

COST Action Distant Reading for European Literary History

European Literary Text Collection (ELTeC)

Distant  *Reading*

Christof Schöch, Carolin Odebrecht, Lou Burnard, Borja
Navarro-Colorado et al.

TS Budapest Track 1, Corpus design and text contribution for ELTeC

TS – Organisation, schedule and data

- ▶ Schedule on <https://github.com/distantreading/WG1/wiki/TS-Budapest>
- ▶ Slides on <https://github.com/distantreading/WG1/tree/master/Training/2019-09-budapest>
- ▶ Data on <https://github.com/COST-ELTeC>

TS – Track 1

- ▶ Introduction to ELTeC
- ▶ Metadata in the `teiHeader`
- ▶ Optical character recognition
- ▶ Hands-on experience in creating ELTeC TEI-XML versions of source texts
- ▶ Start from scanned page images or from a pre-existing HTML version
- ▶ Our goal: contribute new TEI encoded texts to the ELTeC GitHub repository

Outline

1. COST Action DISTANT READING and Working Group
SCHOLARLY RESOURCES
2. ELTeC Corpus Design
3. WG Workflow

COST² Action Distant Reading

- ▶ Christof Schöch, University of Trier
- ▶ CA16204 will
 - ▶ „create a vibrant and diverse network of researchers jointly developing the resources and methods necessary to change the way European literary history is written“
 - ▶ „contribute to the development and distribution of methods, competencies, data, best practices, standards and tools relevant to Distant Reading research“¹
 - ▶ started autumn 2017
- ▶ Working groups
 - ▶ WG 1: Scholarly Resources
 - ▶ WG 2: Methods and Tools
 - ▶ WG 3: Literary Theory and History
 - ▶ WG 4: Dissemination

¹www.distant-reading.net

²European Cooperation in Science and Technology = COST³

Working Group 1: SCHOLARLY RESOURCES

- ▶ Creating an open source multi-lingual benchmark corpus for European literature of the 19th century (novels): European Literary Text Collection (ELTeC)⁴
- ▶ 34 Members from 22 countries
- ▶ Main tasks are
 - ▶ defining corpus design
 - ▶ developing basic encoding schemas
 - ▶ developing workflows

⁴<https://www.distant-reading.net/wg-1/>

Working Group 1: SCHOLARLY RESOURCES

Context of the WG:

- ▶ closely related to WG2 TOOLS AND METHODS and WG3 THEORY, cf. parallel tracks of the TS
- ▶ aims to support many different perspectives on distant reading
- ▶ many members with different scientific background and experience in data creation
- ▶ enabling distant collaborative work on creating ELTEC

Outline

1. COST Action DISTANT READING and Working Group
SCHOLARLY RESOURCES

2. ELTeC Corpus Design

3. WG Workflow

Corpus design

Corpus design defines two things (cf. Hunston 2008; Lüdeling et al. 2016):

- ▶ candidates → sampling
 - ▶ Which text(s) can be included in the corpus? Which can't?
- ▶ proportion → balancing
 - ▶ How many texts with which characteristics should the corpus contain?

Corpus design – approach of WG1

- ▶ Sampling and balancing criteria⁵ will
 - ▶ not define what a novel is
 - ▶ follow a non-normative but metadata-based approach (not canon-based)⁶
 - ▶ aim to represent the variety of a population⁷
 - ▶ allow for a comparability of texts and individual sub-collections according to different metadata set(s)

⁵ https://github.com/distantreading/WG1/blob/master/sampling_proposal.xml

⁶ Each canon is a result of rating texts from different perspectives: intellectual, economical, or/and reader rating (a.o. Herrmann 2011; Winko 1996).

⁷ Cf. for discussion of representativeness Biber (1993) and canonicity and corpus design Algee-Hewitt and McGurl (2018) and Bode (2018).

Sampling criteria

- ▶ **language:** European languages, no translations
- ▶ **prose:** narrative fictional prose
- ▶ **period:** 1840–1920
- ▶ **length:** min. 10.000 words
- ▶ **publication:** prefer books over novels published in serial publications
- ▶ **access:** only freely available digitizations

Balancing criteria

100 texts per language (language collection)

- ▶ **period:** distribution over time
 - ▶ group T1: 1840-1859
 - ▶ group T2: 1860-1879
 - ▶ group T3: 1880-1899
 - ▶ group T4: 1900-1920
- ▶ **gender:** min. 10% and max. 50% written by female authors
- ▶ **authorship:** 9 - 11 authors with exactly three novels (otherwise, only one text for each author)
- ▶ **length:** min. 20% are short novels (10-50k word tokens), min. 20% are long novels (>100k word tokens).
- ▶ **reprint:** min. 30% are highly canonized novels, min. 30% should be non-canonized novels (rarely or never reprinted), based reprint counts within the period 1970-2009 ("unmarked", information not yet available, "unspecified" information about the number of reprints not yet determined)

(ideal) Composition

Current composition of the language collection

		length.cat			gender.cat
		long	medium	short	
year.cat	T1	5	4	4	F
		4	4	4	M
T2		5	4	4	F
		4	4	4	M
T3		5	4	4	F
		4	4	4	M
T4		5	4	4	F
		4	4	4	M

Amount of texts with balancing categories (year, length, gender)

ELTeC summary <https://distantreading.github.io/ELTeC/>

Outline

1. COST Action DISTANT READING and Working Group
SCHOLARLY RESOURCES

2. ELTeC Corpus Design

3. WG Workflow

WG workflow

- ▶ WG1 on GitHub⁸
 - ▶ organization and documentation
- ▶ ELTeC on GitHub⁹
 - ▶ data and scripts
- ▶ Archiving on Zenodo¹⁰

⁸<https://github.com/distantreading/WG1>

⁹<https://github.com/COST-ELTeC>

¹⁰<https://zenodo.org/communities/eltec/>

References I



Algee-Hewitt, Mark and Mark McGurl (2018). *Between canon and corpus: six perspectives on 20th-century novels*.
URL: <https://litlab.stanford.edu/LiteraryLabPamphlet8.pdf>.



Biber, Douglas (1993). "Representativeness in Corpus Design". In: *Literary and Linguistic Computing* (8), pp. 243–257.



Bode, Katherine (2018). *A World of Fiction - Digital Collections and the Future of Literary History*. eng.
University of Michigan Press.



Herrmann, Leonhard (2011). "System? Kanon? Epoche?" In: *Kanon, Wertung und Vermittlung. Literatur in der Wissensgesellschaft*. Ed. by Claudia Stockinger Matthias Beilein and Simone Winko. Berlin: De Gruyter, pp. 59–75.



Hunston, Susan (2008). "Collection strategies and design decisions". In: *Corpus Linguistics. An International Handbook*. Ed. by Anke Lüdeling and Merja Kytö. Vol. 1. 2 vols. 1. Berlin: De Gruyter, pp. 154–168.



Lüdeling, Anke, Julia Ritz, Manfred Stede, and Amir Zeldes (2016). "Corpus Linguistics". In: *OUP Handbook of Information Structure*. Ed. by Caroline Fery and Shinishiro Ishihara. Oxford: Oxford University Press, pp. 599–617.



Winko, Simone (1996). "Literarische Wertung und Kanonbildung". In: *Grundzüge der Literaturwissenschaft*. Ed. by Heinz Ludwig Arnold and Heinrich Detering. München: Deutscher Taschenbuch Verlag, pp. 585–600.