

Agree to disagree: Modelling co-existing scholarly perspectives on literary text

Elli Bleeker, Bram Buitendijk, and Ronald Haentjens Dekker
Research and Development Group, KNAW Humanities Cluster,
Amsterdam, the Netherlands

Abstract

This essay addresses two open challenges in the domain of digital scholarly editing: (1) formally defining the meaning of markup, and (2) allowing the reuse and exchange of textual data through a distributed editorial workflow that allows the editing of texts from multiple, diverging yet co-existing perspectives. We argue that successfully addressing these issues would promote the distribution and exchange of scholarly knowledge, on a technical as well as a theoretical level. The essay introduces ongoing work on a new data model for text called ‘TAG’ (Text-as-Graph) and its reference implementation ‘Alexandria’. The essay outlines how TAG, based on a hypergraph for text, can improve the modeling of complex literary texts, and how Alexandria supports the exchange of markup files in a way that sustains scholarly discourse. We discuss three components of TAG: first, the markup technology stack allows for the formal definition of the meaning of markup (‘markup semantics’); secondly, users can add multiple layers of markup that each represent an alternative perspective on text; and finally the editorial workflow is set up in a git-like distributed version management system. As a result, the TAG model provides for the synthesis of dispersed scholarly practices and the advancement of academic discourse.

Correspondence:

Elli Bleeker, Research and Development Group, KNAW Humanities Cluster, Amsterdam, the Netherlands.

E-mail:

elli.bleeker@gmail.com

1 Introduction

‘Scholars disagree on everything’, a satirical piece in the *New Yorker* claims.¹ Ironically, most scholars would agree with that statement. Opposing views, dialectic discourse and debate are driving forces in research. Humanities research, in particular, requires room for diverging interpretations: cultural and historical artefacts can often be approached from multiple coexisting perspectives. This is one of the reasons why the academic world welcomed the digital medium with great enthusiasm: its infrastructure and theoretically infinite storage space would be able to accommodate various views on

an artefact without privileging one over the other. What is more, it would facilitate the exchange of findings and promote the reuse of one another’s work. In short, digital humanities (DH) was said to promote the organization and dissemination of scholarly knowledge in a way that leaves ample room for multiple co-existing perspectives on data; data that are both reusable and understandable for our scholarly peers. So far, however, that particular academic ideal has not been entirely realized. On the contrary, some have criticized the application of the digital medium to humanities research for reducing research to ‘problem solving’ and for promoting a neo-positivist take on research

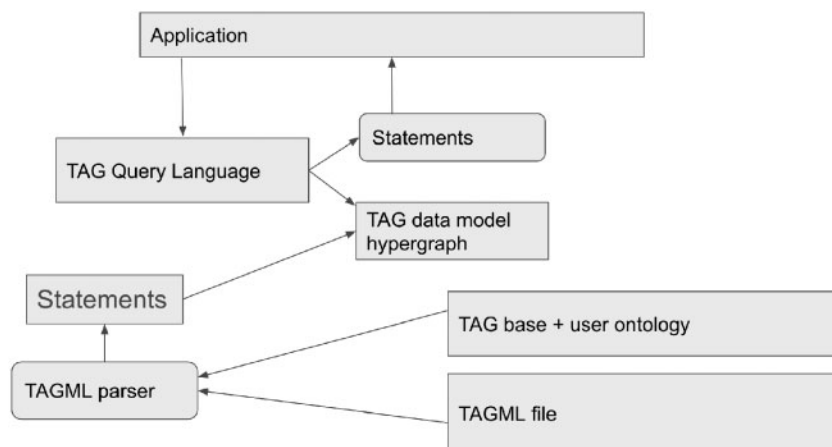


Fig. 1 Schematic representation of the TAGML markup stack

the query language TAG Query Language (TAGQL), applications will be able to query the data in the hypergraph.

By making use of the detailed scholarly knowledge that is stored in the markup and the ontologies (the domain-specific knowledge), applications will be able to carry out more complicated and powerful queries. Furthermore, as others have pointed out, having a formal description of the markup elements used would also make the file's content more accessible. Others may be interested in your transcription and may not share your interpretation of a markup element, but whether you have encoded an added signature with 'fw', 'handNote', or 'add', your interpretation is now described formally and thus readable for humans and machines alike. In short, by complementing a TAGML transcription with a formal description of the markup elements used in that file and a mapping to one or more ontologies, it can be more easily understood, shared, and reused by others.

Another important point regards the human-readability of the markup files. Embedding RDF statements into a transcription results in a verbose and significantly less legible text: a downside that cannot be underestimated in view of sharing and reuse. For that reason, encoding systems based on RDF usually resort to a standoff approach: the RDF statements are stored in an external file and linked to the source text transcription by means of pointers

with identifiers.¹² The source text transcription consists primarily of text, perhaps with some minimal markup, making it easy to read. However, such standoff systems depend on a stable base text: the transcription of the source text is the reference text and cannot change as that will mess up the pointer system. Furthermore, having the markup in a separate file would hinder indexing or searching using the information in the markup. So, whereas standoff approaches have traditionally been commended for their potential to enhance interoperability and interchange (see Schmidt, 2014), we pose that to successfully query, share, and reuse markup files, a better option would be an embedded markup system with (1) the markup directly associated with the source text transcription and (2) the possibility to make changes to the markup as well as the source text transcription. The human-readability of a transcription with multiple layers of embedded markup that most probably overlap remains a challenge, which TAG addresses with its 'layer' functionality.

In TAG, the term 'layer' is used to identify a set of markup nodes that are grouped together for some reason: a layer may express a certain scholarly perspective on text (e.g. a linguistic perspective or a documentary perspective), but a layer can also identify the markup added by a particular user. Layer identifiers can thus be used for both query and display purposes. Furthermore, layers can be used to

