

Development and annotation of a newspaper corpus as part of a doctoral thesis on text structure and cohesion in news items from the 17th and 18th centuries

Katrin Goldschmidt, Universität Bonn

10 March 2017

Content

1. Historical Newspapers
 - Characteristics of historical newspapers
 - Objectives of doctoral thesis
2. A Historical Newspaper Corpus
 - Corpus Development
 - ANNIS
3. Segmentation
 - Typographical Segmentation
 - Functional Segmentation

Historical Newspapers in the 17th and 18th century

Textual structure

- collections of correspondences published by an editor
- single correspondences inform about different news stories
- 3 levels: issue > correspondence > news item(s)

But news stories were written as continuous text...

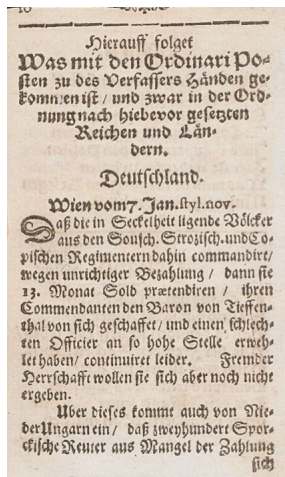
Historical Newspapers in the 17th and 18th century

Textual structure

- collections of correspondences published by an editor
- single correspondences inform about different news stories
- 3 levels: issue > correspondence > news item(s)

But news stories were written as continuous text...

Historical news items in the 17th and 18th centuries



Nordischer Mercurius 1664

Amsterdam den 9. Winterm.
Den Londoner Nachrichten zufolge, denkt
der Großbritannienische Hof auf eine ansehn-
liche Vermehrung der Seemacht. Was
die Gelegenheit hierzu sey, kann aus folgen-
den geschlossen werden; Zwischen den Spa-
niern und Portugiesen soll sich neulich we-
gen der Provinz Nova Colonia, ein ziem-
lich ernsthafter Austritt erdugnet haben.
Hier ist zu wissen, daß die Spanier im letz-
ten Kriege diese Provinz besetzt haben, und
nun die Portugiesen, nach dem Inhalt des
Friedens, da Spanien alles, was es Por-
tugal abgenommen hatte, wieder erstatten
sollte, auch diese Provinz wieder zurück be-
gehrien. Dessen ungeachtet hat Spanien
selbige behalten, unter dem Vorwand, daß
diese Provinz ihm zugehöre, und ihm in
dem Frieden zu Utrecht wäre zuerkannt,
aber von den Portugiesen niemals ab-
getreten worden. Aus dieser Ursache
behielten sie dieselbe nach geschlossenem Frie-
den. Die Portugiesen aber behaupten ihr
Recht mit allem Ernst.

Aus einem Schreiben von Warschau
vom 28. Weinmonat.

Es war am 23. dieses Nachmittags um
2. Uhr, als die russischen Truppen, in Bey-
seyn Sr. Maj. unser Königs, und vieler an-
dern Großen, zwischen Ujazdow und Wola

Wienerisches Diarium 1767

Objectives

Textual Structure

- typographical segmentation of text items
- super- and macrostructural elements as title, headlines, comments, citations

News Structure

- publicistic features as persons, locations, time references
- event-related vs. source-related entities

Coherence Structure

- entities as phrases vs. parts of phrases
- direct and indirect relations within and between news items

Objectives

Textual Structure

- typographical segmentation of text items
- super- and macrostructural elements as title, headlines, comments, citations

News Structure

- publicistic features as persons, locations, time references
- event-related vs. source-related entities

Coherence Structure

- entities as phrases vs. parts of phrases
- direct and indirect relations within and between news items

Objectives

Textual Structure

- typographical segmentation of text items
- super- and macrostructural elements as title, headlines, comments, citations

News Structure

- publicistic features as persons, locations, time references
- event-related vs. source-related entities

Coherence Structure

- entities as phrases vs. parts of phrases
- direct and indirect relations within and between news items

Content

1. Historical Newspapers
 - Characteristics of historical newspapers
 - Objectives of doctoral thesis
2. A Historical Newspaper Corpus
 - Corpus Development
 - ANNIS
3. Segmentation
 - Typographical Segmentation
 - Functional Segmentation

From transcription to annotation

Texts

- 7 german newspaper issues (1609-1767)
- from Berlin, Hamburg (2x), Salzburg, Straßburg, Vienna (2x)
- segmentation of correspondences into news items by 3 raters
=> main corpus
- additional corpus: all January issues of Nordischer Mercurius 1667 with syntactical annotation (Mercurius corpus)

Corpus size

- token: ca. 30.000 [main corpus: 17.784]
- sentences: ca. 1.700 [main corpus: 670]
- news items: ca. 460 [main corpus: 245]
- entities: ca. 5700 [main corpus: 3.384]

From transcription to annotation

Texts

- 7 german newspaper issues (1609-1767)
- from Berlin, Hamburg (2x), Salzburg, Straßburg, Vienna (2x)
- segmentation of correspondences into news items by 3 raters
=> main corpus
- additional corpus: all January issues of Nordischer Mercurius 1667 with syntactical annotation (Mercurius corpus)

Corpus size

- token: ca. 30.000 [main corpus: 17.784]
- sentences: ca. 1.700 [main corpus: 670]
- news items: ca. 460 [main corpus: 245]
- entities: ca. 5700 [main corpus: 3.384]

From transcription to annotation

Texts

- 7 german newspaper issues (1609-1767)
- from Berlin, Hamburg (2x), Salzburg, Straßburg, Vienna (2x)
- segmentation of correspondences into news items by 3 raters
=> main corpus
- additional corpus: all January issues of Nordischer Mercurius 1667 with syntactical annotation (Mercurius corpus)

Corpus size

- token: ca. 30.000 [main corpus: 17.784]
- sentences: ca. 1.700 [main corpus: 670]
- news items: ca. 460 [main corpus: 245]
- entities: ca. 5700 [main corpus: 3.384]

Facsimile and Transcription



Titel: Berlinische Nachrichten von Staats- und gelehrten Sachen
% Erscheinungsort: Berlin
% Herausgeber: Haude & Spener
% Ausgabe: Di, 34
% Gedruckt: 21-03-1741
% Transkription: IDS-Zeitungskorpus, zur Verf. gestellt durch Ins
% Satzsegmentierung und Preprocessing: Katrin Goldschmidt, Univer
% letzteKorrektur: 16-05-2014
% nicht aufgenommene Originalseiten bzw Textstellen: +K Rubriken

Seitel

Ao. 1741. +K Titelkupfer "WAHRHEIT.VND.FREIHEIT." @K No. XXXIV .
Dienstag +K Titelkupfer "WAHRHEIT.VND.FREIHEIT." @K den 21. Merz
+K danach durchgehende Linie @K
&spbr& Berlinische Nachrichten &spbr&
&spbr& von &spbr&
&spbr& Staats- und gelehrten Sachen . &spbr&
+K danach durchgehende Linie @K
+K Beginn linke Spalte @K
&spbr& Berlin , vom 21. Merz . &spbr&
&spbr& +K Schmuckinitiale Adler hält Majuskel A @K *(A*6)m verwicl
Prinz Wilhelmische Regiment von
hier nach Schlesien auf . &spbr& Seine
Königl. Hoheit führten dasselbe in
hoher Person auf und ein jeder
bewunderte die Vollkommenheit
dieses Regiments . &spbr&
&spbr& In wenig Tagen werden Se. Königl. Hoheit
selbst nachfolgen . &spbr&
&spbr& London , vom 9. Merz . &spbr&
&spbr& Am 4. dieses Monats erhielt der Hof Briefe
von dem Ritter Ogle , welche am 4. Jenner geschrie-
ben waren , mit der Nachricht , daß sich der Ritter am
30. December des vorigen Jahres mit seiner Flotte
zu St. Domingo vor Anker gelegt hätte . &spbr& Die
Kriegsschiffe Rippon und York wären zu ihm ge-
+K Beginn rechte Spalte @K
stossen , allein am 12. November 66. Meilen vom
Cap Lezard durch einen Sturm wieder von ihm ge-
trennt worden . Er füget noch hinzu , daß die Schiffe
auf dem Fall einer Trennung , sich auf der Insel zu
St. Christoph wieder zu finden abgeredet hätten ,
weswegen er sie auch daselbst anzutreffen hoffte .
Diesen Briefen zu Folge hat der Ritter am 5. die-

Berlinische Nachrichten 1741

Annotation

dieses Regiments.

In wenig Tagen werden St. Königl. Hoheit selbst nachfolgen.

Londen, vom 9. März.

Am 4. dieses Monats erhielt der Hof Briefe von dem Ritter Ogle, welche am 4. Jenner geschrieben waren, mit der Nachricht, daß sich der Ritter am 30. December des vorigen Jahres mit seiner Flotte zu St. Domingo vor Anker gelegt hätte. Die Kriegsschiffe Rippon und York waren zu ihm ge-

Der Lord Cathcart, commandirender General unserer Truppen in Westindien, ist am 19. Decembr., von einem heftigen Durchfalle angegriffen worden, und am 31. als am andern Tage nach der Ankunft des Ritters Ogle zu St. Domingo, verstorben. So schmerzhaft als seine Krankheit gewesen, so zeigte er sich als ein wahrhafter Engländer und starb eben so großmüthig als er gelebt hatte. Alles was man von Misvergnügen bey ihm fand, rührte nur daher, daß er ohne seinen Landesleuten die

tok	tok_ID	pos	text	entity_p	coref_target_p	coref_entity_p	indref_target	indref_entity	indref_reftype	typosegm	is
seiner	107	PPOSAT					sour_inf		80	poss	
Flotte	108	NN									
zu	109	APPR		loc_H							
St.	110	NE									
Domingo	111	NE									
vor	112	APPR									
Anker	113	NN	bn_b								
(.)	196	\$.	bn_e								
Der	197	ART	bn_b	who_H							
Lord	198	NN									
Cathcart	199	NE									
ist	208	VAFIN									
zu	233	APPR		loc_H	loc_H_bn	2	109				
St.	234	NE									
Domingo	235	NE									
(.)	236	\$.									

Annotation levels

Textual Structure

- **token**
- pos
- **text**
- makrostr
- **typosegm**
- issue
- correspondence
- origin
- meta

News Structure

- **entity_p (phrase)**
 - entity_pp (part of phrase)
 - how_why
 - episode
-
- indref_reftype (func, part, poss)
 - time_reftype (before, during, after)

Coherence Relations

- **coref_target_p**
- coref_target_pp
- indref_target
- lexref_target
- ref_episode_target
- ref_howwhy_target

entity reference view

11/21

HistZeitKorp in local ANNIS

discourse view

Displaying Results 1 - 10 of 18 Result for: entity_p=/loc.* / & node & node #2->coref_entit ...

dieses Monats erhielt der Hof Briefe von dem Ritter Ogle .(,) welche am 4. Jenner geschrieben waren .(,) mit der Nachricht .(,) daß sich der Ritter am 30. December des vorigen Jahres mit seiner Flotte zu St. Domingo vor Anker gelegt hätte .(,) Die Kriegsschiffe Rippon und York wären zu ihm gestossen .(,) allein am 12. November 66. Meilen vom Cap Lezard durch einen Sturm wieder von ihm getrennet worden .(,) Er füget noch hinzu .(,) daß die Schiffe auf dem Fall einer Trennung .(,) sich auf der Insel zu St. Christoph wieder zu finden abgeredet hätten .(,) weswegen er sie auch daselbst anzutreffen hoffete .(,) Diesen Briefen zu Folge hat der Ritter am 5. dieses Monats seinen Weg nach Jamaika fortgesetzt .(,) Der Lord Catheard .(,) commandirender General unserer Truppen in Westindien .(,) ist am 19. Decembr. .(,) von einem heftigen Durchfalle angegriffen worden .(,) und am 31. als am andern Tage nach der Ankunft des Ritters Ogle zu St. Domingo .(,) verstorben .(,) So schmerzhaft als seine Krankheit gewesen .(,) so bezeugte er sich als ein wahrhafter Engländer und starb eben so großmüthig .(,) in Misvergnügen bey ihm fand .(,) rührte nur daher .(,) daß er ohne seinen Landesleuten die Wirkungen seines E .(,) nach seinem Tode hat der General Wentworth seine Stelle übernommen .(,) Paris , vom 8. Merz .(,) Der Marschall von .(,) dieses Monats nach Metz abgereiset .(,) wo sie sich einige Zeit aufhalten werden .(,) bevor sie

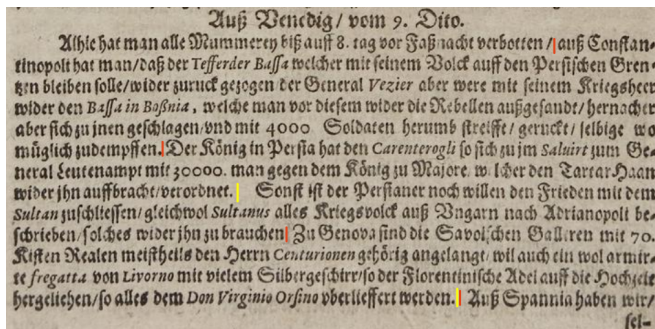
Component: 38, **Type:** coref_entity_p, **outgoing**
Annotations: coref_target p=loc_H_bn

Content

1. Historical Newspapers
 - Characteristics of historical newspapers
 - Objectives of doctoral thesis
2. A Historical Newspaper Corpus
 - Corpus Development
 - ANNIS
3. Segmentation
 - Typographical Segmentation
 - Functional Segmentation

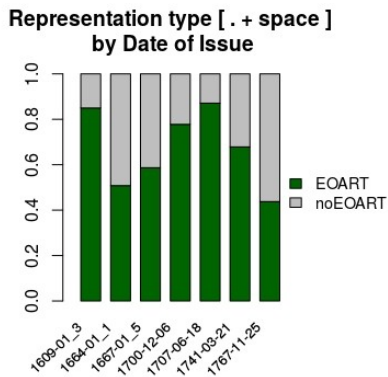
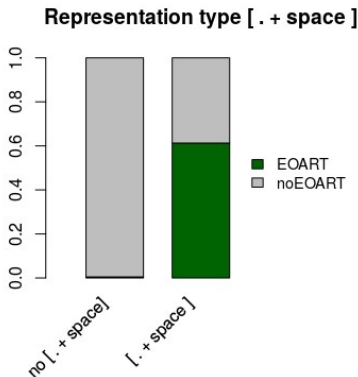
Segmentation and the role of typographical means

- following Lefèvre (2013) shifts between paragraphes are marked by the representation type [. +space(+UpperCase)]
- but: no quantitative analysis



Relation 1609 [boundaries: Lefèvre (2013); HistNewsKorp]

Typographical segmentation is not that easy...



Segmentation and the role of functional means

- following Schröder (1995) shifts between news items are marked by change of theme, person, location and time [Relation 1609, Aviso 1609]
- over 75% of the news items are introduced by entities that respond to journalistic wh-questions (like *who, where, when ...?*)

A first approach: "How many news items begin with entities?"

- ANNIS search for entities as phrases at the direct beginning of a news item: `text=/bn_b/ & entity_p & #2 _#1`
- 167 matches

Segmentation and the role of functional means

- following Schröder (1995) shifts between news items are marked by change of theme, person, location and time [Relation 1609, Aviso 1609]
- over 75% of the news items are introduced by entities that respond to journalistic wh-questions (like *who, where, when ...?*)

A first approach: "How many news items begin with entities?"

- ANNIS search for entities as phrases at the direct beginning of a news item: `text=/bn_b/ & entity_p & #2 _#1`
- 167 matches

Help/Examples

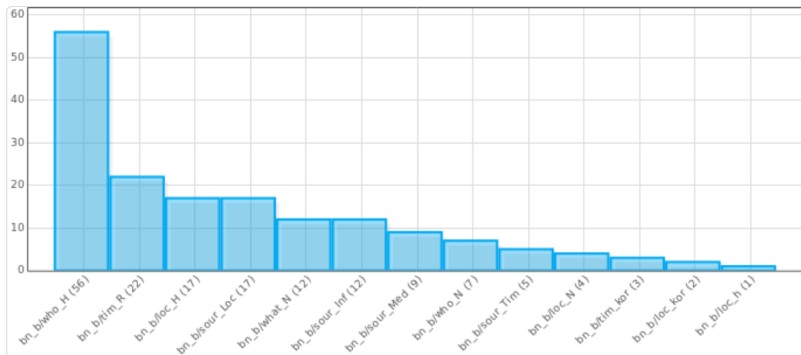
Query Result ×

Frequency Analysis ×

New Analysis

☒ linear scale

☐ log₁₀ scale



Also functional segmentation is not that easy...

- we see that 64% of the news items begin with an entity that responds to the forecited journalistic wh-questions:

	Schröder1995: Relation 1609	HistZeitKorp	
<u>PublInfo</u>	<u>percentage</u>	<u>percentage</u>	<u>N</u>
<u>per_loc_tim_them</u>	75%	64%	158
<u>other</u>	8%	4%	9
<u>none</u>	17%	32%	78
total	100%	100%	245

A second approach: "How many news items begin with entities that refer to entities within former news items?"

- ANNIS search: 11 matches ($11 \div 158 \approx 7.0\%$)

Also functional segmentation is not that easy...

- we see that 64% of the news items begin with an entity that responds to the forecited journalistic wh-questions:

	Schröder1995: Relation 1609	<u>HistZeitKorp</u>	
<u>PublInfs</u>	<u>percentage</u>	<u>percentage</u>	<u>N</u>
<u>per_loc_tim_them</u>	75%	64%	158
<u>other</u>	8%	4%	9
<u>none</u>	17%	32%	78
total	100%	100%	245

A second approach: "How many news items begin with entities that refer to entities within former news items?"

- ANNIS search: 11 matches ($11 \div 158 \approx 7.0\%$)

Conclusions

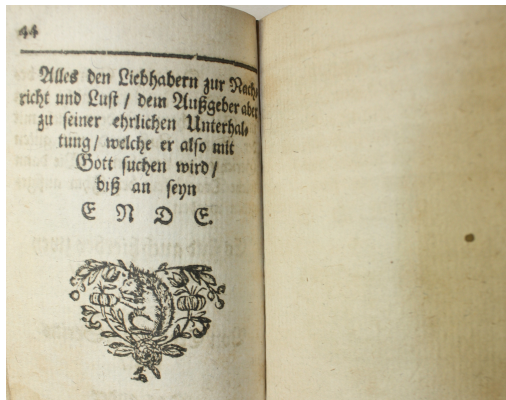
- typographical means alone are not an adequate means of news item segmentation (maybe in 60% of the correspondences they are helpful)
- at least 7% of phrasal entities at the beginning of news items refer to entities in the preceding news item(s)
=> implies no thematic change
- **but:** considering typographical AND functional means could lead to successive segmentation methods

With help of the Historical Newspaper Corpus we can get a better understanding of the textual and publicistic structure of these interesting old newspapers.

Conclusions

- typographical means alone are not an adequate means of news item segmentation (maybe in 60% of the correspondences they are helpful)
- at least 7% of phrasal entities at the beginning of news items refer to entities in the preceding news item(s)
=> implies no thematic change
- **but:** considering typographical AND functional means could lead to successive segmentation methods

With help of the Historical Newspaper Corpus we can get a better understanding of the textual and publicistic structure of these interesting old newspapers.



Nordischer Mercurius 1664

„Everything [serves] to the [news] lovers [as] a message and a pleasure / but to the editor [it serves as] his honest entertainment, which he therefore will seek with God / until his END.“

References

- Fritz, G. & E. Straßner (eds.) (1996): Die Sprache der ersten deutschen Wochenzeitungen im 17. Jahrhundert. Tübingen.
- Demske, U. (2007): Das Mercurius-Projekt: eine Baubank für das Frühneuhochdeutsche. In: G. Zifonun & W. Kallmeyer (eds.): Sprachkorpora - Datenmengen und Erkenntnisfortschritt (= Jahrbuch des Instituts für Deutsche Sprache 2006). Berlin, 91-104.
- Krause, T. & A. Zeldes (2014): ANNIS3: A New Architecture for Generic Corpus Query and Visualization. Literary and Linguistic Computing.
- Lefèvre, M. (2013): Textgestaltung, Äußerungsstruktur und Syntax in deutschen Zeitungen des 17. Jahrhunderts. Zwischen barocker Polyphonie und solistischem Journalismus. Berlin.
- Schröder, T. (1995): Die ersten Zeitungen. Textgestaltung und Nachrichtenauswahl. Tübingen.

Special thanks to all who supported me with the newspaper records (Österreichische Nationalbibliothek, Staatsbibliothek zu Berlin, Staats- und Universitätsbibliothek Bremen, Institut für Deutsche Sprache Mannheim) and Prof. Ulrike Demske for providing the Mercurius Corpus.

Additional corpus :: NM_1667-01_synt

- all January issues of Nordischer Mercurius 1667
- integrated syntactical annotation from the Mercurius corpus

The screenshot displays the exmaralda search interface. On the left, a query builder shows the query: `exmaralda:entity_p="who_N" & #1 _= tiger:cat="PP"`. Below the query builder, search options and a corpus list are visible. The corpus list shows 102 matches in 1 document, with the selected corpus being NM_1667-01_synt (1 text, 12,389 tokens).

The main search results area shows the following information:

- Path: NM_1667-01_synt > NM_1667-01_synt (tokens 6108 - 6121)
- Result for: exmaralda:entity_p="who_N" &
- Left context: 4

The search results table displays the following text and POS tags:

Der	Venetianische	Ambassadeur	hat	bey	ihrer	Königl.	M	eine	extraord.	Audientz	gehabt
ART	ADJA	NN	VAFIN	APPR	PPOSAT	ADJA	NN	ART	ADJA	NN	VVPP

Below the search results, a table shows the correspondence between the search results and the corpus list:

issue	NM_1667-01_5											
correspondence	Paris_1667-01-14											
text	bn_b											
episode	e											
entity_p	who_H	who_N	what_N									
pos	ART	ADJA	NN	VAFIN	APPR	PPOSAT	ADJA	NN	ARTU	ADJA	NN	VVPP
pos_mercproj	ART	ADJA	NN	VAFIN	APPR	PPOSAT	ADJA	NN	ART	ADJA	NN	VVPP
tok_mercproj	Der	Venetianische	Ambassadeur	hat	bey	ihrer	Königl.	M	eine	extraord.	Audientz	gehabt
tok	Der	Venetianische	Ambassadeur	hat	bey	ihrer	Königl.	M	eine	extraord.	Audientz	gehabt

Below the table, a tree diagram (tiger) is shown, illustrating the syntactic structure of the sentence. The tree is rooted at S, which branches into SB (subject) and VP (verb phrase). SB branches into NP (noun phrase) and H (head). NP branches into INK (indefinite article) and N (noun). VP branches into V (verb) and NP (noun phrase). The tree structure is as follows:

```
graph TD
    S --> SB
    S --> VP
    SB --> NP1[NP]
    SB --> H1[H]
    NP1 --> INK1[INK]
    NP1 --> N1[N]
    VP --> V[V]
    VP --> NP2[NP]
    NP2 --> INK2[INK]
    NP2 --> N2[N]
```

The tree diagram shows the syntactic structure of the sentence. The root node S branches into SB (subject) and VP (verb phrase). SB branches into NP (noun phrase) and H (head). NP branches into INK (indefinite article) and N (noun). VP branches into V (verb) and NP (noun phrase). The tree structure is as follows: