), *Meaning, use and interpretation of language* (pp. 302–323). Walter de Gruyter.
- Lyons, J. (1995). *Linguistic semantics*. Cambridge: Cambridge University Press.
- Marslen-Wilson, W. (1975). Sentence perception as an interactive parallel process. *Science*, *189*, 226–228.
- Morris, R. K. (1994). Lexical and message-level sentence context effects on fixation times in reading. *Journal of Experimental Psychology: Learning, Memory and Cognition*, *20*, 92–103.
- Moxey, L. M., Sanford, A. J., Sturt, P., & Morrow, L. I. (2004). Constraints on the formation of plural reference objects: The influence of role, conjunction, and type of description. *Journal of Memory and Language*, *51*, 346–364.
- Müller, C., & Strube, M. (2003). Multi-level annotation in MMAX. In *Proc. of the 4th sigdial* (pp. 198–207).
- Pearson, J., Stevenson, R., & Poesio, M. (2001). The effects of animacy, thematic role, and surface position on the focusing of entities in discourse. In M. Poesio (Ed.), *Proc. of the first workshop on cognitively plausible models of semantic processing (sempro)*.
- Pinkal, M. (1995). *Logic and lexicon*. Dordrecht: D. Reidel.

- Poesio, M. (1999). *Utterance processing and semantic underspecification* (HCRC/RP No. 103). University of Edinburgh, HCRC.
- Poesio, M. (2004, May). The MATE/GNOME scheme for anaphoric annotation, revisited. In *Proc. of sigdial*. Boston.
- Poesio, M., & Artstein, R. (2005, July). Annotating (anaphoric) ambiguity. In *Proc. of the corpus linguistics conference*. Birmingham.
- Poesio, M., Bruneseaux, F., & Romary, L. (1999). The MATE meta-scheme for coreference in dialogues in multiple languages. In M. Walker (Ed.), *Proc. of the ACL workshop on standards and tools for discourse tagging* (pp. 65–74).
- Poesio, M., Reyle, U., & Stevenson, R. (2001). Justified sloppiness in anaphoric reference. In H. Bunt & R. Muskens (Eds.), *Computing meaning 3*. Kluwer. (To appear)
- Prince, E. F. (1992). The ZPG letter: subjects, definiteness, and information status. In S. Thompson & W. Mann (Eds.), *Discourse description: diverse analyses of a fund-raising text* (pp. 295–325). John Benjamins.
- Sanford, A. J., & Sturt, P. (2002). Depth of processing in language comprehension: not noticing the evidence. *Trends in Cognitive Science*, 6, 382–386.
- Sidner, C. L. (1979). *Towards a computational theory of definite anaphora comprehension in english discourse*. Unpublished doctoral dissertation, MIT.
- Simpson, G. B. (1994). Context and the processing of ambiguous words. In M. A. Gernsbacher (Ed.), *Handbook of psycholinguistics* (pp. 359–374). Academic Press.
- Stevenson, R. J., Crawley, R. A., & Kleinman, D. (1994). Thematic roles, focus, and the representation of events. *Language and Cognitive Processes*, 9, 519–548.
- Sturt, P. (2003). The time-course of the application of binding constraints in reference resolution. *Journal of Memory and Language*, 48, 542–562.
- Swinney, D. A. (1979). Lexical access during sentence comprehension: (re)consideration of context effects. *Journal of Verbal Learning and Verbal Behavior*, 18, 545–567.
- Tanenhaus, M. K., Spivey-Knowlton, M., Eberhard, K. M., & Sedivy, J. C. (1995). Integration of visual and linguistic information in spoken language comprehension. *Science*, 268, 1632–1634.