

Design of Chinese HPSG Framework for Data-Driven Parsing

Xiangli Wang^a, Shunya Iwasawa^a, Yusuke Miyao^a, Takuya Matsuzaki^a,
Kun Yu^a and Jun'ichi Tsujii^{a,b,c}

^a Department of Computer Science, University of Tokyo,
Hongo 7-3-1, Bunkyo-ku, Tokyo, 113-0033, Japan
{xiangli, iwasawa, yusuke, matuzaki, kunyu, tsujii}@is.s.u-tokyo.ac.jp

^b School of Computer Science, University of Manchester

^c National Center for Text Mining

Abstract. Data-driven parsing has been a main method for analyzing natural languages. We aim at exploring a data-driven Chinese parser, by basing it on Head-driven Phrase Structure Grammar (HPSG). Unlike for English, there is still no available Chinese HPSG framework. As the first step of our work, we design a Chinese HPSG framework, which can be used as the basis for a practical parser. In this paper, 1) we present a Chinese syntactic structure system and 2) we design a primary Chinese HPSG framework.

Keywords: HPSG, data-driven parsing, Chinese HPSG framework, coverage, consistency.

1 Introduction

Data-driven parsing has been proven to be the most effective approach to development of a practical parser. It can deliver a parser with broad-coverage and high-accuracy. Some English data-driven syntactic parsers have been developed in the past (Charniak and Johnson, 2005; McDonald and Pereira, 2006; Miyao and Tsujii, 2005). Following the success of the research on English data-driven parsing, the same methodology has been applied to Chinese parsing (Levy and Manning, 2003; Wang et al., 2005; Guo et al., 2007).

The goal of our research is to develop a data-driven Chinese parser that is based on Head-driven Phrase Structure Grammar (HPSG) (Sag et al., 2003). Since an English data-driven parser based on the HPSG framework has been developed by our group (Miyao and Tsujii, 2005), we follow the same methodology for developing a Chinese parser. We first convert an existing Chinese treebank into an HPSG treebank, based on which we can obtain a large lexicon and a statistical model for choosing the most plausible interpretation.

Since the HPSG framework for English has been studied comprehensively (Sag et al., 2003), we can rely on the framework developed by linguists to establish the initial framework for a data-driven HPSG parser. For Chinese, however, although some linguistic studies have been conducted in the HPSG framework (Gao, 2000; Wang and Liu, 2007), they have narrowly focused on very specific phenomena, and a HPSG framework for Chinese which is comprehensive and systematic enough to cover wide range of phenomena in Chinese still does not exist.

As the first step towards a data-driven HPSG parser for Chinese, our work aims to design a Chinese HPSG framework, by which we can trigger lexical acquisition from an annotated corpus.

The requirements for such a framework are: 1) it should be versatile enough to cover a wide range of sentences which appear in the real-world text, and at the same time, 2) it should be systematic enough to avoid superfluous ambiguities. To achieve these two goals, the framework should have a disciplined view of syntactic structures of the language and provide a simple but versatile data structure to treat diverse phenomena. Our experience of developing a data-driven parser for English

shows that a set of a restricted number of rule schemas with rich lexical constraints of feature structures in HPSG provide such a disciplined as well as versatile view of linguistic structures.

Our contributions in this paper are as follows:

1. We present a Chinese syntactic structure system.
2. We design a primary Chinese HPSG framework.

In section 2, we present a Chinese syntactic structure system; in section 3, we propose a primary Chinese HPSG framework. In section 4, we conclude and present our future work.

2 Syntactic Structure System of Chinese

We first define a set of *syntactic structures* in Chinese in a way which reflects the traditional view of Chinese linguistics. The set is defined in the form of six graphical schemas in this paper, but these graphical schemas correspond to rules in Sentence Structure Grammar (SSG), a Chinese grammatical rule system, which was designed for and actually used by a rule-based parser (Wang and Miyazaki, 2007). In SSG, we divide structures in Chinese into three levels: *predicative part level*, *simple sentence level*, and *complex sentence level*. Then, we provide six graphical schemas.

Figure 1 provides a schema of the predicative parts. The predicate (P) is the head. A predicate subcategorizes for up to two objects (O_1, O_2). We call the constituents that modify the predicate from the left-side as: *adverbial* (Z). Z includes five types of constituents: temporal phrase, prepositional phrase, auxiliary verb, adverb, and “地/de” phrase.

We refer to the constituent that modifies the predicate from the right-side as: *complement*¹ (C). We divide the complements into three types according to their position. C_1 refers to the complements after the predicate; C_2 refers to the complements after the first object; C_3 refers to the complements after the second object. The constituents in the predicative part are divided into two types: one is the *indispensable constituent*, and the other is the *dispensable constituent*. The predicate (P) and objects (O_1, O_2) are indispensable constituents, and the adverbial (Z) and complement (C) are dispensable constituents. Structures which consist of indispensable constituents are *basic predicative structures*.

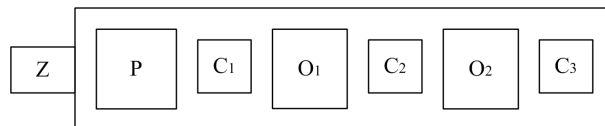


Figure 1: The model of predicative parts

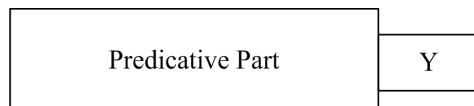


Figure 2: The model of subject-less simple sentence

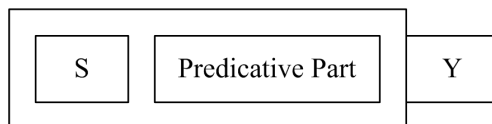


Figure 3: The model of ordinary simple sentence

Basic predicative structures can be divided into three types: 1) P: consists of the predicate; 2) P O_1 : consists of the predicate and one object; and 3) P $O_1 O_2$: consists of the predicate and two objects. 1a, 1b and 1c are examples of the basic predicative structures of the predicative parts. The predicative part framework obtains a high coverage by adding dispensable constituents into each basic predicative structure. For example, 1d, 1e, 1f and 1g are structures that consist of a basic predicative structure P O_1 and some adverbial constituents. “在家/at home”, “偷偷地/secretly”, “今天/today”

¹ The term *complement* in our paper is a different concept from that in the HPSG. *Complement* here refers to a grammatical constituent that appears on the right-side of a predicate, and modifies it.

and “可以/may” are adverbial constituents (Z). These show typical phrases for Z. “在家/at home” is a prepositional phrase; “今天/today” is a temporal phrase; “偷偷地/secretly” is a “地” phrase; and “可以/may” is an auxiliary verb.

- 1a. 吃/eat
eat
- 1b. 吃/eat 苹果/apple
eat apples
- 1c. 送/give 麦克/Mike 苹果/apple
give Mike apples
- 1d. 在/at 家/home 吃/eat 苹果/apple
eat apples at home
- 1e. 今天/today 在/at 家/home 吃/eat 苹果/apple
eat apples at home today
- 1f. 今天/today 在/at 家/home 偷偷地/secretly 吃/eat 苹果/apple
eat apples secretly at home today
- 1g. 今天/today 可以/may 在/at 家/home 偷偷地/secretly 吃/eat 苹果/apple
may eat apples secretly at home today

Figure 2 and 3 present the two schemas for simple sentences. We divide the simple sentences into two types: *subject-less simple sentence* and *ordinary simple sentence*. A subject-less simple sentence consists of one predicative part. An ordinary simple sentence consists of a subject (S) and a predicative part. Both subject-less simple sentences and ordinary simple sentences can be followed by sentence final particles (Y) (Dexi, 1982). 2a and 2b are corresponding examples. The predicate for an ordinary simple sentence subcategorizes for a subject and sentence final particles are optional. The predicate is the syntactic head of the whole sentence, and determines the number and the types of arguments that it can take. For example, “吃/eat” takes one nominal object in 2b, while “学/learn” takes one verb phrase as its object in sentence 2c. Generally, a subject is a noun phrase, but some predicates take a clause or a verb phrase as a subject, as is shown in 2d.

- 2a. 下/fall 雨/rain 了/past-tense
It rained
- 2b. 他/he 吃/eat 苹果/apple
he eats apples
- 2c. 他/he 学/learn 游泳/swimming
he learns swimming
- 2d. 学/learn 游泳/swim 很/very 有趣/interesting
Learning swimming is very interesting

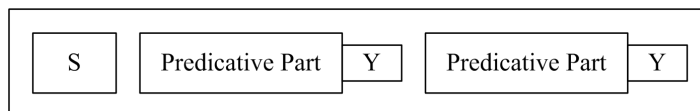


Figure 4: The model of subject-sharing complex sentence

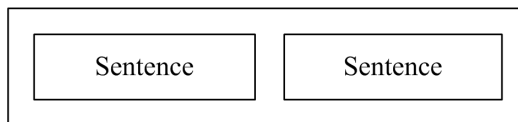


Figure 5: The model of coordination complex sentence

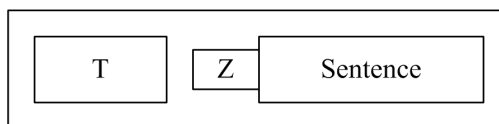


Figure 6: The model of topic sentence

Complex sentences are classified into three types: *subject-sharing complex sentence*, *coordination complex sentence* and *topic sentence*. A subject-sharing complex sentence refers to the sentences that consist of a subject and multiple predicative parts that can be followed by optional sentence-final particles. Figure 4 provides a graphical schema for this type of sentence. 3a is an example for such a sentence. Coordination complex sentence refers to sentences that consist of

multiple simple sentences or complex sentences. Figure 5 presents a schema for this type of sentence. 3b is an example of such a sentence. Topic sentences refer to sentences that have topics (T) in the beginning of a sentence, while adverbials may appear between the topic and the main part of the sentence. Figure 6 presents a schema for this type of sentence. 3c is an example of topic sentence.

- 3a. 他/he 唱歌/sing, 也/also 跳舞/dance
he sings, but also dances
- 3b. 他/he 唱歌/sing, 她/she 跳舞/dance
he sings and she dances
- 3c. 大象/elephant 鼻子/nose 长/long
elephant's nose is long

Six graphical schemas express most of the Chinese syntactic structures at a sentence-level, with the exception of the *Ba*-structure, and *Bei*-structure; these two structures involve dislocation of phrases from their basic positions which these six schemas require. These structures are covered by our Chinese HPSG framework, and we will introduce the details in the next section.

3 The Design of Chinese HPSG Framework

The formalized framework HPSG uses a small number of rule schemas and a large number of lexical entries to describe language. Our basic policy of Chinese HPSG is to exploit rule schemas defined for English with minimum changes. Although a possible solution would be to create an initial grammar with the help of the Grammar Matrix (Bender et al., 2002), we refer to the rule schemas used in an existing HPSG parser (Miyao, 2006), because we intend to apply the technology of this parser to our Chinese parser. This does not only reduce the cost of development of Chinese grammar but also confirm the assumption that, despite surface diversity, human languages share the same organization principles. For example, we do not introduce new rule schemas specific to Chinese unless they are absolutely necessary.

We generalize Chinese syntactic structures into five structures based on the Chinese syntactic structure system that we proposed in the previous section. The five structures are: 1) *predicate argument syntactic structure*, 2) *modifier head structure*, 3) *coordination structure*, 4) *topic structure* and 5) *filler-head structure*. We design signs and schemas for these five structures. Topic-Head Schema is newly introduced to deal with the topic structure, while the others are the same ones for English with some revisions necessary for Chinese. In section 3.1, we introduce our design for signs and schemas. In section 3.2, we introduce ways of dealing with some particular and essential structures in Chinese.

3.1 Design of Signs and Schemas

We define signs and schemas used by our English grammar which basically follow the definition by Sag et al. (2003). Figure 7 shows a lexical sign for the transitive verb, “吃/eat”. PHON is a feature for a surface string of a word. HEAD is a feature that expresses the characteristics of the head word of a constituent. FORM describes the form of a POS tag. MOD, SPR and COMPS represents the selectional constraints of a modifier, left arguments and right arguments. GAP and STOP_GAP express the constraints for moved arguments. TOPIC is a feature that deals with topic structures. INDEX and RESTR express semantic structures. INDEX represents the predicate argument structure of the main predicate, and RESTR provides semantic restrictions to the main predicate.

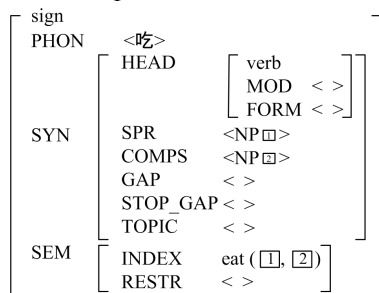


Figure 7: A lexical sign for “吃” (eat)

We define the following schemas to deal with Chinese syntactic structures: Specifier-Head Schema, Head-Object Schema, Head-Modifier Schema (left-head), Modifier-Head Schema (right-head), Filler-Head Schema, Coordination Schema and Topic-Head Schema. The Topic-Head Schema is unique to our Chinese grammar (the details will be discussed in Section 3.2.2), while the other schemas are the same as those for English. We explain how the five structures correspond to HPSG schemas.

The predicate argument syntactic structure is a syntactic concept and is different from the data structure (PAS) for the semantic representation. They include the *argument-head structure* and *head-argument structure*. The argument-head structure refers to structures for which the argument appears at the left-side of the head, and head-argument structure refers to structures for which the argument appears at the right-side of the head. Generally, the subject appears on the left-side of the predicate, and the object appears at the right side of the predicate in Chinese, as shown in Figure 1 and 3. We use the *Specifier-Head Schema* to deal with the argument-head structure, and use the *Head-Object Schema*² to deal with the head-argument structure. In a *Ba*-structure or *Bei*-structure, the object marked by *Ba* or the subject marked by *Bei* appears on the left-side of the predicate. We use Specifier-Head Schema to capture the arguments marked by *Bai* or *Bei*. 4a and 4b are examples for *Ba*-structures and *Bei*-structures.

- 4a. 他/he 把/Ba 书/book 看/read
he reads books.
- 4b. 书/book 被/Bei 他/him 看/read
he reads books

Modifier head structures include the *modifier-head structure* and the *head-modifier structure*. At the sentence level, a modifier-head structure refers to the structures for which modifiers appear on the left-side of the predicate, as is shown by Z in Figure 1 and 6, and the head-modifier structure refers to structures for which modifiers appear at the right-side of the predicate, as in C₁, C₂, C₃ and Y shown by Figure 1, 2, 3, 4 and 5, respectively. We use the *Modifier-Head Schema* and the *Head-Modifier Schema* to deal with the two structures.

Coordination structures include the *predicative-part-coordination structure* and the *sentence-coordination structure*. The two structures are shown by Figure 4, and Figure 5. We use the *Coordination-Schema* to deal with coordination structures.

The topic structure refers to structures for which a topic appears before the subject, as is shown in Figure 6. We use the *Topic-Head Schema* to deal with such kind of structures. The Topic-Head Schema is proposed by Gao (Gao, 2000), as shown in Figure 8. We introduce this concept in detail in the next section, since topic structures are particular and essential structures in Chinese.

The filler-head structure refers to long-dependency structures. It includes the *pre-object-as-subject structure* and *pre-object-as-topic structure*. We introduce the two structures here. The pre-object-as-subject structures refer to structures for which the object appears at the subject position. For example, in 5a “car” is the object of the predicate verb, “repair” in deep meaning; however, it appears in the subject position. The pre-object-as-topic structure refers to structures for which the object appears on the topic position. For example, in 5b, “car” is the object of the verb, “repair.” It appears at the topic position of the sentence. We use the *Filler-Head Schema* to deal with the two structures.

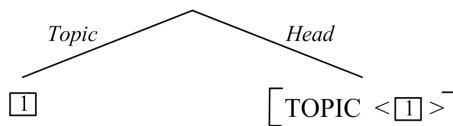


Figure 8: Topic-Head Schema

- 5a. 车/car 修/repair 了/past-tense
a car was repaired.
- 5b. 车/car 他/he 修/repair 了/past-tense
he repaired a car.

² The term *Head-Object Schema* is the same as the Head-Complement Schema, which was proposed by Sag et al. (2003); we use this term for avoiding the confusion between the two concepts in our paper: the object and complement.

predicate are subjects (Ma, 1898; Wang, 1957; Chao, 1968). On the contrary, some linguists consider Chinese to be a topic prominence language, and Chinese does not have a subject. Rather, it only has a topic (Lapolla, 1990; Schachter, 1976). From another perspective, Gao presents a definition of topic: Topic refers to a syntactic component of a sentence that appears to the left of the subject (Gao, 2000).

Gao has examined the topic structure with the theoretical framework of the HPSG. In Gao's study, he suggested that the topic structure should be treated as base-generated, and represented with a new ID schema: the Topic-Head Schema for all the topic structure.

There are two problems in Gao's analysis:

1. Topic is not formally defined.

In Gao's paper, the topic is defined as a constituent that appears before the subject. That is, the definition of topic depends on the definition of subject; however, it is not clear what a subject is.

2. Multiple topics may appear.

According to the definition of Gao's topic, in our observation, there may be multiple topics preceding the subject, and six kinds of constituents that can be a topic: noun phrase (NP), verb phrase (VP), sentence (S), prepositional phrase (PP), locative phrase (LP) and temporal phrase (TP). For example, there are two topics in 7a, three topics in 7b, and four topics in 7c. Then, a large number of lexical entries would be needed by one verb. For example, to cover 7a, 7b and 7c, we have to design three lexical entries for "eat," TOPIC <TP, NP>, TOPIC <TP, PP, NP>, TOPIC <PP, TP, NP>.

- 7a. 昨天/yesterday 约翰/John 他/he 吃/eat 苹果/apple 了/past-tense
John ate apples yesterday
- 7b. 昨天/yesterday 在/at 学校/school 约翰/John 他/he 吃/eat 苹果/apple 了/past-tense
John ate apples at school yesterday
- 7c. 在/at 学校/school 昨天/yesterday 约翰/John 他/he 吃/eat 苹果/apple 了/past-tense
John ate apples at school yesterday
- 7d. 约翰/John 今天/today 吃/eat 苹果/apple
John eats apples today
- 7e. 今天 约翰 吃 苹果
John eats apples today

To resolve the two problems above, we give clear definitions for the subject and topic as listed below:

- 1) The subject should be a NP, VP, or S that appears before the predicate.
- 2) Topics are the constituents that appear before the subject.
- 3) A NP, VP or S is an *ordinary topic*, and TP, PP and LP are *modifier topics*.

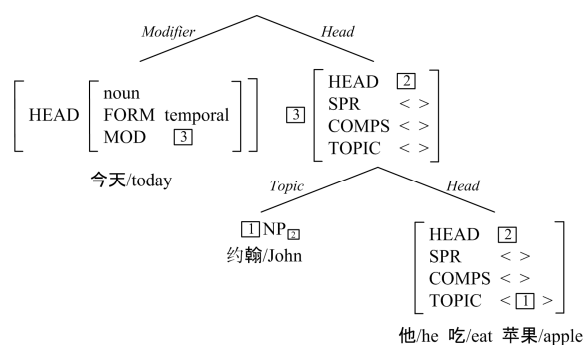


Figure 10: Treatment of a multi-topic structure

Our definitions are stricter than Gao's definition. For example, in Gao's definition, both 7d and 7e may be treated as topic structures, but by our definition, 7d is not a topic structure, and only 7e is treated as a topic structure. The ambiguity can be removed effectively.

In our observation, there is only one ordinary topic in a sentence. We use the Topic-Head Schema (Gao, 2000) to deal with the ordinary topic, and the Modifier-Head Schema to deal with modifier topics. We resolve Gao's multiple topic problem in this way.

Figure 10 presents how we deal with a multi-topic structure “今天/today 约翰/John 他/he 吃/eat 苹果/apple (Today John eats apples)”. “他/he” is the subject of the sentence “他/he 吃/eat 苹果/apple (he eats apples)”. Since “约翰/John” is the ordinary topic, we use the Topic-Head Schema to deal with the relationship between “约翰/John” and the sentence “他/he 吃/eat 苹果/apple (he eats apples)”, while we use the Modifier-Head Schema to deal with “今天/today” and the sentence after it.

4 Conclusion and Future Work

We presented a Chinese syntactic structure system and a primary HPSG framework for Chinese. Our framework covers Chinese structures that frequently appear in Chinese real-world texts. Currently, we are implementing our grammar and evaluating the coverage of the grammatical framework with sentences taken from a Chinese grammar textbook and newswire texts. In future work, we will obtain an HPSG lexicon and a disambiguation model from the Penn Chinese Treebank, using the grammatical framework presented in this paper. A data-driven HPSG parser for Chinese will be explored based on the framework we proposed.

References

- Bender, E. M., D. Flickinger and S. Oepen. 2002. The Grammar Matrix: An Open-Source Starter-Kit for the Rapid Development of Cross-Linguistically Consistent Broad-Coverage Precision Grammars. In *Proceedings of the Workshop on Grammar Engineering and Evaluation, COLING 2002*, Taipei Taiwan. Pp. 8-14.
- Chafe, Wallace L. 1976. Givenness, Contractiveness, Definiteness and Subject. In Charles Li, ed., *Subject and Topic*. New York: Academic Press.
- Chao, Y. 1968. *A Grammar of Spoken Chinese*. Berkeley: University of California Press.
- Charniak, E. and M. Johnson. 2005. Coarse-to-fine n-best parsing and MaxEnt discriminative reranking. In *Proc. ACL*.
- Ding, D. 1994. Ba-Construction as a Causitive Construction. In *Proc. NACCL*.
- Gao, Q. 1992. *Chinese Ba Construction: Its Syntax and Semantics*. Unpublished manuscript, the Ohio State University.
- Gao, Q. 2000. *Argument Structure, HPSG and Chinese Grammar*. PhD thesis, Ohio State University.
- Guo, Y., J. van Genabith and H. Wang. 2007. Treebank-based acquisition of LFG resources for Chinese. In *Proc. LFG*.
- Huang, J. 1991. Verb Movement, (In)definiteness, and the Thematic Hierachy. In *Proc. the second International symposium on Chinese Languages and Literatures*. Taipei: Academia Sinica.
- Kuno, S. 1973. *The structure of the Japanese Language*. Cambridge,MT: MIT Press.
- Levy, R. and C. D. Manning. 2003. Is it harder to parse Chinese, or the Chinese Treebank? In *Proc. ACL*.
- Lopolla R. 1990. *Grammatical Relations in Chinese: Synchronic and Diachronic Considerations*. Doctoral Dissertation, University of California at Berkeley.
- Ma, J. 1898. (reprinted in 1954). *Ma's Grammar*. Beijing: Zhonghua Books.
- McDonald, R. and F. Pereira. 2006. Online learning of approximate dependency parsing algorithms. In *Proc. EACL*.
- Miyao, Y. 2006. *From Linguistic Theory to Syntactic Analysis: Corpus-oriented Grammar Development and Feature Forest Model*. PHD thesis, University of Tokyo.
- Miyao, Y. and J. Tsujii. 2005. Probabilistic disambiguation models for wide-coverage HPSG parsing. In *Proc. ACL*.
- Sag, I., T. Wasow and E. Bender. 2003. *Syntactic Theory: A Formal Introduction*. University of Chicago Press.
- Schachter, P. 1976. The Subject in Philippine Languages: Topic, Actor, Actor-Topic or None of the Above? In Charles Li, ed., *Subject and Topic*. New York: Academic Press.
- Shibatani, M. 1990. *The Language of Japan*. Cambridge: Cambridge University Press.
- Tsao, F. 1987. A Topic-Comment Approach to the ba Construction. *Journal of Chinese Linguistics*, 15:1-55.
- Wang, L. 1957. *Zhongguo Yufa Lilun (Theory of Chinese Syntax)*. Beijing: Zhonghua Press.
- Wang, L. and H. Liu. 2007. A Description of Chinese NPs using Head-Driven Phrase Structure Grammar. *Proc. The HPSG Conference*.
- Wang, Q., D. Schuurmans, and D. Lin. 2005. Strictly lexical dependency parsing. In *Proc. 9th IWPT*.
- Wang, X. and M. Miyazaki. 2007. 文構造文法に基づく中国語構文解析 (Chinese syntactic analysis using sentence structure grammar). *Journal of Natural Language Processing*, vol.14 no.2, pp.69-93.
- Yang, S. 1995. Ba and Bei constructions in Chinese. *Journal of the Chinese Language Teachers Association* vol.30 no.3, pp.1-36.
- Zhu, D. 1982. *语法讲义 (Lectures on Grammar)*, Beijing: Comercial Press.