

From DOCX via TEI to Literature Map

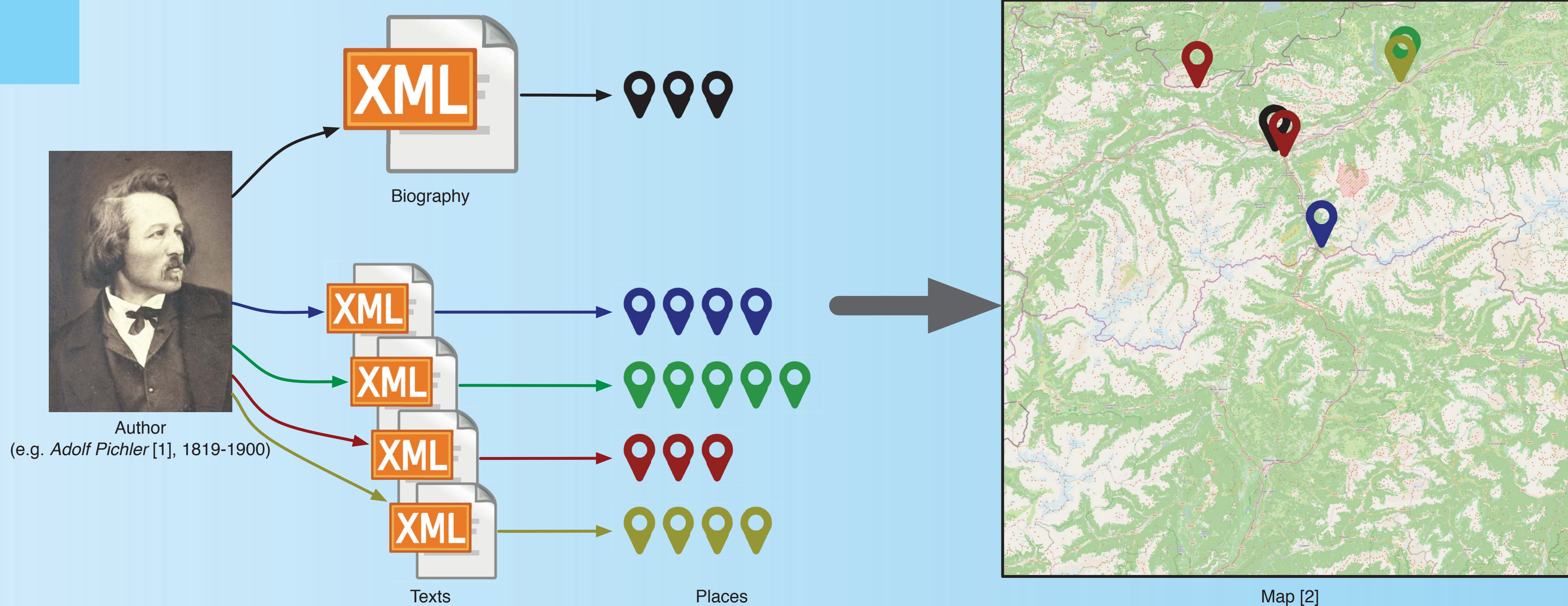
A presentation at the *TEI Conference and Members' Meeting 2016* in Vienna, September 2016, by Joseph Wang, University of Innsbruck

Aim of the Project

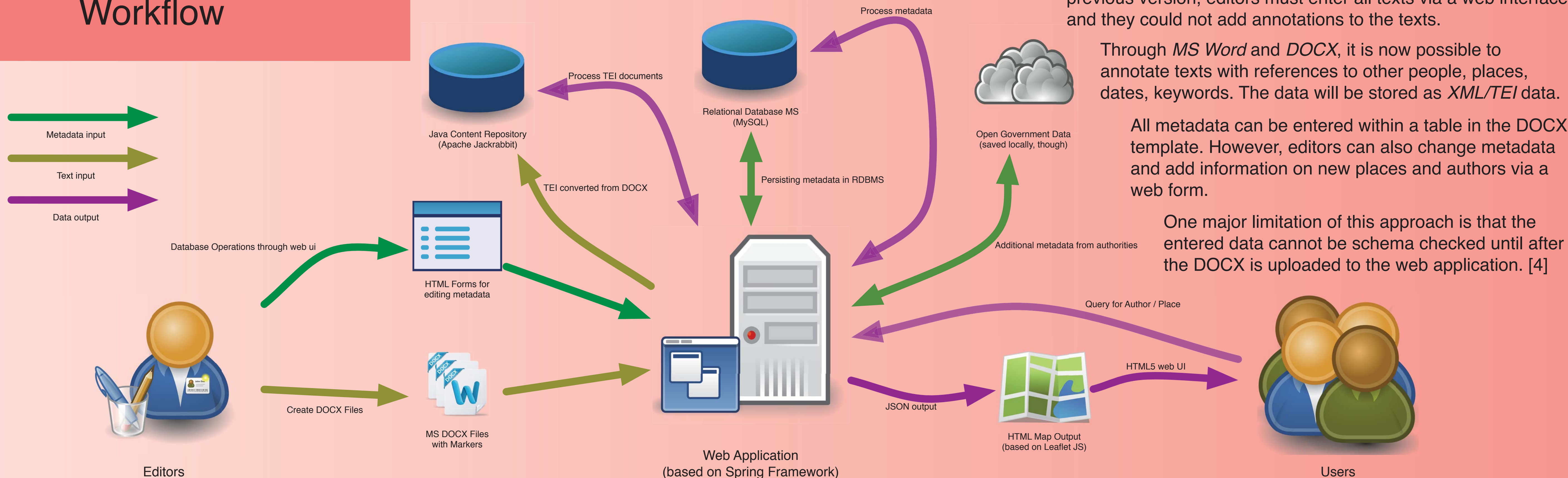
The aim of the project is to create an input system to feed the *Literature Georeferencing Database*, so the software is able to generate Literature Map out of the data.

A *Literature Map Software* must be able to:

1. create a geographical map given an author showing georeferences in his life and in his works.
2. create a geographical map given a work showing georeferences in this work
3. create a geographical map given a geo-location showing other georeferences linked to this location through literary works [3]



Workflow



The workflow was designed to ease the work of editors. In the previous version, editors must enter all texts via a web interface and they could not add annotations to the texts.

Through *MS Word* and *DOCX*, it is now possible to annotate texts with references to other people, places, dates, keywords. The data will be stored as *XML/TEI* data.

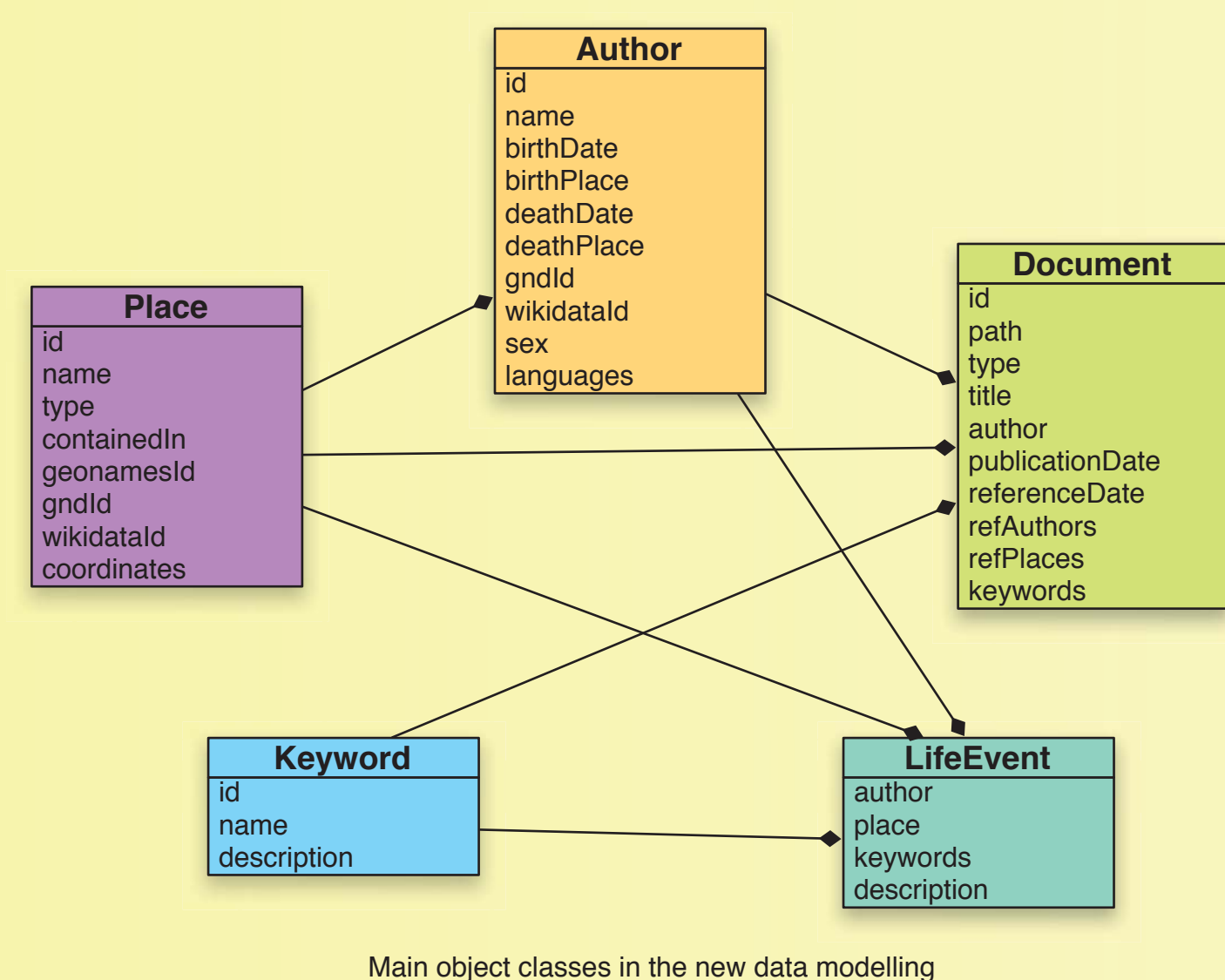
All metadata can be entered within a table in the *DOCX* template. However, editors can also change metadata and add information on new places and authors via a web form.

One major limitation of this approach is that the entered data cannot be schema checked until after the *DOCX* is uploaded to the web application. [4]

Rational behind the data modelling:

1. The researchers should be able to enter large amount of data quickly.
2. We apply DRY-Principle (*Don't Repeat Yourself*) as good as possible.
3. Maps are generated from the database, and not drawn by the scholars using GIS tools.

Thus, we come up with this solution: For the metadata, a web frontend for database operations is created; for the TEI data, we found a solution based upon *DOCX* to *TEI* Conversion. Furthermore, queries on the authority givers (*GND*, *Geonames.org*, and *Wikidata*) is made during the data ingest, thus the database entities are linked to other authorities. However, the same cannot be done with open government data made available by the local government of Tyrol and South Tyrol; the reason: these data can only be downloaded but not queried through a web service. Therefore, the data is stored locally and queried during the data ingest.



Map Generation:

Places have coordinates, they can be plotted in a HTML map. Using *LeafletJS*, one can embed interactive maps within a web page. Different colors and different icons mark different meaning of a place.

Benefits of this modelling:

1. Captures many aspects of use cases (e.g.: queries on objects are just queries on databases.)
2. Database operations on metadata of objects can be done quickly.
3. Inserting and Updating texts can be done easily.
4. Data models reflect their usage in humanities research.

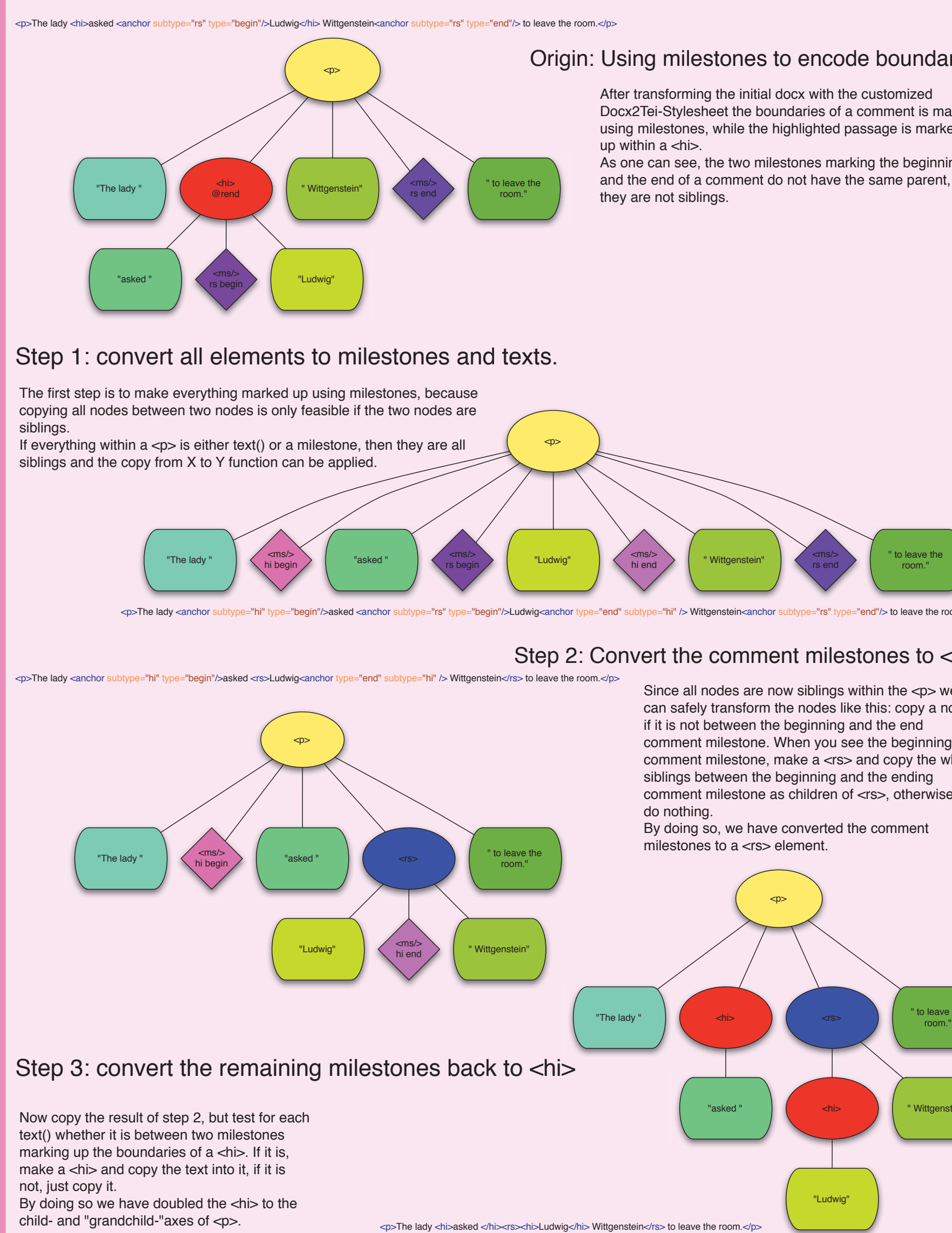
Drawbacks of this modelling:

1. Complicated database scheme is difficult to maintain.
2. Redundant data storage, no back-propagation of changes in RDBMS to *TEI* Document.
3. Dependencies on data of other resources.

Data Modelling

Convert comments marked using milestones to <rs> in 3 steps

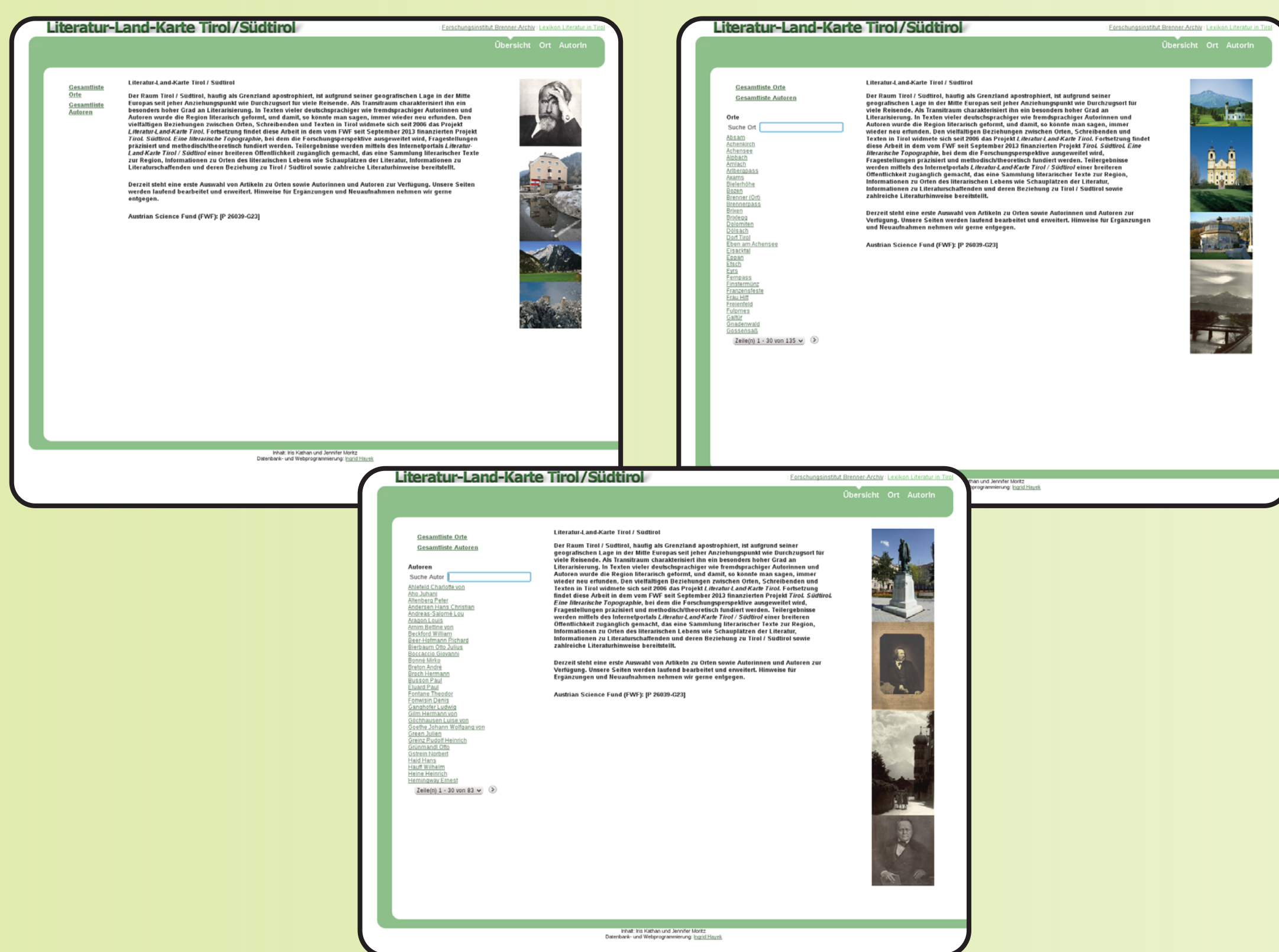
A simplified example without attributes



Having the relationships between authors, works and geographical places in focus, the Project *The Tyrol / The South Tyrol - A literary Topography* (FWF P26039) created a database on these entities.

At the beginning, data modelling only captures basic metadata of the entities and their relationships. However, it was impossible to annotate texts, thus scholars could not tag the actual section of the texts with additional references.

After careful consideration, it was decided to redo the data modelling: While the metadata should still be kept and maintained by a RDBMS (in our case: *MySQL*), the actual texts (primary sources, biographies and description of places) should be modelled as *XML/TEI* data and kept in a repository (in our case: *Apache Jackrabbit*). Furthermore, we want to reuse authority controlled data and open government data.



Screenshot from the Database *Literatur-Land-Karte Tirol/Südtirol* (access on 2016/09/15)

Starting Point