# Mining Paper Catalogues

## A Multilingual Solution to Reduce Verbose Fields to Consistent Terminology

Tim Evans[1],    Felix Kußmaul[2]

23rd Annual Meeting EAA, Maastricht

31 August 2017

[1]Archaeology Data Service, University of York

[2]Archaeological Institute, University of Cologne

ARCHAIDE

ARCHAEOLOGICAL
AUTOMATIC INTERPRETATION
AND DOCUMENTATION
OF CERAMICS

# MOTIVATION

# Data Source

Figure 1: Sample from Ettlinger, *Conspectus*.

# Oh dear!

## 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

# TEXT MINING

# Classification

## Definition: Text Mining

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

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

### Definition: Text Mining

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

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

# Classification

### Definition: Text Mining

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

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

# Classification

## Definition: Text Mining

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

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

# Classification

## Definition: Text Mining

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

- Information retrieval
- Statistical analysis
- **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 Information Extraction:

- Computer scientists have a hard time with IE (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 Information Extraction:

- Computer scientists have a hard time with IE (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 Information Extraction:

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

# Sorry!

## DISCLAIMER

In this presentation, we show **preliminary** results, as this project is still work in progress.
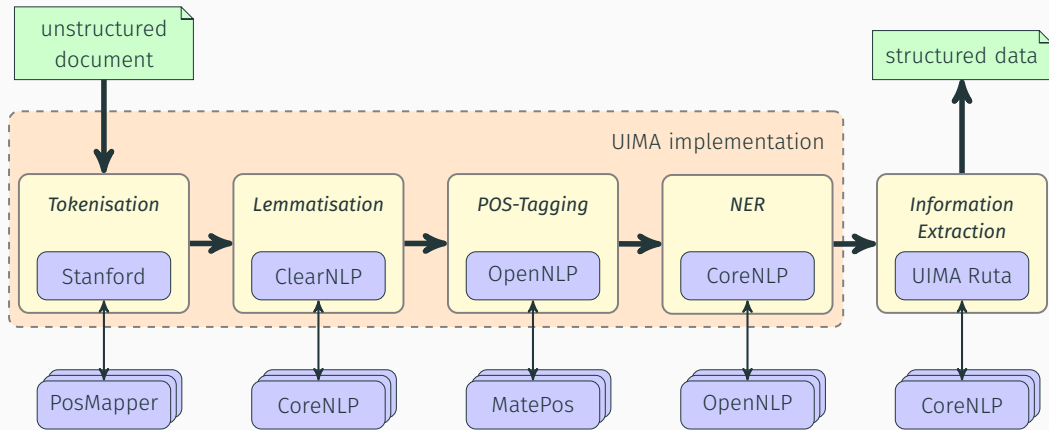
Figure 2: IE Process Pipeline.

# POS-Tagging

The    quick    brown    fox    jump (jumps)    over    the    lazy    dog    .

| DT | JJ | JJ | NN | VBD | IN | DT | JJ | NN | . |

Figure 3: POS-tagging examples after lemmatisation.

jumps

The  quick  brown  fox  jump  over  the  lazy  dog  .
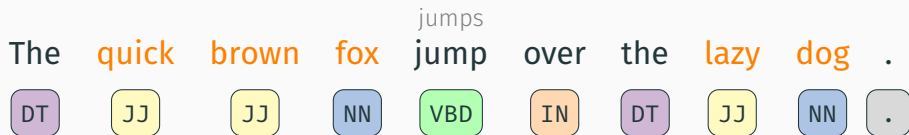
| DT | JJ | JJ | NN | VBD | IN | DT | JJ | NN | . |

**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!

## 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. Schmaler, 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
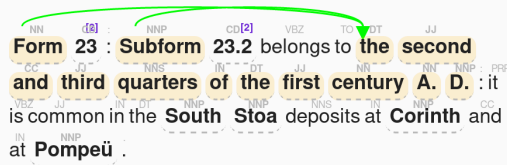M.    Dm 9.5–10 cm
Vk.   Selten
Dat.  Spätes 3.–4. Jh. n. Chr.

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

## Machine-learning approach

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

## 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. Schmaler, 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!

## Rule-based approach

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



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`.

## Temporal Expressions

With *HEIDELTIME* temporal expressions are mapped to TIMEX3 standard

$$
\begin{aligned}
\text{around 140 B.C.} &\longmapsto \text{APPROX BC0140} \\
\text{Spätes 3.–4. Jh. n.Chr.} &\longmapsto \text{END 02; 03} \\
\text{second quarter first century B.C.} &\longmapsto \text{XXXX-Q2 BC00} \\
\text{first half third century A.D.} &\longmapsto \text{XXXX-H1 02}
\end{aligned}
$$

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

## Relation Extraction

| Subject | Relation | Object |
|---|---|---|
| quick brown fox | jump over | lazy dog |
| K 612 | dates | 03[1] |
| Form 23 | dates | XXXX-Q2 00[2] |
| Subform 23.2 | dates | XXXX-Q2 00 |

$\Rightarrow$ e.g.
```
{ "form":    "23.2",
  "dating": "XXXX-Q2 00" }
```

---

[1] "4th century A. D."

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

# Relation Extraction

| Subject | Relation | Object |
|---------|----------|--------|
| quick brown fox | jump over | lazy dog |
| K 612 | dates | 03[1] |
| Form 23 | dates | XXXX-Q2 00[2] |
| Subform 23.2 | dates | XXXX-Q2 00 |

$$\Rightarrow e.\,g. \quad \{ \ \ \texttt{"form":} \quad \texttt{"23.2",}$$
$$\texttt{"dating": "XXXX-Q2 00"} \ \ \}$$

---

[1] "4th century A. D."

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

# Locations with HEIDELPLACE



**Figure 6:** Screenshot of HEIDELPLACE.

# MULTILINGUALISM

# Background

Two problems:
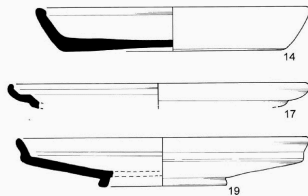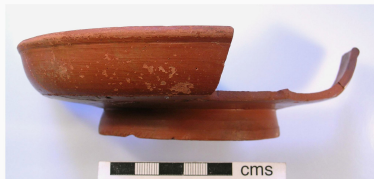
- Linguistic
- Conceptual

Figure 7: Plate, platter or dish?

# Creating controlled vocabularies

- Sherd type (e.g. rim)
- Form (e.g. plate)
- Decoration form (e.g. burnished)
- Decoration color (e.g. yellow)
- Fabric (e.g. bla)

# Lessons from ARAIDNE

Using tools developed for the ARIADNE project by the Hypermedia Research
Group at the University of South Wales

Creation of a neutral spine based on the Getty Institute's Art and Architecture
Thesaurus (AAT)

more bla

more bla

more bla

# OUTLOOK

## Challenges to meet

challenges:

choice of tools, coreferences in text, eloquence of archaeologists, maybe calculating F-value?

HEIDELTIME:

second and **third quarter** of the **first century A. D.** $\longmapsto$ XXXX-Q3; 00

# 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.

📕 Indurkhya, Nitin and Fred J. Damerau. *Handbook of Natural Language Processing.* 2nd. Chapman & Hall/CRC, 2010.

Thank you very much for your attention!

Questions?

# Mining Paper Catalogues

A Multilingual Solution to Reduce Verbose Fields to Consistent Terminology

Tim Evans[1],   Felix Kußmaul[2]

23rd Annual Meeting EAA, Maastricht

31 August 2017

[1]Archaeology Data Service, University of York

[2]Archaeological Institute, University of Cologne

ARCHAIDE

ARCHAEOLOGICAL
AUTOMATIC INTERPRETATION
AND DOCUMENTATION
OF CERAMICS