# Optical Character Recognition with a Neural Network Model for Coptic

#### Kirill Bulert So Miyagawa Marco Buechler December 8, 2017 DH2017 Montreal, Canada Virtual Short Paper







# Coptic

- The final stage of the Ancient Egyptian language (third century)
- Multiple dialects (most important: Bohairic, Sahidic)
- Many manuscripts in Sahidic Coptic



## Alphabet

- Simple: Ca. 30 unique characters
- No upper/lower case in historic texts
- Diacritics add some complexity, but not much
- Based on Greek alphabet

**ΔΒΓΔϾ(5)ΖΗΘΙΚ**λ мизопрстуфхү ω ω ч (b) ε (ε) х с †  $\hat{OY}$  MN  $\bar{N}$  MNT  $\hat{H}$   $\ddot{I}$  , . ` :

- SFB 1136 (Goettingen) https://www.uni-goettingen.de/de/sfb-1136/521113.html
- Digital Edition of the Coptic old testament (Goettingen) http://coptot.manuscriptroom.com/
- Coptic SCRIPTORIUM (Georgetown/Pacific) http://copticscriptorium.org/



## Tesseract vs Ocropy, free OCR frameworks

#### Tesseract

- Font based learning
- Single characters are decomposed
- Decomposed parts are matched against given text

#### Input

• Fonts

#### Problem

• Few Coptic fonts available, not all historic variations covered

#### Ocropy/Ocropus

- Neural nets with online learning
- Model accuracy proportional to ground truth size
- But, size of ground truth is limited

#### Input

 Ground truth and corresponding images

#### Problem

Limited data for learning

## Tesseract vs Ocropy, free OCR frameworks

#### Tesseract

- Font based learning
- Single characters are decomposed
- Decomposed parts are matched against given text

#### Input

• Fonts

#### Problem

• Few Coptic fonts available, not all historic variations covered

#### Ocropy/Ocropus

- Neural nets with online learning
- Model accuracy proportional to ground truth size
- But, size of ground truth is limited

#### Input

• Ground truth and corresponding images

#### Problem

Limited data for learning

# **Pre-Processing**

- Training: Data  $\rightarrow$  training  $\rightarrow$  model
- Garbage in, garbage out
- Cleaner images improve results tremendously
- Dust and stains can be cleaned algorithmically, but ...
- What if text itself is noise the problem?

#### **Original Page**

#### апа внса

[Fragment 35] A DENUNCIATION OF AN ERRING NUN

 $\begin{array}{cccc} & & & & & \\ & & & & & \\ \hline \begin{array}{c} M_{J}\lambda & & & & \\ & & & & \\ \hline \end{array} \begin{array}{c} M_{J}\lambda & & \\ \end{array} \end{array} \begin{array}{c} M_{J}\lambda & & \\ \end{array} \begin{array}{c} M_{J}\lambda & & \\ \end{array} \begin{array}{c} M_{J}\lambda & & \\ \end{array} \end{array}$ 

• Flawless, almost ...

#### **Original Page**

20 ΝΤΕΦΑΙΚΑΙΟCΥΝΗ ΜΝΤΕΦΠΙCΤΙC <sup>11</sup> · 5. ΝΤΦΤΝ ΔΕ ΗΜΕΡΑΤΕ 20 ΝCNHY ΕΤΡ2ΟΤΕ 2ΗΤΫ ΗΠΝΟΥΤΕ ΦΛΗΛ ΤΗΡΤΝ ΧΕΚΑC ΕΡΕΠΝΟΥΤΕ ΝΑ2ΑΡΕ2 ΕΡΟΝ ΕΒΟΛ 2ΗΠΕΦΟΟΥ ΝΙΗ · ΑΥΦ ΝΫΤΟΥ-ΧΟΝ ΕΒΟΛ 2ΗΠΕΊΑΙΦΝ ΗΠΟΝΗΡΟΝ <sup>12</sup>. ΗΝΝΚΙΝΑΥΝΟC ΝΝΕCΝΗΥ ΝΝΟΥΧ <sup>13</sup> : 6. ΟΥΝδΟΜ ΜΕΝ ΗΜΟΝ ΕΤΑΜΦΤΝ ΕΦΕ ΤΗΡΖ ΕΝΤΑΥΑΑC ΝδΙΝΕΤΗΜΑΥ. ΑΛΛΑ ΤΕΝΟΥ ΡΦ ΕΙC2HHTE ΑΝδΦ ΕΝΚΦ ΗΠΜΑ ΗΠΝΟΥΤΕ · ΝΕΤΗΜΑΥ ΔΕ 2ΦΟΥ, ΕΦΦΠΕ ΠΕΥ2HT ΝΑΗΤΟΝ · Η CENAXITOT Ν2HT · CEPΦΦΕ ΗΝΠΝΟΥΤΕ · 7. ΑΝΟΚ ΜΕΝ ΚΑΤΑΤΑΜΝΤΤΑΛΑΙΠΦΡΟC. ΝΤΝΑΥ, ΑΝ, ΕΠΕΦΟΟΥ ΕΑΪΕΙΡΕ ΝΡΦΗΕ · ΑΛΛΑ QCH2 ΧΕΠΕΤΝΑΥ ΑΝ 2ΗΠΕΦΗΙ<sup>14</sup> :

#### апа внса

[Fragment 35] A DENUNCIATION OF AN ERRING NUN

• Easily removable with ScanTailor

#### Foreign language

<sup>20</sup> ΝτεφαικλιοςγΝΗ ΜΝΤεφπιςτις <sup>11</sup> · 5. Ντωτή Δε Ημερλτε <sup>20</sup> ΝζεμικλιοςγΝΗ ΜΝΤεφπιςτις <sup>11</sup> · 5. Ντωτή Δε Ημερλτε <sup>20</sup> Νζεμικ, <sup>20</sup> Νζε

#### апа внса

[Fragment 35] A DENUNCIATION OF AN ERRING NUN

No multilingual OCR models for Coptic

#### Annotations

<sup>20</sup> ΝτεφαικλιοςγΝΗ ΜΝΤεφπιςτις <sup>11</sup> · 5. Ντωτή Δε Ημερατε <sup>20</sup> ΝζεμικαιοςγΝΗ ΜΝΤεφπιςτις <sup>11</sup> · 5. Ντωτή Δε Ημερατε <sup>20</sup> Νζενμα ετρέζοτε 2μτζι Ηπινογτε ωληλ τηρτή χεκας εφεπινογτε Ναζαρές έρον εβολ 2μπεφοογ ΝΙΗ · Αγω Νζτογχον εβολ 2μπειαών Ηπονήρον <sup>12</sup>. ΗΝΝΚΙναγνος Ννεςνηγ Ννογχ<sup>13</sup> : 6. Ογίλομ μεν Ημον εταμωτή εφε τηρς ενταγαας Νδινετήμας. Αλλα τένογ ρω εις2μητε ανόω ενκω Ηπινα Ηπινογτε · Νετήμας Δε 2ωογ, εφωπε πεγ2μτ Ναήτον · η εεναχιτώτ Νζητ · ζερωψε μνπινογτε · 7, ανοκ μεν καταταμήτταλαιπωρος. Νταγ, αν, επεφοογ εαιείρε Νρωμε · Αλλα ζεμεταλά τος Νακιμ αν 2μπεφαί <sup>14</sup> : <sup>30</sup>

#### апа внса

[Fragment 35] A DENUNCIATION OF AN ERRING NUN

 $[\lambda_20]M \in \mathbb{Z}P\lambda^{"} \in \mathbb{Z}^{"} \oplus \mathbb{Z}^{"} \cap \mathbb{Z}^{"}$ 

- Special characters might not be part of any model (, ,)
- Not all annotations wanted

#### Language specific variations

20 ΝΤΕΦΑΙΚΑΙΟCΥΝΗ ΜΝΤΕΦΠΙCΤΙC <sup>11</sup> · 5. ΝΤΦΤΝ ΔΕ ΜΜΕΡΑΤΕ 20 ΝCNHY ΕΤΡ2ΟΤΕ 2ΗΤΫ ΜΠΝΟΥΤΕ ΦΛΗΛ ΤΗΡΤΝ ΧΕΚΑC ΕΡΕΠΝΟΥΤΕ ΝΑ2ΑΡΕ2 ΕΡΟΝ ΕΒΟΛ 2ΜΠΕΦΟΟΥ ΝΙΗ · ΑΥΦ ΝΫΤΟΥ-ΧΟΝ ΕΒΟΛ 2ΜΠΕΊΔΙΦΝ ΜΠΟΝΗΡΟΝ <sup>12</sup>. ΜΝΝΚΙΝΑΥΝΟC ΝΝΕCΝΗΥ ΝΝΟΥΧ <sup>13</sup> : 6. ΟΥΝΤΟΜ ΜΕΝ ΜΜΟΝ ΕΤΑΜΦΤΝ ΕΦΕ ΤΗΡΖ ΕΝΤΑΥΔΑC ΝΤΙΝΕΤΗΜΑΥ. ΑΛΛΑ ΤΕΝΟΥ ΡΦ ΕΙC2HHTE ΑΝΤΦ ΕΝΚΦ ΜΠΜΑ ΜΠΝΟΥΤΕ · ΝΕΤΗΜΑΥ ΔΕ 2ΦΟΥ, ΕΦΦΠΕ ΠΕΥ2HT ΝΑΜΤΟΝ · Η CENAXITOT Ν2HT · CEPUDE ΜΝΠΝΟΥΤΕ · 7. ΑΝΟΚ ΜΕΝ ΚΑΤΑΤΑΜΝΤΤΑΛΑΙΠΦΡΟC. ΝΤΝΑΥ, ΑΝ, ΕΠΕΦΟΟΥ ΕΔΙΈΙΡΕ ΝΡΦΙΕ · ΑΛΛΑ QCH2 ΧΕΠΕΤΝΑΥ ΑΝ 2ΜΠΕΦΟΥ ΑΙΤΗΛΟΥΟΥ.
30 ΠΕΤΝΑΝΟΥΟΥ. ΜΠΕΦΟΟΥ ΝΑΚΙΜ ΑΝ 2ΜΠΕΦΗΙ<sup>14</sup> : 30

#### апа внса

nt 35] A DENUNCIATION OF AN ERRING NUN  $(\dots, [a] Y ω [n], \dots, 0 \cdot H [nim π][ε] τn[aa] φ$ paï ε[x] ω · H nim πε[τη] aktoq nε [ε] γειρηνη ·

Might also not be included in a model

- Coptic models created by Moheb for Tesseract (2013)
- Trained with several Coptic fonts, no non-Coptic letters support
- No support for diacritics
- Non-Coptic letters get replaced with similar Coptic letters

#### The good, the bad, and the problematic

20 ΝΤΕ ΑΙΚΑΙΟ CYNΗ ΜΝΤΕ (ΠΙCTIC <sup>11</sup> · 5. ΝΤΦΤΝ ΔΕ ΜΗΕΡΑΤΕ 20 ΝCNHY ΕΤΡ2ΟΤΕ 2ΗΤΙ ΜΠΝΟΥΤΕ ΦΛΗΛ ΤΗΡΤΝ ΧΕΚΑΟ ΕΡΕΠΝΟΥΤΕ ΝΑ2ΑΡΕ2 ΕΡΟΝ ΕΒΟΛ 2ΗΠΕΘΟΟΥ ΝΙΗ · ΑΥΦ ΝΙΤΟΥ-ΧΟΝ ΕΒΟΛ 2ΗΠΕΊΑΙΦΝ ΗΠΟΝΗΡΟΝ <sup>12</sup>. ΜΝΝΚΙΝΑΥΝΟΟ ΝΝΕΟΝΗΥ ΝΝΟΥΧ <sup>13</sup> · 6. ΟΥΝΤΟΜ ΜΕΝ ΗΜΟΝ ΕΤΑΜΦΤΝ ΕΘΕ ΤΗΡΕ ΕΝΤΑΥΑΑΕ ΝΤΙΝΕΤΗΜΑΥ. ΑΛΛΑ ΤΕΝΟΥ ΡΦ ΕΙC2HHTE ΑΝΤΦ ΕΝΚΦ ΗΠΗΑ ΗΠΝΟΥΤΕ · ΝΕΤΗΜΑΥ ΔΕ 2ΦΟΥ, ΕΦΦΠΕ ΠΕΥ2ΗΤ ΝΑΗΤΟΝ · Η CENAXITOT Ν2ΗΤ · CEPΦΦΕ ΗΝΠΝΟΥΤΕ · 7. ΑΝΟΚ ΜΕΝ ΚΑΤΑΤΑΜΝΤΤΑΛΑΙΠΦΡΟΕ. ΝΤΝΑΥ, ΑΝ, ΕΠΕΘΟΟΥ ΕΑΪΕΙΡΕ ΝΡΦΗΕ · ΑΛΛΑ QCH2 ΧΕΠΕΤΝΑΥ ΑΝ 2ΗΠΕΘΗΙ<sup>14</sup> : 30

#### апа внса

[Fragment 35] A DENUNCIATION OF AN ERRING NUN

Even clean scans still contain noise

# Results

## Without line numbers



- Without time consuming pre-processing
- Ocropy model trained on 10 pages more accurate
- Non-multilingual Tesseract model less accurate

## Without line numbers



#### Accuracy

#### Input

йршмб алла цсне жепе петнаночоч, йпбфооч с

[Fragment 35] A DENUI

.... мја ..... јајуш ћршне алла цснг жепе петнаноуоу йпефооу

- Without non-Coptic letters
- Difference results mostly from diacritics

Μλ λΥΨ

٤

## Without line numbers



#### Accuracy

#### Input

приме алла цсне жепе петнаночоч мпефооч

#### м а а уш Npwme алла qchz жепе Петнаноуоу мпефооу

MA

- Diacritics removed
- Pure Coptic Tesseract model outperforms Ocropys mixed model

16

2

٤

a yw

#### Workload comparison



• Utilising OCR decreases human workload

- Utilisation of OCR beneficial for most clean documents
- **Tesseract** best for monolingual documents with limited fonts and font variations
- Ocropy best for large documents with multiple languages

- Approached by Google for collaboration on OCR for Coptic
- Create data set for Coptic OCR testing
- Transition from typeset to handwritten Coptic texts
- Combination of different models

# Thank you. Questions? Get our Ocropy models at



Presentation Kirill Bulert, So Miyagawa

#### Team (in alphabetical order)

Kirill Bulert, Marco Büchler, So Miyagawa.



- 1. Uwe Springmanns OCR Workshop http://www.cis.uni-muenchen.de/ocrworkshop/program.html
- 2. Scantailor for pre-processing http://scantailor.org/
- Ocropy/Ocropus https://github.com/tmbdev/ocropy
- 4. Kraken an Ocropy fork http://kraken.re/
- 5. Tesseract OCR https://github.com/tesseract-ocr/
- 6. Moheb's Coptic Pages http://www.moheb.de/ocr.html

- 1. SFB 1136 (Goettingen) https://www.uni-goettingen.de/de/sfb-1136/521113.html
- 2. Digital Edition of the Coptic old testament (Goettingen) http://coptot.manuscriptroom.com/
- Coptic SCRIPTORIUM (Georgetown/Pacific) http://copticscriptorium.org/

The LaTeX theme this presentation is based on is licensed under a Creative Commons Attribution-ShareAlike 4.0 International License. Changes to the theme are the work of eTRAP.

