

Annotating “tense” in a tense-less language

Nianwen Xue*, Hua Zhong† and Kai-Yun Chen‡

University of Colorado
The Center for Computational Language and Education Research
Boulder, CO 80309-0594, USA
Nianwen.Xuecolorado.edu*, Hua.Zhongcolorado.edu†, Kaiyun.Chencolorado.edu‡

Abstract

In the context of Natural Language Processing, annotation is about recovering implicit information that is useful for natural language applications. In this paper we describe a “tense” annotation task for Chinese, a language that does not have grammatical tense, that is designed to infer the temporal location of a situation in relation to the temporal deixis, the moment of speech. If successful, this would be a highly rewarding endeavor as it has application in many natural language systems. Our preliminary experiments show that while this is a very challenging annotation task for which high annotation consistency is very difficult but not impossible to achieve. We show that guidelines that provide a conceptually intuitive framework will be crucial to the success of this annotation effort.

1. Introduction

The notion of tense is a subject of lively debate in recent theoretical Chinese linguistics literature. Much of the controversy surrounds the question of whether tense is a useful theoretical construct. Huang (1984) and Li (1990) suggest there is a finite-nonfinite distinction in Chinese syntax which implies the existence of a phonologically empty tense category in the syntactic representation of a clause. Other researchers (Hu et al., 2001; Lin, 2003; Lin, 2006; Smith and Erbaugh, 2005) contend that such an abstract notion of tense is not justified for Chinese. In spite of this controversy, there is a general agreement that Chinese does not have grammatical tense. That is, Chinese does not have tense morphemes that are dedicated to mark the temporal location of situations, like English. However, Chinese speakers, like speakers of other world languages, use the moment of speech as the temporal deixis to temporally locate situations in their communication. One piece of evidence for this is that temporal adverbials like 今天 (“today”), 明天 (“tomorrow”) and 昨天 (“yesterday”) all assume a temporal deixis that is the moment of speech.

Corpus annotation provides a new perspective from which we can look into how Chinese speakers interpret situations temporally. Through linguistic annotation, we can determine whether Chinese speakers consistently identify the temporal locations of situations. To the extent that they do, this information would be highly useful for natural language processing systems that can learn from the human interpretations through statistical and machine learning approaches. For example, this information would be highly valuable to Machine Translation. To translate a language like Chinese into a language like English in which tense is grammatically marked with inflectional morphemes, an MT system will have to infer the necessary temporal information to determine the correct tense for verbs. Statistical MT systems, the currently dominant research paradigm, typically do not address this issue directly. As a result, when evaluated for tense, current MT systems often perform miserably. For example, when a simple sentence like “他/he 明天/tomorrow 返回/return 上海/Shanghai” is given to

Google’s state-of-the-art Machine Translation system¹, it produces the output “He returned to Shanghai tomorrow”, instead of the correct “he will return to Shanghai tomorrow”. The past tense on the verb “returned” contradicts the temporal expression “tomorrow”. Ye et al (2006), for example, reported that the best MT systems only get the tense correct 47% of the time, worse than the 69.5% when the most frequent tense (which is the past tense) is assigned. In this paper we report a preliminary study in which we ask Chinese native speakers to temporally locate situations denoted by verbs in naturally occurring text and assign a label that indicates their temporal relations to the document creation time.

2. Deciding on a tag set

The three temporal points, Speech Time, Reference Time and Situation Time are generally accepted as important to temporal interpretation since Reichenbach (1947) first proposed them. Speech Time and Reference Time are important to the interpretation of tense. Speech Time is the moment of speech, and Reference time is the temporal perspective from which the speaker invites his audience to consider a situation, and Situation Time is the time at which the situation actually occurs. In the “tense” annotation study we report here, we annotate the relation between Speech Time and Situation Time instead of the relation between Speech Time and Reference time, which in our judgment is too subtle to be annotated consistently. In determining the different values of tense, one important question to ask is whether one should use the tense taxonomy of a particular language like English, or use a tag set that abstracts away from the language-specific properties. In a previous pilot tense annotation study motivated by Chinese-English machine translation, Ye (2007) used a tagset that is a direct reflection of the English tense system. She used a set of 11 tags that represent different combinations of the three tenses (present, past and future) and aspect (perfect, progressive). However, there are reasons for not borrowing the English tense taxonomy directly. First, the English system does not consistently encode the relation between Speech

¹http://www.google.com/language_tools

time and Situation time. For example, in “He will call me after he gets here”, while his “getting here” happens at a time in the future, it is assigned the present tense because it is in a clause introduced by “after”. The annotation of the Chinese tense should not be bound by such English idiosyncrasies. Instead it should annotate temporal notions that are more universal. Similarly, in English one can either use an infinitive marker or the future tense to indicate that the temporal location of a situation is in the future. For example, in “He hopes to leave”, one can infer that leaving is something in the future, much like the leaving in “he hopes that she will leave”, where the future tense is explicitly marked. In a language like Chinese that lacks overt tense markers, if the task is to annotate tense as it is used in English, the choice will be indeterminate in such situations. It makes more sense to annotate the abstract relation between Speech Time and Situation Time which is more deterministic. Adopting the English tense taxonomy also limits the pool of potential annotators, who would have to be proficient in both Chinese and English. Such annotators are generally hard to find.

In written text, which is the primary source of data that we are dealing with, the temporal deixis is the document creation time. All situations are temporally related to this document creation time except in direct quotations, where the temporal location is relative to the moment of speech of the speaker who is quoted. For example, in (1), there are several verbs involved. The temporal location of 说 (“say”) assumes that the temporal deixis is the document creation time, while 是 (“be”), 愿意 (“want”) are inside the direct quotations of 索康, the speaker, and their temporal location is relative to the time when he performed that speech act.

- (1) 索康 曾 对他的 母亲 说：“他是我的
Suokang once to he DE mother say：“he be I DE
财产，我愿意把他揉成团，装在口袋
property, I want BA he roll into ball, put in pocket
里，随我；...”
inside, up to me, ...”

”Suokang once said to his mother: ‘He is my property. If I want to roll him into a ball and put it in my pocket, it’s up to me, ...’ ”

Tenses that relate the situation time to the document creation time or the moment of speech are considered to be *absolute tenses*, even though the moment of speech is variable and constantly changes. The term *situation* covers events and states in the broad sense of those terms and they are generally signaled by the presence of a verb. In other cases, the temporal location of a situation cannot be defined in relation to the moment of speech, at least not directly. For example in (2), The temporal location of 有意 (“intend”) cannot be determined independently of the temporal location of 透露 (“reveal”). The temporal location of 有意 is simultaneous with 透露. If the temporal location of 透露 is in the past, then the temporal location of 有意 is also in the past. If the temporal location of 透露 is in the future, then the temporal location of 有意 is also in the future. In this specific case, the situation denoted by the matrix verb 透露 is in the past. Therefore the situation denoted by 有意 is also located in the past.

- (2) 他还 透露 俄罗斯有意 在今后十年 内，
he also reveal Russia intend in next ten years within，
向伊朗提供 武器。
to Iran provide weapons .

“He also revealed that Russia intended to provide weapons to Iran within the next ten years.”

In our Chinese “tense” annotation task, we annotate both *absolute* and *relative* tenses. By “tense”, we mean the temporal relation between the temporal deixis and the situation (in the case of absolute tense) and between situations (in the case of relative tense). The situation time can be anterior to (in the past), simultaneous with (in the present), or posterior to (in the future) the moment of speech. Since there is no grammatical tense that can be directly observed in Chinese text, the tense annotation as defined here is a temporal inference exercise. The “situations” that we are interested in are expressed as clauses centered around a verb, and for the sake of convenience we mark the “tense” on the verb itself instead of the entire clause. However, when inferring the temporal location of the situation, we have to take into consideration the entire clause, because the arguments and modifiers of a verb are just as important as the verb itself when determining the temporal location of the situation. For example, in the absence of additional information such as temporal adverbials, the default temporal location of (3a) is the present while the default temporal location of (3b) is the past, even though the verb is the same in these two clauses. (3a) describes a habitual situation where the person described sleeps a lot repeatedly, even though he is not in a perpetual state of sleep. The habitual interpretation is induced by the adverbial modifier 整天 (“all the time”). In contrast, the duration of three hours quantifies the sleep situation in (3b) and establishes temporal bound for 睡 (“sleep”). By default bounded situations are temporally located in the past (Smith and Erbaugh, 2005), a point we will elaborate further.

- (3) a. 他整天 睡 大觉。
he all the time sleep sleep .
“He sleeps all the time.”
b. 他睡 了 三 小时。
he sleep ASP three hour .
“He slept for three hours.”

3. Specifications

3.1. Inferring Absolute “Tense”

Our annotation task is inspired by the theoretical work of Smith (2005), in which she shows that for both languages with or without grammatical tense, there is a default pattern of temporal interpretation, which is that unbounded situations are located in the Present while bounded situations are located in the past by default. There is also a bounded event constraint which says that bounded situations are not located in the Present. In the absence of explicit temporal expressions that can overwrite this default pattern of temporal interpretation, *boundedness* is thus a crucial semantic concept when inferring the temporal location of a situation. There are many different ways in which a situation

can be bounded (or unbounded, for that matter). Boundedness can be a function of the inherent properties of some verbs. For example, punctual verbs are usually bounded while stative verbs are usually unbounded. Boundedness can also come from the compositional properties of some syntactic constructions. For example, resultative constructions tend to be bounded. In addition, aspectual viewpoints are also good indications of boundedness, when they appear in text. For example, perfective viewpoints are compatible with bounded situations while imperfective viewpoints are compatible with unbounded situations. We do not have space to go into a detailed discussion of boundedness and its relationship to situation types and viewpoints here, and the reader is referred to (Xiao and McEnery, 2004) for a more detailed discussion. In this section, we will describe each temporal location when it is a default or non-default interpretation of a situation.

3.1.1. Present tense

A situation is assigned the present tense if it is true at an interval of time that includes the present moment. Bounded situations is incompatible with the present tense because it can only be viewed “at a distance” and thus cannot be true at the present moment. They can only be in the past or in the future. The present tense is compatible with states and activities. States are temporally located in the present by default. When non-stative situations are temporally located in the present, they either have an imperfective aspect or have a habitual or frequentive reading which makes them look like states, e.g.,

- (4) 他常常参加户外活动。
he often attend outdoors activities .
“He often attends outdoors activities.”

3.1.2. Past tense

Situations that happen before the moment of speech (or the document creation time) are temporally located in the past. Bounded situations are temporally located in the past by default (5a). States and unbounded activities are temporally located in the present by default (5b), but they can be located in the past when they are modified by explicit temporal expressions indicating the past (5c).

- (5) a. 中方人员及侨胞安全撤离乍得。
Chinese personnel and Chinese nationals safely withdraw from Chad .
“Chinese personnel and Chinese nationals safely withdrew from Chad.”
- b. 他喜欢看功夫片。
he like watch kung-fun movie .
“He likes to watch kung-fu movies.”
- c. 他小时候喜欢看功夫片。
he young like watch kung-fu movie .
“He liked to watching kung-fu movies when he was young.”

3.1.3. Future tense

Situations that happen posterior to the moment of speech are temporally located in the future. Future situations are not simply the opposite of past situations. While past situations have already happened by definition, future situations by nature are characterized by uncertainty. That is, future situations may or may not happen. They really just refer to situations that have not happened yet. Therefore, future situations are often linked to possibilities, not just to situations that will definitely happen. No situation type is associated with the future by default, and therefore a situation is temporally located in the future only if it is modified by an explicit temporal expression (6a), or if it is modified by a modal verb (6b) (not all modal verbs are future-oriented), or if the situation is embedded in a conditional clause, as indicated by a subordinate conjunction (6c) or by a noun denoting condition (6d), or if the situation is embedded in a future-oriented verb (6e) or noun (6f), in which case its temporal location is sensitive to the temporal location of the embedding situation (also see Section 3.2.).

- (6) a. 大会明年在新加坡举行。
conference next year in Singapore hold .
“The conference will be held in Singapore next year.”
- b. 法官在判决书中说：“他们如果获释，很可能立即潜逃。”
judge in verdict in say : “ They if release , very likely will immediately escape . ”
“The judge said in the verdict: ‘if they are released, they are very likely to escape immediately.’ ”
- c. 如果油价持续上涨...
if oil price continue rise ...
“If oil prices continue to rise...”
- d. 制定这样一项法律是他1996年在同墨西哥政府谈判时提出的条件之一。
establish like this one CL law be he 1996 in with Mexican government negotiate when put forward DE condition one of .
“Establishing such a law was one on the conditions that he put forward when he negotiated with the Mexican government.”
- e. 美国政府决心将每一位在海外战事中捐躯的官兵遗骸运回美国。
U.S. government determine BA every one CL in oversea battle in die DE troop remains bring back the U.S. .
“The U.S. government is determined to bring back to the U.S. the remains of every troop who died in a battle overseas.”
- f. 这说明他们有随时潜逃的准备。
this indicate they have at any time escape DE preparation .

“This shows that they have made preparations to escape at any time.”

Modal verbs in Chinese include 会 (“will”), 要 (“should, will”), 将 (“will”), 应该 (“should”), 可能 (“may”), 可以 (“may”), 必须 (“must”), etc.). Many modal verbs indicate futurity. The verb complement of these modals tend to be temporally located in the future, that is, they indicate situations that have not yet happened. For example, in (7), the verb 合作 (“cooperate”) is temporally located in the future because of the modal 要. Note that in some cases 要 denotes intension and is not necessarily a modal verb in all context. Even when 要 is not a modal verb, it is still a future-oriented verb. So situations modified by 要 always has a sense of futurity.

- (7) 双方 表示 要 进一步合作。
two sides indicate will further cooperate.

“The two sides indicate that they will cooperate further.”

Not modals are future-oriented. For example, the verb complement of 能 (“can”) does not indicate futurity. The default tense is the present unless overwritten by an overt temporal adverbial:

- (8) 他能 做好 这件工作。
he can do-well this CL job.

“He can do this job well.”

Future-oriented verbs include verbs of future having (e.g., 保证 (“guarantee”), 承诺 (“pledge”), verbs of future situation (e.g., 提议 (“propose”), 计划 (“plan”), 准备 (“prepare”), etc.), verbs of wish and desire (e.g., 要 (“want”), 妄想 (“dream”), 需要 (“need”), 向往 (“look forward to”), verbs of future events (e.g., 预告 (“foretell”), 想 (“want”), 预言 (“predict”) and verbs of future prevention (e.g., 排除 (“preclude”), 阻止 (“prevent”). Future-oriented nouns are often derived from these future-oriented verbs and they often co-occur with a support verb. For example, 准备 (“prepare”) is a future-oriented verb (9a) and its complement, 撤离 (“withdraw”) is located in the future relative to the temporal location of 准备. Even when the 准备 (“preparation) is nominalized, as in (9b), the situation denoted by its complement clause is still located in the future relation to this nominalized verb.

- (9) a. 科索沃独立 可能引发 骚乱, 联合国
Kosovo independence may cause riot . UN
人员 已 准备 撤离。
personnel already prepare withdraw .

“Kosovo independence may cause riot. UN personnel have already prepared to leave.”

- b. 科索沃独立可能引发骚乱联合国人员
已做好[撤离]的准备。

3.2. Inferring Relative tense

The temporal interpretation of one situation is often bound by the temporal location of a reference situation. One common scenario in which this kind of dependency occurs is

when the target situation, the situation we are interested in at the moment, is embedded in the reference situation as a complement. Just as the absolute “tense” represents a temporal relation between the situation time and the moment of speech or document creation time, the relative “tense” represents a relation between the temporal location of a situation and its reference situation. The target situation can be anterior to (10a), simultaneous with (10b), or posterior to the reference situation (10c).

- (10) a. 印尼 警方 透露 前总统 苏哈托
Indonesian police reveal ex-president Suharto
已经 死亡。
already die .

“The Indonesian police revealed that ex-president Suharto had already died.”

- b. 漫游费 听证 代表 透露
roaming charges testimony witness reveal
手机 漫游 成本不到 5 分
mobile phone roaming cost less than five fen
钱。
money .

“The witnesses at the testimony on roaming charges revealed that the cost for roaming was less than 0.05 Chinese yuan.”

- c. 公司 员工 透露 《星际2》测试
company personnel reveal “ Star 2 ” trial
版 即将面世。
version soon face the world .

“The company personnel revealed that ‘Star 2’ trial version would soon face the world.”

Relative “tense” is invoked when the absolute (in relation to moment of speech) temporal location of a situation can only be determined through the temporal location of a reference situation. For example, in (10c), the temporal location of 面世 (“face the world”) is posterior to the reference situation 透露 (“reveal”), and thus the absolute temporal location of 面世 depends on the temporal location of 透露. If the reference situation is located in the present, the target situation is located in the future. If the reference is located in the past, then the temporal location of the target situation is a future-in-past. If the reference situation is located in the future, then the temporal location of the target situation is the future of future. The combination of absolute “tense” of the reference situation and the relative temporal ordering of the reference and target situations yields nine possibilities, listed in Table 1.

ref/tgt	anterior	simultaneous	posterior
present	past	present	future
past	past	past	fip
future	pff	future	future

Table 1: Nine possibilities (Acronyms: ref=reference, tgt=target, pff=past-from-future, fip=future-in-past)

Among the nine possibilities, past-from-future is rarely attested. When the reference situation is located in the

present, the temporal location of the target situation is relative to the moment of speech or document creation time, and thus are not different from the absolute “tense”. When the reference situation is temporally located in the past and the target situation is anterior to or simultaneous with the reference situation, we simply stipulate that the target situation is temporally located in the past without making any further distinctions. The only possibility we explicitly encode is future-in-the-past, which occurs very often in actual text.

3.3. Target situation is posterior to reference situation

One scenario where the target situation is posterior to the reference situation is when the reference situation is denoted by a future-oriented verb. For example, in (11a), the situation denoted by 阻止 (“prevent”) is embedded in the reference situation denoted by 意图 (“intend”), which is a future-oriented verb. 阻止 is something that has not yet happened at the time of 意图 and is thus posterior to 意图. 意图 is a stative verb, and by default is located in the present. The temporal location of 阻止 is thus a simple future. Similarly, in (11b), 乱砍滥伐 (“irresponsibly cut trees”) is also located in the future because of the future-oriented verb 许 (“allow”).

- (11) a. 俄罗斯意图 阻止 美国 独霸太空
Russia want prevent the U. S.
。
monopolize outer space .
“Russian want to prevent the U.S. from monopolizing the outer space”
- b. 二 不 许 乱砍滥伐
second not allow cut trees irresponsibly
“Second, irresponsible logging is not allowed.”

When the embedding future-oriented verb is temporally located in the past, the tense of the complement verb is a future-in-past. For example, in (12), the temporal location of the embedding verb 准备 (“get ready to”) is in the past because it is in the scope of the temporal adverbial 昨天 (“yesterday”). The tense of the complement verb 去 (“go”) is thus a future-in-past.

- (12) 昨天 , 他正 准备 去 机场 时 ,
yesterday , he right prepare go airport when ,
“Yesterday, right when he was preparing to go to the airport...”

3.4. Target situation is simultaneous with reference situation

By “simultaneous” we mean there is some overlap between the intervals that the two situations hold true. For example, in (13a), the temporal interpretation of 寻找 (“seek”) is dependent on the temporal location of 乐于 (“like”). Since 乐于 is a stative verb, by default we assume that the situation it denotes is temporally located in the present. Since 寻找 is simultaneous with 乐于, it is also located in the present. In (13b), the situations denoted by 穿上 (“put on”) are simultaneous with the situation denoted by 看

到 (“see”). Since 看到 is temporally located in the past, the situation denoted by 穿上 is also temporally located in the past.

- (13) a. 不少 中国人 乐于 在 工资 收入 较
many Chinese like in salary income relatively
高 的 外企 寻找 就业
high DE foreign enterprises seek employment
机会 .
opportunity .

“Many Chinese like to seek employment opportunities in foreign enterprises where the salary income is higher.”

- b. 记者 今天 在 火车站 广场 看到 ,
reporter today at railway station square see ,
便衣 警察 也 全部 穿上 制服
plain-clothes police also all put on uniform
参与 维持 秩序 .
participate maintain order .

“This reporter saw at the railway station square today that all plain-clothes policemen also put on uniforms and participated in maintaining order.”

3.5. Target situation is anterior to reference situation

A third situation is when the target situation is temporally anterior to the reference situation. Verbs such as 庆幸 (“feel fortunate”), 后悔 (“regret”) are opposite of future-oriented verbs in that they imply something has already happened. For example, in (14), the temporal location of 要 (“ask for”) is anterior to 后悔 (“regret”). Since 后悔 is a stative verb and by default is located in the present, the situation denoted by 要 occurs before the moment of speech or document creation time. Thus it should have a simple past tense.

- (14) 陕西 官员 后悔 未 向 拍摄者 多
Shanxi official regret not from photographer more
要 几 张 华南虎 照片 .
ask for a few CL South China tiger picture .

“Shanxi officials regret that they didn’t ask the photographer for a few more pictures of the South China tiger.”

4. Experiments

We performed two experiments with two annotators that are native speakers of Chinese who are also fluent in English. In the first experiment, we replicated the annotation condition reported in (Ye, 2007) under which the annotators were asked to choose a tense they would use to translate a Chinese verb into English without providing them with a set of guidelines. The experiment was done on 1076 verb instances from 6 articles randomly selected from the Chinese Treebank. The annotators were presented with each verb in its context in the article. The results are presented in Table 2. The agreement is higher than the 43% is reported in (Ye, 2007), but on a reduced tag set that is only relevant to tense. In the second experiment, we provided the annotators with a simple set of guidelines that are summarized in Section 3.. The guidelines summarize factors that the annotators

Tag	Ann1	Ann2	Agree	f-measure
Present	386	269	185	.565
Past	302	269	216	.757
Future	123	111	52	.444
N/A	265	427	174	.503
Overall	1076	1076	627	.583

Table 2: Inter-annotator agreement without guidelines

should consider when making their judgment on the tense selection. The guidelines are presented as a set of rules that can be overwritten when “hard evidence” such as temporal adverbials are present. The second experiment was done on a different set of verbs taken from six different articles from the first experiment. The results (Table 3) show that even with this set of very simple guidelines, there has been substantial improvement (17%) in inter-annotator agreement, from the 58% in Experiment 1 to the 75% in Experiment 2. One interesting observation is that the past tense is annotated with very high inter-annotator agreement (86%) while the future and present tenses are much harder to annotate. This is something that we will try to explain in future work.

Tag	Ann1	Ann2	Agree	f-measure
Present	276	123	106	.531
Past	501	480	424	.864
Future	180	320	171	.684
N/A	134	168	122	.61
Overall	1091	1091	823	.754

Table 3: Inter-annotator agreement with guidelines

5. Conclusion and future work

In the context of natural language processing, annotation is about recovering implicit information that is useful for natural language applications. In this sense annotation of Chinese tense would be a highly rewarding endeavor if successful. Our preliminary experiments show that providing a set of intuitive of guidelines is crucial to Chinese tense annotation. Although the inter-annotator agreement we report here has not yet reached satisfactory levels, we are confident that with refined guidelines the inter-annotator agreement can be further improved.

6. References

- Jianhua Hu, Haihua Pan, and Liejiong Xu. 2001. Is there a finite-nonfinite distinction in Chinese. *Linguistics*, 39:1117–1148.
- James C.T. Huang. 1984. On the distribution and reference of empty pronouns. *Linguistics Inquiry*, 15:531–574.
- Audrey Yen-hui Li. 1990. *Order and Constituency in Mandarin Chinese*. Dordrecht: Kluwer Academic Publishers.
- Jo-Wang Lin. 2003. Temporal Reference in Mandarin Chinese. *Journal of East Asian Linguistics*, 9:259–311.
- Jo-Wang Lin. 2006. Time in a language without tense: The case of Chinese. *Journal of Semantics*, 23:1–53.

- Carlota S. Smith and Mary Erbaugh. 2005. Temporal interpretation in Mandarin Chinese. *Linguistics*, 43(4):713–756.
- Z. Xiao and A. McEnery. 2004. *Aspect in Mandarin Chinese: A corpus-based study*. Amsterdam, John Benjamins.
- Yang Ye, Victoria Li Fossum, and Steven Abney. 2006. Latent features in automatic tense translation between Chinese and English. In *The Proceedings of the 5th SIGHAN Workshop on Chinese Language Processing*, Sydney, Australia.
- Yang Ye. 2007. *Automatica Tense and Aspect Translation between Chinese and English*. Ph.D. thesis, University of Michigan.