



# Datathon Miniproject

## Linking XML corpora

Tutors: Maxim Ionov, Christian Chiarcos

Participants: Sabine Tittel, ...



# Corpora and edited works in the Humanities

- TEI/XML as a de facto standard
- Linked Data in Digital Philology
  - prosopography ([social] network analysis)
  - managing cross references
- Little to no application of LOD for resource modeling
  - the lack of convincing prototypes
  - *in this field*, XML is likely to stay

# TEI and LOD

- the relevance of LOD has also been recognized in the TEI community
  - yet, there is no convenient formalization of semantically typed URI references in the TEI
    - beyond special-purpose applications (prosopography)
- here, experiment with using RDFa for the purpose
  - preserve TEI spirit and formalization
  - add LOD

# Sample Data: Edition (printed book)

[14v<sup>o</sup>a] Ou nom de Dieu misericord. Cy commence le premier traictier de ceste oevre qui parle de l'anathomie et contient deux doctrines: la premiere doctrine parle de l'anathomie des membres *communs*, universelz et simples. La seconde sera des membres propres, particuliers [14v<sup>o</sup>b] et compost. La premiere doctrine  
5 ne contient .v. chappitres: le premier, c'est ung chappitre universel qui parle de l'anathomie et de la nature des membres du corps.

[14v<sup>o</sup>a] POUR CE QUE, selon Galien, lumiere des mediciens, ou .xvij.<sup>e</sup> [14v<sup>o</sup>b] livre qui se intitule «De utilité des parties», ou penultime chappitre, y sont quatre utilités de la science de anathomie: l'une, qui est la tres grande, pour amiracion  
10 de la puissance de Dieu; la seconde, pour cognoistre les parties des paciens; la tierce, pour pronostiquer des [15r<sup>o</sup>a] dispositions du corps qui doivent avenir; la quarte si est pour curer les maladies. Et pour ce, c'est chose necessaire et prouffitable a ung *chacum* medecin de savoir la anathomie. Et c'est ce *que* disoit Galien ou livre qui se intitule «*Liber scienciarum sive interiorum medicorum*», ou il dit  
15 ainsi: les jeunes clers, et les anciens aussi, estudient a cognoistre les parties et les passions d'icelles, car selon la difference d'icelles. Et ja soit ce que les parties qui appparent ou sens soient cogneues appertement, toutesvoies, celles qui sont en parfont occultes, elles ont mestier de homme qui soit excercités en l'anathomie et es accions et utilités d'icelles. Et de ce lieu ci est prins le principes de tout le  
20 continent. Et dit qu'il est escript ou premier du livre des membres que le medecin hardi doit estre sage en la cognoissance des membres qui viennent en *chacum* lieu. Et se c'est chose prouffitable aux phisiciens, elle est plus necessaire aux

# Sample Data: Edition (LaTeX source)

```
\ProvidesFile{EDITION.tex}[2017/06/13]

\begin{edition}{Kritischer Text}
\markboth{Die \emph{Anathomie}}{Kritischer Text}

\beginnumbering
\pstart
%
[14v\hoch{o}a] Ou nom\wdx{nom}{m. `mot
servant à désigner les êtres, les choses
qui appartiennent à une même catégorie
logique'}{\textbf{au nom de} \emph{`en vertu de'}}
de Dieu misericord\wdx{misericort}{adj. `qui a de la
miséricorde; miséricordieux'}{misericord
\emph{m.sg.}}. Cy
co\emph{m}mence\wdx{*comencier}{v.intr. `entrer
dans son commencement'}{commence
\emph{3.p.sg.
ind.prés.}} le premier\wdx{premier}{adj. `qui vient
avant les autres, dans un ordre; premier'}{}
traictier\wdx{*traitié}{m. `ouvrage
didactique, où est exposé d'une manière
systématique un sujet ou un ensemble de
```

# Sample Data: Edition

- (snippet of the) edition of the Middle French *Grande Chirurgie* by *Gui de Chauliac* (15th c.)
  - text with glossary (originally typeset in LaTeX)

1060 ryr et engendrer le corps. Les non naturelles sont devisees: en sont envoyés ens certains [32r<sup>o</sup>b] lieux pour aucuns aides, ou elles sont expellees du corps. Et saches **que la cole est envoyee ou cestim du fiel**, et la merancolie a l'esplein, et le fleume aux jointures, et la superfluité aigouse aux reins et a la vecie. Celles qui sont expellees du corps, elles vont avec le sang et se putrefient aucunes fois et

edited text

**cestim du fiel** m. terme d'anat. 'vésicule biliaire' 1061

Isolierte Form. Cp. AmphYpL<sup>2</sup> f<sup>o</sup>150a *la ciste ou bourse du fiel*. Cp. DEAF H 676,40ss. *huche du fiel* mit gleicher Bedeutung und den Kommentar ib. 673,42–674,3. Cp. *cistis fellis*.

sample gloss (\index)

# Text edition (Mid. French, 15th c.)

```
Et
saches
que
la
<wdx>
  <orth>cole</orth>
  <lemma>cole</lemma>
  <gloss>f. terme de méd. `bile jaune (l'une des quat
</wdx>
est
envoyee
ou
cestim
du
<wdx>
  <orth>fiel</orth>
  <lemma>cestim du fiel</lemma>
  <gloss>m. terme d'anat. `vésicule biliaire'</gloss>
</wdx>
,
et
la
<wdx>
  <orth>merancolie</orth>
  <lemma>melancolie</lemma>
  <gloss>f. terme de méd. `bile noire (l'une des quat
  <var>merancolie</var>
</wdx>
a
<wdx>
  <orth>l'esplein</orth>
  <lemma>esplein</lemma>
  <gloss>m. terme d'anat. `organe lymphoïde situé sou
</wdx>
```

- philological edition
  - primary text
  - plus in-line glossary
- XML
  - converted from LaTeX sources
  - one-word-per-line formatting for viewing convenience only

# Sample Data: Dictionary

- <http://deaf-server.adw.uni-heidelberg.de/lemme/fiel>

**FIEL** m.



[Étymologie]



(*fiel*, *fel* ca. 1000, *feel*, *fele*, *feil*, *feil*, *fius*)

1° t. de méd. "liquide verdâtre et amer qui est contenu dans la vésicule biliaire, bile" (dep. ca. 1160, EneasS<sup>2</sup> 8221 [*el cors m'as mis une amartume Peor que sui*]; MoamT II 48,17 [*burse dou fiel dedens, et sont .6.: l'estomach, les b*]; l'esplain, la ciste ou bourse du fiel]; Hu 2,160a [CERTEAU]; Li 1,1666b; D

◆ dans des collocations *huche de* foie et qui emmagasine la bile, vésicule [*fel*]; MoamT II 48,17 [*burse dou fiel dedens, et sont .6.: l'estomach, les b*]; l'esplain, la ciste ou bourse du fiel]; Hu 2,160a [CERTEAU]; Li 1,1666b; D

```

7 <part type="mainpart">
8   <title>
9     <lemma developed="false" language="afr.">fiel</lemma>
10    <pos>m.</pos>
11  </title>
12  <etymology>
189 <variants>
1062 <senses>
1063   <sense>
1064     <description>&#x2060;<m:terminology type="medecine">t. de
1065       m&#xE9;d.</m:terminology>
1066     <m:definition>liquide verd&#xE2;tre et amer qui est contenu dans la
1067       v&#xE9;sicule biliaire,
1068       bile</m:definition>&#x2060;</description>
1069   <datings>
1070     <dating type="since">
1071       <date numeric="1160">ca.1160</date>
1072     </dating>
1073   </datings>
1074   <references type="primary">
1075     <reference id="10698">
1076       <m:date numeric="1160">ca.1160</m:date>
1077       <m:siglum>EneasS<m:sup>2</m:sup></m:siglum>
1078       <m:text-reference>8221</m:text-reference>
1079     <context>
1080       <m:quotation>el cors m&#x2019;as mis une amartume Peor que sui
1081         ne que fiel</m:quotation>
1082     </context>
1083   </reference>
1084 </senses>

```



# Miniproject

Data:

- edition
  - text
  - glosses
- dictionary

Ultimate goal:

- dictionary and glossary as lemon
- edited text as TEI/XML
- cross-references between edited text (XML) and lemon dictionary

Here:

- model for referencing lemon data from the edited text
- RDFa+XML

for this and similar datasets  
provided by participants