

THE RESOLUTION OF QUANTIFICATIONAL AMBIGUITY IN THE TENDUM SYSTEM

Harry Bunt
Computational Linguistics Research Unit
Dept. of Language and Literature, Tilburg University
P.O.Box 90153, 5000 LE Tilburg
The Netherlands

ABSTRACT

A method is described for handling the ambiguity and vagueness that is often found in quantifications - the semantically complex relations between nominal and verbal constituents. In natural language certain aspects of quantification are often left open; it is argued that the analysis of quantification in a model-theoretic framework should use semantic representations in which this may also be done. This paper shows a form for such a representation and how "ambiguous" representations are used in an elegant and efficient procedure for semantic analysis, incorporated in the TENDUM dialogue system.

The quantification ambiguity explosion problem

Quantification is a complex phenomenon that occurs whenever a nominal and a verbal constituent are combined in such a way that the denotation of the verbal constituent is predicated of arguments supplied by the (denotation of the) nominal constituent. This gives rise to a number of questions such as (1) What objects serve as predicate arguments? (2) Of how many objects is the predicate true? (3) How many objects are considered as potential arguments of the predicate?

When we consider these questions for a sentence with a few noun phrases, we readily see that the sentence has a multitude of possible interpretations. Even a sentence with only one NP such as

- (1) Five boats were lifted

has a variety of possible readings, depending on whether the boats were lifted individually, collectively, or in groups of five, and on whether the total number of boats involved is exactly five or at least five. For a sentence with two numerically quantified NPs, such as 'Three Russians visited five Frenchmen', Partee (1975) distinguished 8 readings depending on whether the Russians and the Frenchmen visited each other individually or collectively and on the relative scopes of the quantifiers. Partee's analysis is in fact still rather crude; a somewhat more refined analysis, which distinguishes group readings and readings with equally wide

scope of the quantifiers, leads to 30 interpretations (Bunt, in press).

This presents a problem for any attempt at a precise and systematic description of semantic structures in natural language. On the one hand an articulate analysis of quantification is needed for obtaining the desired interpretations of every sentence, while on the other hand we do not want to end up with dozens of interpretations for every sentence.

To some extent this "ambiguity explosion problem" is an artefact of the usual method of formal semantic analysis. In this method sentences are translated into formulae of a logical language, the truth conditions of which are determined by model-theoretic interpretation rules. Now one might want to consider a sentence like (1) not as ambiguous, but only as saying that five boats were lifted, without specifying how they were lifted. But translation of the sentence into a logical representation forces one to be specific. That is, the logical representation language requires distinction between such interpretations as represented by (2) (individual reading) and (3) (group reading):

$$(2) \#(\{x \in \text{BOATS} : \text{LIFTED}(x)\}) = 5$$

$$(3) \exists x \in \{y \subseteq \text{BOATS} : \#(y) = 5\} : \text{LIFTED}(x)$$

In other words, the analysis framework forces us to make distinctions which we might not always want to make.

To tackle this problem, I have devised a method of representing quantified expressions in a logical language with the possibility of leaving certain quantification aspects open. This method has been implemented in the TENDUM dialogue system, developed jointly at the Institute for Perception Research in Eindhoven and the Computational Linguistics Research Unit at Tilburg University, Department of Linguistics (Bunt, 1982; 1983; Bunt & thoe Schwartzberg, 1982;). This method is not only of theoretical interest, but also provides a computationally efficient treatment of quantification.

Ambiguity resolution

In a semantic analysis system which translates natural language expressions into formal representations, all disambiguation takes place during this translation.

This applies both to purely lexical ambiguities and to structural ambiguities. For lexical disambiguation this means that a lexical item has several translations in the representation language (RL), which are all produced by a dictionary lookup at the beginning of the analysis. The generation of semantic representations for sentences that display both lexical and structural ambiguity thus takes place as depicted in Fig. 1:

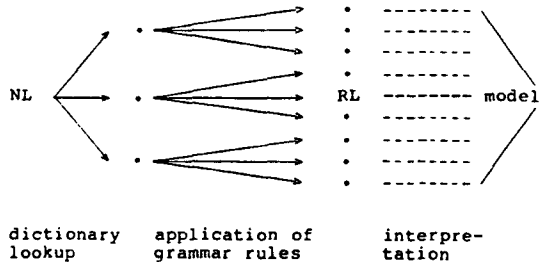


Fig. 1 Longer arrows indicate larger amount of processing.

Since the lexical ambiguities considered here are purely semantic, the same grammar rules will be applicable to all the lexical interpretations (assuming that the grammar does not contain world knowledge to filter out those interpretations that are meaningless in the discourse domain under consideration). Since the amount of processing involved in the application of grammar rules is very large compared to that of translating a lexical item to its RL instances, this set-up is not very efficient. In the PHLIQAL question-answering system (Bronnenberg et al., 1980) the syntactic/semantic and lexical processing stages were therefore reversed, so that disambiguation takes place as depicted in Fig. 2:

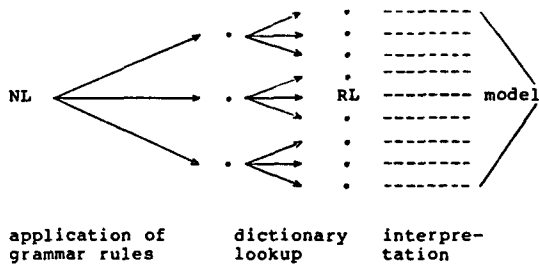


Fig. 2 Longer arrows indicate larger amount of processing.

In this setup an intermediate representation language is used which is identical to RL except that it has an ambiguous constant for every content word of the natural language.

It turns out that semantic analysis along these lines can be formulated entirely in terms of the traditional model-theoretic framework (Bunt, in press), therefore this method is appropriately called two-level model-theoretic semantics. This method has been implemented in the TENDUM system, with an intermediate representation language that

contains ambiguous constants corresponding to quantification aspects, in addition to ambiguous constants corresponding to nouns, verbs, etc.

Quantification aspects

The different aspects of quantification are closely related to the semantic functions of determiners. These functions depend on their syntactic position in a determiner sequence. A full-fledged basic noun phrase has the layout:

- (4) pre- + central + post- + head
determiner determiner determiner noun

(see Quirk et al., 1972, p.146). For example, in the NP

- (5) All my four children

the central determiner 'my' restricts the range of reference of the head noun 'children' to the set of my children; the predeterminer 'all' indicates that a predicate, combined with the noun phrase to form a proposition, is associated with all the members of that set, and the postdeterminer 'four' expresses the presupposition that the set consists of four elements. This set is determined by the central determiner plus the denotation of the head noun; I will call it the source of the quantification. In the case of an NP without central determiner the source is the denotation of the head noun. For the indication of the quantity or fraction of that part of the source that is involved in a predication I will use the term source involvement.

Quantification owes its name to the fact that source involvement is often made explicit by means of quantitative (pre-)determiners like 'five', 'many', 'all', or 'two liters of'. Obviously, source involvement is a central aspect of quantification.

Another important aspect of quantification is illustrated by the following sentences:

- (6a) The chairs were lifted by all the boys
(6b) The chairs were lifted by each of the boys

These sentences differ in that (6b) says unambiguously that every one of the boys lifted the chairs, whereas (6a) is unspecific as to what each individual boy did: it only says that the chairs were lifted and that all the boys were involved in the lifting, but it does not specify, for instance, whether every one of the boys lifted the chairs or all the boys together lifted the chairs. The quantifiers 'all' and 'each (of)' thus both indicate complete involvement of the source, but differ in their determination of how a predicate ('lifted the chairs') is applied to the source. 'Each' indicates that the predicate is applied to the individual members of the source; 'all' leaves open whether the predicate is applied to individual members, to groups of members, or to the sources as a whole. To designate the way in which a predicate is applied to, or "distributed over", the source of a quantification, I use the term distribution. A way of expressing the distribution of a quantification is by specifying the class of objects that the predicate is applied to, and how this class is related to the source. In the distributive case this class is precisely the :

source; in the collective case it is the set having the source as its only element. I will refer to the class of objects that the predicate is applied to as the domain of the quantification. The distribution of a quantification over an NP denotation can be viewed as specifying how the domain can be computed from the source. Where domain = source I will speak of individual distribution, where domain = {source} of collective distribution.

Individual and collective are not the only possible distributions. Consider the sentence

(7) All these machines assemble 12 parts.

This sentence may describe a situation in which certain machines assemble sets of twelve parts, i.e. a relation between individual machines and groups of twelve parts. If PARTS is the set denoted by 'parts', the direct object quantification domain is \mathcal{P}_{12} (PARTS), the subset of \mathcal{P} (PARTS) containing only those subsets of PARTS that have twelve members. I call this type of distribution group distribution. In this case the numerical quantifier indicates group size.

A slightly different form of "group quantification" is found in the sentence

(8) Twelve men conspired.

In view of the collective nature of conspiring, it would seem that 'twelve' should again be interpreted as indicating group size, so that the sentence may be represented by

(9) $\exists x \in \mathcal{P}_{12}$ (MEN): CONSPIRE(x)

However, as the existential quantifier brings out clearly, this interpretation would leave open the possibility that several groups of 12 men conspired, which is probably not what was intended. The more plausible interpretation, where exactly one group of 12 men conspired, I will call the strong group reading of the sentence, and the other one the weak group reading. On the strong group reading the quantifier 'twelve' has a double function: it indicates both source involvement and group size.

In a sentence like

(10) The crane lifted the tubes

there is no indication as to whether the tubes were lifted one by one (individual distribution), two by two (weak group distribution with group size 2), one-or-two by one-or-two (weak group distribution with group size 1-2), ..., or all in one go (collective distribution). The quantification is unspecific in this respect. In such a case I will say that the distribution is unspecific. If S is the source of the quantification, the domain is in this case the set consisting of the elements of S and the plural subsets of S.

Distribution and source involvement are the two central aspects of quantification that I will focus on here.

Quantification in two-level model-theoretic semantics

Consider a non-intensional verb, denoting a one-place predicate P (a function from individuals to truth values), which is combined with a noun phrase with associated source S (a set of indivi-

duals). The quantification then predicates the source involvement of the set of those elements of the quantification domain, defined by S and the distribution, for which P is true. This can be represented by a formula of the following form:

(11) S-INVOLVEMENT({x ∈ QUANT.DOMAIN: P(x) })

For example, consider the representation of the readings of sentence (1) 'Five boats were lifted', with individual, collective, and weak and strong group distribution:

(12a) $(\lambda z: \#(z)=5) (\{x \in \text{BOATS: LIFTED}(x)\})$

(12b) $(\lambda z: \#(z) \geq 1) (\{x \in \mathcal{P}_5(\text{BOATS}): \text{LIFTED}(x)\})$

(12c) $(\lambda z: \#(z)=1) (\{x \in \mathcal{P}_5^+(\text{BOATS}): \text{LIFTED}(x)\})$

(12d) $(\lambda z: \#(z)=5) (\cup_{\text{BOATS}} (\{x \in \text{BOATS} \cup \mathcal{P}^+(\text{BOATS}): \text{LIFTED}(x)\}))$

where $\mathcal{P}^+(S)$ denotes the set of plural subsets of S. The notation $\cup_S(D)$ is used to represent the set of those members of S "occurring in D"; the precise definition is:

(13) $\cup_S(D) = \{x \in S: x \in D \vee (\exists y \in D: x \in y)\}$

Note that in all cases the quantification domain is closely related to the source in a way determined by the distribution. I have claimed above that the distribution can be construed as a function that computes the quantification domain, given the source. Indeed, this can be accomplished by means of a function of two arguments, one being the source and the other the group size, in the case of a group distribution. A little bit of formula manipulation readily shows that all the formulas (12a-d) can be cast in the form

(14) $(\lambda z: N(\cup_S(z))) (\{x \in d(k,S): P(x)\})$

where S represents the quantification source, $(\lambda z: N(\cup_S(z)))$ the source involvement, k the group size, and d the "distribution function" computing the quantification domain. (For technical details of this representation see Bunt, in press). The most interesting point to note about this representation is that the distribution of the quantification, which in other treatments is always reflected in the syntactic structure of the representation, corresponds to a term of the representation language here. For this term we substitute expressions like $(\lambda k,s: \mathcal{P}_k^+(s))$ to obtain a particular interpretation.

I will now indicate how representations of the form (14) are constructed in the TENDUM system.

The construction of quantification representation in the TENDUM system

The TENDUM system uses a grammar consisting of phrase-structure rules augmented with semantic rules that construct a representation of a rewritten phrase from those of its constituents (see Bunt, 1983). For the sentence 'Five boats were lifted' this works as follows.

The number 'five' is represented in the lexicon as an item of syntactic category 'number' with representation '5'. To this item, a rule applies that constructs a syntactic structure of category 'numeral' with representation

($\lambda y: \#(y)=5$), which I abbreviate as FIVE. To this structure a rule applies that constructs a syntactic structure of category 'determiner' with representation

(15) $(\lambda x: (\lambda P: FIVE(\cup_x \{ \{x \in d(FIVE, x): P(x) \} \})))$

A rule constructing a syntactic structure of category 'noun phrase' from a determiner and a nominal (in the simplest case: a noun) applies to 'five' and 'boats', combining their representations by applying (15) as a function to the noun representation BOATS. After λ -conversion, this results in

(16) $(\lambda P: FIVE(\cup_{BOATS} \{ \{x \in d(FIVE, BOATS): P(x) \} \}))$

A rule constructing a sentence from a noun phrase and a verb applies to 'five boats' and 'were lifted', combining their representations by applying (16) as a function to the verb representation LIFTED. After λ -conversion, this results in

(17) $FIVE(\cup_{BOATS} \{ \{x \in d(FIVE, BOATS): P(x) \} \})$

Now suppose the sentence is interpreted relative to a domain of discourse where we have such boats and lifting facilities that it is impossible for more than one boat to be lifted at the same time. This is reflected in the fact that the RL predicate $LIFTED_x$ is of such a type that it can only apply to individual boats. Assuming that the ambiguous constant BOATS has the single instance $BOATS_x$ and that LIFTED has the single instance $(\lambda z: LIFTED(z))$, the instantiation rules, constrained by the type restrictions of RL, will produce the representation:

(18) $FIVE(\cup_{BOATS_x} \{ \{x \in BOATS_x: LIFTED(x) \} \})$

(For the instantiation process see Bunt, in press, chapter 7.) This is readily seen to be equivalent to the more familiar form:

(19) $\#(\{x \in BOATS_x: LIFTED(x)\}) = 5$

If, in addition to, or instead of the distributive reading we want to generate another reading of the sentence, then we extend or modify the instantiation function for LIFTED accordingly.

This shows how the analysis method generates the representations of only those interpretations which are relevant in a given domain of discourse, and does so without generating intermediate representations as artefacts of the use of a logical representation language.

References

- Bronnenberg, W.J., Bunt, H.C., Landsbergen, S.P.J., Scha, R.J.H., Schoenmakers, W.J., van Utteren, E.P.C. (1979) The question answering system PHLIQAI. In L. Bolc (ed.), Natural communication with computers, McMillan, London; Hanser Verlag, München.
- Bunt, H.C. (1982) The IPO Dialogue Project. SIGART Newsletter 80.
- Bunt, H.C. (1983) A grammar formalism with augmented phrase-construction rules. IPO Annual Progress Report 18.
- Bunt, H.C. (in press) Mass terms and model-theoretic semantics. Cambridge University Press.
- Bunt, H.C. and thoe Schwartzenberg, G.O. (1982) Syntactic, semantic and pragmatic parsing for a natural language dialogue system. IPO Annual Progress Report 17.
- Partee, B. (1975) Comments on C.J. Fillmore's and N. Chomsky's papers. In: D. Austerlitz (ed) The scope of American linguistics. De Ridder Press, Lisse.
- Quirk, R., Greenbaum, S., Leech, G., and Svartvik, J. (1972) A grammar of contemporary English. Longman, London.