



Compilation of a Large Ground-Truth Data Set Using Transkribus

Matthias Boenig & Kay-Michael Würzner
{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

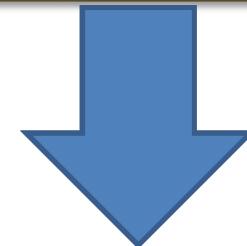
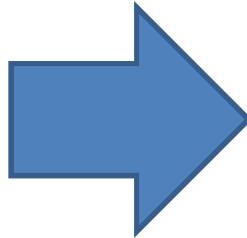
- a. Funding by the *Deutsche Forschungsgemeinschaft*
→ Improvement of OCR-D tools for **historical printings** (i.e. VD 16, 17, 18)
- b. Coordination project
 - Identify to-dos, desiderata and improvement options
 - Development of a call for proposals
 - Merge (sub-)project results into a **productive 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 shows a digital library interface with a central document ingest/upload dialog box overlaid on a main workspace.

Main Workspace:

- Header:** Search current document... /11, various icons for search, navigation, and tools.
- Left Sidebar:** Server Overview, Layout, Metadata, Tools; Logout boenig@bbaw.de; Document..., Jobs, Versions, User activity; Recent documents...; Collections: GT3-Sample (5120, Owner).
- Table View:** Shows a list of documents with columns: ID, Title, Pages, and Upload status (boeni). Examples include:
 - 17737 a_gehema_feldapotheke_1688 11
 - 17736 a_gehema_feldapotheke_1688_1 9
 - 17735 a_gehema_feldapotheke_1688_10
 - 17734 a_gehema_feldapotheke_1688_10
 - 16012 braendl_thaumatographia_1692 11
 - 16011 braendl_thaumatographia_1692 11
 - 15965 carrichter_speiszkammer_1610_11
 - 15964 carrichter_speiszkammer_1610_10
 - 15034 hoeftler_bienenkunst_1614_2 11
 - 15033 hoeftler_bienenkunst_1614_1 11
- Right Sidebar:** Various icons for file operations like TR, L, BL, W, etc.

Document Ingest / Upload Dialog:

Options:

- Upload via private FTP
- Upload single document
- Upload via URL of DFG Viewer METS
- Extract and upload images from pdf

FTP Location: ftp://transkribus.eu

File List:

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

Buttons:

- Add to collection: GT3-Sample (5120, Owner)
- Cancel
- Upload

Workflows

The screenshot shows a digital library interface with a sidebar on the left and a main content area on the right.

Left Sidebar:

- Server Overview Layout Metadata Tools
- Logout boenig@bbaw.de
- Document... Jobs
- Versions User activity
- Recent documents...
- Collections: GT3-Sample (5120, Owner)

Main Content Area:

Title Page:

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

Text Below Title:

Aus Chur-Fürstl. Pfalz Verordnung
kürzlich erklärt / und mit Zeugnissen der

Table (List of Items):

ID	Title	Pages	Upload
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_speiszkammer_1610_	11	boeni
15964	carrichter_speiszkammer_1610_	10	boeni
15034	hoefler_bienenkunst_1614_2	11	boeni
15033	hoefler_bienenkunst_1614_1	11	boeni

Workflows

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

-
-
-
-



Workflows

The screenshot shows the Transkribus application interface. On the left, there is a sidebar with various tools and options:

- Layout Analysis:** Method set to "NCSR". Options include "Current transcript" (selected), "Pages (22) 1-22", "Find Text Regions" (checked), "Find Lines in Text Regions" (checked), and "Find Words in Lines (experimental)". A "Run" button is present.
- Text Recognition:** Method set to "OCR (Abbyy FineReader)". Options include "Models..." and a "Run..." button.
- Compute Accuracy:** Reference set to "(Correct Text)" and Hypothesis set to "(HTR Text)". Buttons for "Compare" and "Compare Versions in Textfile" are available.
- Other Tools:** Current transcript selected. Options include "Add Baselines to Polygons" and "Add Polygons to Baselines".

A central window displays a page from a historical document in Gothic script. Overlaid on this is an "Optical Character Recognition" dialog box:

- Radio button: "On this page" (unchecked)
- Radio button: "Pages:" (checked) with value "1-3" (unchecked)
- Text Face: "Gothic" (selected)
- Languages:
 - Basque
 - Bulgarian
 - Catalan
 - Croatian
 - Czech
 - Danish
 - Dutch
 - English
 - Estonian
 - Finnish
 - French
 - GaelicScottish
 - Galician
 - German** (checked)
 - Greek
 - Hungarian
 - Irish
 - Italian
 - Latin
 - Latvian
 - LatvianGothic

At the bottom of the dialog box are "Cancel" and "OK" buttons.



Workflows

The screenshot shows the Transkribus software interface. On the left, there is a sidebar with a tree view of the document structure. The main area displays a document in Gothic script. A specific section of the text is highlighted with a red box and a green line underneath it. The sidebar also contains a toolbar with various icons for editing and navigating the document.

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

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

The screenshot shows the Transkribus digital edition interface. On the left, the 'Metadata' tab is selected in the top navigation bar. Below it, the 'Structural' tab is active. The main workspace displays a page from a historical book. The title 'CATECHISMUS, Oder Bürger Unterricht' is visible, with the subtitle 'Christlicher Lehr / wie der in Kirchen und Schulen der Churfürstlichen Pfalz getrieben wird.' The text is presented in large, decorative Gothic script. A green rectangular selection box highlights a portion of the subtitle text. The left sidebar contains a tree view of the document structure, with nodes labeled TR, L, BL, W, and ... under the root node 'Document'. A list of element types is shown, including paragraph, heading, header, footer, drop-capital, floating, catch-word, footnote, endnote, other, and more, with 'TextRegion' selected. At the bottom of the interface, there are various toolbars for navigating the document and applying changes.

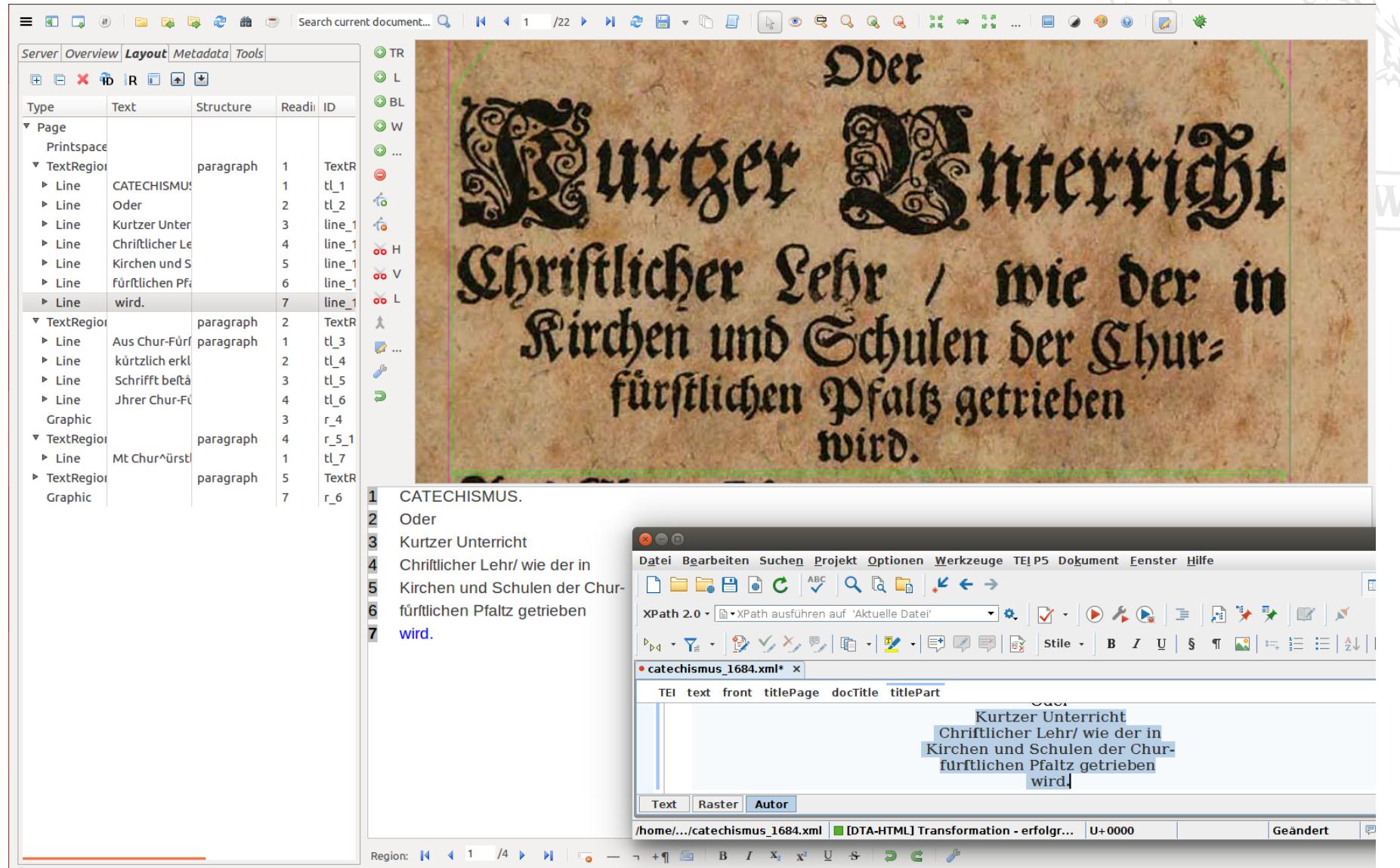
Workflows

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
Copy and paste text region by region	

-
-
-
-



Workflows



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
Copy and paste text region by region	
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
Copy and paste text region by region	
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: <http://kaskade.dwds.de/~matthias/ocr-d/>

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



Collaboration



- 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.

