

# Utilization of Common OCR Tools for Typeset Coptic Texts



\*KELLIA\*

EN LY ISBOR



So Miyagawa Kirill Bulert Marco Büchler



#### Coptic

- The final stage of the Ancient Egyptian language used in Egypt from ca. the third century
- A new writing system based on the Greek alphabet and several letters from Demotic (a language stage and writing system used in Egypt from ~ 700 BCE)
- A language transmitted in several regional forms (dialects) with a large production of manuscripts in Sahidic Coptic, the dialect at the basis of our OCR work

#### **Coptic alphabet**

- Ca. 30 letters.
- Several diacritics such as tremas, circumflexes, supralinear strokes etc.
- Several punctuation marks such as dots, commas, and colons
- Editorial marks in editions

#### Why is Coptic OCR needed?

- OCR for Coptic is not well-developed.
- Almost all the Coptic texts in past publications were not OCRed.
- OCR for Coptic is needed by many DH projects in Coptic.
- There is a small amount of human power in Coptology compared with the large amount of unOCRed Coptic editions.

#### Coptic DH projects (selected)

- SFB 1136 (Göttingen)
  - Creates a text corpus of selected monastic works in Coptic
- Digital Edition of the Coptic Old Testament (Göttingen)
  - Creates a digital edition of the Coptic translation of the Old **Testament**
- Coptic SCRIPTORIUM (Georgetown/Pacific)
  - > Creates a linguistically annotated Coptic corpus

#### **Existing Coptic OCR**

Tesseract (developed by Ray Smith) for Coptic, trained by Moheb Mekhaiel (http://www.moheb.de/ocr.html)

#### **New method: Ocropy**

- Python-based OCR package
- Using recurrent neural networks
- Originally developed by Thomas Breuel
- Available at https://github.com/tmbdev/ocropy

George Ghaly, Marwan Kilani, Rebecca Krawiec, and Nicola Denzey Lewis.

Trained for Coptic by our group and our collaborator Eliese-Sophia Lincke (Berlin)

## Upper: the Coptic alphabet

Lower: diacritics and punctuation

авгде (5) хнөікх мизопрстуфхү 4 9 x (3) 8 (6) b 60 a

OY MIN IN MINT HI, . ':

Optical Character Recognition for Coptic (OCR)

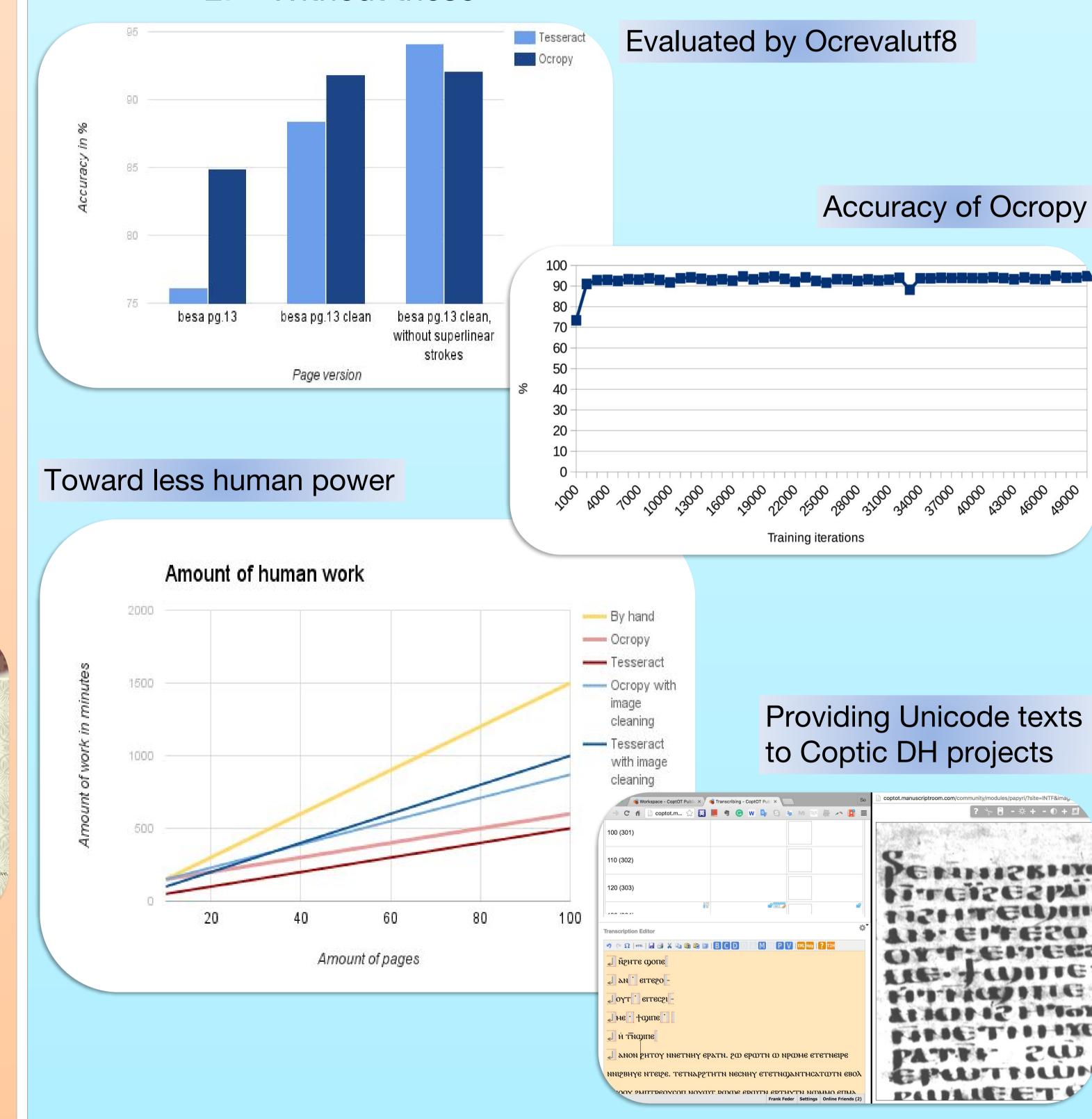
### Coptic diacritics, punctuation, and editorial signs

- Create problems with Tesseract
- Better processed by Ocropy

летнаноуоу. Мпефооу накім ан 2Мпечнї 14 : апа внса [Fragment 35] A DENUNCIATION OF AN ERRING NUN ....  $M_1A$  .....  $A_1YW$   $M_1$  ....  $O \cdot H$  NIM  $\Pi$   $E_1TN[AA]$  W-AZO]M EZPAÏ E[X]W· H NIM ПE[TN]AKTOY NE [E]YEIPHNH•

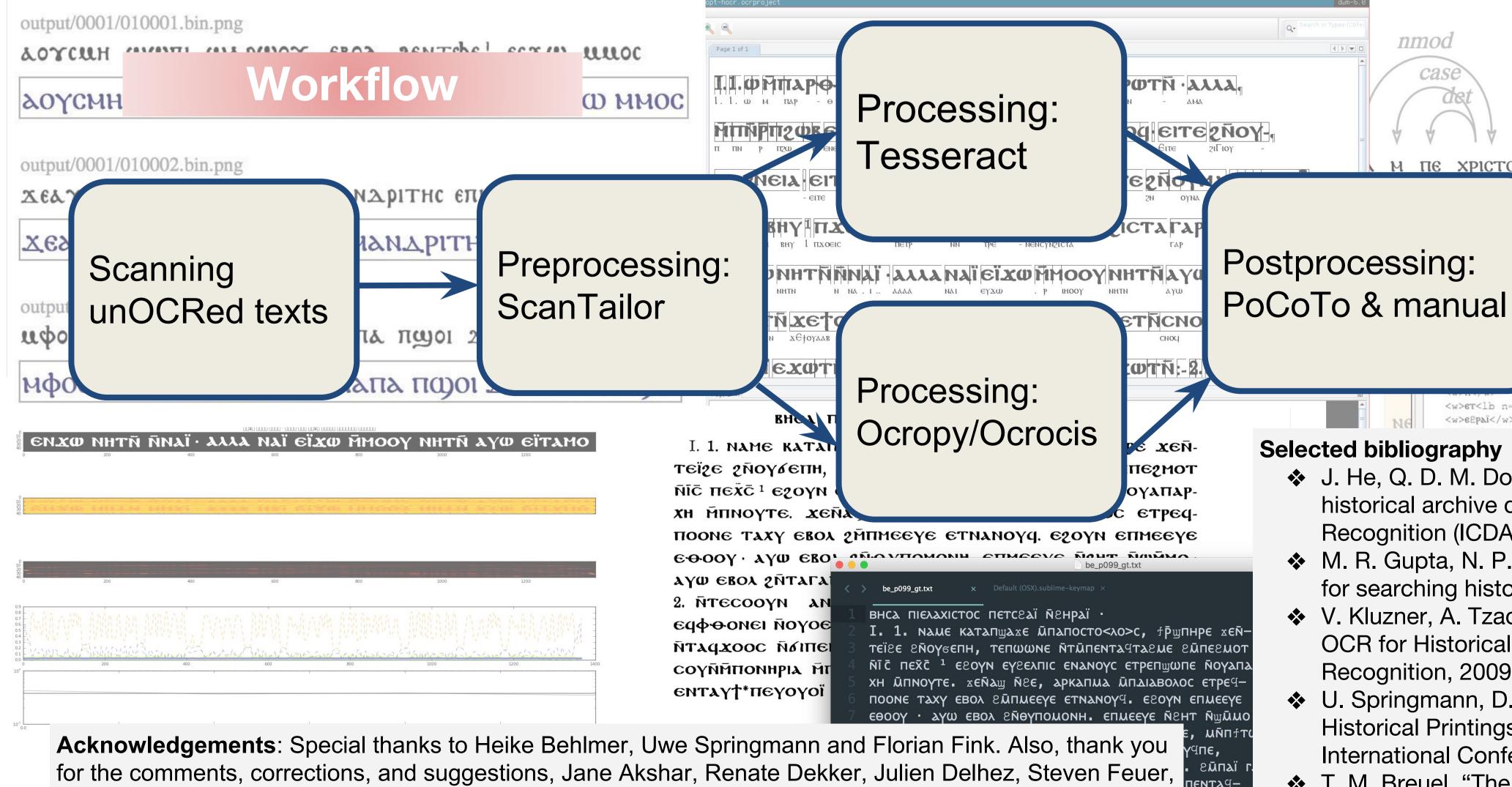
#### **Comparison between Tesseract and Ocropy**

- Our test case: 1956 edition of the works of Besa, a fifth-century abbot in Upper Egypt, written in the Sahidic dialect of Coptic
- Mekhaiel's model: aimed at the Bohairic dialect (different type face, but containing all the diacritics and letters of Sahidic)
  - With diacritics, punctuation, and editorial signs
  - Without these



ccomp

nsubi



Selected bibliography

nmod

case

❖ J. He, Q. D. M. Do, A. C. Downton, and J. H. Kim, "A comparison of binarization methods for historical archive documents," in Eighth International Conference on Document Analysis and Recognition (ICDAR'05), 2005, p. 538-542 Vol. 1.

nmod

THUNDAUT

**Providing Unicode Coptic texts** 

to the Coptological community

Enter Coptic text in UTF-8 (XML markup is also allowed, 10,000 characters max).

COPTIC

**SCRIPTORIUM** 

- M. R. Gupta, N. P. Jacobson, and E. K. Garcia, "{OCR} binarization and image pre-processing for searching historical documents," Pattern Recognit., vol. 40, no. 2, pp. 389–397, 2007.
- ❖ V. Kluzner, A. Tzadok, Y. Shimony, E. Walach, and A. Antonacopoulos, "Word-Based Adaptiv OCR for Historical Books," in 2009 10th International Conference on Document Analysis and Recognition, 2009, pp. 501-505.
- U. Springmann, D. Najock, H. Morgenroth, H. Schmid, A. Gotscharek, and F. Fink, "OCR of Historical Printings of Latin Texts: Problems, Prospects, Progress," in Proceedings of the First International Conference on Digital Access to Textual Cultural Heritage, 2014, pp. 71–75.
- ❖ T. M. Breuel, "The OCRopus open source OCR system," Proc. SPIE 6815, Doc. Recognit. Retr. XV, 2008.