

## A List of heuristic rules

Table 1 lists the methods used to identify and locate the signal tokens based on the annotated signal labels and the gold lexical and syntactic parse provided in Penn Treebank (Marcus et al., 1993). Certain signal labels are supplemented with optional *meta comment* that specifies the signal in the particular relation.

Signals of the relation are identified from the *elementary discourse units* (EDUs) related by the relation. As a result, a word-level annotation of the discourse signals is obtained. For example:

### Tokenized text

In the first year , the bank eliminated 800 jobs .  
Now it says it will trim more in the next year .

### RST Treebank

Relation			
Index	Nucleus	Satellite	Sense
R1	word 1-11	word 12-23	Temporal
...	...	...	...

### RST Signaling Corpus

Signal Index	Relation Index	Signal label	Meta comment
S1	R1	(13) discourse marker	now
S2	R1	(12) tense	
S3	R1	(25) lexical chain	first year – next year
...	...	...	...

Using the heuristics in Table 1, the signal tokens are identified (the underlined words below):

In the first year , the bank eliminated 800 jobs.  
Now it says it will trim more in the next year .

Each of the signal tokens is tagged by the corresponding signal index, leading to below word-level annotation:

-- 3 3 --- 2 --- 1 - 2 - 2 2 --- 3 3 -

and the relation boundaries and senses can be retrieved with reference to the RST Signaling Corpus, and, in turn, the RST Discourse Treebank.

The converted annotation is available on <http://www.coli.uni-saarland.de/~frances/contents/rstsignal/rstsignal.html>.

## References

Mitchell P Marcus, Mary Ann Marcinkiewicz, and Beatrice Santorini. 1993. Building a large annotated corpus of english: The penn treebank. *Computational linguistics*, 19(2):313–330.

Category	Signal label	Method	
Syntactic	1	relative clause	identify POS/syntactic patterns: onset of SBAR, e.g. <i>that</i>
	2	infinitive clause	POS patterns: TO VB
	3	present participial clause	first VBG in VP
	4	past participial clause	first VBN in VP
	5	imperative clause	manually identify imperative verb
	6	nominal modifier	POS patterns: TO VB
	7	adjective modifier	POS patterns: TO VB
	8*	parrallel syntactic construction	–
	9	subject-auxillary inversion	manually identify subject-auxillary inversion
	10*	interrupted matrix	–
	11	reported speech	identify a list of verbs, e.g. <i>said, according to</i>
Morphological	12	tense	identify verbs in different tenses
Discoure marker	13	one of 201 defined markers	identify labelled discourse marker e.g. <i>when, in addition to</i>
Reference	14	personal reference	identify tokens/phrases specified in meta comment: e.g. <i>Johnson – she</i>
	15	comparative reference	e.g. <i>equal</i>
	16	demonstrative reference	e.g. <i>these</i>
	17	propositional reference	e.g. <i>naming a candidate – it</i>
Lexical	18	alternative expression	identify tokens/phrases specified in meta comment: e.g. <i>what’s more</i>
	19	indicative word	e.g. <i>compared with</i>
Semantic	20	synonymy	identify tokens/phrases specified in meta comment: e.g. <i>United Airlines – UAL</i>
	21	antonymy	e.g. <i>short–long</i>
	22	meronymy	e.g. <i>Californians–Johnson</i>
	23	repetition	the word / phrase that is repeated
	24	indicative word pair	e.g. <i>asked–replied</i>
	25	lexical chain	e.g. <i>company–sold–shares–holdings</i> unspecific comments, e.g. <i>a few lexical chains</i> , are excluded
	26	general word	e.g. <i>issues</i>
Numerical	27	same count	identify CD
Graphical	28	colon	identify :
	29	semi-colon	identify ;
	30	dash	identify - or –
	31	comma	identify ,
	32	parentheses	identify -LRB-, -RRB-, -LSB-, -RSB-, -LCB-, -RCB-
	33*	items in sequence	–
Genre	34*	newspaper style attribution	–
	35*	newspaper style contrast	–
	36*	newspaper style elaboration	–
	37*	newspaper layout	–
	38*	newspaper style definition	–
	39*	inverted pyramid cheme	–
Unsure	40*	no signals identified	–

Table 1: The heuristics used to identify and locate the signal tokens based on the annotated signal labels. Signals marked with \* are excluded in analysis.