



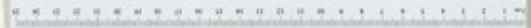
# Efficient Algorithms and User Interaction for Metadata Extraction from Historical Maps

Benedikt Budig, Universität Würzburg

# Overview

- Historical Maps: what and why?
- Sketch of a Pipeline
  - from bitmap image to georeferenced metadata
- Module 1: Locate Map Elements
- Module 2: Match Labels and Place Markers
- Open Questions & Future Work





*Ortschaft der Gebürteten GRAFFSCHAFT HENNEBERG, mit dem angrenzenden FÜRSTENTH. COBURG u. andern GRAENZLÄNDERN, nach authent. alten Documenten, und Nachrichten, verzeichnet und in der Ordnung alphabet. Anzettelten von Anzettelten Orten in Preuss.*

La CARTE du Comté de HENNEBERG avec les parts confins de Principauté de COBURG. Le tout subdivisé en ses Baillages et droits, selon les mémoires les plus authentiques.  
*Les Mémoires tirés des papiers d'un Prince de Saxe.*  
*A Nuremberg chez les Libraires de l'Empire, M.D.C.C.XX.*

ut subdivisé en ses B  
Homañ. A. 1743.



Nach Selig  
Leuffenorth  
Marb am Mayn  
Leutendorff  
an  
au  
ach El  
M

S. R. I. COMITATVS  
**HENNEBERG**  
secundum Præfectu-  
ras & modernas Dy-  
nastias, una cum con-  
fini PR:COBURGENSI geo-  
graphice consignatus  
& in hac Tabula editus.











## Significatio notarum



Fortalitia

2 Stahlwerck.



Urbes prae-  
cipuae.

⊖ Salzwerck.



Oppida  
Municipia.

3 Eisen Gruben.

⊕ Vitriol.



Pagi.

Δ Glasz. und  
Schmelz hütte



Villae.

⊙ Alaun.



Arces.

K Kobalt.



Coenobia.

Gr. Granit. od. Eg.



Ruinæ.

ypt: Marmor.



W. Wüstung.

Lf. Stein Kohlen.



Gold.

/ Gewehr Fabrique



Silber.

MB. Marm. Bruch.



Kupfer.

U. Umbra.

7 Bergwercke.

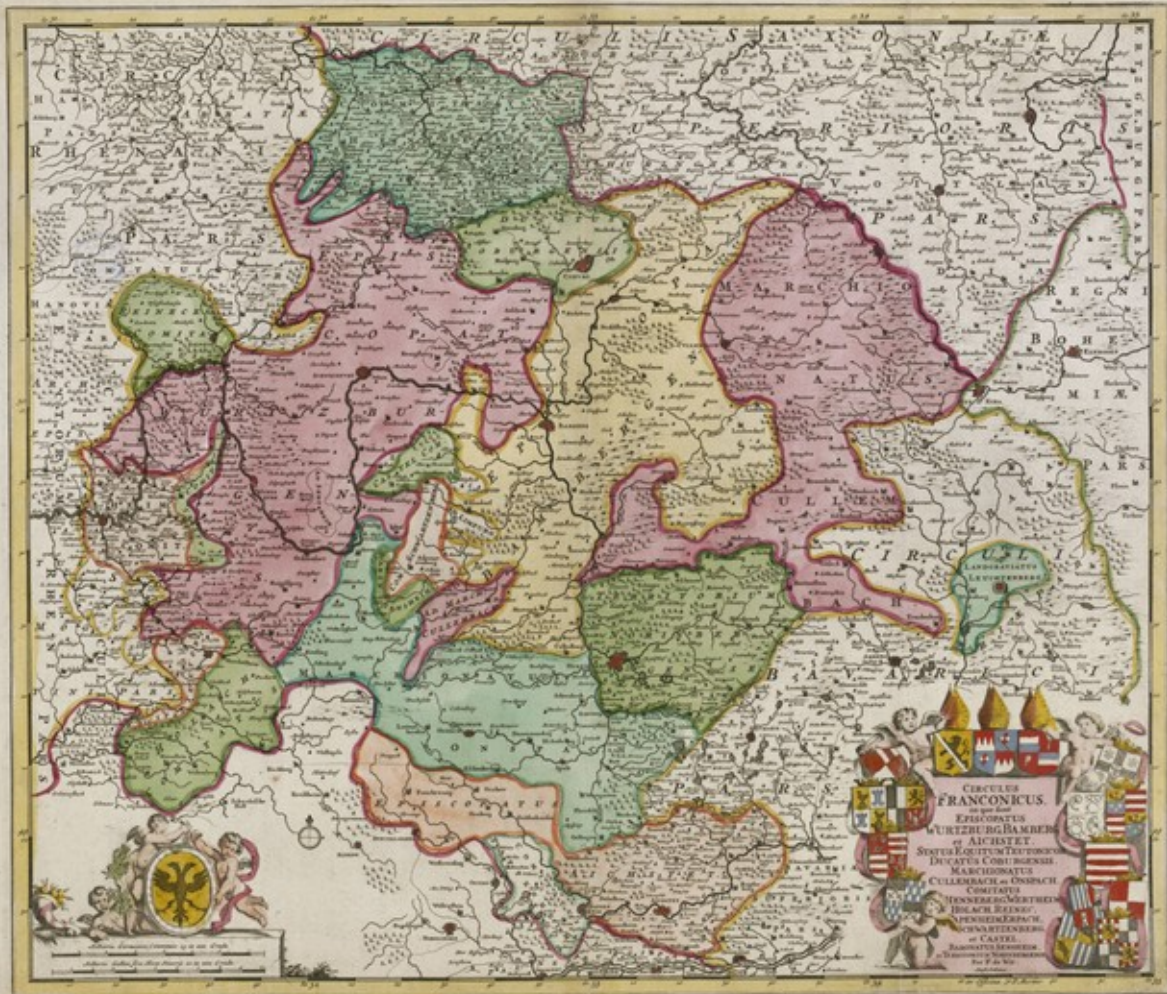


Altendamm  
achbach  
Friedrichs  
Hirschbach  
Erlau  
Leupoldshof  
Wilcken  
hoff  
Silbach  
S. Kilian  
Rajen  
Winternahe  
K...  
Vesser  
Stuten  
haus  
Breitenbach  
Amt  
Schleus  
Neundor  
Steinbach  
Schleusingen  
Lanobach  
Schönau  
Frauen  
Ster  
Ahlstätt  
Fischbach  
Ziegehof  
Neuhof  
S. Kilian  
Rajen  
Winternahe  
K...









**CIRCULUS  
FRANCONICUS.**  
EPISCOPATUS  
WURZBURG-HAMBERG  
et ARCHIEP.  
SANTA EQUITUM TUTORUM  
DUCATUS COMITATUS  
MARCHIONATUS  
CULEMBACH et OBERACH  
COMITATUS  
HONNERGEMERTHEIM  
HILACH REGENS  
HESSENKREFTZ  
SCHNARTENBERG  
et CASTEL  
RABENSTEIN  
et SCHNEIDERHOF  
et F. de N. de





Haslach

1

Schlusselfeldt

Ellendorf

AD LIMPURG

Speckfeldt

Ober Schainfeldt

1

Marck Pibrach

Marckschaintfeldt

Lamach

Erensdorf

Ebrack Flu

Schlusselfeldt

Pommersfeldt

Aisch Flu

Muthausen

Adelsdorf

Hallersdorf

Neuhaus

Kaufen

Hochstett

Clepach

Hanburg

Morendorf

Rotenbach

Puttenhaid

Hirshaid

Neuses

Forchheim

Kerspach

Paierdorf

Vite

Erlang

Buck

Dachsbach

Wesent Flu



# Das Francken Land Chorographi Franckia Dre



Dem Hochfürstlichen in Westphalen Erbprinzen  
Johann Christianen: Vorgetragen von  
Wilhelm Friderich Schöner,  
Königlichen Hofrath.

Es hat hochfürstliche erucht, In dem Erbprinzen  
Johann Christianen, den 14. Junii 1617, zu  
Frankfurt am Main, zu sein, dass er die  
Karte von Francken, welche er von dem  
Hochfürstlichen Rathe, Herrn  
Johann Christianen, erhalten hat, in  
seinem Reichthum, zu vergrößern, und  
zu verbessern, befohlen hat.



Die Beschreibung  
des Francken  
Landes  
von  
Johann  
Christianen  
1617







# Study historical maps: why?

- Many libraries have large collections of historical maps
- Relevant for the (digital) humanities
  - History of cartography
  - General history
  - Specific example: onomastics



# What happens with historical maps?

- Stored in a library basement
  - Retrievable by bibliographic information
- High-quality bitmap scans, online catalogue
  - Browsable by bibliographic information
- Useful queries?
  - In actual research practice
  - By interested laypeople



# Metadata: what?

- Contained settlements
  - Landscape topography
  - Geopolitical features
  - ...
- 
- Ideally: georeferenced



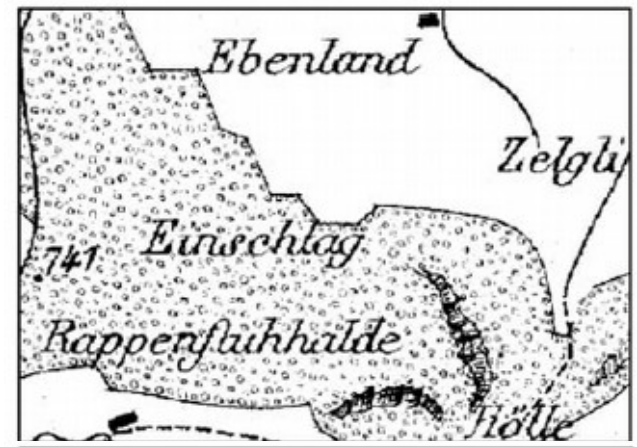
# Metadata: how?

- Do it by hand
- Software: usability improvements
  - Gains in efficiency are limited
- Software: computer vision
  - No panacea, but can work well for restricted corpora
  - Significant custom R&D effort every time



# For example...

- Forest-cover analysis of the "Siegfried Map"  
[Leyk, Boesch, Weibel]
- 6000 sheets, produced 1870 to 1922



(f)



# Our scope

- We consider maps from early modern period forward
- Unique graphical styles
- Different fonts, handwriting
- Different cartographic conventions

# So what now?

- Split problem into smaller goals
- Design a modular pipeline



## Segmentation



# Clustering and Matching



# Optical Character Recognition



## Georeferencing





Segmentation



Clustering  
Matching

# Segmentation

- Smaller goals
- Look for one particular element on one map

*[Budig and van Dijk 2015]*





Pollich

Leibes

**Damburg**

五

Stadelhousen

Güspach

Remerci

Оберг

Scheßlig

સાતેલપોર

260  
arr

ary

Afterdeck

Esapfendorf

使biffelt

2017年

४८३

Quinta

Quinta

[Demo]



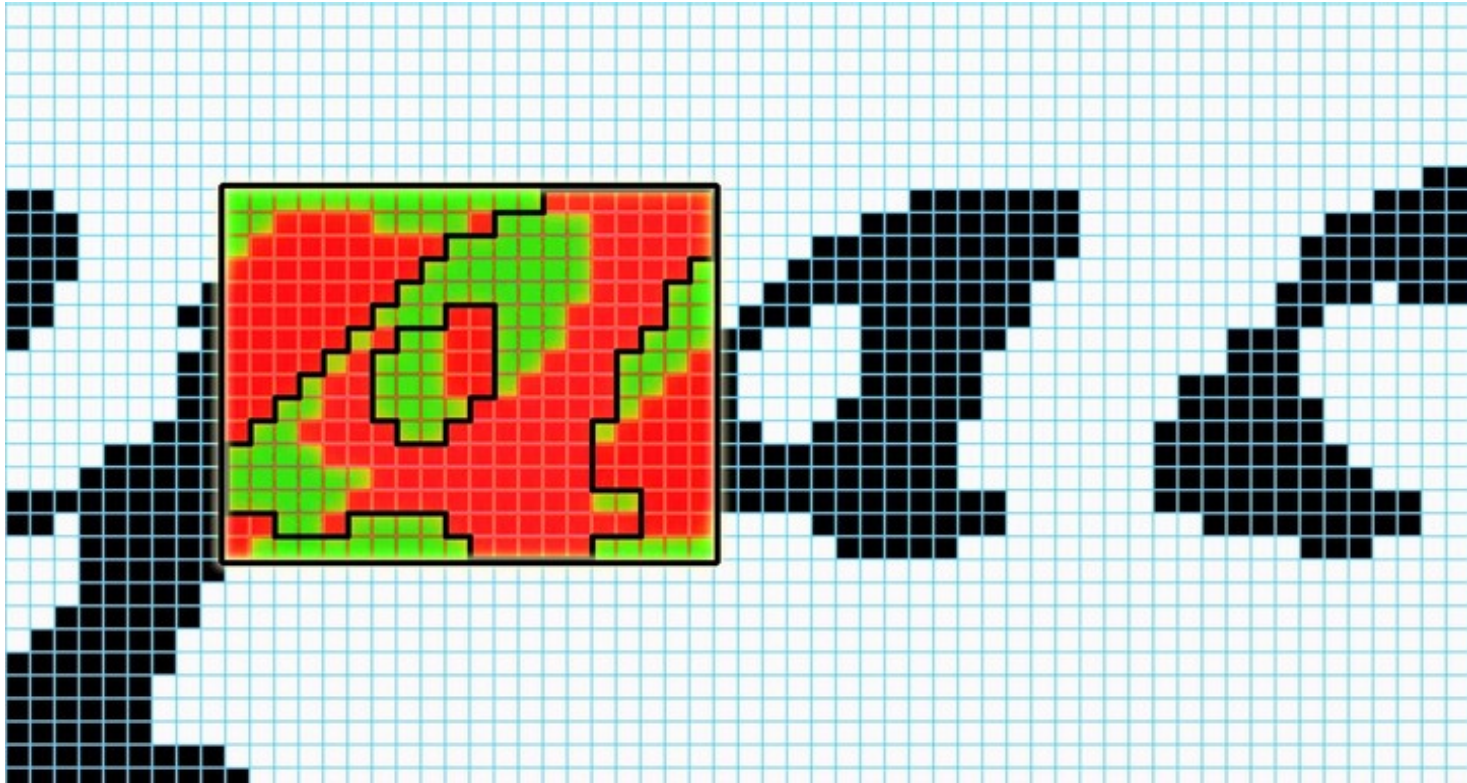
# Ingredient 1: Template Matching

# Template matching

- Find approximate repeat-occurrences of an example image
- Here:
  - 1-Bit black-and-white
  - Only considers translation
  - Percentage-correct score









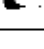

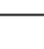
# Matching score



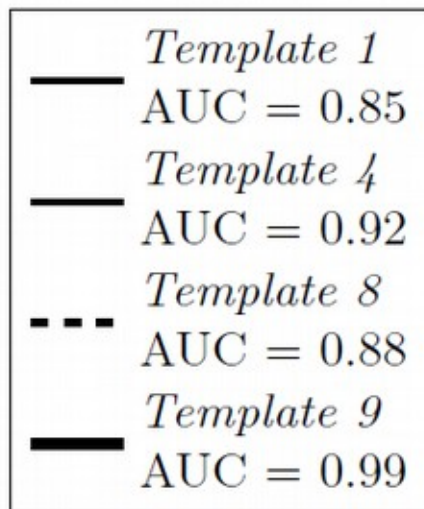
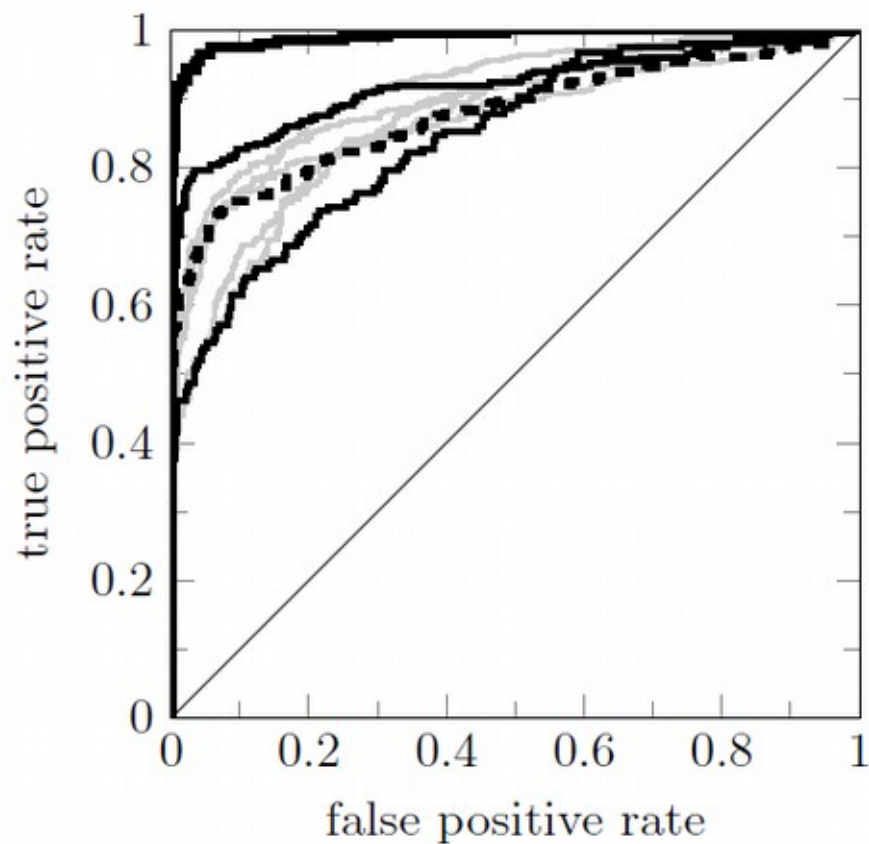
## Ingredient 2: Active Learning



# Data sets

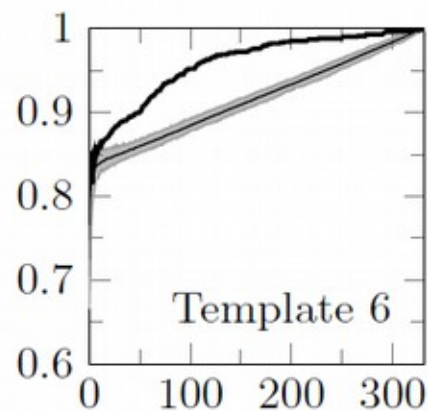
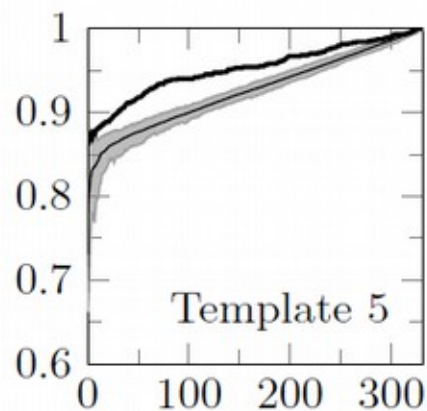
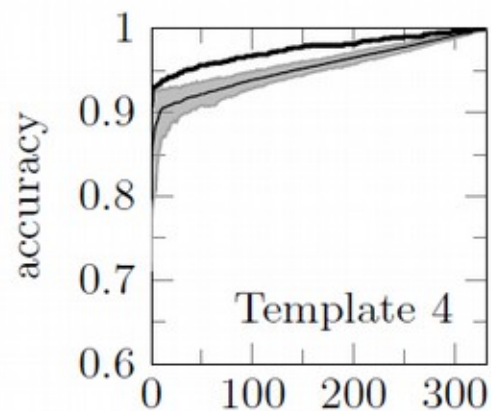
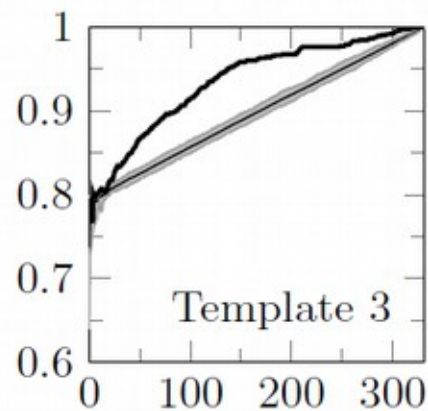
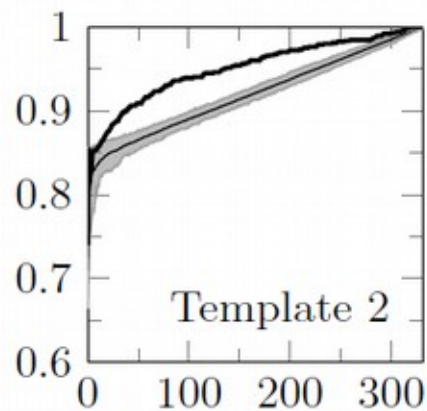
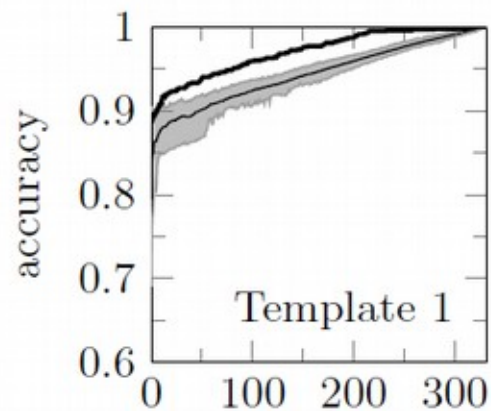
	Historical Map	Template	Accepted
1	<i>Carte Topo. D'Allemagne</i> (1787)		b, h
2	<i>Franciae Orientalis</i> (1570)		a, g, d
3	<i>Franciae Orientalis</i> (1570)		e
4	<i>Circulus Franconicus</i> , De Wit (1706)		a, g, d
5	<i>Das Franckenlandt</i> (1533)		a, g
6	<i>SRI Comitatus Henneberg</i> (1743)		n, m, h
7	<i>SRI Comitatus Henneberg</i> (1743)		e
8	<i>Circulus Franconicus</i> , De Wit (1706)		
9	<i>Circulus Franconicus</i> , Seutter (1731)		

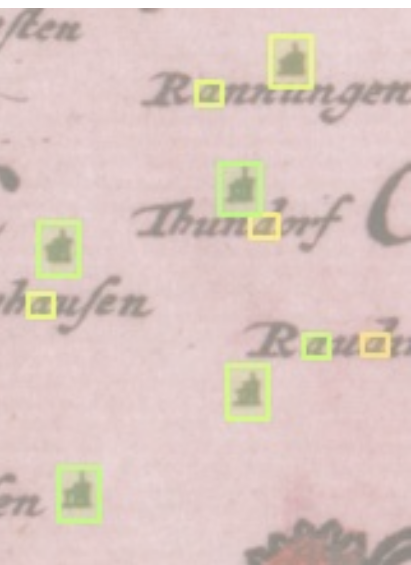
# ROC





# Learning curves





Segmentation



Clustering and  
Matching



Optical Character  
Recognition



# Matching Labels and Place Markers

- Assumption: labels and markers already detected
- Match the corresponding ones *[Budig, van Dijk, Wolff, 2014]*



# Wanted: a Matching

- Find a *matching* of labels and place markers
- No 1-to-1 assignment possible
- Basic assumption: labels are *near* their corresponding markers
- Greedy strategy?
  - does not work well!



# Experimental Results

- Franckenlandt (1533)

- 539 markers, 524 labels
- our algorithm: error rate 3.5%
- greedy algorithm: error rate 17.8%



- Circulus Franconicus (1706)

- 1663 markers, 1669 labels
- our algorithm: error rate 1.3%
- greedy algorithm: error rate 5.9%





# What now?

- Error rates in experiments: 1.3% and 3.5%
- Unclear situations:



- Manual verification or correction necessary

# Sensitivity

- Only show assignments our algorithm is *uncertain* about
- Calculate sensitivity analysis for the matching
  - For each assignment,
    - calculate the best matching that does *not* use it
    - and see how much worse it is.



# Conclusion

- Historical maps are cool, but hard to search
- Modular pipeline is reasonable
- Human effort is necessary → smart interactions!
- Template matching & active learning work well
- Sensitivity analysis for efficient interactions



# Open Questions & Future Work

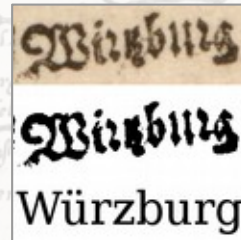
- Solve more small goals from the pipeline, then integrate
  - Cluster template matches (e.g. into labels)
  - Use already collected information for OCR
  - Georeferencing, ...
- Should the pipeline really be sequential?
- Crowdsourcing?



Segmentation



Clustering and  
Matching



Optical Character  
Recognition



Georeferencing

# Smartphone





# Open Questions & Future Work

- Develop remaining modules in extraction pipeline
  - Cluster template matches (e.g. into labels)
  - Use already collected information for OCR
  - Georeferencing, ...
- Should the pipeline really be sequential?
- Crowdsourcing! Yes, but how exactly?
- What other algorithmically-guided user interactions?