

Preparing a Text Corpus with **Transkribus**

Pre-Processing, Training, Post-Processing

Dr. Markus Müller, Milena Ramirez
Leibniz-Institut für Europäische Geschichte, Mainz

19. November 2019

Initial problem:

In the 16th century, the inquisition in Spain, Italy and France censored the books of the Mainz cathedral preacher Johann Wild OFM.

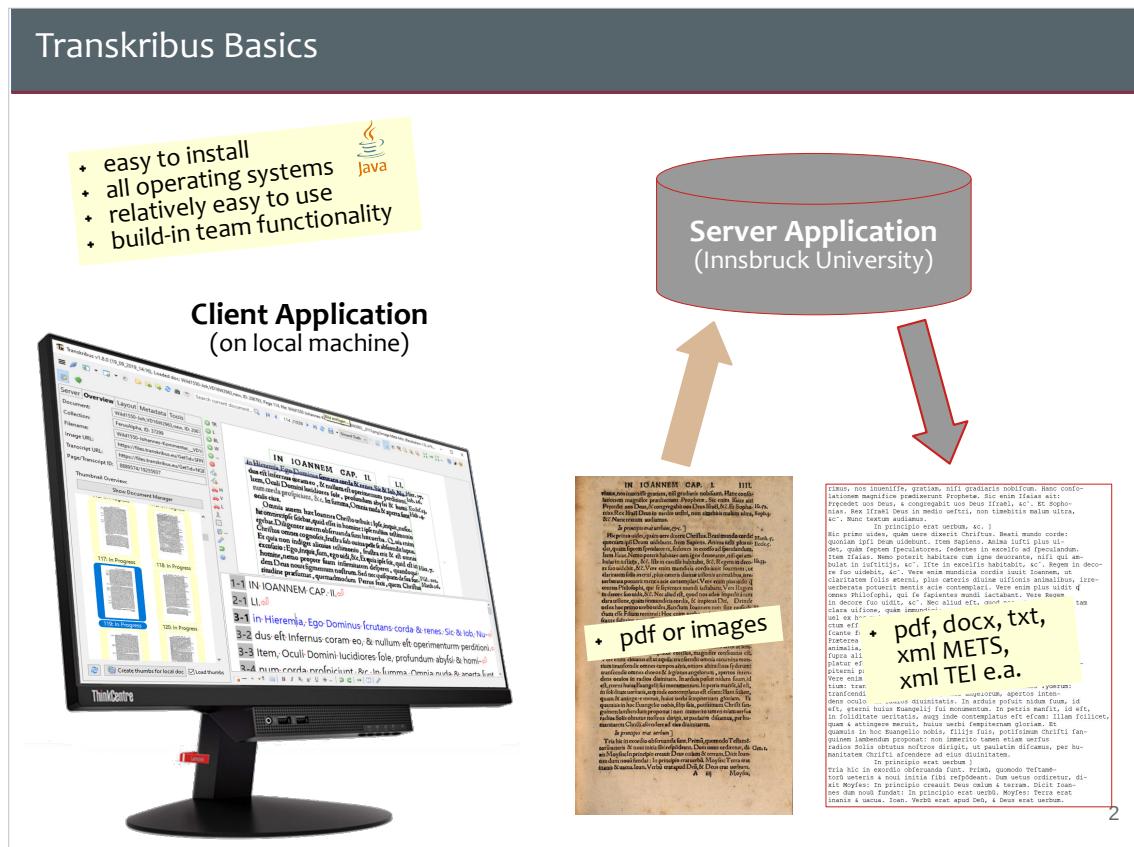
Interestingly, the actual words and sentences that were deleted or modified or inserted, differed between the different inquisitions, thus various censored versions of the same book were printed.

What exactly did the different censors change in the text? What were exactly are the differences between Spanish, Italian and French censorship?

In 2018, I began to search manually for these differences, but very soon, it became clear that I would never come to a reliable result because manual comparision of texts is extremely time-consuming and error-prone.

Therefore, I decided to invest my time into up-to-date OCR technology instead of manual search.

Transkribus Basics



What is Transkribus?

- a server application for handwritten text recognition (HTR) running on servers at Innsbruck University.
- a very nice and easy to use client software for all steps of the OCR process.
- client software needs Java Runtime Environment therefore, it works on all operating systems without installation.
- Originally made for handwritten text recognition (HTR), i.e. manuscripts from the archives, it also works very well with printed texts.
- Basic facts about the Transkribus project as such on <https://de.wikipedia.org/wiki/Transkribus>

What do you need to get started?

- a bunch of images of a handwritten or printed text.
- a lot of time and patience.

Our Workflow

(A) Preparation & Training

get pdfs or images of your sources

- color images work better than pure black & white images
- resolution is not too important (it works quite well with 3.000 x 2.000 and 10.000 x 6.000 pixels)

upload images/pdfs to Transkribus

produce training data (“ground truth”)

- a) manual segmentation (**4+ min**/page)
 - define text regions
 - detect baselines (semi-automatically)
- b) manual transcription (**30+ min**/page)
 - 50 pages (for printed texts)
 - 100 pages (for manuscripts)
 - special characters (unicode)
- c) manual correction (**20+ min**/page)
In total: 2 people, 4 months = 75 pages

train the neural network

- ca. **6 h** for 75 pages, 200 epochs
 - character error rate of ca. 3% (manuscripts)
 - CER of ca. 0,25% (printed texts)

(B) Application

run automatic layout analysis

- ca. **20 s**/page
- manual correction

run OCR

- ca. **30 s**/page
- manual correction
(if you need 100% perfect result)

export the result

- Excel, Word
- plain text
- xml formats: TEI, METS
- pdf

(C) Post-Processing

3

The whole process took a lot of time and energy:

April–August 2019: build a workflow and manual transcription of 76 pages in a team of two with and mutual revision of the transcriptions

30 August 2019: training of a first model: >> character error rate: 0,23%

September/October 2019: transcription of another 20 pages of a second book and struggling with the relatively poor layout analysis

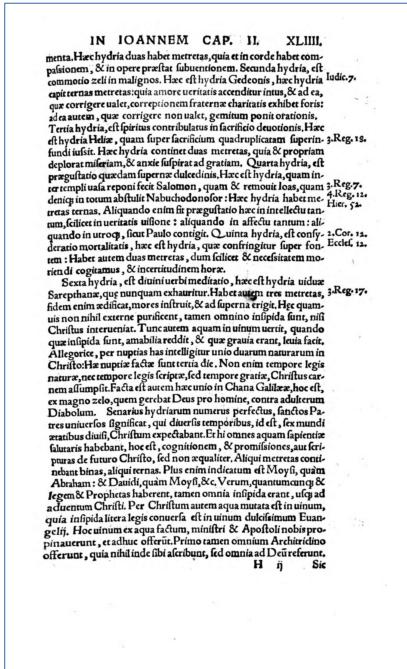
November 2019: second training with 20 pages: >> character error rate: 1,60%

Roadmap for the next months: transcribe and train more books of the same author and build a stable post-processing workflow, i.e. invest a lot of time into programming Python scripts.

„character error rate of xy%“ = of all the characters on the page, xy% are wrong.

In 2018 , the same training still needed ca. 20 h and produced an error rate of ca. 10%. Thanks to better hardware and an improved architecture of the underlying neural network, the Innsbruck Transkribus team could reduce the time needed for training to ca. 6 h and improve the character error rate to under ca. 3% for manuscripts and under 1% for printed texts (2019).

Images



- Download pdfs from Google Books (only available for copyright-free stuff)

(Hover with the cursor over the „e books“ button on the left. A context dialog opens up. Click on „Download PDF“ in the context dialog. The pdf will be downloaded immediately and opens in your browser or in a local pdf reader.)

Pro: super easy.
Contra: only black & white, sometimes bad image quality.

- Download high quality pdfs or images from library websites

(e.g. <https://www.digitale-sammlungen.de/>)

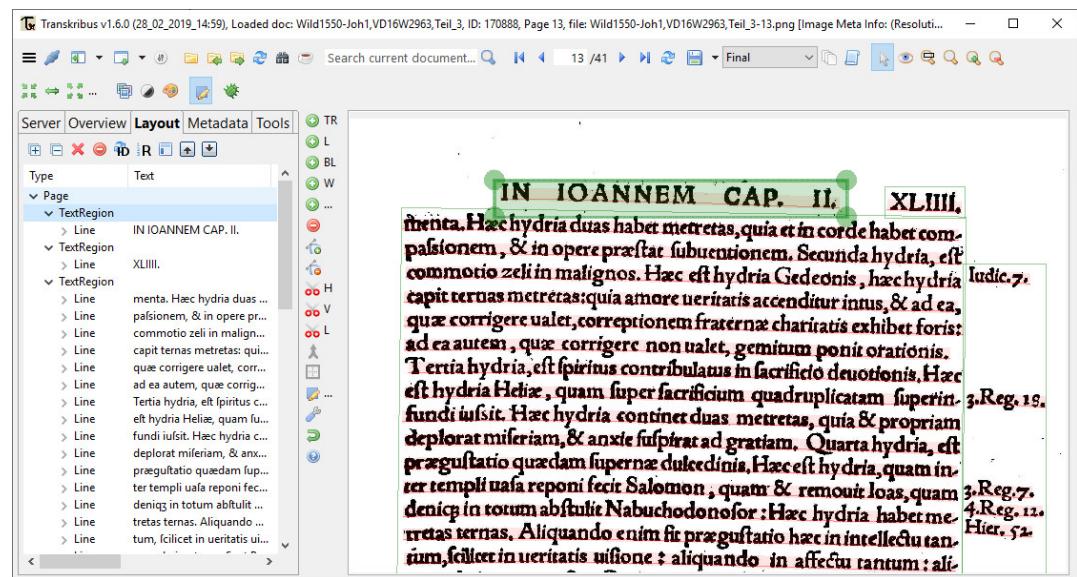
Pro: best image quality available
Contra: often sophisticated user interfaces, sometimes the same images as in Google Books.

- Produce your own digital images.

(e.g. with your smartphone or even with specialized tools like <https://scantent.eu/en/>)

Pro: maximum control etc.
Contra: you need the technical skills at the time to do it well; may violate copyright rules or personal rights.

Segmentation Mode

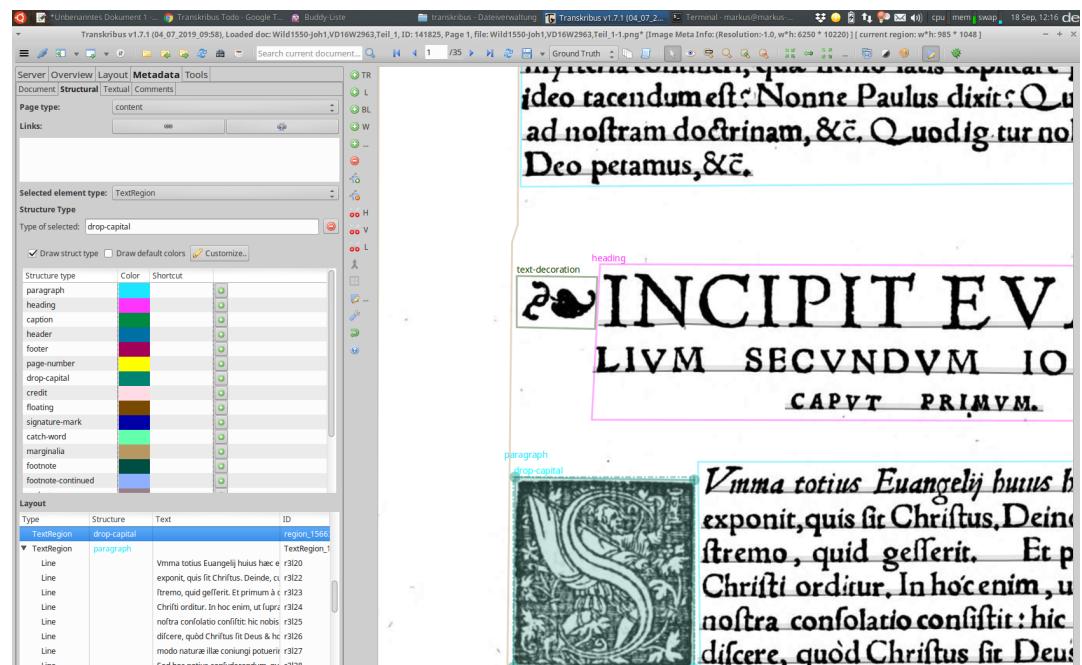


5

There are two major „modes“ in Transkribus:

(1) Segmentation mode: By drawing rectangles and lines, you tell the software where the different text blocks are situated on the page.

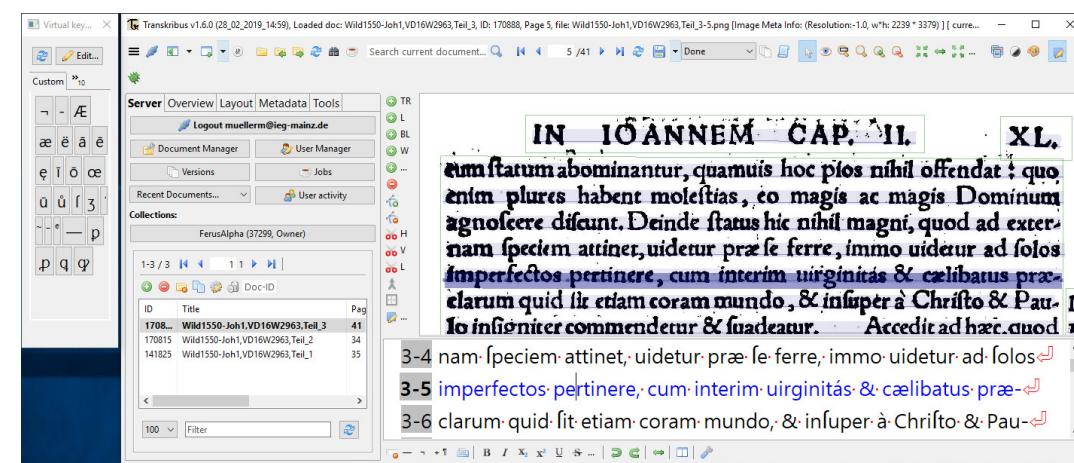
Structural Metadata: Text Regions



6

You can add metadata to the different text regions (as well as to many other parts of the dataset, e.g. metadata of the document itself, metadata for documentation of the transcription standards etc.).

Transcription Mode

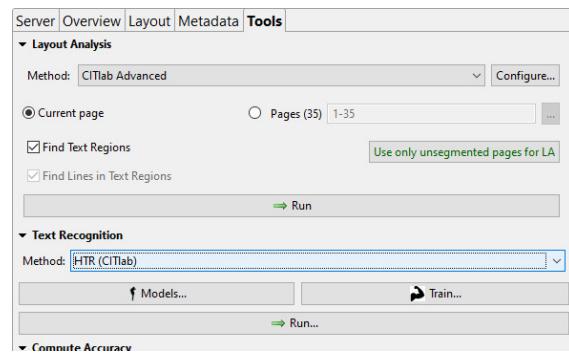


7

(2) Transcription mode: before the software can recognize text, it has to be trained with sample data. That means, you have to transcribe a lot of pages manually (ca. 50 pages for printed texts, 100 pages for manuscripts).

The Transkribus interface makes transcription very easy: The corresponding line is highlighted on the photo, you can define shortcuts for special characters etc.

Tools for Layout Analysis, Training and Text Recognition



8

As soon as the transcription phase is completed, you can **train a „model“**, i.e. the software will try to figure out how the images and the given text („ground truth“) are connected. As soon as the software has found a reliable and stable way to come from the image to the desired text, the „training“ process is completed. The set of rules that the software produced is called „model“.

If you provide the software with another image of a similar kind that it has never „seen“ before, it should be able to „read“ the text on the image without any „help“ or intervention by the user.

Having finished the training, you can „run“ the OCR with the pre-trained model on the rest of the book and Transkribus will automatically transcribe the text (in my case: another 1500 pages).

Unfortunately, the automatic layout analysis does not recognize the marginalia as an separate text region, thus still have to invest 3–5 min per page to correct the layout before you can acutally run the OCR.

Result: A (nearly) Perfect Text!

The screenshot shows a comparison between a handwritten Latin manuscript and its digital transcription. The manuscript on the left is from the Codex Bezae Cantabrigiensis, folio 49v, featuring dense Gothic script. The transcription on the right is in modern Latin, with red highlights indicating errors or specific features. The transcription reads:

1-1 IN-IOANNEM-CAP. I.
2-1 4.
3-1 iustitia nobis autē confusio faciei noſtræ; &c. Et infra: Non
enim in iuſtificationibus noſtris proterminus preces noſtræ coram te, ſed in miferationibus tuis multis. Præterea
ſubdit, ac clarius ſeipſum exponit, dicens: Dirige uiam Do-
mini, hoc eft appareat, complanare, impedimenta remouete.
In uia naque eft Dominus uiam uite prafentis ingeffus
eſt, ad grande opus proficisciuit. Ecce enim tanquam ſpon-
ſus procedit de thalamo ſuo, exultauit ut gigas ad curren-
dam viam a ſummo celo egrefio eius, &c. Nolite igiur ob-
ſtare, nolite offendicula ponere. Nam ſuper quem ceci-
derit lapis hic, comminuet eum. Vae enim illis, quorum ſce-
lere in mortem caſurus eft. Deinde: Dirige uias Domini
hoc eſt, recipie aduenientem Chriftum, date locum Spir-
iti ſancto. Pænitentia, & agnoscite peccata. Rectas faci-
te ſemitas, id eft, abiuite peccata, implete legem, faci-
te fructus dignos pænitentia. In hanc modum Ioannes uiā
Domino parat, quia quanto plus ad pænitentiam & imple-
tionem mandatorum nos hortatur, tanto magis ad Chriftū
nos compellit, ſine quo nec pænitentia fieri, nec lex imple-
ri potefit. Tale igiur erat officium Ioannis, quale in libro
Genes̄is legitur: quod cum Iofeph a Pharaone exaltatus cir-
cunduceretur, prece clamabat, ut omnes corā eo genu fler-
etarent, & prepoſitum ſcirent uniuersitatem Aegypti. Si-
gnanter autem dicit: Dirige uias Domini, non uelras. In
Iſaiā enim dicit Dominus: Non ſunt uite mee, uite uelras,
&c. Vix enim hominum ſunt ambitione, hypocritis, auaricia,

The resulting text produced by Transkribus has only very few errors. The model used here has been trained with 30 pages of the same and 75 pages of a similar book.

Pitfalls During Transcription

non credat, Christum esse uerum Deum, quemadmodum & Paulus
a'c: Pr̄edicamus Christū, Iudeis scandalum, Gentibus stultitiam, &c.
Deinde etiam quidam hæretici ē nobis egressi, quamuis non erant

a't: Pr̄edicamus Chirſtū, Iudeis fcandalum, Gentibus stultitiam, &c.

How many errors can you spot in this transcription?

non credat, Christum esse uerum Deum, quemadmodum & Paulus
a'c; Pr̄edicamus Christū, Iudeis scandalum, Gentibus stultitiam, &c.
Deinde etiam quidam hæretici ē nobis egressi, quamuis non erant

X a't: Pr̄edicamus Chirſtū, Iudeis scandalum, Gentibus stultitiam, &c.

✓ ait: Pr̄edicamus Christū, Iudeis scandalum, Gentibus stultitiam, &c.

Nota bene: Errors in the transcription will scale up!

- the neural net will ‘learn’ every single mistake you make
- OCR will produce erroneous output

It would be so nice to have a Latin spell-checker (for early modern Latin, including all the abbreviations ;-) ...

Why OCR alone is not enough...

Version A printed in Mainz, 1550

non credat, Christum esse verum Deum, quemadmodum & Paulus
a*t*: Pr*ed*icamus Christ*ū*, Iudeis scandalum, Gentibus stultitiam, &c.
Deinde etiam quidam hæretici e nobis egressi, quamvis non erant

ait: Pr*ed*icamus Christ*ū*, Iudeis scandalum, Gentibus stultitiam, &c.

ait: Pr*æ*
dicamus Christum, Iudæis scandalum, Gentibus stulti-
tiam, &c.

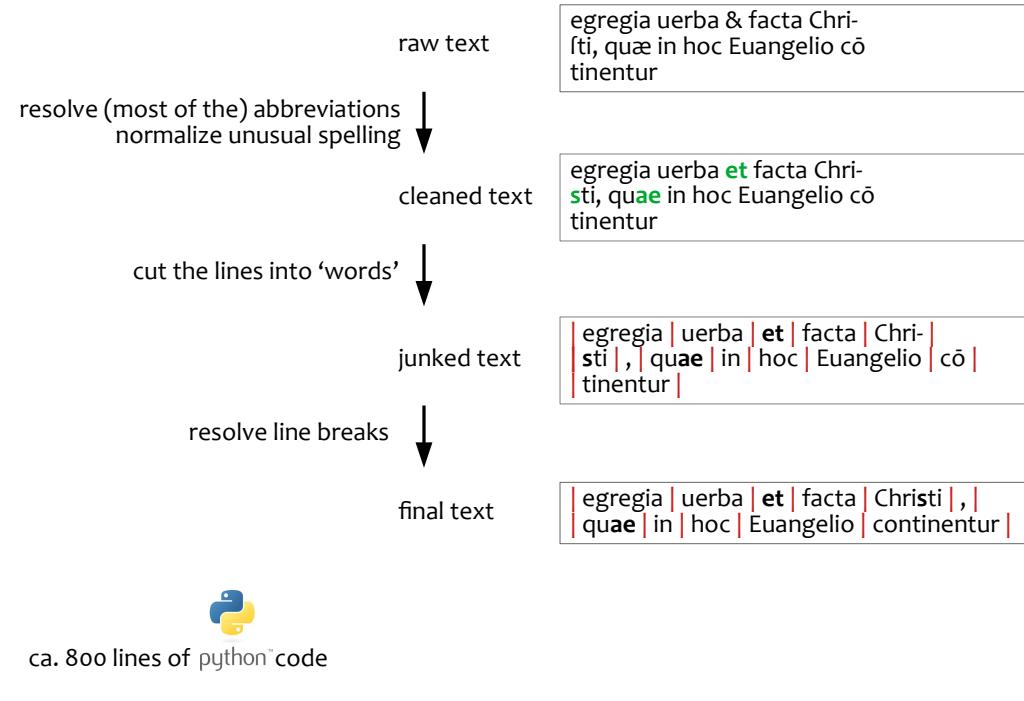
stum esse verum Deum, quemadmodum & Paulus a*t*: Pr*ed*
dicamus Christum, Iudeis scandalum, Gentibus stulti-
tiam, &c. Deinde etiam quidam hæretici e nobis egressi,

Version B printed in Rome, 1577

12

I want to compare one version of the text with another version. Due to the different usage of abbreviations and different spelling, the text produced by Transkribus has be post-processed in order to be comparable.

Post-Processing to make the texts comparable



13

Post-processing cannot be done with Transkribus. At the moment, I try to tackle the problem with a Python script.

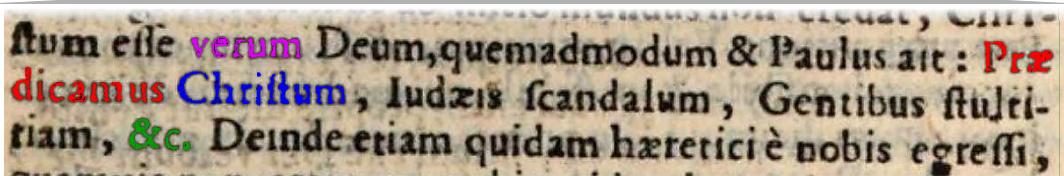
After post-processing: perfectly searchable & comparable texts

Version A printed in Mainz, 1550

**non credat, Christum esse uerum Deum, quemadmodum & Paulus
ait: Praedicamus Christum, Iudeis scandalum, Gentibus stultitiam, &c.
Deinde etiam quidam haeretici e nobis egressi, quamuis non erant**

non credat, Christum esse **uerum** Deum, quemadmodum et Paulus
ait: **Praedicamus Christum**, Iudeis scandalum, Gentibus stultitiam, **etc.**
Deinde etiam quidam haeretici e nobis egressi, quamuis non erant

esse **uerum** Deum, quemadmodum et Paulus ait: **Praedicamus Christum**,
Iudeis scandalum, Gentibus stultitiam,
etc. Deinde etiam quidam haeretici e nobis egressi,



stum esse verum Deum, quemadmodum & Paulus ait : Prae dicamus Christum , Iudæis scandalum , Gentibus stulti tam , &c. Deinde etiam quidam haeretici e nobis egressi ,

Version B printed in Rome, 1577

14

The Python script resolves the abbreviations and normalizes the text and, finally, the data for text comparison is ready.

Use Case: Compare Text Versions to Find Censorship

Juxta - welcome.jxt

File Edit View Collation Help

Comparison Set

Wild1577_cleaned.txt Wild1550_cleaned.txt

Epistolis nihil non agunt, ut ad bona opera nos extimulent: Maximeque eos detestantur, qui libertatem uertunt in occasionem carnis, et uelamen malitiae. Ideo aeternus Christus doctrinam suam laetum **nuncius** nominat, quod ea non, ut, lex, poena et minis **ad bonum** cogit, sed donis et promissis afficit. sic igitur uetus Testamentum, tametsi plures habent **adjunctas** promises, ut patet in Propterea et Psalmis: quia tamen haec promises nondum impletas sunt, erant, nomen legis obtinuit, unde et Christus etiam Psalmos legem nominat. Nonne, inquit, in lege uestra scriptum est: Ego dixidi estis, sic doctrina Christi, tametsi multa doceat de operibus, quia tamen promises nunc per Christum impletatae sunt, et implentur, Euangeli nomen obtinuit. **Sicut igitur impie cum patribus veteris Testamenti egisset, qui omissionis promises, solam eis legem praedicasset;** per hoc enim eos ad desperationem adegisset, sic impie fact, qui Christiculis, neglecta doctrina operum, solas promises inculcat: per hoc enim eos in uanum securitatem, immo leuitatem inducit, quemadmodum uidemus. summa, sicut lex non excludit promises, immo multas adjunctas habuit: sic Euangelium nonexcludit opera, immo summe ea commendat. Haec ideo pluribus exposui, propter eos, qui de Euangeliis gloriabantur, cum nesciant, quid sit Euangelium. Praeter hos sunt ali, qui obgaudent: Quid, inquit, necesse fuit, scribere Euangeli, cum Christus non dixit: scribite Euangeli, sed praeificate Euangeli. Quod propter omnes non in tabulis lapideis, sed in cordibus hominum scripturam legem suam. Atque hinc quidam in tantam insaniam peruererunt, ut scripturam in terra reciant, ac nouas uitiosas reuelationes dicant, et cetera. Quod immundissimus diabolus, statim post Euangeli luis transfiguravit, eosque mirabiliter decipit, ut uidemus. Audi igitur, uerum quidem est, melius esse, Euangelium Christi in corde habere scriptum, quam in membrana. Deinde oportuerat quidem Christianos non indigere auxilio litterarum, sed tam mundum habere cor, spiru Dei purificatum, ut illud

Wild1577_cleaned.txt - Source

opus domini inserviatis, derelinquere sine compunctione incertum **uero**, quod uideam uos non similliter affectos, nec eadem appetere. Quod si eam

Source Images Notes Moves Search

Comparision of two versions of the same text produced with Transkribus to find differences, especially censorship. (Using plain text files with the free [Juxta](#) tool.)

15

As you can see, there are still a lot of „false positives“, i.e. words that are marked as different in both texts not because they were actually censored but because they were spelled differently, or due to an error in the transcription/OCR process, or due to an error made by the printer in the 16th century (believe me, they made quite a lot of errors!). There is still a lot of work to be done to improve the post-processing as well as the comparison itself...

Further Information

Information

- Transkribus website: <https://transkribus.eu/Transkribus/>
- Transkribus wiki: English: https://transkribus.eu/wiki/index.php/Main_Page
German: <https://transkribus.eu/wikiDe/index.php/Hauptseite>
- Article on medieval manuscript OCR with Transkribus (in German):
<https://mittelalter.hypotheses.org/21828>

Ideas

- Random selection of the pages to be transcribed will cover more ‘special cases’ and exceptions than transcribing whole chapters.
- Use [state-of-the-art spell checking technology](#) to reduce errors during transcription and to speed up the post-processing of the OCR output of Transkribus.
- Further automation: Use the [Transkribus REST api](#) to download OCR data for post-processing and to upload the results of automated post-processing.

