

Discourse as a Knowledge Resource for Sentence Disambiguation

Tetsuya Nasukawa Naohiko Uramoto

IBM Research, Tokyo Research Laboratory

1G23-14, Shimotsuruma, Yamato-shi Kanagawa-keu 242 Japan

Phone 81 462 73-4574 FAX 81 462 73-7413

{nasukawa, uramoto} Ctrl net lbm com

Abstract

A consistent text contains rich resources of information such as collocation patterns that can be used to resolve ambiguities within its sentences. For example, attachment ambiguities in a sentence can be resolved by selecting a candidate attachment that matches attachments found in other sentences in the same discourse. Thus, discourse can be regarded as a valuable knowledge resource for sentence analysis. In this paper, we examine some features of discourse as a knowledge resource and propose a framework for natural language processing that provides a simple algorithm for using information extracted from discourse together with information stored in knowledge bases.

The experimental results of using our framework to disambiguate sentences in technical documents offer good prospects for improving the accuracy of a broad coverage natural language processing system that handles various texts without constructing knowledge bases for each text in advance. Some noteworthy features of discourse information are also deduced from the results of our experiments.

1 Introduction

In natural language processing, knowledge resources are indispensable for disambiguating input sentences and a large quantity of such resources is necessary to improve the accuracy of sentence analysis. However, constructing knowledge resources from hand-coded data requires enormous amounts of time and human resources and thus leads to a knowledge acquisition bottleneck. Much research has been devoted to developing methods for resolving ambiguities by using large quantities of knowledge resources that can be acquired semi-automatically from on the texts [Brown *et al.* 1991, Gale *et al.* 1992, Hindle and Rooth, 1993, Jensen and Binot, 1987, Nagao, 1990, Utsuro *et al.* 1992, Uramoto 1991]. Yet in the day-to-day operation of a practical natural language processing system such as a machine translation system, it is inevitable that a new text will contain various ambiguities that cannot be resolved with the knowledge

contained in the system. Moreover, a knowledge resource tends to be less effective for a text in a new domain.

On the other hand, a consistent text itself contains rich resources of information that can be used to resolve ambiguities within its sentences. As shown by Gale *et al.* [Gale *et al.* 1992] and Nasukawa [Nasukawa, 1993], polysemous words within the same discourse tend to have the same word sense with a probability of over 95% and the accuracy of word sense disambiguation can be improved by applying discourse constraints in such a manner that each polysemous word in a discourse takes the same word sense. Rinoshita *et al.* [Rinoshita *et al.* 1993] and Nasukawa [Nasukawa 1993] have reported that co-occurrence between words extracted from a text improves the accuracy with which attachment ambiguities in the text are resolved. Furthermore, information extracted from discourse by means of a simple algorithm makes it easier to select the correct antecedent of a pronoun from among candidate noun phrases, as shown by Nasukawa [Nasukawa, 1994].

In this paper, we study the potential usefulness of information extracted from discourse in contrast to information stored in knowledgebases, and develop a method based on a simple algorithm that resolves ambiguities in sentences by using information extracted from discourse together with information stored in knowledge bases.

In the next section we examine a couple of features of discourse information that allow more accurate sentence analysis than is possible with information stored in knowledge bases. Then, in the third section, we propose a framework of natural language processing that extracts information from a discourse and disambiguates each sentence in the discourse by using the discourse information together with information stored in knowledge bases. Finally, in the fourth section we give the results of experiments on the use of our framework to resolve structural ambiguities.

2 Information extracted from discourse

In this section we examine two features in discourse that are applicable to sentence disambiguation and compare the information extracted from discourse with the information stored in knowledge bases, in terms of its effectiveness for sentence disambiguation.

2 1 Discourse constraint and discourse preference

The first feature we discuss in this section is the "discourse constraint" that polysemous words within a discourse have the same word sense. According to Gale et al [Gale et al., 1992] and Nasukawa [Nasukawa, 1993] there is a high probability (98% according to Gale et al.) that polysemous words within the *same* discourse actually carry the same word sense. If this discourse constraint is assumed to hold, a result of word sense disambiguation applied in one sentence can be shared with all the morphologically identical words within the discourse.

Since a valid word sense disambiguation cannot always be made for every word appearing in a text, this procedure of sharing the results of word sense disambiguation improves the total accuracy of word-sense disambiguation. For example, if a polysemous word appears alone in a title, a heading, or a list item, there is no way to disambiguate its word sense except by referring to information in other sentences.

The second feature is "discourse preference," namely, a tendency for each word to modify or be modified by similar words within a discourse. During our analysis of computer manuals, we noticed that the same collocation patterns appeared frequently. For example, a collocation pattern in which the noun *cursor* occurred as the subject of the verb *move* appeared more than the times on one page of a computer manual. By assuming discourse preference, we obtain another hint on how to improve the accuracy of sentence analysis. For example, consider the following sentences extracted from a computer manual:

- (2 1) *Type your user name on the user line*
- (2 2) *You probably typed it on the line correctly, but*

Sentence (2 2) actually comes immediately after the sentence (2 1) in the manual. In sentence (2 1), the modifier of the prepositional phrase *on the user line* can be either *Type* or *your user name*, but in sentence (2 2), *on the line* uniquely modifies *typed*, therefore, by applying the discourse preference that the noun *line* tends to modify the verb *type* with the preposition *on*, we can resolve the structural ambiguity of sentence (2 1). Furthermore, the candidates for the antecedent of the pronoun *it* in sentence (2 2) are *your user name* and *the user line* in sentence (2 1), since both of them satisfy syntactic constraints. To select the actual antecedent, we can apply discourse preference, and thus *your user name* is preferred as an object of the verb *type* since a collocation pattern in which *your user name* is the object of the verb *type* appears in sentence (2 1).

2 2 Comparison with information stored in knowledge bases

In this paper, we assume that

- A general natural language processing system holds one or more knowledge bases that contain information for sentence disambiguation.
- Data are initially stored in the knowledge bases and updated through day-to-day operations.

- Multiple knowledge bases may be provided, to improve the accuracy of disambiguation by supplying appropriate information for various domains.

With the aim of relieving the knowledge acquisition bottleneck, we consider an example-based approach to sentence disambiguation in which disambiguation is based on information extracted from on-line texts. For example, a sentence analyzer called SENA [Uramoto 1991] analyzes a sentence and extracts word-to-word relationships that basically consist of both a modifier word and a modifiee word along with a marker of the relation between them (such as preposition, case marker, and so on) and prefers an attachment candidate that matches word-to-word relationships in its knowledge base (in this case, an example-base). Word-to-word relationships can be extracted from on-line texts semi-automatically with a little human intervention to eliminate ambiguities in the data [Nagao 1990, Uramoto 1991]. Within the framework of the example-based approach, it is easy to correct a specific disambiguation error by adding hand-coded word-to-word relationships, and an additional set of hand-coded word-to-word relationships may easily improve (the accuracy with which sentences in a specific text can be disambiguated). However, a knowledge base that contains only word-to-word relationships has a disadvantage in its coverage [Uramoto 1994]. Since the vocabulary as well as the usage of specific words varies from text to text depending on the writer and the topic, improving the coverage of word-to-word relationships in a practical domain is an important issue in the area of example-based systems. However, it may reduce the particularity of some word-to-word relationships such as those extracted from an idiomatic expression in which the modifier-modifiee relationship may not be general.

On the other hand, information extracted from a discourse contains some ideal features that can be applied to each sentence in the discourse. A sentence to be disambiguated and the discourse share the same writer and the same topic, thus information extracted from a discourse is more reliable for disambiguating a sentence inside the discourse than information in a knowledge base. However, the discourse cannot be the only resource for sentence disambiguation because

- Not all ambiguities can be covered by information extracted from the discourse.
- Information extracted from the discourse may itself contain ambiguities.
- Correct information may not always be extracted from the discourse because of ellipses and analysis failures and other reasons.

Examples of these reasons are provided in Section 4 along with the results of experiments on computer manuals. To improve the accuracy of sentence disambiguation, we propose to use information extracted from discourse together with information stored in knowledge bases. The next section describes the framework of our method.

3 Framework for using discourse information

In this paper we focus on structural ambiguities rather than any of the other ambiguities in sentences such as word sense ambiguities since information extracted from discourse can resolve structural ambiguities without being supplemented by other information and the effects of discourse information can be evaluated separately from those of information stored in knowledge bases. Resolution of other ambiguities such as selection of a word sense for a polysemous word, requires information initially stored in knowledge bases and discourse information is applied in such a manner that information extracted from knowledge bases is circulated within the discourse. In this section therefore we describe a framework for using information extracted from discourse to disambiguate a sentence by focusing on its structural ambiguities.

3.1 Extraction of discourse information

The first step of our framework is to extract information from a discourse while the second step is to resolve the ambiguities in each sentence of the discourse.

In the first step each sentence in the whole text given as a discourse is processed by a parser and the position of each instance of every lemma and its modifier-modifier relationships with other content words in the parser output are stored as discourse information. In accumulation of discourse information preference scores are given for each definite modifier-modifier relationship, such as the modification of type *on the line* in sentence (2.2) and lower scores are given for each ambiguous modifier-modifier relationship such as the possible modification of *Type B on the user line* in sentence (2.1).

Table 1 shows an example of such information extracted from the text of a computer manual consisting of sentences such as

A server shares resource with other workstations on the LAN*

In this table, scores of 10 for each definite modifier-modifier relationship and 3 for each ambiguous modifier-modifier relationship are assigned. Thus the sentence above yielded preference scores of 10 for *share* modifying *share* as a subject, 3 for *workstation* modifying *share* with the preposition *with* and so on.

3.2 Structural disambiguation

To resolve structural ambiguities, the discourse information is referred to for each modifier candidate. If the same modifier-modifier relationship that is the same collocation pattern, is found in the discourse information, the score of the collocation is extracted and added to the candidate's preference value. In order to improve the coverage of the discourse information, an online synonym dictionary [Collins 1984] is consulted so that a collocation pattern can be checked with a synonym in one side of the collocation (either the modifier or the modifier). The score after being multiplied by a weighting value is then added to its preference value of the

candidate. For example to select the verb *see* as the modifier of the propositional phrase, *for information*, in the sentence

Set Where to Look for Additional Task, Information in "Up and Running" for information on which books to use

a collocation pattern in which *for instruction* modifies *see* extracted from a later sentence in the same discourse (11 sentences after the sentence above)

For instructions see "Up and Running"

is referred to since *instruction* is defined LS a synonym of *information* in the online synonym dictionary [Collins, 1984]

In addition to the dictionary information coordinate structures in discourse information are referred to, and the coordinated words are also treated as synonyms in examining collocation patterns. For example in the sentence

In the MVS environment the product runs under TSO/E CICS and IMS

TSO/E CICS and *IMS* are coordinated and thus treated as synonyms in the same discourse.

Information on word-to-word relationships stored in knowledge bases also provides preference stores, in addition to those given by the discourse information and a store is added to the preference value of each candidate modifier. After evaluation of the preference value for each candidate modifier, the candidate with the highest preference is selected.

4 Results and discussion

After implementing our framework on an English-to-Japanese machine translation system called Shalt2 [Takeda et al, 1992] we processed a large number of technical documents such as patents, technical letters and in particular computer manuals. The results were encouraging and we found some interesting phenomena that should be taken into account when information extracted from discourse is used along with information stored in knowledge bases.

4.1 Coverage

The coverage of information used for structural disambiguation is shown in Figure 1. The figure reflects the results of disambiguating all the sentences in a computer manual consisting of 791 consecutive sentences by changing the size of the discourse from 10 to 791 sentences. Except for the result of referring to the whole 791 sentences as a discourse all the results indicate the averages of the results obtained by referring to each of several sample areas as a discourse. For example, to obtain data for the case in which the size of a discourse is 20 sentences, we examined 32 areas each consisting of 20 sentences, such as the 1st sentence to the 20th, the 51st to the 70th, and the 701st to the 720th.

In this experiment, we disambiguated the modifiers of structurally ambiguous phrases such as prepositional phrases, present and past participle clauses, relative

Table 1 Example of discourse information on the verb 'share'

Modifiers	POS	Relation	Word (preference value)
	Noun	from	
with			user (30) workstation (3)
OBJ			printer (50) directory (50) resource (20)
on			network (3) LAN (3)
SUBJ			server (10)
Verb	when		complete (10)
	if		exist (10)
Adverb	DIRECT		for example (10) also (10)

Modifiees	POS	Relation	Word (preference value)
	Noun	PASTPART	
for			procedure (3)
preinf			subdirectory (3) procedure (3)
Verb	preinf		want (10) install (10) create (3) follow (3)
	PASTPART		identify (3) type (3)
	for		go (3)

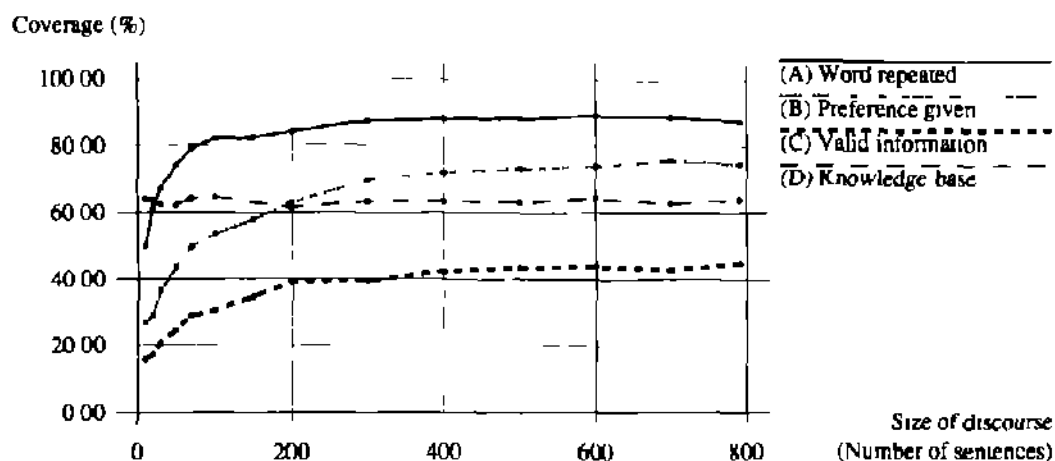


Figure 1 Relation between the coverage of discourse information and the size of the discourse

clauses and to-infinitive clauses. In Figure 1 (A) indicates the percentage of cases in which either the modifier Word or out of the modifiee candidate words was found in the discourse information. In other words, at least one of the words related to a resolution of a structural ambiguity was repeated in the discourse in a percentage of cases*, indicated by (A), this shows that main of the words were repeated within a relative small area of the text. To be precise, more than 80% of the words were used more than once within 100 sentences. (B) shows the percentage of cases in which some preference values were given to any of the candidate modifiers on the basis of information extracted from the discourse. However, the assignment of preference values does not imply that a valid disambiguation has been earned out since identical values may be assigned to different candidates. (C) indicates the percentage of cases in which

discourse information allowed a preferable candidate to be disambiguated. The percentage of cases in which the information stored in the knowledge base provided preference stores is shown by (D). Our knowledge base contains more than 57,000 unique word-to-word relationships extracted from almost 10,000 sentences in various computer manuals and technical documents and some broader coverage rules such as the one that prefers a prepositional phrase with *to* to modify a verb *go* is shown in the figure about 40% of structural ambiguities can be resolved with information extracted from the discourse provided the discourse contains more than 200 consecutive sentences.

To determine whether an input text should consist of (consecutive sentences) or combined 100 consecutive sentences from three different parts of the manual examined in the previous experiments on coverage - say the