



# Compilation of a Large Ground-Truth Data Set Using Transkribus

Matthias Boenig & Kay-Michael Würzner

[{boenig|wuerzner}@bbaw.de](mailto:{boenig|wuerzner}@bbaw.de)

Transkribus User Conference

Vienna, 2nd November 2017



# Overview

**Goal:** Compilation of a large, **homogeneous** Ground Truth (GT) data set

- Various **heterogeneous** sources
- Annotation on the **textual** and/or **structural** level

**Background:** OCR-D initiative

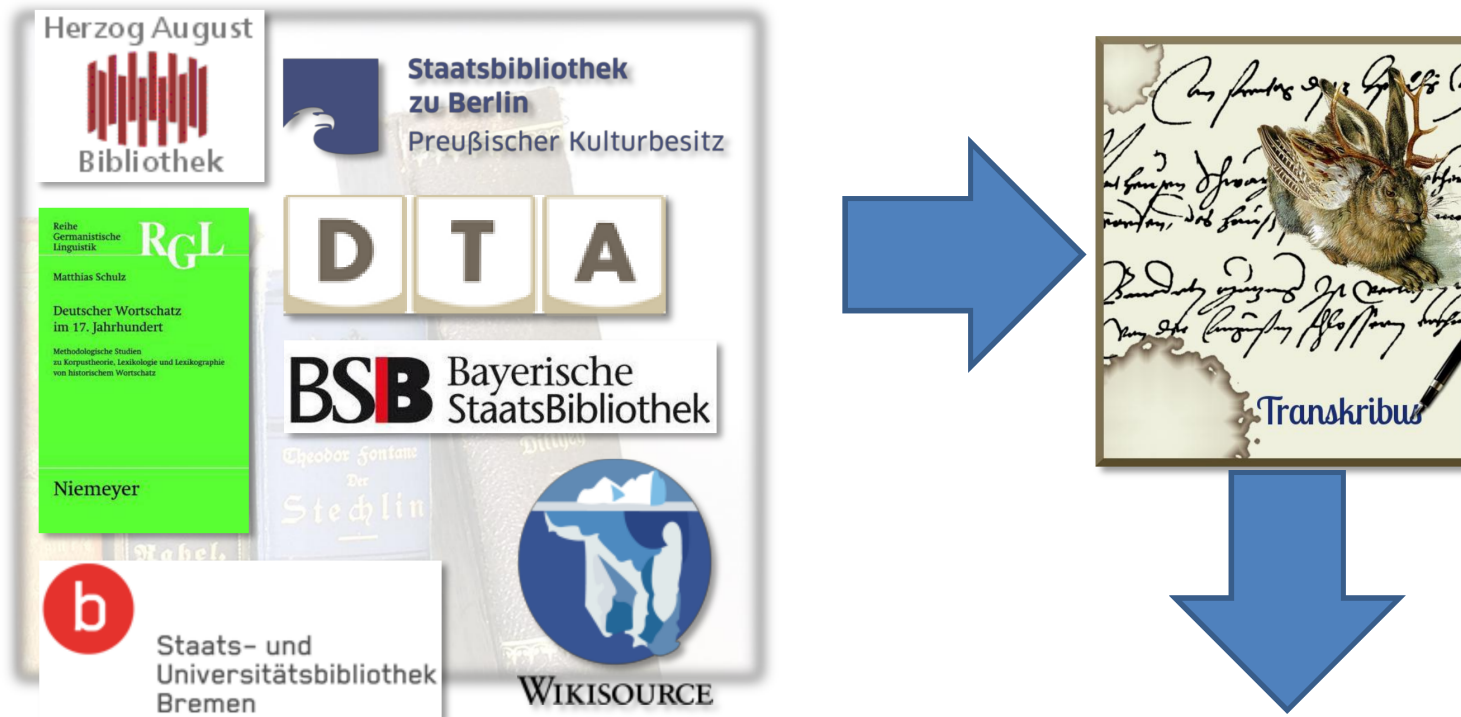
- Funding by the *Deutsche Forschungsgemeinschaft*  
→ Improvement of OCR tools for **historical printings** (i.e. VD 16, 17, 18)
- Coordination project
  - Identify to-dos, desiderata and improvement options
  - Development of a call for proposals
  - Merge (sub-)project results into a **productive workflow**

**Procedure:** Annotation with Transkribus

1. Import images and existing text and/or structural information
2. **Harmonization** and completion within Transkribus



# Overview



- Various GT sources
- Containing either text or structural annotations in differing quality
- By now,  $\approx 130$  documents with  $\approx 500$  pages
- A lot more to come!



# Workflows

Existing text	Existing structure





# Workflows

Existing text	Existing structure
Import images	



# Workflows

The screenshot displays the Transkribus web interface. On the left, a sidebar contains navigation tabs (Server, Overview, Layout, Metadata, Tools) and a list of collections, including 'GT3-Sample (5120, Owner)'. Below this is a table of documents with columns for ID, Title, Pages, and Upload status. The main area shows a 'Document ingest / upload' dialog box with options for upload methods (private FTP, URL of DFG Viewer METS, single document, or extract images from PDF). The 'Upload via private FTP' option is selected. A table lists documents on private FTP, including 'nn\_catechismus\_1684' and 'Archiv'. The dialog also includes a field for 'Add to collection:' set to 'GT3-Sample (5120, Owner)' and 'Cancel' and 'Upload' buttons.

ID	Title	Pages	Upload
17737	a_gehema_feldapothke_1688_11	11	boeni
17736	a_gehema_feldapothke_1688_9	9	boeni
17735	a_gehema_feldapothke_1688_10	10	boeni
17734	a_gehema_feldapothke_1688_10	10	boeni
16012	braendl_thaumatoграфия_1692	11	boeni
16011	braendl_thaumatoграфия_1692	11	boeni
15965	carrichter_speiskammer_1610_11	11	boeni
15964	carrichter_speiskammer_1610_10	10	boeni
15034	hoefer_bienenkunst_1614_2	11	boeni
15033	hoefer_bienenkunst_1614_1	11	boeni

Directory	Title	Nr. of Files	Last modified
nn_catechismus_1684	nn_catechismus	1	Thu Oct 26 14:07:45 CEST 2017
Archiv	Archiv	28	Thu Oct 26 14:01:35 CEST 2017



# Workflows

The screenshot displays the Transkribus web interface. On the left, a sidebar contains navigation tabs: Server, Overview, Layout, Metadata, and Tools. Below these are buttons for 'Logout boenig@bbaw.de', 'Document...', 'Jobs', 'Versions', and 'User activity'. A 'Recent documents...' section is also present. The 'Collections:' section shows 'GT3-Sample (5120, Owner)'. Below this is a table listing documents:

ID	Title	Pages	Uploa
28396	nn_catechismus_1684	22	boeni
17737	a_gehema_feldapotheke_1688_	11	boeni
17736	a_gehema_feldapotheke_1688_	9	boeni
17735	a_gehema_feldapotheke_1688_	10	boeni
17734	a_gehema_feldapotheke_1688_	10	boeni
16012	braendl_thaumatographia_1692	11	boeni
16011	braendl_thaumatographia_1692	11	boeni
15965	carrichter_speiskammer_1610_	11	boeni
15964	carrichter_speiskammer_1610_	10	boeni
15034	hoefler_bienenkunst_1614_2	11	boeni
15033	hoefler_bienenkunst_1614_1	11	boeni

The main area shows a large image of a historical manuscript page. The text on the page is:

CATECHISMUS,  
Oder  
Kurzer Unterricht  
Christlicher Lehr / wie der in  
Kirchen und Schulen der Chur-  
fürstlichen Pfalz getrieben  
wird.  
Aus Chur-Fürstl. Pfalz Verordnung  
kürzlich erklärt / und mit Zeugnissen der

The interface includes a search bar at the top, a toolbar with various icons, and a footer with a 'Region:' label and a progress bar.



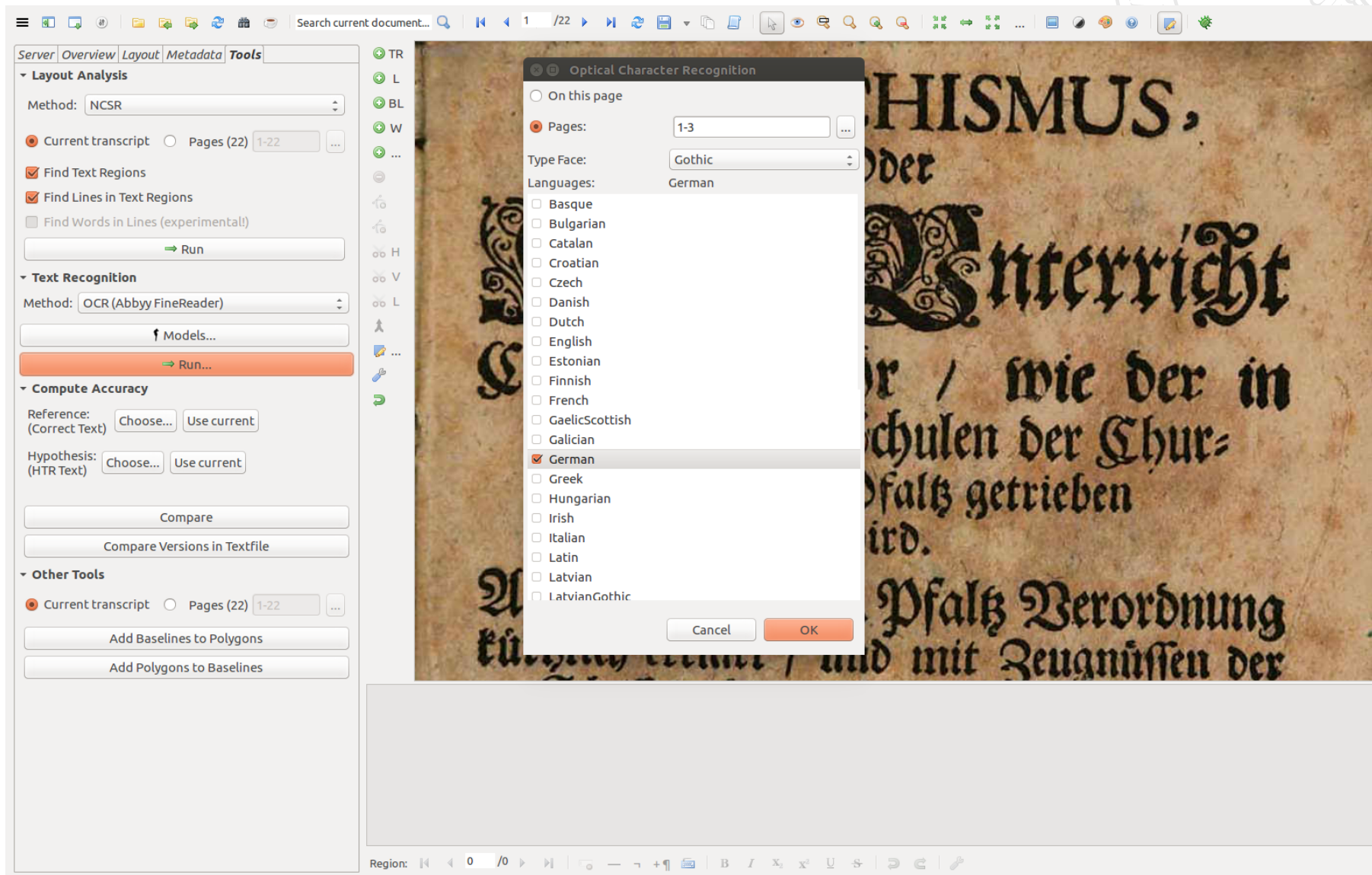
# Workflows

Existing text	Existing structure
Import images	
Run FineReader for initial layout version	Import Page XML

- 
- 
- 
- 



# Workflows





# Workflows

The screenshot displays the Transkribus software interface. The main window shows a manuscript page with German text in a historical script. The text is organized into paragraphs and lines, with green bounding boxes indicating the current selection. The sidebar on the left provides a detailed structure of the document, including page, text regions, and lines, with corresponding read IDs and line IDs.

**Text on the manuscript page:**

Kirchen und Schulen der Chur=  
fürstlichen Pfalz getrieben  
wird.

Aus Chur=Fürstl. Pfalz Verordnung  
fürstlich erklärt / und mit Zeugnissen der  
Schrift bestätigt: Für die Schul=Jugend in  
Ihrer Chur=Fürstl. Durchl. Landen.

**Structure Table (from sidebar):**

Type	Text	Structure	Read ID	ID
Page				
Printspace				
TextRegion		paragraph	1	r_1_1
Line	CATECHISMU		1	tl_1
TextRegion		paragraph	2	r_1_2
Line	Oder		1	tl_2
Graphic			3	r_2
TextRegion		paragraph	4	r_3_1
Line	Alls Lhur-Fürs		1	tl_3
TextRegion		paragraph	5	r_3_2
Line	kürtzlich erklä		1	tl_4
TextRegion		paragraph	6	r_3_3
Line	Schnsfr bestä		1	tl_5
Line	JbrerKhur'Fü		2	tl_6
Graphic			7	r_4
TextRegion		paragraph	8	r_5_1
Line	Mt Chur^ürstl		1	tl_7
TextRegion		paragraph	9	r_5_2
Line	Gedruckt und		1	tl_8
TextRegion		paragraph	10	r_5_3
Line	Bty denen Wa		1	tl_9
TextRegion		paragraph	11	r_5_4
Line	Dero Uoiverlt		1	tl_10
TextRegion		paragraph	12	r_5_5
Line	/' — "		1	tl_11
Graphic			13	r_6

**Transcription Area:**

1 Alls Lhur-Fürstl. Pfaltz Verordnung



# Workflows

Existing text	Existing structure
Import images	
Run FineReader for initial layout version	Import Page XML
Manually correct layout	Run external OCR for initial text version

- 
- 
- 
- 



# Workflows

Server Overview Layout **Metadata** Tools

Document **Structural** Textstyle Tagging Comments

Edit status: In Progress

Page type:

Links:

Selected element type: TextRegion

Structure Type

- ☐ paragraph ☐ heading
- ☐ caption ☐ header
- ☐ footer ☐ page-number
- ☐ drop-capital ☐ credit
- ☐ floating ☐ signature-mark
- ☐ catch-word ☐ marginalia
- ☐ footnote ☐ footnote-continued
- ☐ endnote ☐ TOC-entry
- ☐ other

Apply Apply down

1  
2  
3  
4  
5

CATECHISMUS,  
Oder  
Kurtzer Unterricht  
Christlicher Lehr / wie der in  
Kirchen und Schulen der Chur-  
fürstlichen Pfalz getrieben  
wird.

Region: 3 / 11





# Workflows

Existing text	Existing structure
Import images	
Run FineReader for initial layout version	Import Page XML
Manually correct layout	Run external OCR for initial text version
<b>Copy and paste text region by region</b>	

- 
- 
- 
- 



# Workflows

The screenshot displays the Transkribus software interface, which is used for transcribing historical documents. The main window shows a manuscript image with the title "Oder Kurtzer Unterricht Christlicher Lehr / wie der in Kirchen und Schulen der Churfürstlichen Pfaltz getrieben wird." The left sidebar shows a tree view of the document structure, including "Page", "TextRegion", and "Line" elements. The bottom panel shows the "CATECHISMUS." title and the main text. A smaller window in the foreground shows the "catechismus\_1684.xml" file with the text transcribed into XML format.

Server Overview Layout Metadata Tools

Type Text Structure Readli ID

▼ Page

Printspace

▼ TextRegion

Line CATECHISMUS paragraph 1 TextR tl\_1

Line Oder 2 tl\_2

Line Kurtzer Unter 3 line\_1

Line Christlicher Le 4 line\_1

Line Kirchen und S 5 line\_1

Line fürstlichen Pf 6 line\_1

Line wird. 7 line\_1

▼ TextRegion

Line Aus Chur-Fürf paragraph 2 TextR tl\_3

Line kürztlich ekl paragraph 1 tl\_4

Line Schriftt beftä 2 tl\_5

Line Jhrer Chur-F 3 tl\_6

Graphic 3 r\_4

▼ TextRegion

Line Mt Chur^ürst paragraph 4 r\_5\_1

TextRegion paragraph 5 TextR tl\_7

Graphic 7 r\_6

1 CATECHISMUS.

2 Oder

3 Kurtzer Unterricht

4 Christlicher Lehr/ wie der in

5 Kirchen und Schulen der Chur-

6 fürstlichen Pfaltz getrieben

7 wird.

File Bearbeiten Suchen Projekt Optionen Werkzeuge TEI P5 Dokument Fenster Hilfe

XPath 2.0 XPath ausführen auf 'Aktuelle Datei'

catechismus\_1684.xml

TEI text front titlePage docTitle titlePart

Kurtzer Unterricht

Christlicher Lehr/ wie der in

Kirchen und Schulen der Chur-

fürstlichen Pfaltz getrieben

wird.

Text Raster Autor

/home/.../catechismus\_1684.xml [DTA-HTML] Transformation - erfolgr... U+0000 Geändert



# Workflows

Existing text	Existing structure
Import images	
Run FineReader for initial layout version	Import Page XML
Manually correct layout	Run external OCR for initial text version
<b>Copy and paste text region by region</b>	
	Manually correct text

- 
- 
- 
- 



# Workflows

Existing text	Existing structure
Import images	
Run FineReader for initial layout version	Import Page XML
Manually correct layout	Run external OCR for initial text version
<b>Copy and paste text region by region</b>	
	Manually correct text

- Somewhat naïve approach
- External Page XML creation or
- Intermediate export and (re-)import as alternative options
- **Not very comfortable**



# Desiderata

- Transkribus is a wonderful tool!
  - ▶ Support for **polygonal** regions
  - ▶ **Multiple** OCR options
  - ▶ **Collaborative** working environment with basic version control
  - ▶ TEI export
- For **GT creation**, we would welcome
  - ▶ OCR application on **specific regions** also for FineReader
  - ▶ Dedicated **text import** functionalities (e.g. on paragraph level)
  - ▶ METS import which accounts for **existing structural annotations** and linked ALTO
  - ▶ **Automatic support** during manual post correction
  - ▶ TEI import



# Collaboration

## ■ OCR-D GT Guidelines

- ▶ Documentation of existing OCR-D GT
- ▶ Instructions for GT creation
  - Already used within the OCR-D project
  - Perspectively also used in a **broader context** (community use)
- ▶ Automatic **validation** of GT data
- ▶ (Semi-)automatic **conversion** of existing GT data sets
- ▶ Plans for setting up a **GT repository** for print publications and handwritten documents

## ■ Availability

**View:** <https://kaskade.dwds.de/~matthias/ocr-d/>

**Sources:** <https://github.com/OCR-D/>



- Transkribus User Documentation: A proposal
  1. Step: Change the documentation format from Wiki to DITA
    - ▶ XML-based documentation format
    - ▶ **Topic-oriented** internal and “external” structure (i.e. presentation)
    - ▶ Various automatically generated **presentation modes**
  2. Step: Build and organize a documentation **source repository** (e.g. on github)
  3. Step: Involve the user community into the documentation process
    - ▶ Non-developer view point
    - ▶ Recipes for frequent tasks





**Many thanks for your attention.**

