

Mining Paper Catalogues

A Multilingual Solution to Reduce Verbose Fields to Consistent Terminology

Felix Kußmaul¹

Digital Humanities Colloquium, University of Cologne

22 November 2017

¹Archaeological Institute, University of Cologne



Form 23 Conical cup with smooth vertical rim
Konische Schale mit glattem Steilrand
Coppa troncoconica con orlo verticale
Coupe tronconique à rebord vertical lisse

Conical cup representing the further evolution of Form 22. The floor is now always flat or biconical (meeting the wall at a sharp angle on the inside), usually with a low foot.

23.1: Plain tapering rim, inclined slightly inwards, sometimes bearing applied decoration.

23.2: Rim with flat outer face bearing applied decoration bounded above and below by simple convex mouldings; inner face plain or with a groove at lip.

Production

Subform 23.1 is probably made in many parts of Italy; examples in Padana ware do not show applied decoration. Subform 23.2 is made in Italy but apparently not in the Padana region.

Date

Subform 23.2 belongs to the second and third quarters of the first century A. D.; it is common in the South Stoa deposits at Corinth and at Pompeii. Subform 23.1 is less readily datable as it may occur as a simplified version of Form 22 or Form 23; other features of the vessel (e.g. foot-profile, decoration) may provide a clearer indication of date than the shape of the rim.

Distribution

Subform 23.2 is very common throughout the Mediterranean region, with sporadic examples found in the North and in North-Italy; Subform 23.1 is relatively uncommon.

References

- 23.1.1 Karthago K 78/172a, unpublished. Stamp L.M.A., O.-C. -, Italy.
- 23.1.2 Berenice B210.2. Anepigraphic stamp. Italy.
- 23.2.1 Corinth 1973 pl.84,70. Stamp CAMVRI, O.-C. 397. Arezzo.
- 23.2.2 Berenice B216.2. Italy.

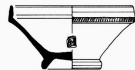
Other findspots

- 23.1 Asciburgium, Bologna, Bolsena, Conimbriga, Köln, Luni, Magdalensberg, Ordona, Pollentia, Roma.
- 23.2 Not separately listed.

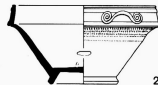
Concordance

- 23.1: Goudineau 20c; 25a; 37a. – Berenice B210.2.
- 23.2: Goudineau 40. – Barocelli 11. – Berenice B216. – Hayes 23.
- Pieces described as Haltern 9 sometimes belong to this form.

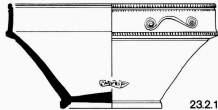
P. M. K.



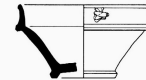
23.1.1



23.1.2



23.2.1



23.2.2

Figure 1: Sample from Ettlinger, *Conspectus*.

Problem

Running texts contain a lot of *irrelevant information* (for machine processing).
This makes database lookups without keywords **extremely inefficient**.

What we have:

Production

Subform 23.1 is probably made in many parts of Italy; examples in Padana ware do not show applied decoration. Subform 23.2 is made in Italy but apparently not in the Padana region.

Distribution

Subform 23.2 is very common throughout the Mediterranean region, with sporadic examples found in the North and in North-Italy; Subform 23.1 is relatively uncommon.

UNSTRUCTURED DATA

What we want:

```
{  
  "form": "23.1",  
  "origin": "Italy",  
  "decoration": "none",  
  "occurs": "uncommon"  
},  
{  
  "form": "23.2",  
  "origin": "Italy, not Padana",  
  "occurs": "Mediterranean region;  
             North-Italy"  
}
```

STRUCTURED DATA

What we have:

Production

Subform 23.1 is probably made in many parts of Italy; examples in Padana ware do not show applied decoration. Subform 23.2 is made in Italy but apparently not in the Padana region.

Distribution

Subform 23.2 is very common throughout the Mediterranean region, with sporadic examples found in the North and in North-Italy; Subform 23.1 is relatively uncommon.

UNSTRUCTURED DATA

What we want:

```
{  
  "form": "23.1",  
  "origin": "Italy",  
  "decoration": "none",  
  "occurs": "uncommon"  
},  
{  
  "form": "23.2",  
  "origin": "Italy, not Padana",  
  "occurs": "Mediterranean region;  
             North-Italy"  
}
```

STRUCTURED DATA

What we have:

Production

Subform 23.1 is probably made in many parts of Italy; examples in Padana ware do not show applied decoration. Subform 23.2 is made in Italy but apparently not in the Padana region.

Distribution

Subform 23.2 is very common throughout the Mediterranean region, with sporadic examples found in the North and in North-Italy; Subform 23.1 is relatively uncommon.

UNSTRUCTURED DATA

What we want:

```
{
  "form": "23.1",
  "origin": "Italy",
  "decoration": "none",
  "occurs": "uncommon"
},
{
  "form": "23.2",
  "origin": "Italy, not Padana",
  "occurs": "Mediterranean region;
             North-Italy"
}
```

STRUCTURED DATA

Definition: Text Mining

Text Mining is a **general term** covering several different ideas, e. g.:

- Statistical analysis
- Information retrieval
- Information extraction

Definition: Text Mining

Text Mining is a **general term** covering several different ideas, e. g.:

- Statistical analysis
- Information retrieval
- Information extraction
- ...

Definition: Text Mining

Text Mining is a **general term** covering several different ideas, e. g.:

- Statistical analysis
- Information retrieval
- Information extraction
- ...

Definition: Text Mining

Text Mining is a **general term** covering several different ideas, e. g.:

- Statistical analysis
- Information retrieval
- **Information extraction**

• ...

Definition: Text Mining

Text Mining is a **general term** covering several different ideas, e. g.:

- Statistical analysis
- Information retrieval
- **Information extraction**
- ...

Definition: Information Extraction (IE)

*“[IE refers to] the **identification** and **extraction** of instances of a particular class of events or relationships in a natural language text and their **transformation** into a structured representation.”*

– Grishman 1997, Eikvil 1999

Some other facts about IE:

- Computer scientists have a hard time with it (for over 30 years now!)
- IE is **really super difficult** and **often inaccurate**.

Definition: Information Extraction (IE)

*“[IE refers to] the **identification** and **extraction** of instances of a particular class of events or relationships in a natural language text and their **transformation** into a structured representation.”*

– Grishman 1997, Eikvil 1999

Some other facts about IE:

- Computer scientists have a hard time with it (for over 30 years now!)
- IE is really super difficult and often inaccurate.

Definition: Information Extraction (IE)

*“[IE refers to] the **identification** and **extraction** of instances of a particular class of events or relationships in a natural language text and their **transformation** into a structured representation.”*

– Grishman 1997, Eikvil 1999

Some other facts about IE:

- Computer scientists have a hard time with it (for over 30 years now!)
- IE is **really super difficult** and **often inaccurate**.

IE Process Pipeline

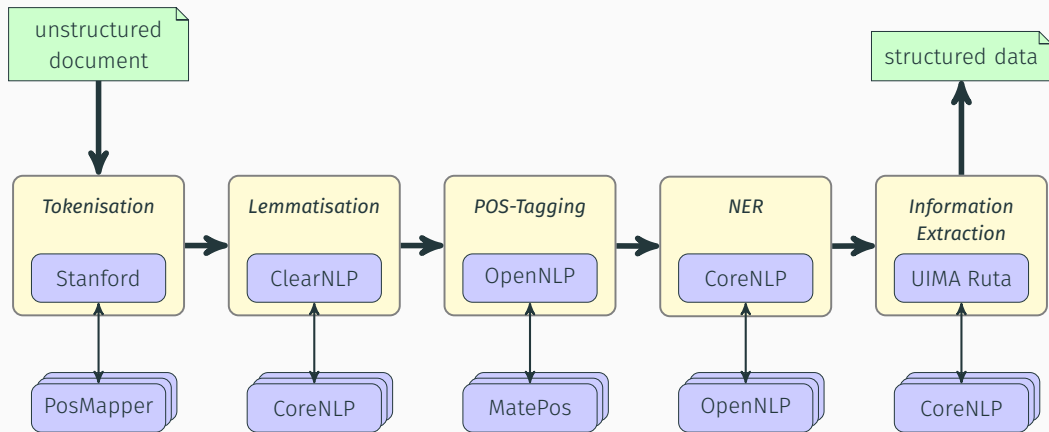


Figure 2: IE process pipeline.

Named Entity Recognition



Figure 3: POS-tagging examples after lemmatisation.

Named Entity Recognition



Figure 3: POS-tagging examples after lemmatisation.

Most NERs (e.g. **Stanford CoreNLP**) only recognise 8 entities types:

PERSON

DATE

ORGANIZATION

TIME

LOCATION

MONEY

PERCENT

MISC

So we have to add the **custom entity type FORM**.

Two approaches for NER

Rule-based approach

- High precision, but lower recall
⇒ Many many rules?!

Machine-learning approach

- Lower precision, but high recall
- Needs to be trained!

Two approaches for NER

Rule-based approach

- High precision, but lower recall
⇒ Many many rules?!

K 612

Abb. 117,2

Amphoriskos mit breitem, annähernd zylindrischem Hals und davon abgesetzter Schulter. Schmäler, aussen vorkragender Wulstrand mit an der Innenseite umlaufender breiter Riefe. Zwei Bandhenkel vom Mündungsrand zur Schulter. Form des Gefäßkörpers unbekannt.

Ton I, B mit rotem Überfang

M. Dm 9,5–10 cm

Vk. Selten

Dat. Spätes 3.–4. Jh. n. Chr.

Figure 4: Excerpt from Gempeler, *Elephantine X*.

Machine-learning approach

- Lower precision, but high recall
- Needs to be trained!

Two approaches for NER

Rule-based approach

- High precision, but lower recall
⇒ Many many rules?!

K 612

Abb. 117,2

Amphoriskos mit breitem, annähernd zylindrischem Hals und davon abgesetzter Schulter. Schmäler, aussen vorkragender Wulstrand mit an der Innenseite umlaufender breiter Riefe. Zwei Bandhenkel vom Mündungsrand zur Schulter. Form des Gefäßkörpers unbekannt.

Ton I, B mit rotem Überfang

M. Dm 9,5–10 cm

Vk. Selten

Dat. Spätes 3.–4. Jh. n. Chr.

Figure 4: Excerpt from Gempeler, *Elephantine X*.

Machine-learning approach

- Lower precision, but high recall
- Needs to be trained!

Two approaches for NER

Rule-based approach

- High precision, but lower recall
⇒ Many many rules?!

K 612

Abb. 117,2
Amphoriskos mit breitem, annähernd zylindrischem Hals und davon abgesetzter Schulter. Schmäler, aussen vorkragender Wulstrand mit an der Innenseite umlaufender breiter Riefe. Zwei Bandhenkel vom Mündungsrand zur Schulter. Form des Gefäßkörpers unbekannt.

Ton I,B mit rotem Überfang

M. Dm 9,5–10 cm

Vk. Selten

Dat. Spätes 3.–4. Jh. n. Chr.

Figure 4: Excerpt from Gempeler, *Elephantine X*.

Machine-learning approach

- Lower precision, but high recall
- Needs to be trained!

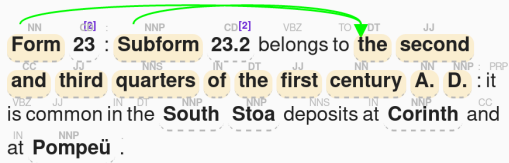


Figure 5: Manually annotated sentence from Ettlinger, *Conspectus in iepy*.

With *HEIDELTIME* temporal expressions are mapped to TIMEX3 standard

around 140 B.C.	⟶	APPROX BC0140
Spätes 3.-4. Jh. n.Chr.	⟶	END 02; 03
second quarter first century B.C.	⟶	XXXX-Q2 BC00
first half third century A.D.	⟶	XXXX-H1 02

HEIDELTIME supports many other languages, e.g. German, Italian, French, ...

Locations with HEIDELPLACE

Text document:

My name is Ludwig. This has nothing to do with the town Ludwig.
I study at Ruprechts-Karl-Universität in Heidelberg, Germany.
New York Times is a newspaper. It is 5:30 p.m.
U.S. is an abbreviation for the United States of America.
The Simpsons live in Springfield.
What about the place Steffies Hostel Heidelberg?

Recognition Modules:
STANFORD_NER, GAZETTEER_...
Linking Modules:
GAZETTEER_LOOKUP_EXACT
Disambiguation Modules:
HIGHEST_ADMIN_LEVEL_DISA...

Geoparse

Find Toponyms

Link Toponyms

Disambiguate Toponyms

Cancel

Reset

Running geoparsing step ALL took 5768ms

STANFORD_NER, GAZETTEER_LOOKUP_EXACT, HIGHEST_ADMIN_LEVEL_DISAMBIGUATION
DISAMBIGUATION took 2473ms

My name is **Ludwig** . This has nothing to do with the town **Ludwig** .
I study at Ruprechts-Karl-Universität in **Heidelberg** , **Germany** .
New York Times is a newspaper. It is 5:30 p.m.
U.S. is an abbreviation for the **United States of America** .
The **Simpsons** live in **Springfield** .
What about the place Steffies Hostel **Heidelberg** .

Named Entity	Type	Linked Places
Ludwig	PERSON	3468
Ludwig	PERSON	
Heidelberg	LOCATION	
Germany	LOCATION	
New York Times	ORGANIZATION	
U.S.	LOCATION	

Figure 6: Screenshot of HEIDELPLACE.

Relation Extraction

Subject	Relation	Object
quick brown fox	jump over	lazy dog
K 612	dates	03 ¹
Subform 23.2	occurs	North Italy
Subform 23.2	dates	XXXX-Q2 00 ²

⇒ e.g.

```
{  "form":    "23.2",  
  "dating":  "XXXX-Q2 00"  
  "occurs":  "North Italy" }
```

¹“4th century A. D.”

²“second and third quarters of the first century A. D.”

Relation Extraction

Subject	Relation	Object
quick brown fox K 612	jump over dates	lazy dog 03 ¹
Subform 23.2	occurs	North Italy
Subform 23.2	dates	XXXX-Q2 00 ²

⇒ e.g. `{ "form": "23.2",
"dating": "XXXX-Q2 00"
"occurs": "North Italy" }`

¹"4th century A. D."

²"second and third quarters of the first century A. D."

Multilingualism

Two problems:

- Linguistic
- Conceptual

Different languages



Different traditions

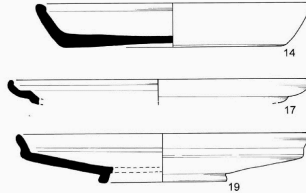
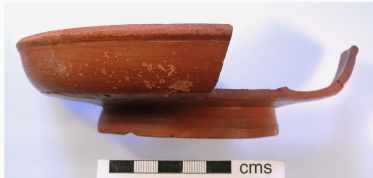


Figure 7: Plate, platter or dish?

Creating controlled vocabularies

Creating wordlists that project team would be most useful to describe the key features of a vessel or sherd

- Sherd type (e.g. rim or handle)
- Form (e.g. plate or bowl)
- Decoration form (e.g. burnished)
- Decoration color (e.g. yellow)
- Fabric (e.g. Dressel 28 fabric)



Used tools and methodology developed for the ARIADNE project by the Hypermedia Research Group at the University of South Wales

- Created a neutral spine based on the Getty Institute's Art and Architecture Thesaurus (AAT)
- This spine was populated by members from partner organisations, identifying common terms and concepts within it
- Project partners then mapped terms in their language to this neutral spine
- French terms supplied courtesy of a 2001 Masters thesis by Caroline SOURZAT (thanks to Eleni Schindler Kaudelka for identifying this on the ArchAIDE blog!)

Mapping terms and concepts (part 1)

Often this was very straightforward, for example:

- The Italian terms *graffita*, *graffita a punta*, *graffita a stecca* = “sgraffito” (<http://vocab.getty.edu/aat/300266416>)
- The Spanish term *Cántaro* = “jars” (<http://vocab.getty.edu/aat/300195348>)
- The German terms *gebogener Henkel*, *Ohrförmiger Henkel*, *langer Vertikalhenkel* = “handles” (<http://vocab.getty.edu/aat/300266416>)

Mapping terms and concepts (part 2)

Often this was more complicated, with partners having differing perceptions on what to call something (e.g. “plate” versus “platter”)

In truth, this confusion may also be reflected by what has come out of the ground!

An advantage of using the AAT (a “SKOS’d” thesaurus), is that ambiguity or difference in nomenclature can be resolved by a broader term or concept, so for example ...

Mapping terms and concepts (part 3)

Looking at the hierarchies for plate and platter in the AAT we can see that both are “dishes (vessels for food)”, or even broader “culinary containers”. So whole we can retain our original classifications (and this is essential for text mining), we can agree at a fundamental level *what these fundamentally are*



Figure 8: AAT Hierarchies for Plate and Platter

- Recognize reigns of emperors as **DATE** entities
- Coreferences in general
- HEIDELTIME:
second and **third quarter** of the **first century A. D.** \mapsto XXXX-Q3; 00
- Returning to difference in ceramic recording details
- Fabric names often contain locations, e. g. *Magdalensberg xyz*
- Location sometimes narrow, sometimes whole regions
- In many cases the form is not named in particular but just described

References

-  Ettlinger, Elisabeth. *Conspectus formarum terrae sigillatae Italico modo confectae*. Ed. by Deutsches Archäologisches Institut zu Frankfurt and Römisch-Germanische Kommission. Materialien zur römisch-germanischen Keramik. Bonn: Habelt, 1990.
-  Gempeler, Robert D. *Elephantine X. Die Keramik römischer bis früharabischer Zeit*. Mainz: Von Zabern, 1992.
-  Indurkha, Nitin and Fred J. Damerau. *Handbook of Natural Language Processing*. 2nd. Chapman & Hall/CRC, 2010.
-  Richter, Ludwig et al. “HeidelPlace: An Extensible Framework for Geoparsing”. In: *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. Copenhagen, Denmark: Association for Computational Linguistics, 2017, pp. 85–90.
-  Strötgen, Jannik and Michael Gertz. “HeidelTime: High Quality Rule-Based Extraction and Normalization of Temporal Expressions”. In: *Proceedings of the 5th International Workshop on Semantic Evaluation*. Uppsala, Sweden: Association for Computational Linguistics, July 2010, pp. 321–324.

<http://www.archaide.eu/blog>

Thank you very much for your attention!

Questions?

This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement № 693548



Mining Paper Catalogues

A Multilingual Solution to Reduce Verbose Fields to Consistent Terminology

Felix Kußmaul¹

Digital Humanities Colloquium, University of Cologne

22 November 2017

¹Archaeological Institute, University of Cologne

