

Chapter 16: Discourse

Pierre Nugues

Pierre.Nugues@cs.lth.se
http://cs.lth.se/pierre_nugues/

October 14, 2021



A Definition of Discourse

A discourse is a sequence of sentences: a text or a conversation

A discourse is made of words or phrases that refer to things: the **discourse entities**

A discourse normally links the entities together to address topics

Within a single sentence, grammatical structures provide with a model of relations between entities.

Discourse models extend relations to more sentences



Reference

Discourse entities – or discourse referents – are the real, abstract, or imaginary objects introduced by the discourse. **Referring expressions** are mentions of the discourse entities through the text

- ① Susan drives a Ferrari
- ② She drives too fast
- ③ Lyn races her on weekends
- ④ She often beats her
- ⑤ She wins a lot of trophies



Discourse Entities

Mentions (or referring expressions)	Discourse entities (or referents)	Logic properties
<i>Susan, she, her</i>	'Susan'	'Susan'
<i>Lyn, she</i>	'Lyn'	'Lyn'
<i>A Ferrari</i>	X	ferrari(X)
<i>A lot of trophies</i>	E	$E \subset \{X \mid \text{trophy}(X)\}$



Reference and Named Entities

Named entities are entities uniquely identifiable by their name.

Some definitions/
clarifications:

- Named entity recognition (NER): a partial parsing task, see Chap. 10;
- Reference resolution for named entities: find the entity behind a mention, here a name.

Words	POS	Groups	Named entities
U.N.	NNP	I-NP	I-ORG
official	NN	I-NP	O
Ekeus	NNP	I-NP	I-PER
heads	VBZ	I-VP	O
for	IN	I-PP	O
Baghdad	NNP	I-NP	I-LOC
.	.	O	O

As it is impossible to set a physical link between a real-life object and its mention, we use unique identifiers or tags in the form of URIs instead (from Wikidata, DBpedia, Yago).



Mentions of Named Entities are Ambiguous

Cambridge: England, Massachusetts, or Ontario?

Given the text (from Wikipedia):

*One of his translators, Roy Harris, summarized **Saussure**'s contribution to linguistics and the study of language in the following way...*

Which Saussure? *Saussure* has 11 entries in Wikipedia:

- *Ferdinand de Saussure*:
 - Wikidata: <http://www.wikidata.org/wiki/Q13230>
 - DBpedia: http://dbpedia.org/resource/Ferdinand_de_Saussure
- *Henri de Saussure*: <http://www.wikidata.org/wiki/Q123776>
- *René de Saussure*: <http://www.wikidata.org/wiki/Q13237>



Collecting Entity-Mention Pairs from Wikipedia

Wikipedia has a mark up that enables an editor to link a word or phrase to a page:

- `[[Ferdinand_de_Saussure|Saussure]]` or
- `[[target or link|text or label or anchor]]`

In our case, it is an association between a mention and an entity:

`[[Entity|Mention]]`

All the links can be extracted from a wikipedia dump to derive two probabilities:

- The probability of a mention given an entity, how we name things:
 $P(M|E)$
- The probability of a entity given a mention, the ambiguity of a mention: $P(E|M)$



Göran Persson in Swedish




In Wikipedia, at least four entities can be linked to the name *Göran Persson*:

























- ❶ **Göran Persson** (född 1949), socialdemokratisk partiledare och svensk statsminister 1996–2006 (Q53747)
- ❷ **Göran Persson** (född 1960), socialdemokratisk politiker från Skåne (Q5626648)
- ❸ Göran Persson (militär), svensk överste av 1:a graden
- ❹ **Göran Persson** (musiker), svensk proggmusiker (Q6042900)
- ❺ Göran Persson (litterär figur), överkonstapel i 1930-talets Lysekil
- ❻ Göran Persson (skulptör) (född 1956), konstnär representerad i bl.a. Karlskoga
- ❼ **Jöran Persson**, svensk ämbetsman på 1500-talet (Q2625664)



$P(\text{Mention}|\text{Entity})$, An Exemple

Mentions of *Göran Persson*, Q53747, in Swedish:
How do we name Q53747?

Mentions




   Göran Persson	 <input type="checkbox"/> <input type="checkbox"/> Göran Persson (s)
   Göran Perssons	 <input type="checkbox"/> <input type="checkbox"/> Göran Persson (statsminister)
 <input type="checkbox"/>  Persson	 <input type="checkbox"/> <input type="checkbox"/> Göran Persson i Stjärnhov
 <input type="checkbox"/>  (Hans) Göran Persson	 <input type="checkbox"/> <input type="checkbox"/> Hans G. Persson
 <input type="checkbox"/>  Han Som Bestämmer	 <input type="checkbox"/> <input type="checkbox"/> Hans Göran Persson
 <input type="checkbox"/>  Perssonplanen	 <input type="checkbox"/> <input type="checkbox"/> Persson, Göran
 <input type="checkbox"/>  Perssons	 <input type="checkbox"/> <input type="checkbox"/> Påven vid Båven
 <input type="checkbox"/> <input type="checkbox"/> Goran Persson	
 <input type="checkbox"/> <input type="checkbox"/> Göran Person	

From <http://klang.cs.lth.se:8888/en/data/wiki>



$P(\text{Entity}|\text{Mention})$, An Exemple

Entities linked to the mention *Göran Persson* in Swedish:
The things behind *Göran Persson*

sv mention:Göran Persson 🔍

▶ 6 hits in 0,003 s

- ▶ **Göran Persson (född 1960)**
wikidata:Q5626648
- ▶ **Göran Persson (musiker)**
wikidata:Q6042900
- ▶ **Lars Göran Persson**
wikidata:Q6043257
- ▶ **Göran Persson**
wikidata:Q53747
- ▶ **Jöran Persson**
wikidata:Q2625684
- ▶ **Regeringen Persson**
wikidata:Q4570330

From <http://klang.cs.lth.se:8888/en/data/wiki>



Disambiguation of Named Entities

Given:

*One of his translators, Roy Harris, summarized **Saussure**'s contribution to linguistics and the study of language...*

Disambiguation is a classification problem dealing with mention-entity pairs:

Mention	Entity	Q number	T/F
Saussure	Ferdinand de Saussure	Q13230	1
Saussure	Henri de Saussure	Q123776	0
Saussure	René de Saussure	Q13237	0
...			

Feature vectors represent pair of mentions and entities:

- Cosine similarity between the mention context and the named entity page in Wikipedia and bag-of-word vectors of the mention context
- Training set built from Wikipedia markup:
[[Ferdinand_de_Saussure|Saussure]]



Named Entities and Linked Data

Graph databases are popular devices used to represent named entities, especially the resource description framework (RDF).

Entities are assigned unique resource identifiers (URIs) similar to URLs (as in HTTP addresses) and can be linked to other data sources (Linked data)

Examples of databases using the RDF format:

DBpedia: A database of persons, organizations, locations, etc.
DBpedia is automatically extracted from Wikipedia
semi-structured data (info boxes)

Geonames: A database of geographical names (a gazetteer).

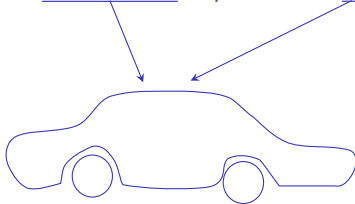
SPARQL is a database query language that enables a programmer to extract data from a graph database (similar to Prolog or SQL)



Coreference

[entity1 Garcia Alvarado], 56, was killed when [entity2 a bomb] placed by [entity3 urban guerrillas] on [entity4 his vehicle] exploded as [entity5 it] came to [entity6 a halt] at [entity7 an intersection] in [entity8 downtown] [entity9 San Salvador].

on his vehicle exploded as it came to a halt



Anaphora

Anaphora, often pronouns

Pronouns: it, she, he, this, that

Cataphora

*I just wanted to touch **it**, this stupid animal.*

***They** have stolen my bicycle.*

Antecedents

Ellipsis is the absence of certain referents

I want to have information on caterpillars. And also on hedgehogs.



Coreference Annotation

The MU Conferences have defined a standard annotation for noun phrases. It uses the COREF element with five possible attributes: ID, REF, TYPE, MIN, and STAT.

- `<COREF ID="100">Lawson Mardon Group Ltd.</COREF> said`
`<COREF ID="101" TYPE="IDENT" REF="100">it</COREF>`
- `<COREF ID="100" MIN="Haden MacLellan PLC">Haden`
`MacLellan PLC of Surrey, England</COREF> ...`
`<COREF ID="101" TYPE="IDENT" REF="100">Haden`
`MacLellan</COREF>`



Coreference Annotation: CoNLL 2011 simplified

0		"	"	...	-
1	Vandenberg	NNP			(8 (0)
2	and	CC			-
3	Rayburn	NNP			(23) 8)
4	are	VBP			-
5	heroes	NNS			-
6	of	IN			-
7	mine	NN			(15)
8	,	,			-
9	"	"			-
10	Mr.	NNP			(15
11	Boren	NNP			15)
12	says	VBZ			-
13	,	,			-
14	referring	VBG			-
15	as	RB			-
16	well	RB			-
17	to	IN			-
18	Sam	NNP			(23
19	Rayburn	NNP			-
20	,	,			-
21	the	DT			-
22	Democratic	JJ			-
23	House	NNP			-
24	speaker	NN			-
25	who	WP			-
26	cooperated	VBD			-
27	with	IN			-
28	President	NNP			-
29	Eisenhower	NNP			23)
30	.	.			-

Entities and mentions:

$e_0 = \{Vandenberg\}$

$e_8 = \{Vandenberg \text{ and } Rayburn\}$

$e_{15} = \{mine, Mr. Boren\}$

$e_{23} =$

$\{Rayburn, Sam Rayburn, 'the Democratic House speaker who cooperated with President Eisenhower\}$



Coreference Chains

In the MUC competitions, coreference is defined as symmetric and transitive:

- If A is coreferential with B, the reverse is also true.
- If A is coreferential with B, and B is coreferential with C, then A is coreferential with C.

It forms an equivalence class called a **coreference chain**.

The TYPE attribute specifies the link between the anaphor and its antecedent.

IDENT is the only possible value of the attribute

Other types are possible such as part, subset, etc.



Solving Coreferences

Coreferences define a class of equivalent references

Backward search with a compatible gender and number

98% of the antecedents are in the current or previous sentence

Focus: an integer attached to all objects, incremented when:

- It is mentioned: subject, object, adjunct
- It is visible or pointed at.

The focus is decremented over time

Constraints are also applied: subject \neq object, grammatical role

Anaphora is resolved by taking the highest focus



A Simplistic Method

*Garcia Alvarado, 56, was killed when **a bomb** placed by urban guerrillas
on **his vehicle** exploded as **it** came to a halt at an intersection in
downtown San Salvador*

Diagram illustrating a simplistic method for discourse analysis. The text is segmented into two parts, labeled 1 and 2, with arrows indicating the flow of information.

Part 1 (labeled 1) covers the first line: *Garcia Alvarado, 56, was killed when*

Part 2 (labeled 2) covers the second line: *a bomb placed by urban guerrillas on his vehicle exploded as it came to a halt at an intersection in downtown San Salvador*



Machine Learning to Solve Coreferences

Instead of manually engineered rules, machine learning uses an annotated corpus and trains the rules automatically.

The coreference solver is a decision tree. It considers pairs of noun phrases (NP_i, NP_j).

Each pair is represented by a feature vector of 12 parameters.

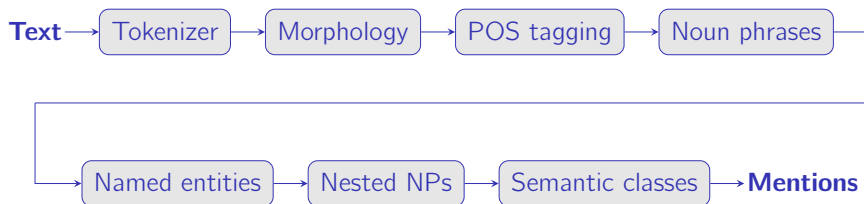
The tree takes the set of NP pairs as input and decides for each pair whether it corefers or not.

Using the transitivity property, it identifies all the coreference chains in the text.

The ID3 learning algorithm automatically induces the decision tree from texts annotated with the MUC annotation standard.



Architecture



The coreference engine takes a pair of extracted noun phrases (NP_i, NP_j) For a given index j , the engine considers from right to left, NP_i as a potential antecedent and NP_j as an anaphor. It classifies the pair as positive if both NPs corefer or negative if they don't.



Some Features

- Positional feature:
 1. Distance (DIST): This feature is the distance between the two noun phrases measured in sentences: 0, 1, 2, 3, ... The distance is 0 when the noun phrases are in the same sentence.
- Grammatical features:
 2. *i*-Pronoun (I_PRONOUN): Is NP_i a pronoun i.e. personal, reflexive, or possessive pronoun? Possible values are true or false.
 3. *j*-Pronoun (J_PRONOUN): Is NP_j a pronoun? Possible values are true or false.
- Lexical feature:
 12. String match (STR_MATCH): Are NP_i and NP_j equal after removing articles and demonstratives from both noun phrases? Possible values are true or false.



Training Examples: The Positive Examples

The classifier can be a decision tree or logistic regression.

It is trained from positive and negative examples extracted from the annotated corpus

The positive examples use pairs of adjacent coreferring noun phrases.

If $NP_{a1} - NP_{a2} - NP_{a3} - NP_{a4}$ is a coreference chain in a text, we have

Noun phrases	Coreference chains
NP_{a1}	Chain 22
...	
NP_{a2}	Chain 22
...	
NP_{a3}	Chain 22
...	
NP_{a4}	Chain 22
...	

The positive examples correspond to the pairs: (NP_{a1}, NP_{a2}) , (NP_{a2}, NP_{a3}) , (NP_{a3}, NP_{a4})



Training Examples: The Negative Examples

The negative examples consider the noun phrases NP_{i+1} , NP_{i+2}, \dots , NP_{j-1} intervening between adjacent pairs (NP_i, NP_j) .

Noun phrases	Coreference chains	Relation
NP_i	Chain 22	Antecedent
NP_{i+1}	Not part of Chain 22	
NP_{i+2}	Not part of Chain 22	
...		
NP_{j-1}	Not part of Chain 22	Anaphor
NP_j	Chain 22	

For each positive pair (NP_i, NP_j) , the training procedure generates negative pairs:

- They consist of one intervening NP and the anaphor NP_j : (NP_{i+1}, NP_j) , (NP_{i+2}, NP_j) , \dots , and (NP_{j-1}, NP_j) .
- The intervening noun phrases can either be part of another coreference chain or not.



Performances

At this point, it is useful to have the current performances in mind

- Morphological parsing can parse correctly 99 % of the words in many languages (Koskenniemi 1984)

Bilolyckorna "bil#olycka" N UTR DEF PL NOM

- Part-of-tagging reached and exceeded 97% as early as Church (1991)

En bilolycka med tre bilar

En/dt_utr_sin_ind bilolycka/nn_utr_sin_ind_nom med/pp
tre/rg_nom bilar/nn_utr_plu_ind_nom

- Sentence parsing reaches ~89% in Swedish (CoNLL 2018) – labeled dependencies.



Performances (II)

- Semantic parsing (extraction of predicate–argument structures.)
The F-measure reaches about 85.5 in 2019 (CONLL 2009).

[Judge **She**] **blames** [Evalued **the Government**] [Reason **for failing to do enough to help**]
`blames(judge, evaluatee, reason)`

`blames('She', 'The Government', 'for failing to do enough to help')`.

- Coreference solving reaches a F-measure of ~ 80 in 2019 using the CoNLL 2012 CoNLL script up from 60 in Pradhan et al. (2011)
- A site to have up-to-date figures:
<https://github.com/sebastianruder/NLP-progress>



Discourse Theories and Models

Discourse theories are used to develop organization models of texts. They have three objectives: **represent**, **parse automatically**, and **generate** a discourse.

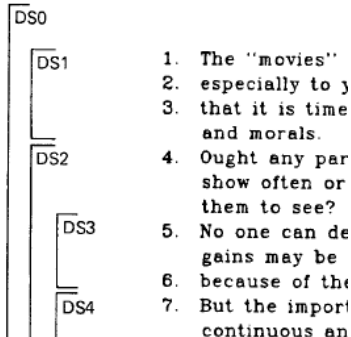
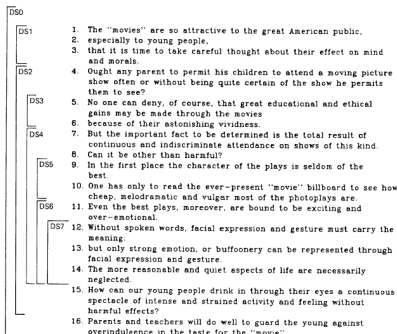
There are many ways to represent a text and competing theories. In 1992, Mann and Thompson compared 12 different representations obtained from experts in the field. The most significant are:

- Grosz and Sidner's theory (1986) and Centering (1995)
- Rhetorical structure theory (RST) (Mann and Thompson 1988)



Grosz and Sidner's Theory

Discourse describes a hierarchical tree



Centers

Centers are entities that link one a sentence to another one.

Grosz divides centers in a unique **backward-looking center** that is the most important entity in the segment and others **forward-looking** centers.

Two relations link segments: dominance and satisfaction-precedence.



Rhetoric

- Invention (*Inventio*).
- Arrangement (*Dispositio*): introduction (*exordium*), a narrative (*narratio*), a proposition (*propositio*), a refutation (*refutatio*), a confirmation (*confirmatio*), and finally a conclusion (*peroratio*).
- Style (*Elocutio*): emote (*movere*), explain (*docere*), or please (*delectare*).
- Memory (*Memoria*)
- Delivery (*Actio*).



Rhetorical Structure Theory

The rhetorical structure theory is a text grammar that analyzes argumentation: A text consists of:

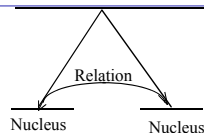
- **Text spans** that can be sentences or clauses
- **Rhetorical relations** that link the text spans

Relations are richer than with Grosz and Sidner.

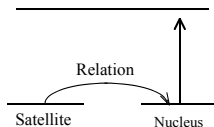


Relations

Relations between segments can be symmetrical when spans have the same importance: Both spans are **nuclei**.



When relations are asymmetrical, we have a **nucleus** and a **satellite** where the nucleus is the most important



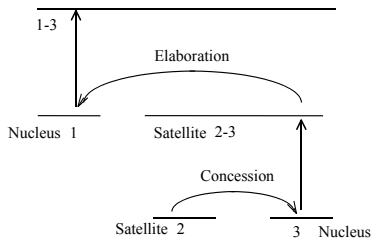
The text analysis produces a tree of text spans that are linked by different relation types.



Graphical Representation

Example cited by Mann and Thompson (1987):

- ① *Concern that this material is harmful to health or the environment may be misplaced.*
- ② *Although it is toxic to certain animals,*
- ③ *evidence is lacking that it has any serious long-term effect on human beings.*

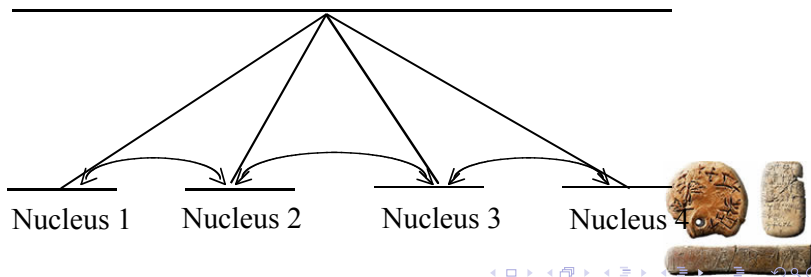


Links Between Nuclei

Spans can have a same importance and are linked by a sequence relation:

- ① *Napoleon met defeat in 1814 by a coalition of major powers, notably Prussia, Russia, Great Britain, and Austria.*
- ② *Napoleon was then deposed*
- ③ *and exiled to the island of Elba*
- ④ *and Louis XVIII was made ruler of France.*

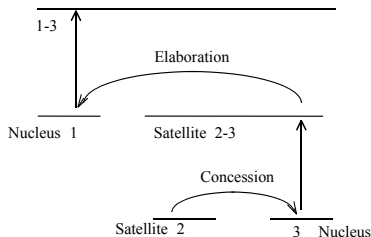
Microsoft Encarta, cited from Simon Corston-Oliver (1998)



Attempt to Formalize Structure

Mann and Thompson gave a formal structure to the graph that correspond to a parse tree:

- 1 The tree extends over the whole text;
- 2 Each text span part of the text analysis is either a terminal symbol or a node constituent;
- 3 A span has a unique parent;
- 4 Relations bind adjacent spans.



RST Relations

The original relations in RST are:

Nucleus-satellite relations

Circumstance	Evidence	Otherwise
Solutionhood	Justify	Interpretation
Elaboration	Cause	Evaluation
Background	Antithesis	Restatement
Enablement	Concession	Summary
Motivation	Condition	

Multi-nucleus relations

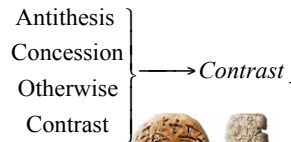
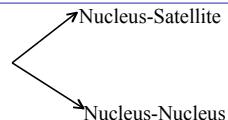
Sequence	Contrast	Joint
----------	----------	-------



Relation Number

The number of relations is somewhat arbitrary.
Mann and Thompson first proposed 15 relations, then 23.
It is possible to group and simplify them.

Symmetrical (nucleus-nucleus) and asymmetrical
relations (nucleus-satellite)



Group classes in a superclass



Definition of the Relations

The following text corresponds to an **evidence** relation that links a nucleus (segment 1) and a satellite (segment 2):

- ① *The program as published for calendar year 1980 really works.*
- ② *In only a few minutes, I entered all the figures from my 1980 tax return and got a result which agreed with my hand calculations to the penny.*

Mann and Thompson defined each relation in the RST model using a set of “constraints”.



Definition of the Relations (II)

Relation name	EVIDENCE
Constraints on the nucleus N	The reader R might not believe to a degree satisfactory to the writer W
Constraints on the satellite S	The reader believes S or will find it credible
Constraints on the $N + S$ combination	R 's comprehending S increases R 's belief of N
The effect	R 's belief of N is increased
Locus of the effect	N



Automatic Processing of Discourse

Is it possible to process automatically texts with these definitions?

And how can we do?

The description of an evidence relation is:

Reader believes Satellite or finds it credible

How can we measure this?



Cues in Text

The idea is to map a certain relation to certain words.

Words like *and*, *so*, *but*, *although*, and commas denote frontiers and ideas in a text.

The automatic text analysis uses these signs, *cues*, *cue phrases*, to segment a text and recognize relations



Ambiguity

Cues are often be ambiguous. Example:

*[Karl **and** Jan came to the lecture] [**and** asked questions]*

The first *and* has a syntactic role only. The second one defines a sequence

We must use supplementary constraints like position constraints between spans to carry out the analysis



Events

Research on the representation of time, events, and temporal relations dates back the beginning of logic.

It resulted in an impressive number of formulations and models.

A possible approach is to **reify** events: turn them into objects, quantify them existentially, and connect them using predicates

John saw Mary in London on Tuesday

$$\exists \epsilon [\text{saw}(\epsilon, \text{John}, \text{Mary}) \wedge \text{place}(\epsilon, \text{London}) \wedge \text{time}(\epsilon, \text{Tuesday})],$$

where ϵ represents the event.



Event Types

Events are closely related to sentence's main verbs

Different classifications have been proposed to associate a verb with a type of event, Vendler (1967):

- A state – a permanent property or a usual situation (e.g. *be, have, know, think*);
- An achievement – a state change, a transition, occurring at single moment (e.g. *find, realize, learn*);
- An activity – a continuous process taking place over a period of time (e.g. *work, read, sleep*). In English, activities often use the present perfect *-ing*;
- An accomplishment – an activity with a definite endpoint completed by a result (e.g. *write a book, eat an apple*).



Temporal Representation of Events (Allen 1983)

#	Relations	#	Inverse relations	Graphical representations
1.	before(a, b)	2.	after(b, a)	
3.	meets(a, b)	4.	met_by(b, a)	
5.	overlaps(a, b)	6.	overlapped_by(b, a)	
7.	starts(a, b)	8.	started_by(b, a)	
9.	during(b, a)	10.	contains(a, b)	
11.	finishes(b, a)	12.	finished_by(a, b)	
13.	equals(a, b)			

TimeML, an Annotation Scheme for Time and Events

TimeML is an effort to unify temporal annotation, based on Allen's (1984) relations and inspired by Vendler's (1967) classification.

TimeML defines the XML elements:

- TIMEX3 to annotate time expressions (at four o'clock),
- EVENT, to annotate the events (he slept),
- "signals".

The SIGNAL tag marks words or phrases indicating a temporal relation.



TimeML, an Annotation Scheme for Time and Events (II)

TimeML connects entities using different types of links

Temporal links, TLINKs, describe the temporal relation holding between events or between an event and a time.

TimeML elements have attributes. For instance, events have a tense, an aspect, and a class.

The 7 possible classes denote the type of event, whether it is a STATE, an instantaneous event (OCCURRENCE), etc.



TimeML Example

All 75 people on board the Aeroflot Airbus died when it ploughed into a Siberian mountain in March 1994

(Ingria and Pustejovsky 2004):

All 75 people

```
<EVENT eid="e7" class="STATE">on board</EVENT>
```

```
<MAKEINSTANCE eiid="ei7" eventID="e7" tense="NONE" aspect="NONE"/>
```

```
<TLINK eventInstanceID="ei7" relatedToEvent="ei5"  
relType="INCLUDES"/>
```

the Aeroflot Airbus

```
<EVENT eid="e5" class="OCCURRENCE" >died</EVENT>
```

```
<MAKEINSTANCE eiid="ei5" eventID="e5" tense="PAST" aspect="NONE"/>
```

```
<TLINK eventInstanceID="ei5" signalID="s2" relatedToEvent="ei6"  
relType="IAFTER"/>
```



TimeML Example

All 75 people on board the Aeroflot Airbus died when it ploughed into a Siberian mountain in March 1994

(Ingria and Pustejovsky 2004):

```
<SIGNAL sid="s2">when</SIGNAL>
it
<EVENT eid="e6" class="OCCURRENCE">ploughed</EVENT>
<MAKEINSTANCE eiid="ei6" eventID="e6" tense="PAST" aspect="NONE"/>
<TLINK eventInstanceID="ei6" signalID="s3" relatedToTime="t2"
relType="IS_INCLUDED"/>
<TLINK eventInstanceID="ei6" relatedToEvent="ei4"
relType="IDENTITY"/>
into a Siberian mountain
<SIGNAL sid="s3">in</SIGNAL>
<TIMEX3 tid="t2" type="DATE" value="1994-04">March 1994</TIMEX3>
```

