

# Guided ontology-backed semantic annotation with neonion

**Bachelors thesis:** Using Wikidata as a referential knowledge source for entity linking and interactive controlled vocabulary alignment

Jakob Höper  
hoepfer@zedat.fu-berlin.de

Institute of Computer Science  
Freie Universität Berlin

Thesis Advisor:  
Prof. Dr. Claudia Müller-Birn  
clmb@inf.fu-berlin.de

Human-Centered Computing  
Institute of Computer Science  
Freie Universität Berlin

**neonion** is a light-weight collaborative workbench for scholarly reading and annotation of documents [1]. Apart from basic forms of annotation, like highlighting and commenting selected portions of text, the software provides a model for semantic information extraction tasks like concept tagging and annotating relations [2]. It also contains modules for automation of such tasks, and even comes with a triple store engine for retrieval of semantic statements.

Although its technical capabilities and conceptual complexity are no unique features among the various existing approaches toward semantic annotation and enhancement, neonion aims at optimizing usability and follows paradigms for synergistic interaction between users and underlying automated processes. This focus on usability and mixed-initiative user interfaces [5] is a consequence of the observation that researchers who don't come from a technical background won't adopt semantic technologies as commonly as their colleagues in life sciences or technical domains. Scholars from the humanities, who could make valuable contributions as subject-matter experts on their respective fields of research, often hesitate to use available tools, expecting them not to be designed for their needs or to require too much prior knowledge about semantic web technologies.

Development of features and extensions of the neonion semantic annotation software aims at meeting the needs of subject-matter experts who are interested in generating explicit semantic data from their research and publications. In order to offer usable tooling while at the same time enabling researchers, authors and editors to generate structured extensions to their content in a way that yields interoperable and schematically sound semantics, we propose a semi-automated approach which interacts with the user based on their input and lets them choose from recommendations for vocabulary alignment.

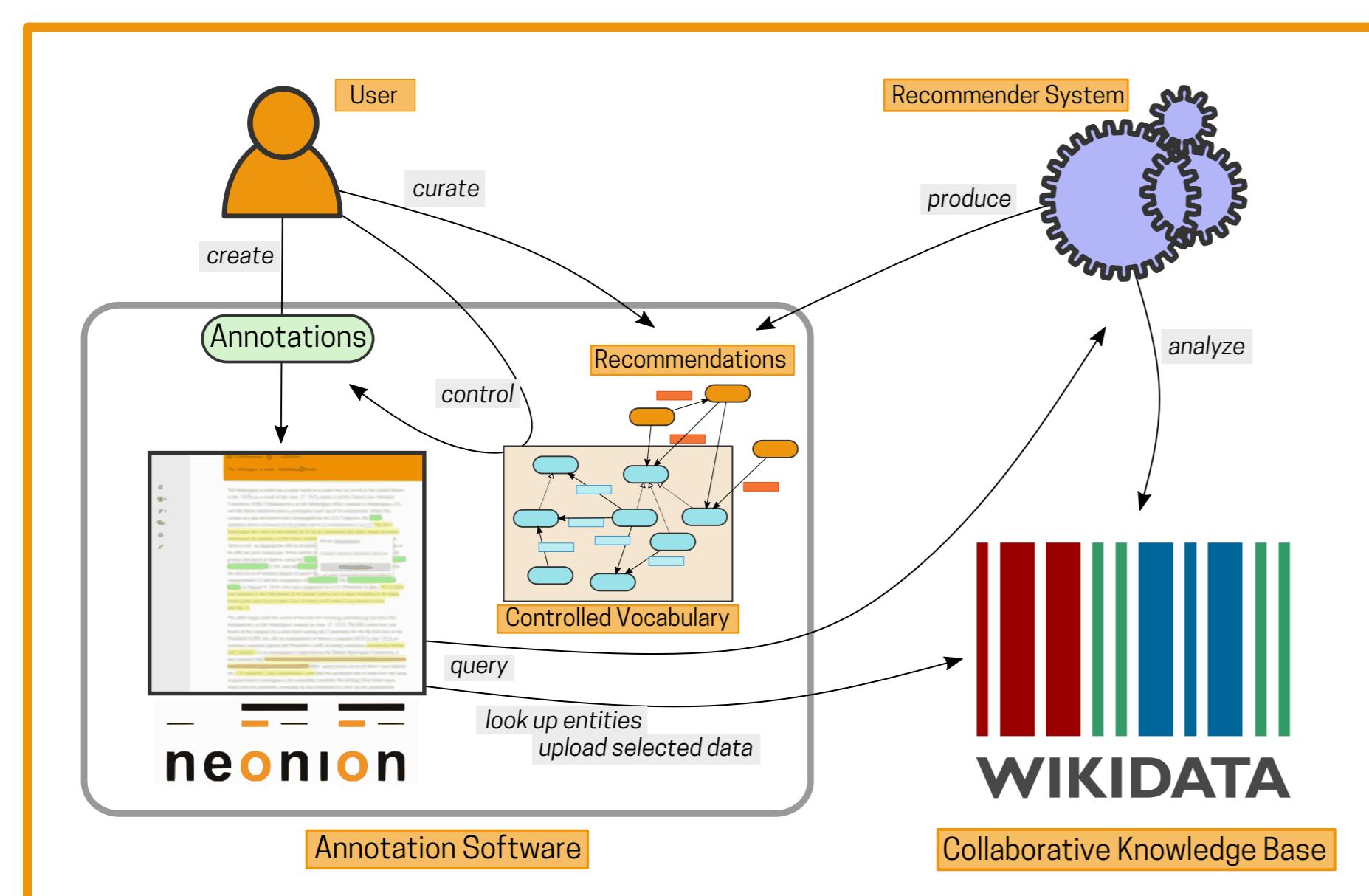
... and recording new ones (mostly the aforementioned historians), them in audio(visual) counterparts to the footage itself. A late but not was added to the "ghetto" puzzle by Claude Lanzmann in his 220 *Le Dernier des justes / The Last of the Unjust* (2013). This edits down and 24 minute long 1975 interview with the last "Jewish Elder" and adds t.5 Benjamin Murmelstein was not only a crown-witness of the end of the Holocaust but later also became a witness for the prosecution commandant of Theresienstadt, Karl Rahm.

the central provocation of *Le Dernier des justes* was that in the 1970s interview material – originally intended for *Shoah* Lanzmann broke his old rule of banning perpetrator footage from extracts from the historical Theresienstadt "ghetto" footage can be two ways: Either Lanzmann no longer believes his own central credo – probable – or he tacitly condones the behaviour of Geron and the otherants in the "ghetto" Theresienstadt in 1944. Moreover, that would also be



**Entity linking in neonion**, using Wikidata item page identifiers for contextualization. When users select some text in annotation mode, a search dialog appears next to the document and Wikidata is queried for matching items. Results are restricted to matching instances of appropriate classes in order to ensure valid input.

The proposed functionality is being added to neonion as a prototypical implementation targeting a specific workflow involving scholarly texts from an scientific online journal. In order to assess whether Wikidata as a community repository for common knowledge can be used to extract domain ontologies based on subjects covered in scholarly texts, the **open access online journal Apparatus** [3] was chosen as a specialized textual resource, with cinematographic and historical studies with focus on central and eastern Europe as the specific knowledge domain to be covered. In order to generate candidates for vocabulary amendments, association rule mining **property recommender system Snoopy** is used [4].



Without need for prior configuration, new users can start on annotating their documents with a very basic controlled vocabulary. Thanks to its mapping to Wikidata terms, semantic recommenders can easily exploit Wikidata's structured contents and determine related class items and properties which can be suggested to the user as additional concept and property elements of their vocabulary. Continuous creation of annotation generates a content-aligned vocabulary curated by the user.

"Ghetto Film" Puzzles and The Last of the Unjust

In Cheminot remained SS Sonderkommando Kullmann in April 1942, as the deputy of camp commander Herbert Lange, one of the experts in mobile genocidal technologies (Montague 2012: 532). If Otto had not died – or rather gone missing on May 6, 1945 in Prague – he would have most probably been convicted of war crimes and hanged either in Czechoslovakia or Poland – similar to one of his Kulmhof successors.<sup>4</sup>

For several decades the mysteries of the "ghetto film" have occupied Holocaust, film and military historians. At the same time until recently, footage from the "ghetto" shot by Nazi cameramen has been used in documentary films and TV programmes to illustrate the history of the European Jews – and not only during the Holocaust. Documentary filmmakers have been eager to offer solutions to the "ghetto" puzzle, regrouping survivors' voices and recording new ones (mostly the aforementioned historians), contextualising them in audio(visual) counterparts to the footage itself. A late but not conclusive piece was added to the "ghetto" puzzle by Claude Lanzmann in his 220 minute long film *Le Dernier des justes / The Last of the Unjust* (2013). This edits down his eleven hour and 24 minute long 1975 interview with the last "Jewish Elder" and adds new footage to it. Benjamin Murmelstein was not only a crown-witness of the beginning and end of the Holocaust but later also became a witness for the prosecution against the last commandant of Theresienstadt, Karl Rahm.

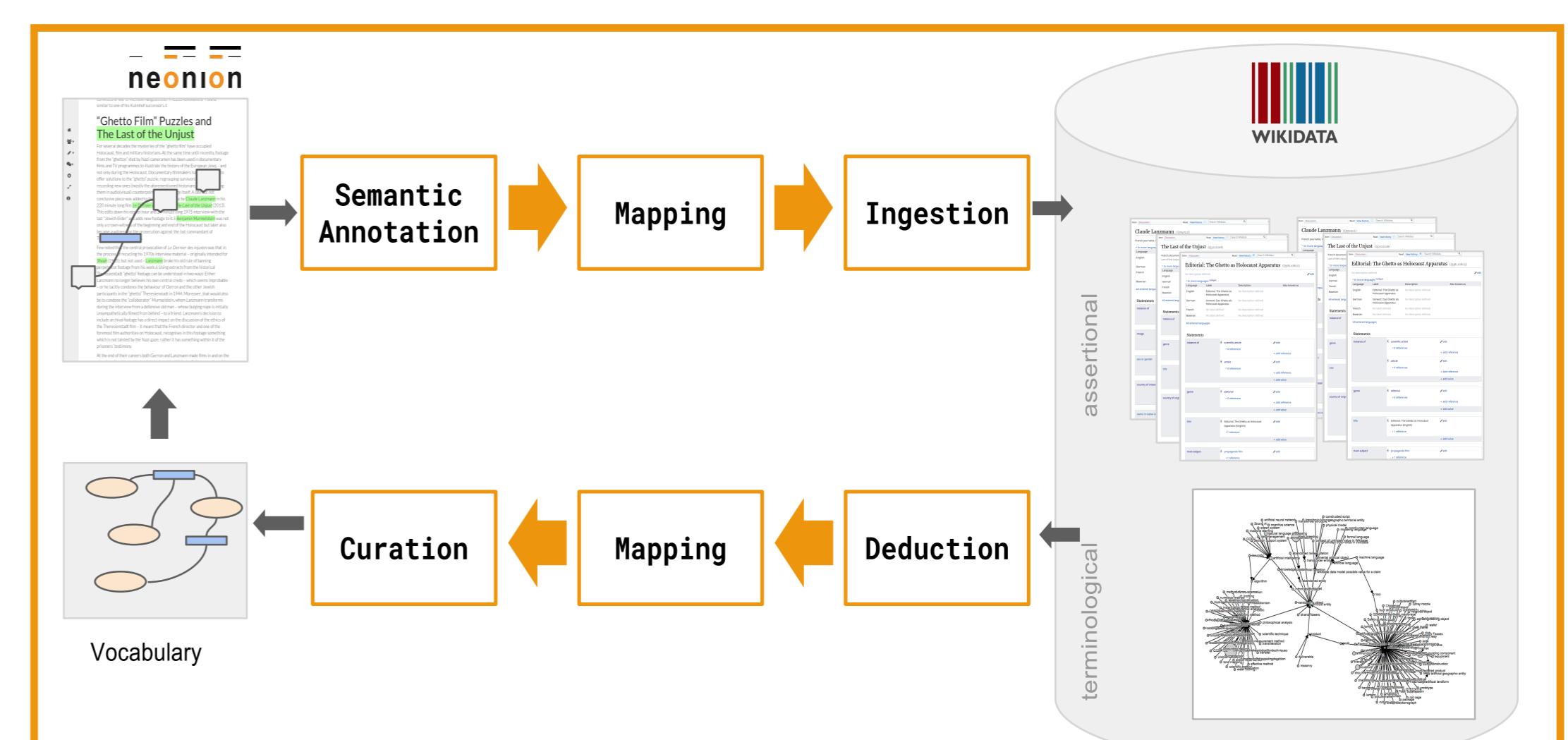
Few noted that the central provocation of *Le Dernier des justes* was that in the process of recycling his 1970s interview material – originally intended for *Shoah* (1985), but not used – Lanzmann broke his old rule of banning perpetrator footage from his work.<sup>6</sup> Using extracts from the historical Theresienstadt "ghetto" footage can be understood in two ways: Either Lanzmann no longer believes his own central credo – which seems improbable – or he tacitly condones the behaviour of Geron and the other Jewish participants in the "ghetto" Theresienstadt in 1944. Moreover, that would also to condone the "collaborator" Murmelstein, whom Lanzmann transforms during the interview from a defensive old man – whose bulging nape is initially unsympathetically filmed from behind – to a friend. Lanzmann's decision to include archival footage has a direct impact on the discussion of the ethics of the Theresienstadt film – it means that

neonion's user interface for **annotating relations**. Occurrences of Named Entities are classified by annotations assigning concepts from a controlled vocabulary. Available properties for description of relationships are visually indicated as hinted connections between entities of applicable class membership.

**Free collaborative knowledge base Wikidata** serves as a semantic backbone for structured data in Wikimedia projects. It is meant to cover of all public knowledge rather than selected scholarly and scientific domains and its large community of currently almost 20 000 active users and workflows for import of other knowledge bases such as Freebase seem to promise continuing progress. neonion is genuinely designed to allow connection to external knowledge bases in its annotation model and the internal knowledge organisation system it uses as a vocabulary for annotation. However, designing an ontology to formally define terminology and semantics for a specific research field requires not only subject-matter knowledge and labour time, but also comprehension of and skill in employed technologies.

A closer connection of neonion and Wikidata is expected to achieve mutual benefits.

- Use terminological knowledge gathered from Wikidata to extend annotation vocabulary: Align vocabulary with domain of annotated content.
- Use mapping of vocabulary elements to extract Wikidata statements from neonion annotations: Add new factual claims to Wikidata or confirm facts known to Wikidata with bibliographic references to annotated publication.



**Proposed Workflow** of mutually beneficial employ of Wikidata as an authoritative ontological reference. Users with specialized research contents and subject-matter expertise use a basic vocabulary for annotation of structured information from texts (upper left corner). A mapping from neonion's knowledge organisation system to Wikidata allows for optional contribution of extracted statements to the Wikidata knowledge base directly from within the neonion workbench. On the other hand, existing descriptions of identified entities are retrieved from Wikidata and schematic knowledge is deduced (bottom right corner). In combination with an association rule mining-based property recommender, neonion is able to suggest candidates for vocabulary extension to the user, who can thereby align new annotation choices with terminological semantics based on the actual content they work with.

## References

- [1] Müller-Birn, Claudia and Klüwer, Tina and Breitenfeld, André and Schlegel, Alexa and Benedix, Lukas. neonion – combining human and machine intelligence. In Proceedings of the 18th ACM Conference Companion on Computer Supported Cooperative Work & Social Computing, pages 223–226, 2015.
- [2] André Breitenfeld, Maximilian Mackeprang, Ming-Tung Hong, and Claudia Müller-Birn. Enabling structured data generation by nontechnical experts. In Manuel Burghardt, Raphael Wimmer, Christian Wolff, and Christa Womser-Hacker, editors, Mensch und Computer 2017 - Tagungsband, pages 181–192, Regensburg, 2017. Gesellschaft für Informatik e.V.
- [3] Natascha Drubek, Anke Hennig, and Irina Sandomirskaia. Apparatus: Zur Einführung. In Apparatus. Film, Media and Digital Cultures in Central and Eastern Europe 1, 2015. <http://www.apparatusjournal.net>
- [4] Wolfgang Gassler, Eva Zangerle, and Günther Specht. Guided curation of semistructured data in collaboratively-built knowledge bases. Future Generation Computer Systems, 31:111–119, February 2014.
- [5] Eric Horvitz. Principles of mixed-initiative user interfaces. In Proceedings of the SIGCHI conference on Human Factors in Computing Systems, pages 159–166, ACM, 1999.

