# Comprehensive annotation of cross-linguistic variation in the category of tense

Zymla, Mark-Matthias
University of Konstanz
Mark-Matthias.Zymla@uni-konstanz.de

### Abstract

In this paper, we present part of a new, cross-linguistically valid annotation scheme for annotating tense and aspect information on a syntactic and semantic level, focussing on the category of tense. Primarily, the annotation maps morphosyntactic information to representations of eventualities. We specify mapping conventions which are represented as inference rules expressing language specific variations of the syntax/semantics interface. Eventualities are expressed in terms of a cluster of features whose values can each be mapped to a formal description based on insights from the tense and aspect semantics literature. The annotation is integrated into a broader effort of achieving cross-linguistically viable temporal annotation and combines recent efforts of bringing together computational and formal approaches to temporal semantics. This allows for an overall comprehensive representation of syntax, semantics and the syntax/semantics interface regarding tense and aspect. The annotation scheme is especially well suited for computational research seeking to understand and extract tense/aspect information across languages.

## 1   Introduction

Annotation of temporal information as a whole has made steady progress in recent years following the TimeML ISO standard (Derczynski et al., 2013; Pustejovsky et al., 2003, 2002) and going beyond it (Bethard and Parker, 2016; Gast et al., 2016, 2015). However, there are only limited possibilities readily available to use this information to gain a better understanding of the relation between explicitly stated temporal information such as temporal expressions and and the category of tense. This has two implications: First, the fine nuances in meaning encoded in the interplay between verbal form and meaning cannot be captured, and second, existing efforts cannot provide a cross-linguistically adequate representation of tense and aspect categories. In this paper we address these issues by providing a tripartite annotation of tense. Thereby, we bring together deep linguistic parsing and a novel semantic annotation of eventualities. The system may be supplemented by elements from existing temporal annotation schemes ultimately providing a comprehensive description for temporal information and its relation with tense syntax and semantics.

Current temporal annotation schemes are well equipped with expressive power to describe the overall temporal situation within a text or sentence. However, the current state of the art focuses more on the annotation of explicit temporal expressions rather than on the meaning that remains implicit in the morphosyntactic form of verbal predicates. Thus, a great deal of the information with regard to the mapping from syntax to semantics is difficult to access. For example, Chinese timeML annotations in the 2010 TempEval shared tasks (UzZaman et al., 2012) are not marked up for tense and aspect at all, although, several morphosyntactic as well as semantic and pragmatic factors for analyzing tenses have been identified in the formal literature. This is shown for example in Bittner (2014); Smith (2006), among many others.

This makes it difficult to use existing efforts to get a better understanding of the interactions between temporal structure as a whole and different instantiations of tense and aspect categories cross-linguistically. Following approaches such as Gast et al. (2015) or Bethard and Parker (2016), we propose an annotation scheme that is faithful to the interaction between syntax and semantics in this paper.

The paper is structured as follows: We begin with identifying the main issues we want to address wrt. the annotation of tense. Section 3 illustrates the syntactic and semantic annotation of tense features and their realization in terms of the semantics we propose for eventualities. We illustrate the annotation in terms of the concrete annotation of the category of tense. In 4 we apply the entire annotation scheme to data stemming from the formal semantic literature illustrating how the scheme deals with cross-linguistic variation. Section 5 wraps up the paper with some remarks on implementability of the proposed annotation scheme.

## 2   Research questions

In this paper we focus on answering two questions: First, how to map morphosyntax to meaning in an implementable way, and second, which properties of eventualities are required (and of interest) for a linguistically rich representation of tense and aspect features. To address the first question, we propose separate levels for syntactic vs. semantic annotation, complemented by an alignment system that links syntactic features to semantic features by applying a set of inference rules. Such a system allows us to maintain parallelism in terms of syntactic and semantic features of languages while still explicitly capturing cross-linguistic variation. In terms of the abstract semantic foundations, the annotation is independent of a syntactic system. Nevertheless, our approach entails that the richer the syntactic representation, the more detailed the interactions between syntax and semantics can be modeled.

The alignment system is based on two different linking operations between syntactic and semantic features: a compatibility relation and an implication relation, such that for example the following simple rules are possible:

(1)  a.  The people killed    the king.
         the  people kill.**Past** the king

   b.  logoN=ne   baadshaah=ko maar daalaa.
       people=Erg king=Acc       hit   **put.Perf**

   c.  Das Volk    hat           den König getötet.
       the  people have.Pres.3.Sg the  king  **kill.Perf**

   d.  Orang membunuh raja
       people **AV.kill**     king

(2)  a.  *indicative mood ∧ syntactic past*
         *→ semantic past*

   b.  *indicative mood ∧ syntactic perfect*
       *∘ semantic past*

   c.  *indicative mood ∧ syntactic present ∧*
       *syntactic perfect → semantic past*

   d.  *indicative mood ∘ semantic past*

Example (1) and (2) illustrate how tense might be annotated in four different languages. The semantic interpretation might be a obligatory inference of a syntactic tense marker as in (2a) (English), a non-obligatory byproduct of a syntactic marker that does not necessarily express semantic past tense as in (2b) (Urdu), an obligatory inference of a combination of a syntactic tense with a syntactic perfect marker as in (2c) (German), or simply not inferable from the syntactic structure but merely a feature that overlaps with it as in (2d) (Indonesian). [1]  Overall, the semantics in (2) are parallel; mapping the syntax to the semantics, however, captures cross-linguistic variation, even in tense systems that are more on par with one another than the ones described above.

The second aim of this paper is to define different types of variation in the syntax/semantics interface across languages and to illustrate how this variation is encoded in our annotation scheme. We introduce a crucial distinction of different form-to-meaning mappings that allows us to identify primary and secondary meaning features of certain syntactic constructions. Primary meaning features are those that follow directly from morphosyntactic grammaticalization (and are thus readily available from the syntactic analysis), while secondary meaning features describe semantic properties that are based on complex semantic and pragmatic processes, such as implications or alternative, semantically-constructed meanings of morphosyntactic material. Thus, the system presented here allows us to explicitly capture variations of the form-to-meaning mapping within a language, such as for example the so-called Sequence-of-tense phenomenon. Sequence-of-tense (hence: SOT) describes an unexpected pattern in the interpretation of

---

[1]In cases such as the latter, the semantic interpretation is usually contextual, a problem that is discussed in depth in Zymla (to appear)

past tense morphology in embedded contexts. This is illustrated in (3). SOT sentences result in two different possible paraphrases, undermining the assumption that syntactic past tense is always synonymous with a temporal back-shift of the respective eventuality (denoted in (3) by the two predicates E1 and E2).[2]

(3)  Tom $\underline{said}_{E1}$ that Karen $\underline{was\ sick}_{E2}$.
   a.  Tom said: Karen is sick.
   b.  Tom said: Karen was sick.

Furthermore the annotation presented here allows us to capture cross-linguistic variation of tense categories. In the upcoming sections, we discuss syntactic variation of tense categories across languages such as the ones already introduced in (2). Moreover, we present a case of cross-linguistic semantic variation based on the SOT-phenomenon sketched above. Concretely, we illustrate the explicit annotation of syntax/semantics interface processes that vary between languages that express the SOT phenomenon on the one hand and languages that do not express the SOT phenomenon on the other hand. By explicitly annotating syntax, semantics and the interface, we can capture formal linguistic insights that at best remain implicit in existing annotation schemes.

# 3   Syntactic and semantic annotation of tense and aspect

The main innovation of our annotation system are the semantic features representing various properties of (verbally expressed) eventualities. The term eventuality thereby covers both events as well as states, following the classical ontology of Bach (1986). We treat eventualities as semantic objects similar to event variables in event semantics with certain semantic features. Relating these features to a sufficiently rich syntactic representation is the important first step for a comprehensive annotation of tense (and aspect). For this purpose we make use of the computational grammars developed within the ParGram project (Butt et al., 2002; Sulger et al., 2013) as syntactic input. This is done for two reasons: first, parsers are already available for a wide number of languages (for a full list and more description, see Sulger et al. (2013)) and second, the syntactic annotation of tense and aspect features is sufficiently exhaustive and more importantly parallel across languages. This allows us to clearly define parameters of syntactic variation.

The ParGram grammars are all developed using the XLE parser (Crouch et al., 2017) and are couched within the syntactic theory of Lexical Functional Grammar (LFG, Dalrymple (2001); Bresnan (2001)) which describes syntax in terms of two distinct representations: the c(onstituent)-structure and the f(unctional)-structure. The c-structure in LFG is a tree representation of the surface structure of a given sentence loosely based on X'-theory (Bresnan, 2001); since it is entirely language-specific and does not serve for encoding tense and aspect categories, our approach does not make use of it. The f-structure is a flat, quasi-logical (Crouch and King, 2006) representation of syntactic relations in terms of attribute-value matrices (Butt et al., 1999). F-structures contain information about syntactic dependencies (predicate-argument structures) as well as further morphosyntactic information such as number, gender, person or tense and aspect. At f-structure, ParGram grammars encode a language-universal level of syntactic analysis, allowing for crosslinguistic parallelism at this level of abstraction.

---

[2]It is well known that past tense markers do not always express a temporal backshift in non-indicative constructions (e.g., conditionals). However, the variation we are concerned with here is situated within constructions that fall under indicative mood.

(4) Er        schrieb  Brief-e
    pron.3sg.m write.past letter-pl
    'He wrote letters.'

$$\begin{bmatrix} \text{PRED} & \text{'schreiben} < \boxed{1}\text{:Er}, \boxed{2}\text{:Brief}> \text{'} \\ \text{SUBJ} & \boxed{1} \begin{bmatrix} \text{PRED} & \text{'pro'} \\ \text{NUM} & \text{sg} \\ \text{CASE} & \text{nom} \end{bmatrix} \\ \text{OBJ} & \boxed{2} \begin{bmatrix} \text{PRED} & \text{'Brief'} \\ \text{NUM} & \text{pl} \\ \text{CASE} & \text{acc} \end{bmatrix} \\ \text{TNS-ASP} & \begin{bmatrix} \text{TENSE} & \text{past} \\ \text{MOOD} & \text{indicative} \end{bmatrix} \\ \text{VTYPE} & \text{main} \\ \text{VFORM} & \text{fin} \\ \text{CLAUSE-TYPE} & \text{decl} \end{bmatrix}$$

Figure 1: f-structure: *Er schrieb Briefe*

As (3) illustrates, the syntactic features are categorized in terms of grammatical functions (the LFG term for syntactic arguments such as SUBJ and OBJ), as well as tense and aspect (TNS-ASP) features of the main verb *schreiben* and other relevant morphosyntactic features. The f-structure is a projection alongside the c-structure that makes accessible morphosyntactic information that is semantically relevant. The information is rendered in terms of feature-value pairs (e.g. [*TENSE past*]). The f-structure can thus readily serve as input to a syntax-semantics interface. The syntactic features for tense and aspect used in the ParGram grammars specifically are the same cross-linguistically. Thus a wide variety of features is available and distinctly describable. They are carefully carved out with the consideration of several different languages and language families in mind (Butt et al., 2002; Sulger et al., 2013). If a language does not express a feature, then the feature is not expressed at f-structure level. Thus, only overtly realized morphosyntactic properties can be read off the f-structure. This allows for a clear-cut differentiation between structural, morphosyntactic and semantic as well as pragmatic features.

## 3.1   Annotation of eventualities at the interface

As stated before, semantic and pragmatic information is encoded in terms of properties of eventualities. On the one hand, there are properties that are concerned with the internal structure of an eventuality, such as telicity or verbal number. In this paper we omit an analysis of these properties, since they are not crucial to the argument defended in this paper. On the other hand, there are properties that relate eventualities to temporal intervals, i.e. tense and grammatical aspect. In what follows we focus on the category of tense which is traditionally understood as an operator relating two time intervals putting them in sequence on a time line (Goranko and Galton, 2015). This notion will be refined in the forthcoming pages. In general, the annotation of a semantic feature always requires the formalization of an inference rule that leads to the annotation. The resulting meaning features are of the form in (5).

(5) ⟨attribute ::= [...],tier ::= [t1 | t2]⟩

Effectively, the annotation is flat, i.e. a set that consists of tuples of the form presented above. Thereby each attribute has a designated feature space as illustrated in figure 2. For the sake of visualization we group certain features under a governing *main feature* presented in angular brackets. The main feature is such that it generalizes the semantic features it subsumes onto a cross-linguistically valid feature space. The main feature should be a cross-linguistically parallel, sufficiently general annotation of the respective *semantic category* it classifies making it an important point of meta data for the comparison of semantic categories. In other words, the main feature is a meta label for a semantic category, while the semantic properties it governs represent the instantiation of the category in a specific language. Consider as a concrete example the category of tense: the logical possibilities of temporal reference are universal to all languages by the nature of its definition, however, some languages further restrict temporal reference as is shown in (6). In this example a temporal remoteness morpheme glossed with *IMM* marks that the sentence is about a time in the immediate past rather than about a time far back in the past which

would be marked with a different temporal remoteness morpheme. This restriction is encoded in the *restr* feature and is a non-necessary modification of a tense category. As such *ref 'past'* $\wedge$ *restr 'immediate'* or *ref 'past'* $\wedge$ *restr 'unspec'* are different realizations of the semantic, cross linguistic category (or main feature) *past*.

(6) **Temporal remoteness in Gĩkũyũ (Cable, 2013)**:

    a. Nĩ-ma-∅-gũr-ire          TV njeru
        ASRT-3pl-**IMM**-buy-<u>PST.PRV</u> TV new
        'They bought a new TV **(today)**'

    b. $\left[ \text{TEMP-REF} <\text{'past'}> \quad \begin{bmatrix} \textbf{ref} ::= \text{'past'} \\ \textbf{restr} ::= \text{'imm'} \end{bmatrix} \right]$

Overall, the mapping from form to meaning is straight forward. We can map tense markers on values for the feature *ref* (short for (temporal) reference) and temporal remoteness markers to further restrict certain values of *ref*. The full spectrum of semantic tense categories and their specific configurations is shown in Figure 2 below. In comparison to the existing state of the art, we provide a more fine-grained set of features that is necessary to annotate cross-linguistic variations of the category of tense.

$$\left[ \text{TEMP-REF} <\text{'past'} \mid \text{'present'} \mid \text{'future'} \mid ... > \quad \begin{bmatrix} \textbf{ref} ::= \text{'past'} \mid \text{'present'} \mid \text{'future'} \mid \\ \text{'non-past'} \mid \text{'non-present'} \mid \text{'non-future'} \mid \text{'unspec'} \\ \textbf{restr} ::= \text{'immediate'} \mid \text{'non-recent'} \mid \text{'remote'} \mid \text{'unspec'} \end{bmatrix} \right]$$

Figure 2: Possible annotations for temporal reference

The mapping from morphosyntax to semantics is not always as straight forward as in the example above. Recall example (2) where we claimed that at least two different types of relations between syntax and semantics obtain. We model these relations in so-called inference rules consisting of a source and a target. The source, represented by the premises of the inference rules, may either be a syntactic or a semantic feature or a set of features. The target is a semantic feature. The most simple inference rules represent a mapping between a meaning feature and its syntactic exponent. However, more complex rules are possible. The basic syntax of inference rules is illustrated below.

(7) $\phi$, $\psi$ are semantic feature/value pairs; $\alpha$, $\beta$, $\gamma$ are morpho-syntactic features, such that the following types of rules are possible:

    a. $\alpha \rightarrow \phi$, $\psi \rightarrow \phi$

    b. $\alpha \wedge \beta \wedge ... \wedge \gamma \rightarrow \phi$

    c. $\alpha \circ \phi$, $\psi \circ \phi$

In (7) $\rightarrow$ describes the implication relation and $\circ$ describes the compatibility relation. A feature or set of features implies a semantic feature, iff there are no two equally strong rules that generate the respective feature. The compatibility relation holds if there are two or more equally strong rules that generate a value for the same feature or if the feature is optional, i.e. an implicature. This means, for each syntactic feature or feature complex, there might be multiple rules that target it, such that we need to formalize certain principles — principles of strength — according to which these rules operate:

Firstly, the implication is stronger than the compatibility relation. This means if we have two rules $\beta \rightarrow \alpha$ and $\gamma \circ \alpha'$, where $\beta$ and $\gamma$ are semantic or syntactic features and $\alpha$ and $\alpha'$ are two different annotations of the same feature, then the attribute/value pair $\alpha$ is generated. Secondly, a rule is stronger if it requires more premises. This means, if there are two rules $\beta \rightarrow \alpha'$ and $\beta \wedge \gamma \rightarrow \alpha$, where $\beta$ and $\gamma$ are semantic or syntactic features and $\alpha$ and $\alpha'$ are two different annotations of the same feature, then, again, the attribute/value pair $\alpha$ is generated.

Based on these rules we now can define primary and secondary meanings. The primary meaning of any syntactic element is the meaning that is generated by the weakest implication rule that exists for this element. The corresponding meaning is labeled the *tier-1* meaning (t1 above) of its syntactic exponent. All other meanings are labeled t2 (*tier-2*), although the range of possible t2 meanings is less coherent than the range of t1 meanings. Tier 2 covers both semantic and pragmatic processes that generate

semantic features, while tier 1 only covers the direct mapping from syntax to semantics. For illustrative purposes assume that $\alpha, \beta, \gamma$ are syntactic features, $\phi, \psi$ are attribute-value pairs(avps) describing semantic features. $\phi, \phi'...$ are alternative annotations of the respective semantic attribute.

(8) **Inference rules for syntactic feature $\alpha$:**

$\alpha, \beta, \gamma \rightarrow \phi \rightarrow$ **tier 2**
$\alpha, \psi \rightarrow \phi' \quad \rightarrow$ **tier 2**
$\alpha \rightarrow \phi'' \qquad \rightarrow$ **tier 1:** $\phi''$ **is the primary meaning of** $\alpha$

The alignment system introduced above is complemented by cross-linguistically universal assumptions about hierarchical relations between features. We illustrate this in terms of the semantic feature of temporal reference that is formalized such, that we can sort the domain of time intervals $D_i$ so that each time interval in $D_i$ belongs to a set representing a possible value for tense. This is done using the temporal precedence relation $\prec$ and the temporal overlap relation $\otimes$ wrt. a given evaluation time $t_0$. The resulting sets can be realized as set that is partially ordered in terms of the inclusion relation (omitting the empty set). In formal semantics a prevalent assumption is that the tense feature is determined by the relation of the so called topic time to the evaluation time Klein (1994). This theory has also been incorporated within the existing TimeML standard by Gast et al. (2016). Keeping future compatibility in mind we also presuppose topic times although we will not explicitly illustrate the system in this paper. This means we will not concern ourselves with complex tense constructions such as the perfect at this point. Depending on the membership of the topic time wrt. the sets introduced in figure 3 the tense feature is then determined. For example if a verbal predicate is restricted to topic times that are included within yesterday, then the verbal predicate is automatically past tense (see section 3.2 and following).

$$\{t|t \in D_i\}$$

$$\{t|t \prec t_0 \vee t \otimes t_0\} \quad \{t|t \prec t_0 \vee t_0 \prec t\} \quad \{t|t \otimes t_0 \vee t_0 \prec t\}$$

$$\{t|t \prec t_0\} \qquad \{t|t \otimes t_0\} \qquad \{t|t_0 \prec t\}$$

Figure 3: Formalization of tense features

The order by inclusion visualizes the hierarchical structure of the different values for tense. For example, the value past $\{t|t \prec t_0\}$ provides a stronger restriction than non-future $\{t|t \prec t_0 \vee t \otimes t_0\}$, thus, if two rules of equal strength compete such that one rule attributes the value 'past' to the feature tense and one rule attributes the value 'non-future' to the feature tense, then the former is applied since it makes a stronger claim.

The formal system introduced above covers a wide array of syntactic and semantic variations in the category of tense and a similar point can be made about aspect. It should be made clear again, that the rules introduced above are language specific. In parallel corpora we assume that the semantics are fairly comparable between languages, however, the inference rules that describe the processes that occur during the mapping from syntax to semantics and within the semantics express the actual variation between different languages.

## 3.2 Advanced annotations – The sequence-of-tense phenomenon

We initially illustrated the annotation in terms of the category of tense. It is not surprising that the form to meaning mapping varies between languages, however, even within languages the mapping from form to meaning is not always completely clear. A famous example that has been widely discussed in the formal semantic literature is the Sequence-of-tense (SOT) phenomenon. A case where syntactic tense is not distinctly mappable to a specific meaning. It is illustrated in (9) below.

(9) Tom said that Karen was sick.

a. Tom said: Karen is sick.

b. Tom said: Karen was sick.

As shown above the SOT-phenomenon allows for two different interpretations of the sentence *Tom said that Karen was sick*. Following the guiding principles introduced in the last section we can provide the following basic tense rule for English. However, it does not allow us to capture the two readings that arise in the SOT sentence.

(10)    a.  *TENSE past $\wedge$ MOOD indicative $\rightarrow$ ref ::= past*

The rule above basically says that the syntactic past tense marker and indicative mood produce a semantic past tense. If the topic times of both eventualities are linked to the speech time of the sentence via a past operator, the resulting semantics are not quite right. Assume that $t^m$ is the topic time of the propositional attitude verb and $t^c$ is the topic time of the verb embedded in the complement, then the the rules above would give us the following three logical possibilities: $(t^m \prec t^c) \prec t^0, (t^c \prec t^m) \prec t^0$ and $(t^c \otimes t^c) \prec t^0$. However, only two of these logical possibilities fit the SOT data, namely the latter two. Either, the eventualities overlap at some point in the past, or the embedded eventuality precedes the matrix eventuality. This also means that we cannot infer the appropriate temporal sequence of eventualities directly from the syntax without assuming some intermediate semantic processes. The virtue of the annotation scheme presented here is that we can explicitly formalize these processes.

We have seen above that the feature past relates the topic time of an eventuality such that it is prior to some evaluation time. Above, we hooked both topic times to the speech time of the sentence. However, tenses may be relative, rather than absolute. In formal semantics absolute tenses are always interpreted with respect to the speech time while relative tenses are (possibly) relative to topic times provided by other tense markers. Usually these tense markers govern the tense in question syntactically (Kusumoto, 2005). However, simply assuming relative tenses does not suffice to reach the proper semantics, because if we interpret the embedded tense as past wrt. the matrix tense, then we only get the so-called past under past reading: $(t^c \prec t^m) \prec t^0$. We have to provide a more elaborate semantic system to account for SOT readings.

Following the formal semantic literature we treat SOT as an ambiguity (Grønn and von Stechow, 2010; Kusumoto, 2005). Seemingly, this ambiguity arises only in specific syntactic contexts as shown in the f-structure in Figure 4. The crucial point in the f-structure template for prototypical SOT sentences above is that two verbal PREDs (the *matrix* verb and the embedded *comp* verb) marked for past tense stand in a syntactic complement relation (COMP), i.e. a past-under-past structure. Thereby, the matrix verb has to be a propositional attitude verb, e.g. *say*, *believe*, *think*. Since this information is not part of the annotation of temporal properties of verbal predicates we follow the annotation guidelines of the TimeML standard differentiating between eventualities and instances of eventualities. For the feature *class* we use the dummy feature *propositional attitude*.[3]

---

[3]This information is annotated in the EVENT tag in the TimeML architecture as the attribute `class` Pustejovsky et al. (2003). At this point we do not roll out the discussion of this feature. For this paper we assume that propositional attitude verbs are certain I_STATE (e.g. believe, think) or REPORTING (e.g. say, report) that occur in (counter-)factive or (negative) evidential subordination links (SLINKS). Concretely, *CLASS(E1) propositional attitude $\wedge$ COMP(E1,E2)* equates such an SLINK. Without going into detail, these SLINKS can be derived from syntactic information and the annotation of the event attribute `class` in terms of inference rules that are coherent with the formal system introduced above. In terms of implementability Crouch and King (2006) provide a semantic system that centers around these types of verbs and the links they invoke. Thus, their system readily provides the means to infer such links inside our implementation.

$$
\begin{bmatrix}
\text{PRED} & \text{matrix} \\
\text{SUBJ} & [...] \\
\text{COMP} & \begin{bmatrix}
\text{PRED} & \text{comp} \\
\text{SUBJ} & [...] \\
\text{TNS-ASP} & \begin{bmatrix} \text{MOOD indicative, ...,} \\ \textbf{TENSE past} \end{bmatrix} \\
\text{COMP-FORM} & \text{that} \\
\text{CLAUSE-TYPE} & \text{decl}
\end{bmatrix} \\
\text{TNS-ASP} & \begin{bmatrix} \text{MOOD indicative, ....,} \\ \textbf{TENSE past} \end{bmatrix} \\
\text{CLAUSE-TYPE} & \text{decl}
\end{bmatrix}
$$

(11) **Additional annotation steps:**

   a. PRED matrix $\rightarrow$ E1

   b. PRED comp $\rightarrow$ E2

   c. CLASS(E1)
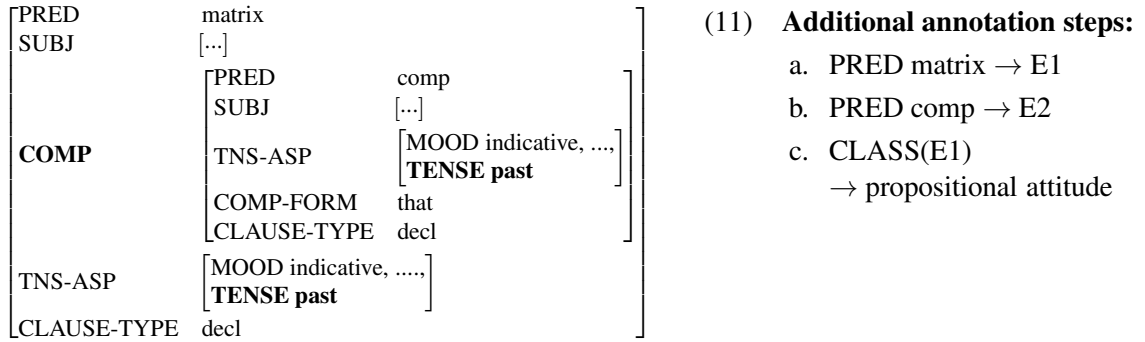      $\rightarrow$ propositional attitude

Figure 4: f-structure skeleton for SOT

With this information in place we can provide an additional rule that fulfills the SOT requirements in English. As (12b) shows the rule is incomparably more complex than the basic tense rule we presupposed before. It results in non-future temporal reference. Due to the definition of tense given above the non-future value for the embedded tense subsumes both cases: simultaneous readings and back-shifted readings. At the core of this rule are *propositional attitude* feature of the matrix predicate and the syntactic complement construction. In short, if an event is subordinated under an event via this rule and both of these events are annotated as past tense, then the embedded tense receives the value *non-future* or *non-successive* for temporal reference. This feature is interpreted with respect to the topic time of the matrix tense rather than the speech time. In this respect English is relative in this annotation scheme. The respective rules are illustrated below in a simplified manner:

(12) **Inference rules for (past) tense in English:**

   a. *tier 1:*
    *TENSE past $\wedge$ MOOD indicative $\rightarrow$ temp-ref 'past'*

   b. *tier 2:*
    *CLASS(E1) propositional attitude $\wedge$ COMP(E1,E2) $\wedge$ MOOD indicative $\wedge$*
    *TENSE(E1) past $\wedge$ TENSE(E2) past $\rightarrow$ temp-ref(E2) 'non-future'*

In TimeML the annotation of a SOT sentence may be obtained, as is illustrated below (omitting optional/unnecessary tags):

(13)   a.
```
EVENT : eid ::= e1, class ::= reporting
EVENT : eid ::= e2, class ::= state
```
   b.
```
MAKEINSTANCE : eiid ::= ei1, eventID ::= e1
              pos ::= 'VERB'
              tense ::= 'PAST'
              aspect ::= 'NONE'

MAKEINSTANCE : eiid ::= ei2, eventID ::= e2
              pos ::= 'VERB'
              tense ::= 'PAST'
              aspect ::= 'NONE'
```
   c.
```
SLINK : lid ::= 1, eventInstanceID ::= ei1
        subordinatedEventInstance ::= ei2
        relType ::= 'Evidential'
```
   d.  i.
```
TLINK : lid ::= 2, eventInstanceID ::= ei1
        relatedToEventInstance ::= ei2
        relType ::= 'BEFORE'
```
      ii.
```
TLINK : lid ::= 2, eventInstanceID ::= ei1
        relatedToEventInstance ::= ei2
        relType ::= { 'SIMULTANEOUS' | 'IS_INCLUDED' }
```

As example (13d) shows, it is difficult to find a unique annotation for the SOT sentence wrt. the relation between the matrix event and the embedded event (TLINKs to the speech time have been omitted for reasons of space). In the annotation presented in this paper we obtain a semantic restriction at sentence level for the embedded event, namely, that is has to be non-successive (or non-future) wrt. the matrix event. This seems like a small gain, however, it captures that, cross-linguistically, languages behave differently wrt. the SOT-parameter: They build a different semantic frame for interpreting temporal relations. The importance of such a frame becomes more apparent when the SOT phenomenon is explored cross-linguistically, as shown in, for example, Bochnak (2016); Mucha (2015), where the role of a purely morphosyntactic frame is challenged. Furthermore, the next section will show, that an annotation along the lines of timeML sketched above is not sufficient to describe cross-linguistic variation of tense categories.

# 4 Cross-linguistic variation in the category tense

In this section we apply the overall annotation scheme to different linguistic expressions. We define two types of variations which mark cornerstones in the spectrum of variation within syntax/semantics interface: Syntactic variation and semantic variation. Syntactic variation occurs when two eventualities that have the same semantic features are realized differently on a morphosyntactic level. Semantic variation occurs when syntactically similar constructions result in different annotations of the corresponding eventualities. We illustrate the two types of similarity in terms of minimal pairs so as to keep variation within a reasonable degree.

## 4.1 Syntactic variation

Syntactic variation is arguably the simpler of the two types of variation. We define it in terms of similarity on f-structure level, that is, two sentences are syntactically similar if they express alignable f-structures, i.e. they are non-contradicting (Sulger et al., 2013). We focus on variation in terms of the TNS-ASP grammatical functions in this paper (However, not all syntactic variation relevant for annotation of eventualities is necessarily confined within the TNS-ASP node). Thus, the two f-structures below are syntactically identical except for the lexical instantiation of the PRED.

(14) $\begin{bmatrix} \text{PRED 'schreiben} < ... > \text{'} \\ ... \\ \text{TNS-ASP} \quad [\text{TENSE past, MOOD indicative}] \end{bmatrix}$ (15) $\begin{bmatrix} \text{PRED 'write} < ... > \text{'} \\ ... \\ \text{TNS-ASP} \quad [\text{TENSE past, MOOD indicative}] \end{bmatrix}$

Assume now a language that does not express tense overtly such as for example Indonesian (Arka et al., 2013). A corresponding TNS-ASP matrix is shown in (16).

(16) $\begin{bmatrix} \text{TNS-ASP} \quad [\text{MOOD indicative}] \end{bmatrix}$

The minimal difference between the TNS-ASP matrix in (16) and the matrices in (14) and (15) is that there is no syntactic tense. Consider the sentences in (17), where this would be the minimal difference (in terms of tense) between the two sentences. Out of the blue the two sentences would be different in their compatibility with semantic features. However, they occur in a specific context such that both of the sentences are about an event in the past. Thus, in a context-driven annotation both sentences are semantically past. However, the explicit past should be annotated as conveying past time reference while the contextual past in Indonesian is only compatible with past time reference.

(17) Q: What did the farmer do (yesterday)?

    a. The farmer groaned.

    b. Petani itu mengaduh
       Farmer that groan

Since the past time reference in Indonesian has no overt syntactic exponent as reflected in the partial f-structure in (16), we have to provide inference rules that exhibit the difference between English and Indonesian. The first step is to annotate syntax and semantics. The second step is to annotate the corresponding inference rules for English and Indonesian as shown in (18). The resulting representations of the two different instantiations of the category past tense are illustrated in (19).

(18)  a.  English:
          *TENSE past ∧ MOOD indicative → ref ::= past*

      b.  Indonesian:
          *MOOD indicative ∘ ref ::= past*

(19)  a.  I met Peter (at the market).

| F-Structure: | ParTMA Temporal reference: |
|---|---|
| $\big[$TNS-ASP $\;$ [TENSE **past**, MOOD indicative]$\big]$ | $\boxed{101}\;\big[$TEMP-REF <'past'> $\big[$ref ::= **'past,t1'** $\;$ restr ::= 'unspec'$\big]\big]$ |

      b.  Saya bertemu Peter (di pasar (itu)).

| F-Structure: | ParTMA Temporal reference: |
|---|---|
| $\big[$TNS-ASP $\;$ [MOOD indicative]$\big]$ | $\boxed{101}\;\big[$TEMP-REF <'past'> $\big[$ref ::= **'past,t2'** $\;$ restr ::= 'unspec'$\big]\big]$ |

These rules capture the relationship between syntax and semantics in these two languages. In English there is a very direct mapping from syntax to semantics. In Indonesian the connection is only apparent. The compatibility relation between indicative mood and past time reference does not entail that every sentence that is marked in indicative mood has past time reference. Where does the past time reference come from? Intuitively, the answer is that the past time reference is a result of contextual inference based on an available, salient time interval. In the context above this time interval is denoted by *yesterday*. We claimed before that tense sorts the available temporal intervals respective to the evaluation time. However, this does not mean that a past tense annotation denotes every interval that qualifies as that tense. Rather a tense annotation is appropriate if the temporal interval that a sentence is about, the Kleinian topic time (Klein, 1994), is included in the set described by the respective feature. With Gast et al. (2015) and following Klein (1994) we claim that finite verbs introduce topic times. In formal semantics topic times are treated as pronominal elements that point to a specific time interval (Partee, 1973). This leads to the following annotation process for English:

(20)  **Annotation of:**
      Q: What did the farmer do yesterday?
      A: The farmer groaned.

      a.  Syntactic parsing

      b.  did → E1, groaned → E2

      c.  temporal expressions:
          → [*yesterday*, *evaluation time*]

      d.  Relate events to temporal intervals
          → [*E1, E2 ⊆ yesterday*]

      e.  time reference:
          $E2 ⊆ yesterday ∧ \underline{TENSE(E2)\ past} → ref ::= 'past',t1$

In Indonesian the process would be similar for the most part. The only difference would be the underlined part, which would be void. Thus, the inference of past time reference would be purely semantic(/pragmatic). This would also mean that the Indonesian annotation would receive the label *tier 2* since we can't relate the semantics to any syntactic exponent via an implication relation as illustrated in (18).

The example above illustrates the interaction between the module of syntax, semantics and a contextual component currently carried out through manual disambiguation which regulates the interactions between topic times and (salient) temporal variables (step d. above). For current annotations we employ the dummy label ctx for elements outside of the syntactic and the eventuality module, such that *ctx(past)* ∧ *MOOD indicative ∘ ref ::= 'past',t2* describes a rule that, given some context that requires a past interpretation combined with a sentence in indicative mood determines semantic past tense such as in the Indonesian example in (17). We, thus, slightly revise the rule given in (18).

## 4.2 Semantic variation

In section 3.2 we introduced the Sequence-of-tense phenomenon and illustrated how our annotation scheme captures shifts in meaning within a language. However, this anomaly is not cross-linguistically robust. There seems to be a parameter that distinguishes two groups of languages. Those that express the phenomenon such as English and those that do not express it, such as for example, Japanese.

(21) Tom said that Karen was sick.
   a. Tom said: Karen is sick.
   b. Tom said: Karen was sick.

(22) Jon wa Karen ga byōkida to itta
   John top Karen subj be-sick.past comp say.past
   'Tom said: Karen <u>was</u> sick.

Table 1: SOT vs NON-SOT language

| Readings | SOT | NON-SOT |
|---|---|---|
| $(t^c \prec t^m) \prec t^0$ | + | + |
| $(t^c \otimes t^c) \prec t^0$ | + | - |
| $t^m \prec t^c) \prec t^0$ | - | - |

The annotation system provided here provides a simple solution to capture the difference in the semantics – namely a variation in terms of the syntax/semantics interface. If we treat SOT as a parameter that allows for specific rules, there is simply no SOT rule in Japanese. Concretely, applying the simple past tense rule we provided initially, we will get the required result for Japanese. Thereby, both English and Japanese tenses are treated as relative tenses rather than distinguishing between absolute and relative tense systems.

Overall, we have shown that the annotation scheme presented here can handle different cases of variations in terms of tense categories easily. Thus, the annotation scheme provides a valuable tool for qualitative analysis of tense categories while pertaining integratability in the broader picture of temporal annotation. A similar point can be made about various aspectual features as well, but we leave this endeavor for elsewhere (see e.g. Zymla (to appear)).

## 5 Summary

In this paper we presented a novel annotation for tense using as example different variations related to the category of tense. The research presented here is situated within a broader effort to provide a linguistically and formally sound annotation for tense and aspect for NLP applications. We focused on small data examples from the formal semantic literature to illustrate the overall architecture. We bring together syntactic resources as well as existing resources in the realm of temporal tagging to provide annotations that encode the intrinsic meaning variations that arise from the morphosyntactic marking of tense or the absence thereof.

We did not go into detail with regards to the implementability of the system presented here with regards to the language specific inference rules. Complemented with a suitable database for lexical semantics the presented annotation could be automated to a large degree relying on human supervision for disambiguation and resolving contextual inferences the system cannot make due to the lack of a

formalized pragmatic module. The biggest semantic resource specifically for XLE grammars is the semantic system employed in the ParcBridge Q & A system. Crouch (2005); Crouch and King (2006). We work on an extension of this system to incorporate the annotation presented in this paper. Languages without the respective semantic resources have to be annotated manually or at least require substantially more manual labor.

To summarize, we provide an annotation scheme for primarily qualitative research on the syntax/semantics interface cross-linguistically. Thereby we use the parallel grammars of the ParGram project as a foundation. We integrated the annotation scheme into the broader effort of temporal annotation and thus provide means to bring temporal annotation and deep linguistic parsing closer together. This promises to allow us to better test formal linguistic insights concerning the mapping from syntax to semantics. Furthermore, we have shown, how the separation of syntax, semantics and the syntax/semantics interface provides a more clear representation of various tense categories in general.

## Acknowledgments

## References

I Wayan Arka et al. 2013. On the typology and syntax of tam in indonesian. *Tense, aspect, mood and evidentiality in languages of Indonesia* pages 23–40.

Emmon Bach. 1986. The Algebra of Events. *Linguistics and Philosophy* 9(1):5–16.

Steven Bethard and Jonathan L Parker. 2016. A semantically compositional annotation scheme for time normalization. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC 2016)*.

Maria Bittner. 2014. *Temporality: Universals and variation*. John Wiley & Sons.

M Ryan Bochnak. 2016. Past time reference in a language with optional tense. *Linguistics and Philosophy* 39(4):247–294.

Joan Bresnan. 2001. *Lexical-functional Syntax*, volume 16 of *Blackwell textbooks in linguistics*. Blackwell Publishing.

Miriam Butt, Helge Dyvik, Tracy Holloway King, Hiroshi Masuichi, and Christian Rohrer. 2002. The parallel grammar project. In *Proceedings of the 2002 workshop on Grammar engineering and evaluation*. Association for Computational Linguistics, volume 15, pages 1–7.

Miriam Butt, Tracy Holloway King, María-Eugenia Niño, and Frédérique Segond. 1999. *A Grammar Writer's Cookbook*. CSLI Publications.

Seth Cable. 2013. Beyond the past, present, and future: towards the semantics of 'graded tense'in Gĩkũyũ. *Natural Language Semantics* 21(3):219–276.

Dick Crouch, Mary Dalrymple, Ronald M. Kaplan, Tracy Holloway King, John T. Maxwell III, and Paula Newman. 2017. *XLE Documentation*. Palo Alto Research Center.

Richard Crouch. 2005. Packed rewriting for mapping semantics to KR. In *Proceedings of the Sixth International Workshop on Computational Semantics (IWCS-6)*. Tilburg, pages 103–114.

Richard Crouch and Tracy Holloway King. 2006. Semantics via F-Structure Rewriting. In Miriam Butt and Tracy Holloway King, editors, *Proceedings of the LFG06 Conference*. CSLI Publications, Stanford, CA, pages 145–165.

Mary Dalrymple. 2001. *Lexical Functional Grammar*, volume 34 of *Syntax and Semantics*. Academic Press, New York.

Leon Derczynski, Hector Llorens, and Naushad UzZaman. 2013. Timeml-strict: clarifying temporal annotation. *arXiv preprint arXiv:1304.7289* .

Volker Gast, Lennart Bierkandt, Stephan Druskat, and Christoph Rzymski. 2016. Enriching timebank: Towards a more precise annotation of temporal relations in a text. In Nicoletta Calzolari (Conference Chair), Khalid Choukri, Thierry Declerck, Sara Goggi, Marko Grobelnik, Bente Maegaard, Joseph Mariani, Helene Mazo, Asuncion Moreno, Jan Odijk, and Stelios Piperidis, editors, *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC 2016)*. European Language Resources Association (ELRA), Paris, France.

Volker Gast, Lennart Bierkandt, and Christoph Rzymski. 2015. Creating and retrieving tense and aspect annotations with graphanno, a lightweight tool for multi-level annotation. In *Proceedings 11th Joint ACL-ISO Workshop on Interoperable Semantic Annotation*. page 23.

Valentin Goranko and Antony Galton. 2015. Temporal logic. In Edward N. Zalta, editor, *The Stanford Encyclopedia of Philosophy*, Metaphysics Research Lab, Stanford University. Winter 2015 edition.

Atle Grønn and Arnim von Stechow. 2010. Complement tense in contrast: the sot parameter in russian and english. *Oslo Studies in Language* 2(1).

Wolfgang Klein. 1994. *Time in language*. Psychology Press.

Kiyomi Kusumoto. 2005. On the quantification over times in natural language. *Natural language semantics* 13(4):317–357.

Anne Mucha. 2015. *Temporal interpretation and cross-linguistic variation*. Ph.D. thesis, PhD thesis, University of Potsdam, Potsdam, Germany.

Barbara Hall Partee. 1973. Some structural analogies between tenses and pronouns in english. *The Journal of Philosophy* 70(18):601–609.

James Pustejovsky, José M Castano, Robert Ingria, Roser Sauri, Robert J Gaizauskas, Andrea Setzer, Graham Katz, and Dragomir R Radev. 2003. Timeml: Robust specification of event and temporal expressions in text. *New directions in question answering* 3:28–34.

James Pustejovsky, Roser Saurí, Andrea Setzer, Rob Gaizauskas, and Bob Ingria. 2002. Timeml annotation guidelines. *TERQAS Annotation Working Group* 23.

Carlota Smith. 2006. The pragmatics and semantics of temporal meaning. In *Proceedings, Texas Linguistics Forum*.

Sebastian Sulger, Miriam Butt, Tracy Holloway King, Paul Meurer, Tibor Laczkó, György Rákosi, Cheikh M Bamba Dione, Helge Dyvik, Victoria Rosén, Koenraad De Smedt, Agnieszka Patejuk, Özlem Çetinŏglu, I Wayan Arka, and Meladel Mistica. 2013. ParGramBank: The ParGram Parallel Treebank. In *ACL*. pages 550–560.

Naushad UzZaman, Hector Llorens, James Allen, Leon Derczynski, Marc Verhagen, and James Pustejovsky. 2012. Tempeval-3: Evaluating events, time expressions, and temporal relations. *arXiv preprint arXiv:1206.5333* .

Mark-Matthias Zymla. to appear. Cross-Linguistically Viable Treatment of Tense and Aspect in Parallel Grammar Development. In *Proceedings of the LFG17 Conference*. CSLI Publications.