# Empirical Study of Mandarin Chinese Discourse Analysis:
# An event-based approach

Yuan-Kai Wang, Yi-Shiou Chen, and Wen-Lian Hsu

Institute of Information Science, Academia Sinica,

115, Nankang, Taipei, Taiwan, R.O.C.

{ykwag,ysc,hsu}@iis.sinica.edu.tw

## Abstract

*Discourse analysis plays an important role in natural language understanding. Mandarin Chinese discourse, which has many different properties compared with English discourse, is still far behind in the construction of a basic computational model. In this paper, we propose an event model to elucidate anaphora and ellipsis in Mandarin Chinese. An event-based approach (EBA) based on the model is designed to resolve anaphora and ellipsis in Mandarin Chinese discourse. In this approach, we provide an event-based partial parser and an event-based reasoning mechanism. This approach is applied to the mathematics word problems of elementary school (MWES). Our results provide empirical evidence that the EBA can resolve many difficult problems in Mandarin Chinese discourse.*

## 1. Introduction

Among a variety of researches in natural language understanding, discourse analysis is one of the most important topics. Discourse is a sequence of utterances, for example, reports, books, dialogue, and conversations, that expresses certain ideas in writing or speech. In discourse analysis, the treatment of pronoun resolution, the finding of elliptic expressions, the derivation of implications, and the recognition of intentions touch the heart of understanding: What is the meaning of a discourse? Is there a unified computational model of meaning that is appropriate for the analysis of all kinds of linguistic expressions, from words to sentences, and for extracting information that is realized in an utterance?

Discourse analysis can be regarded as a basic research topic in search of answers for these questions. Approaches for discourse analysis are abundant in English. Many studies have focused on different issues of it. One important issue is the complex linguistic phenomenon: anaphora. Ellipsis is another critical issue, though it occurs less frequently in English.

Elliptic sentences appear overwhelmingly in Chinese. The use of ellipsis in Chinese discourse produces very versatile expressions. This situation also makes Chinese texts less formal in general. Lin and Soo [19] claimed that a highly accurate Chinese syntax parser may not be achieved without considering anaphora and ellipsis. Some people believe that the study of Chinese should emphasize more on semantics than that of English. Thus, to develop a comprehensive model for Chinese discourse understanding, one needs to focus on semantics and pragmatics.

How can we understand an utterance and make a meaningful reply? The major claim of this paper is that utterances are produced from actions (speech act theory) in order to have some effect on the hearer. We believe the underlying structure of a discourse is an event sequence. An event is an activity or a state expressed by a sequence of words. We believe events are discourse entities crucial for finding explanations of linguistic behavior into representation and computational issues. Our discourse analysis is modeled as a process of generating a set of events that can be resolved, deleted, and merged.

In this paper, we provide a computational linguistic model for Mandarin Chinese discourse analysis. In this model, an event-based partial parser is adopted to deal with ill-formed sentences in Mandarin Chinese. In

particular, the output of the parser is event sequence rather than parse tree.

Our study in discourse processing is conducted in the context of a realistic application domain. A system based on the work in this paper has been implemented in a *microworld*: the mathematics word problems of elementary school (MWES). Normally, a word problem uses several natural language sentences to form a coherent discourse that describes an event sequence. The discourse contains a question sentence that asks the student to compute an answer by performing some algebraic operations. While the word problem is normally simple and restricted in respect to its special domain on mathematical computation, it provides a wide range of linguistics phenomena and can be an excellent test-bed for our model. Consider the following example:

**Example 1**

```
    7 ,13        ?
John is1 7 years old this year. How old
will he be after 13 years?
```

The answer can be obtained by a simple calculation 7+13 = 20. The first clause in Chinese does not have a verb, which makes a classical parser difficult to generate a parse tree. The pronoun in the second clause is referred to the object *John* in the preceding sentence. How to find it and then know that the two digits 7 and 13 are associated by an addition relation? Empirical study on this microworld has another benefit that the performance of a discourse model can be easily verified through a numerical answer.

The remainder of the paper is organized as follows: In Section 2, we discuss the phenomena motivating the development of event model. Section 3 briefly describes our architecture of discourse processing. Sections 4 and 5 explain an event-based partial parser and reasoning mechanism. Section 6 introduces the application of elementary mathematics word problems. Finally, Section 7 makes conclusions.

## 2. Studies of Mandarin Chinese Discourse

There are four levels of ambiguity in Chinese discourse: lexical, syntactic, semantic (referential ambiguity), and pragmatic ambiguity. Anaphora, ellipsis, and metaphor are everywhere. Moreover, ellipsis is very important in Chinese discourse to express thoughts concisely.

### 2.1 Anaphora and Ellipsis

Anaphora is the use of a word as a regular grammatical substitute for a preceding word or a group of words. Considerable attention of research involving natural language processing, linguistics, and cognitive science has been paid to study anaphora. Many anaphora resolution methods have been proposed [9][16][18][22]. There are also many discussions or extensions of Chinese anaphora [3][4][15][19][23]. However, most attempts to model anaphor resolution in Chinese considered intrasentential anaphora [3][15][16][18][23]. Only a few considered intersentential anaphora [9][11][22], which is an important problem in discourse-level analysis. In [19] Lin and Soo resolved consecutive-sentential Chinese anaphora by a theta-grid chart parser. In this paper, an event-based approach for intersentential Chinese anaphora resolution is proposed.

To discover the problem of anaphora in Chinese discourse, consider the example 2 for pronominal anaphora. An underlined word in English sentences represents an anaphor. An italic word in Chinese example represents an elliptic word. Only underlined word is appeared in Chinese example.

**Example 2**

```
5  ,      — 5   ,——
?
John has 5 books. Mary has five more books
than him. How many books do they have?
```

To understand the anaphor " (he)" in example 2, the words " (John)" should be bound as the *antecedent* of " (he)". Another anaphor " (they)" is bound to the two objects: " (John)" and " (Mary)", that appear in different clauses. To resolve these anaphors needs to take account of intersentential structure.

Ellipsis is the omission from a sentence or other construction of one or more words understandable from the context that would complete or clarify the construction. It refers to the "hole" where an NP is understood and would have to be present in a complete sentence. Ellipsis produces ill-formed sentences that can not be recovered by syntactic interpretation. It also produces ambiguity. Consider example 3 modified from example 2:

**Example 3**

```
   5  ,       5 ,——      ?
John has 5 books. Mary has five more books
than him. How many books do they have?
```

Example 3 exhibits an elliptic form of the example 1 in that two phrases " (than him)" and " (books)" in

the second clause is disappeared, and the noun phrase " (books)" is still disappeared in the third clause.

## 2.2 Related Research in Discourse Analysis

Anaphora and ellipsis are the two main problems in discourse analysis. There are many approaches to deal with the above linguistic phenomena.

Coherence theory is one that is widely accepted. Many computational models have been proposed based on it [1][21]. In coherence theory, a discourse is composed of **discourse segments**. The utterances of a segment play a particular role in the discourse. Many coherence relations exist, such as evaluation, causal, elaboration, explanation, sequence, and so on. Hobbs [**12**] outlined the theory of discourse coherence. A local coherence of a segment is the coherence among utterances within a segment, and global coherence, which Grosz and Sidner argued that it depends on the intentional structure [6], is the coherence with other segments. In addition to cohesion, they addressed the role of attention and intention in discourse. Their theory includes linguistic structure, intentional structure, and a pushdown stack of focus spaces. Utterances cause the focus to shift by pushing or popping elements off the stack.

In contrast to the coherence theory, studies of knowledge representation theory [8][17] generally start from representing linguistic knowledge in logic form. It does not emphasize on pragmatics and world knowledge. Much work is influenced by this theory. Many researchers have found cue phrase to be an important structuring element for discourse [6][10].

Work in centering theory has been addressed by Grosz, Joshi, and Weinstein [7]. They believed that certain discourse entities were centers of an utterance that relates the utterance to other utterances. The same sentence may have different centers in different discourses. They argued that the coherence of discourse was affected by the compatibility between centering properties of an utterance and choice of referring expression. Hsu, Chen, and Wang [25] proposed a context sensitive model to interpret the understanding of concept by human beings.

There is some work done on Chinese discourse, most of them from linguistics. A few of them provided simple computational models to resolve part of phenomena of anaphora [2][3][19], but no basic theory for computing Mandarin Chinese discourse was proposed. In this paper we propose an event model as a model for Chinese discourse. The model examines phenomena of Chinese discourse, and considers cohesion at a local context level with some guidelines from linguistic theories.

# 3. Our Framework

This section will address the event model, and then briefly elucidate three components of our framework: knowledge, parser, and reasoning mechanism.

## 3.1 Event Model

Center is an important feature in discourse. However, it is inherently dynamic [7]. It generally changes across a sequence of discourse entities. This change will also produce expectation in other utterances. We use event to represent change and expectation. We treat a sequence of utterances in a virtual time line. Discourse entities are constituents across the time line. The change of center and the shift of focus are state transitions according to eventuality. Constructing discourse structure is to infer new events that generate state transitions.

An event is an activity that may be static (e.g., state, possession, description) or dynamic (e.g., action, transition). Utterances within a discourse are decomposed into events as discourse entities that simplify the analysis and understanding of discourse. A set of events within a sentence is called an event list. By specifying the basic units of discourse and decomposing discourses into structures of the basic units across a sequence of sentences, a discourse is transferred into events sequentially placed in a virtual time line.

A proper account of discourse can be processed by deleting events, merging events, and creating expected events. An expected event is an event generated after reasoning on event lists. It is a related event that is induced by existing events across different event lists. Therefore, expected events can be used beyond intrasentential level to resolve word sense disambiguation and anaphora. Besides, dynamic property of discourse (e.g., centering, focus) can also be represented by events.

## 3.2 Knowledge

Knowledge is important in natural language understanding. In our framework, four types of knowledge come into play in arriving at an understanding: 1. General knowledge about syntax, semantic, and linguistic knowledge, 2. general knowledge about discourse, 3. general knowledge about the world, 4. specific knowledge about the domain being discussed.

Syntactic, semantic, and linguistic knowledge are collected in dictionaries. There are two dictionaries in our framework. One is lexicon dictionary, and another is measure word dictionary. Lexicon dictionary contains more than forty thousand lexicons plus their syntactic category, word senses, and semantic category. The

lexicons, syntactic taxonomy, and semantic taxonomy are collected and analyzed mainly from our prior research [13]. Part of semantic taxonomy comes from *Tong2yi4ci2 ci2lin2*. Measure word dictionary comes from [14]. It has 427 measure words that is classified into 7 categories. It has 12352 noun phrases that each corresponds to one or several measure words. Measure word dictionary can serve as constraints for noun phrases. It is useful in parsing and ellipsis resolution.

General knowledge about discourse, general knowledge about the world, and special knowledge about the domain are represented by an expectation database. The expectation database is exploited by a general reasoning mechanism to produce expected events.

### 3.3 Parsing and Reasoning

We address an event-based approach that can interpret sentences. There are two subsystems in the approach. The first is called *event-based partial parser*, or EP parser. The second is a reasoning mechanism called *programmable reasoning with expectation* (PRE). Event-based representation provides the basis for computing the meaning of discourse in both subsystems.

EP parser is a parsing algorithm that takes a sentence, parses it, and generates an event list. It is an intrasentential processing algorithm. PRE resolves the problems of intersentential understanding. It receives event lists and generates expected events. Two parts of PRE: general reasoning and specific reasoning are performed iteratively until a special ending event is issued. Architecture of our approach is shown in Fig. 1. Details of EP parser is addressed in section 4. The PRE performing pragmatics analysis is addressed particularly in section 5.
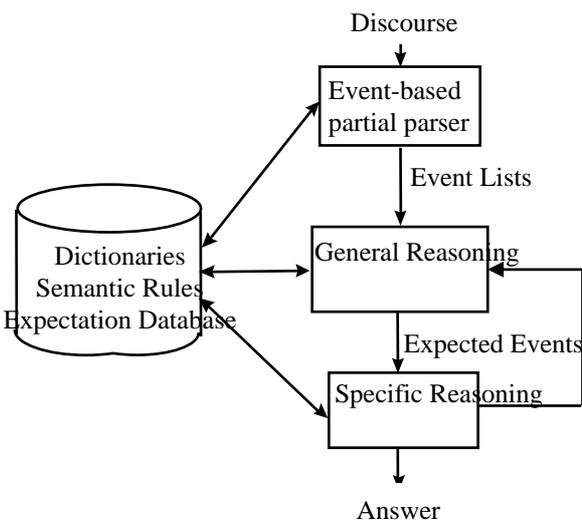


Figure 1    The framework of event-based Mandarin Chinese discourse analysis.

There is a tool that provides a visual knowledge engineering environment to assist the acquisition of knowledge. It is a front-end interface for collecting four types of knowledge. Event sequence is constructed by a visual programming environment as shown in Figs. 2, 3, and 4.



Figure 2    Interface for acquiring lexicons.



Figure 3. Interface for acquiring templates in EP parser.

## 4. Event-Based Partial Parser

Classical parser generally accepts a sentence, resolves syntactic ambiguity within the sentence, and generates syntactic parse trees of the sentence. However, our intention is for discourse analysis that crosses sentences. This will induce the problems of semantic processing, anaphora, and ellipsis. These problems cannot be resolved by classical parsers, such as chart parser. Even worse, these problems can degenerate. Therefore, our

4

parser should not just parse a sentence, but provide enough information for other sentences.



Figure 4    Interface for acquiring expectations in PRE.

An EP parser is proposed to perform partial parsing and generate event information for latter advanced processing. It is refined work of our previous research [13] that is extended to event-based approach. It considers both syntactic and semantics analyses. EP parser includes three stages: *template matching*, *consistency checking*, and *event generation*. Our templates are linguistic structures combined with semantic and syntactic information. The template matching stage identifies the template that matches the sentence and places them in the conflict set. There are semantic preference scores that allows semantic preferences to influence parse priorities, avoiding a combinatorial explosion of possible parse trees. EP parser solves syntax and semantics at the same time.

Generally, more than one acceptable meaning for an input expression can be obtained after the template matching stage. It occurs wherever a single syntactic constituent has more than one semantic interpretation. However, each semantic interpretation of phrases within a sentence is given a preference score. The consistency checking stage uses dynamic programming to choose the most probable interpretation.

The third stage generates events according to the set of optimal semantic interpretations. This stage uses event-generation templates to generate events. Event-generation template has similar representation with semantic rules. In other words, a sentence can produce many noun phrases, and many events will also be produced to form an event list.

One important characteristic of the EP parser is that it can carry out robust parsing, that does partial matching in template matching stage between phrase templates and a sentence. Partial parsing is necessary for ill-formed

sentences, which are very prevalent expressions in Chinese due to the use of ellipsis.

Another characteristic of EP parser is that some ambiguity is resolved by the use of constraints between phrases. For example, there is a common constraint in Chinese between measure words and noun phrases. The two legal sentences " 3 (I have three)" and " 3 (apple has three)" have different syntax and semantics though they have the same sentential form. The verb in the first sentence represents the *possession* of its subject, and the one in the second represents the *existence* of its subject. EP parser will resolve this type of ambiguity at event generation stage by means of constraints between subject and measure words to generate different events. It will check that the measure word " " can not be used for pronoun, but for a list of legal noun phrases recorded in measure word dictionary.

## 5. Reasoning Mechanism

A reasoning mechanism called *programmable reasoning with expectation*, or *PRE*, accepts event lists generated by EP parser. Programmable reasoning with expectation means that the reasoning process is guided by activating appropriate expectation to evoke other expectations. PRE produces expected events, which are the events generated in the reasoning process of PRE, and reasons in an iterative fashion.

An expectation includes three parts: score, category, and procedure construct. Score is preference priority of an expectation. Category indicates the classification of an expectation. A procedure construct includes a sequence of operations.

There are four parts of an operation in a procedure construct: s*tep number*, *operator*, *operand*, and *fail*. *Step number* is the order number assigned to the construct. *Operator* is an execution code performing a single task such as going to previous sentence or generating an expected event. *Operand* is the parameter of operator. It has two parts: an example template, and a set of events. *Fail* includes three directives: abort, break, and continue. When an operation is failed, the abort directive will escape the procedure construct. The break directive will stop a Repeat-EndRepeat loop. The continue directive will skip those operations between this operation and the nearest EndRepeat operation.

There are 20 operators that are classified into four categories: *condition*, *move*, *report* and *other*. Condition operators include **Match**, **Peak**, **ContextMatch**, **Check**, and **CheckList**. Move operators contain **Reverse**, **Backward**, **First**, **NextS**, **PrevS**, **Pointer**, and **Goto**. Report operators involve **Insert**, **WriteDown**. **Calculate**, **Answer**, and **ContextInsert**. Other operators comprise

**Ask**, **Repeat**, **EndRepeat**, and **End**.

PRE will reason as follows: It matches an expectation to an event, where the match condition is necessary to be described in the first operation of each expectation. This is called the general reasoning process. A matched expectation will be evoked and its procedure construct is executed. Some expected events may be generated, or answer may be calculated. When an ask operator is being executed, PRE will go into a deeper expectation resolution process. The original procedure construct will be resumed after the expectation temporarily evoked by the ask operator is completed. If there is no ask operator in the procedure construct, general reasoning is performed again to match other expectations to unmatched or new events.

Anaphora has proved to be a very interesting but difficult linguistic phenomenon that involves a large number of problems among lexical, syntactic, semantic, and pragmatic levels. PRE can handle anaphor resolution (including pronouns, reflexives, and definite reference), temporal analysis, and ellipsis. Here we give an example showing the pronoun resolution strategy of PRE, that is, how to solve the structure that anaphor/antecedent pair does not need to occur in the same clause or sentence.

We treat all information useful for resolving pronoun as constraints, and implement an idea of recency constraint satisfaction by PRE, which states that the antecedent should be the most recently mentioned object that satisfies all the constraints. The implementation will check the most recent local event for an antecedent that matches all the constraints related to the pronoun. If no antecedent is found in the local context, then move forward along the virtual time line to the next most recently local context and check constraints for it. It acts as backward-looking recency constraint satisfaction.

Constraints come from many sources. Pronoun resolution may use gender and number to eliminate some objects that are not possible to be antecedents, and other constraints derived from imposing symptoms may introduce further restrictions.

Fig. 5 is an example procedure construct for the third-person pronoun " (he)." This procedure construct will be executed when its first operation is true. Next line issues a control for the analysis to backward direction. PrevS operator then put a starting sentence pointer in previous sentence. The operations with step numbers from 4 to 8 enter loop control. The operation 6 examines the constraints of gender between pronoun and objects by an ask operation. When the examination is passed, an expected event is produced by an insert operator. Be aware that the word " (father)" appeared in steps 5, 6, and 7 is a variable. PRE will check if the semantic category of the variable is compatible with that of " (father)",

which is *consanguinity* in our lexicon.

| Step number | Operator | Operand | Fail |
|---|---|---|---|
| 1 | Match | | Abort |
| 2 | Backward | | Abort |
| 3 | PrevS | | Abort |
| 4 | Repeat | | Abort |
| 5 | Match | | Break |
| 6 | Ask | -? -# - | Continue |
| 7 | Insert | -# - | Abort |
| 8 | EndRepeat | | Abort |
| 9 | End | | Abort |

Figure 5. An example of execution construct for pronoun resolution.

## 6. Application System

Corpus of MWES is taken from many sources. One of the major sources is *Textbook of Mathematics of Elementary School* in 1996. Our corpus includes about 300 word problems based on mathematics at the third and fourth grades of elementary school in Taiwan. Its difficulty is at the level of hybrid operations of four elementary algebraic operators: addition, subtraction, multiplication, and division. About 90 percentage of word problems in our corpus can be represented and solved by the event-based approach. The application interface of MWES is illustrated in Fig. 6.

Here we will take example 1 to illustrate the processing of word problems by our event model. The example discourse will produce four events as shown in Fig. 7. The pronoun " (he)" will be resolved by the procedure construct of the expectation in Fig. 5. The anaphor is bound to " (John)" because the gender and semantic category are compatible between them. The third event in Fig. 7 will match an expectation that issues an expected event to infer that "13 (after 13 years)" means *add* 13 to John's current age. A simple computation 7+13=20 is then obtained by performing an Answer operator in an expectation.

## 7. Conclusions

In this paper, we have sketched a new discourse processing approach. It is an initial attempt to develop a model that uses event as discourse entity of Mandarin
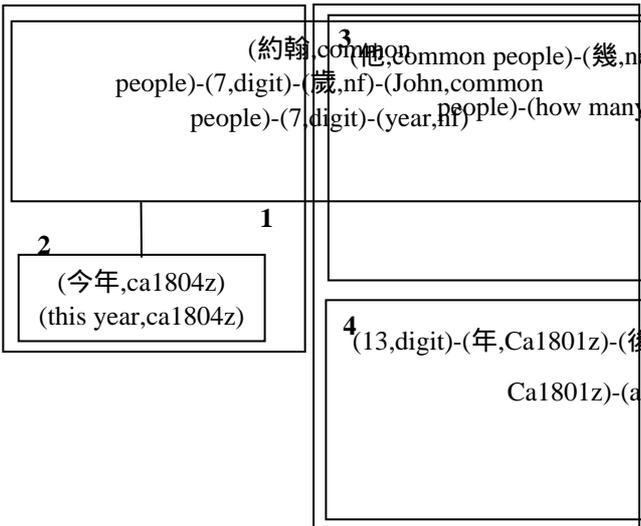
Chinese. Nonetheless, we believe the event model can handle both Chinese and English discourses because anaphora and ellipsis are common phenomena within both natural languages.

Event model is a more expressive representation than logic. It provides procedural execution constructs to describe constraints, retain memory, and generate outputs. Information encoded by procedural constructs includes both representation and operations. That is, it is not a static representation of discourse, but an operational representation of discourse, which does not use a complex proposition to represent discourse structures only.

In this paper, we presented a mechanism and tools to analyze Mandarin Chinese discourse. Although this is only a pilot study, the preliminary results in the word problems of elementary mathematics are quite encouraging. We hope this will lay down the ground work for further development into many practical discourse application systems, such as informational retrieval, natural language generation, and machine translation.



Figure 6. Application interface of word problems.



Figure 7. Four events of the event list for example 1.

## References

[1] Allen, J. *Natural Language Understanding*, 2nd, Benjamin/Cummings, 1995.

[2] Chen, B. L. and V. W. Soo, "An acquisition model for both choosing and resolving anaphora in conjoined Mandarin Chinese sentences," *COLING-92*, pp. 274-280, August, 1992.

[3] Chen, H. H. "Anaphora resolution algorithms for Mandarin Chinese," *Communications of COLIPS*, vol. 3, no. 2, pp. 59-67, 1993.

[4] Chen, H. H. "The transfer of anaphors in translation," *Literary and Linguistic Computing*, vol. 7, no. 4, pp. 231-238, 1992.

[5] Dahlgren, K., "Knowledge representation for commonsense reasoning with text," *Computational Linguistics*, vol. 15, no. 3, pp. 149-170, 1989.

[6] Grosz, B. J. and C. L. Sidner, "Attention, intentions, and the structure of discourse," *Computational Linguistics*, vol. 12, no. 3, pp. 175-204, 1986.

[7] Grosz, B. J., A. K. Joshi, and S. Weinstein, "Centering: a framework for modeling the local coherence of discourse," *Computational Linguistics*, vol. 21, no. 2, pp. 203-225, 1995.

[8] Guenthner, F., H. Lehmann, and W. Schonfeld, "A theory for the representation of knowledge," *IBM Journal of Research and Development*, vol. 30, no. 1, pp. 39-56, January, 1986.

[9] Heeman, P. A. and G. Hirst, "Collaborating on referring expressions," *Computational Linguistics*, vol. 21, no. 3, pp. 351-382, 1995.

[10] Hirschberg, J. and D. Litman, "Empirical studies on the disambiguation of cue phrases," *Computational Linguistics*, vol. 19, no. 3, pp. 501-530, 1993.

[11] Hirst, G. "Discourse-oriented anaphora resolution in natural language understanding: a review," *American Journal of Computational Linguistics*, vol. 7, no. 2, pp. 85-98, 1980.

[12] Hobbs, J. R. *Literature and Cognition*, CSLI Press, Stanford, California, 1990.

[13] Hsu, W. L. "Chinese parsing in a phoneme-to-character conversion system based on semantic pattern matching," *International Journal on Computer Processing of Chinese and Oriental Languages*, vol. 40, pp. 227-236, 1995.

[14] Huang, C. R., K. J. Chen, and C. S. Lai, *Guo2Yu3Ri4Bao4 Liang2C2Dian3*, Taiwan, 1997.

[15] Huang, S. "Getting to know referring expressions: anaphor and accessibility in Mandarin Chinese," *Proceedings of ROCLING V*, pp. 27-51, 1992.

[16] Huls, C., E. Bos, and W. Claassen, "Automatic referent resolution of deictic and anaphoric expressions," *Computational Linguistics*, vol. 21, no. 1, pp. 59-79, 1995.

[17] Kamp, H. "A theory of truth and semantic representation," *Formal Methods in the Study of Language*, MC TRACT 135, J. A. G. Groenendijk, T. M. V. Janssen, and M. B. J. Stokhof (Eds.), Amsterdam, p. 277, 1981.

[18] Lappin, S. and H. J. Leass, "An algorithm for pronominal anaphora resolution," *Computational Linguistics*, vol. 20, no. 4, pp. 535-561, 1994.

[19] Lin, K. H. C. and V. W. Soo, "Toward discourse-guided theta-grid parsing for Mandarin Chinese -- a preliminary

report," *Proceedings of ROCLING II*, pp. 259-270, 1993.

[20] Nakhimovsky, A. "Aspect, aspectual class, and the temporal structure of narrative," *Computational Linguistics*, vol. 14, no. 2, pp. 29-43, 1988.

[21] Russell, S. and P. Norvig, *Artificial Intelligence,* Prentice Hall, Chapter 22, 1995.

[22] Sidner, C. L. "Focusing for interpretation of pronouns," *American Journal of Computational Linguistics*, vol. 7, no. 4, pp. 217-231, 1981.

[23] Tang, C. C. J. "Chinese reflexives," *Natural Language and Linguistic Theory*, vol. 7, 1989, pp. 93-121.

[24] Webber, B. L. "Tense as discourse anaphora," *Computational Linguistics*, vol. 14, no. 2, pp. 61-73, 1988.

[25] Hsu, W. L., Y. S. Chen, and Y. K. Wang, "A context sensitive model for concept understanding," *Proceedings of Third Int. Conf. on Information-Theoretic Approaches to Logic, Language, and Computation*, 1998.