# *Coptic OCR:*
# Even better models and improvements on user-friendliness

Eliese-Sophia Lincke (Humboldt-Universität zu Berlin)
eslincke@staff.hu-berlin.de

*Digital Coptic 3*, July 12-13, 2020 (online)

# Overview

- more fonts

- *Calamari:* new and better OCR software (as compared to OCRopus)

- *OCR4all*: a graphical user interface (GUI)

# *Coptic OCR*: Chronological overview since 2016

2016: presentation of OCRopus results at the Coptic Congress, Claremont

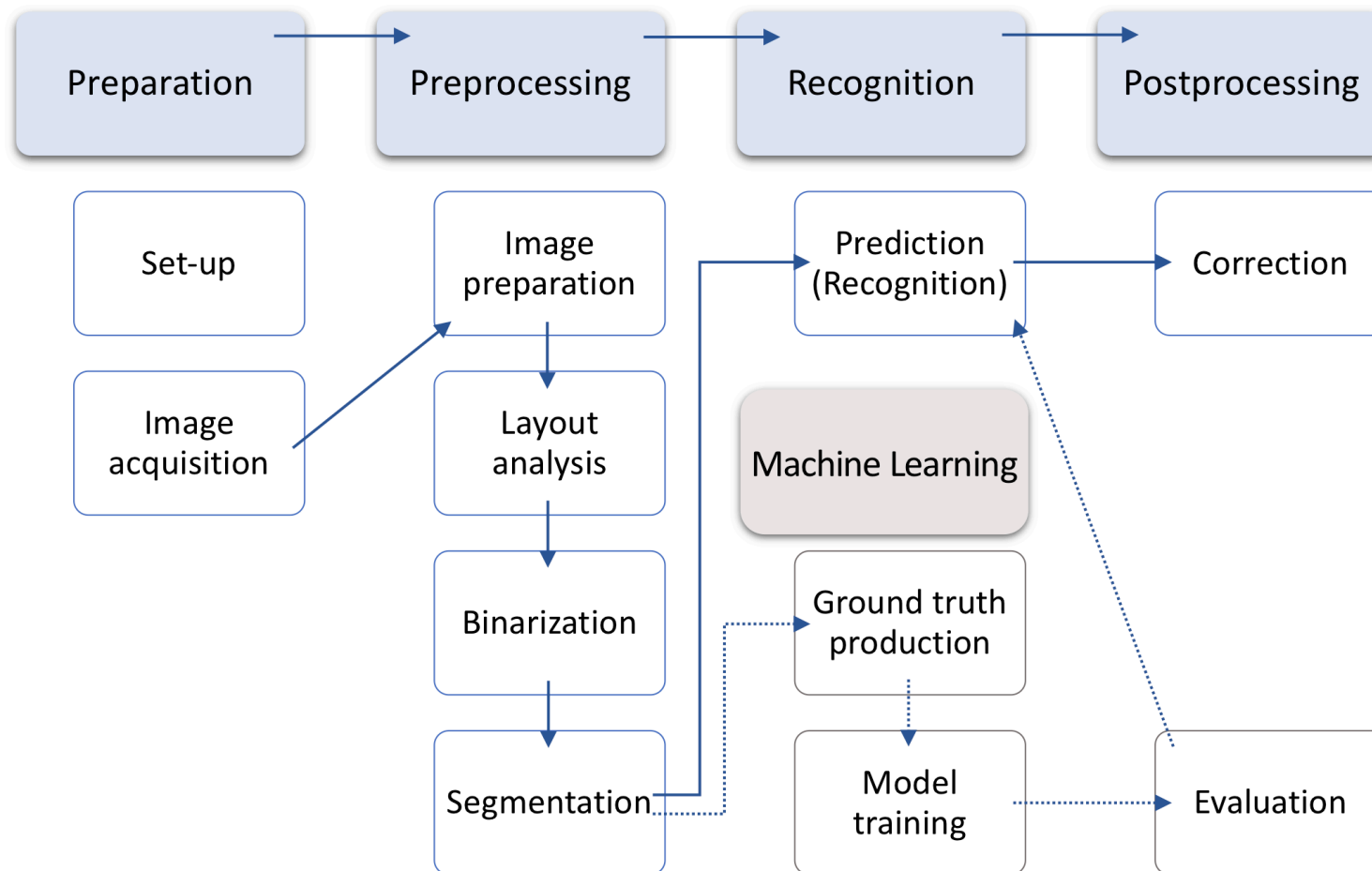2018: *Coptic OCR* becomes an official project

- 6 months PostDoc fellowship at the Göttingen Centre for Digital Humanities, CampusLab *Digitization and Computational Analytics* (DCA)
  PIs: Heike Behlmer, Marco Büchler, Camilla Di Biase-Dyson

2019: publication of the results of the fellowship and new tools

- presentation at the 3[rd] international conference on *Digital Access to Textual Cultural Heritage* (DATeCH) in Brussels (May)
- publication (open access) in the conference proceedings: Lincke, Bulert & Büchler (2019)
- testing new tool – Calamari
- training data of more fonts, model training (OCRopus, Calamari)
- presentation of new results in the *Berlin Digital Classicist Seminar* (November)

2020: testing new tool – OCR4all

# OCR workflow

Preparation → Preprocessing → Recognition → Postprocessing

Preparation:
- Set-up
- Image acquisition

Preprocessing:
- Image preparation
- Layout analysis
- Binarization
- Segmentation

Recognition:
- Prediction (Recognition)

Machine Learning:
- Ground truth production
- Model training
- Evaluation

Postprocessing:
- Correction

# Font overview



CSSC_1

ⲙ̄ⲡⲉⲛⲧⲁⲁⲛⲧⲱⲛⲓⲟⲥ ⲭⲟⲟϥ ⲉⲁϥϫⲱⲕ ⲉⲃⲟⲗ ⲉϫⲱϥ
ϩⲛⲟⲩϭⲉⲡⲏ.

87. Ⲛ̄ⲧⲉⲓϩⲉ ⲟⲛ ⲛⲉϥⲧⲥⲃⲱ *ⲛ̄ⲣⲱⲙⲉ ⲛⲓⲙ ⲉⲧⲣⲉⲩ-
ⲉⲣⲟⲕⲡⲉ· ϩⲉⲛⲕⲟⲟⲩⲉ ⲇⲉ ⲟⲛ ⲉⲁⲩⲉⲓ
ϩⲱⲥⲧⲉ[1] ⲉⲧⲣⲉⲩⲣⲡⲱⲃϣ[2] ⲛⲥⲉϫⲓ ...
ⲧⲏⲣⲓⲟⲛ ⲁⲩⲱ ⲛ̄ⲥⲉⲙⲁⲕⲁⲣⲓⲍⲉ
Ⲛⲉϥⲡⲣⲟϩⲓⲥⲧⲁ ⲇⲉ ⲉⲛⲉⲧⲟⲩϫⲓ ...
ⲛ̄ⲧⲉϩⲟⲓⲛⲉ ⲙⲉⲉⲩⲉ ϩⲙ̄ⲡⲉⲩϩⲏⲧ[4]
ϭⲟⲛⲥ̄ ⲉⲣⲟϥ. Ⲛⲉϥⲟ ⲇⲉ ⲛ̄ϩⲓⲕⲁⲛ
ⲛ̄ⲛⲟⲩⲟⲛ ⲛⲓⲙ ϩⲱⲥⲧⲉ ⲟⲩⲙⲏⲏϣ
ⲛ̄ϩⲁϩ ⲛ̄ⲛ̄ⲕⲁ ⲁϥⲧⲣⲉⲩⲕⲁⲁⲩ ...
Ⲛⲉϥⲟ ⲇⲉ ⲛ̄ⲛⲟⲩⲥⲁⲉⲓⲛ[7] ⲉⲁⲡⲛⲟⲩ ...
Ⲛⲓⲙ ⲅⲁⲣ ⲉⲛⲉϩ ⲡⲉⲛⲧⲁϥⲉⲓ ϣⲁ ...
ⲉⲙⲡⲉϥⲕⲧⲟϥ ⲉⲡⲉϥⲏⲓ ⲉϥⲣⲁϣⲉ·
ⲉⲧⲃⲉⲛⲉϥⲣⲱⲙⲉ ⲛ̄ⲧⲁⲩⲙⲟⲩ ⲙ̄ⲡⲉϥⲕ ...
ⲛⲓⲙ ⲛ̄ϩⲏⲕⲉ[10] ⲡⲉⲛⲧⲁϥⲉⲓ ϣⲁⲣⲟϥ
ϩⲛⲧⲉϥⲙⲛ̄ⲧϩⲏⲕⲉ · ⲛⲓⲙ ⲙ̄ⲙⲟⲛ ...
ⲛ̄ϩⲏⲧ ⲁⲩⲱ ⲛϥⲉⲓ ϣⲁⲣⲟϥ ⲉⲙⲡⲉϥ ...
ϣⲏⲙ ⲡⲉⲛⲧⲁϥⲉⲓ ⲉⲡⲧⲟⲟⲩ ⲁϥⲛⲁⲩ ...
ϣⲟⲟⲩⲉ ⲛ̄ⲧⲉⲩⲛⲟⲩ ⲉⲃⲟⲗ ϩⲛ̄ⲛⲁ ...

CSSC_2

ⲱⲛϩ̄ ⲉⲃⲟⲗ ⲙ̄ⲡⲧⲁⲓⲟ ⲛ̄ⲛⲉⲧⲉⲣϩⲏⲓⲃⲉ · ⲉϥϫⲱ ⲙ̄ⲙⲟⲥ : ϫⲉⲛⲁⲓ̈ⲁⲧⲟⲩ
ⲛ̄ⲛⲉⲧⲣ̄ϩⲏⲓ̈ⲃⲉ · ϫⲉⲛ̄ⲧⲟⲟⲩ ⲛⲉⲧⲟⲩⲛⲁⲥⲉⲡⲥⲱⲡⲟⲩ[1] :- 3. ⲛⲉⲧⲉⲣ-
ϩⲏⲓ̈ⲃⲉ · ϣⲁⲩⲉⲓ̈ⲛⲉ ⲛⲁⲩ ⲛⲟⲩⲙ ...
ⲙ̄ⲛⲟⲩⲙ̄ⲛ̄ⲧⲣⲉϥϯⲕⲁⲣⲡⲟⲥ ϩⲛ̄
ⲛⲉⲧⲉⲣϩⲏⲓ̈ⲃⲉ ϣⲁⲩϥⲓ ⲉϩⲣⲁⲓ̈ ...
ⲛ̄ⲧⲙⲛ̄ⲧϫⲁⲥⲓ̈ϩⲏⲧ :- ⲛⲉⲧⲉⲣ ...
ⲛ̄ⲛⲁⲩ ⲛⲓⲙ · ϫⲉⲕⲁⲥ ⲉⲩⲉ̄ⲥⲟ ...
ϩⲛ̄ⲥⲱⲙⲁ ⲛ̄ⲥⲉⲃⲱⲕ ϣⲁⲡϫⲟⲉ ...
ⲙⲁⲓ̈ⲣⲱⲙⲉ ⲙ̄ⲡⲉⲛⲭⲟⲉⲓⲥ · ⲁⲩ ...
ⲡⲁⲓ̈ ⲉⲃⲟⲗ ϩⲓ̈ⲧⲟⲟⲧϥ̄ · ⲉⲣⲉⲉ̄ ...
ⲛⲁⲅⲁⲑⲟⲥ · ⲙ̄ⲡⲉⲡⲛ̄ⲁ ⲉⲧ ...
ⲛ̄ϩⲟⲙⲟⲟⲩⲥⲓⲟⲛ · ⲧⲉⲛⲟⲩ · ⲁ ...
ⲛ̄ⲛⲁⲓ̈ⲱⲛ · ϩⲁⲙⲏⲛ :· -

☦ ⲡⲓⲱⲧ ⲙⲛ̄ⲡϣⲏⲣⲉ ⲙⲛ̄
ⲃ[ⲓⲡ]ⲣⲁϣ ⲉⲡⲓⲕⲩⲫⲁⲗⲓⲟⲛ
ⲉⲡⲁⲣⲭⲁⲅⲅⲉⲗⲟⲥ ⲙⲓⲭⲁⲏⲗ
ϩⲁⲙⲏⲛ ⲉⲥⲛ̄ϣⲱⲡⲓ[3] -

CSSC_3

ⲙⲡⲟⲩⲙⲟⲩ ⁿ· ⲁⲩϣⲱⲃⲉ ⲉⲛⲉⲩⲙⲁ ⲛⲉⲣⲙⲏ · ⁿ· ⲁⲩⲱ ⲡⲉϩⲣⲟⲟⲩ
ⲛⲧⲃⲁⲥⲁⲛⲟⲥ ⲛ̄ⲧⲡⲟⲗⲓⲥ ⲁϥⲃⲱⲕ ⲉϩⲣⲁⲓ ⲉⲧⲡⲉ .

VI. 1 ⲧⲕⲟⲓ̄ⲃⲱⲧⲟⲥ ⲇⲉ ⲛⲉ ...
ⲛⲥⲁϥϥ ⲛⲉⲃⲟⲧ · ⲁⲩⲱ ⲡⲉⲩⲕ ...
2 ⲛⲁⲗⲗⲟⲫⲩⲗⲟⲥ ⲇⲉ ⲁⲩⲙⲟⲩⲧⲉ ...
ⲁⲩⲱ ⲛⲉⲩⲣⲉϥⲙⲟⲩⲧⲉ · ⲉⲩⲭ ...
□ ⲛ̄ⲧⲕⲟⲓ̄ⲃⲱⲧⲟⲥ ⲙⲡⲭⲟⲉⲓⲥ · ⲧⲟ ...
ⲉϩⲣⲁⲓ ⲉⲡⲉⲥⲙⲁ ϩⲛ̄ ⲟⲩ̄ · 3 ⲡⲉⲭⲁ ...
ⲛ̄ⲧⲱⲧⲛ̄ ⲛⲧⲕⲟⲓⲃⲱⲧⲟⲥ ⲙⲡⲭⲟⲉ ...
ⲉⲃⲟⲗ ⲉⲥϣⲟⲩⲉⲓⲧ · ⲁⲗⲗⲁ ϩⲛ̄ ⲟⲩ ...
ⲁⲩⲱ ⲧⲉⲧⲛ̄ⲛⲁⲉⲙⲧⲟⲛ ⲛⲧⲉ ⲡⲭ ...
ⲙⲙⲟⲛ ⲛⲧⲃ̄ϭⲓϫ ⲙⲡⲭⲟⲉⲓⲥ ⲛⲁⲗ ...
ϫⲉ ⲟⲩ ⲡⲉ ⲡⲧⲱϣ ⲛⲧⲃⲁⲥⲁⲛⲟ ...
ⲛ̄ⲣⲉϥϣⲓⲛⲉ ⲛⲁⲩ · ϫⲉ ⲕⲁⲧⲁ ...
ⲫⲩⲗⲟⲥ · ☦ ⲛⲁⲥ ⲛ̄ϯⲟⲩ ⲙⲙⲁ ...
ϣⲱⲡⲉ ⲛϩⲏⲧⲧⲏⲩⲧⲛ̄ · ⲙⲛ ⲛⲉⲧ ...
ϩⲛ̄ⲕⲉⲡⲓⲛ ⲛⲛⲟⲩⲃ · ⲛⲧⲉⲧⲛ *...
ⲛⲁⲓ ⲉⲧⲧⲁⲕⲟ ⲙⲡⲉⲧⲛ̄ⲕⲁϩ · ⲛⲧ ...

Aubert

ⲉⲃⲟⲗ ⲁⲩⲉⲓ ⲉⲧⲉⲕⲕⲗⲏⲥⲓⲁ ⲛ̄ⲧⲕⲩⲣⲓⲁⲕⲏ ⲁⲩϩⲙⲟⲟⲥ ⲙ̄ ...
ⲟⲩⲣⲱⲙ ⲛ̄ⲟⲩⲱⲧ ⲛⲉⲣⲉ ⲡϩⲗ̄ⲗⲟ ⲇⲉ ϩⲛ̄ ⲧⲉⲩⲙⲏⲧⲉ · ...
ⲉⲧϩⲓϩⲟⲩⲛ ⲁⲩⲱ ⲛ̄ⲧⲉⲣⲟⲩⲕⲱ ⲉϩⲣⲁⲓ ⲙ̄ⲡⲟ(ⲥⲗ̄ⲍ ⲡ. ⲁ37 ⲁ)...
ⲉⲧⲟⲩⲁⲁⲃ ⲁϥⲟⲩⲱⲛϩ ⲉⲃⲟⲗ ⲙ̄ⲡϣⲟⲙⲛ̄ⲧ ⲙⲁⲩⲁⲁⲩ ⲛ̄ⲑⲉ ...
ⲛ̄ⲧⲉⲣⲉ ⲡⲉⲡⲣⲉⲥⲃⲩⲧⲉⲣⲟⲥ ⲥⲟⲟⲩⲧⲛ̄ ⲉⲃⲟⲗ ⲛ̄ⲧⲉϥϭⲓϫ ...
ⲉⲓⲥ ⲟⲩⲁⲅⲅⲉⲗⲟⲥ ⲁϥⲉⲓ ⲉⲃⲟⲗϩⲛ̄ ⲛ̄ⲙⲡⲏⲩⲉ ⲉⲟⲩⲛ̄ ⲟⲩϭ ...
ⲁϥϣⲱⲧ ⲙ̄ⲡⲕⲟⲩⲓ ⲛ̄ϣⲏⲣⲉ ⲁϥⲡⲱⲧ ⲙ̄ⲡⲉϥⲥⲛⲟϥ ⲉⲡ ...
ⲡⲉⲡⲣⲉⲥⲃⲩⲧⲉⲣⲟⲥ ⲇⲉ ⲉⲣ ⲡⲟⲉⲓⲕ ⲛ̄ⲅⲗⲁⲥⲙⲁ ⲕⲗⲁⲥⲙⲁ ⲛ ...
ⲡⲱϣ ⲙ̄ⲡϣⲏⲣⲉⲕⲟⲩⲓ ϣⲏⲙϣⲏⲙ · ⲁⲩⲱ ⲛ̄ⲧⲉⲣⲟⲩϯ ⲙ̄ⲡ ...
ⲛⲉⲧⲟⲩⲁⲁⲃ ⲁϥϫⲓ ⲛ̄ϭⲓ ⲡϩⲗ̄ⲗⲟ ⲛ̄ⲟⲩⲕⲗⲁⲥⲙⲁ ⲛⲁϥ ⲉϥⲡⲏ ...
ⲣⲉϥⲛⲁⲩ ⲁϥⲣ̄ϩⲟⲧⲉ ⲁϥϫⲓϣⲕⲁⲕ ⲉⲃⲟⲗ ϫⲉ ☦ⲡⲓⲥⲧⲉⲩⲉ ...
ⲡⲉ ⲡⲉⲕⲥⲱⲙⲁ ⲁⲩⲱ ⲡⲡⲟⲧⲏⲣⲓⲟⲛ ⲡⲉ ⲡⲉⲕⲥⲛⲟϥ · ⲁⲩ ...
ⲉⲧϩⲛ̄ ⲧⲉϥϭⲓϫ ϩⲟⲉⲓⲕ ⲕⲁⲧⲁ ⲡⲉⲟⲟⲩ ⲙ̄ⲡⲙⲩⲥⲧⲏⲣⲓⲟⲛ ...
(ⲡ. ⲁ37 ⲃ) ⲁⲩⲱ ⲁϥϫⲓ ⲉϥⲉⲩⲭⲁⲣⲓⲥⲧⲓ ⲙ̄ⲡⲭⲟⲉⲓⲥ · ⲡⲉ ...
ϫⲉ ⲡⲛⲟⲩⲧⲉ ⲥⲟⲟⲩⲛ ⲛ̄ⲧⲉⲫⲩⲥⲓⲥ ⲛ̄ⲛ̄ⲣⲱⲙⲉ ϫⲉ ⲙⲛ̄ⲟ ...
ⲁⲃ ⲉϥⲟⲩⲱⲧ ⲉⲧⲃⲉ ⲡⲁⲓ ϣⲁϥⲧⲣⲉ ⲡⲉϥⲥⲱⲙⲁ ϣⲱⲡⲉ ⲙ̄ ...
ⲛ̄ⲏⲣⲡ̄ ⲛ̄ⲛⲉⲧϫⲓ ⲙⲙⲟϥ ϩⲛ̄ ⲟⲩⲡⲓⲥⲧⲓⲥ ⲁⲩⲱ ⲁⲩϣⲛ̄ϩⲙⲟⲧ ...
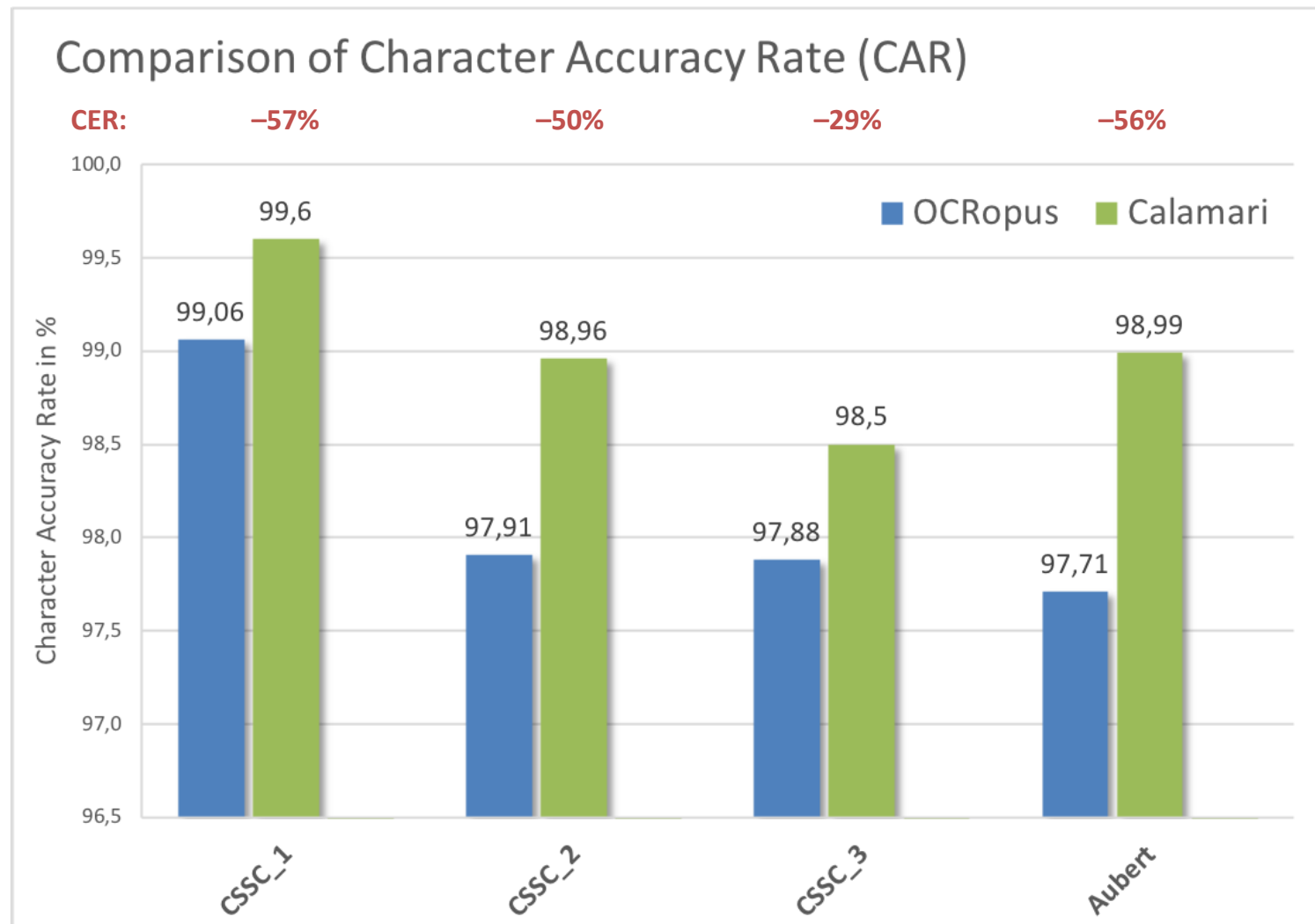ⲡⲉⲛⲧⲁϥϣⲱⲡⲉ ϫⲉ ⲙ̄ⲡⲉϥⲕⲁ ⲡϩⲗ̄ⲗⲟ ⲛ̄ⲣⲱⲙⲉ ⲉϯⲟⲥⲉ ⲙ̄ ...

# Criteria for training (and test) data selection

- several printed text editions per font

- inhomogeneous scans

- Sahidic and Bohairic dialects  (if available)

- several layouts (one column vs. two columns)

- no normalization, simplification or cleaning

- comprehensive character set for each font
(lower case, upper case, supralinear strokes, various punctuation marks, footnote signs and other philological markup etc.)

# Best Models

# Best Models

## Best Models

*Calamari* uses *Cross Fold Training* and *Confidence Voting* to improve its predictions
cf. Reul, Springmann, Wick & Puppe (2018: Fig. 1; Table II)



|      | c       | e       |
|------|---------|---------|
| M1   | 66.83%  | 38.40%  |
| M2   | 93.27%  | 19.77%  |
| M3   | -       | 99.91%  |
| M4   | 7.56%   | 98.02%  |
| M5   | 90.31%  | 50.07%  |
| $\sum$ Rec   | **250.41%** | 197.93%  |
| $+\sum$ Alt  | 257.97%  | **306.17%** |

# OCR4all: Workflow overview

cf. Reul et al. (2019)

- Project at the University of Würzburg
- originally designed for historical (German) prints from the Early Modern Period
- can be used for Coptic printed texts too

https://www.uni-wuerzburg.de/en/zpd/ocr4all/

# OCR4all: Workflow overview



https://www.uni-wuerzburg.de/en/zpd/ocr4all/

# OCR4all: Workflow overview

cf. Reul et al. (2019)

- Project at the University of Würzburg
- originally designed for historical (German) prints from the Early Modern Period
- can be used for Coptic printed texts too
- **Interface** that runs in a browser (using Docker) or in a virtual machine (VirtualBox)



https://www.uni-wuerzburg.de/en/zpd/ocr4all/

# OCR4all: Workflow overview

# OCR4all: Model selection

# OCR4all: Recognition (Calamari)



Calamari:
Wick et al. (to appear 2020)

# OCR4all: Postcorrection (Larex)



Larex:
Reul et al. (2017)

# OCR4all: Output (xml)

```xml
<PcGts xmlns="http://schema.primaresearch.org/PAGE/gts/pagecontent/2017-07-15" xmlns:xsi="
http://www.w3.org/2001/XMLSchema-instance" xsi:schemaLocation="http://schema.primaresearch.org/PAGE/gts/
pagecontent/2017-07-15 http://schema.primaresearch.org/PAGE/gts/pagecontent/2017-07-15/pagecontent.xsd">
    <Metadata>
        <Creator>User123</Creator>
        <Created>2020-07-11T23:46:38</Created>
        <LastChange>2020-07-11T23:46:38</LastChange>
    </Metadata>
<Page imageFilename="0001.png" imageHeight="4852" imageWidth="3944">
    <TextRegion id="r0" type="paragraph"><Coords points="1,1 3942,1 3942,4850 1,4850"/>
    <TextLine id="r0_l001">
        <Coords points="400,120 3569,120 3569,207 400,207"/>
        <TextEquiv index="1">
            <Unicode>ⲉⲑⲃⲃⲓⲉ ⲡⲉϥⲙⲉⲉⲩⲉ ⲡⲉⲭⲁϥ ⲛⲁϥ ϫⲉ ⲃⲱⲕ ϣⲁ ⲙⲉϣⲉⲛⲓⲙ ⲛ̄ⲁⲣⲭⲓⲙⲁⲇⲣⲓⲧⲏⲥ</Unicode></TextEquiv></TextLine>
    <TextLine id="r0_l002">
        <Coords points="398,245 3574,245 3574,334 398,334"/>
        <TextEquiv index="1">
            <Unicode>ⲁⲩⲱ ⲡⲉⲧⲉϥⲛⲁϫⲟⲟϥ ⲛⲁⲕ ⲁⲣⲓϥ · ⲁ ⲡⲛⲟⲩⲧⲉ ⲇⲉ ϭⲱⲗⲛ̄ ⲉⲃⲟⲗ ⲙ̄ⲡⲓⲁⲣⲭⲓ-</Unicode></TextEquiv></TextLine>
    <TextLine id="r0_l003">
        <Coords points="403,369 3569,369 3569,459 403,459"/>
        <TextEquiv index="1">
            <Unicode>ⲙⲁⲇⲣⲓⲧⲏⲥ ⲉϥϫⲱ ⲙ̄ⲙⲟⲥ ϫⲉ ⲉⲓⲥ ⲙⲉϣⲉⲛⲓⲙ ⲛ̄ⲁⲛⲁⲭⲱⲣⲓⲧⲏⲥ ⲛⲏⲩ ϣⲁⲣⲟⲕ</Unicode></TextEquiv></TextLine>
    <TextLine id="r0_l004">
        <Coords points="398,486 3572,486 3572,585 398,585"/>
        <TextEquiv index="1">
            <Unicode>ϯ ⲟⲩϥⲣⲁⲅⲉⲗⲗⲓⲟⲛ ⲛⲁϥ ⲛ̄ϯⲣⲉϥⲙⲟⲟⲛⲉ ⲛ̄ⲛ̄ⲣⲓⲣ · ⲁϥⲉⲓ ⲇⲉ ⲛ̄ϭⲓ ⲡⲉⲗⲗⲟ ⲁϥ-</Unicode></TextEquiv></TextLine>
    <TextLine id="r0_l005">
        <Coords points="401,617 3568,617 3568,707 401,707"/>
        <TextEquiv index="1">
            <Unicode>ⲕⲱⲗⲏ̄ ⲉⲡⲣⲟ ⲁⲩⲱ ⲁϥⲃⲱⲕ ⲉϩⲟⲩⲛ ϣⲁ ⲡⲁⲡⲉ ⲛ̄ⲧⲥⲟⲟⲩⲉϩ̄ ⲁⲩⲁⲥⲡⲁⲍⲉ ⲛ̄ⲛⲉⲩⲉ-</Unicode></TextEquiv></TextLine>
    <TextLine id="r0_l006">
        <Coords points="397,736 3564,736 3564,833 397,833"/>
        <TextEquiv index="1">
            <Unicode>ⲣⲏⲩ ⲁⲩϩⲙⲟⲟⲥ ⲁⲩⲱ ⲡⲉⲭⲁϥ ⲛ̄ϭⲓ ⲡⲁⲛⲁⲭⲱⲣⲓⲧⲏⲥ ϫⲉ ⲟⲩ ⲡⲉϥⲛⲁⲁⲁϥ ϫⲉ</Unicode></TextEquiv></TextLine>
    <TextLine id="r0_l007">
```

# Hopes (Plans?) for the future

- a Digital Coptic repositorium
- an infrastructure in which OCR4all can be run on a (remote) server – no need for local installation

- better documentation for Coptic OCR
- hands-on workshops (like the one planned for this years' Coptic Congress …): offline and online

## Coptic OCR data

Data repository for the Coptic OCR project ("working GitLab repository")

(At the moment, the repository, is accessible without registration. This may, however, change in the future due to copyright concerns etc. In that case, please, contact us.)

DOI: 21.11101/0000-0007-C9D1-A

https://vcs.etrap.eu/Coptic-OCR/datasets

eslincke@staff.hu-berlin.de

**Thank you …**

… for having contributed to and/or supported *Coptic OCR*:

Heike Behlmer

Marco Büchler

Kirill Bulert

Camilla Di Biase-Dyson

Frank Feder

Florian Finck

Jürgen Knauth

So Miyagawa

Tobias Paul

Christian Reul

Malte Rosenau

Caroline Sporleder

Uwe Springmann

Ronnie Vuine

# References

Bourcellier, Laurent. 2006. *Création d'une typographie numérique copte adaptée aux usages éditoriaux*, Diplôme supérieur d'arts appliqués, arts et techniques de communication option création typographique, Livret d'accompagnement, Ecole Estienne, Paris.

Lincke, Eliese-Sophia, Kirill Bulert and Marco Büchler. 2019. Optical Character Recognition for Coptic fonts: A multi-source approach for scholarly editions, in: *DATeCH2019 – Proceedings of the 3rd International Conference on Digital Access to Textual Cultural Heritage,* 87-91.
DOI: 10.1145/3322905.3322931

Reul, Christian, Uwe Springmann, Christoph Wick and Frank Puppe. 2018. Improving OCR Accuracy on Early Printed Books by Utilizing Cross Fold Training and Voting, in: *Proceedings of the 13th IAPR International Workshop on Document Analysis Systems (DAS)*, Vienna, 24-27 April 2018, IEEE: 423-428.
DOI: 10.1109/DAS.2018.30

Reul, Christian, Dennis Christ, Alexander Hartelt, Nico Balbach, Maximilian Wehner, Uwe Springmann, Christoph Wick, Christine Grundig, Andreas Büttner and Frank Puppe. 2019. OCR4all---An Open-Source Tool Providing a (Semi-)Automatic OCR Workflow for Historical Printings, in: *Applied Sciences* 9(22), No. 4853.
DOI: 10.3390/app9224853

Reul, Christian, Uwe Springmann and Frank Puppe. 2017. LAREX: A semi-automatic open-source Tool for Layout Analysis and Region Extraction on Early Printed Books, in: *Proceedings of the 2nd International Conference on Digital Access to Textual Cultural Heritage – DATeCH2017*, ACM: 137-142.
DOI: 10.1145/3078081.3078097

Wick, Christoph, Christian Reul and Frank Puppe. to appear (2020). Calamari - A High-Performance Tensorflow-based Deep Learning Package for Optical Character Recognition, in: *Digital Humanities Quarterly* 14(2).