

Text Encoder and Annotator: an all-in-one editor for transcribing and annotating manuscripts with RDF

Fabio Valsecchi¹, Matteo Abrate¹, Clara Bacciu¹,
Silvia Piccini², and Andrea Marchetti¹

¹ Institute of Informatics and Telematics (IIT) {name.surname}@iit.cnr.it,

² Institute for Computational Linguistics (ILC) {name.surname}@ilc.cnr.it,
National Research Council (CNR), Via G. Moruzzi 1, Pisa, Italy

Abstract. In the context of the digitization of manuscripts, transcription and annotation are often distinct, sequential steps. This could lead to difficulties in improving the transcribed text when annotations have already been defined. In order to avoid this, we devised an approach which merges the two steps into the same process. Text Encoder and Annotator (TEA) is a prototype application embracing this concept. TEA is based on a lightweight language syntax which annotates text using Semantic Web technologies. Our approach is currently being developed within the *Clavius on the Web* project, devoted to studying the manuscripts of Christophorus Clavius, an influential 16th century mathematician and astronomer.

Keywords: Manuscript Transcription, Annotation, RDF, Semantic Web

1 Introduction

Within the field of Digital Humanities, several projects are devoted to preserving, analyzing and studying the large amounts of manuscripts, books, newspapers, maps, photos and paintings stored in archives, museums and libraries around the world.

The *Clavius On the Web*³ project [2] aims to restore and enrich the manuscripts written by Christophorus Clavius (1538-1612), one of the most respected and influential mathematicians and astronomers of his time. Among the works preserved by the Historical Archives of the Pontifical Gregorian University (APUG) there is the autograph manuscript, used for the printed edition of 1574, entitled *Euclidis Elementorum Libri XV. Accessit XVI de solidorum regularium comparatione. Omnes perspicuis demonstrationibus, accuratissimeque scholiis illustrati*. It is an annotated translation from Greek into Latin of Euclid's Elements, the famous text of arithmetic and geometry from the 3rd century BC. The text was considered one of the most comprehensive and authoritative of the 16th century

³ <http://claviusontheweb.it/>

so that personalities such as René Descartes, Marin Mersenne and Johannes Kepler built their knowledge on it.

In this context, Semantic Web technologies, such as the Resource Description Framework (RDF), make it possible to approach text annotation in an innovative way. RDF-based annotations provide a method for enriching texts using structured data already described in details and maintained by the Semantic Web community (i.e., Linked Data sets). In addition, due to the interlinked structure of Linked Data, RDF-based annotations produce valuable annotated documents, characterized by a strong connection with external resources.

1.1 Background

Text annotation consists in attaching additional information such as comments, tags or links, to specific portions of a text. Annotations can be mainly performed using two methods: *inline* and *standoff*. The former directly includes annotations within the text, while the latter defines them in a different location. Inline markup keeps the annotations and the annotated text close together, but it has the drawback of weighing the document down. Moreover, depending on the complexity of the markup language, the text could become hard to read. This aspect is crucial and must be taken into account when developing manual annotation tools, as users need to be able to read the annotated text with ease. Last but not least, a complex and heavyweight markup language could make the manual annotation process even more difficult since users have to firstly know all the syntax rules and secondly write a considerable amount of additional markup.

In contrast, the standoff approach does not have markup overloading problems due to its total independence from the resource text. Annotations are in fact separately defined in a different location where the relative text offsets are stored and kept up to date. In addition, standoff markup has the advantage of allowing overlapping annotations. Nevertheless, this approach has some drawbacks related to the sequential process of transcribing and annotating. Typically, the available tools of this type separate the transcription and the annotation phases. However, if a transcription error is found during the subsequent annotation phase, it is necessary to recompute the offsets in order to reflect the changes in the transcribed text. This implies an automatic recomputation of the offsets, a process that could be complex and costly. The logical conclusion is therefore to make transcription and annotation a joint process.

1.2 Our Approach

The core idea of this work is to combine transcription and annotation of text, thus streamlining the workflow process. Hence, we devised a lightweight language that enables this continuous and mixed process. Another key-point is that we propose to treat every textual phenomenon as an annotation independently from its specific type (e.g., semantic, syntactic, lexical). Every portion of text is treated in the same way, and RDF-based annotations are used to describe their content. RDF supports our purpose by allowing any possible annotation to be specified

using the enormous amount of ontologies, vocabularies and Linked Data sets available on the Web.

To the best of our knowledge, none of the existing tools include this approach. For instance, Pundit [4] and Refer.cx [10] provide RDF-based annotation features but limited only to web pages and without the possibility of transcribing text. Brat [9] allows for text annotation supported by Natural Language Processing technology, however it does not provide a transcription feature. RDFaCE [6] is a text editor for annotating text using a graphical UI and displaying results with different views such as WYSIWYG and WYSIWYM (i.e., What You See Is What You Get/Mean).

2 Text Encoder and Annotator

In the light of the above, this article presents the Text Encoder and Annotator, a Web application, which provides an editor to transcribe texts and a lightweight language for annotating them with RDF (all within the same environment). It was envisaged for linguists, historians and more generally for scholars and students. We devised a layout composed of three main views, horizontally placed along the interface of the application (Figure 1):

1. *Image box*: it displays the digitized image of a specific manuscript to be transcribed and annotated.
2. *Editor box*: it is the main component of the tool. It is used to write text and enrich it with RDF-based annotations, which follow the specific language syntax described below. Text highlighting identifies the markup and improves its legibility. Moreover, a top bar contains shortcut buttons for inserting some basic annotations, characterized by common RDF predicates, such as *rdfs:seeAlso*, which can be used to link to external resources, *rdfs:comment* to specify text comments and *foaf:page* to include hyperlinks related to the topic the annotation is about.
3. *Diagram box*: it is an optional view that can be activated on-demand if the user needs a summary of their annotations. The node-link diagram contains a white node representing the whole text of the document, blue nodes identifying the portion of text referred to in an annotation, orange nodes displaying the identifiers of the annotation and gray nodes describing the objects of the RDF triples specified. The edges between orange and gray nodes represent the predicates of the triples defined in the annotations.

2.1 Lightweight language

Considering the pros and cons of inline and standoff annotations discussed above, we think that a hybrid lightweight syntax is a suitable solution to fulfill our requirement of having a simultaneous workflow of transcription and annotation.

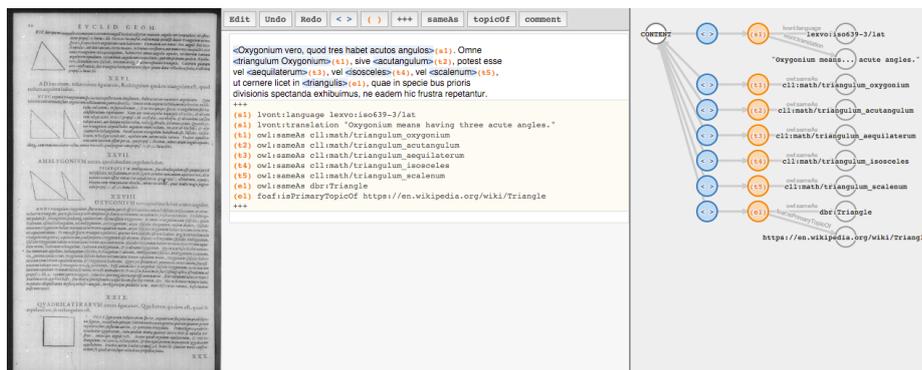


Fig. 1. The interface of the application is mainly composed by three views. From left to right there is a box containing the image of a document, an editor allowing the transcription and annotation and a diagram displaying a visual summary of the annotations. The prototype tool is available on github at <http://github.com/nitaku/TEA>.

We propose a Lightweight Markup Language⁴ (LML) employed only within the interface of TEA, in order to provide an easy and quick way of transcribing and annotating text using RDF. The main reason behind the adoption of an LML is that common markup languages based on XML (e.g., TEI) are not easy to write and read in their raw form, due to their complex syntax. Moreover, the use of LMLs has already proven to be beneficial in other systems such as the Leiden-plus⁵ language employed in the papyri.info editor. It is worth clarifying that we do not propose our approach as a format for the representation and interchange of texts, consequently it cannot be compared to standards such as the Text Encoding Initiative (TEI) [5]. Furthermore our language markup should be considered as distinct from semantic markup languages like Microformat [7], RDFa [3] and Microdata [8] since it has not been not devised for annotating HTML and XML documents.

Our language is used for marking portions of text and assigning them identifiers that will be used in a different, reserved and “standoff-like” section of the text where the details of the annotations are specified (Figure 2). More precisely, a portion of text can be annotated by enclosing it in a *span* using angle brackets, while round brackets specify a string identifying the annotation (i.e., *<annotated portion of text>(identifier)*). This inline syntax uses a very limited amount of characters and does not weigh the text down too much, keeping it easy-to-read. Identifiers are then used in a distinct part of the text, called *directive section*, where the annotation body is specified. Three plus signs (i.e., *+++*) are used both to open and close this section that can be repeated within the text more than once. Inside this block, annotations can be specified as RDF triples with the identifier of a certain span of text as subject. The choice of predicates and objects

⁴ http://en.wikipedia.org/wiki/Lightweight_markup_language

⁵ <http://papyri.info/editor/documentation?docotype=text>

is totally free, although some default predicates are suggested in order to perform the most common and basic annotations (e.g., `rdfs:seeAlso`, `rdfs:comment`, `foaf:page`). Currently, the syntax used in the directive section defines triples as three text values separated by a space.

2.2 Example

We here provide an example of annotation performed on a portion of text extracted from the *Euclidis Elementorum Libri XV. Accessit XVI* [1] (paragraph 30) by Clavius. A free English translation of the Latin text fragment, shown in Figure 2, follows: "It is called Oxigonium as it has three acute angles. Every Oxigonium triangle, or Acutangulum triangle, could be either Equilateral, or Isosceles or Scalene as you can see from the classification provided above and not reported here". Figure 2 shows the code resulting from the text encoding and annotation process. Portions of text are marked using spans while the body of the annotations is specified within the directive section. The annotation identifiers (e.g., `s1`, `t1`, `t2`) are used as the subjects of the triples while predicates and objects are freely chosen by annotators. The Latin language (i.e., `lexvo:iso639-3/lat`) and the translation have been specified in the first annotation `s1` using the Lexvo ontology⁶. Lexical entries of the mathematical lexicon of Clavius⁷ (e.g., `cll:math/triangulum_oxygonium`) and the DBpedia Triangle resource (i.e., `dbr:Triangle`) have been linked through the *seeAlso* predicate of the RDF Schema⁸. The triangle entry of Wikipedia has been specified as an interesting web page (i.e., `foaf:page`) for the annotation `e1` using the FOAF vocabulary.

```
<Oxygonium vero, quod tres habet acutos angulos>(s1). Omne
<triangulum Oxygonium>(t1), sive <acutangulum>(t2), potest esse
vel <aequilaterum>(t3), vel <isosceles>(t4), vel <scalenum>(t5),
ut cernere licet in <triangulis>(e1), quae in speciebus prioris
divisionis spectanda exhibuimus, ne eadem hic frustra repetantur.
+++
(s1) lvont:language lexvo:iso639-3/lat
(s1) lvont:translation "...
(t1) rdfs:seeAlso cll:math/triangulum_oxygonium
(t2) rdfs:seeAlso cll:math/triangulum_acutangulum
(t3) rdfs:seeAlso cll:math/triangulum_aequilaterum
(t4) rdfs:seeAlso cll:math/triangulum_isosceles
(t5) rdfs:seeAlso cll:math/triangulum_scalenum
(e1) rdfs:seeAlso dbr:Triangle
(e1) foaf:page https://en.wikipedia.org/wiki/Triangle
+++
```

Fig. 2. An annotated fragment of the *Euclidis Elementorum Libri XV. Accessit XVI*. Text is marked with spans highlighted in blue. Identifiers in orange are used in the directive section where they correspond to the subjects of RDF triples.

⁶ <http://lexvo.org/ontology>

⁷ <http://claviusontheweb.it/lexicon/math/>

⁸ <http://www.w3.org/TR/rdf-schema/>

3 Conclusion & Future Works

This article describes an approach for merging the distinct steps of transcription and annotation as a single process. We implemented a tool based on a lightweight syntax language that allows RDF annotations to be performed. We conducted some preliminary tests, which involved 50 students, who were asked to use the prototype, and provide feedback. Future works will consist in developing an improved language with a syntax capable of handling nested as well as overlapping (i.e., not hierarchically nested) annotations. New syntax elements will be introduced: *milestone* elements, for annotating a single location in the text (e.g., a gap), *partition* elements, for the identification of phenomena such as line, page or sentence breaks. Additional syntax will be introduced to provide shortcuts ("syntactic sugar") to the most common annotations. The Turtle syntax⁹ will also be taken into account for the RDF triples specification. Finally, different formats (e.g., turtle, json, csv, xml) will be chosen for exporting annotations according to various data models (e.g., Open Annotation, NLP Interchange Format (NIF)).

References

1. *Euclidis Elementorum Libri XV: Accessit XVI De Solidorum Regularium Comparatione. Omnes perspicuis demonstrationibus, accuratissime scholiis illustrati.*
2. Matteo Abrate, Angelo Mario Del Grosso, Emiliano Giovannetti, Angelica Lo Duca, Damiana Luzzi, Lorenzo Mancini, Andrea Marchetti, Irene Pedretti, and Silvia Piccini. Sharing cultural heritage: the clavus on the web project. In *Language Resources and Evaluation Conference*, 2014.
3. Ben Adida, Mark Birbeck, Shane McCarron, and Steven Pemberton. Rdfa in xhtml: Syntax and processing. *Recommendation, W3C*, 2008.
4. Marco Grassi, Christian Morbidoni, Michele Nucci, Simone Fonda, and Francesco Piazza. Pundit: augmenting web contents with semantics. *Literary and linguistic computing*, 2013.
5. Nancy Ide and Jean Véronis. *Text encoding initiative: Background and contexts.* Springer Science & Business Media, 1995.
6. Ali Khalili, Sören Auer, and Daniel Hladky. The rdfa content editor-from wysiwyg to wysiwym. In *Computer Software and Applications Conference.* IEEE, 2012.
7. Rohit Khare and Tantek Çelik. Microformats: a pragmatic path to the semantic web. In *Proceedings of the 15th international conference on World Wide Web*, 2006.
8. Steven Ruggles, Matthew Sobek, Catherine A Fitch, Patricia Kelly Hall, and Chad Ronnander. *Integrated public use microdata series: Version 2.0.* Historical Census Projects, Department of History, University of Minnesota, 1997.
9. Pontus Stenetorp, Sampo Pyysalo, Goran Topić, Tomoko Ohta, Sophia Ananiadou, and Jun'ichi Tsujii. Brat: a web-based tool for nlp-assisted text annotation. In *Conference of the European Chapter of the Association for Computational Linguistics.* Association for Computational Linguistics, 2012.
10. Tabea Tietz, Jörg Waitelonis, Joscha Jäger, and Harald Sack. Smart media navigator: Visualizing recommendations based on linked data. In *13th International Semantic Web Conference, Industry Track*, 2014.

⁹ <https://www.w3.org/TR/turtle/>