

# The Annotation of Event Schema in Chinese

Hongjian Zou<sup>1</sup>, Erhong Yang<sup>1</sup>, Yan Gao<sup>2</sup>, Qingqing Zeng<sup>1</sup>

<sup>1</sup>Institute of Applied Linguistics, Beijing Language and Culture University

<sup>2</sup>Zhanjiang Normal University

hongjianzou@gmail.com, yerhong@blcu.edu.cn

## Abstract

We present a strategy for revealing event schema in Chinese based on the manual annotation of texts. The overall event information is divided into three levels and events are chosen as the elementary units in annotation. Event-level annotation content and the obtaining of events patterns are explored in detail. The discourse-level annotation, annotation of relations between events and annotation of the functional attributes provide a simple way to represent event schema.

## 1 Introduction

When we want to understand a report on occurrences, we need to catch the following information: the categorization of events, the relationships between them, the participants and the attributes of the events such as polarity and modality, the attitudes towards the events and the following actions or consequences. Only the information above cannot be the precisely described. Furthermore, we need to form a schema which incorporates all of the above, that is, to compile all this information together to get the integral structure about the report.

The available annotated corpora concerning the different types of information mentioned above include: the event-annotated corpora such as ACE corpora, the corpora annotating temporal information such as TimeBank, the corpora annotating event factuality such as FactBank, the corpora annotating various types of discourse relations such as RST corpus and Penn Discourse TreeBank. Meanwhile, we lack the

annotation of event schema, which is important for providing the integral meaning of the reports.

Currently for Chinese language, the annotation of event information corpora is just beginning and still far from being sufficient, when compared with English, hence it needs further exploration.

## 2 Related Work

The work and theories concerning event schema annotation can be divided into three categories. The first kind is focused on annotation of the event argument structure, such as in ACE. The second kind is focused on annotation of the temporal information and event factuality. The last is focused on the annotation of the relations among different discourse units such as RST corpus and Penn Discourse TreeBank.

ACE(2005) is an in-depth study of research oriented annotated corpus for the purpose of textual information extraction. The annotation task includes event annotation besides the annotation of entities, values and relations between entities. The event annotation is limited to certain types and subtypes of events, that is, *Life, Movement, Transaction, Business, Conflict, Contact, Personnel*, and *Justice*. The argument structure of events including participants and other components such as time and place are predefined and tagged. Besides these, four kinds of attributes of events, *polarity, tense, genericity* and *modality*, are tagged. The expression characters of events, including the extent and the triggers, are also tagged.

TimeML(Pustejovsky et al., 2003; TimeML, 2005) is a system for representing not only all events but also temporal information. The events tagged are not limited to certain types as in ACE, but are classified in a different way. Event tokens and event instances are distin-

guished and tagged respectively. For each event instance, four kinds of attributes, namely, *tense*, *aspect*, *polarity* and *modality* are tagged. TimeML defines three kinds of links between events and times. *TLINK* represents temporal relationships, including *simultaneous*, *before* and *after*. *SLINK* represents subordinative relationships. And *ALINK* represents relationships between an aspectual event and its argument event. Several TimeML corpora have been created now, including TimeBank and AQUAINT TimeML Corpus.

FactBank(Roser and Pustejovsky, 2008, 2009; Roser, 2008) is a corpus that adds factuality information to TimeBank. The factual value of events under certain sources is represented by two kinds of attributes, *modality* and *polarity*.

Besides the annotation of events and their temporal relationships or factuality information, there are various types of discourse annotation, which can be divided into two trends: one under the guidance of a certain discourse theory(such as RST) and the one independent of any specific theory(such as PDTB).

RST (Mann and Thompson, 1987; Taboada and Mann, 2006) was originally developed as part of studies on computer-based text generation by William Mann and Sandra Thompson in 1980s. In the RST framework, the discourse structure of a text can be represented as a tree. The leaves of the tree correspond to text fragments that represent the minimal units of the discourse; the internal nodes of the tree correspond to contiguous text spans; each node is characterized by its *nuclearity* and by a rhetorical relation that holds between two or more non-overlapping, adjacent text spans. RST chooses the clause as the elementary unit of discourse. All units are also spans, and spans may be composed of more than one unit. RST relations are defined in terms of four fields: (1) Constraints on the nucleus; (2) Constraints on the satellite; (3) Constraints on the combination of the nucleus and the satellite; and (4) Effects. The number and the types of relations are not fixed. It can be reduced or extended. Carlson et al. (2003) describes the experience of developing a discourse-annotated corpus grounded in the framework of Rhetorical Structure Theory. The resulting corpus contains 385 documents selected from the Penn Treebank.

Penn Discourse TreeBank(Miltsakaki et al., 2004; Webber et al., 2005) is to annotate the million-word WSJ corpus in the Penn TreeBank with a layer of discourse information. Although the idea of annotating connectives and their arguments comes from the theoretical work on discourse connectives in the framework of lexicalized grammar, the corpus itself is not tied to any particular theory. Discourse connectives were treated as discourse-level predicates of binary discourse relations that take two abstract objects such as events, states, and propositions. The two arguments to a discourse connective were simply labeled Arg1 and Arg2.

### 3 The Levels and Elementary Unit of Event Schema Annotation

#### 3.1 The Elementary Unit of Event Schema Annotation

What counts as an elementary unit of Event Schema annotation in Chinese?

It is common to set sentences or clause as the basic units in discourse annotation such as RST corpus. However, there will be certain limitations if we choose sentences or clauses as the elementary units of Chinese event schema annotation:

First, a Chinese **sentence** is generally defined as a grammatical unit that has pauses before and after it, a certain intonation, and expresses a complete idea. But the definition is not exact or operational. The only notable borders of Chinese sentences in writings are the punctuations at the end of the sentences. The same is true of clauses in Chinese.

Second, there is generally more than one event in a sentence or a clause in Chinese. Hence, if we choose sentences or clauses as the basic units of event schema annotation, the relations between the events in one sentence/clause cannot be described in detail. For example:

1. 不到 24 小时, 俄罗斯南部黑海边一个老年人之家又燃起熊熊大火, 至少 62 人葬身火海。(In less than 24 hours, a **fire** swept through an old people's home in the Black Sea coast of southern Russia and **killed** at least 62 people.)
2. 智利南部艾森大区自 22 日以来频繁发生地震, 该地区政府已宣布该地区进入“早期警报”状态。(Earthquakes have hit the Aysen

region in southern Chile frequently since the 22nd. The government has **declared** the region to be a state of "early warning".)

In example 1, there are two events in bold type: the fire and the death in one sentence. In example 2, there are also two events in a single sentence: the earthquake and the declaration.

The "event" in this paper covers the same meaning defined by ACE(2005), which refers to "a specific occurrence involving participants".

Zou and Yang(2007) shows that an average of 2.3 times events per sentence are reported in Chinese texts and hence chose events as the basic discourse unit in their annotation. This consideration also fits the elementary unit of event schema annotation.

### 3.2 Three Levels of Event Schema Annotation

The overall event information in a report is complex and consists of different levels. In order to simplify the annotation task, we first divide the total event information into three levels, that is, the discourse level, the event level, and the entity level, choosing the event as the elementary unit of the event schema annotation.

**The event level** is defined as the level relating to atomic events. A report of occurrences always has many related events that are very easy to recognize. The events are atomic, which means the events are divided into small and minimal events. For example, when reading a report about an *earthquake* that happened in Haiti, the reader will not only know about the *earthquake* itself, but also other relating happenings such as the number of *casualty* or the following *search* and *rescue*. These things are divided into different atomic events, though they are still linked closely.

**The entity level** means the entities, times, and locations that are involved in events. For example, in "China rescues 115 from a flooded mine", "China" is the agent of the rescue; "115(miners)" are the recipients; "a flooded mine" is the location. These three entities are the arguments of the *rescue* event and should be annotated before tagging them as the arguments of the *rescue* event.

**The discourse level** is the level above the event level which creates the integral meaning of the event schema. For example, the report concerning the rescue of miners from a flooded

mine involves the *rescue*, the *coalmine accident* and possibly *injuries*. These events are linked together but have different significances within the report. So it is necessary to annotate the different significances of the events, as well as relations between events.

The following passages discuss in detail the event-level and the discourse-level annotation, while the entity-level annotation will not be discussed considering its relative simplicity.

## 4 Event-level Annotation

### 4.1 Definition of Events

ACE(2005) defines an **event** as follows: An event is a specific occurrence involving participants. An event is something that happens. An event can frequently be described as a change of state. According to ACE's definition, we define **event** as the following: An Event is an occurrence that catches somebody's attention and a change of state.

### 4.2 Obtainment of Event Patterns

The event patterns are the argument structures of certain types of events, which are the directors of argument annotation. They are extracted from large-scale texts category by category. The above categories are based on the classification of sudden events. In other words, sudden events are divided into 4 categories: *natural disasters*, *accidental disasters*, *public health incidents*, and *social security incidents*, and each category includes different types of events, for example, the *natural disasters* includes *earthquakes*, *tsunamis*, *debris flows* and so on. In dealing with a specific kind of texts, only the closely related events that appear frequently are annotated. For example, when annotating the events of *earthquake*, only *earthquake* itself and closely related events such as *loss*, *rescue*, etc, are annotated.

The event patterns are manually extracted from real texts as follows, taking *earthquake* for instance:

- A search engine is used to obtain the reports whose titles and main bodies contain the key word 'earthquake', and then manually filter out those texts whose topics are not;
- The remaining texts are then split into sentences and only the sentences that narrate an earthquake or are closely relate to the earthquake are selected;

- Specific entities in these sentences are replaced with general tags such as ‘<TIME>’, ‘<PER>’ and ‘<LOC>’ to get the patterns for earthquake type events;
- Frequently used patterns for earthquake events are extracted from the descriptions;
- The arguments of the event are numbered in sequence, and given corresponding explanations;
- The arguments are appended to event patterns when new roles are found.

The following principles should be abided by when extracting event patterns:

- Event triggers are the words or expressions that indicate existence of an event or events. If there is an event trigger in a sentence, we consider that there exists a corresponding event;
- Event triggers of different categories indicate different kinds of events;
- Some arguments of an event can be indistinct in a sentence. In other words, the different roles of the same event need to be merged into different patterns to get the complete argument structure of a certain event.

Some arguments are common roles in many events, such as *time*, *location*, and some arguments are specific to some events, such as *the magnitude*, and *the focus of an earthquake*. After the extraction of a certain amount of patterns, we can then merge the similar events. So far, we have obtained 31 categories of event patterns for 4 topics of news events.

Here is the event pattern corresponding to the *earthquake* event type extracted:

arg0	Time
arg1	Location
arg2	Magnitude
arg3	Epicenter
arg4	Focus
arg5	Focal depth
arg6	Quake-feeling locations
arg7	Frequency

Table 1. The earthquake event pattern.

### 4.3 Annotation of Types and Arguments

After obtaining the event patterns, we can annotate the types and the arguments of events according to the predefined types and patterns. If a certain event is not yet defined, the annotator should tag the event as “Other” and retag it later after obtaining the pattern of that category pro-

vided that the category is not too rare in similar reports.

The annotation of arguments consists of two steps. Firstly, we locate the entities and other expressions that belong to the arguments of a certain event. Then, we locate the roles of fixed arguments according to the corresponding event pattern. The arguments of an event are sought in the scope of the sentence in which the event trigger appears.

For example, according to the earthquake event pattern listed before, the annotation of the following sentence would be as follows:

美国地质勘探局称，这起地震发生在当地时间 12 日下午 4 时 53 分，震中位于海地首都太子港西南方向 16 公里处，震源深度为 10 公里，强度达到里氏 7.0 级。(The earthquake, with a magnitude estimated at 7.0, struck Haiti at 4:53 p.m. local time and was centered about 16 kilometers southwest of Port-au-Prince, at a depth of 10 km, the U.S. Geological Survey reported. )

arg0 Time	当地时间 12 日下午 4 时 53 分(about 4:53 p.m. local time)
arg1 Location	
arg2 Magnitude	里氏 7.0 级 (7.0)
arg3 Epicenter	海地首都太子港西南方向 16 公里处 (16 kilometers southwest of Port-au-Prince)
arg4 Focus	
arg5 Focal depth	10 公里 (10 km)
arg6 Quake feeling locations	
arg7 Frequency	1

Table 2. The annotation of the Haiti Earthquake.

### 4.4 Annotation of Event Attributes

Besides the types and arguments, the attributes of events are also tagged, which is necessary for a comprehensive description of events. Based on the analysis of various attributes in the reports, we decided to annotate the following: *Polarity*, *Modality*, *Tense*, *Aspect*, *Level*, *Frequency*, *Source*, and *Fulfillment*. Among these attributes, *Polarity*, *Modality* and *Tense* are adopted by both ACE and TimeML. *Aspect*, *Frequency* and *Source* are adopted by TimeML. The primary reason for annotating these attributes is that they have an important role in



Paris on the morning of the 10th, **killing** at least 2 women, seriously **injuring two firemen** and causing huge property **damage**.)

The Event Words “火灾”(fire), “死亡”(killing), “重伤”(injuring) and “损失”(damage) in the sentence above indicate four events respectively.

Besides annotating Event Words for events, the annotator also needs annotating indicators from texts to help to locate the attributes of the events. The attributes annotated should be clearly indicated by some linguistic hints, so the value of a certain attribute will not be specified if the hints are not so clear.

## 5 Discourse-level Annotation

The purpose of discourse level annotation is to integrate the information from the event-level into a structure. We annotate two kinds of discourse information, the relationships among events as annotated before and the functional attributes of events, to represent the event schema.

### 5.1 Annotation of Relations among Events

The events in the same report are not self-sufficient or independent, but are linked by various relationships, such as the causal relationships between an earthquake and an injury.

Taking into account of both the frequency of relationships between events and the ease and accuracy of distinguishing them, we have decided to focus on the following: *causality*, *co-reference*, *sequential*, *purpose*, *part-whole*, *juxtaposition* and *contrast*.

**Causality** is very common in reports. If event A is responsible for the happening of event B, then there exists a causal relationship between A and B. For example, in

“海地首都太子港附近发生里氏 7.0 级强烈地震，造成一家医院倒塌，另有多座政府建筑损毁。” (A magnitude 7.0 earthquake hit Haiti, causing a hospital to collapse and damaging government buildings in the capital city of Port-au-Prince.)

there are three events, called “地震”(earthquake), “倒塌”(collapsing) and “损毁”(damaging), and a causal relationship between “地震” and “倒塌”/“损毁”.

**Co-reference** is not the relationship between two different events but the relationship between two expressions of events that refer to the same object.

**Sequential** is the relation between A and B such that B follows A chronologically but there is not necessarily a causal relationship between them. For example, in

“尼日利亚南部经济中心拉各斯一名 22 岁妇女 1 月 17 日因病死亡，后经尼卫生部门检测，该妇女死于高致病性禽流感。” (A 22-year-old woman died of illness on Jan. 17 in Lagos, Nigeria's southern economic hub. After being tested by the Nigerian health sector, it was found that the woman had died of bird flu.)

the events “死亡”(death) and “检测”(testing) have sequential relationship.

**Purpose** is the relation between A and B that A happened for B. For example, in

“尼日利亚政府目前已经在全国范围内加大了卫生监管力度，以控制高致病性禽流感的扩散。” (The Nigerian government has already strengthened hygienic supervision and regulation nationwide to control the spread of the highly pathogenic avian influenza.)

the purpose of the event “监管”(supervision) is to “控制”(control).

**Part-whole** relationship between A and B is when B is part of A. For example, in

“台风“桑美”给福鼎市带来重大人员伤亡，截至目前已有 138 人死亡，其中海上遇难人员 116 人，陆上遇难人员 22 人，还有 86 人失踪。” (Saomai caused significant casualties in Fuding: at least 138 people have been killed so far, including 116 at sea, and 22 were on land, with 86 missing.)

the event “遇难”(killed) appeared first and is part of the event “伤亡”(casualties).

**Juxtaposition** relationship means that A and B are caused by the same thing, or that A and B are simultaneous. For example, in

“大同市、左云县有关部门已对被困矿工家属进行了妥善安置。同时，环境部门正在对水质进行监测。” (Datong, Zuoyun authorities have made proper arrangements for the families of trapped miners. Meanwhile, the department for environmental protection has been monitoring water quality.)

the “安置”(arrangement) and “监测”(monitoring) are simultaneous.

**Contrast** relationship is when A would usually cause B, but here A happened and didn't in fact cause B. For example, in

“萨尔瓦多中部地区 2 日发生里氏 5.3 级地震，但没有造成人员伤亡和财产损失。” (A 5.3 magnitude earthquake hit the central region of Salvador on the 2nd, but caused no casualties or property losses.)

the “地震”(earthquake) usually causes “伤亡”(casualties), but here there is no “伤亡”.

The contrast relationship between A and B is not equal to the negation of a causal relationship, because in a contrast relationship A is positive and B is negative, while in the negation of causal relationship, the A is negative.

Besides those relationships between events described above, the annotator could tag the relation as “Underspecified” if he/she feels that relationship belongs to a new kind and deserves to be annotated.

These relations are also annotated with the attributes similar to those of events, but only including **Polarity**, **Modality**, **Tense**, **Aspect** and **Source**.

## 5.2 Annotation of Functional Attributes

The annotation of relations among events only represents the local discourse structure of the report. To represent the overall information it is necessary to integrate the event-level information globally. We find that the events annotated in one text are not owning equal significance, and they can be divided into at least two basic kinds according to their role in expressing the highlight of the text. The two basic kinds of role we decide to tag are “**core**” and “**related**”. We call this the **functional attribute** of the events.

The **core events** are the events that are the topics of the reports. Other events are the **related events**. If core events were removed, the elementary topics would change and the remaining events could not easily be organized together well. For example, in a report concerning the *earthquake* that happened in Haiti several months ago, the report's core events are the events representing the *earthquake*. The other events such as the *rescue* or the *injuries* are not integral and cannot be meaningful alone. But if the other events were removed, the topic and

logic of the report would still be clear, though the details might be somewhat incomplete.

After annotating the relationships among events and functional attributes of these events, we can represent a report about an earthquake which happened in Kyrgyzstan as follow:

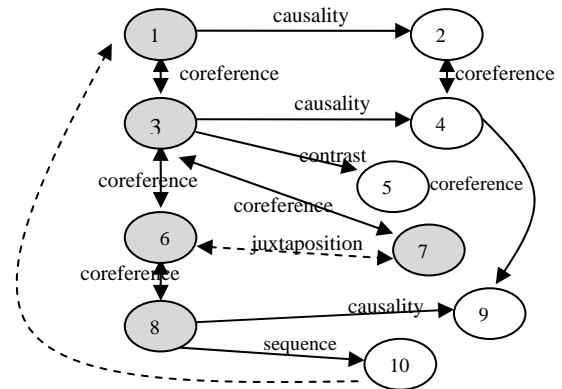


Figure 2. Event schema of Kyrgyzstan earthquake. Nodes of 1, 3, 6, 7 and 8 represent earthquakes; Nodes of 2, 4, and 9 represent damage; Node 5 represents casualty; Node 10 represents investigation.

In the graph above, the nodes represent the events, and the edges represent the relationships between events. The gray nodes represent the core events, while the white nodes represent the related events. As can be seen from the graph, the core events are at the center of the text and the related events are attached to the core events.

## 6 Preliminary Results and Discussion

In order to check the taggability of the annotation strategy mentioned above, three graduate students manually annotated about 60 news reports in 3 categories, including *earthquake*, *fire* and *attack*, using *sina* search engine, according to the method and principles above. Each text was annotated by two annotators and discussed jointly if the annotation results were inconsistent or not proper.

As can be seen from Table 3 below, 1) the event patterns extracted can cover the texts well because up to 78% sentences have been annotated. 2) There are 1.6 times more annotated events than annotated sentences. This shows that there is generally more than one event in a sentence. So, it is reasonable to assume that the annotation method can accomplish the task of a detailed description of relationships between events. 3) The relevant events are more numer-

ous than the core events. This shows that it is necessary to distinguish the core events from the relevant events.

C	T	S	NS	EV	CE	RE	AR
C1	20	277	45	361	191	170	588
C2	20	309	66	394	183	211	515
C3	20	356	93	401	121	280	605
C4	60	942	204	1156	495	661	1708

Table 3. The annotation of EVENTS

C: Sub-category; C1: earthquake; C2: fire;  
 C3: terrorist attacks; C4: total  
 T: the number of texts; S: the number of sentences  
 NS: the number of sentences not annotated  
 EV: the number of EVENTS  
 CE: the number of core EVENTS  
 RE: the number of relevant EVENTS  
 AR: the number of arguments

We have also analyzed the event attributes in detail (Zou and Yang, 2010). An interesting event attribute is Fulfillment, which is only applicable to those events with intentions whose result is often emphasized. Sometimes, readers care about the intended results or outcomes as much as or more than the events themselves. Therefore it would be useful to explore the notion of Fulfillment, and investigate which linguistic categories could play a role in deciding the value of Fulfillment. We plan to create a Fulfillment corpus in the next stage.

The annotation of event schema is time-consuming, partly because it needs to annotate all three levels of event information of every text, and partly because of the difficulties to identify the event information from trivial descriptions, in other words, one question we often discuss is whether it deserves to annotate certain parts of a text. Also, we often need to make a balance between obtaining enough event patterns to cover various types of related events well and omitting low frequent event types to simply the obtainment of event patterns. In discourse-level annotation, the main difficulty is the identification of relations between events without lexical hints. This discourse-level annotation is only just underway. We also plan to give detailed analysis in the next stage.

**Acknowledgements.** This paper is sponsored by National Philosophy and Social Sciences Fund projects of China (No. 06YY047).

## References

- ACE. 2005. *ACE (Automatic Content Extraction) Chinese Annotation Guidelines for Events*. [http://www ldc.upenn.edu/Projects/ACE/docs/Chinese-Events-Guidelines\\_v5.5.1.pdf](http://www ldc.upenn.edu/Projects/ACE/docs/Chinese-Events-Guidelines_v5.5.1.pdf)
- Carlson L., D. Marcu, M. E. Okurowski. 2003. *Building a Discourse-Tagged Corpus in the Framework of Rhetorical Structure Theory*. Current Directions in Discourse and Dialogue, Jan van Kuppevelt and Ronnie Smith eds., Kluwer Academic Publishers.
- Mann W. and S. Thompson, 1987. *Rhetorical Structure Theory: A Theory of Text Organization* (No. ISI/RS-87-190). Marina del Rey, CA, Information Sciences Institute.
- Miltsakaki E., R. Prasad, A. Joshi, and B. Webber. 2004. *Annotating Discourse Connectives and their Arguments*. Proceedings of the HLT/NAACL Workshop on Frontiers in Corpus Annotation. Boston, MA.
- Pustejovsky J., J. Castaño, R. Ingria, S. Roser R. Gaizauskas, A. Setzer and G. Katz. 2003. *TimeML: Robust Specification of Event and Temporal Expressions in Text*. Fifth International Workshop on Computational Semantics.
- Taboada M. and W. Mann. 2006. *Rhetorical Structure Theory: Looking Back and Moving Ahead*. Discourse Studies 8(3): 423-459.
- TimeML. 2005. *Annotation Guidelines Version 1.2*. [http://www.timeml.org/site/publications/timeMLdocs/anguide\\_1.2.1.pdf](http://www.timeml.org/site/publications/timeMLdocs/anguide_1.2.1.pdf).
- Webber B., A. Joshi, E. Miltsakaki, et al. 2005. *A Short Introduction to the Penn Discourse Tree-Bank*. Copenhagen Working Papers in Language and Speech Processing.
- Roser S. 2008. *A Factuality Profiler for Eventualities in Text*. Ph.D. Thesis. Brandeis University.
- Roser S. and J. Pustejovsky. 2008. *From Structure to Interpretation: A Double-layered Annotation for Event Factuality*. Proceedings of the 2nd Linguistic Annotation Workshop.
- Roser S. and J. Pustejovsky. 2009. *FactBank: A Corpus Annotated with Event Factuality*. Language Resources and Evaluation.
- Zou H.J. and E.H. Yang. 2007. *Event Counts as Elementary Unit in Discourse Annotation*. International Conference on Chinese Computing 2007.
- Zou H.J. and E.H. Yang. 2010. *Annotation of Event Attributes*. The 11th Chinese Lexical Semantics Workshop.