**Universität Stuttgart**

Institut für Maschinelle Sprachverarbeitung

Janis Pagel
Nils Reiter
Ina Rösiger
Sarah Schulz

# A Unified Text Annotation Workflow for Diverse Goals

# Why do we annotate?

- Empirical validation of theories

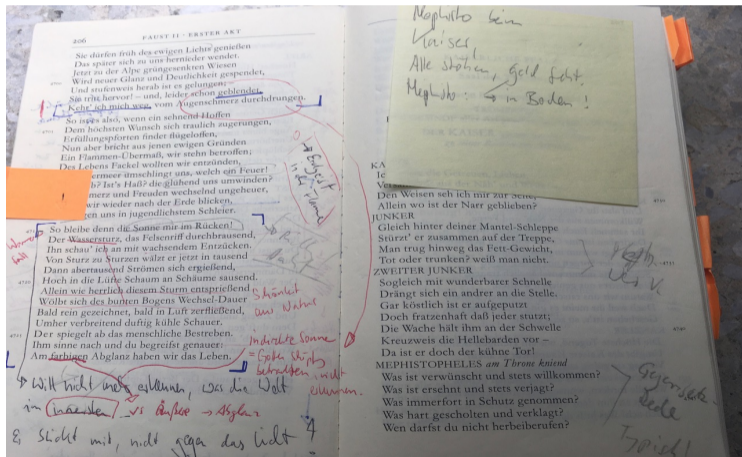- Data creation

# Why do we annotate?

- Empirical validation of theories
  - Discovering phenomena not covered by a theory
  - Strengthening definitions in a theory
    - Often confused categories might be overlapping or at least unclear
  - Uncovering implicit assumptions
- Data creation

**Why do we annotate?**

- Empirical validation of theories
  - Discovering phenomena not covered by a theory
  - Strengthening definitions in a theory
    - Often confused categories might be overlapping or at least unclear
  - Uncovering implicit assumptions
- Data creation
  - Manually annotated data can be analysed
    - Which categories are how frequent in what context?
  - Automatic tools can be evaluated
    - How well do machines do this task?
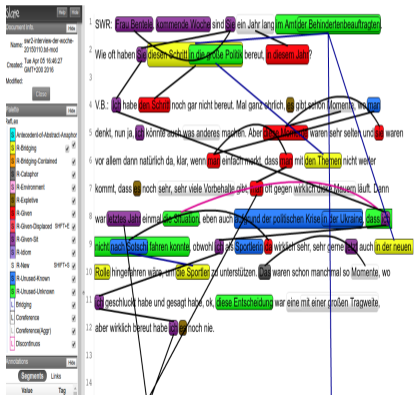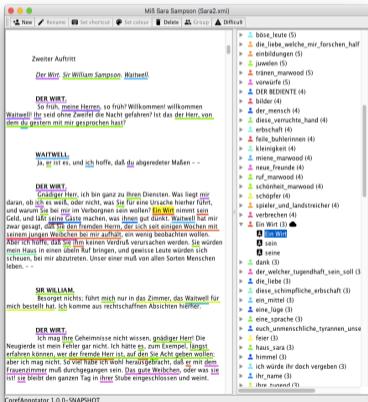  - Supervised tools can be trained

# Analog annotation …

- Ideas attached to spans of text
  - Sometimes fuzzy text spans



Janis Pagel, Nils Reiter, Ina Rösiger, Sarah Schulz, Institute for Natural Language Processing (IMS), University of Stuttgart: Unified Annotation Workflow

3

# …and digital annotation

- Explicit assignment of categories to text spans
  - Text spans are explicitly bounded (begin, end)



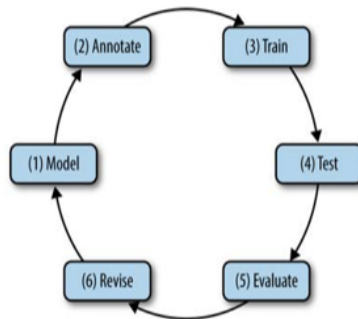Janis Pagel, Nils Reiter, Ina Rösiger, Sarah Schulz, Institute for Natural Language Processing (IMS), University of Stuttgart: Unified Annotation Workflow

4

# Circles

- **Annotation (Circle)**
  - Well known in Computational Linguistics
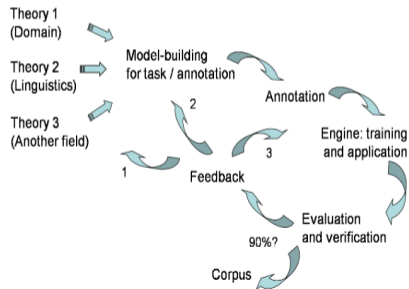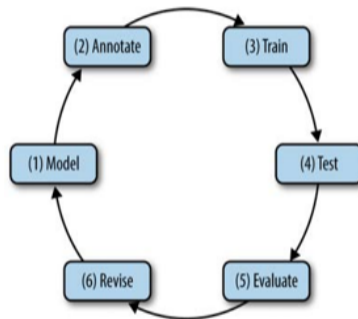  - Automation-centered
  - "Enrich with knowledge"

# Circles

- **Annotation (Circle)**
  - Well known in Computational Linguistics
  - Automation-centered
  - "Enrich with knowledge"



Pustejovsky and Stubbs (2012)

# Circles

- **Annotation (Circle)**
  - Well known in Computational Linguistics
  - Automation-centered
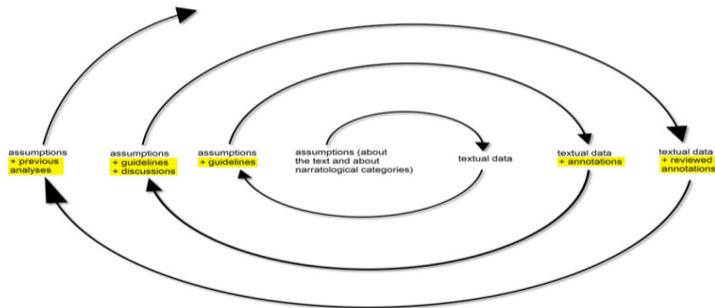  - "Enrich with knowledge"



Hovy and Lavid (2010)



Pustejovsky and Stubbs (2012)

# Circles

- **Hermeneutic Circle**
  - Well known in Humanities
  - Interpretation-centered
  - "Retrieve knowledge"

# Circles

- **Hermeneutic Circle**
  - Well known in Humanities
  - Interpretation-centered
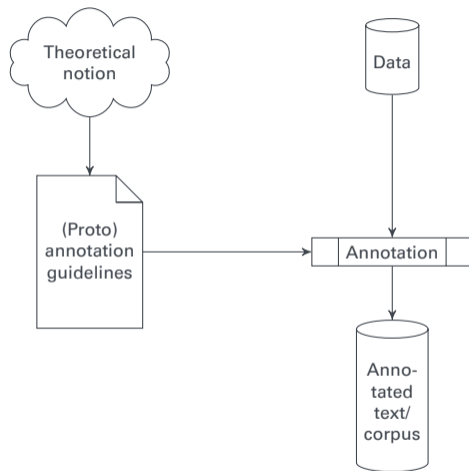  - "Retrieve knowledge"



Gius and Jacke (2017)

# Workflow

Theoretical notion

# Workflow



Theoretical notion
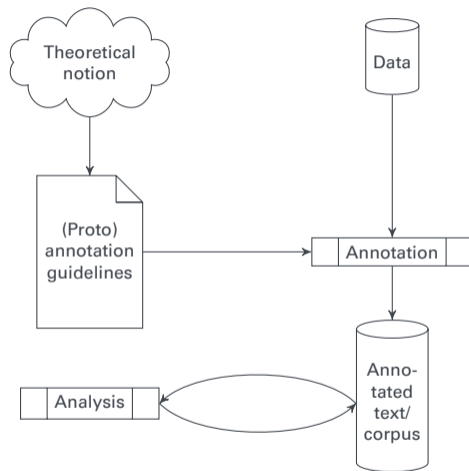
(Proto) annotation guidelines

# Workflow



Theoretical notion

Data

(Proto) annotation guidelines

# Workflow

# Workflow

# Workflow



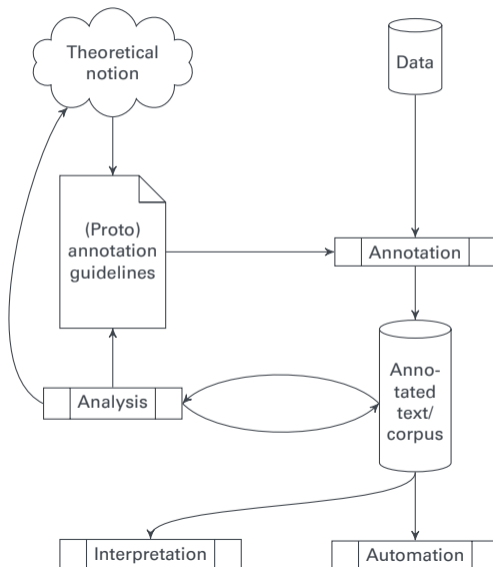Janis Pagel, Nils Reiter, Ina Rösiger, Sarah Schulz, Institute for Natural Language Processing (IMS), University of Stuttgart: Unified Annotation Workflow

7

# Workflow

# Workflow

Janis Pagel, Nils Reiter, Ina Rösiger, Sarah Schulz, Institute for Natural Language Processing (IMS), University of Stuttgart: Unified Annotation Workflow

7

# Workflow

# Workflow

# Goals of Annotation

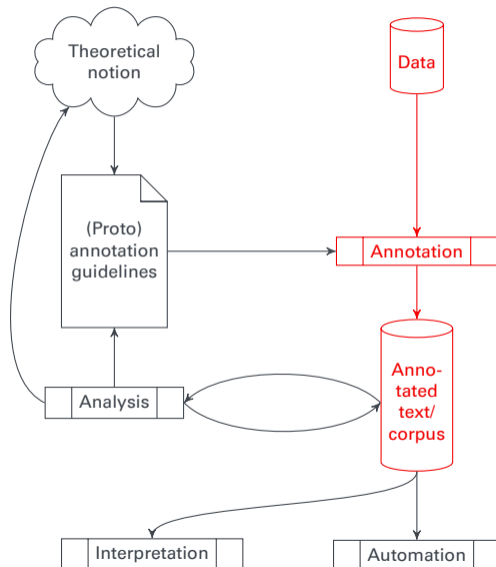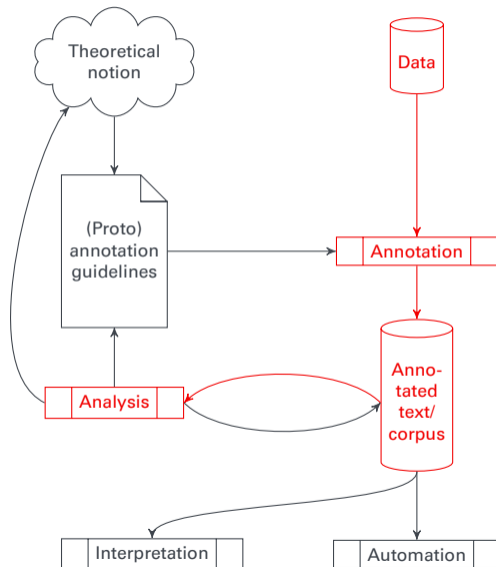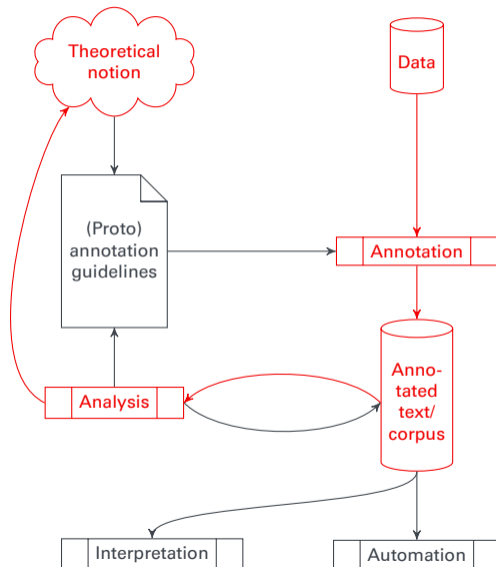- **Exploratory**
  - Mostly note-taking
  - Semi-organized
  - Humanities centered
  - Ideally completely free of presuppositions

# Goals of Annotation

- **Exploratory**
  - Mostly note-taking
  - Semi-organized
  - Humanities centered
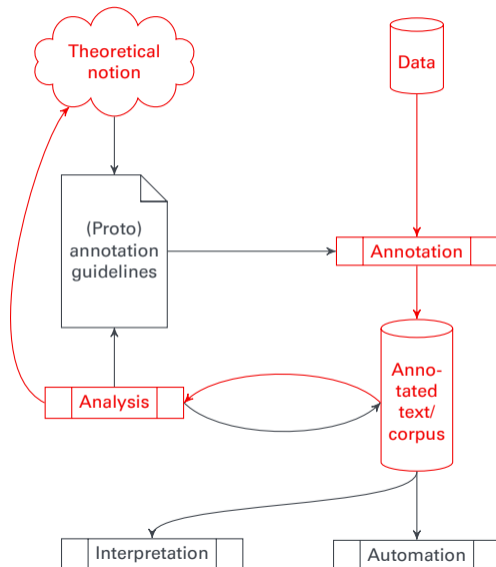  - Ideally completely free of presuppositions



Janis Pagel, Nils Reiter, Ina Rösiger, Sarah Schulz, Institute for Natural Language Processing (IMS), University of Stuttgart: Unified Annotation Workflow

8

# Goals of Annotation



- **Exploratory**
  - Mostly note-taking
  - Semi-organized
  - Humanities centered
  - Ideally completely free of presuppositions

Diagram elements: Theoretical notion, Data, (Proto) annotation guidelines, Annotation, Analysis, Annotated text/corpus, Interpretation, Automation

Janis Pagel, Nils Reiter, Ina Rösiger, Sarah Schulz, Institute for Natural Language Processing (IMS), University of Stuttgart: Unified Annotation Workflow

8

# Goals of Annotation

- **Exploratory**
  - Mostly note-taking
  - Semi-organized
  - Humanities centered
  - Ideally completely free of presuppositions

# Goals of Annotation

- **Exploratory**
  - Mostly note-taking
  - Semi-organized
  - Humanities centered
  - Ideally completely free of presuppositions

# Goals of Annotation

- **Exploratory**
  - Mostly note-taking
  - Semi-organized
  - Humanities centered
  - Ideally completely free of presuppositions

# Goals of Annotation

- **Exploratory**
  - Mostly note-taking
  - Semi-organized
  - Humanities centered
  - Ideally completely free of presuppositions
- Examples
  - *Pliny* (Bradley, 2008)
  - *3DH* (Kleymann, Meister, and Stange, 2018)



Janis Pagel, Nils Reiter, Ina Rösiger, Sarah Schulz, Institute for Natural Language Processing (IMS), University of Stuttgart: Unified Annotation Workflow

8

# Pliny



Janis Pagel, Nils Reiter, Ina Rösiger, Sarah Schulz, Institute for Natural Language Processing (IMS), University of Stuttgart: Unified Annotation Workflow

9

# 3DH

# Goals of Annotation

- **Conceptualizing**
  - Sharpen theoretic notions
  - Find disagreements of theoretical claims on concrete texts
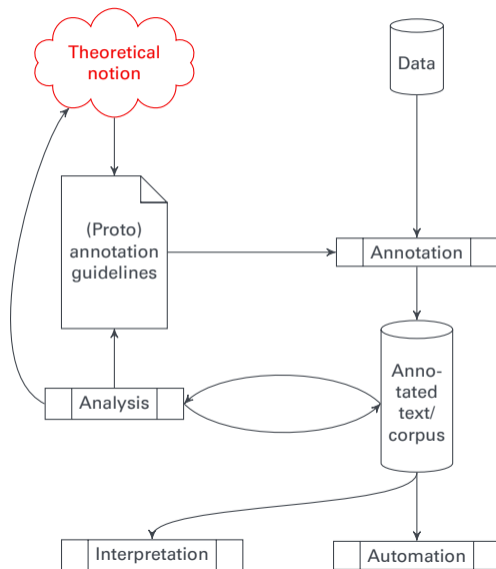
# Goals of Annotation

- **Conceptualizing**
  - Sharpen theoretic notions
  - Find disagreements of theoretical claims on concrete texts

Janis Pagel, Nils Reiter, Ina Rösiger, Sarah Schulz, Institute for Natural Language Processing (IMS), University of Stuttgart: Unified Annotation Workflow
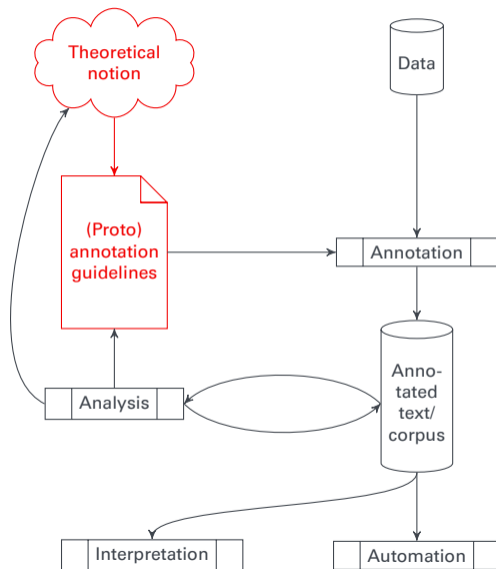
11

# Goals of Annotation
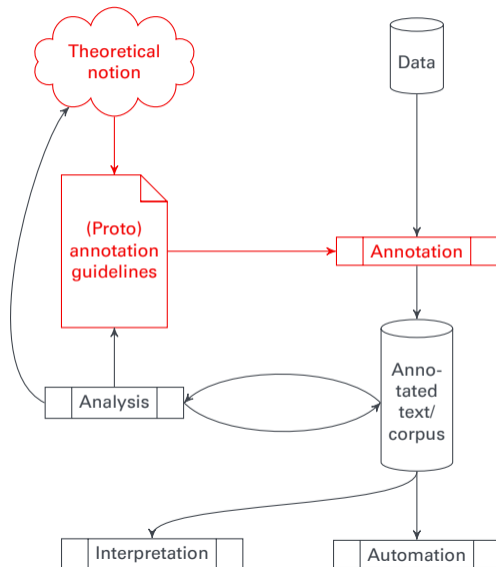
- **Conceptualizing**
  - Sharpen theoretic notions
  - Find disagreements of theoretical claims on concrete texts

# Goals of Annotation

- **Conceptualizing**
  - Sharpen theoretic notions
  - Find disagreements of theoretical claims on concrete texts
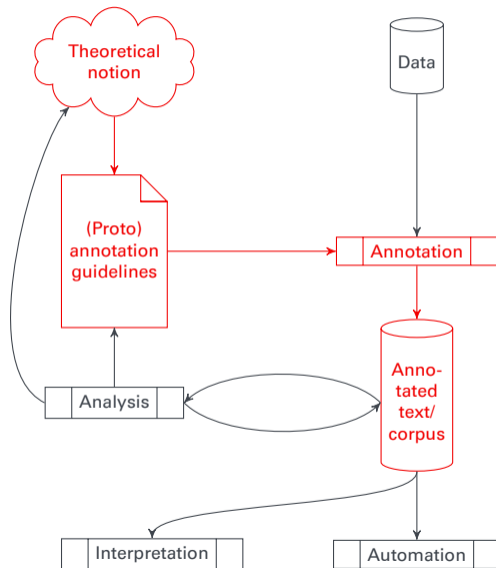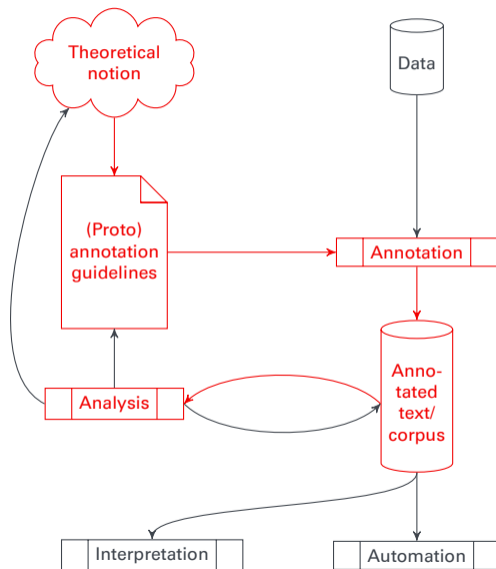
# Goals of Annotation

- **Conceptualizing**
  - Sharpen theoretic notions
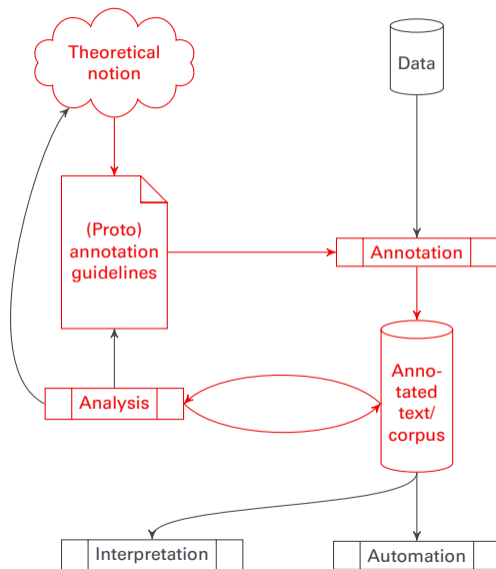  - Find disagreements of theoretical claims on concrete texts

# Goals of Annotation

- **Conceptualizing**
  - Sharpen theoretic notions
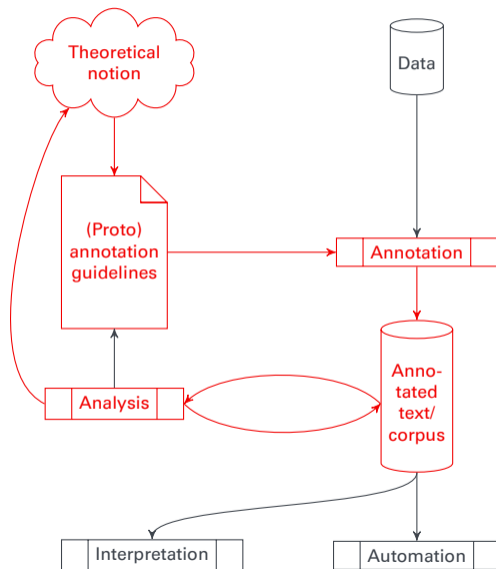  - Find disagreements of theoretical claims on concrete texts

# Goals of Annotation

- **Conceptualizing**
  - Sharpen theoretic notions
  - Find disagreements of theoretical claims on concrete texts

# Goals of Annotation

- **Conceptualizing**
  - Sharpen theoretic notions
  - Find disagreements of theoretical claims on concrete texts

# Goals of Annotation

- **Conceptualizing**
  - Sharpen theoretic notions
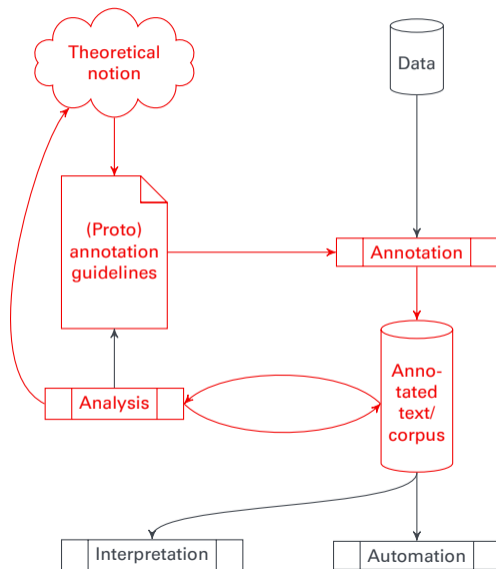  - Find disagreements of theoretical claims on concrete texts

- Examples
  - *heureCLÉA* (Bögel et al., 2015)
  - *QuaDramA* (Rösiger, Schulz, and Reiter, 2018)

# heureCLÉA (Catma)

# QuaDramA

# Goals of Annotation

- **Explicating**
  - Structure text in an observable way
  - Helps interpretation

# Goals of Annotation

- **Explicating**
  - Structure text in an observable way
  - Helps interpretation

# Goals of Annotation

- **Explicating**
  - Structure text in an observable way
  - Helps interpretation



Janis Pagel, Nils Reiter, Ina Rösiger, Sarah Schulz, Institute for Natural Language Processing (IMS), University of Stuttgart: Unified Annotation Workflow

14

# Goals of Annotation

- **Explicating**
  - Structure text in an observable way
  - Helps interpretation

# Goals of Annotation

- **Explicating**
  - Structure text in an observable way
  - Helps interpretation

Janis Pagel, Nils Reiter, Ina Rösiger, Sarah Schulz, Institute for Natural Language Processing (IMS), University of Stuttgart: Unified Annotation Workflow

14

# Goals of Annotation

- **Explicating**
  - Structure text in an observable way
  - Helps interpretation
- Examples
  - Nantke and Schlupkothen (2018)

Beispiel: Das *Rheinweinlied* und seine Intertexte

# Nantke and Schlupkothen: Modeling Complex Intertextual Relations

Janis Pagel, Nils Reiter, Ina Rösiger, Sarah Schulz, Institute for Natural Language Processing (IMS), University of Stuttgart: Unified Annotation Workflow

16

# Goals of Annotation

- **Automation-oriented**
  - Create input data for automatic learning procedures
  - Usually uses existing theoretic notions
  - NLP/CL centered

# Goals of Annotation

- **Automation-oriented**
  - Create input data for automatic learning procedures
  - Usually uses existing theoretic notions
  - NLP/CL centered

# Goals of Annotation

- **Automation-oriented**
  - Create input data for automatic learning procedures
  - Usually uses existing theoretic notions
  - NLP/CL centered

# Goals of Annotation

- **Automation-oriented**
  - Create input data for automatic learning procedures
  - Usually uses existing theoretic notions
  - NLP/CL centered

# Goals of Annotation

- **Automation-oriented**
  - Create input data for automatic learning procedures
  - Usually uses existing theoretic notions
  - NLP/CL centered

# Goals of Annotation

- **Automation-oriented**
  - Create input data for automatic learning procedures
  - Usually uses existing theoretic notions
  - NLP/CL centered
- Examples
  - GRAIN (Schweitzer et al., 2018)
  - Schulz and Kuhn (2016)

Janis Pagel, Nils Reiter, Ina Rösiger, Sarah Schulz, Institute for Natural Language Processing (IMS), University of Stuttgart: Unified Annotation Workflow

17

# POS



Janis Pagel, Nils Reiter, Ina Rösiger, Sarah Schulz, Institute for Natural Language Processing (IMS), University of Stuttgart: Unified Annotation Workflow

18

# Discussion

- Annotation in DH
  - Helps disciplines with high theoretic focus
  - Supports intersubjectivity

## Discussion

- Annotation in DH
  - Helps disciplines with high theoretic focus
  - Supports intersubjectivity
- Different goals enforce different tools
  - Not entirely clear how tools' functionality and annotation goals interrelate

## Discussion

- Annotation in DH
  - Helps disciplines with high theoretic focus
  - Supports intersubjectivity
- Different goals enforce different tools
  - Not entirely clear how tools' functionality and annotation goals interrelate
- How to treat true disagreements? (Gius and Jacke, 2017)
  - CL usually requires unambiguous annotation data
  - Humanities deal with ambiguity of concepts (different interpretations)
  - How to measure these wanted disagreements?
  - Representation problems
  - Machine learning problems

## Conclusion

- Workflow
- Goals
  - Exploratory
  - Conceptualizing
  - Explicating
  - Automation-oriented
- Example projects and areas
  - Pliny (Tool)
  - 3DH (Tool)
  - heureCLÉA (Narratology)
  - QuaDramA (Coreference)
  - Nantke and Schlupkothen (Intertextuality)
  - GRAIN (Part-of-Speech tagging)

**Universität Stuttgart**

Janis Pagel, Nils Reiter, Ina Rösiger, Sarah Schulz
Institute for Natural Language Processing (IMS), University of Stuttgart

eMail      {pageljs,reiterns,roesigia,schulzsh}@ims.uni-stuttgart.de
Phone      +49-711-685 813 {89,54,30,94}
Website    http://www.ims.uni-stuttgart.de/

# References I

Bögel, Thomas et al. (2015). "Collaborative Text Annotation Meets Machine Learning: heureCLÉA, a Digital Heuristic of Narrative". In: *DHCommons* 1.

Bradley, John (2008). "Thinking about interpretation: Pliny and scholarship in the humanities". In: *Literary and Linguistic Computing* 23.3, pp. 263–279. DOI: 10.1093/llc/fqn021. URL: http://dx.doi.org/10.1093/llc/fqn021.

Gius, Evelyn and Janina Jacke (2017). "The Hermeneutic Profit of Annotation: On Preventing and Fostering Disagreement in Literary Analysis". In: *International Journal of Humanities and Arts Computing* 11.2, pp. 233–254. DOI: 10.3366/ijhac.2017.0194. URL: https://doi.org/10.3366/ijhac.2017.0194.

Hovy, Eduard and Julia Lavid (2010). "Towards a 'Science' of Corpus Annotation: A New Methodological Challenge for Corpus Linguistics". In: *International Journal of Translation Studies* 22.1, pp. 13–36.

Kleymann, Rabea, Jan Christoph Meister, and Jan-Erik Stange (Feb. 2018). "Perspektiven kritischer Interfaces für die Digital Humanities im 3DH-Projekt". In: *Book of Abstracts of DHd 2018*. Cologne, Germany.

Nantke, Julia and Frederik Schlupkothen (2018). "Zwischen Polysemie und Formalisierung: Mehrstufige Modellierung komplexer intertextueller Relationen als Annäherung an ein 'literarisches' Semantic Web". In: *Proceedings of DHd*.

Pustejovsky, James and Amber Stubbs (2012). *Natural Language Annotation for Machine Learning: A Guide to Corpus-Building for Applications*. Sebastopol, Boston, Farnham: O'Reilly Media. ISBN: 9781449306663.

Rösiger, Ina, Sarah Schulz, and Nils Reiter (2018). "Towards Coreference for Literary Text: Analyzing Domain-Specific Phenomena". In: *Proceedings of the Joint SIGHUM Workshop on Computational Linguistics for Cultural Heritage, Social Sciences, Humanities and Literature*. Santa Fe, USA.

Schulz, Sarah and Jonas Kuhn (2016). "Learning from Within? Comparing PoS Tagging Approaches for Historical Text". In: *LREC*. Ed. by Nicoletta Calzolari et al. European Language Resources Association (ELRA). URL: http://dblp.uni-trier.de/db/conf/lrec/lrec2016.html#SchulzK16.

Schweitzer, Katrin et al. (2018). "German Radio Interviews: The GRAIN Release of the SFB732 Silver Standard Collection". In: *Proceedings of the 11th International Conference on Language Resources and Evaluation*. LREC 2018.

Janis Pagel, Nils Reiter, Ina Rösiger, Sarah Schulz, Institute for Natural Language Processing (IMS), University of Stuttgart: Unified Annotation Workflow

22