# POSTDATA:
# Towards publishing European Poetry as Linked Open Data

**Mariana Curado Malta (mariana@iscap.ipp.pt; mariana.malta@linhd.uned.es),** ISCAP, IPP | LINHD – UNED

**Paloma Centenera (paloma.centenera@linhd.uned.es), LINHD - UNED**

**Elena González-Blanco (egonzalezblanco@flog.uned.es), LINHD - UNED**

# Outline

- The Context

- The Problem

- The Approach

- Where are we now?

- Future Work

- Conclusions

# Context

Project **POSTDATA**
Poetry Standardization and Linked Open Data

ERC Grant
PI: Elena Gonzalez-Blanco

# Context – Where?



*A digital humanities center is an entity where new media and technologies are used for humanities-based research, teaching, and intellectual engagement and experimentation. The goals of the center are to further humanities scholarship, create new forms of knowledge, and explore technology's impact on humanities based disciplines".*

*Diane M. Zorich, A Survey of Digital Humanities Centers in the United States, 2008*

# Context – Where?

- LINHD: a bridge between different fields of knowledge

- LINHD has:

    ➔ Philologists
    ➔ Software Developers
    ➔ Natural Language Processing Experts
    ➔ Ontologists & LOD technologists

# DH in the World

# Context - What

- The metrics on European Poetry

Estuans intrinsecus     ira vehementi
in amaritudine     loquor mee menti.
factus de materia     levis elementi
4    folio sum similis,     de quo ludunt venti.

Cum sit enim proprium     viro sapienti,
supra petram ponere     sedem fondamenti,
stultus ego comparor     fluvio labenti,
8    sub eodem aere     numquam permanenti.

Feror ego veluti     sine nauta navis,
ut per vias aeris     vaga fertur avis;
non me tenent vincula,     non me tenet clavis,
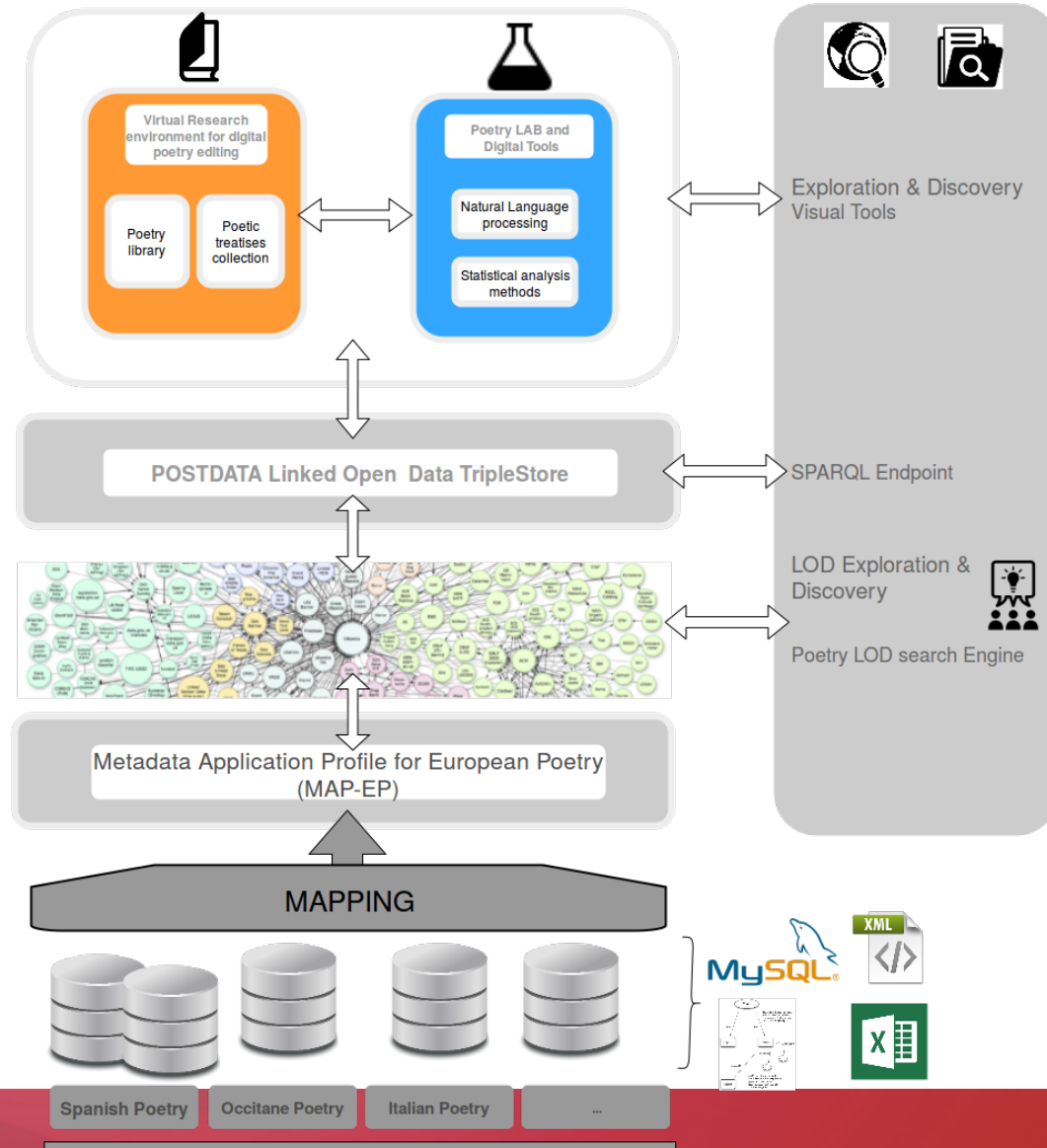12    quero mei similes     et adiungor pravis.
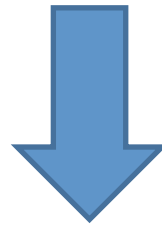
*Carmina Burana, 191*

# The problem

- At least 21 repertoires on Poetry metrics & other information (in the Web of Documents)

- This community wants to share all  the data among  repertoires

- ....to enhance  its  research

- And  more...

# The problem

- First issue: standardize poetic features

  ➔ Different languages

  ➔ Different cultures/traditions

# Philologists take care of this issue!

# Philological barriers: different ways of conceptualization

Estuans intrinsecus     ira vehementi

in amaritudine     loquor mee menti.

factus de materia     levis elementi

4     folio sum similis,     de quo ludunt venti.

Cum sit enim proprium     viro sapienti,

supra petram ponere     sedem fondamenti,

stultus ego comparor     fluvio labenti,

8     sub eodem aere     numquam permanenti.

Feror ego veluti     sine nauta navis,

ut per vias aeris     vaga fertur avis;

non me tenent vincula,     non me tenet clavis,

12     quero mei similes     et adiungor pravis.

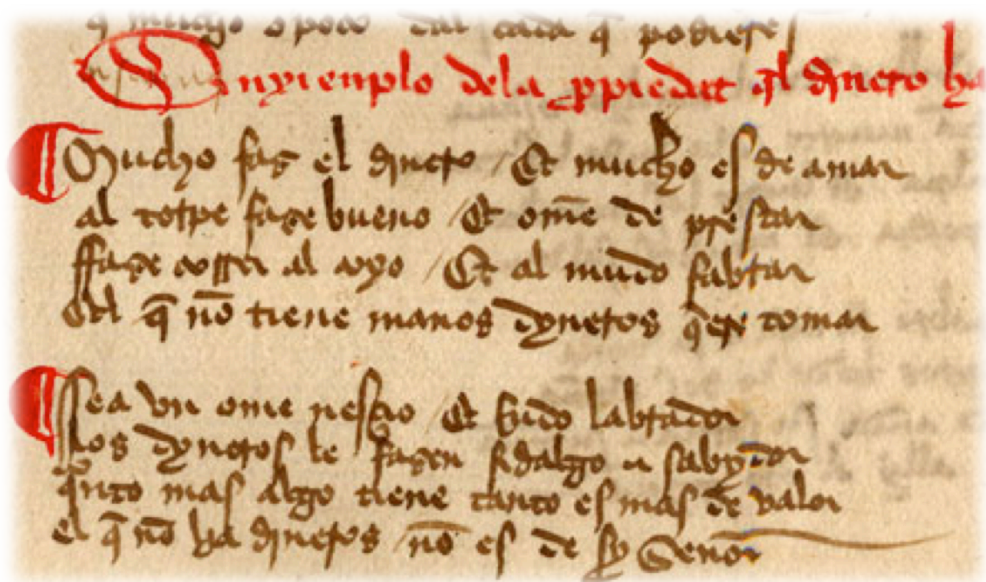*Carmina Burana, 191*

Alexandrines

Goliardic

12A12A12A12A
(Romance)

4x(7pp+7p)
(Classic Latin)

# Philological standardization: starting point



Author
Title
Incipit
Manuscript
Post quem
Ante quem
Language
Topics
Edition
Online edition
Work

Isometrism
Isostrophism
Metrical scheme
Rhyme scheme
Rhyme
Musical notation
Number of stanzas
Number of lines
Poetic form

# The problem

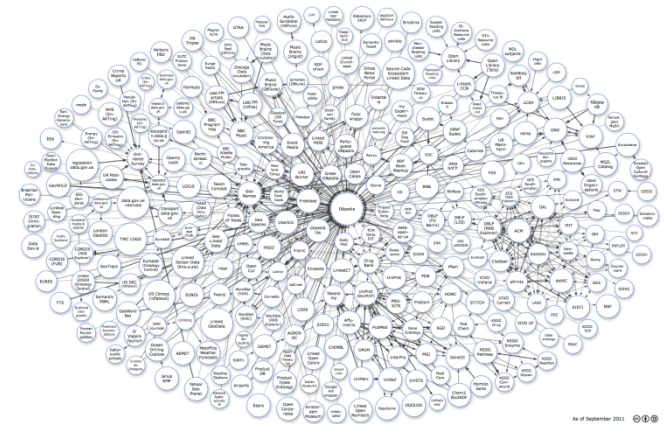Second issue: repertoires locked in their silos of information:

➔ Different paradigms: Local and Web of Docs

➔ Different technologies: XML, Excel, Access, MySQL, SQL, Data stored in Perl Objects *(so far)*
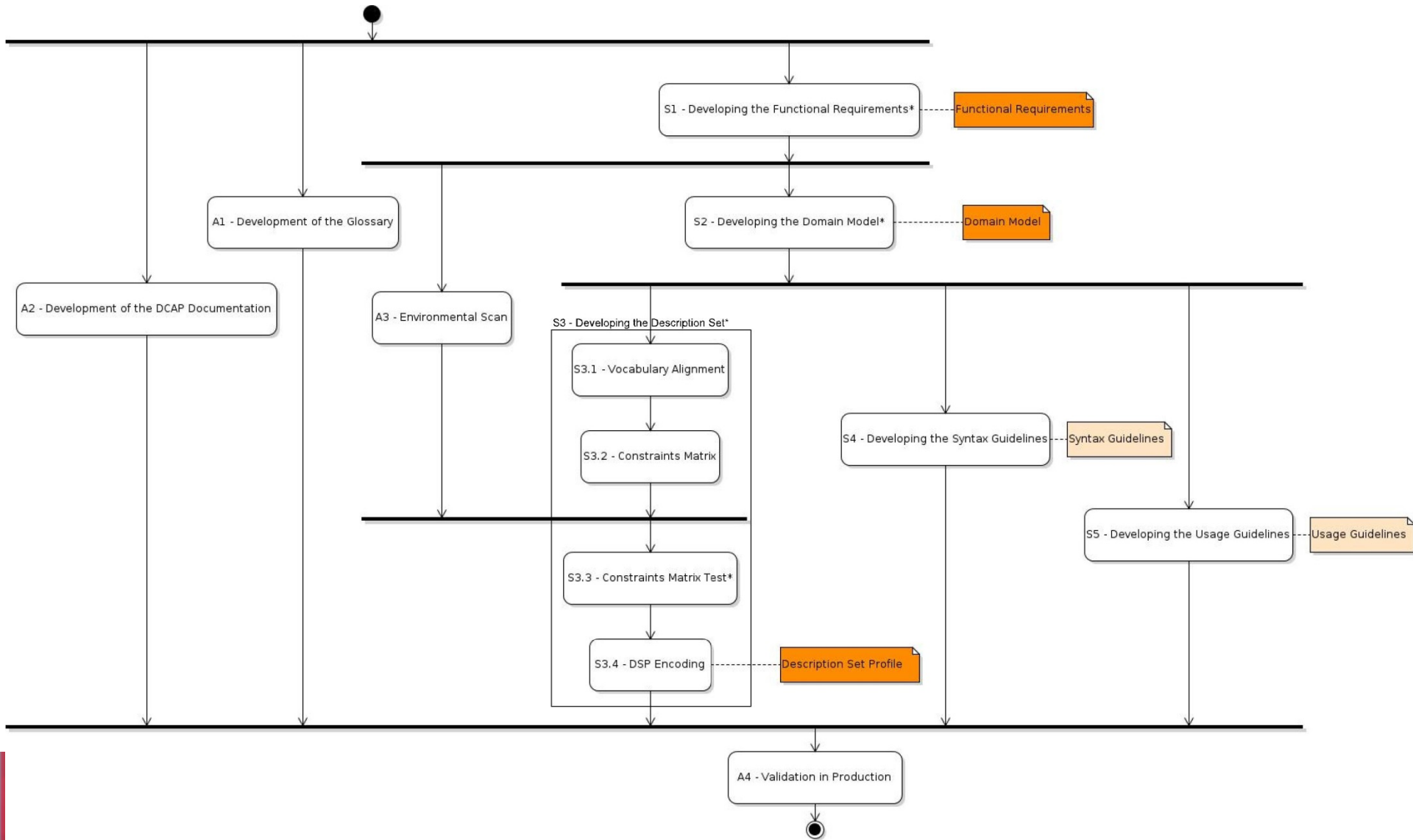
➔ Different data models

# The problem

- How to overcome these diferences?

- LOD technology

- Development of a Metadata **Application Profile** for the European Poetry community

# The Approach

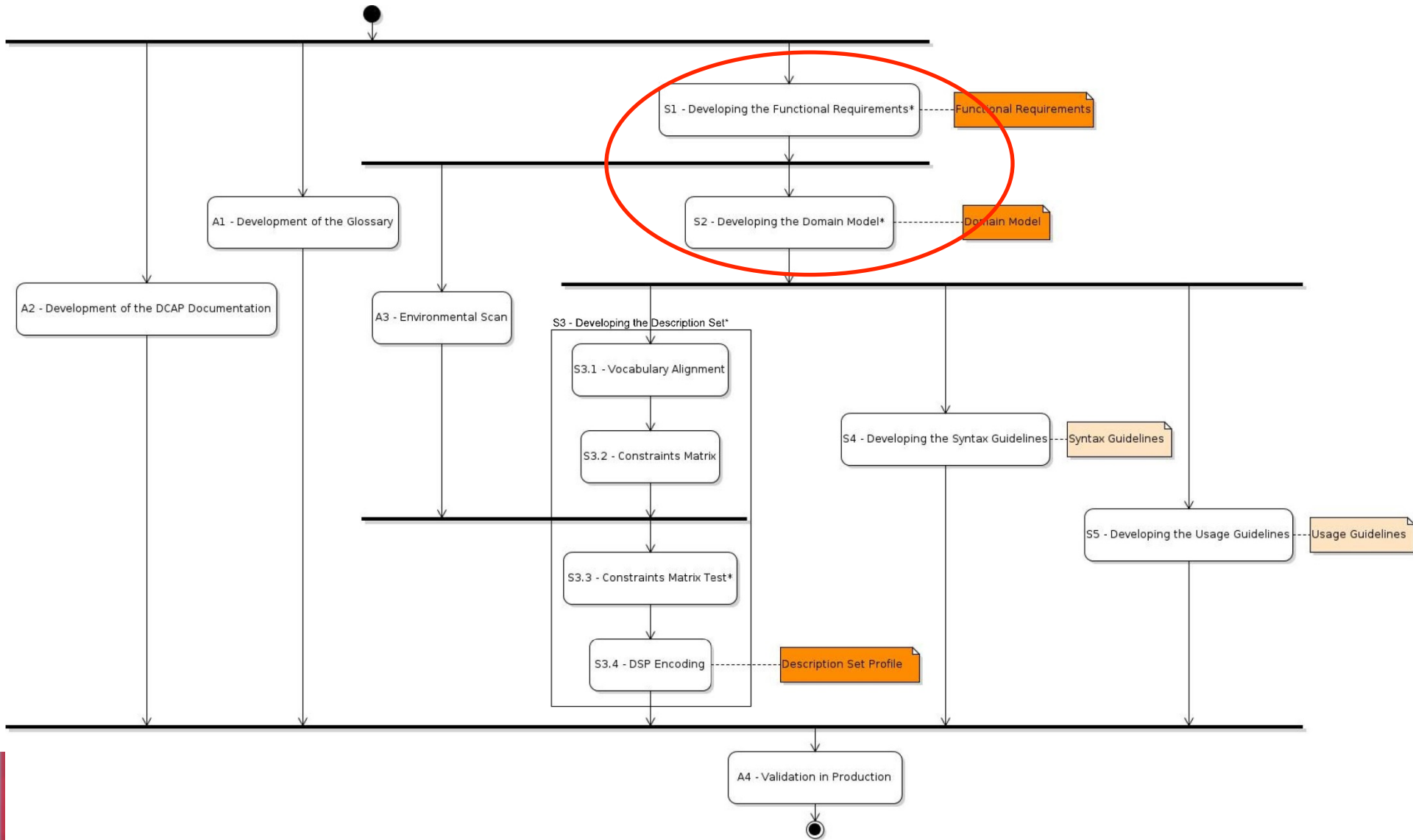- Method for the development of Metadata Application Profiles (Me4MAP)

- Me4MAP establishes a well defined process for the development of a MAP:

  ➢ defines activities
  ➢ when those activities should take place
  ➢ how they interconnect
  ➢ and their resulting deliverables

# The Approach



S1 - Developing the Functional Requirements* ----- Functional Requirements

A1 - Development of the Glossary

S2 - Developing the Domain Model* ----- Domain Model

A2 - Development of the DCAP Documentation

A3 - Environmental Scan

S3 - Developing the Description Set*

S3.1 - Vocabulary Alignment

S3.2 - Constraints Matrix

S4 - Developing the Syntax Guidelines --- Syntax Guidelines

S5 - Developing the Usage Guidelines --- Usage Guidelines

S3.3 - Constraints Matrix Test*

S3.4 - DSP Encoding ----- Description Set Profile

A4 - Validation in Production

# Were are we now?

# Where are we now ?

- S1: Defining the Functional Requirements

  - Analysing the Websites' functionalities and the Logical Models of the databases (when possible)

- S2: Defining the Domain Model

# Where are we now ?

- Reverse engineering process eliminates all the details that have to do with the implementation/ representation

- We have followed the process:
  - ID keys deleted

  - Separate different concepts that are represented in the same table

  - Tables that enumerate terms deleted → become properties that can be repeated

  - When models have conceptual problems → fix problems

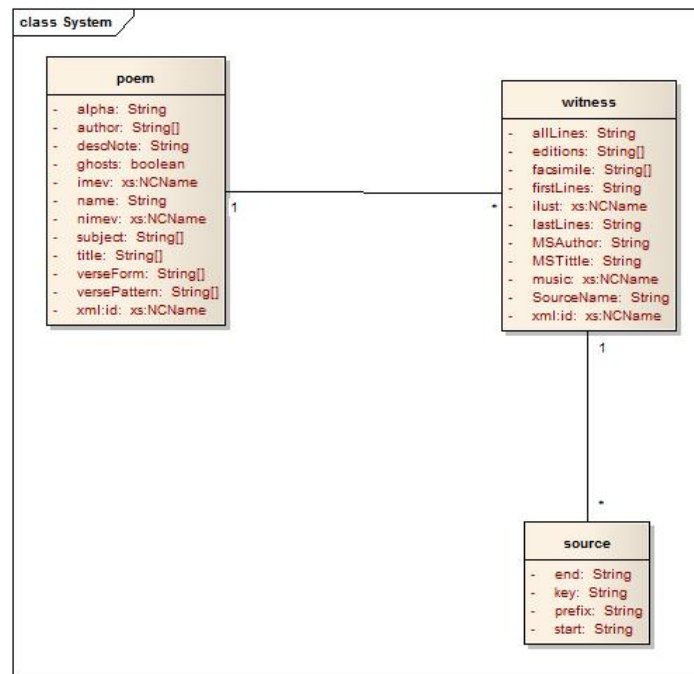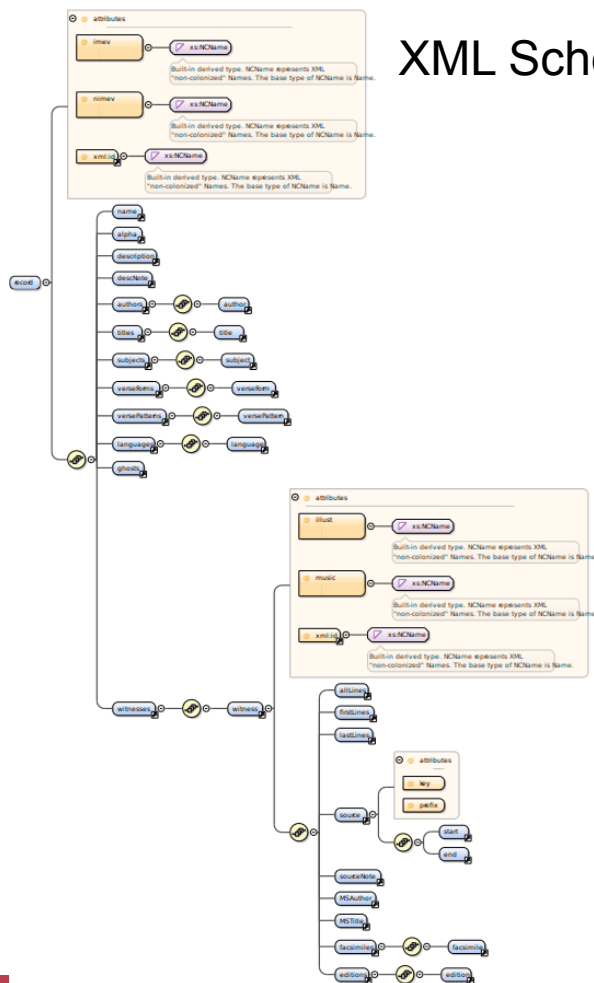# Analysing Data Models

Relational Database ➡️ Conceptual Model

# Analysing Data Models

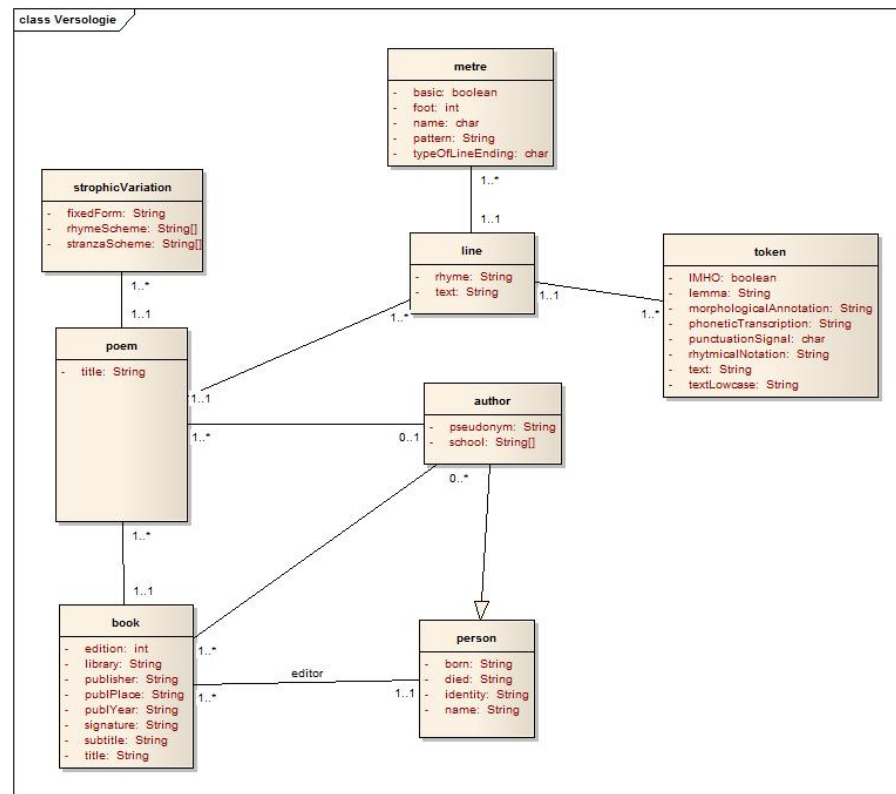XML Schema Model ➡ Conceptual Model

# Analysing Data Models

Perl Scritp structure → Conceptual Model

# Analysing Data Models

- During the process of reverse engineering we standardize, i.e.

  ➢ Call the same concepts by the same name (working together with the philologist)

  ➢ Try to call the same names to tables or properties as classes or terms that already exist in A3: Environmental Scan

# Where are we now ?

- If the repertoire's responsibles did not provide database definition →analyse the functionalities of the Website


- Study the controlled vocabularies and standardize them: 1) ask for them, 2) collect them in the Websites

# Where are we now ?

- Study other communities/projects for interoperability, ex:

  ➢ Biblioteca Nacional de España - http://www.bne.es

  ➢ Biblissima - http://www.biblissima-condorcet.fr

  ➢ Pelagios - http://pelagios.org/

  ➢ Claros - http://www.clarosnet.org/

# Where are we now ?

- A3: Environmental Scan:

  - A report

  - contains a review schemas available in any serialization of the Semantic Web

  - that may serve the needs of the Domain Model

# Future Work

- At the end of S1 & S2 processes we will have Functional requirements & Domain Model defined

- We will validate the Domain Model in two meetings (Jan/Feb 2017) with:
  1. The repertoire's responsibles (circa 20)
  2. Semantic Modelers experts (3)

# Next step ?

- Defined the Domain Model & the Environmental Scan → develop the S2.1: Vocabulary Alignment

  ➢ To match the terms of the metadata schemas identified in the Environmental Scan (A3) with the needs of the Domain Model.

# Conclusions

- POSTDATA aims to put poetry metrics data in LOD

- There are at least 21 repertoires on the Web of Documents

- To achieve that we need to: 1) standardize the way poetry metrics is defined, 2) create a Metadata Application Profile (MAP) for the European Poetry community

- We are following Me4MAP to create this MAP

# Europeana Poetry?

**Mariana  Curado  Malta**
mariana.malta@linhd.uned.es
**Paloma  Centenera**
paloma.centendera@linhd.uned.es
**Elena  González-Blanco**
gonzalezblanco@flog.uned.es

linhd.uned.es
postdata.linhd.es
@linhduned