# Layout Detection and Table Recognition

Recent Challenges in Digitizing Historical Documents and Handwritten Tabular Data

# About us

**Constantin Lehenmeier**

IT EMPLOYEE / UNIVERSITY LIBRARY OF REGENSBURG

PHD STUDENT / CHAIR OF MEDIA INFORMATICS / UNIVERSITY OF REGENSBURG

**Junior Prof. Dr. Manuel Burghardt**

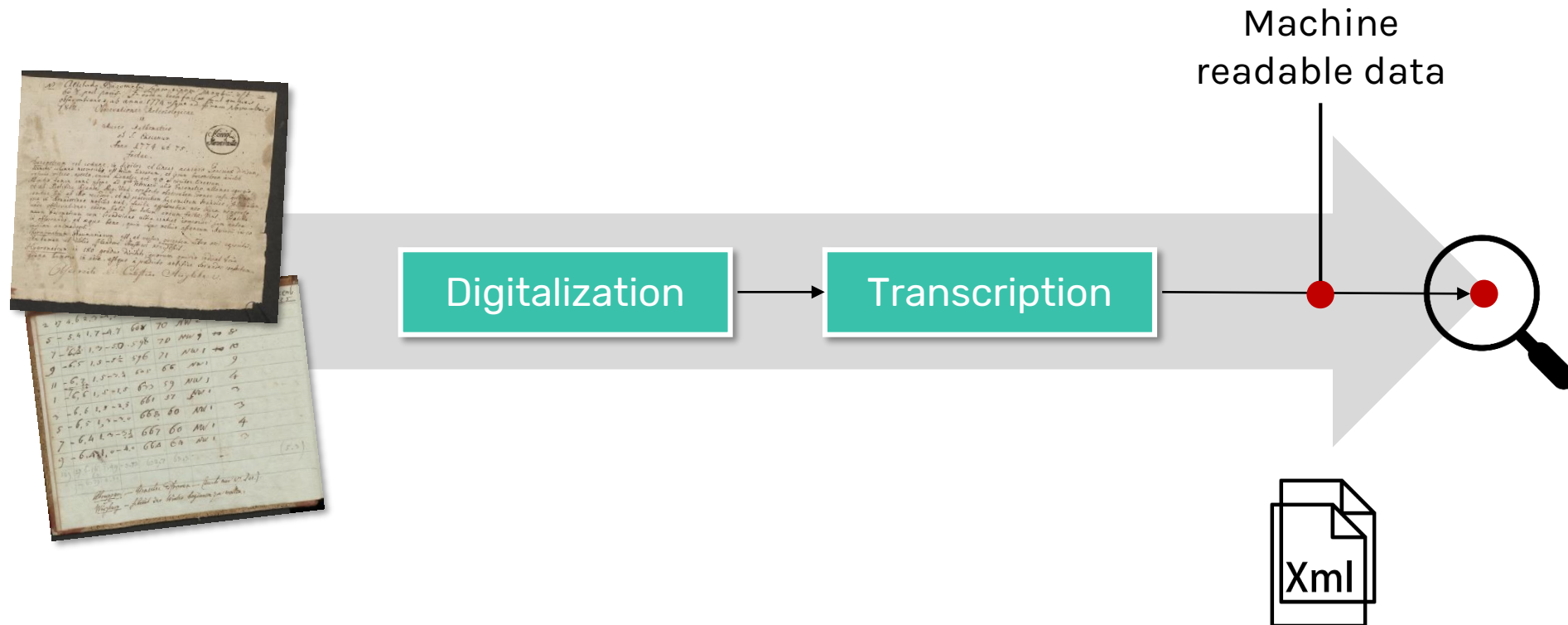CHAIR HOLDER COMPUTATIONAL HUMANITIES / UNIVERSITY OF LEIPZIG

**Bernadette Mischka**

PHD STUDENT / CHAIR OF EUROPEAN HISTORY/ UNIVERSITY OF REGENSBURG

# Libraries and the digitzed humanities

The digitization of cultural artifacts is part of the "digitized humanities"



Machine readable data

Digitalization → Transcription

Xml

... there are billions of documents

# Automated document recognition

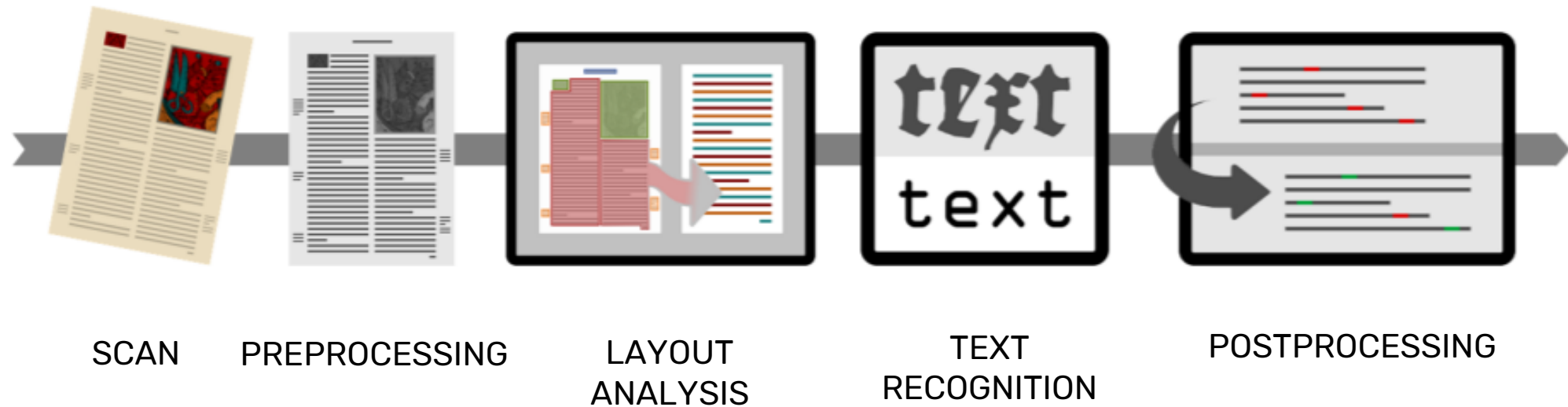OCR is used to automatically convert visible language to a searchable digital format



| SCAN | PREPROCESSING | LAYOUT ANALYSIS | TEXT RECOGNITION | POSTPROCESSING |

Figure 1: The main steps of the OCR process

# State-of-the-art

CER under 1% on modern printed books without specific training

CER under 2% on early printed books with specific training

Character Error Rate (CER) is the edit distance between two sequences

# State-of-the-art

CER under 1% on modern printed books without specific training

CER under 2% on early printed books with specific training

Character Error Rate (CER) is the edit distance between two sequences

❗Historical documents pose special challenges

**DEGRADATION**

**INCONSISTENT LAYOUT**

**NON-STANDARDIZED TYPOGRAPHY**

Figure 2: Main challenges of the automatic
recognition of historical documents

# Observationes meteorologicae

Continious weather reports for over 53 years. The university library owns 55 volumes for the years 1774 – 1827.

# There is something about the weather…

Reconstruction of past climate conditions

Historical climate impact research

Scientific history of the climate

Asses futues climate fluctuations

9

# Problems

**Tables** + **Handwritten text**

# The authors

The reports were written by three scientists



**Coelestin Steiglehner**

*1771 - 1778*

**Placidus Heinrich**

*1771 - 1825*

**Ferdinand von Schmöger**

*1825 - 1827*

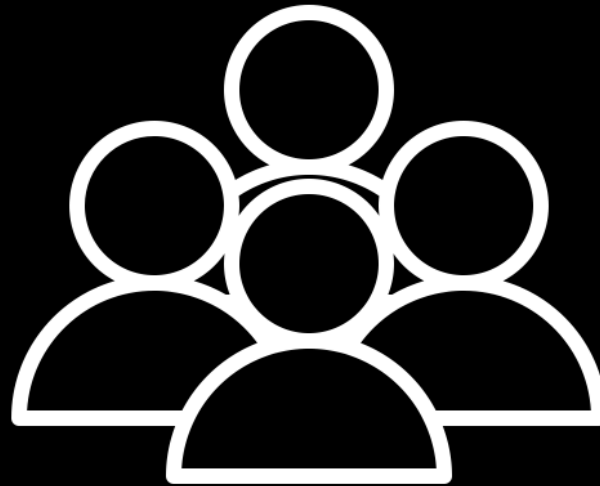Figure 3: The main authors of the observationes meteorologicae

# The authors

The reports were written by three scientists + various students



**Coelestin Steiglehner**

*1771 - 1778*

**Students**

*1791 - 1798*

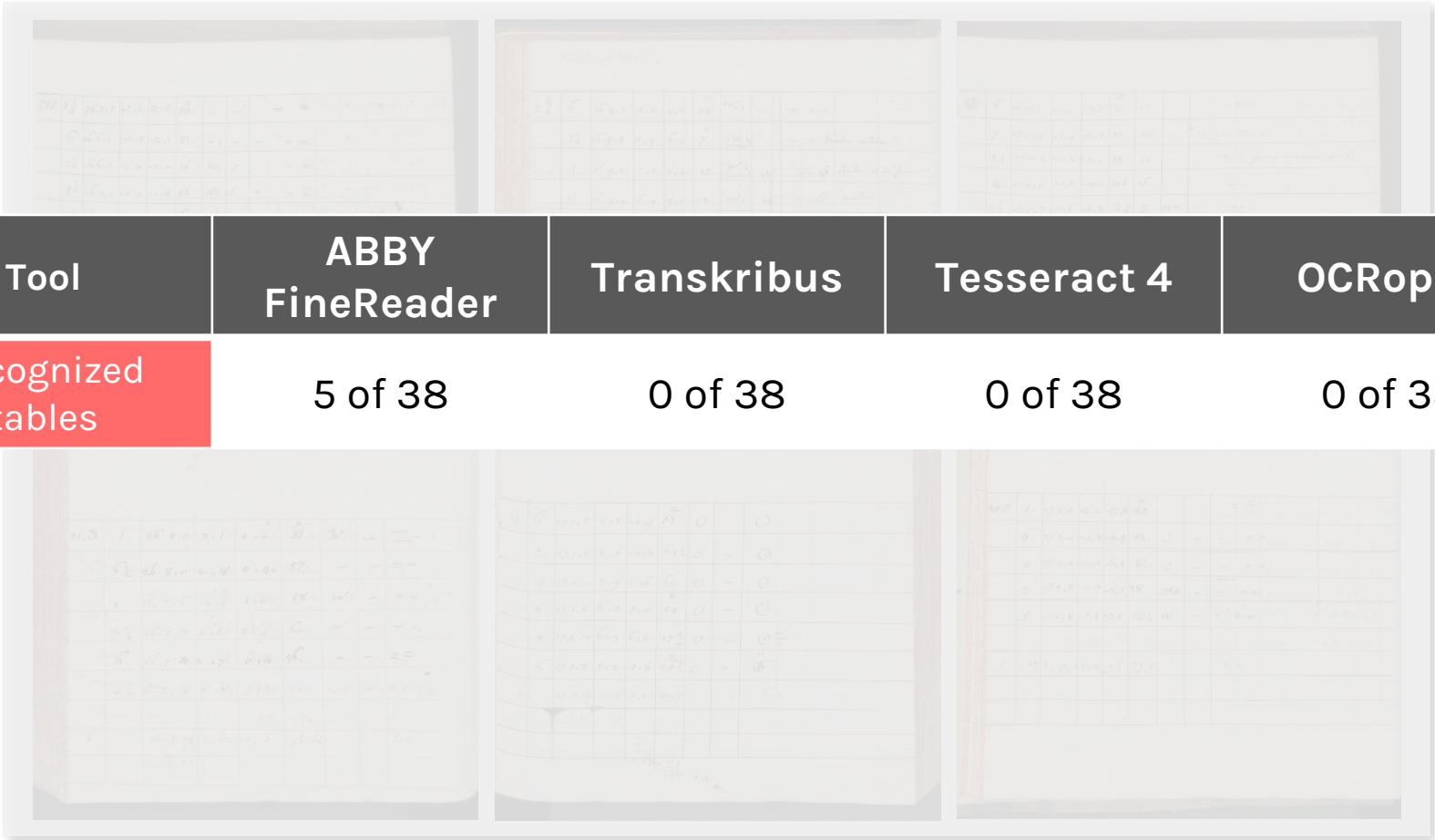**Ferdinand von Schmöger**

*1825 - 1827*

# Results of out-of-the-box tools

A test with four existing OCR engines was conducted on 38 pages of one volume (1793)
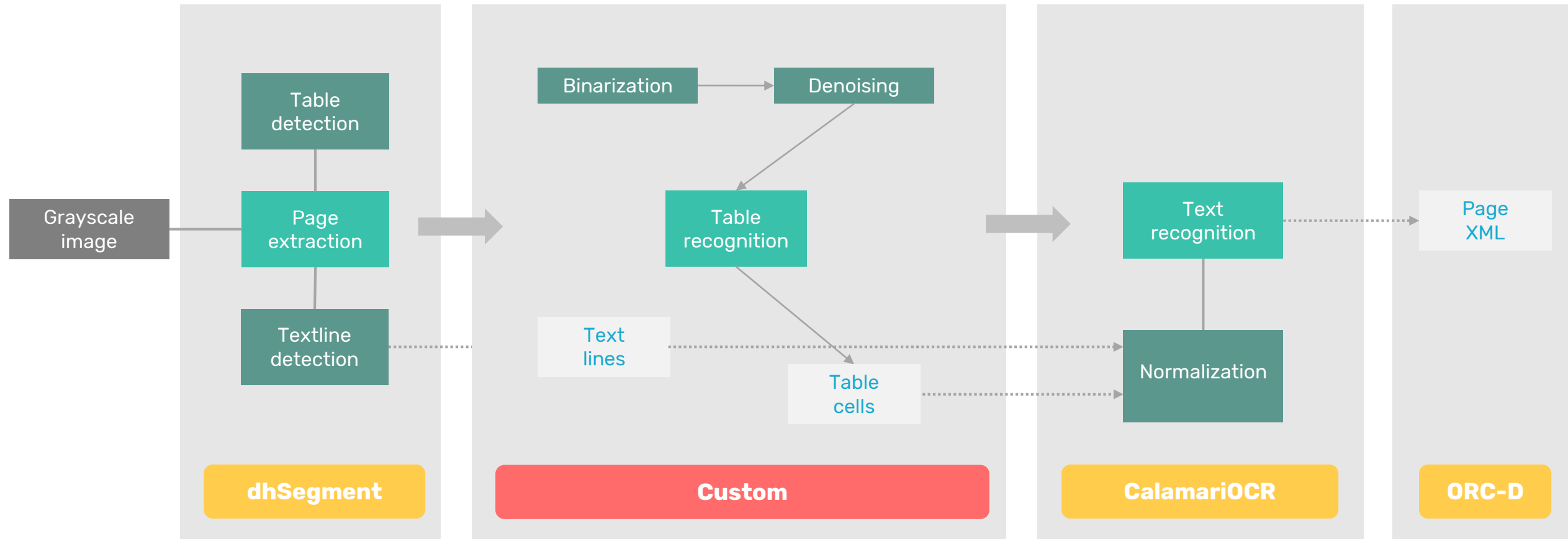
# Results of out-of-the-box tools

A test with four existing OCR engines was conducted on 38 pages of one volume

| Tool | ABBY FineReader | Transkribus | Tesseract 4 | OCRopus |
|---|---|---|---|---|
| Recognized tables | 5 of 38 | 0 of 38 | 0 of 38 | 0 of 38 |

# Current workflow

# Layout analysis and table recognition

# Layout analysis and table recognition

# Evaluation of the table recognition

| | |
|---|---|
| Training set | 40 pages |
| Evaluation set | 38 pages |
| Table detection | 100% |
| Table recognition | 87% |

Table detection results were evaluated by using the Jaccard index

Table recognition results were evaluated by converting the structure to HTML and computing the BLEU metric

# Text recognition

| | |
|---|---|
| Model | C, Mp(2x2), C, Mp(2x2), LSTM(200) |
| Training set | 32 pages (~1.800 text lines) |
| Evaluation set | 4 pages |
| CER | 25,86% |

# Recent challenges

## HISTORIC METEOROLOGICAL SYMBOLS

| HISTORIC | SEMANTIC | UNICODE | VISUAL |
|---|---|---|---|
|  | Coelum totum serenum et fine nube est. | ☉ | ☉ |
|  | Nubes et coeruleum aeque diuisum. | 〰 | = |
|  | Nubeculae paucae plerum que albicantes. | ⚡ | ⁛ |
|  | Nubes et nubeculae minori coeli parte. | ⚡ | ⚡ |

### How to map certain symbols?

- not every symbol has a semantic or visuel Unicode pendant

## BIG GT OR SMALLER SETS?



### How to generate a sufficient GT data set for the entire collection?

- not every volume covers on year
- layout and writing style changes mainly with autor

# Future directions

**Student projects**

(virtual exhibition on topic)

**Creating ground truth data**

Mapping all historic symbols

Trying to include every author

**Further evaluation of the table recognition**

Comparing various methods

Choose workflow for final indexing

**Developing GUI tool**

Specialized on handwritten tables

Focus on UCD and HCI

Current

Next

Future

# Conclusion

OCR engines do not perform optimally out-of-the-box

Layout segmentation can not be done completely automatically

No real guidelines exist when creating ground truth data

# Thank you for your attention!

# References

Glaser, R., Hagedorn, H.: Klimageschichte - Antworten auf die Verfinderlichkeit yon Wetter, Witterung und Klima? In: Naturwissenschaften 81, pp. 97-107 (1994)

Reul, C., Christ, D., Hartelt, A., Balbach, N., Wehner, M., Springmann, U., Wick, C., Grundig, C., Büttner, A., Puppe, F.: OCR4all – An Open-Source Tool Providing a (Semi-)Automatic OCR Workflow for Historical Printings, (2019)

# Figures

| Figure # | Creator |
|---|---|
| 1 | Reul, C., Christ, D., Hartelt, A., Balbach, N., Wehner, M., Springmann, U., Wick, C., Grundig, C., Büttner, A., Puppe, F.: OCR4all – An Open-Source Tool Providing a (Semi-)Automatic OCR Workflow for Historical Printings, (2019) |
| 2 | Dhali, M., Wit, JW. & Schomaker, L.: BiNet: Degraded-Manuscript Binarization in Diverse Document Textures and Layouts using Deep Encoder-Decoder Networks, (2019) |
| | Binmakhashen, G. & Mhamoud, S.: Document Layout Analysis: A Comprehensive Survey. ACM Comput. Surv. 52(6), (2020) |
| 3 | https://de.wikipedia.org/wiki/Coelestin_II._Steiglehner#/media/Datei:Coelestin_II._Steiglehner.JPG https://de.wikipedia.org/wiki/Placidus_Heinrich#/media/Datei:Placidus_Heinrich.jpg https://rzbvm050.uni-regensburg.de/meteorologie/schmoeger.htm |