

The Problem of Scale Handling (More) Data with the *WikiLexicographica*



Bruno Bon
Institut de Recherche et
d'Histoire des Textes, CNRS
Paris, France

Krzysztof Nowak
Institute of Polish Language
Polish Academy of Sciences
Kraków, Poland



(Pre)History of the Project (1)



ca. 1920

Union Académique Internationale launches an initiative aiming at replacing Du Cange's *Glossarium* with a new scholarly dictionary – *Novum Glossarium Mediae Latinitatis*

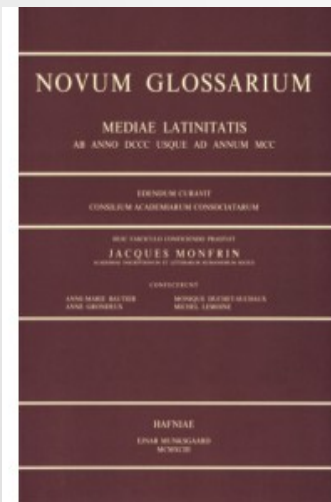
soon after

several national editorial committees decide to create separate national Medieval Latin dictionaries

21st century

several dictionaries varying in their

- geographical and chronological scope
- work progress *etc.*



(Pre)History of the Project (2)

2011-2015

Cost Action *Medioevo Europeo*



- main goal

VCMS = „Virtual Center of Medieval Studies”

- WG 3 goal

unified access to the corpora and dictionaries

2013 -

Prototype of the *WikiLexicographica*, a tool of integrated access to the dispersed lexicographic content



Outline

- ▶ *WikiLexicographica* – prototype and data import process
- ▶ macrostructural organisation
- ▶ issues with data quantity & quality
- ▶ gardening the *WikiLexicographica*
 - ▶ templates
 - ▶ redirections
 - ▶ concept pages

WikiLexicographica. Software



this word's headword is
mandragora



this word's headword is
[[headword: :mandragora]]

WikiLexicographica. Advantages (1)

- separation of basic and advanced view

EU:Mandragora

Summarium **Articulus plenus** Alia

Lemma
mandragora

Formae
mandragoras

Grammatica

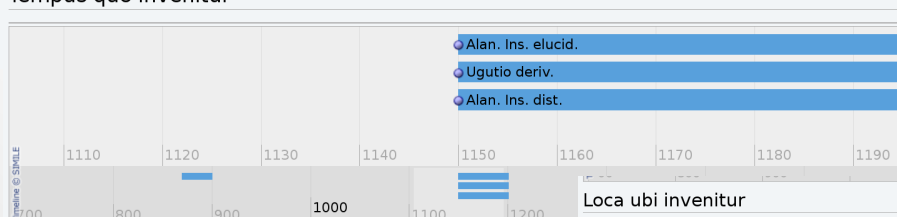
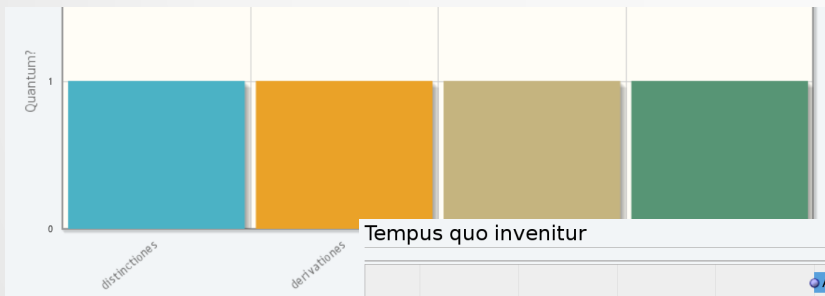
Paradigma
-e

Summarium **Articulus plenus** Alia

mandragora
, -e f. *sive* mandragoras, -e f.

- **Ugutio deriv.** : -a, genus pomi, simili: *mandragore, plante* :
- **Anon. de navig. et agric.** (MGH, Poe
- **Alan. Ins. dist.** col. 848^B : -a, proprie
- **Alan. Ins. elucid.** col. 103^A : per -as, l intelligitur.

- lexicographic data visualisation

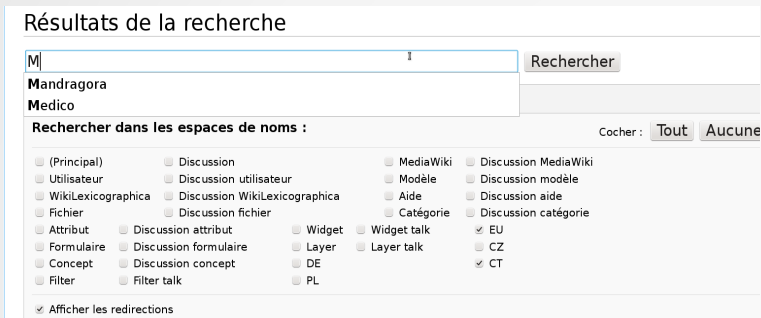


Loca ubi invenitur

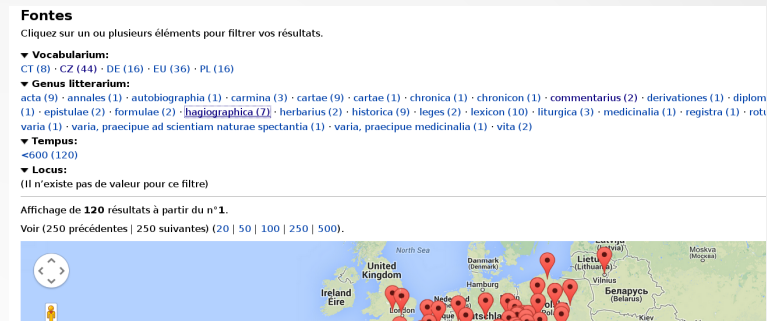


WikiLexicographica. Advantages (2)

- efficient and varied data access



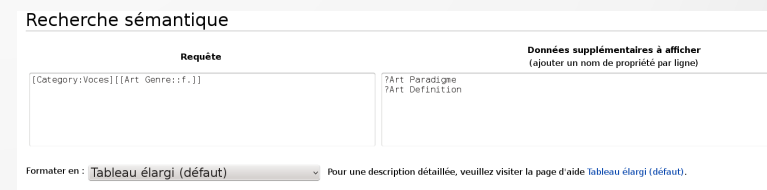
simple search



Semantic Drilldown



Factbox browsing



advanced search

- collaboration support

Importing (More) Data

imported data

NGML

NGML sources

LMILP

LMILP sources

Latin Wiktionary

DuCange

import tools

Import pages

Upload XML data

Please export the file from the source wiki using the [export utility](#). Save it to your computer and upload it here.

Filename: Nie wybrano pliku.

Comment:

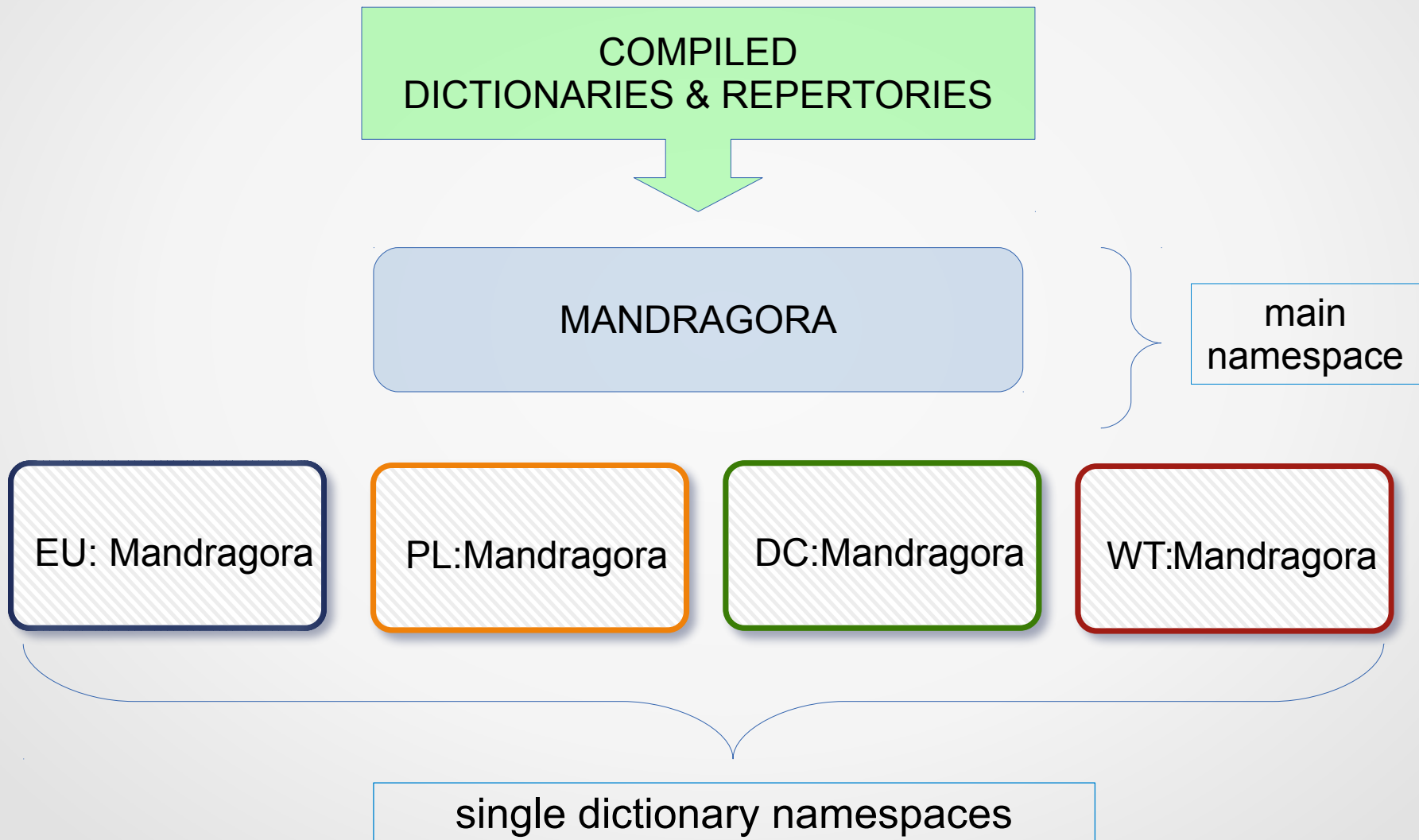
Destination root page (optional):

```
1775 php importDump.php --conf ../LocalSettings.php < ../../work/WikiLexicogra  
phica/export/A-Z export.xml  
1779 php importDump.php --conf ../LocalSettings.php < ../../work/WikiLexicogra  
phica/export/A-Z/A-Z export.xml  
1785 php importDump.php --conf ../LocalSettings.php < ../../work/WikiLexicogra  
phica/export/A-Z/A-Z_main export.xml  
1786 php importDump.php --conf ../LocalSettings.php < ../../work/WikiLexicogra  
phica/export/A-Z/A-Z main export.xml -v
```

import format

```
1 <?xml version="1.0" encoding="UTF-8"?>  
2 <mediawiki version="0.8" xml:lang="">  
3 <page>  
4 <title>a</title>  
5 <ns>0</ns>  
6 <revision>  
7 <text>{{Art|{{Art lemma|a|silent}}|{{Art form|a|silent}}|{{Art form|ab|silent}}|{{Art  
form|abs|silent}}|{{Art inv|a|silent}}|{{Art inv|ba|silent}}|{{Art inv|sba|silent}}|{{Art  
iType|-|silent}}|{{Art iType|-|silent}}|{{Art iType|-|silent}}|{{Art pos|PRE|silent}}|{{Art def|( + abl.) à  
partir de, loin de, du côté de, par suite de, par|silent}}|silent}}</text>  
8 </revision>  
9 </page>  
10 <page>  
11 <title>abacinus</title>  
12 <ns>0</ns>
```

Macrostructure. Namespaces



Data Quantity Issues. Timelines

too much data

The screenshot shows a web interface with a navigation menu on the left and a main content area. The main content area has four tabs: "Fontes enumerati", "Praecipue citati", "Geographice ordinati", and "Chronologica ordinati". The "Chronologica ordinati" tab is selected, displaying a list of data items represented by blue horizontal bars. A tooltip is visible over one of the items, containing the following text:

EU:Abbo Flor. circ. decenn.
De: 945
Ad: 1004
Fri, 01 Jan 0945 00:00:00 GMT
Mon, 31 Dec 1004 00:00:00
GMT

Below the tooltip, the following text is visible in the list:

- EU:Abbo Flor. circ. mund.
- EU:Abbo Flor. circ. decenn.
- EU:Abbo Flor. can.
- EU:Abbo Flor. apol.
- EU:Abbo Flor. Rom. pont.
- EU:Abbo Flor. Eadm.
- EU:Abbo Flor. (?) carm.

limiting data display

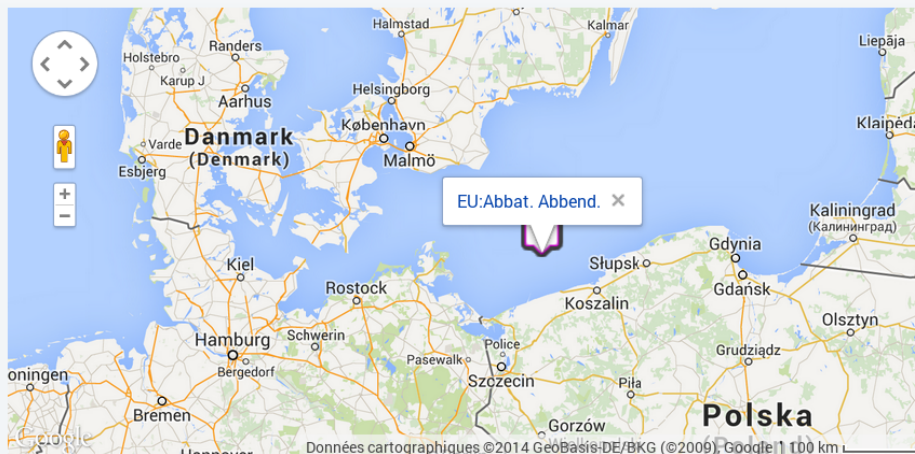
Data Quality Issues. Maps (1)

low granularity of the geographical data

eliminating uncertain data?

Catégorie:Fontes

Fontes enumerati Praecipue citati Geographice ordinati Chronologice ordinati



Kraków, Poland, Europe



Poland, Europe

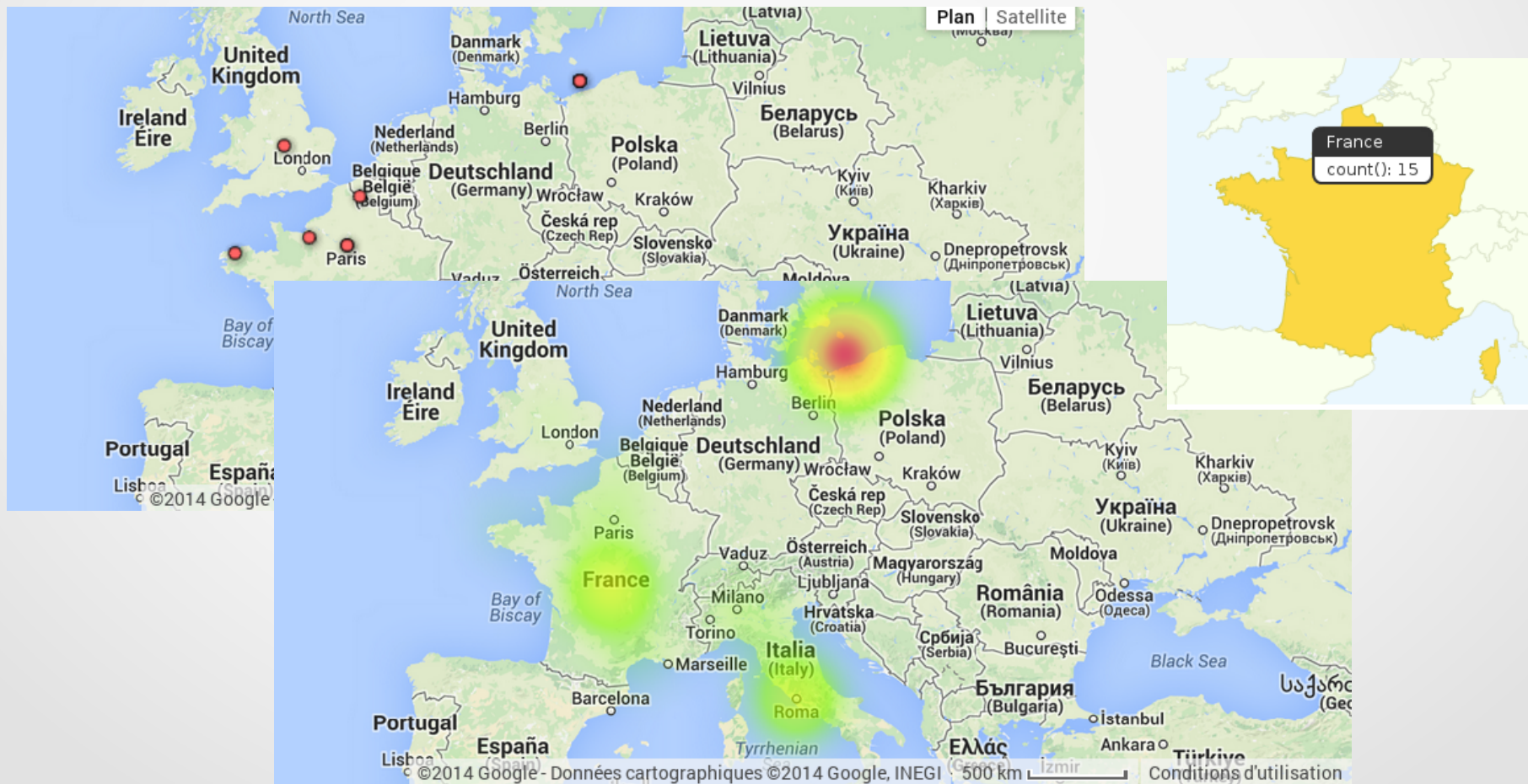


Europe



Data Quality Issues. Maps (2)

implementing static choropleths with Google Fusion Tables



Microstructural Problems

poor formatting

difficult to reflect nested structures
(senses, notes *etc.*)

data extraction approach

different lexicographic conventions, differing encoding

copyright fears

Gardening the WikiLexicographica. Templates

automatisation of repetitive tasks

```
{{#set:Lemma= {{#regex:
{{{1}}}|/([0-9.]+\s*)+/}}
}}
```

1. pinna



[[Lemma::Pinna]]

cuts **homonym number** from the headword's form and sets
[[Lemma]] property

adding redundant annotation

```
{{#set: POS= {{#if: ... [[iType::~~1*]] ... VBE ... }} }}
```

pinno ... 1. → [[POS::VBE]]

checks if word's **inflectional type** begins with **number**;
if this is the case, sets [[POS]] property to **VBE**

Gardening the WikiLexicographica. Redirection Pages

dealing with
reference entries

**pinguesso v.
pinguesco**

EU:Pinguesso

```
#REDIRECTION  
[[EU:pinguesco]]
```

dealing with
varying labels

**Lemma
Lemme**

Lemme

```
#REDIRECTION  
[[Lemma]]
```

Gardening the WikiLexicographica. Concept Pages

adding hierarchy, exposing data to the users

1st Conjugation Markers

Lemma ending with **-o(r)**

EU → conjugation number
pinguesco 1. ...

&
||

PL → inflectional endings
-ar(e|i)

1st Conjugation Concept Page

Concept:Coniugatio_I

`[[Lemma::~*o||~*or]] [[iType::~~1*||~-ar*]]`

More?

E-mail

- Bruno Bon bruno.bon@irht.cnrs.fr
- Krzysztof Nowak krzysztofnowak@ijp-pan.krakow.pl

WWW

<http://scriptores.pl/wiki>