

Intercoder agreement in discourse analysis

Seminar in methodology and statistics

13th May 2009

Ildikó Berzlánovich, Myrthe Faber

University of Groningen

Center for Language and Cognition Groningen

Outline

PART 1: Intercoder agreement for the study of discourse structure (Ildikó)

PART 2: Intercoder agreement in conversation analysis (Myrthe)

Outline – Discourse structure

- Own research
- Annotation problems
- Intercoder agreement measures (percentage agreement, Cohen's kappa)
- Practice so far
- Decisions for own research

Own research

Annotation problems
Measures

Practice so far
Decisions for own research

General aim:

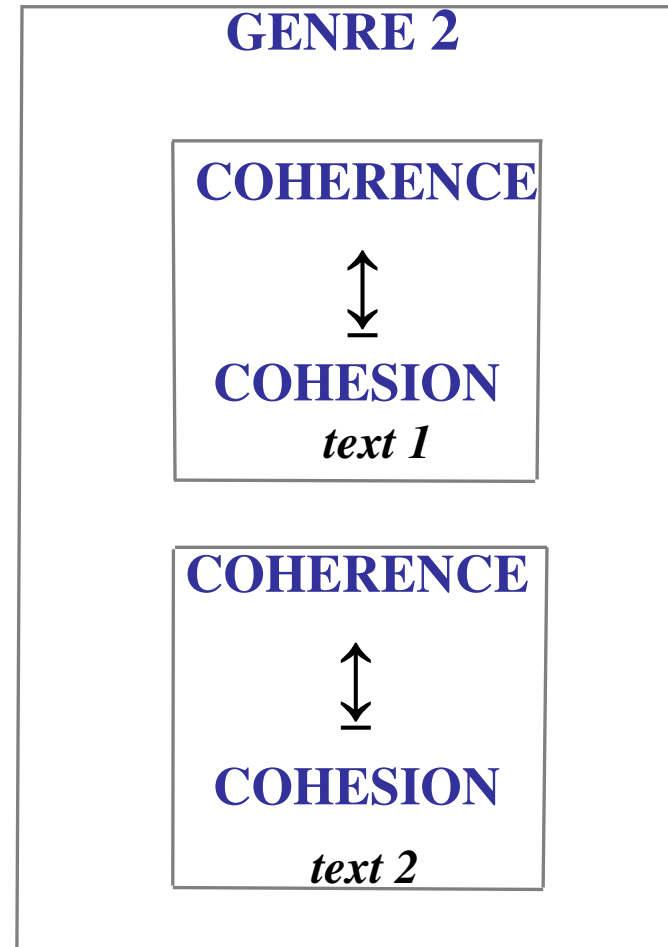
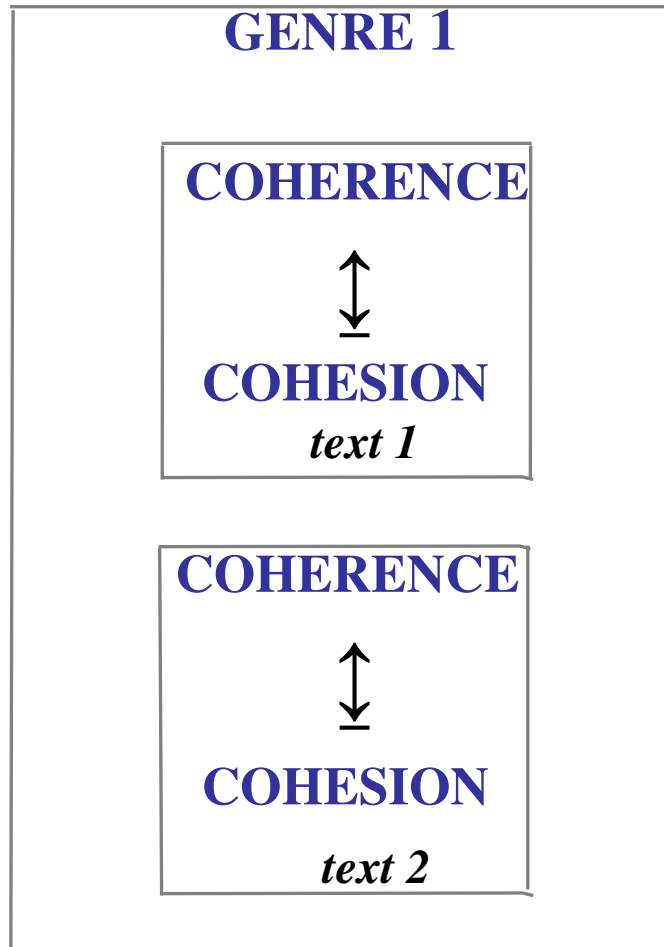
interaction between coherence and lexical cohesion across genres

Alignment hypothesis:

Lexical cohesion is more closely aligned with coherence in thematically organized texts than in intentionally organized texts.

Specifically:

- › close alignment in expository texts
- › less or no alignment in persuasive texts



Own research

Annotation problems
Measures

Practice so far
Decisions for own research

Genre

- class of communicative events with common communicative purposes shared in a discourse community (Swales 1990)
- genre-specific move structure

Coherence

- underlying relations between discourse units in text

Cohesion

- semantic relations between surface elements in text

Own research

Annotation problems

Measures

Practice so far
Decisions for own research

Texts

- expository texts: encyclopedia entries (EE01, EE02)
- persuasive texts: fundraising letters (FL01, FL02)

Own research

Annotation problems
Measures

Practice so far
Decisions for own research

Encyclopedia entries

1. name the object
2. define the object
3. describe in general (e.g., size, age, category)
4. describe details (e.g., surface, past/future development, discovery)

Fundraising letters - seven moves (Upton 2002)

1. get attention
2. introduce the cause and/or establish credentials of organization
3. solicit response
4. offer incentive
5. reference insert
6. express gratitude
7. conclude with pleasantries

Fragment from EE01 (*De Zon*)

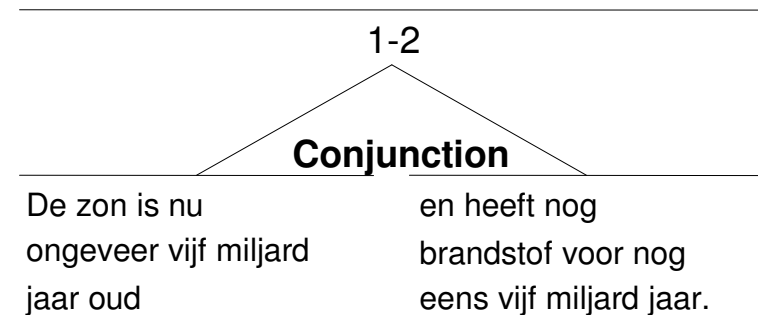
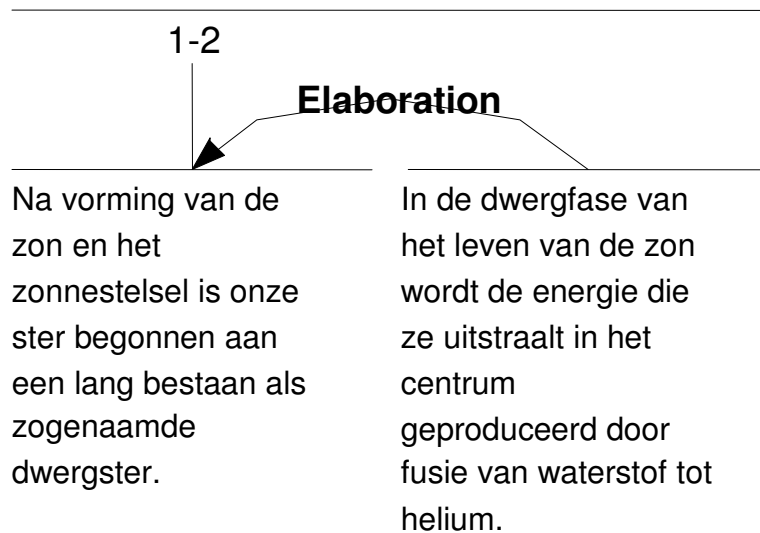
Na vorming van de zon en het zonnestelsel is onze ster begonnen aan een lang bestaan als zogenaamde dwergster. In de dwergfase van het leven van de zon wordt de energie die ze uitstraalt in het centrum geproduceerd door fusie van waterstof tot helium. De zon is nu ongeveer vijf miljard jaar oud en heeft nog brandstof voor nog eens vijf miljard jaar.

After the forming of the sun and the solar system, our star began its long existence as a so-called dwarf star. In the dwarf phase of its life, the energy that the sun gives off is generated in its core through the fusion of hydrogen into helium. The sun is about five billion years old now and it still has enough fuel for another five billion years.

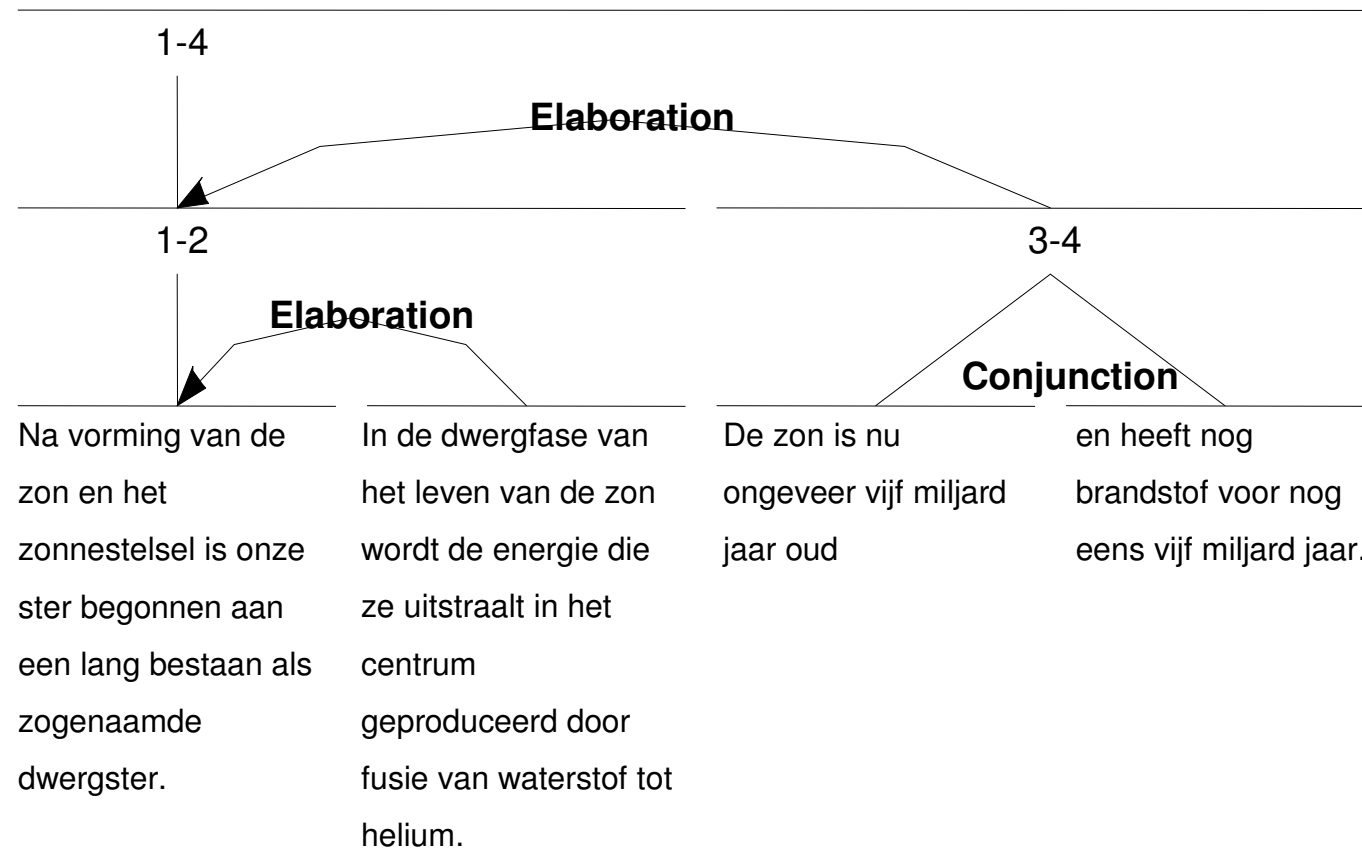
Rhetorical Structure Theory (Mann & Thompson 1988)

- functional relations between propositions
- reconstruction of writer's purposes
- subject matter vs. presentational relations
- mononuclear relations

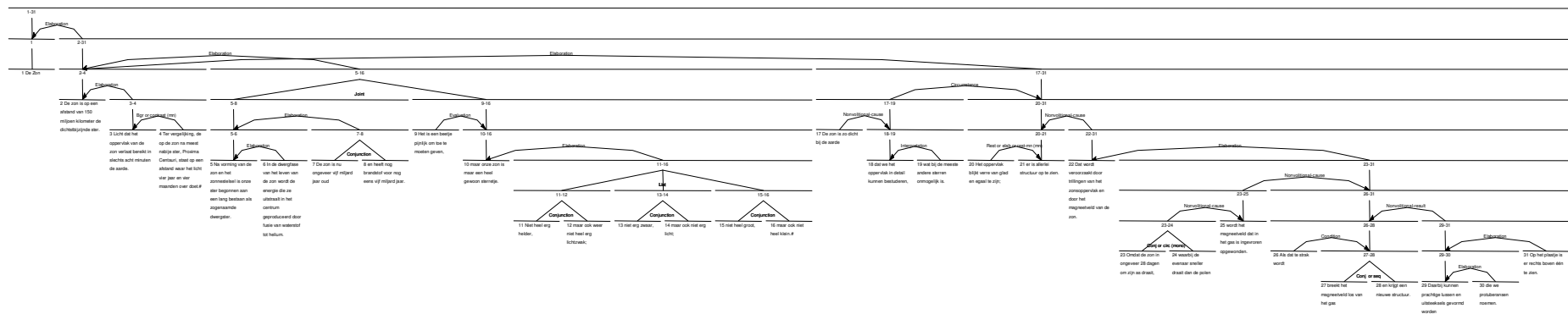
multinuclear relations



hierarchy



hierarchy



Own research

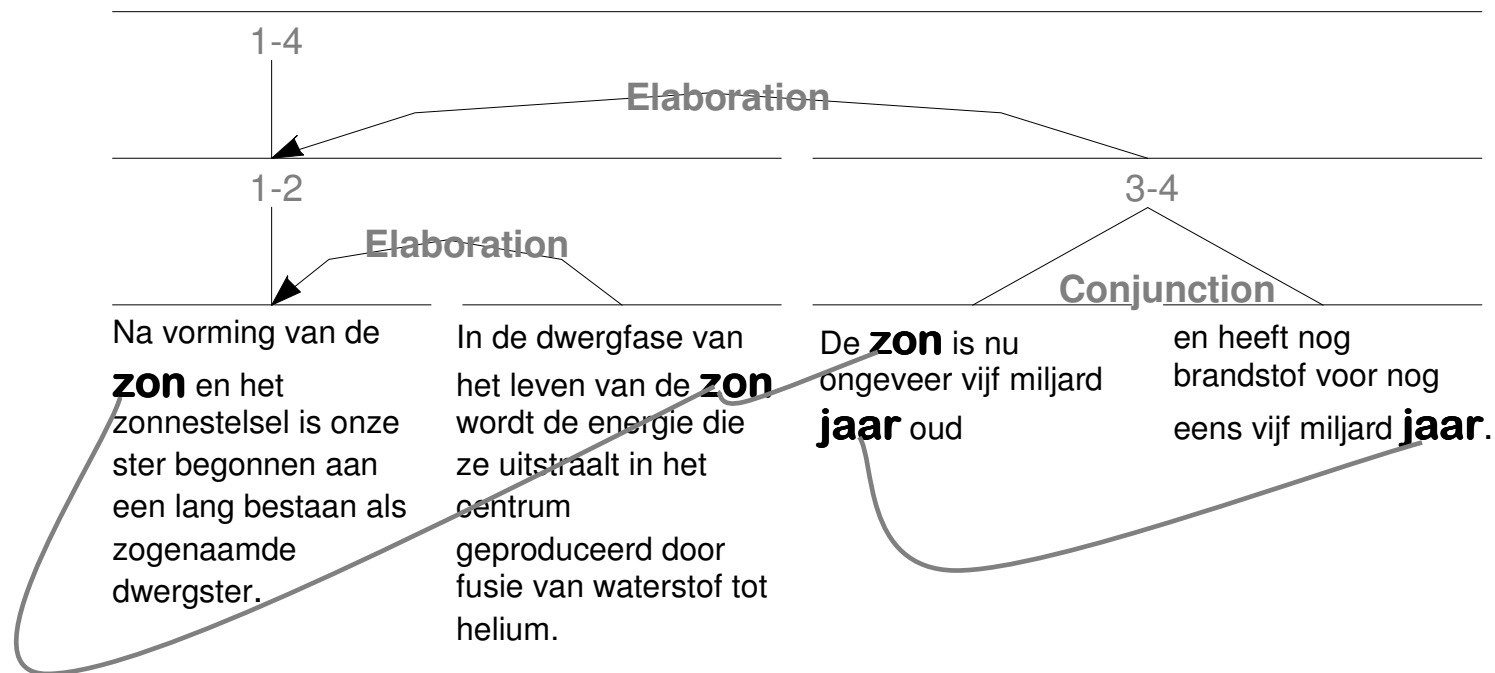
Annotation problems

Measures

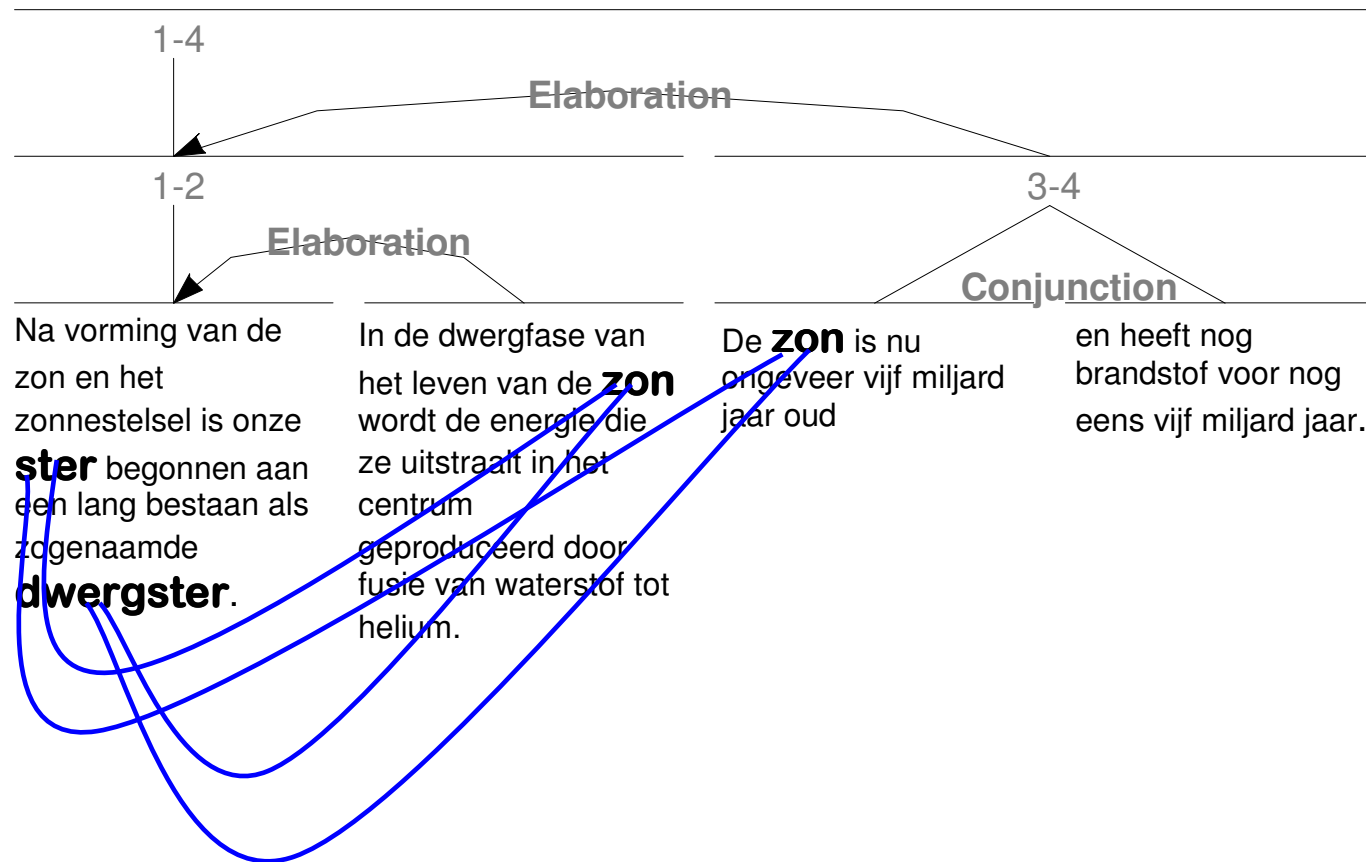
Practice so far
Decisions for own research

- lexical cohesion
- network of relations
- lexical cohesive relations
 - repetition
 - systematic semantic relations
 - hyponymy, hyperonymy, co-hyponymy
 - meronymy, holonymy, co-meronymy
 - synonymy
 - antonymy
 - collocation

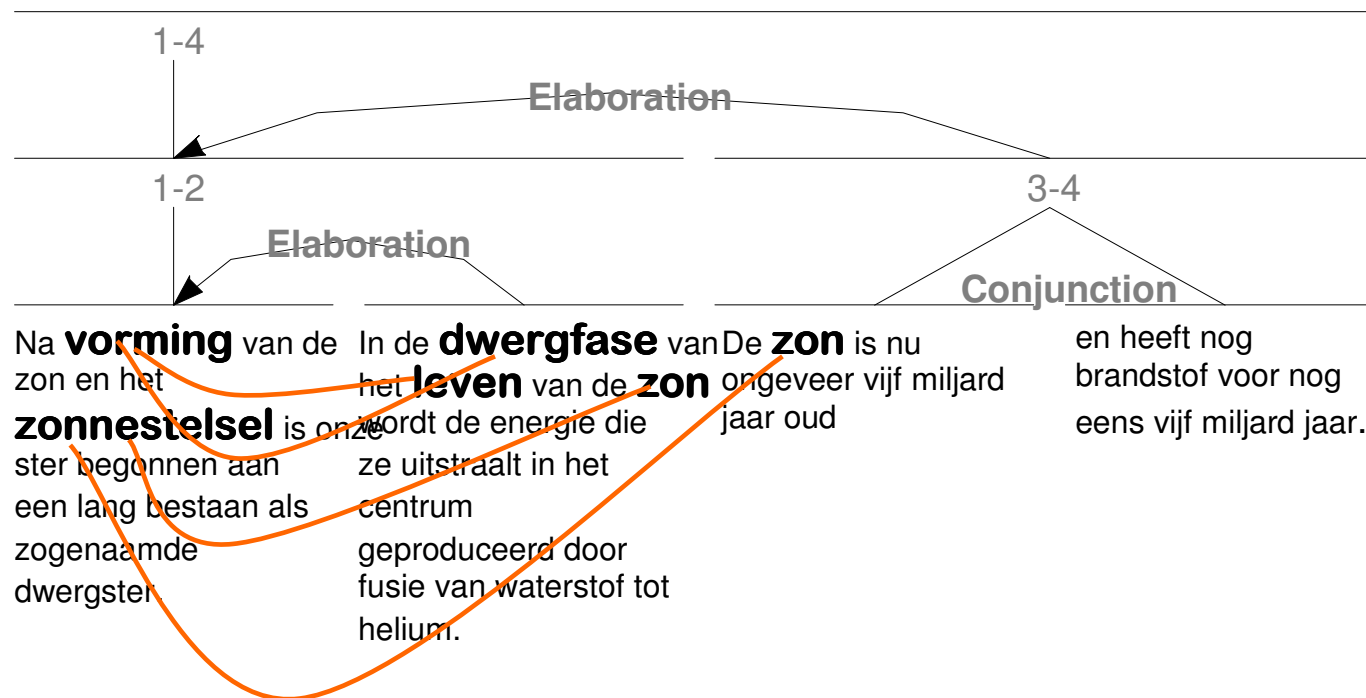
Repetition



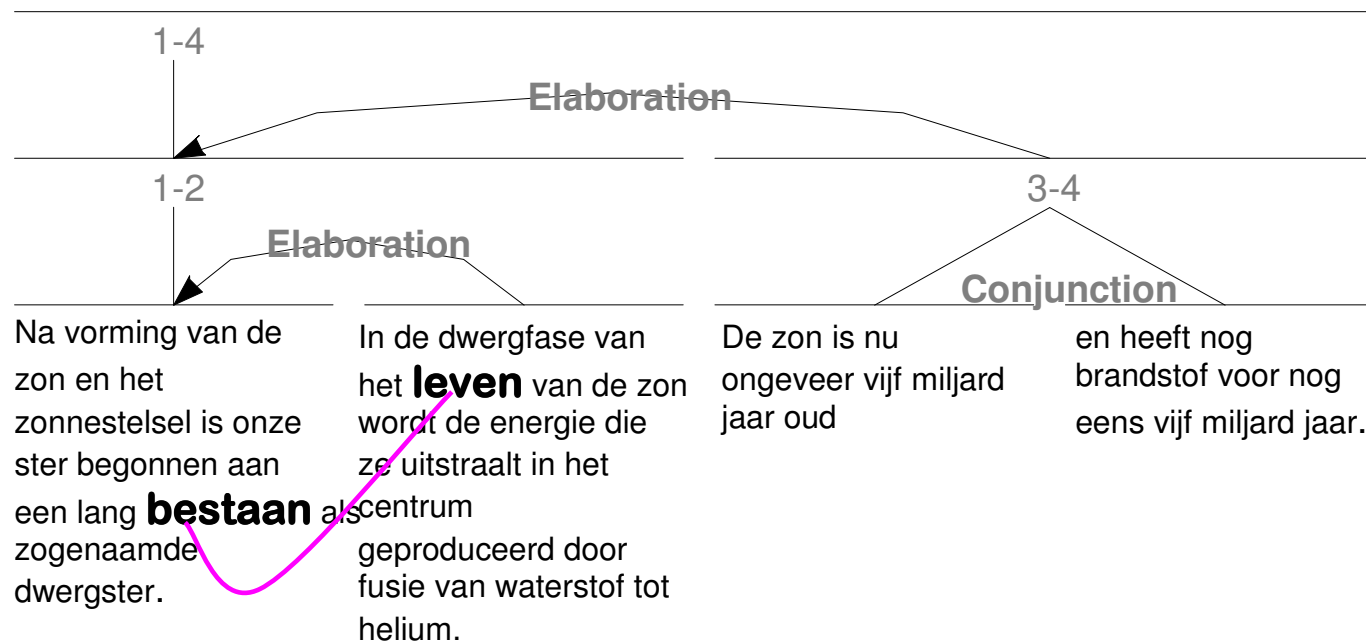
Hyponymy



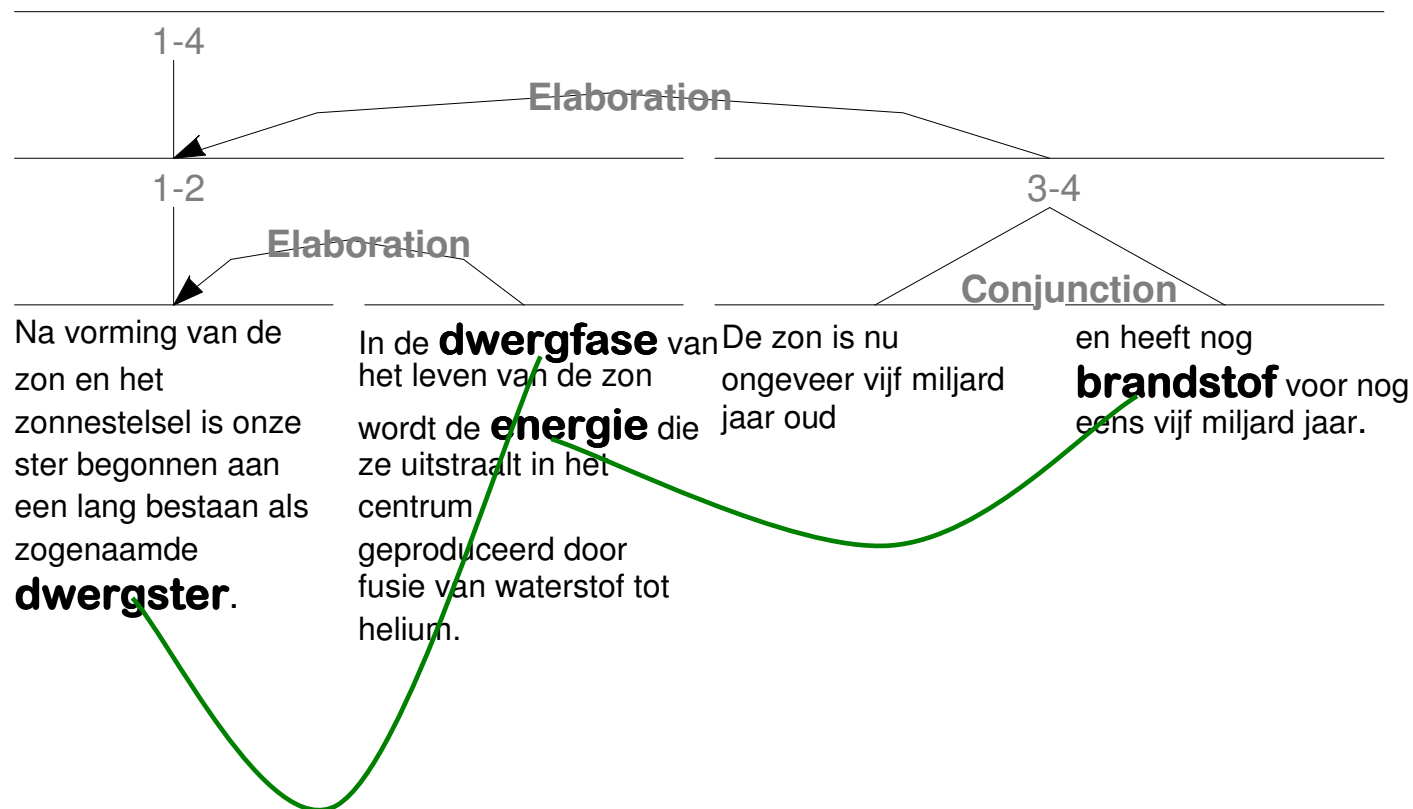
Meronymy

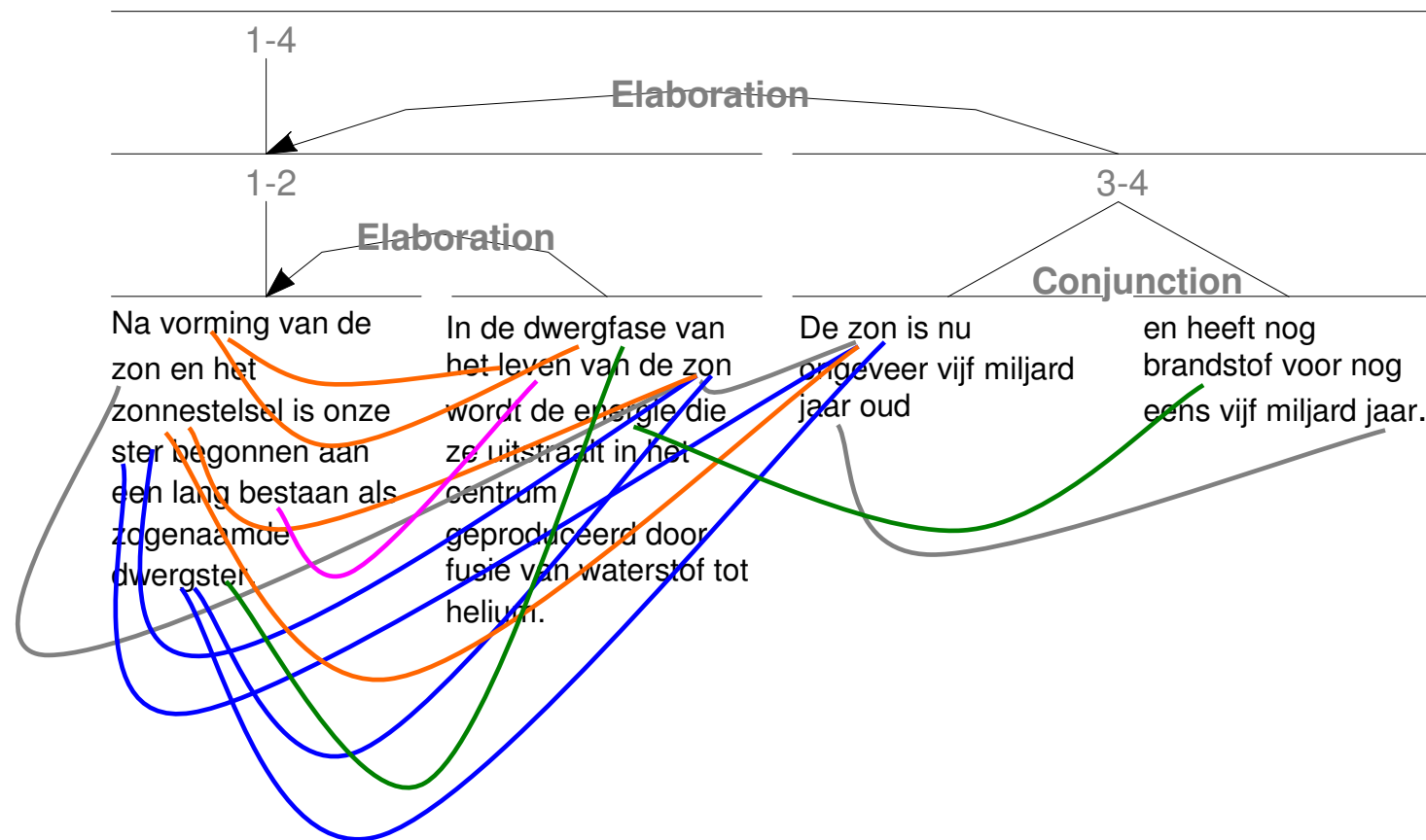


Synonymy



Collocation





Own research
Annotation problems
Measures

Practice so far
Decisions for own research

Centrality of discourse units

EDUs

Moves

*in coherence
structure*

smallest units
in RST tree

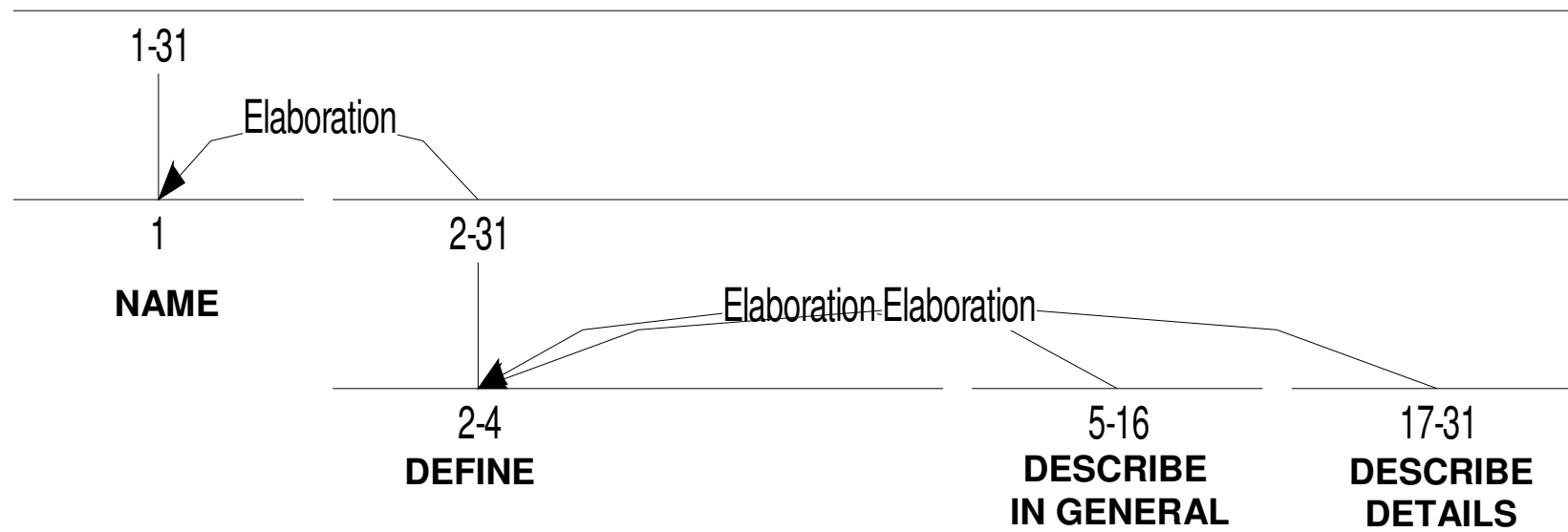
top level
of RST tree

*in lexical
cohesion*

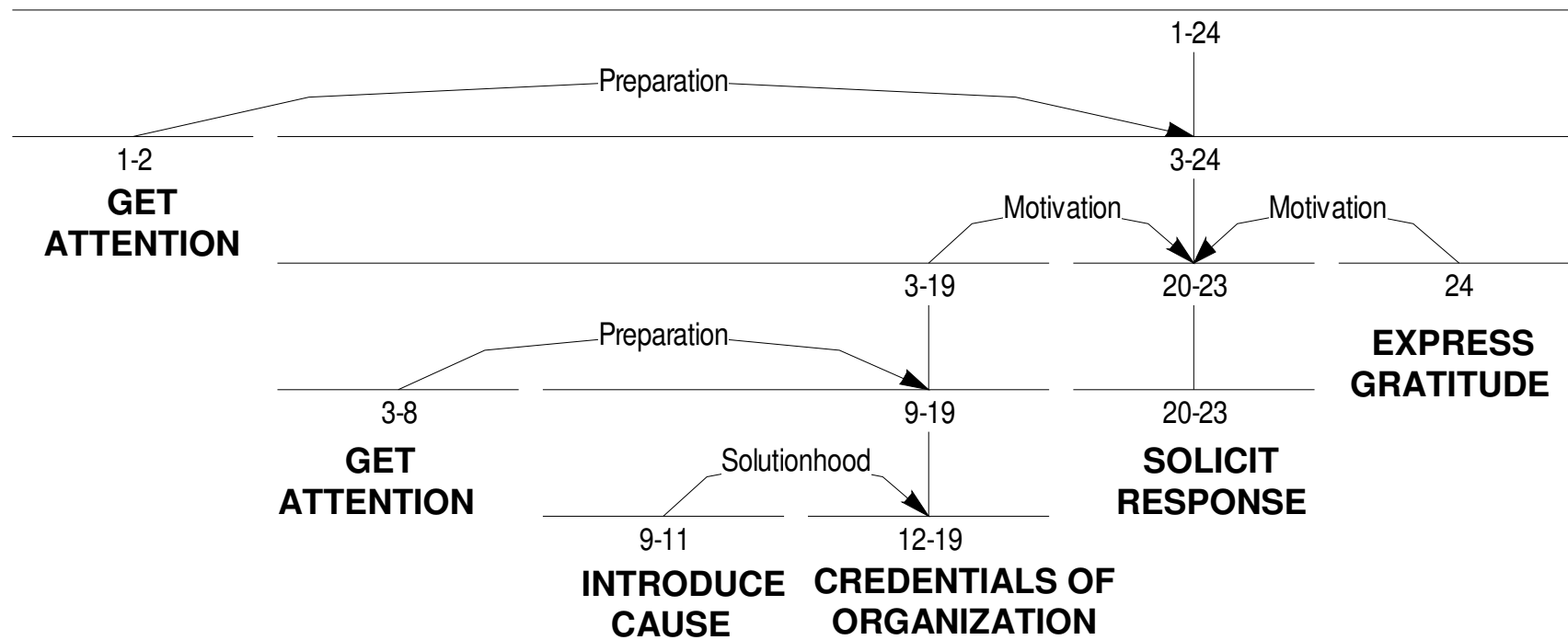
lexical cohesive links
per EDU

average lexical
cohesion density

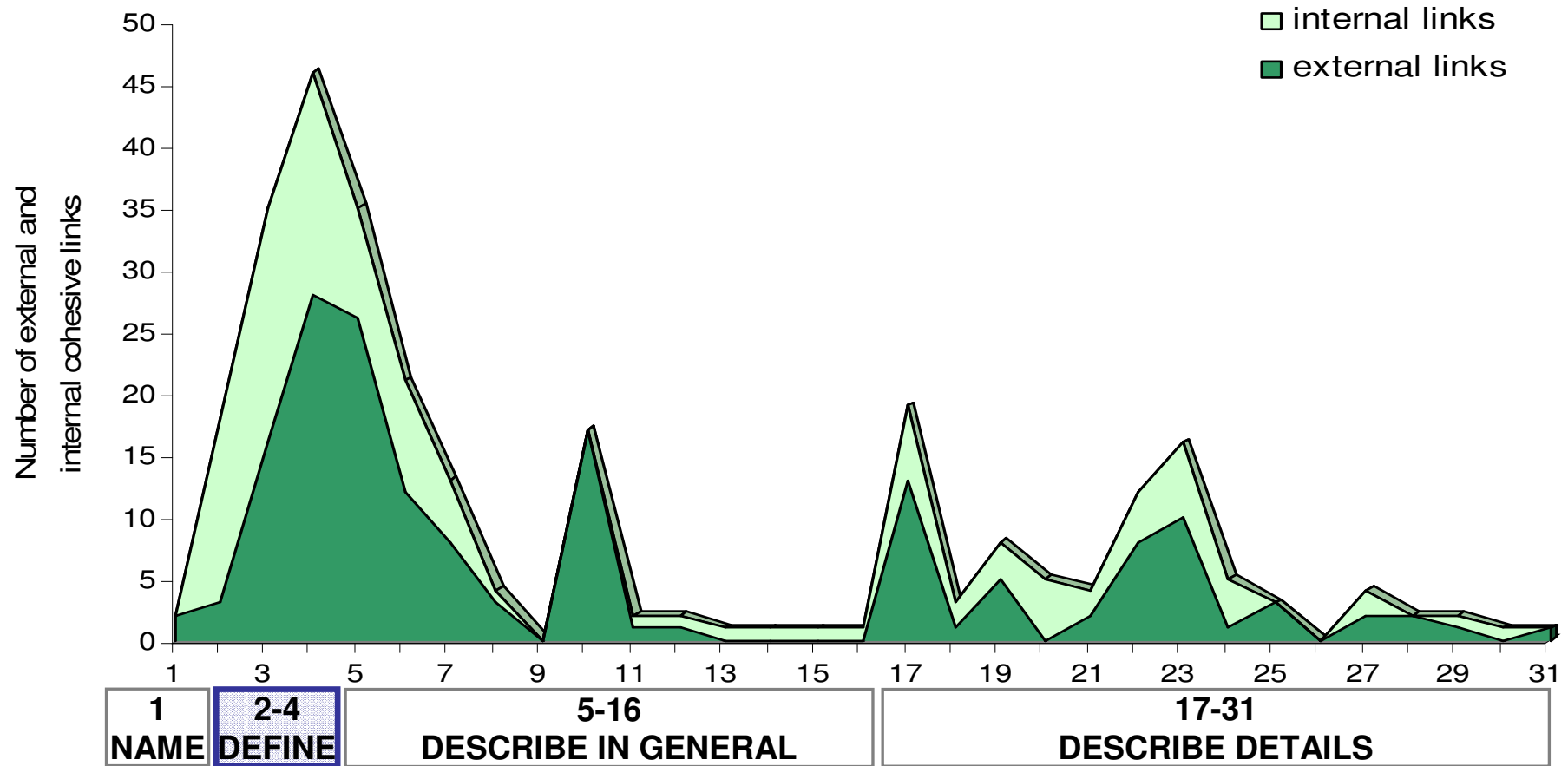
EE01



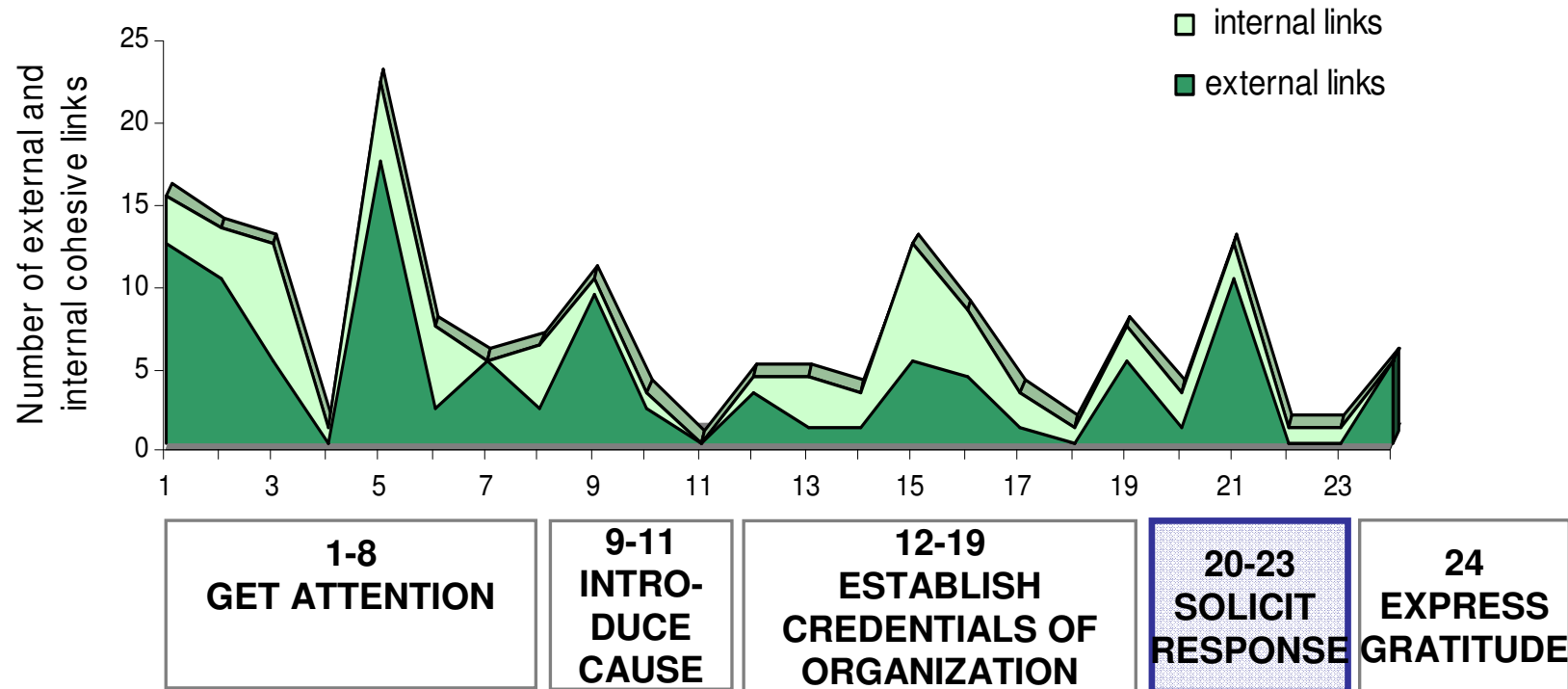
FL01



EE01



FL01



Own research

Annotation problems

Measures

Practice so far

Decisions for own research

- EE: *Define* is the central move both in coherence and in lexical cohesion
- FL: *Solicit response* is the central move in coherence, but no central move was found for lexical cohesion

→ alignment much closer for EE than for FL

Own research

Annotation problems

Measures

Practice so far

Decisions for own research

- Move analysis
- Segmentation
- Cohesion analysis
- Coherence analysis

Structural segments + moves

- [Met vriendelijke groet] [en alvast heel hartelijk dank,]

Overlapping moves

- *Dankzij uw donaties* aan de Nierstichting kunnen wij ons inzetten voor een betere kwaliteit van leven voor nierpatienten.

Two functions in one segment

- De *bijgevoegde* kaarten zijn een *bedankje* voor het lezen van mijn brief over de malaria-epidemie in Afrika.

EDU (elementary discourse unit) = clause

[Na vorming van de zon en het zonnestelsel is onze ster begonnen aan een lang bestaan als zogenaamde dwergster.] [In de dwergfase van het leven van de zon wordt de energie die ze uitstraalt in het centrum geproduceerd door fusie van waterstof tot helium.] [De zon is nu ongeveer vijf miljard jaar oud] [en heeft nog brandstof voor nog eens vijf miljard jaar.]

EDU = clause

[Na vorming van de zon en het zonnestelsel is onze ster begonnen aan een lang bestaan als zogenaamde dwergster.] [In de dwergfase van het leven van de zon wordt de energie die ze uitstraalt in het centrum geproduceerd door fusie van waterstof tot helium.] [De zon is nu ongeveer vijf miljard jaar oud] [en heeft nog brandstof voor nog eens vijf miljard jaar.]

EDU = clause

[Na vorming van de zon en het zonnestelsel is onze ster begonnen aan een lang bestaan als zogenaamde dwergster.] [In de dwergfase van het leven van de zon wordt de energie die ze uitstraalt in het centrum geproduceerd door fusie van waterstof tot helium.] [De zon is nu ongeveer vijf miljard jaar oud] [en heeft nog brandstof voor nog eens vijf miljard jaar.]

EDU = clause

[Na vorming van de zon en het zonnestelsel is onze ster begonnen aan een lang bestaan als zogenaamde dwergster.] [In de dwergfase van het leven van de zon wordt de energie die ze uitstraalt in het centrum geproduceerd door fusie van waterstof tot helium.] [De zon is nu ongeveer vijf miljard jaar oud] [en heeft nog brandstof voor nog eens vijf miljard jaar.]

[Na vorming van de zon en het zonnestelsel] [is onze ster begonnen aan een lang bestaan als zogenaamde dwergster.] [In de dwergfase van het leven van de zon wordt de energie die ze uitstraalt in het centrum geproduceerd] [door fusie van waterstof tot helium.] [De zon is nu ongeveer vijf miljard jaar oud] [en heeft nog brandstof voor nog eens vijf miljard jaar.]

Annotation problems – segmentation

Decisions for own research

Measures

Comparisons

- [De zon is zo dicht bij de aarde] [dat we het oppervlak in detail kunnen bestuderen.]
- [Een neutronenster is ongeveer anderhalf keer zo zwaar als de zon,]
- [Echter, Saturnus produceert meer licht] [dan hij van de zon ontvangt.]
- [De atmosferische druk is op het oppervlak zo'n 90 keer groter dan op Aarde.]

Embedded EDUs

- 12 [In kraters nabij de polen van Mercurius, [...13...] bestaat misschien zelfs ijs. / 13 [waar nooit zonlicht komt,]

Parentheticals

- 14 De binnenste maan [...15...] beweegt iets sneller dan de buitenste / 15 [(van 2002 tot 2005 is dat Epimetheus)]

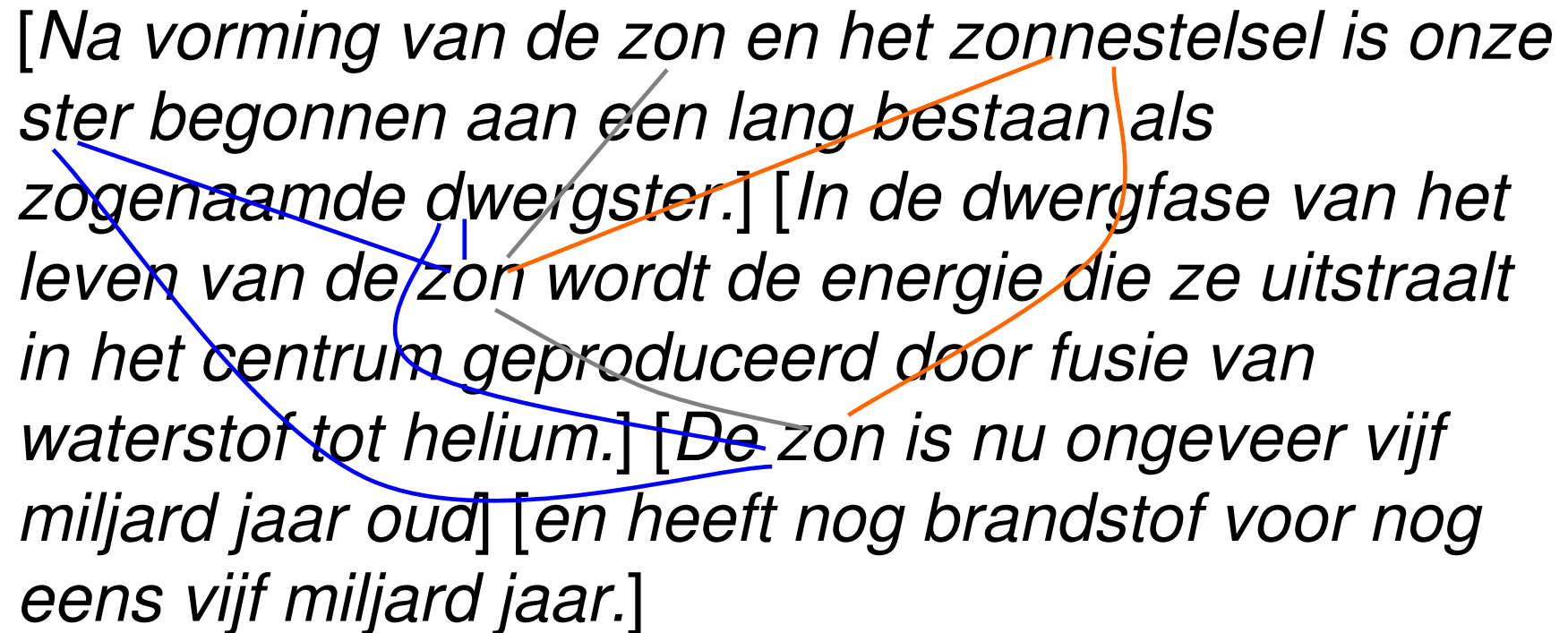
[Na vorming van de zon en het zonnestelsel is onze ster begonnen aan een lang bestaan als zogenaamde dwergster.] [In de dwergfase van het leven van de zon wordt de energie die ze uitstraalt in het centrum geproduceerd door fusie van waterstof tot helium.] [De zon is nu ongeveer vijf miljard jaar oud] [en heeft nog brandstof voor nog eens vijf miljard jaar.]

→ link only with the closest preceding item

[Na vorming van de zon en het zonnestelsel is onze ster begonnen aan een lang bestaan als zogenaamde dwergster.] [In de dwergfase van het leven van de zon wordt de energie die ze uitstraalt in het centrum geproduceerd door fusie van waterstof tot helium.] [De zon is nu ongeveer vijf miljard jaar oud] [en heeft nog brandstof voor nog eens vijf miljard jaar.]

→ all links identified (same item to more preceding items)

[Na vorming van de zon en het zonnestelsel is onze ster begonnen aan een lang bestaan als zogenaamde dwergster.] [In de dwergfase van het leven van de zon wordt de energie die ze uitstraalt in het centrum geproduceerd door fusie van waterstof tot helium.] [De zon is nu ongeveer vijf miljard jaar oud] [en heeft nog brandstof voor nog eens vijf miljard jaar.]



→ all links identified (same preceding item to more succeeding items)

Annotation problems – lexical cohesion

Decisions for own research

Measures

Multiple relations

- *gasplaneet – Aarde*: co-meronymy OR co-hyponymy

Context

- *Mercurius – as / polen / krater*

Word forms

- *Aards – Aarde*: repetition

Abbreviations

- *€ – euro*: repetition; *€ – bedrag*: collocation; *H₂ – waterstof*: synonymy

Multi-word units

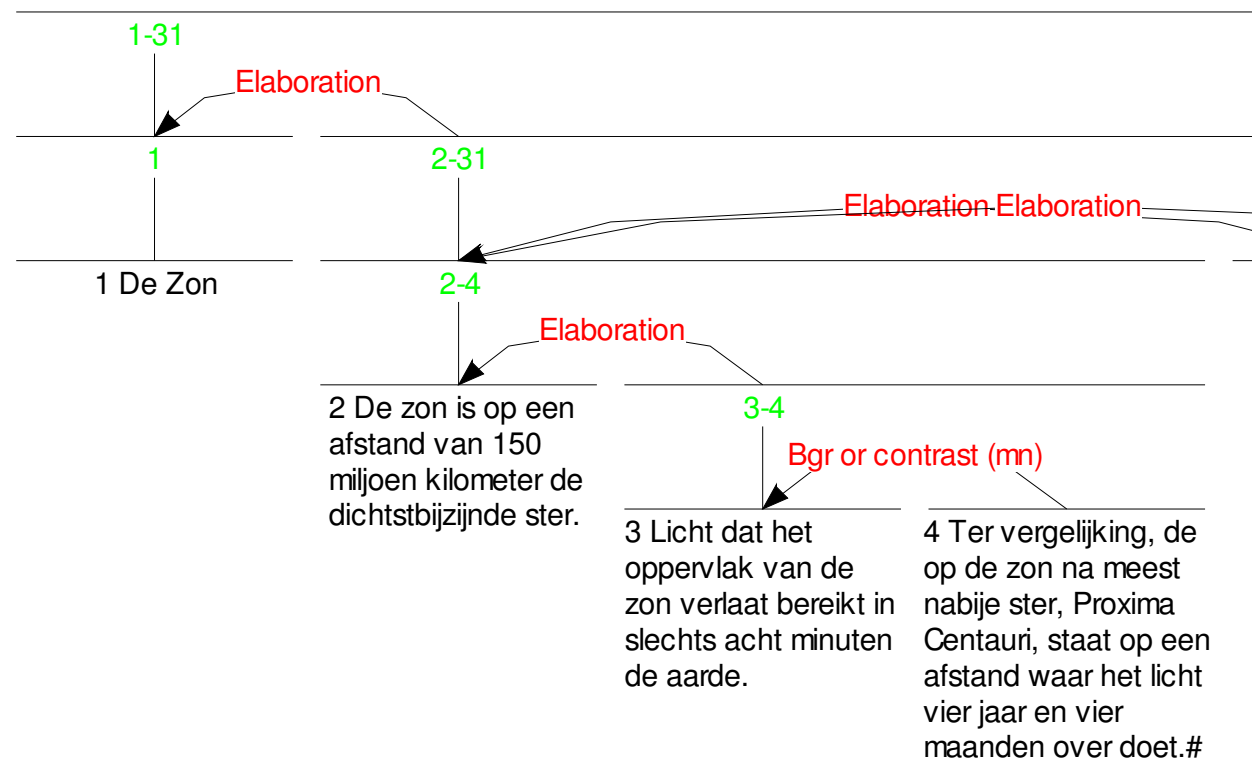
- *dwergster – dwergfase*: collocation; *Proxima Centauri*;
- [De meest spectaculaire structuur op het oppervlak van Jupiter is een *grote rode vlek* iets ten zuiden van de evenaar.] [Deze structuur wordt de *Grote Rode Vlek* genoemd.]

Annotation problems – coherence

Decisions for own research

Measures

- RST allows parallel analyses
- hierarchical structure (decisions of the annotator at a given step affects the decisions made at subsequent steps)



Measures

- for annotation purposes
 - reliability of data
 - validity of coding scheme
 - reproducibility of coding
- detailed reference manual
- no tradition to report agreement
 - about 10% of the corpus
 - 2 coders / more than 2 coders
 - naive vs. trained coders

Measures – percentage agreement

= observed agreement

		CODER A		
		STAT	IREQ	TOTAL
CODER B	STAT	20	20	40
	IREQ	10	50	60
	TOTAL	30	70	100

$$A_o = (20 + 50) / 100 = 0.7$$

- not correct for chance agreement (no comparability, biased)
- not correct for distribution of items among categories

Own research
Annotation problems

Practice so far
Decisions for own research

Measures – percentage agreement

Time Period	Ratings by Psychologist 1	Ratings by Psychologist 2
1	0	0
2	0	0
3	0	0
4	0	0
5	1	0
6	0	0
7	0	0
8	0	0
9	0	0
10	0	0
11	0	0
12	0	1
13	0	0
14	0	0
15	0	0
16	0	0
17	0	0
18	0	0
19	0	0
20	0	0

$$A_o = 18/20 = 90\%$$

$$A_e = 1/20 \times 1/20 + 19/20 \times 19/20 = \\ .0025 + .9025 = \\ .9050$$

Measures – Cohen’s kappa

$$\kappa = \frac{A_o - A_e}{1 - A_e}$$

$$A_o = (20+50) / 100 = 0.70$$

$$A_e = (30/100) \times (40/100) + (70/100) \times (60/100) = 0.54$$

$$\kappa = (0.70 - 0.54) / (1 - 0.54) = 0.348$$

$$P(k|c_i) = \hat{P}(k|c_i) = \frac{n_{c_i k}}{i}$$

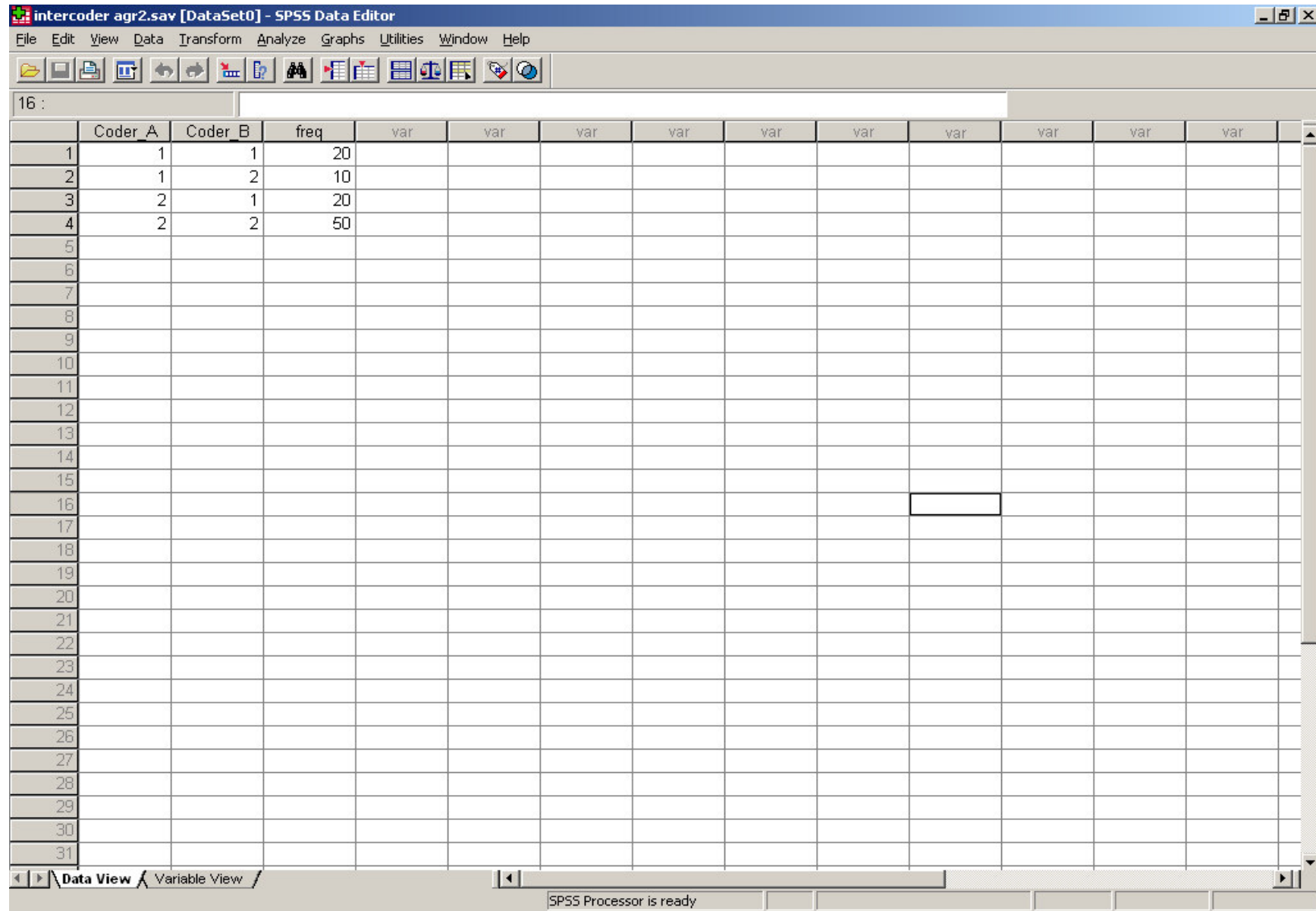
		CODER A		
		STAT	IREQ	TOTAL
CODER B	STAT	20	20	40
	IREQ	10	50	60
	TOTAL	30	70	100

- category judgments
- chance-corrected
- individual coder distributions
- $-1 < K < 1$

Own research
Annotation problems

Practice so far
Decisions for own research

Measures – Cohen's kappa



The screenshot shows the SPSS Data Editor window for a file named 'intercoder agr2.sav [DataSet0]'. The window title bar includes standard OS controls (minimize, maximize, close) and the menu bar contains 'File', 'Edit', 'View', 'Data', 'Transform', 'Analyze', 'Graphs', 'Utilities', 'Window', and 'Help'. Below the menu bar is a toolbar with various icons for file operations and data manipulation. The main area displays a data grid with 31 rows and 14 columns. The first four columns are labeled 'Coder_A', 'Coder_B', and 'freq', with the remaining columns labeled 'var'. The data is as follows:

	Coder_A	Coder_B	freq	var	var	var	var	var	var	var	var	var	var
1	1	1	20										
2	1	2	10										
3	2	1	20										
4	2	2	50										
5													
6													
7													
8													
9													
10													
11													
12													
13													
14													
15													
16													
17													
18													
19													
20													
21													
22													
23													
24													
25													
26													
27													
28													
29													
30													
31													

The status bar at the bottom indicates 'SPSS Processor is ready'.

Own research
 Annotation problems

Practice so far
 Decisions for own research

Measures – Cohen’s kappa

Coder_A * Coder_B Crosstabulation

Count

		Coder_B		Total
		statement	info-req	
Coder_A	statement	20	10	30
	info-req	20	50	70
Total		40	60	100

Symmetric Measures

		Value	Asymp. Std. Error(a)	Approx. T(b)	Approx. Sig.
Measure of Agreement	Kappa	,348	,095	3,563	,000
N of Valid Cases		100			

a Not assuming the null hypothesis.

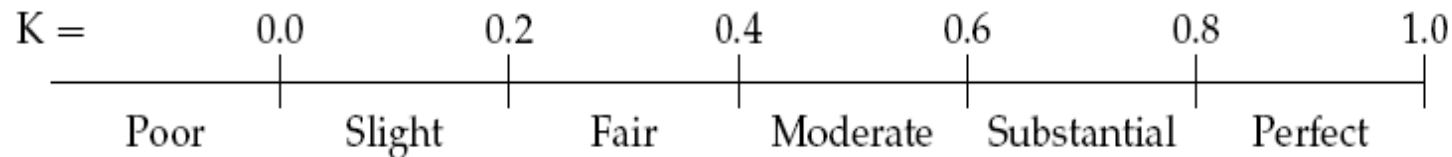
b Using the asymptotic standard error assuming the null hypothesis.

Measures – Cohen's kappa

Interpretation

Krippendorff (1980): $.67 < K < .80$ tentative; $.80 < K$ definite conclusions

Rietveld & van Hout (1993):



Craggs & Wood (2005): no general threshold

Artstein & Poesio (2008): .80 (.70)

Spooren & Degand (2009): .70 (for coherence relations)

Measures – Cohen’s kappa

Category prevalence

		Coder 2		
		Accept	Ack	
Coder 1	Accept	90	5	95
	Ack	5	0	5
		95	5	100

$$P(A) = 0.90, P(E) = 0.905$$

$$\kappa_{Co} = -0.048$$

		Coder 2		
		Accept	Ack	
Coder 1	Accept	45	5	50
	Ack	5	45	50
		50	50	100

$$P(A) = 0.90, P(E) = 0.5$$

$$\kappa_{Co} = 0.80$$

Bias

		Coder 2		
		Accept	Ack	
Coder 1	Accept	40	15	55
	Ack	20	25	45
		60	40	100

$$P(A) = 0.65, P(E) = 0.52$$

$$\kappa_{Co} = 0.27$$

		Coder 2		
		Accept	Ack	
Coder 1	Accept	40	35	75
	Ack	0	25	25
		40	60	100

$$P(A) = 0.65, P(E) = 0.45$$

$$\kappa_{Co} = 0.418$$

Measures

Other measures

- same distribution of coders → Scott's π
- more than two coders → multi- π , multi-K
- types of disagreements differentiated → weighted coefficients:
weighted K, Krippendorff's α

Own research
Annotation problems
Measures

Practice so far
Decisions for own research

- annotation of discourse structure
- complexity of annotation tasks
 - segmentation
 - move analysis
 - lexical cohesion analysis
 - coherence analysis

- proposition / sentence / clause / turn
- CL: automatic segmentation
- early studies: percentage agreement; later: K
- agreement on "bulk" of segments, but disagreement on exact boundaries
- broad vs. finer segments

Fewer boundaries, higher expected agreement.

Case 1: Broad segments
 $A_o = 0.96, A_e = 0.89, K = 0.65$

		CODER A		
		BOUNDARY	NO BOUNDARY	TOTAL
CODER B	BOUNDARY	2	1	3
	NO BOUNDARY	1	46	47
	TOTAL	3	47	50

Case 2: Fine discourse units
 $A_o = 0.88, A_e = 0.53, K = 0.75$

		CODER A		
		BOUNDARY	NO BOUNDARY	TOTAL
CODER B	BOUNDARY	16	3	19
	NO BOUNDARY	3	28	31
	TOTAL	19	31	50

Marcu et al. (1999)

- EDU boundaries for RST
- K calculated in two ways:
K_w: boundary can be after any word
K_u: boundary where at least one coder identified a boundary

$$K_w > K_u$$

- few studies: percentage agreement (“interrater reliability was calculated at 84%”; “the two raters had an agreement rate of 92% in identifying and categorizing the moves”)
- **de Groot (2008)**
 10% of corpus, 2 trained coders, Cohen’s kappa

Table 6.4: Evaluation of intercoder agreement for combinations of moves+strategies in the written texts and combinations of moves+strategies in the photographs.

<i>Section/ text type</i>	<i>Qualification</i>	<i>Mean Kappa</i>	<i>SD</i>	<i>Min. – max.</i>	<i>Agreement %</i>
Management statements/TEXT	Moderate	.58	0.04	.50 – .66	60.63
Management statements/PHOTO	Almost perfect	.92	0.17	.33 – 1.00	93.75
Profiles/ TEXT	Substantial	.73	0.14	.55 – .84	77.21
Profiles/PHOTO	Slight	.13	–	.13 – .13	14.29
Operational reviews/TEXT	Substantial	.67	0.09	.55 – .77	68.27
Operational reviews/PHOTO	Fair	.25	0.45	-.25 – .61	61.48

Own research
Annotation problems
Measures

Practice so far – lexical cohesion

Decisions for own research

- agreement not reported for lexical cohesion
- percentage agreement for word pairs
- word-sense tagging: rely on dictionaries, hierarchical tagsets (e.g., WordNet)
- ? corpus annotated for lexical cohesion

Spooren & Degand (2009)

- category labels for coherence relations bw two fragments
- two expert coders
- change in coding manual, fewer variables, more fragments
→ K improved (.60)

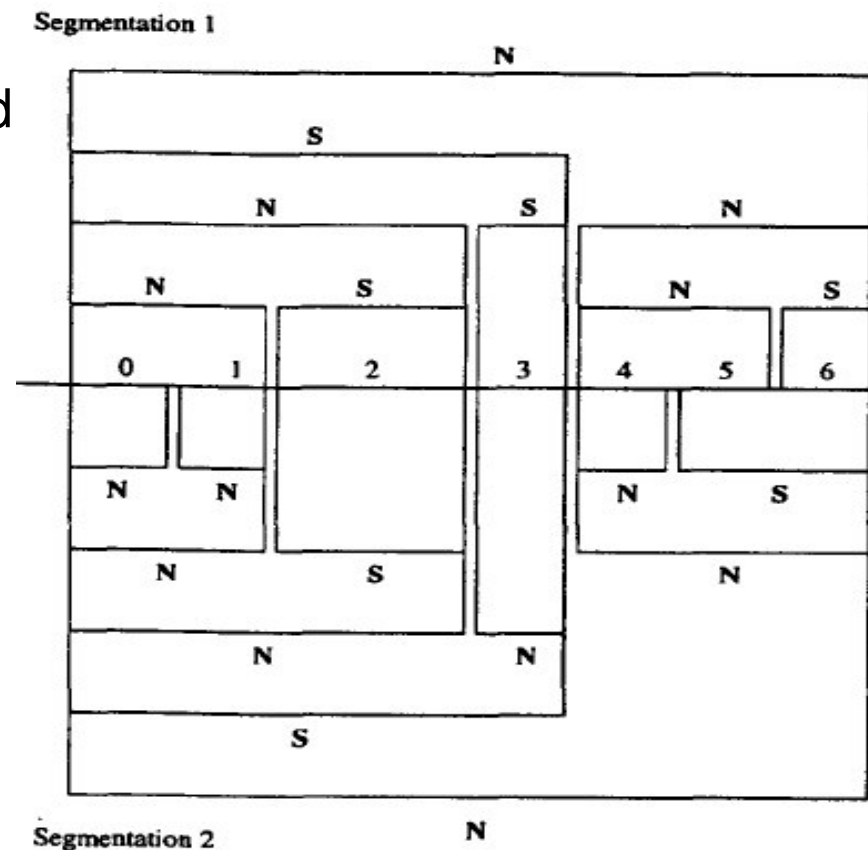
RST: more complex (N vs. S; relation labels; hierarchical tree structure)

Marcu et al. (1999), Carlson et al. (2003)

- for RST
- 2-3 expert coders
- long training phase
- method: mapping hierarchical structures into sets of units that are labeled with category judgments
- label for each EDU which is identified by at least one coder

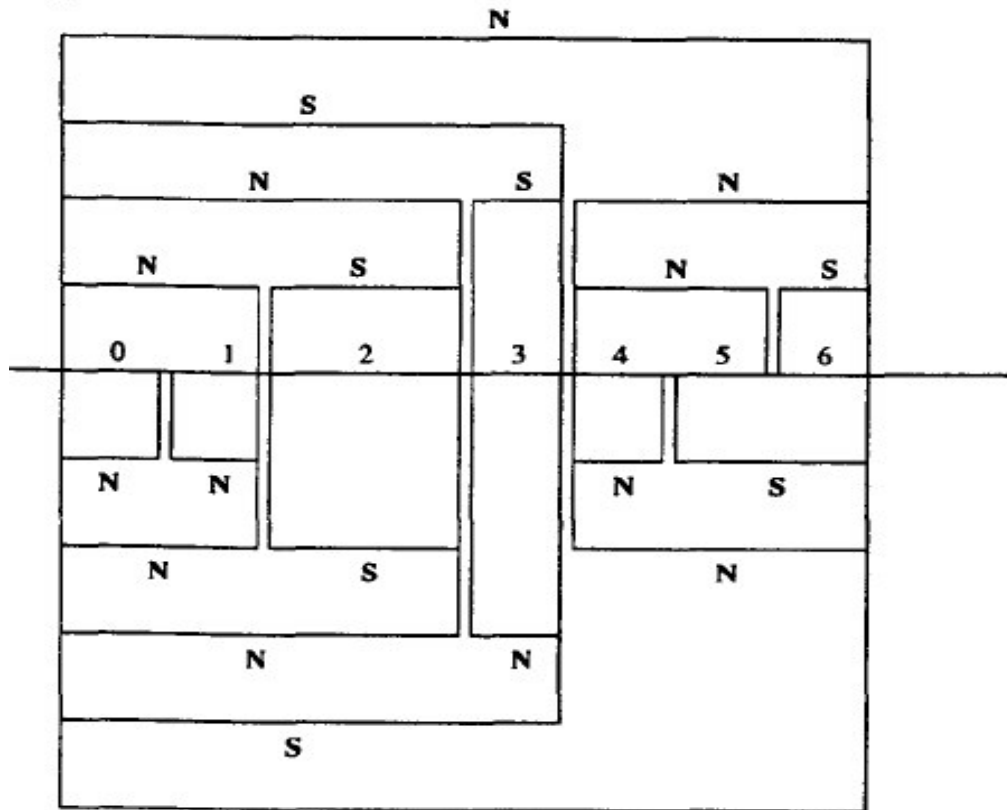
coder 1: [0,1] [2,2] [3,3] [4,5] [6,6]

coder 2: [0,0] [1,1] [2,2] [3,3] [4,4] [5,6]



- mapping for active and non-active EDUs → NONE label

Segmentation 1



Segment	Segmentation 1	Segmentation 2
[0,0]	none	N
[0,1]	N	N
[0,2]	N	N
[0,3]	S	S
[0,4]	none	none
[0,5]	none	none
[0,6]	N	N
[1,1]	none	N
[1,2]	none	none
[1,3]	none	none
[1,4]	none	none
...		
[4,4]	none	N
[4,5]	N	none
[4,6]	N	N
[5,5]	none	none
[5,6]	none	S
[6,6]	S	none

Segmentation 2

- K_n for nuclearity
- K_r for coherence relations
- K_{rr} for a reduced set of coherence relations

Problems:

- violation of independence assumption bw categorical judgments
- NONE agreements make K artificially high
- agreements of different importance
- not for diagnosing disagreements

Q: how to calculate agreement on hierarchical annotation?

- corpus: 150 texts (25 texts per genre)
to annotate for lexical cohesion, RST
- two coders
- expert coders
- detailed coding manual
- coding manual refined as corpus grows
- 20% of the corpus (5 texts per genre)
- training phase: orientation (principles, tools); independent codings (compare → revise manual); final phase (reduce differences)
- presegmented texts for RST analysis and cohesion analysis
- agreement measured per genre as corpus grows
- be explicit (process, scores)

- two coders: GR + IB
- Cohen’s kappa: calculate (1) for EDU boundary, (2) for word boundary (following Marcu et al. 1999)

YES / NO categories for each possible location

example: segmentation of nine encyclopedia entries (EEs)

	EDU level			Word level		
	A _o	A _e	K	A _o	A _e	K
EE17	0.95	0.82	0.72	0.99	0.82	0.94
EE18	0.86	0.74	0.46	0.98	0.83	0.88
EE19	0.84	0.76	0.33	0.98	0.80	0.90
EE20	0.92	0.79	0.62	0.99	0.84	0.94
EE21	0.93	0.92	0.14	0.99	0.84	0.94
EE22	0.89	0.88	0.08	0.99	0.84	0.94
EE23	0.96	0.88	0.67	0.99	0.83	0.93
EE24	0.86	0.78	0.36	0.99	0.91	0.89
EE25	0.91	0.73	0.67	0.99	0.82	0.94

- human errors
- new rules added to coding manual

OR

- percentage agreement (following Carletta et al. 1997)

- two coders: IB + MR
- training the co-coder
- six different move structures for six different genres → time-consuming
- Cohen's kappa (following de Groot 2008)
? two-step procedure:
 - 1: identifying move boundaries – percentage agreement OR K
 - 2: labeling moves identified by both coders – K

- two coders: IB + MR
- training the co-coder
- texts presegmented, potential lexical items preselected
- MMAX2 tool, CORNETTO
- ? Cohen's kappa

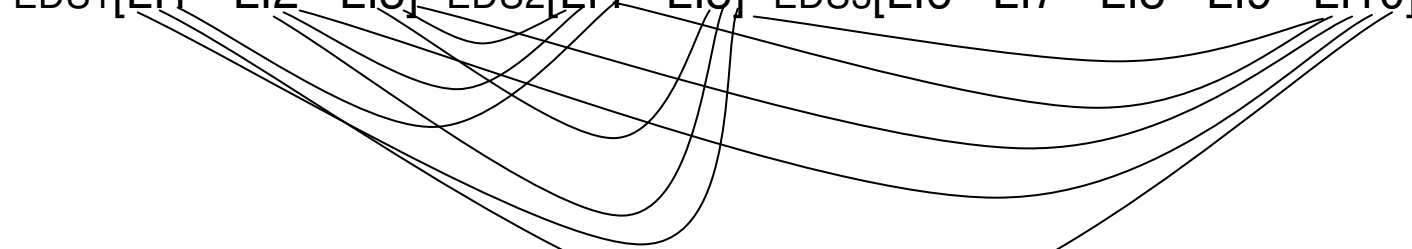
mutually exclusive categories

? equally distinct categories (? weighted coefficient)

lexical cohesive relations across EDU boundaries

multiple relations allowed (= multi-tagged lexical items)

EDU1[LI1 LI2 LI3] EDU2[LI4 LI5] EDU3[LI6 LI7 LI8 LI9 LI10]



LI = lexical item

Q: how to calculate agreement for graph structures?

Own research
Annotation problems
Measures

Practice so far
Decisions – coherence

- two coders: IB + GR
- texts from early training phase not included in the corpus
- ? following Marcu et al. (1999)

Reading

Artstein, R. & Poesio, M. (2008). Inter-coder agreement for computational linguistics. *Computational Linguistics* 34(4), 555-596.

References (intercoder agreement)

Carletta, J. et al. (1997). The reliability of a dialogue structure coding scheme. *Computational Linguistics* 23(1): 13-31.

Carlson, L., Marcu, D. & Okurowski, M.E. (2003). Building a discourse-tagged corpus in the framework of rhetorical structure theory. In J. Kuppevelt & R.W. Smith (Eds.). *Current and new directions in discourse and dialogue*. (pp. 85-112). Dordrecht: Kluwer.

Craggs, Richard & Wood, Mary McGee (2005). Evaluating discourse and dialogue coding schemes. *Computational Linguistics* 31(3): 289-295.

de Groot, E. (2008). *English annual reports in Europe*. Utrecht: LOT.

Krippendorff, K. (1980). *Content analysis: an introduction to its methodology*. Chapter 12. Beverly Hills, CA: Sage.

Marcu, D, Amorrortu, E. & Romera, M. (1999). Experiments in constructing a corpus of discourse trees. In *Proceedings of the ACL Workshop on Standards and Tools for Discourse Tagging*. (pp. 48-57). College Park, MD.

Rietveld, T. & van Hout, R. (1993). *Statistical techniques for the study of language and language behaviour*. Berlin: de Gruyter.

Spooren, W. & Degand, L. (2009). Coding coherence relations: reliability and validity. Unpublished manuscript.

References (IB's research)

Mann, W.C. & Thompson, S.A. (1988). Rhetorical structure theory: Toward a functional theory of text organization. *Text* 8(3), 243-281.

Swales, John (1990). *Genre analysis. English in academic and research settings*. Cambridge: Cambridge University Press.

Upton, Thomas A. (2002). Understanding direct mail letters as a genre. *International Journal of Corpus Linguistics* 7(1), 65-85.