

Data and Annotation in Literary Studies

Some Methodological Remarks



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Evelyn Gius

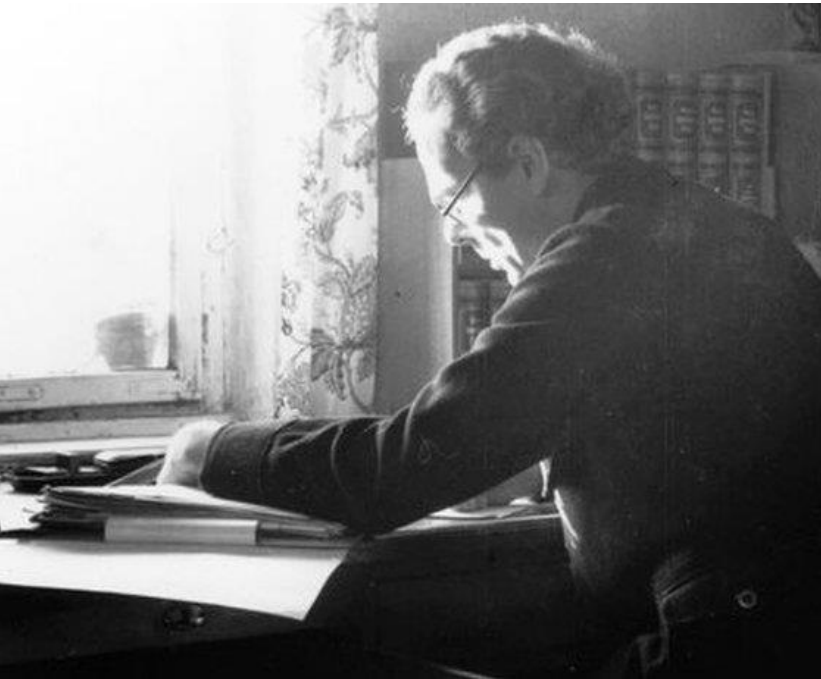
Digital Philology, Technical University of Darmstadt



Literary Studies



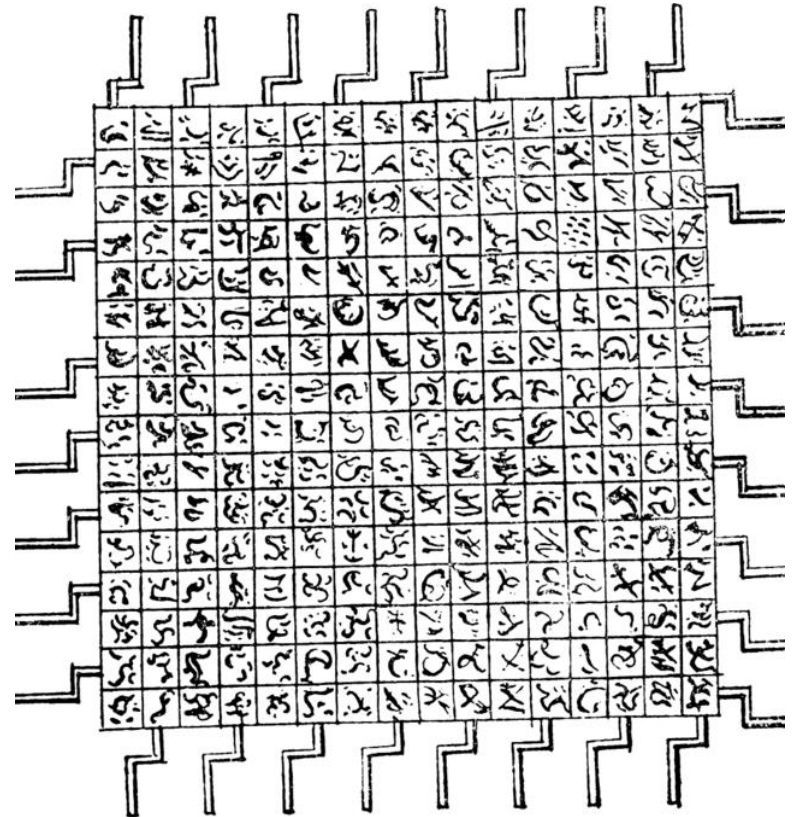
TECHNISCHE
UNIVERSITÄT
DARMSTADT



Computational Literary Studies

.. Every one knew how laborious the usual method is of attaining to arts and sciences; whereas, by his contrivance, the most ignorant person, at a reasonable charge, and with a little bodily labour, might write books in philosophy, poetry, politics, laws, mathematics, and theology, without the least assistance from genius or study.

[...] The pupils, at his command, took each of them hold of an iron handle, whereof there were forty fixed round the edges of the frame; and giving them a sudden turn, the whole disposition of the words was entirely changed. He then commanded six-and-thirty of the lads, to read the several lines softly, as they appeared upon the frame; and where they found three or four words together that might make part of a sentence, they dictated to the four remaining boys, who were scribes.



Swift, Jonathan (1726). [Gulliver's Travels](#). p. Part 3, Chapter 5.

https://upload.wikimedia.org/wikipedia/commons/6/6d/The_Engine_%28Gulliver%29.png

Literary Text Analysis as Data Analysis

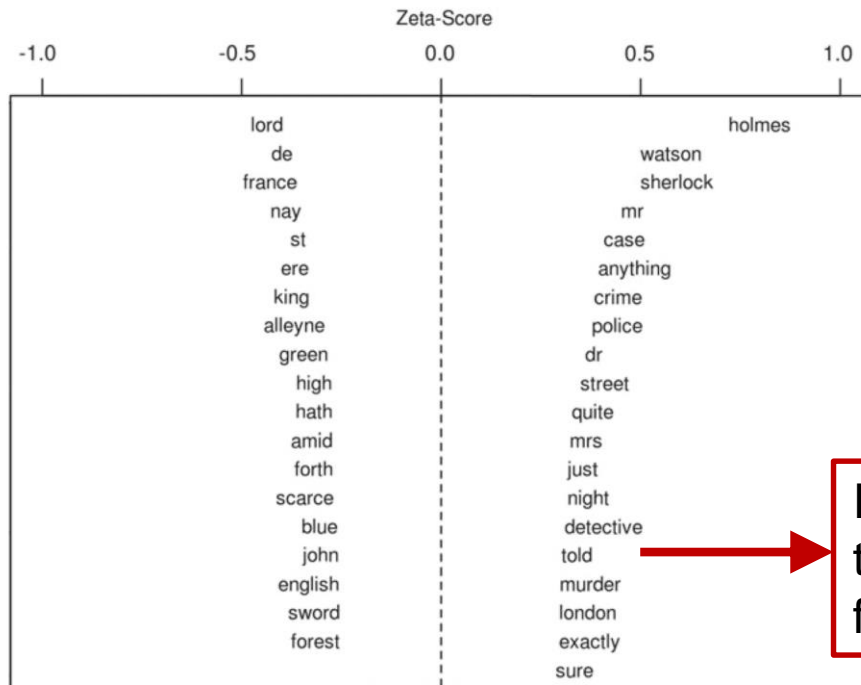


Abb. 59 Zeta-Scores. Überrepräsentierte (rechts, positives Zeta) und unterrepräsentierte (links, negatives Zeta) Types für die Kriminalromane im Vergleich zu den übrigen Romanen in der Doyle-Sammlung

In the computational approach
the very object of analysis shifts
from literary text(s) to data

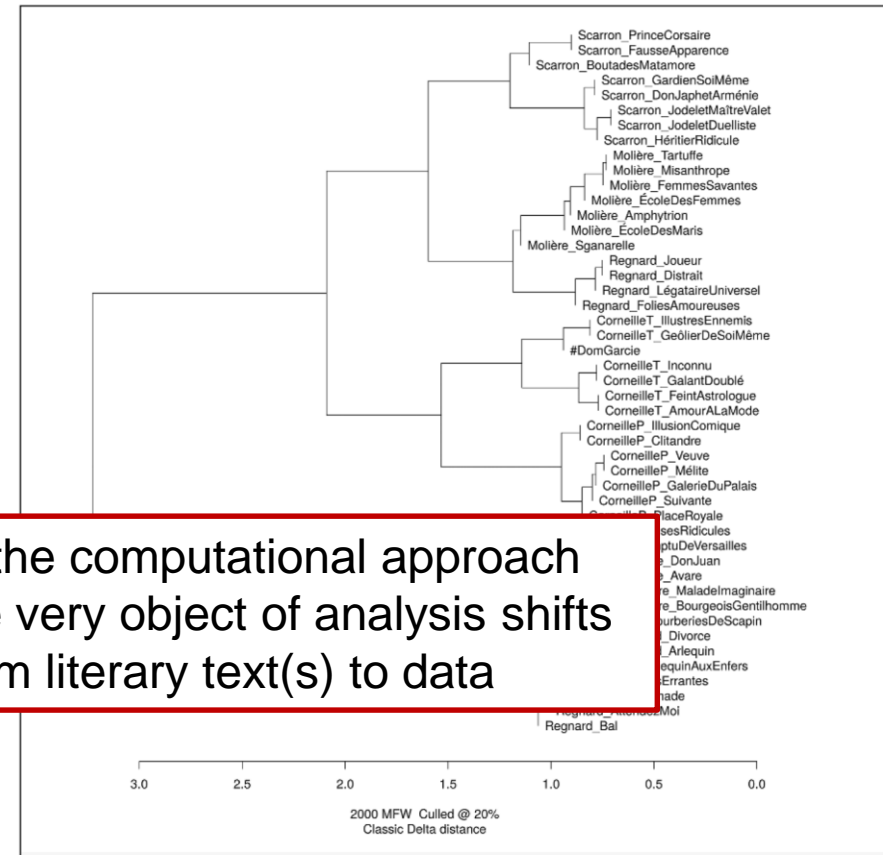
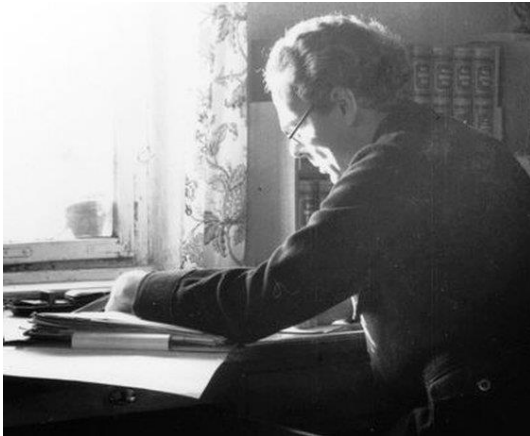


Abb. 61 Stilometrisches Dendrogramm der Ähnlichkeitsbeziehungen zwischen 47 französischen Komödien des 17. Jahrhunderts, darunter ein Stück umstrittener Autorschaft (Label: #DomGarcie)

cf. Schöch, Christof. "Quantitative Analyse." *Digital Humanities: eine Einführung*, edited by Fotis Jannidis et al., J.B. Metzler Verlag, 2017, pp. 279–298.

Literary Text Analysis as Data Analysis

text analysis: text(s) → analysis → findings/theory
≈ data analysis (?): input → analysis → output



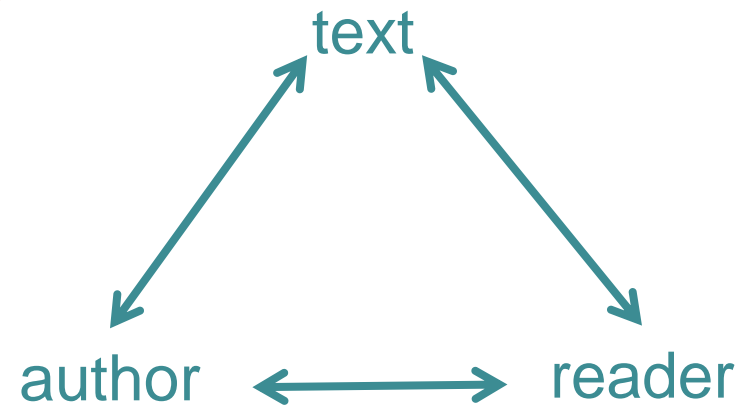
Data in Literary Studies: Input



... is more than analysis of text

Reason 1: Literature as Communication

- Literary Communication includes:
 - author(s)
 - text(s)
 - reader(s)
- + narrator(s), implied author(s), implied reader(s)...

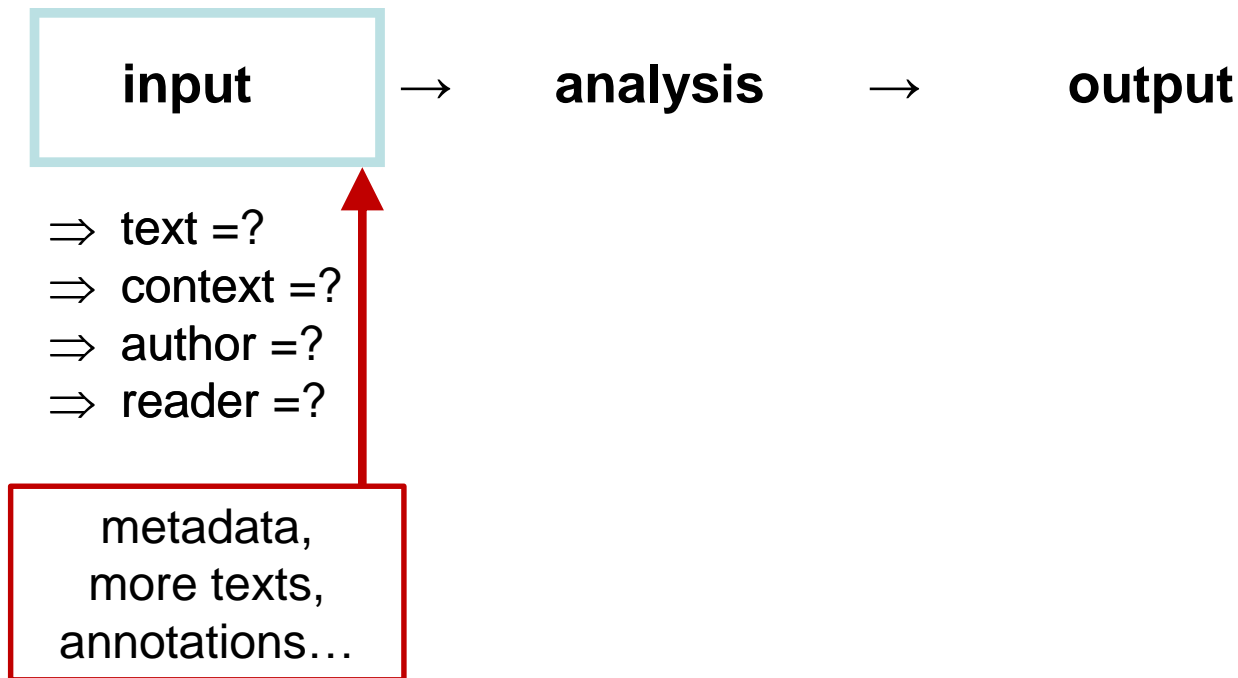


... is more than analysis of text

Reason 2: Theories of Interpretation

in literary studies focus on:

- text (e.g., structuralism, deconstruction),
- context (e.g., discourse analysis, gender studies, system theory),
- author (e.g. hermeneutics, psychoanalytic approaches), and/or
- recipient (e.g. reception aesthetics).



Data in Literary Studies: Analysis



- text analysis in (non-digital) literary studies
 - focusses on different aspects
 - includes intra and extra textual aspects
 - is—necessarily—subjective (cf. esp. extra textual aspects, focus on reader)
- (professional) interpretation: ↑ **intersubjective** understanding of texts

Data in Literary Studies: Input and Analysis

Example: annotation based analysis

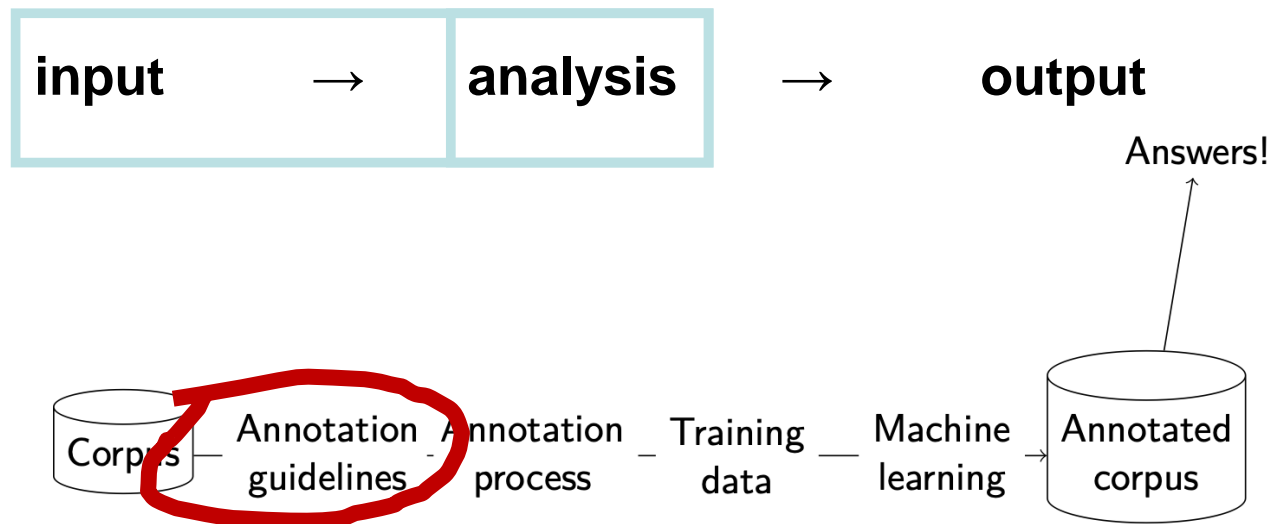


Figure: Towards large scale literature analysis

Figure from presentation „Shared Tasks for the Development of Annotation Guidelines in the Digital Humanities Annotation in Scholarly Editions and Research“, Evelyn Gius, Nils Reiter and Marcus Willand, Workshop *Annotation in scholarly editions and research*, Bergische Universität Wuppertal, February 20-22, 2019

Guideline Creation as Shared Task

cf. Gius, Evelyn, Nils Reiter, and Marcus Willand. 2019. *A Shared Task for the Digital Humanities – Special Issue of Journal of Cultural Analytics*.

Systematic Analysis of Narrative Texts through Annotation

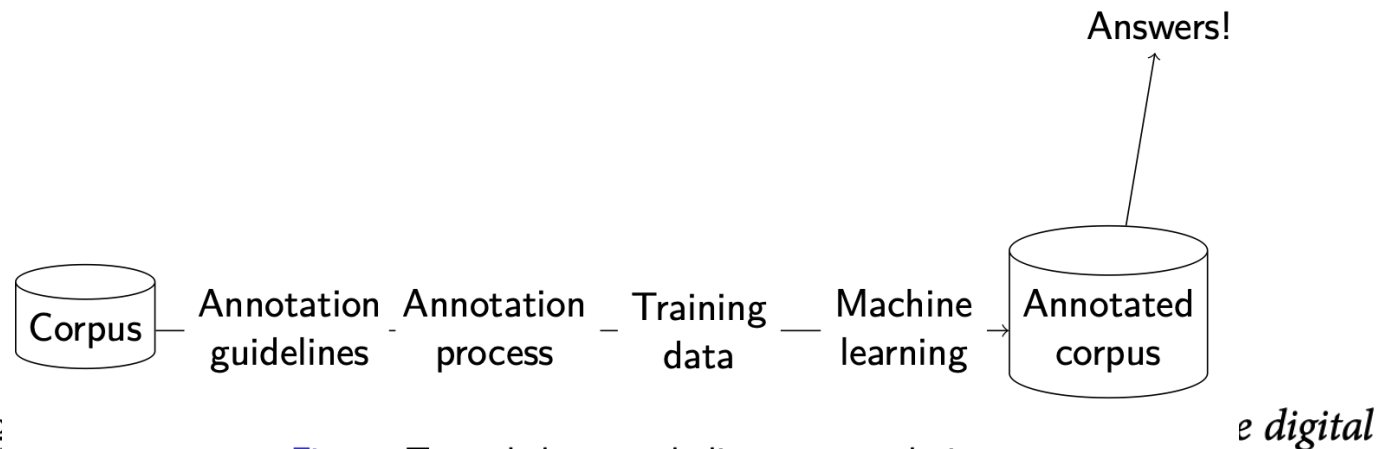


Fig
hu

Figure: Towards large scale literature analysis

Figure from Gius, Evelyn, Nils Reiter, and Marcus Willand. 2019. "A Shared Task for the Digital Humanities Chapter 2: Evaluating Annotation Guidelines." *Journal of Cultural Analytics*. <https://doi.org/10.22148/16.049>.

Data in Literary Studies: Input and Analysis

Example: annotation based analysis

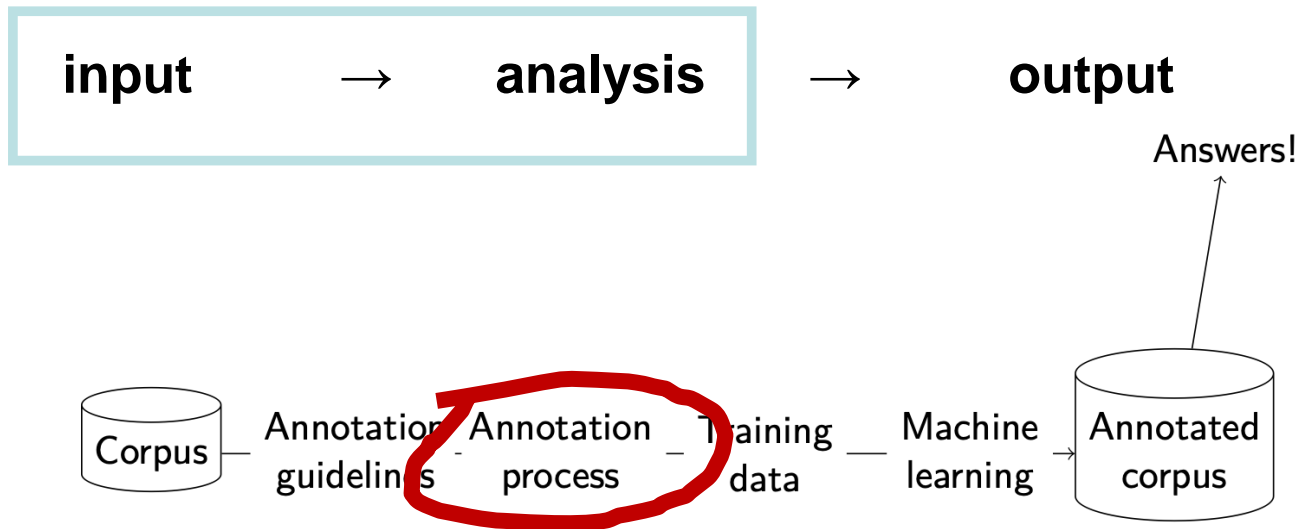
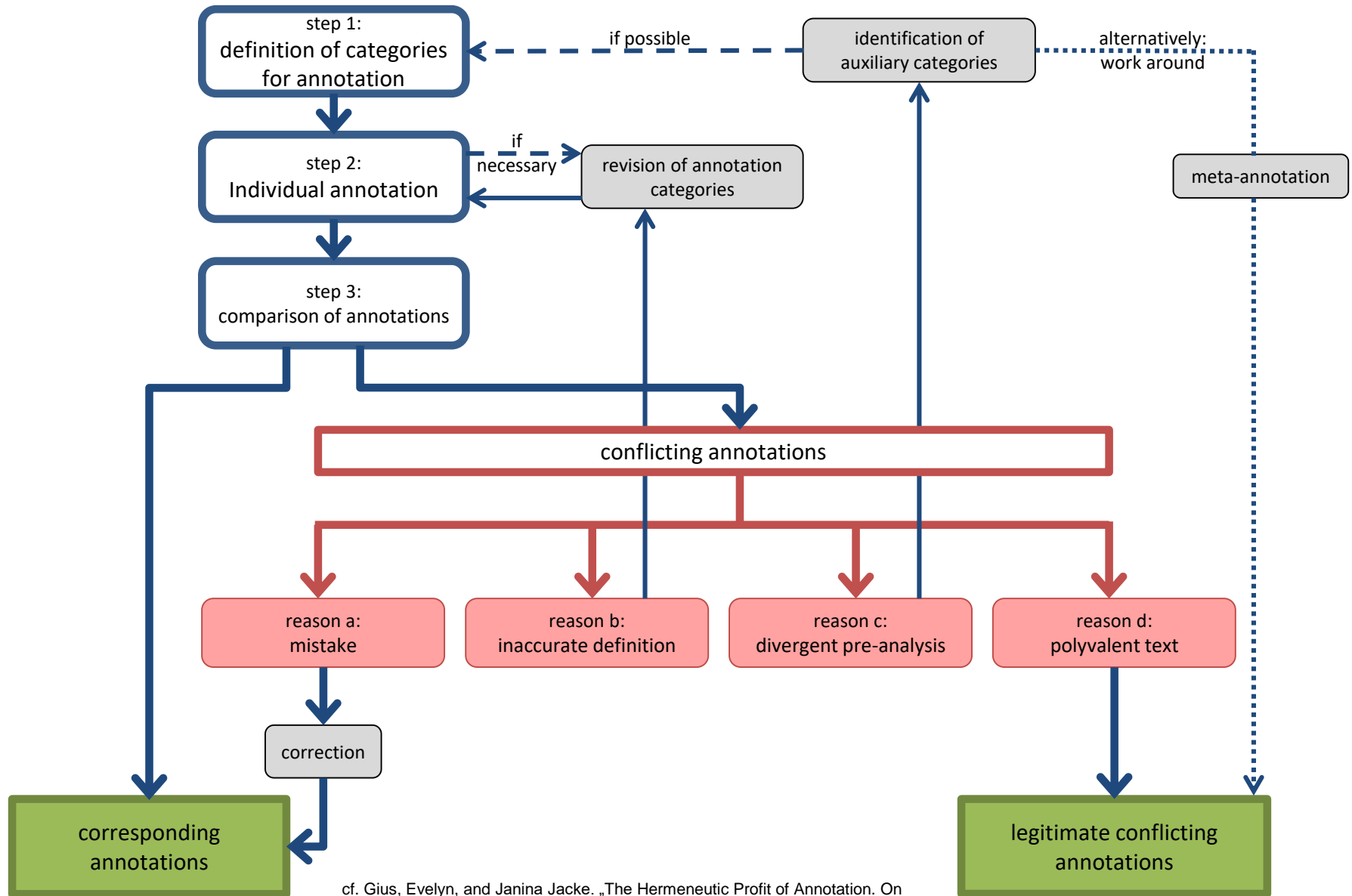


Figure: Towards large scale literature analysis

Figure from presentation „Shared Tasks for the Development of Annotation Guidelines in the Digital Humanities Annotation in Scholarly Editions and Research“, Evelyn Gius, Nils Reiter and Marcus Willand, Workshop *Annotation in scholarly editions and research*, Bergische Universität Wuppertal, February 20-22, 2019

The Annotation Process



cf. Gius, Evelyn, and Janina Jacke. „The Hermeneutic Profit of Annotation. On preventing and fostering disagreement in literary text analysis“. *International Journal of Humanities and Arts Computing* 11, Nr. 2 (2017): 233–54.
<https://doi.org/10.3366/ijhac.2017.0194>.

Data in Literary Studies: Analysis



Analysis between Explanation and Understanding

The traditionally important distinction between explanation and understanding is hardly meaningful in the context of interpretation theory: every interpretation is concerned in one way or another with explaining certain textual findings, and it aims to promote a better understanding of the work (depending on the respective theory of meaning).*

Köppe, Tilmann, Winko, Simone (2013): Theorien und Methoden der Literaturwissenschaft, in: Methoden und Theorien, ed. by. Thomas Anz, Bd. 2, Stuttgart 2013 Handbuch Literaturwissenschaft, 285–371, p.287

* Die traditionell bedeutsame Unterscheidung zwischen Erklären und Verstehen ist im Rahmen interpretationstheoretischer Überlegungen [...] kaum aussagekräftig: Jede Interpretation ist in der einen oder anderen Weise damit befasst, bestimmte Textbefunde zu erklären, und sie zielt darauf, (in Abhängigkeit von der jeweiligen Bedeutungstheorie) ein besseres Verständnis des Werkes zu befördern.

In reality, there is no such thing as a context-independent meaning for a word. As argued by Firth [1935], “the complete meaning of a word is always contextual, and no study of meaning apart from context can be taken seriously”. An obvious manifestation of this is the case of polysemy: some words have obvious multiple senses- [...]. Using a single vector for all forms is problematic. In addition to the multiple senses problem, there are also much subtler context-dependent variations in word meaning.

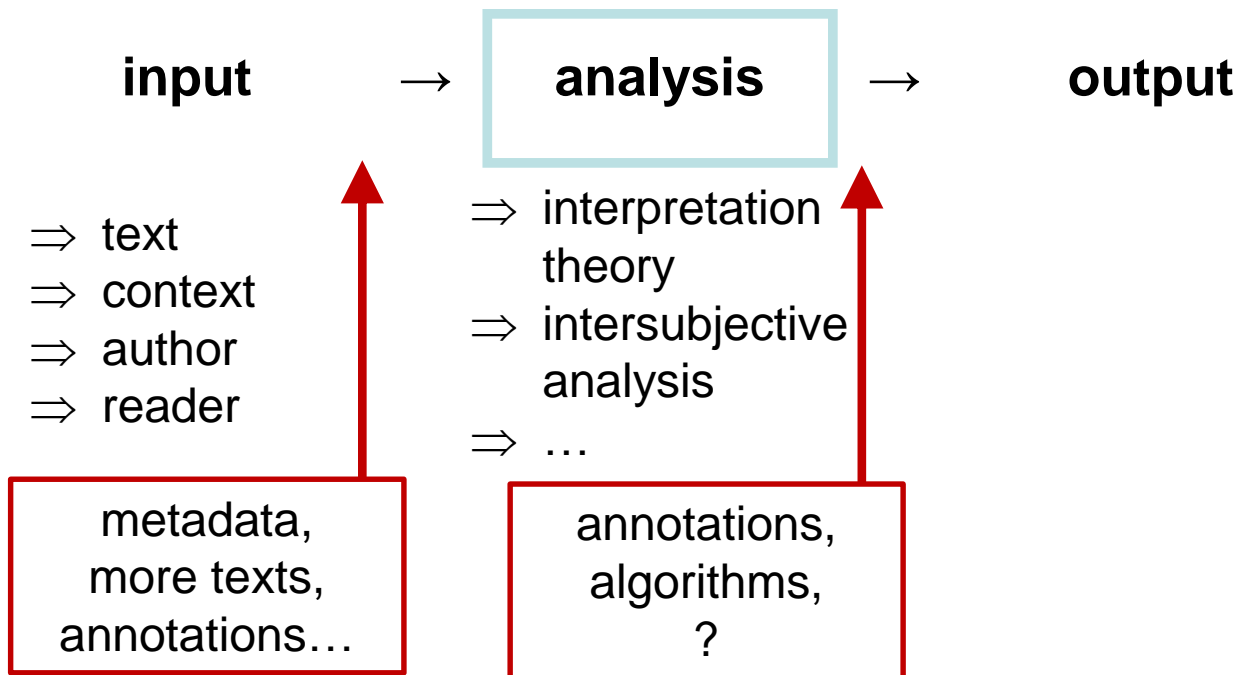
Goldberg, Yoav (2017): Neural network methods for natural language processing, San Rafael 2017 Synthesis lectures on human language technologies 37, p. 134

Analysis: Intersubjectivity and Traceability of Results

Cf. Computer Bias and Ethics

- Critique of the *kNN*-Algorithm (cf. Chun 2017, Dobson 2019)
- IEEE (Institute of Electrical and Electronics Engineers): „The IEEE Global Initiative on Ethics of Autonomous and Intelligent Systems“ in 2019. Principles 5 and 6:
 - Transparency–The basis of a particular A/IS decision should always be discoverable
 - Accountability–A/IS shall be created and operated to provide an unambiguous rationale for all decisions made

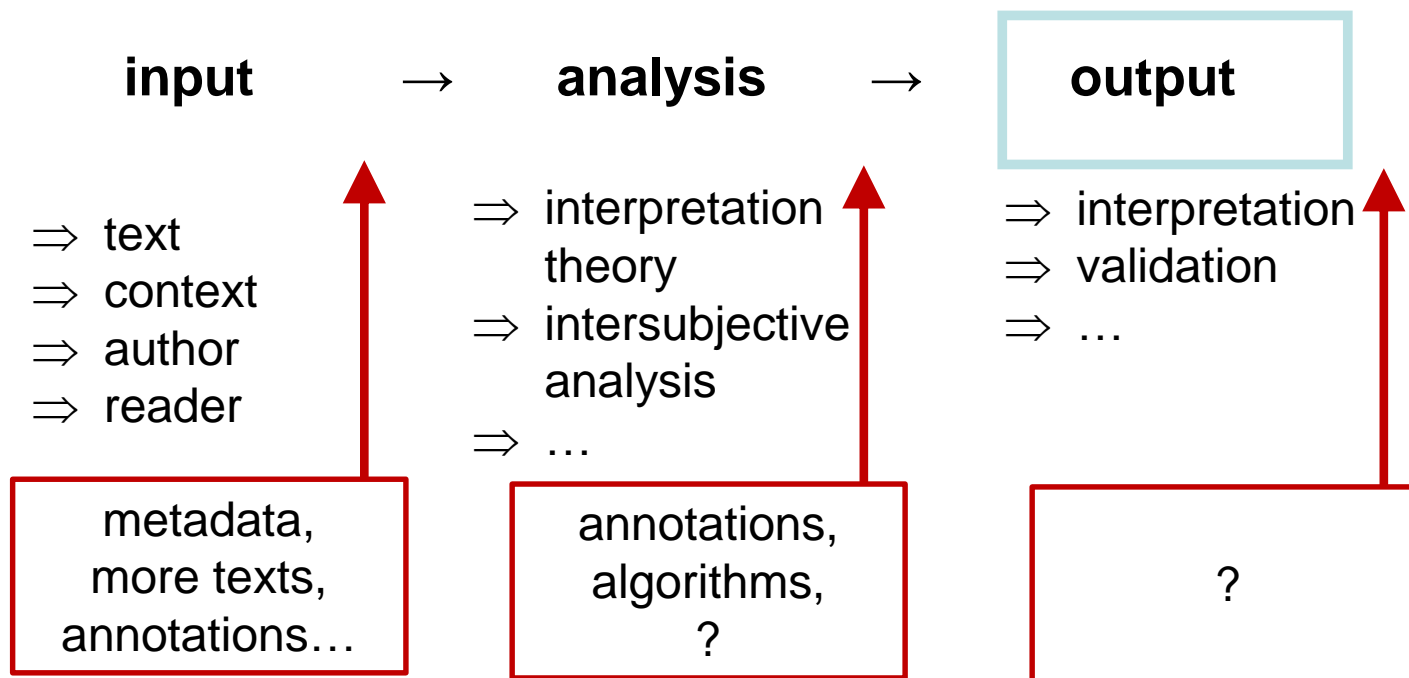
Data in Literary Studies: Input, Analysis



Data in Literary Studies: Input, Analysis, Output



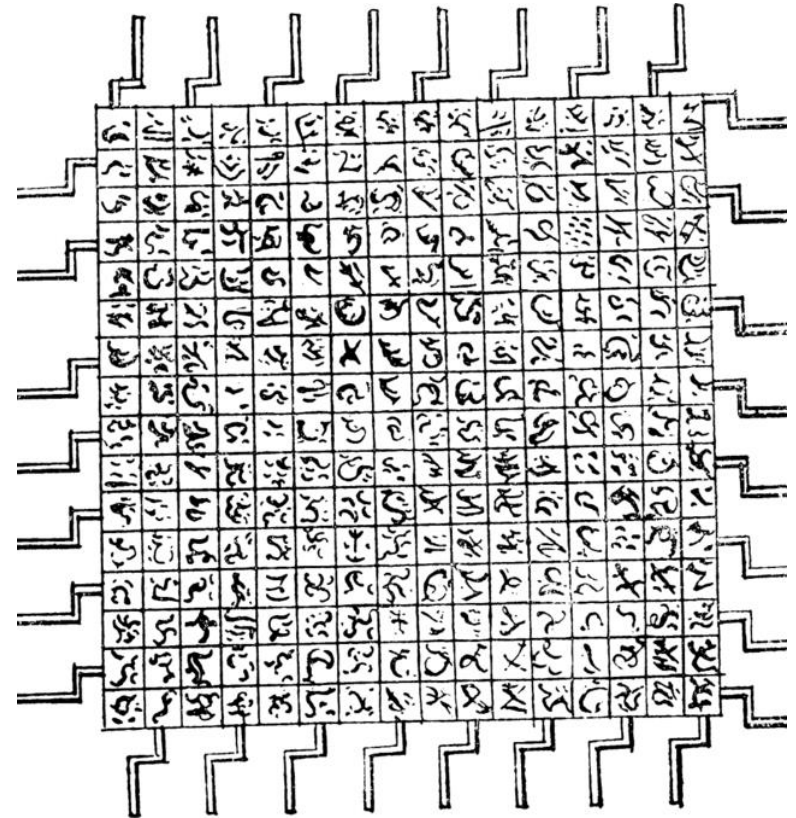
TECHNISCHE
UNIVERSITÄT
DARMSTADT



Computational Literary Studies

.. Every one knew how laborious the usual method is of attaining to arts and sciences; whereas, by his contrivance, the most ignorant person, at a reasonable charge, and with a little bodily labour, might write books in philosophy, poetry, politics, laws, mathematics, and theology, without the least assistance from genius or study.

[...] The pupils, at his command, took each of them hold of an iron handle, whereof there were forty fixed round the edges of the frame; and giving them a sudden turn, the whole disposition of the words was entirely changed. He then commanded six-and-thirty of the lads, to read the several lines softly, as they appeared upon the frame; and where they found three or four words together **that might make part of a sentence**, they dictated to the four remaining boys, who were scribes.



Swift, Jonathan (1726). [Gulliver's Travels](https://upload.wikimedia.org/wikipedia/commons/6/6d/The_Engine_%28Gulliver%29.png). p. Part 3, Chapter 5.

https://upload.wikimedia.org/wikipedia/commons/6/6d/The_Engine_%28Gulliver%29.png

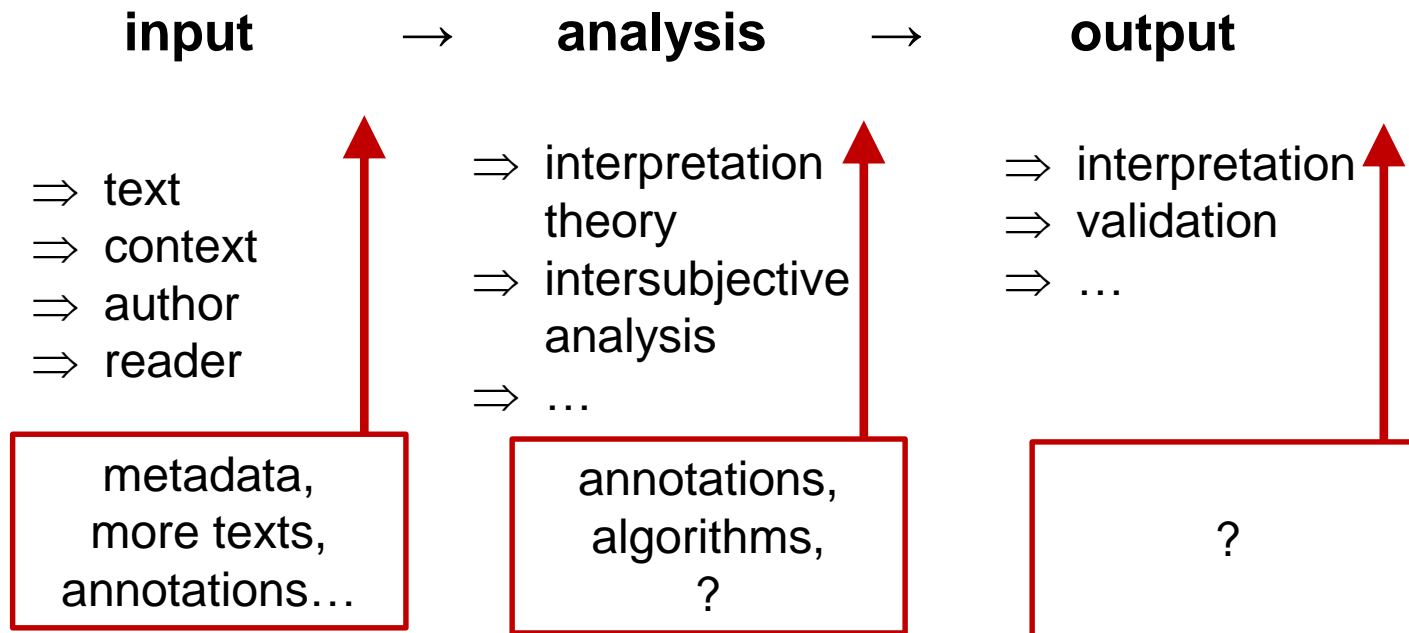
Data in Literary Studies: Input, Analysis, Output

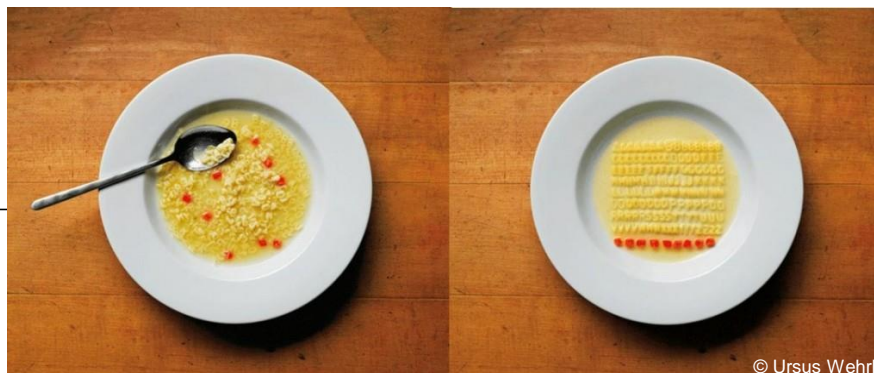


TECHNISCHE
UNIVERSITÄT
DARMSTADT

Minimum Requirements for Computational Literary Studies

1. Text Concept
2. Interpretation Theory
3. Interpretation Method





References

- Chatila, Raja, and John C. Havens. 2019. "Ethically Aligned Design. A Vision for Prioritizing Human Well-Being with Autonomous and Intelligent Systems." In *Robotics and Well-Being*, edited by Maria Isabel Aldinhas Ferreira, João Silva Sequeira, Gurvinder Singh Virk, Mohammad Osman Tokhi, and Endre E. Kadar, 95:11–16. Cham: Springer International Publishing. https://doi.org/10.1007/978-3-030-12524-0_2.
- Chun, Wendy Jui Kyong. 2017. "We're All Living in Virtually Gated Communities and Our Real-Life Relationships Are Suffering." *Wired UK*, April 13, 2017. <https://www.wired.co.uk/article/virtual-segregation-narrows-our-real-life-relationships>.
- Dobson, James E. 2019. "The Cultural Significance of K-NN." In *Critical Digital Humanities*, 101–30. The Search for a Methodology. University of Illinois Press. <http://www.jstor.org/stable/10.5406/j.ctvfjd0mf.8>.
- Firth, J. R. 1935. "The Technique of Semantics." *Transactions of the Philological Society* 34 (1): 36–73. <https://doi.org/10.1111/j.1467-968X.1935.tb01254.x>.
- Gius, Evelyn, and Janina Jacke. 2017. "The Hermeneutic Profit of Annotation. On Preventing and Fostering Disagreement in Literary Text Analysis." *International Journal of Humanities and Arts Computing* 11 (2): 233–54. <https://doi.org/10.3366/ijhac.2017.0194>.
- Gius, Evelyn, Nils Reiter, and Marcus Willand. 2019a. *A Shared Task for the Digital Humanities - Special Issue (Journal of Cultural Analytics)*.
- . 2019b. "A Shared Task for the Digital Humanities Chapter 2: Evaluating Annotation Guidelines." *Journal of Cultural Analytics*. <https://doi.org/10.22148/16.049>.
- Goldberg, Yoav. 2017. *Neural Network Methods for Natural Language Processing*. Synthesis Lectures on Human Language Technologies 37. San Rafael: Morgan & Claypool Publishers.
- Köppe, Tilmann, and Simone Winko. 2013. "Theorien und Methoden der Literaturwissenschaft." In *Methoden und Theorien*, edited by Thomas Anz, 2:285–371. Handbuch Literaturwissenschaft. Stuttgart: Metzler.
- Meister, Jan Christoph. 1999. "Jenseits von Lagado. Literaturwissenschaftliches Programmieren Am Beispiel Des Kodierungsprogramms Move Parser 3.1." *Jahrbuch Für Computerphilologie* 1: 71–82.
- Nowvskie, Bethany. 2004. "Speculative Computing: Instruments for Interpretive Scholarship." University of Virginia. <http://nowvskie.org/dissertation.pdf>.
- Schöch, Christof. 2017. "Quantitative Analyse." In *Digital Humanities: eine Einführung*, edited by Fotis Jannidis, Hubertus Kohle, and Malte Rehbein, 279–298. Stuttgart: J.B. Metzler Verlag.

Conceptual Coverage

Is the narrative level concept explicitly described?

Explanation: Narrative levels can be described or defined. This depends on the narratology used; some of them are structuralist, others are post-structuralist. Regardless of the mode, is the description/definition understandable and clear?

- 1: I did not understand what the guideline describes as “narrative level”.
- 4: I fully understood the concept described in the guideline.

Is the narrative level concept based on existing concepts?

Explanation: The level concepts can be self-designed, oriented on existing narratologies or copied from an existing level definition

- 1: The theory relation of the used level concept is not clear.
- 4: It is clearly mentioned whether the level concept is made up or (partially) based on a theory.

How comprehensive is the guideline with respect to aspects of the theory? Does it omit something?

Explanation: If the guideline is based on a theory or multiple theories, does it include the whole theory or only parts of it? Are there reasons mentioned why aspects are in-/excluded?

- 1: The guideline does not clearly state the extension of its dependence on theory/ies.
- 4: The guideline unambiguously states the scope of its theory-dependance.

How adequately is the narrative level concept implemented by this guideline in respect to narrative levels?

Explanation: Narratologies differ in their complexity. Firstly, you have to decide whether complexity or simplicity (in relation to x) is desirable, then you have to answer:

- 1: The guideline is too simple or too complex for narrative levels and thus not adequate.
- 4: The guideline's complexity is adequate.

1. How easy is it to apply the guideline for researchers with a narratological background?

Explanation: The question asks for an assessment of the ease of use of the guideline for an annotator with some narratological background. Indicators can be: Complexity of the concepts, length of the guideline, clarity of examples, clear structure, difficulty of finding special cases, etc.

- 1: Even as a narratology expert, I needed to read the guideline multiple times and/or read additional literature.
- 4: The guideline is very easy to apply, and I always knew what to do.

2. How easy is it to apply the guideline for researchers without a narratological background?

Explanation: The question asks for an assessment of the ease of use of the guideline if we assume an annotator who doesn't have a narratological background (e.g., an undergraduate student). Indicators can be: Complexity of the concepts, length of the guideline, use of terminology, clarity of examples, reference to examples only by citation, clear structure, difficulty of finding special cases, etc.

- 1: Non-experts have no chance to use this guideline.
- 4: The guideline is very easy to apply, and non-experts can use them straight away.

3./4. Inter-annotator agreement: gamma scores

Thought experiment: Assuming that the narrative levels defined in the annotation guideline can be detected automatically on a huge corpus. How helpful are these narrative levels for an interesting corpus analysis?

Explanation: This question focuses on the relevance of the narrative level annotations for textual analysis of large amounts of texts, e.g., for the analysis of developments over time with regard to narrative levels or a classification of texts with regards to genre, based on narrative levels.

- 1: The narrative levels annotations are irrelevant for corpus analysis.
- 4: The annotations provide interesting data for corpus analysis.

How helpful are they as an input layer for subsequent corpus or single text analysis steps (that depend on narrative levels)?

Explanation: The analysis of some other textual phenomena depends on narrative levels, e.g., chronology should be analyzed within each narrative level before analyzing it for the whole text. This question asks whether the analysis of such phenomena is possible or even better when based on the narrative level annotations.

- 1: The usage of the narrative levels annotations makes no difference for subsequent analyses.
- 4: Subsequent analyses are possible only because of the narrative level annotations.

Do you gain new insights about narrative levels in texts by applying the foreign guideline, compared to the application of your own guideline?

Explanation: In most cases annotating a text in accordance to a guideline changes the evaluation of textual phenomena in the text, e.g., the quality (or quantity) of narrative levels in the text.

- 1: It doesn't make a difference—I get no additional insights with the foreign guideline.
- 4: I gain a lot of new insights about narrative levels in texts based on this guideline.

Does the application of this guideline influence your interpretation of a text?

Explanation: Interpretations are normally based on the analysis of a text and thus on the observation of the presence (or absence) of certain textual phenomena. Therefore, the application of the guidelines may result in annotations that are relevant for your interpretation, e.g. the detection of a narrative level of a certain type may influence your interpretation of the reliability of a narrator.

- 1: My interpretation is independent from the annotations based on the guideline.
- 4: My interpretation is based primarily on the annotations based on the guideline.

A Digital Humanities Problem Feature



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Humanist Discussion Group, Vol. 33, No. 603.
Department of Digital Humanities, King's College London
Hosted by King's Digital Lab
www.dhhumanist.org
Submit to: humanist@dhhumanist.org

Date: 2020-02-13 16:37:10+00:00
From: Willard McCarty <willard.mccarty@mccarty.org.uk>
Subject: missed opportunities

As you're probably aware, not a few messages posted on Humanist are taken from other, more specialised lists. No problem here, except that the overall tendency of activity's migration out into the older disciplines marks lost opportunities for digital humanities.

For sure, this migration is a good sign of acceptance, for which we have fought long and hard against indifference and determined resistance. (Us older ones could tell many stories...) But these other disciplines are not primarily methodological; for me, the heart and soul of digital humanities is. On the positive side, this means that everything that happens with computing in whatever discipline is (potentially) relevant, even of prime interest. But the tendency to think only within the blinkered scope of a single discipline seems to be preventing our colleagues from seeing this.

Mind you, I am not advocating the breaking down of disciplinary boundaries, rather expanding outward from one's discipline of origin. Digital humanities provides a language, as it were, for all -- one that shows another way of seeing things.

Consider, then, if you would, encouraging our colleagues to keep one eye on the methodology of their work wherever it happens and reporting what's happening here. Digital classicists, digital historians et al., you have much of great interest to say to those who are not classicists, historians et al.

Many thanks.

Yours,
WM

--
Willard McCarty (www.mccarty.org.uk/),
Professor emeritus, Department of Digital Humanities, King's College
London; Editor, Interdisciplinary Science Reviews
(www.tandfonline.com/loi/yisr20) and Humanist (www.dhhumanist.org)

Unsubscribe at: <http://dhhumanist.org/Restricted>
List posts to: humanist@dhhumanist.org
List info and archives at: <http://dhhumanist.org>
Listmember interface at: <http://dhhumanist.org/Restricted/>
Subscribe at: http://dhhumanist.org/membership_form.php

Literary Text Analysis as Data Analysis

input → analysis → output

what data is needed to
represent the text(s)?

which qualities in the
data should be
analysed?
how?

How can the results be
interpreted?



Kempfert, I., Anwar, S., Friedrich, A., Biemann, C. (2020):
Digital
History of Concepts: Sense Clustering over Time. 42.
Jahrestagung der
Deutschen Gesellschaft für Sprachwissenschaft (DGfS),
Hamburg, Germany.

<https://www.inf.uni-hamburg.de/en/inst/ab/lt/publications/2020-kempfertetal-dgfs-scot.pdf>

