



How to use Transkribus eXpert (deprecated)

March 2026



1. Downloading and Installing Transkribus eXpert (deprecated).....	3
2. Troubleshooting Transkribus eXpert (deprecated).....	4
Windows.....	4
Mac.....	4
Linux.....	5
Run Transkribus via command line.....	5
Using a proxy server.....	6
Known Issues.....	6
3. Creating a Collection.....	9
4. Uploading.....	10
Images (Single document).....	10
PDF.....	10
Private FTP.....	11
IIIF manifest.....	11
DFG Viewer METS.....	11
5. Downloading.....	13
Transkribus Document.....	14
PDF.....	15
TEI.....	15
DOCX.....	16
TXT.....	17
Tag Export (Excel).....	17
Table Export into Excel.....	17
Page metadata into Excel.....	17
6. Automatically Transcribing your documents.....	19
7. Choosing a Model.....	21
8. Automatic Layout Recognition.....	22
9. Advanced Layout Configuration Settings.....	24
Layout Model.....	24
Baseline Detection Settings.....	24
Region Detection Settings.....	25
10. Editing Layout Manually.....	27



11. P2PaLA.....	28
12. Transcribing Manually.....	28
13. Manually Drawing Tables.....	30
14. Training Baselines Models.....	32
15. Training Text Recognition Models.....	35
Conventions.....	36
Abbreviations.....	37
Retraining with PyLaia.....	37
16. Model Setup and Training.....	39
Model Setup.....	39
Training Set.....	39
Validation Set.....	40
Engine.....	41
Advanced Parameters.....	41
17. Character Error Rate.....	48
Learning Curve.....	48
18. Computing Accuracy.....	50
Compare Text Versions:.....	50
Compare and measure the CER:.....	51
Compare Samples.....	52
19. Structural Tags.....	55
20. Textual Tags.....	58
21. Searching.....	63
Fulltext Search.....	63
Fuzzy Search.....	64
22. Workspace Management.....	65
Managing Collections.....	65
Managing Documents.....	65
Managing Models.....	68
Sharing a Model.....	68



1. Downloading and Installing Transkribus eXpert (deprecated)

The Transkribus eXpert is the desktop version of Transkribus.

Please note that Transkribus eXpert (desktop software) is no longer being updated, and all new features will be exclusively available on the Transkribus web app.

Transkribus eXpert is written in Java. You need a Java JDK version 11 or higher installed on your computer for Transkribus eXpert to work.

We recommend installing the latest JDK from [here](#). (If you need to check your Java version, [read more here](#)).

After checking you have the correct Java version installed, download the ZIP file of Transkribus eXpert [here](#).

After download, you see the ZIP File in the download directory of your computer. Unzip the file before you try to start Transkribus.

Open the Transkribus directory. There you will find the executable files for your operating system. Start Transkribus from your user interface via doubleclick:

- Windows: **Transkribus.bat** or use Transkribus.exe
- Mac OS – Apple: **Transkribus.command**
- Linux: **Transkribus.sh**



2. Troubleshooting Transkribus eXpert (deprecated)

Windows

- If you do not have “Administrator” rights, Windows will produce a warning message, such as: “Your Computer is Protected by Windows.” Do not confirm, but go to “More Information”. There you can agree that this is not malware and that you want to run Transkribus on your computer.
- Transkribus eXpert requires Windows 64 bit. If Transkribus does not start, it may be because a 32 bit operating system is installed. Then a new 64 bit operating system would have to be installed.
To find out your operating system, go to Windows Settings – About or run Transkribus from the console to get a more detailed error message. To do this, type ‘cmd’ in the Windows search, open the ‘Command Prompt’. Then change to the Transkribus folder with cd ‘Transkribus directory’ (e.g. cd C:\Users\Username\Transkribus-1.20.0) and run ‘java -jar Transkribus.exe’.
- Make sure you have a 64-bit version of Java installed on your computer that is not older than Java 11.
- To find out the Java version on your own computer, start the console again and type java -version. Then the current version will be displayed. If no java version is displayed, then the operating system cannot (yet) find the installed Java. Then a reboot or setting the ‘PATH’ variable can help: <https://michster.de/wie-setze-ich-die-path-umgebungsvariablen-unter-windows-10/>
- In case of Java problems, the Java package can also be copied directly into the Transkribus folder. The Java package must be renamed to ‘jre’, so that Transkribus then uses this Java.

Mac

- If you run the program the first time, it may not start because it is a non-signed application (“... can’t be opened because it is from an unidentified developer” message).
In this case, right-click (or control-click) the application and choose “Open”. In the appearing dialog box, click “Open” again.

Alternatively, right click the Track Pad to open the Context Menu and add a security exception for Transkribus.



Another option: right-click on the program icon -> Open (in the context menu) -> via Terminal.

- If the app won't start at all: you can try to move the application out of the Download folder, i.e. copy or move the package to another destination folder like the Desktop or the Application folder.
- Error message when trying to start the app from a terminal using 'open -a Transkribus.app' is: LSOpenURLsWithRole() failed for the application /Users/xxx/Desktop/Transkribus.app with error -10810.

A workaround for starting the program anyhow until then is to start a new terminal (search for 'terminal' after hitting cmd + space), then cd into the directory where Transkribus was unpacked, e.g. the 'Downloads' directory. Then start the program directly from the start script that is included in the Transkribus.app package:

```
./Transkribus.app/Contents/MacOS/Transkribus
```

Linux

- If the OS is based on Ubuntu 17.04, installing libwebkit is necessary: `sudo apt install libwebkitgtk-1.0-0`

Run Transkribus via command line

Transkribus is contained in the main jar file `Transkribus-<version>.jar`. To run the program from command line type: `java -jar Transkribus-<version>.jar`. Note: Java 11 or one of the following versions is needed. Some problems (mainly Java heap space) occur because a Java 32 bit version is installed on a 64 bit operating system. To run the scripts in Mac (or Linux) you may have to make them executable from the command line: (any version before 0.6.8)

- [Mac console basics](#)
- change into the program folder using 'cd' commands
- `chmod +x Transkribus.command` (or `chmod +x Transkribus.sh` for Linux!)

Furthermore, you will find several files in the Transkribus package copied to your computer:

- `config.properties` can be modified to adjust simple appearance properties
- `virtualKeyboards.xml` can be used to specify a set of virtual keyboards



- logback.xml can be modified to adjust logging properties (for expert users only)

The 'libs' subfolder contains the necessary libraries for all platforms. Currently supported are:

- Windows 32/64 bit
- Linux 32/64 bit
- OSX 64 bit

Using a proxy server

If you keep getting the error message "Login failed: already connected" the problem can be the proxy server.

When the program has started, click on the home menu button on the top left and select "Proxy settings...". In the following dialog you can set the proxy host, port, user name (optional) and password (optional). This is the recommended method for using a proxy server.

Alternatively, you can edit the start script (e.g. Transkribus.bat on Windows, Transkribus.sh on Linux) to include the environment variables for the proxy server:

```
java -Dhttps.proxyHost=<proxyserver>  
  
-Dhttps.proxyPort=<proxyPort>  
  
-Dhttps.proxyUser=<user name for proxy>  
  
-Dhttps.proxyPassword=<password for proxy>  
  
-jar Transkribus-0.14.0.jar
```

However, editing this file will be necessary on each update of Transkribus.

Known Issues

Logging in to the server is not possible via Transkribus eXpert, but on the website it works.

- A reason for the error message "Already connected" is that your Java installation may be outdated and can't establish a secure connection to the server. You can check your installed version by opening a terminal/command line and entering "java -version".
If you encounter this problem, try updating Java on your machine. We



recommend a current version not older than Java 11 (Oracle or OpenJdk).

The Mac version of Transkribus includes a Java runtime. If you encounter this problem on a Mac please download a new package from

transkribus.eu/TranskribusX/releases/ and update your installation.

If the error persists please contact info@readcoop.eu, ideally including the log file of your installation (from the Transkribus directory: logs/TrpGui.log) and/or information on your Java version and operating system.

- Note: the Mac version of the expert client comes with a java shipped within the application. If this java version is outdated, you can try to delete or replace it with an updated version. To find the files in Mac finder, right-click (or cmd-click) on the Transkribus application in your programs view, click “show packages contents” in the context menu, then go to the subfolder “Contents/MacOS”. There, the subfolder “jre” contains that java version. If you delete this folder, the application starter will try to find java on your system.
- You may have to configure a proxy server via ‘main menu’ – ‘Proxy settings’.

Wrong JAVA Version on Mac

- After opening the command file on the Mac, Transkribus says that there is a wrong Java version installed. However, there is the most current version of Java installed. The problem is that an older version may still be installed and set as default Java. You can check the default version by opening the terminal and typing ‘java -version’.
- To solve the problem you can either download the latest Java versions as a .tar.gz package from here: <https://www.oracle.com/java/technologies/downloads/> and unpack it into the Transkribus folder – the Transkribus.command file will automatically check for java installations in its sub directories.
- Alternatively you can try to set your Java 11 (or one of the following versions) installation as default one in the command line following e.g. the instructions here: <http://myshittycode.com/2014/03/17/mac-os-x-setting-default-java-version/>
- Or you just try to re-install the latest Java JDK from the installer at <https://www.oracle.com/java/technologies/downloads/> and cross fingers that it replaces the old one.

Java Heap space / No more handles



- 32 bit Java on a 64 bit Windows OS -> install 64 bit Java from here: <https://www.oracle.com/java/technologies/downloads/>
- Too little RAM: Try to allocate more main memory by opening Transkribus.bat and set e.g. java -Xmx2048m -jar Transkribus-1.14.0.jar
Start Transkribus with this bat file

Logging in is prevented by the Firewall of your Internet Provider

- Some IT departments are blocking the SSL port 443 and/or unknown applications via a firewall. Please check with your IT department if that might be the case.

Avira or Norton Antivirus detects a threat and is blocking the zip file from being unpacked.

- Solution: This is a false alarm which Norton gives when encountering software it is not familiar with (WS.Reputation.1). You should be able to restore the file from quarantine by following the instructions from this [resource](#).

Transkribus does not start on (Fedora) Linux - 'MOZILLA_FIVE_HOME not set' error message

- The package "libwebkitgtk" may not be installed. On Fedora you can install the package using dnf on the command line (use "yum" instead of "dnf" in older versions of Fedora): sudo dnf install webkitgtk

Dark mode on Mac

- The dark mode on a Mac can sometimes cause problems with Transkribus, so if Transkribus does not run properly on your Mac and you have the dark mode switched on, please try to turn it off. You might need to reinstall Transkribus after turning off the dark mode in order to take the change into effect.



3. Creating a Collection

A collection in Transkribus eXpert is a folder that contains all the documents pertaining to a project. Before uploading the documents, you need to create the collection.

To create a new collection in Transkribus, go to Server Tab in the Managing & Tools Bar. By clicking on the “Collections” button, a window will open: all your collections are listed there. Click on the green plus “Create” button at the bottom left corner of the window, type the title of the collection and click OK.

Double-click on the collection already created, and it will open. In the same window, you can delete the collection and edit the metadata by selecting, respectively, the “Delete” and “Modify” buttons at the bottom of the window.



4. Uploading

A document in Transkribus is a group of images that belongs to the same organised unit (file, manuscript, book...).

To upload a document to Transkribus eXpert, click the “Import document(s)” button in the Menu Bar at the top. The Import button icon is a folder with a green arrow pointing to the left. A window with all the import options will open up. The collection to which the documents will be uploaded is, by default, the one you have opened in the Managing and Tools Bar, but you can change it in this window by clicking on “Add to collection.”

The five upload options in Transkribus eXpert are:

Images (Single document)

Use this option to upload single folders with images. You can upload a folder at a time. The accepted image file formats are JPEG, PNG and TIFF. JP2 files are not supported.

In this option, you can only choose folders, not the image files within them. Therefore, before starting the upload, make sure that all your images are in the folder you are choosing and that the folder does not contain pdfs or subfolders.

In Transkribus, click the folder icon in the Upload window to search for your local folder with the images. If you open the local folder, it will seem as if the folder would be empty. That is just a matter of displaying. Please select the folder one hierarchy above (not the images themselves) and confirm.

By default, the document title will be the folder name. Type the document title near “Title on server” if you want to rename it (you can also do it later by modifying the Document Metadata).

To start the upload process, click on “Upload”.

PDF

Click the folder icon in the Upload window to search for the PDF file on your computer. Choose the PDF file you want to upload and then click “Upload”.

Each page of the PDF will be extracted and uploaded as an image of the Transkribus document.

You can only upload one PDF file at a time with this option. If you have a couple of PDF files, you might want to use a private FTP client (see below).



Private FTP

The FTP option enables you to upload more than one folder/PDF file at a time. It works with PDF, JPEG, PNG and TIFF files.

Before selecting this option, make sure you have a FTP client installed on your computer.

In the "Upload" window, choose the FTP option. Click the link next to Location: ftp://transkribus.eu. Your standard FTP Client pops up and asks for your password. Please use your login credentials you also use with Transkribus.

In the FTP Client, navigate the current local directory and the local directory tree displayed on the left side of the main window. You can either upload folders with files in them or choose single files. Just drag and drop them from the left to the right side of the window. After everything is uploaded correctly, you can use the FTP client.

After refreshing the upload window with the blue circling arrows in Transkribus, you can see the files that you just have uploaded. Choose the files you need in your collection and press "Upload".

You can see the upload in your "Jobs" page. It might take a while, depending on how many files you upload. Afterwards, you will see the files in your collection, ready to work with.

IIIF manifest

If the archive/library digitised the document you are interested in and has used the IIIF standard to display it online, you don't need to download the images and upload them to Transkribus. Just search for the URL of the IIIF manifest in the online catalogue of the archive/library, copy it and insert it in the provided field in Transkribus; then click "Upload".

For now, Transkribus works only with the IIIF Presentation API version 2.1. If your upload fails, it could be because the archive/library uses another version to display its digital objects.

DFG Viewer METS

Simply insert the URL of the DFG Viewer METS in the provided field and press "Upload".

All documents uploaded to Transkribus are private by default. They are stored on the servers of READ-COOP SCE (i.e. the company that develops and maintains the software). The servers are all located in Innsbruck, Austria, in a GDPR-compliant



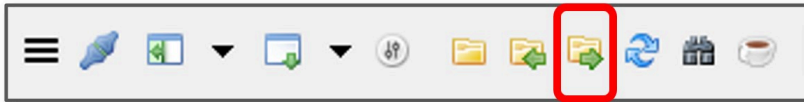
manner, and the data may be processed according to the [terms & conditions](#) on the READ-COOP SCE website.



5. Downloading

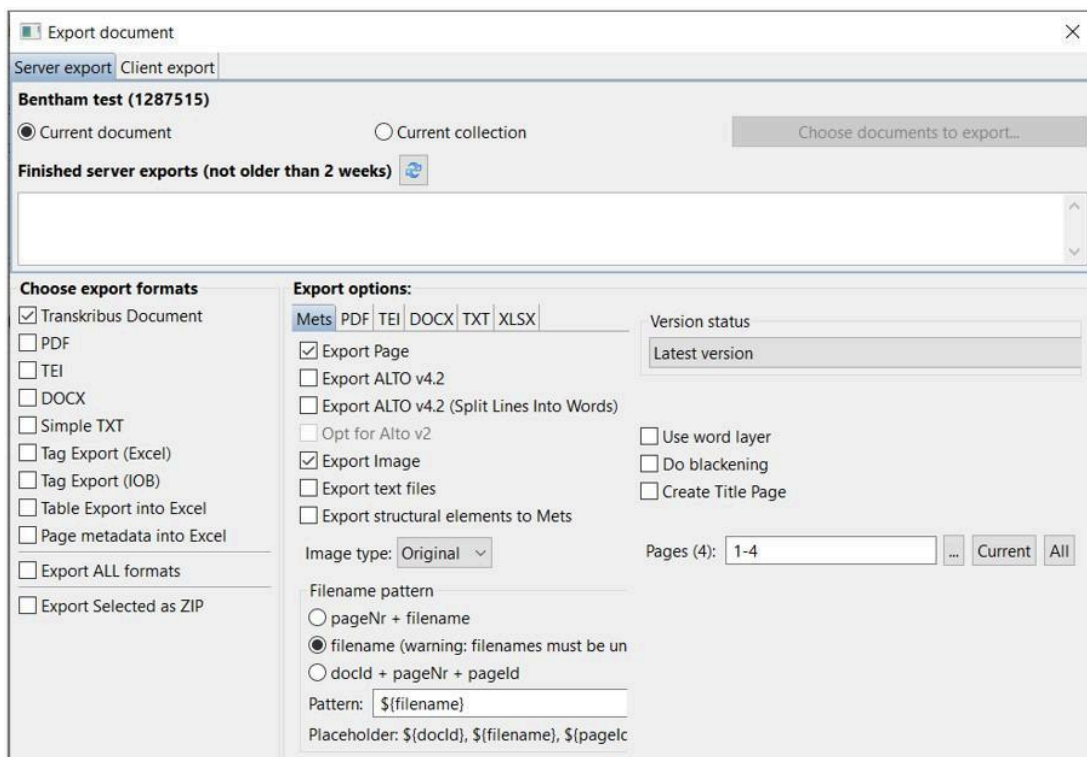
If you want to work with your images and transcriptions outside of Transkribus, you can export your documents from the platform. Different export formats and features are available to suit your needs.

To open the Export window, click on the folder icon with the green arrow pointing to the right in the Main Bar:



The Export document that opens up has two tabs/options between which you have to choose:

- **Server export:** the export will be processed on the Transkribus server, and you will receive a link to download your files. The export will not slow your computer down, and the process will not be interrupted if you switch your computer off. After starting the download, you can check the progress of your export by clicking the “Jobs” button in the “Server” tab.
- **Client export:** the files will be saved directly to your computer. Please choose where you would like to save the exported files: type the file location in the “Base folder” box at the top of the window.





These are the available export formats:

Transkribus Document

If you export your transcription as a Transkribus Document, you will produce a METS (Metadata Encoding and Transmission Standard) file containing the links to PAGE, XMLs, ALTO and/or image files, depending on which options you choose.

A METS file is like a container which includes all the background information about a file. More detailed information about METS can be found at:

<http://www.loc.gov/standards/mets/>

In conjunction with the METS file, you can export your document in these formats:

- **PAGE:** is an XML-based page image representation framework that records information on image characteristics in addition to layout structure and page content. The complete format definition used in Transkribus can be accessed [here](#).
- **ALTO:** is a special output format which allows you to input the exported document into other programs working with this format. The format is similar to XML and works for OCR, for example. It is often used in combination with METS for the description of the whole digitized object and the creation of references across the ALTO files, e.g. description of the reading sequence.

More information about ALTO can be found at:

<http://www.loc.gov/standards/alto/> With the “Split Lines Into Words” option Transkribus will divide the lines into words. The program does this by analysing the spaces between words, even if no word segmentation has been performed previously.

- **Images:** choose this option to download the image file of each page of the document/the selected pages.

In Image type, you can choose to download the Original (the image you uploaded) or the JPEG compressed version of the image (the one you see in the Transkribus Image Window).

Under “Filename pattern”, you can choose how the filename will be composed. The second option, “filename”, is the standard one. With this option, the exported file will have the same name as the document you imported. This is important if you want to match local transcripts with the images in Transkribus. So if you export a document, then adjust it externally, and after that upload it to



Transkribus again, the program will need to have two similar filenames in order to recognise the file properly.

PDF

When you export a PDF file, you can choose between these options:

- “Images plus text layer”: you will see two layers in the exported PDF document: OCR (the transcribed text) and image (image of the document).
- “Images only”: you will produce a PDF file with the document as an image. This means that you will not see the transcribed text.
- “Extra text page”: the transcribed text will be added to the PDF as an extra page after each image.
- “Highlight tags”: select these options to highlight the tags in the exported PDF file. The tags will be shown in the same colours used in Transkribus. At the end of the document, there will also be a symbol legend to explain the signification of the different colours.
- “Highlight article”: the articles will be highlighted with different colours in the exported PDF.
- “PDF/A”: for long-term preservation.

You can also choose the font and image type to use in the PDF.

TEI

This option is for people working with the Text Encoding Initiative (TEI). The Text Encoding Initiative is a text-centric community of practice in the academic field of digital humanities, operating continuously since the 1980s. More information at: <http://www.tei-c.org/index.xml>

You can choose to create the TEI XML file using the the XLS from Dario Kampaskar (available here: <https://github.com/dariok/page2tei>) or the “Client Export” format. You can try both and decide which one best suits your needs.

With the “Client Export” format, you can flag the option to export the predefined tags and attributes only: it creates a valid TEI but note that all the tags and attributes created by you will be ignored.

Transkribus enables you to choose the zones you need (no zones; zone per region; zone per line; zone per word). Furthermore, you can choose between line tags and line breaks to tag lines.

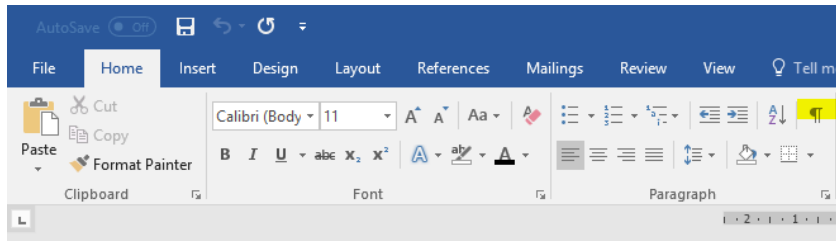


DOCX

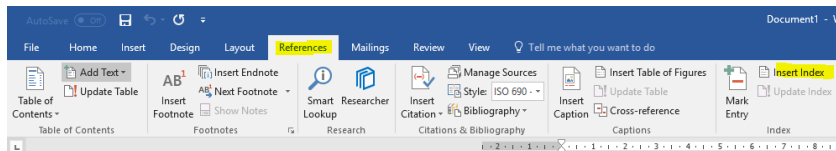
By choosing this option, you will get your transcriptions in Word files. You can select options relating to line breaks, abbreviations and more according to your needs.

Select “Export selected tags” to make the tags visible in the exported DOCX file. After the export of the Word document, please open it and do the following:

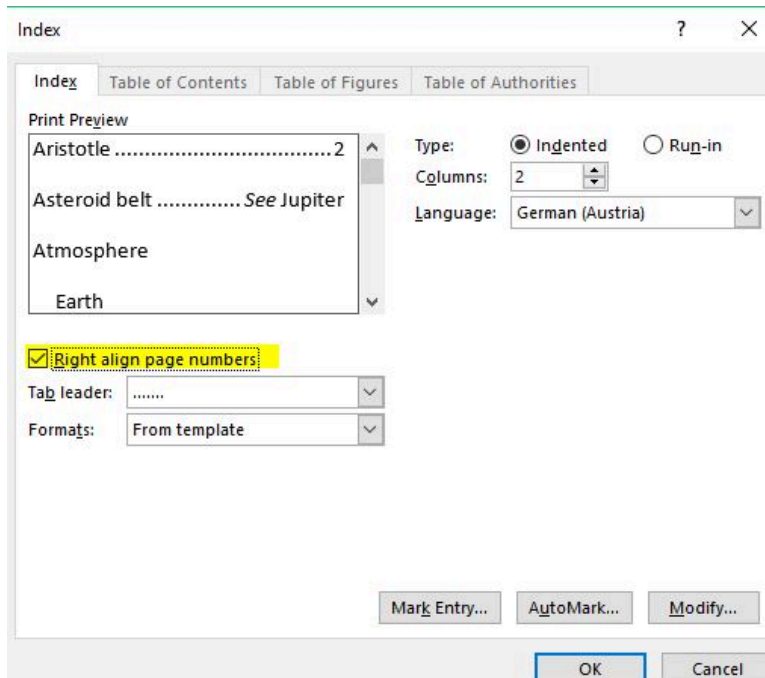
- a. Click on the paragraph button in the Home menu of Word



- b. Go to “References” and choose “Insert Index”



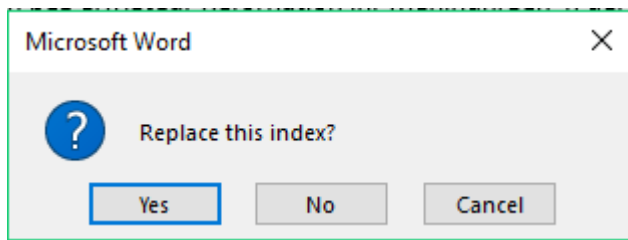
- c. The following Office window will open up



- d. Select “Right align page numbers” and press “OK”



- e. A confirmation window will pop up: click “Yes”



- f. An overview of the tags should now appear at the end of the document. If the overview of the tags does not appear, click “Update Index”. This should solve the problem.

TXT

If you do not usually work with Microsoft Word, it is possible to export your transcription as a simple TXT file.

You can choose to split the text into text files from a start tag to an end tag and create several text files; those files can be named according to one or more attributes of the tag.

Tag Export (Excel)

If you would like to export the tags you assigned to your transcription, select this option to produce an Excel file with individual tabs for each tag category and one tab with an overview of all the tags.

As described above, you can also export the tags in PDF and DOCX files.

Table Export into Excel

If your document presents tables, use this option to export them in Excel format. Each table will be exported as a separate sheet of the Excel file. However, you can check the option “Create one large table” if you prefer to have all the tables of the selected pages in one table in one Excel sheet. You can also choose to export only one column of your table with the cells' image snippets.

Page metadata into Excel

To export the page-related metadata in the Metadata-page tab in the Excel format.

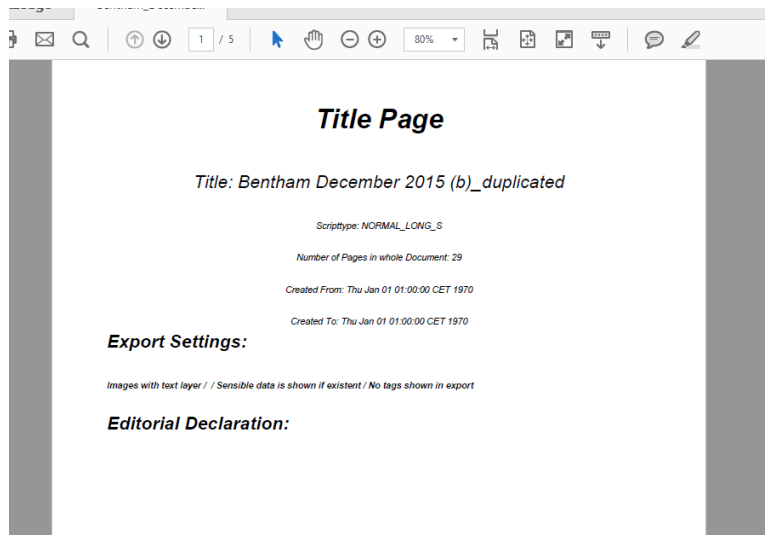
In addition to the export format, these options are selectable during the export:

- **Version Status:** to export a particular version of the document. If you select “Ground Truth”, for example, Transkribus will export only those pages of the document which you have marked as “Ground Truth.”



For the export, the program consults previous versions of your document. This means that if you choose to export all “In Progress” pages, the program will export all pages which have been marked as “In Progress”, even if their status is now updated. The program will export the latest “In Progress” version of your document. If you would like to export a specific former “In Progress” version of your page, open this version of the page in Transkribus. Open the Export window and select “Loaded version for current page” (available only for the Client Export). In the “Pages” option, select “Current” before confirming.

- **Word Layer:** if checked, the text from word layer segmentation will be exported (it works only if you have previously selected the “Add estimated word coordinates” during the Text Recognition).
- **Blackening:** If you have blacked out sensitive sections of your transcription, these words or phrases can also be hidden in the exported files. To do this, select “Do blackening” in the export options. Note: this option only works for Word, PDF and METS files.
- **Create Title Page:** with this option, a title page based on the information added in the “Document” tab within the “Metadata” tab is created. In the “Document” tab, you can add information about the title, author, language and date of your document. You can also create an Editorial Declaration to explain how exactly your document has been transcribed.



- **Pages to be exported:** select the number of pages you wish to export. You can export all the pages in your document or just the current page.
- **All tags/chosen tags:** to choose which tags you want to export.



6. Automatically Transcribing your documents

To automatically transcribe your documents, go to the “Tools” tab, under the “Text Recognition” section and click on the “Run” button. In the pop-up window, choose the page(s)/document(s) to process and then click on “Select HTR model.” Here you can choose the most appropriate text model for your documents.

A text model is the AI algorithm trained on a certain number of data (images and transcriptions), able to detect the most probable sequence of characters for each segmented text line. There does not exist a general model for all the handwritings, so you need to choose the most appropriate one for the script and language of your documents.

You can select both the public models made available by the Transkribus community and team and the private models trained by yourself. You can filter your search by engine, language and name.

Advanced settings you can select are:

- **Use existing line polygons:** use this option if you have corrected the line polygons manually because the computation of polygons from the baselines did not perform well on your documents.
- **Do polygons simplification:** to reduce the number of points of the line polygons.
- **Add estimated word coordinates:** add approximate bounding boxes for each word in the line (you can then decide to show/hide the word boxes with the eye-icon in the Main bar at the top).
- **Restrict on structure tag:** limit the Text Recognition only to the text regions tagged with the selected structural tag. You can decide if you want to keep or delete the text in the other regions.

After having selected the model, click "OK" to launch the recognition. You can check the status by clicking on the “Jobs” button in the top Main Bar. When the recognition is finished, reload the page: the automatically generated transcript will appear in the text editor.

When you launch the Text Recognition, first, the images are automatically segmented into text regions and lines. This step, called Layout Recognition, connects the text and the image. If your documents have a complex layout (e.g. tables, newspapers, postcards, marginalia, multiple columns...), it could be convenient to run the Layout



Recognition as a separate step in order to check and correct it before the Text Recognition. If this is your case, take a look at the section on Layout Recognition.



7. Choosing a Model

The most important thing for good transcripts is to select a model suitable for your documents. There does not exist a general model for all the handwritings and for the next few years, it is expected that specialized models will remain necessary.

When you click "Select HTR model", a window opens: on the left side of the window, you can see an overview of the available models; on the top right side of the window, the details of the model are shown.

When choosing a text model, you need to consider the following:

- the type of material, handwritten or printed;
- the language;
- the period;
- the type of script;
- the Character Error Rate.

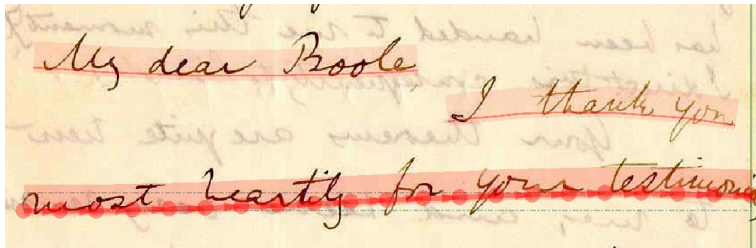
All the models that can be used for Text Recognition have been trained with PyLaia, which is the Handwritten Text Recognition engine currently available within Transkribus. It has been developed by UPVLC (Universitat Politècnica de València) and is open-source.



8. Automatic Layout Recognition

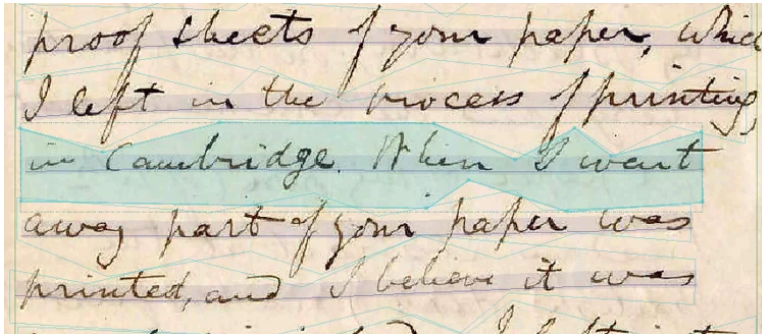
Layout Recognition is the segmentation of the image into text regions, lines and baselines to connect the text and the image.

The text region is a rectangle, encasing all of the handwritten text contained in the image/page.



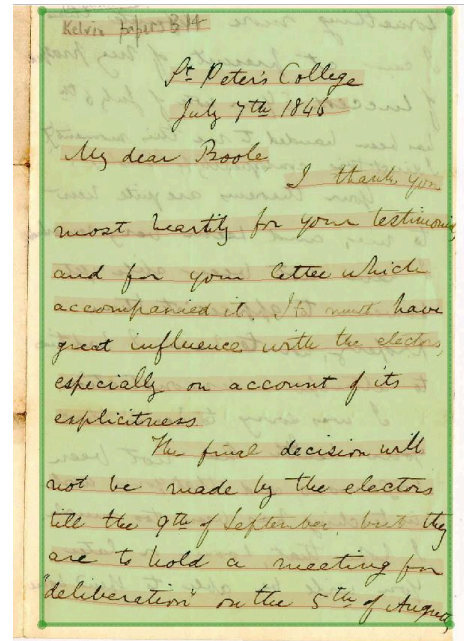
The baseline is a polyline, running along the bottom of the handwritten text line, and is the most important reference point for text recognition.

The lines are regions located within a text region and can be described as polygons, encasing all of the handwritten text in a line.



The Layout Recognition is performed automatically when you start a Text Recognition job, but it can also be run as a separate step.

To run the Layout Analysis as a separate step, go to the "Tools" tab in the Managing & Tools Bar (on the left side of the screen). The section we are interested in is named "Layout Analysis".





Layout Analysis

Method:

Current page

Pages (24):

Document Selection

Document 2.pdf (1371029)

Current collection

Find Text-Regions Find Lines

Model: Horizontal Text Line Orientatio

Select the current page, the pages or the document(s) you want to process and then click "Run" to launch the layout analysis. The Layout Analysis will be performed with the default settings (Horizontal Text Line Orientation model; General region detection method).

To check the progress of the job, click on the "Jobs" button. When the job is finished, reload the page(s) and the text regions, lines and baselines will appear in the Image Window. You can also see the layout structure in the "Layout" tab, in the Managing&Tools Bar.

If the automatic Layout Recognition has performed poorly (e.g. has missed some lines or the clustering of lines in text regions is not correct), you can change the advanced configuration settings.



9. Advanced Layout Configuration Settings

We recommend first trying the automatic layout recognition with the default settings. However, depending on your material, you may need to change some of them to get optimal results.

To access the Layout Configuration Settings, go to the "Layout Analysis" section in the "Tools" tab in the Managing & Tools Bar (on the left side of the screen).

For a complex layout, it could be a trial-and-error process. Try to modify some of the settings depending on the issues with your type of material and run the Layout Analysis on a few pages to test the settings out.

Click on "Configure", to the right of "Method", to open the Layout Analysis Configuration window. The settings you can configure are:

Layout Model

The first thing to choose is the *Layout model* to use:

- **Horizontal Text Line Orientation** (selected by default): when your documents have a homogeneous layout, i.e. only horizontal and vertical lines.
- **Mixed Text Line Orientation**: when your documents have a heterogeneous layout, i.e. lines in all directions.
- **Custom Baselines model**: when you have trained a Baselines model specific to your document typology.

Click on the model's name to open the window with the list of all the public models as well as the private baselines models you have trained.

Baseline Detection Settings

These options offer you the possibility of setting parameters for baseline detection. They are particularly useful if too few/too many baselines have been recognised or if they have been joined or separated when they should not. For each parameter, you can choose one of the three suggested values (Low, Medium, High) or customise the value manually.

- **Minimal baseline length**: it indicates the minimum length of lines in pixels. Lines shorter than this length will not be detected.
- **Baseline accuracy threshold**: in the first stage of the layout recognition, each pixel is labelled as baseline, separator or other. The baseline accuracy threshold



applies to the baseline labelling at this stage. It ranges between 0 and 255, and higher values enforce higher accuracy in the detected baselines.

- **Use trained separators:** separators are small vertical lines drawn beside each baseline; they mark the beginning and end of each baseline (do not confuse them with actual separators in printed document images). As for the baseline accuracy threshold, the separator threshold refers to the first stage, when pixels are labelled.

The separator threshold ranges between 0 and 255: 0 means that separators are not used at all; with a higher value, separators are used, thus, nearby baselines tend not to be merged.

Usually, low values are sufficient to prevent a connection between nearby baselines.

- **Max-dist for merging:** in the second stage, the algorithm tries to merge nearby baselines but only when their distance is smaller than a set value.

Set it to: "Low" to merge only the closest lines (closer than 0.5% of the image width); to "Medium" to merge lines closer than 1% of the image width; to "High" to merge lines that are quite far apart, but closer than 5% of the image width.

"Medium" should work well in most cases.

- **Image scaling:** You can decide to Upscale low-resolution images or Downscale high-resolution images.

We suggest trying this feature only when Layout Recognition does not work with the default setting (for instance, it detects no or few lines).

Region Detection Settings

After the baselines are detected, they are clustered in text regions. There are two **clustering** methods available:

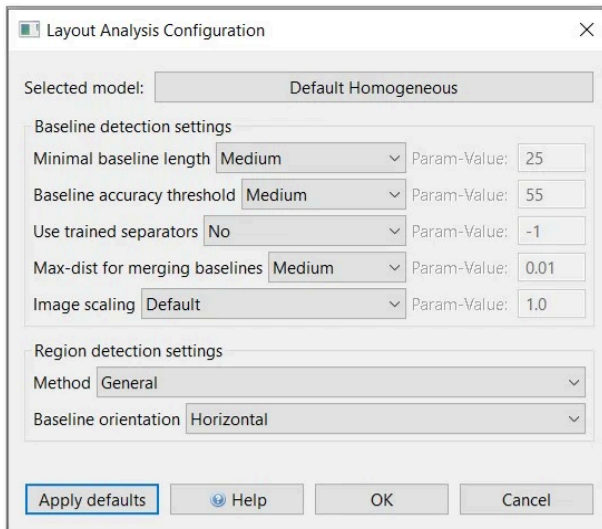
- **General** (default): it clusters the lines from left to right.

To improve the clustering, set the **text line orientation** to "**Horizontal**" if your documents have only horizontal lines or to "**Mixed**" to make the algorithm assume that lines are rotated by 0, 90, 180, and 270 degrees.

- **Custom:** it is a simple agglomerative clustering based on the leftmost point of each line. It clusters lines based on their distance.



You can choose to have one text region, few, medium or many text regions per image.



It may be the case that you have drawn the text regions of your interest by hand (e.g. for tables) and you just want to detect the baselines while keeping the existing text regions. In the "Layout Analysis" section in the "Tools" tab, untick the "Find Text Regions"-option before starting the layout analysis.

When keeping the existing text regions, it is also possible to:

- **restrict on structure tags:** limit the baseline detection to the text regions tagged with certain structure tags.
- **split lines on regions:** make the lines strictly obey the region border, preventing lines close to each other but belonging to different regions from being merged together as one long line.

You can increase the minimum overlap fraction between a detected line and an existing text region in case you want the line to extend slightly beyond the region border. If a line overlaps with multiple regions, however, the region with the most overlap is chosen.

For more information about the Transkribus LA algorithm and setting, consult [this page](#).



10. Editing Layout Manually

Sometimes it may be necessary to make some manual corrections to automate layout recognition, especially when the layout of your documents is complex.

When manually transcribing or correcting the automatic transcription, you may notice some empty regions, some lines too short, or the automatic recognition has detected two lines instead of one.

To correct the layout, use the Canvas Menu (to the left of the image). Most of the common edits are:

- Modify the shape of a text region: click inside the region so that it is highlighted, and drag the border of the text region as needed.
 - Split a shape into two: you can use the Split Shapes buttons in the Canvas Menu. The “H-button” splits the shape horizontally. The “V” button splits it vertically. The “L-button” allows you to split it with a customisable line.
- Delete a shape: click on the shape you wish to delete, and click the red “Remove a shape” button in the Canvas Menu.
- Merge two shapes: sometimes, the automatic layout analysis creates two shapes where only one is needed. In this case, you can easily merge the two together. Hold down the “CTRL” button on your keyboard and select both shapes (they need to be the same type of shape). Click the “Merges the selected shapes” button in the Canvas Menu.
- Add a text region: click the “+TR” button in the Canvas Menu. Click once to start drawing your region and double-click to finish it.
- Add a line: click the “+L” button in the Canvas Menu. Click once to start drawing your line and double-click to finish it.
- Add a baseline: click the “+BL” button in the Canvas Menu. Click once to start drawing your baseline and double-click to finish it.
- Correct baselines: click on one point of the baseline and drag it to the new position. Remember that baselines are most important for HTR and it is important that they run at the bottom of the words; lines do not need to be corrected.



11. P2PaLA

This functionality has been replaced by Field Models, which require less training data and provide more precise outcomes than P2PaLA.

Instead of using this function you can now train and use Field Models in the Transkribus Web App: <https://help.transkribus.org/field-models>.

12. Transcribing Manually

The Transkribus interface also facilitates the manual transcription of documents. With the text editor, you can easily transcribe your documents with the image side by side.

The automatic Layout Analysis creates a correspondence between the lines in the image and the lines in the text editor, so while you are transcribing, the image is automatically centred on the line you are working on.

You can also enrich your transcriptions with tags and download them in different formats (Docx, PDF, TEI, XML...).

There are various reasons why you might want to transcribe documents manually. One of them is the creation of accurate transcriptions to train a Text Model to transcribe new pages automatically.

After uploading your documents, select the page(s) you want to transcribe and the automatic Layout Recognition. Open the first page to transcribe and start transcribing line by line the text in the text editor.

To facilitate the manual transcription, remember that you can:

- enable the Transcription mode with the “Profiles” button in the Main Menu;
- change the position of the text editor (text and image side by side) with the “Change position of transcription widget” button in the Formatting Bar;
- change the font and font size of the text by clicking on “Transcription settings” in the Formatting Bar;
- modify the image viewing settings with Main Menu;
- use the virtual keyboard in the Formatting Bar to add special characters (you can customise the special characters of the virtual keyboard: select the tag “Custom” and then click “Edit”);



- press Enter to move to the next line, and use the arrow keys to move up and down in the text.

If you need to change the reading order, open the “Layout” tab in the Managing&Tools Bar: the line you are editing is highlighted. Within the layout tree, move the highlighted line to the correct position. By clicking the "R" button, you can reorder all the lines of an element according to their coordinates.

Another way is to enable the display of the line reading order in the image window: to do so, click on the “Shape visibility” button in the Main Menu. At the beginning of each line, the reading order value will appear; double-click on the value to change it.



13. Manually Drawing Tables

Please note that table models can only be trained in the Transkribus Web App and we recommend drawing tables in the Web App as you might otherwise encounter incompatibility issues, for more information on Table Models see: <https://help.transkribus.org/tables>.

Segmenting printed or hand-drawn tables using the Table Editor in Transkribus will add graphical lines to your image and assign a tabular structure to the layout of your documents.

Currently, tables must be manually drawn using the Table Editor in Transkribus. But if multiple pages follow the same table template, the table markup can be done on the first page and then copied to the remaining pages.

First, create text regions for any information not belonging to the table.

This refers to information at the top, bottom or sides of the page which is clearly not part of the table, such as page numbers, line numbers, dates and any other markings or annotations.

Then, you can create the table. In the Canvas Menu, select the "Add other item" button and then click "Add a table." Click on the top left corner of the table in the image and then click on the bottom right corner

You can now segment your table into rows and columns. To begin, make sure you are in "Selection mode:" press the "ESC" key on your keyboard or click the "Selection mode" button in the Main Menu. Click on the table that you have created.

To create rows, click the H-button in the Canvas Menu: move your cursor across the page and click wherever you want to create a horizontal line.

To create columns, click the V-button in the Canvas Menu: move your cursor across the page and click wherever you want to create a vertical line. Continue until all table cells are marked.

In some cases, it may be necessary to merge cells together in order to reflect cells spanning multiple rows or columns. To select cells to merge, hold down the "CTRL/CMD" key on your keyboard, click on the relevant cells in your table and then click the "Merges the selected shapes" button in the Canvas Menu. Please note that tables with merged cells should not be used as Ground Truth for table training.

If you focus on having the perfect table segmentation, correcting the shapes of some of the cells in your table may also be necessary. The segmented green lines should then correspond to the lines of your table as far as possible. In order to do so, select



the table cell you wish to edit, click and drag the green dots to move the position of the lines.

Depending on the layout of your table, you might want to treat the spine of the book like an extra column. You can also mark up this column on the table cell level using the "book-binding" tag in the "Metadata/Structural" tab.

If the table layout of several pages is similar, it is possible to transfer the table format from one page to others. To do this, open "Other segmentation tools" in the Canvas Menu; choose "Copy regions (texts or tables) to other pages;" define the pages the layout should be copied to in the appearing window and confirm with "OK". The table layout will be copied to the indicated pages. To definitely run the tool, unselect "Dry run". It might be that the position of the table on the new pages will need to be correct. To do so, select the whole table and then move it by holding the CTRL + SHIFT on your keyboard.

Before manually or automatically transcribing the table, the next step is adding baselines. The baselines should reflect the logical flow of text and can therefore run over the cell borders if necessary. You can either draw the baselines by hand or use the automatic Layout Analysis tool.

You may find that the automatic layout tool on table cells strictly obeys the cell borders. Baselines stretching multiple cells are divided. You can use the merging tool to combine those partial baselines. In case you want to merge baselines stretching more than one cell, you need to move them first to the same cell, select them and use the merging tool. In more detail, open the "Layout" tab in the Tools&Managing Bar, and select, in the image, the line that belongs to the wrong cell: automatically, it will highlight the corresponding line in the layout tree. Within the layout tree, move the highlighted line to the right cell (probably the previous or following cell). Now that both the lines belong to the same cell, you can select both and click the "Merges the selected shapes" button in the Canvas Menu. Customise your table models to suit your documents by adjusting how tables are recognised and data is extracted. With customised training, Table Models can recognise columns, rows, headers and data cells, and more.



14. Training Baselines Models

The default Layout Recognition tool (Preset Model) works well for most document typologies but may not be as accurate with documents with complex layouts, such as newspapers, postcards, registers, annotated documents, etc.

When the default Layout Analysis is unsatisfactory for your documents, you can train a Baselines model specific to your document typology. After the training, you can apply your customised Baselines model to your documents, which will be segmented following the examples you provided for training.

The first step is to prepare the pages on which to train the Baselines model. A good number to start with is 50 pages, but the model efficiency depends on the complexity of the layout. After the first training with 50 pages, you could decide if the Baselines model is good enough or if it needs more training material.

To prepare the pages, it is only necessary to segment, automatically or manually, the text regions and the baselines. Depending on the layout complexity, there are three options to segment the pages:

1. Run the default automatic Layout Analysis that you find under the “Tools” tab and then correct it manually using the Canvas Menu.
2. Draw the Text Regions manually using the “+TR” button in the Canvas menu. Then, under the “Tools” tab, run the automatic Layout Analysis to detect the baselines: before running it, remember to uncheck the “Find Text Regions” option. Finally, go through the pages and correct them manually using the Canvas Menu.
3. Draw both the Text Regions and the Baselines manually, using respectively the “+TR” button and the “+BL” button in the Canvas Menu.

Which option to choose depends on the document type and how poorly the default automatic Layout Analysis recognition performs.

No transcription is required to be added to the pages before the Baselines model training since it focuses only on the baselines, and the presence of transcribed text is irrelevant.

Once the 50 or more pages are segmented, it is time to train the Baselines model. Click on the “Tools” tab. Under the “Model Training” section, click on “Train a new model”. The Model Training window pops up, and on the right, you can choose which engine to train: for the Baseline model, please select “Baselines.”



Before starting training, enter the name and the description of your model. You can also modify the training parameters, i.e. the number of epochs and the learning rate. For the first training and if you are not familiar with machine learning, do not change these parameters.

You need then to select the pages you want to use to train the model, i.e. the pages you previously segmented into text regions and baselines. On the left, select the whole collection or the relevant pages. Click the Training button in the centre to add the selected pages to the Training Set. If you want to consider only the pages with Ground Truth status, select “Ground Truth only” in the drop-down menu on the right, under “Overview”.

Do the same for the Validation Set. The Validation Set should be around 10% of the Training Set, so we suggest, for the first training, including 45 pages in the Training Set and 5 pages in the Validation Set. If you want to automatically assign a percentage of the Training Set to the Validation Set, tick a percentage in the “automatic selection of validation set” option before clicking the “Training” button.

After completing this phase, you can start training the Baselines model by clicking on the “Train” button. Depending on the amount of training material, your training might take a while. Click on the “Jobs” button to check the Job status or your position in the queue (i.e. the number of trainings ahead of yours). You can perform other jobs in Transkribus or close the platform during the training process. If the Job status is “created” or “running”, please don’t start a new training, but just be patient and wait.

When the training is finished, the Baseline model will appear in the “Server” tab, under “Model Data”. To see it, please select “layout” instead of “text” as model output type in the second drop-down menu.

Double-clicking on the Baseline model name, you will see all the details and its learning curve. The “Learning Curve” graph shows the Baseline model’s accuracy. The x-axis indicates the number of Epochs, i.e. the number of times that the training data is evaluated. The y-axis measures the Loss, i.e. the percentage amount of pixels classified incorrectly. The program trains itself first on the Training Set; then, it tests itself on the pages of the Validation Set. For this reason, there are two lines in the graph. The blue line indicated the progress of the training; the red line indicated the progress of the evaluation on the Validation Set. Note that it is important that the two curves do not differ too much. If the two curves diverge, it is most likely that the Training Set differs too much from the Validation Set and the resulting model is not effective.

Underneath the graph, the two percentages indicate how the Baseline model performs on the Training Set and the Validation Set in terms of Loss. The Loss on the Validation



Set is the most significant value because it indicates how the Baseline model performs on new pages that it has not been trained on. Results with a Loss of 10% or below mean that the Baseline model is effective.

To apply the trained Baseline model to your documents, go to the “Tools” tab. Under the “Layout Analysis” top section, click “Configure”. The “Layout Analysis Configuration” window pops up: here, you can choose the Baselines model you trained.

In combination with the Baselines model, it is also possible to change the Layout Analysis settings (minimal baseline length; baseline accuracy threshold; use trained separators; max-dist for merging baselines; number of text regions).

Finally, click the “OK” button at the bottom of the “Layout Analysis Configuration” window. Your trained model has now been selected.

In the “Tools” tab, choose the pages on which to segment and click the “Run” button: the Layout Analysis job will now start. You can check its progress by clicking on the “Jobs” button under the “Server” tab. Once the job is finished, reload the page/pages and the text regions and baselines will appear in the images. No credit will be used to apply the Baseline model to your documents.



15. Training Text Recognition Models

Before starting the training of a Text Recognition model, you need to prepare the Ground Truth data, i.e. the images and the corresponding accurate transcriptions on which the model will learn.

Ground Truth is a term used in Machine Learning. In Transkribus, it is used to indicate the images and the corresponding transcriptions used to train the Artificial Intelligence. The transcriptions should be as accurate as possible because any mistake in the Ground Truth will train the model to learn something wrong.

Depending on the type of material and the number of hands, between 5,000 and 15,000 words (around 25-75 pages) of transcribed material are required to start. In general, the neural networks of the Handwritten Text Recognition engine learn quickly: the more training data they have, the better the results will be.

If you are working on printed material, 5,000 words should be sufficient to achieve a good Character Error Rate.

In the case of handwritten documents, our advice is to train the model on at least 10,000 words for each hand. Models trained on a large training data (more than 100,000 words) comprising many hands from the same period and region should be capable of recognising hands not seen in any way during the training: the results, however, will probably be somewhat worse than the Character Error Rate (which is measured on the Validation Data).

The Ground Truth should include examples of all the scripts that you want your model to be able to transcribe. It is possible to train models capable of recognising two or more hands, languages, types of writing or alphabets at the same time: however, all these variants must be present in a representative manner in the Ground Truth.

The pages to include in the Ground Truth are, therefore, important because they will affect the effectiveness of the model. If you want to train a model that recognises the hands of three different writers, you will have to transcribe about 10,000 words for each writer. In the case of a writer whose handwriting changed over time, the Ground Truth should comprise pages written over various years that are representative of the changes.

To create the Ground Truth, there are two ways:

1. **Manually:**

Run the Layout Recognition on the pages to be included in the Ground Truth; transcribe them accurately in the Text editor and save them as Ground Truth.



2. **Partly automatically, partly manually:**

If there is a Text Recognition model that works sufficiently well on your documents, but you would like to train a more accurate one, you can first run the model on your pages. Correct, then, manually the automatically generated transcriptions and save them as Ground Truth.

In both cases, it is important that the Ground Truth transcriptions are as accurate and correct as possible and that you are consistent with your editorial choices.

Conventions

The most common approach is creating a consistent transcript that accurately represents what you read in your document, including errors and punctuation. This is the case of a diplomatic transcription: combining words, upper and lower cases, superscripts and subscripts, and punctuation marks are all transcribed as they appear in the document. The advantage of this approach is a strong model that will exactly transcribe what is shown in the image.

However, the neural networks could learn, to a certain extent, to apply our transcription conventions. If the conventions are consistently adopted in all our transcriptions, and the Ground Truth is large enough, the model could learn to separate words that appear combined in the document, normalise historical spelling, transcribe superscripts and subscripts as in line with the rest of the text, solve abbreviations (see the next point).

In particular:

- **Diacritical characters** (e.g. accents, circumflexes, cedillas, hyphens, tildes): it depends on you, if you want the Text Recognition model to make a diplomatic transcription or normalise words according to modern orthography. Both approaches are fine; you just need to choose one and be consistent.
- **i/j and I/J**: the letters “i” and “j” were often used interchangeably. You can decide to transcribe the letters as they appear in the document or to follow the spelling in use today.
- **u/v and U/V**: historical documents often use “v” at the beginning of words and “u” in the middle and end. You can decide to transcribe the letters as they appear in the document or to follow the spelling in use today.
- **Ligatures**: are common combinations of letters to form a new character. They can be transcribed in full, using the single characters composing the ligature (e.g. “præs” becomes “praes”).



- **S-characters:** the letter “s” can appear in different forms. Normal and long “s” (with descender) can both get transcribed as normal “s” or according to their shape as “s” or “ſ” (U-017F). Double “s” or “ß” (sharp “s” or “Eszett”) are transcribed according to the original text.
- **Hyphenated words:** when hyphenated words appear at the end of the line, they should be transcribed and broken up according to the original text. Add a “-” at the end of the line only if present.
- **Text styles:** with the Formatting Bar at the bottom of the Text Editor, you can tag words or portions of words as bold, italic, subscript, superscript, underline, and strikethrough. If you train these tags when training the model, the tags will be automatically added when recognising new pages.
- **Fonts:** different fonts like Kurrent or Antiqua are not specially marked.

Each user can use the conventions best suited to their needs. What is important is to be consistent: we recommend taking note of your decisions while transcribing the pages and writing down the conventions you used in the Details field of the model.

Abbreviations

According to your needs, you can decide to train the model to:

- **Keep the abbreviated form:** transcribe the abbreviations as they appear in the documents, using the base characters or the special characters most similar to the characters written by the writer.
- **Transcribe the expanded form:** the neural networks are often able to learn to recognise and use expansions, especially if they appear frequently. You just need to write the expansions of the abbreviations in the transcriptions, paying attention to always solving them in the same way.
- **Tag the abbreviation and add the corresponding expansion as a property:** in the Ground Truth, transcribe the abbreviations as they appear, tag them and add the expanded form in the “expansion” field (property of the Abbreviation Tag). When training the model, check the option to train the Abbreviation tags with expansions as well.

Retraining with PyLaia

Until November 2022, Transkribus supported another text recognition engine called HTR+. Because of its discontinuation within the platform, HTR+ models can neither be trained nor be applied anymore.

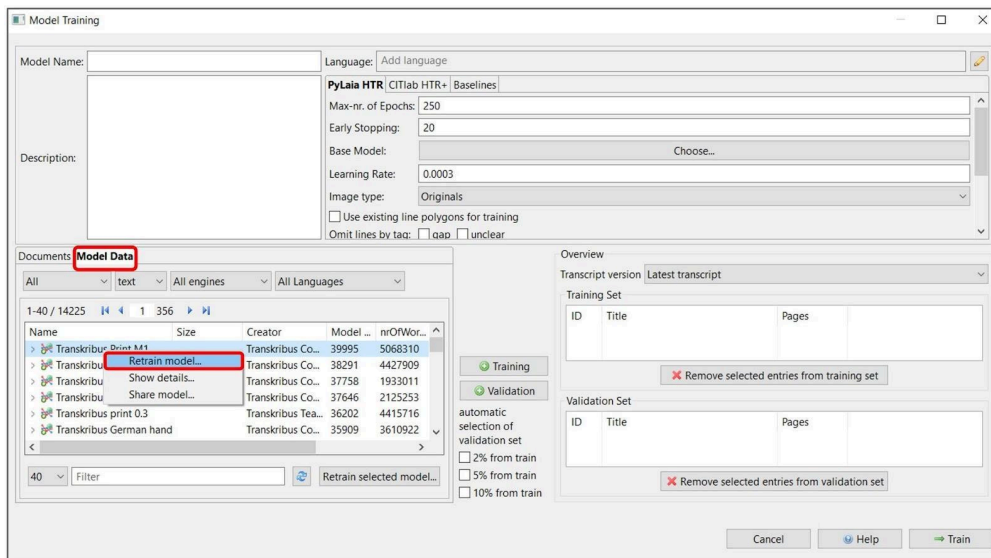


If you have trained an HTR+ model in the past, you can easily retrain it with PyLaia.

Click “Train a new model” in the Tools Tab. Instead of selecting the pages to assign to the Training and Validation Set, click on Model Data: the list of all the models will appear.

Right-click on the model you want to retrain with PyLaia and select “Retrain model...”: all the metadata, the Training Set and Validation Set will be automatically copied into the model setup.

Then, click “Train” to start the training: the PyLaia model will be trained on the same Training and Validation Data used for the HTR+ model. You can also change the advanced parameters and select a base model: just remember that only PyLaia models can be selected as base models.

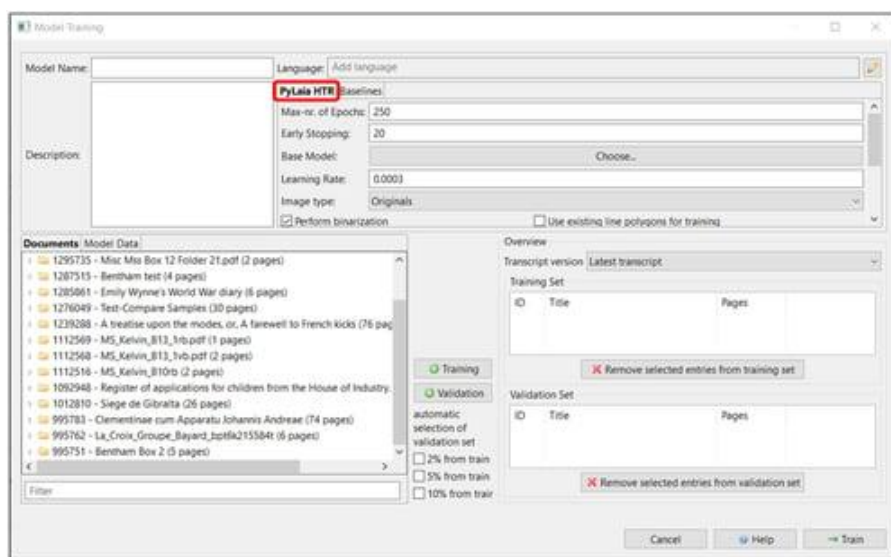




16. Model Setup and Training

Once you have your Ground Truth pages, it is time to train the Text Recognition model.

Click on the “Tools” tab. Under the “Model Training” section, click on “Train a new model”. The Model Training window pops up. By default, “PyLaia HTR”, the engine we are interested in to train the Text Recognition model, is selected, as shown in the figure below.



Model Setup

In the upper section, you will need to add details about your model:

- Model Name (chosen by you);
- Language (of your documents);
- Description (of your model and the documents on which it is trained).

Training Set

During the training, the Ground Truth pages are divided into two groups:

- **Training Set:** set of examples used to fit the parameters of the model, i.e. the data on which the knowledge in the net is based. The model is trained on those pages.
- **Validation Set:** set of examples that provides an unbiased evaluation of a model, used to tune the model's parameters during the training. In other words, the pages of the Validation Set are set aside during the training and are used to assess its accuracy.



In the lower part of the window, first, select the pages that you would like to include in your Training Set.

To add all the pages of your document, select the folder and click “+Training”.

To add a specific sequence of pages from your document to the Training Set, double-click on the folder, click on the first page you wish to include, hold down SHIFT on your keyboard and then click the last page. Then click “+Training”.

To add individual pages from your document to the Training Set, double-click on the folder, hold down CTRL on your keyboard and select the pages you would like to use as training data. Then click “+Training”.

The pages you have selected will appear in the “Training Set” space on the right.

Above it, you can decide which transcription version you want to use for both the Training Set and the Validation Set: latest transcript, Ground Truth only or initial transcript. With the first option, all the latest transcripts, regardless of how they were saved, are displayed and can be selected for training. If you choose “Ground Truth only”, only the pages saved as Ground Truth are selectable.

The screenshot shows a software interface with a 'Documents' pane on the left and an 'Overview' pane on the right. The 'Documents' pane lists several folders, with '995762 - La_Croix_Groupe_Bayard_bpt6k215584t (6 pages)' selected. Below it, pages 1 through 6 are listed with their respective line counts. The 'Overview' pane shows the 'Transcript version' set to 'Latest transcript'. Under 'Training Set', a table lists one entry: '995... La_Croix_Groupe_Bayard_bpt6k21...' with pages '1-2,4-6'. A red box highlights this entry. Below the table is a button 'Remove selected entries from training set'. Under 'Validation Set', an empty table is shown with a button 'Remove selected entries from validation set'. In the center, there are 'Training' and 'Validation' buttons, with 'Training' highlighted in red. Below these are options for 'automatic selection of validation set' with checkboxes for '2% from train', '5% from train', and '10% from train'.

Validation Set

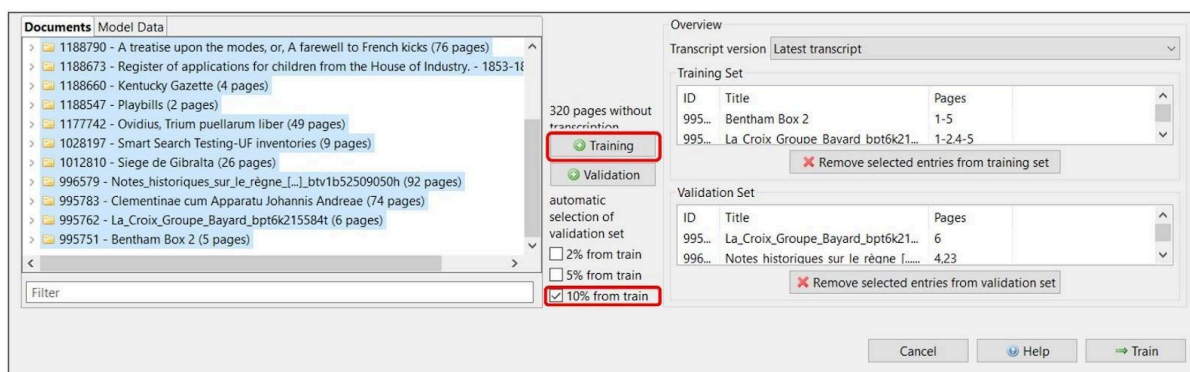
Remember that the Validation Set should be representative of the Ground Truth and comprise all examples (scripts, hands, languages...) included in the Training Set. Otherwise, if the Validation Set is too little varied, the measurement of the model's performance could be biased.

We recommend not to save effort at the Validation Set and assign around 10% of your Ground Truth transcriptions to it.

To add pages to the Validation Set, follow the same process explained above for the Training Set but click the “+Validation” button.

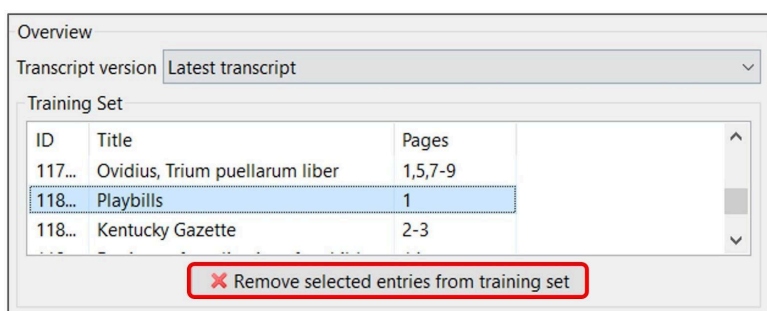


You can also automatically assign 2%, 5% or 10% of the Training Set to the Validation Set. Select the pages to add to the Training Set, flag the percentage you want to assign



to the Validation Set and click the “+Training” button.

To remove pages from the “Training Set” or “Validation Set”, click on the page and then click the red cross button.



Engine

PyLaia is the Text Recognition engine available within Transkribus. When you train a Text Recognition model, it is trained with PyLaia.

In the future, we hope to implement new architectures in Transkribus and offer more text recognition solutions to the users.

Advanced Parameters

The parameters for the PyLaia engine in Transkribus eXpert are divided into two groups:

- “Standard Parameters” (in the upper right section of the window, under “Language”);
- “Advanced Parameters” (accessible by clicking on the button at the bottom of the Parameters section).



Standard Parameters:

PyLala HTR Baselines

Max-nr. of Epochs: 250

Early Stopping: 20

Base Model: Choose...

Learning Rate: 0.0003

Image type: Originals

Perform binarization Use existing line polygons for training

Omit lines by tag: gap unclear

Reverse Text Exclude digits Tag exceptions for reversion: []

Train Abbrevs with expansion Train Tags Include Properties Tag choice: []

Advanced parameters... Reset to defaults

In detail, they are:

- **Max-nr. of Epochs:**

The number of epochs refers to the number of times the learning algorithm will work through the entire Training Data and evaluate itself on both the Training and the Validation Data.

You can increase the maximum number of epochs, but be aware that the training process will take longer. Furthermore, note that the training will be stopped automatically when the model no longer improves (i.e. it has reached the lowest possible CER). To begin with, it makes sense to stick to the default setting of 250.

- **Early Stopping:**

This value defines the minimum number of epochs for the training: the model will at least run this many epochs.

Completed these epochs, if the CER of the Validation Data continues to decrease, the training will continue and end automatically when the model no longer improves. On the other hand, if the CER does not go down anymore, the training will be stopped.

In other words, with the early stopping value, you force the model to train at least that number of epochs.

For most of the models, the default setting of 20 epochs works well, and we recommend leaving it as it is for your first training.

When there is no or little variation in the Validation Set, the model may stop too early. For this reason, we recommend creating a varied Validation Set containing all types of hands and document typologies present in the Training Set.

Only if your Validation Set is rather small, please increase the “Early Stopping”-value in order to avoid the training from stopping before it has



seen all the Training Set. But bear in mind that increasing the value, the training will take longer.

- **Base Model:**

Existing models can be used as starting point to train new models. When you select a Base model, the training does not start from scratch but from the knowledge already learnt by the base model. With the help of a base model, it is possible to reduce the amount of new Ground Truth pages, thus speeding the training process. Likely the base model will also improve the quality of the recognition results.

The benefit of a base model, however, is not always guaranteed and has to be tested for the specific case.

To use a base model, simply choose the desired one with the “Choose...” button next to “Base Model”. You can select one of the PyLaia public models or a PyLaia model of yours.

Remember that, to be beneficial, the base model must have been trained on a writing style similar to the one of your model.

- **Learning Rate:**

The “Learning Rate” defines the increment from one epoch to another, so how fast the training will proceed. With a higher value, the Character Error Rate will go down faster. But, the higher the value, the higher the risk that details are overlooked.

This value is adaptive and will be adjusted automatically. The training is influenced though by the value it is started with. You can go with the default setting here.

- **Image Type:**

We have had some cases where the pre-processing took too much time. If this happens to you, you can switch the “Image Type” to “Compressed”. You can proceed in the following way: start the training with “Original”. When the training has started (“Running” status), every now and then check the progress of the pre-processing with the “Jobs”-button. In case it gets stuck, you can cancel the job and restart it with the “Compressed”-setting.

- **Perform binarization:**

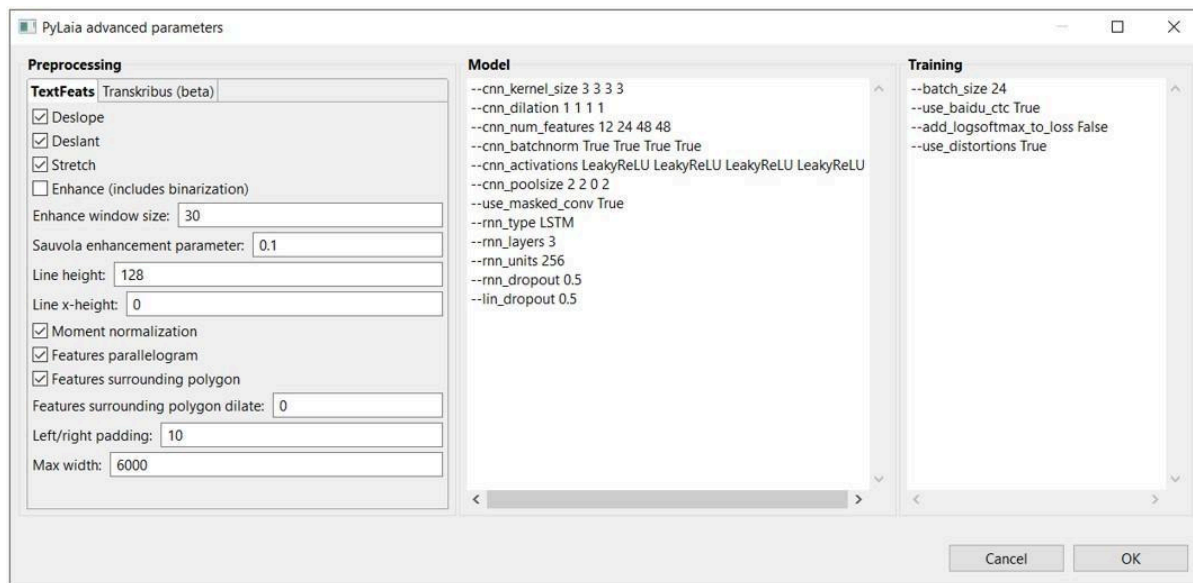
This option is selected by default. Unflag it and use no binarization only if you have homogenous training data, i.e. pages with the same background colour. Only in this case no binarization can lead to better results.



- **Use existing line polygons for training:**
If you flag this option, line polygons will not be computed during the training (as it happens by default), but the existing ones will be used. During the recognition, then, similar line polygons should be used for the best performance of the resulting trained model.
- **Omit lines by tag:**
With this option, you can omit lines containing words tagged as “gap” and/or “unclear” from the training. Please note that the whole line will be ignored during the training, not only the unclear word: this happens because the training happens on the line level.
- **Reverse Text:**
Use this option to reverse text during the training when the writing direction in the image is opposite to that of the transcription, e.g. the text was written right-to-left and transcribed left-to-right. You can also decide to exclude digits or tagged text from reversion.
- **Training Tags:**
It is possible to train textual tags and their properties if they are present in the Ground Truth, so that the model will automatically generate tags during the recognition.
This feature works well with abbreviations and text styles and brings the best results for tags which are repeated in the same way (i.e. the same word) very often.
Select “Train Abbreviations with tags” to train the tagged abbreviations (“Abbrev” tag) and the respective “expansion” properties present in your Ground Truth.
For other tags, select the “Train Tags” option and click “Include Properties”. Use the green plus button to enter the list of tags that should be trained. Be aware that here you would be able to select only the tags that have been added to the “Tag Specification” group.

Advanced Parameters:

Users can set several advanced parameters of PyLaia on their own. You can open the advanced parameters by clicking on the “Advanced parameters”-button at the bottom of the standard parameters.



They are divided into three groups: Preprocessing, Model and Training.

- **Preprocessing Parameters:**

- **Deslant:** choose this option with cursive writing in order to straighten it. Leave out this option with printed documents because if printed documents contain cursive passages in addition to the normal print characters, the effect can be upside down.
- **Deslope:** allows more variation at the baselines, e.g. more tolerance at baselines that are not exactly horizontal but slanting.
- **Stretch:** this option is for narrow writing in order to uncompress it.
- **Enhance:** that is a window which goes over the baselines in order to optimize passages which are difficult to read. This is useful if you have “noise” in the document.
- **Enhance window size:** this setting refers to the option just explained and therefore only needs to be set if you would like to use “Enhance”. This setting defines the size of the window.
- **Sauvola enhancement parameter:** please stick to the default setting here.
- **Line height:** value in pixels; if you need to increase the pixels of the line, you can do this here. 100 is a good value to go for. Attention: if the value is too high it might lead to an “out of memory order”. You can bypass this error in turn by lowering the value of the



“batch size” (top left in the advanced parameters window), e.g. by half. Please be aware that the lower this value, the slower the training will be. The slow-down of the training relating to the batch size should be improved with the new version of PyLaia, which will set the batch size automatically.

- **Line x-height:** this setting applies to the descenders and ascenders. If you put this value, the "Line height" parameter will be ignored.
- Please don't change the following four parameters:
 - Moment normalization
 - Features parallelogram
 - Features surrounding polygon
 - Features surrounding polygon dilate
- **Left/right padding:** 10 (default) means that 10 pixels will be added. This is useful if you are worried that parts of the line could be cut off.
- **Max width:** maximum width that a line can reach; the rest will be cut off. 6000 (default) is already a high value. If you have huge pages, you can further increase this value.
- **Model parameters:**

For all those who are familiar with machine learning and the modification of neural nets. Therefore, these parameters are not further explained here.
- **Training parameters:**
 - **Batch size:** number of pages which are processed at once in the GPU. You can change this value by putting in another number.
 - **Use_distortions True:** the training set is artificially extended in order to increase the variation of the training set and in this way make the model more robust. If you are working on even writing and good scans, you don't need this option. To deactivate it, please write „False“ instead of „True“.
 - The net structure of PyLaia can also be changed – a playground for people who are familiar with machine learning. Modifications on the neural net can be done via the [Github repository](#).

At this point, you can start the training by clicking the “Train” button.



You can follow the progress of the training by clicking the “Jobs” button in the “Server” tab. Depending on the traffic on the servers and the amount of material, your training might take a while. In the “Jobs” window, you can check your position in the queue (i.e. the number of trainings ahead of yours). You can perform other jobs in Transkribus or close the programme during the training process. If the Job status is “created” or “running”, please don’t start a new training, but just be patient and wait.

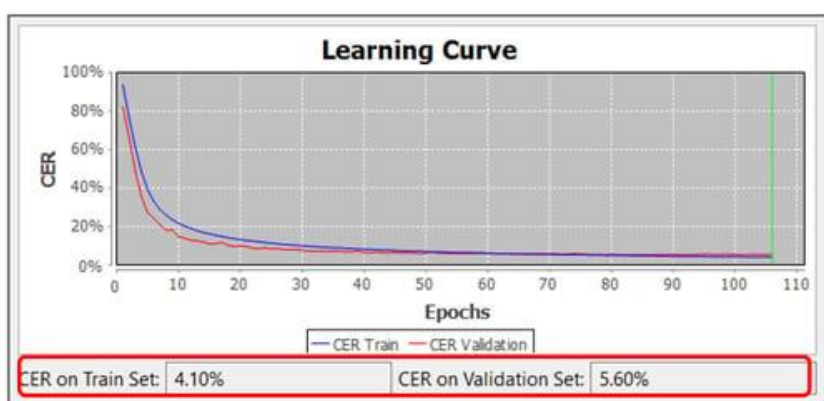
After the training has started, the completion of every epoch will be shown in the Job description, and you will receive an email when the training process is completed.



17. Character Error Rate

The performance of a model is determined based on the “distance” between a perfect transcription (your Ground Truth) and the automatically recognised text. It is measured by the Character Error Rate (CER), i.e. the percentage of characters that have been transcribed incorrectly by the Text Recognition model.

In the Details of the model, you see the CER on both the Training and the Validation Set.



The most representative CER is the one measured on the Validation Set because it shows how the model performs on pages on which it has not been trained.

Results with a CER of 10% or below can be seen as very efficient for automated transcription. However, if applied on hands not seen during the training or scribble notes, the model may perform worse. Results with a CER of 20-30% are sufficient to work with powerful searching tools like Smart Search.

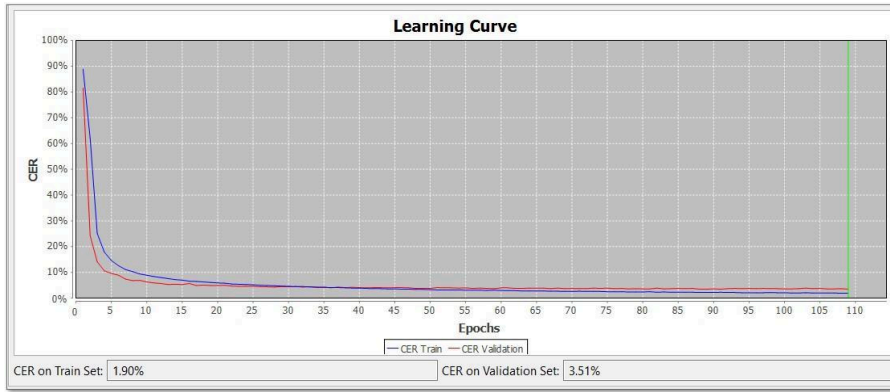
Learning Curve

The “Learning Curve” graph signifies the accuracy of your model.

The y-axis represents the Character Error Rate. The curve goes down as the training progresses and the model improves.

The x-axis represents the Epochs, i.e. the training progress. During the training process, Transkribus makes an evaluation after every epoch.

In the figure below, 109 epochs were trained. In this case, the maximum number of epochs was set to 250, but the training was automatically stopped at 109 because the model no longer improved.



The graph shows two lines, one in blue and one in red. The blue line represents the progress of the training. The red line represents the progress of evaluations on the Validation Set.



18. Computing Accuracy

You can measure the accuracy of your model on specific pages with the “Compute Accuracy” feature in the “Tools” tab.

To do so, you first need two versions of the transcription: the Ground Truth (i.e. the manual transcription as close to the original text as possible) and the transcript generated by the Text Recognition model.

To get the most significant value from the accuracy measurement, it would be best to use pages which have not been used during the training process and therefore are new to the model. Using pages from the Validation Set is also an option but not as ideal as completely new pages. Using pages from the Training Set is not a good idea because this will output CER-values that are lower than they actually are.

To make the comparison, you have to select two transcription versions of the same page. Choose as “Reference” the Ground Truth page version and as “Hypothesis” the version that was automatically generated with the Text Recognition model and on which you would like to test how good the result is.

You can change the versions to be compared by clicking on the grey button next to “Reference” and “Hypothesis”. Double-click to choose the desired version of the document in the appearing window. The versions that can be selected for “Reference” and “Hypothesis” are different versions of your document, which have been created after running a new job or saving transcriptions.

▼ Compute Accuracy...

Reference: (Correct Text) 09.08.22 13:11:56 - s.mansutti@readcoop.eu - Ground Truth Use current

Hypothesis: (HTR Text) 26.10.22 13:17:13 - s.mansutti@readcoop.eu - In Progress - PyLaia decoding 0.7.8 - Model: 37646, Transkribus English Handwriting M3, LM: none Use current

Compare Text Versions...

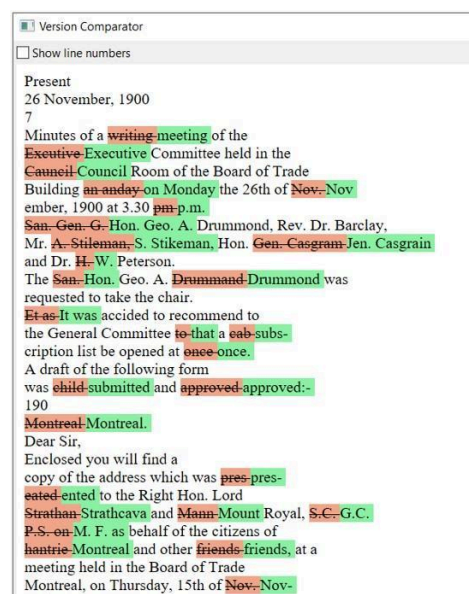
Compare...

Compare Samples...

There are three comparison options available:

Compare Text Versions:

This feature gives you a visual representation of what the HTR model transcribed correctly and incorrectly. After selecting the Reference and Hypothesis versions, click on “Compare Text Versions...” The words marked in red are the



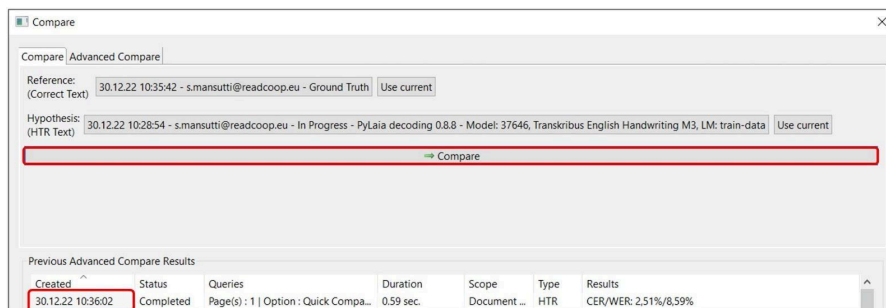


ones transcribed incorrectly by the Text Recognition model; in green, the words are shown as written in the Ground Truth transcription. Please note that even if only one character is wrong, the whole word is marked in red. In the passages without colour, the recognised text is identical to the Ground Truth.

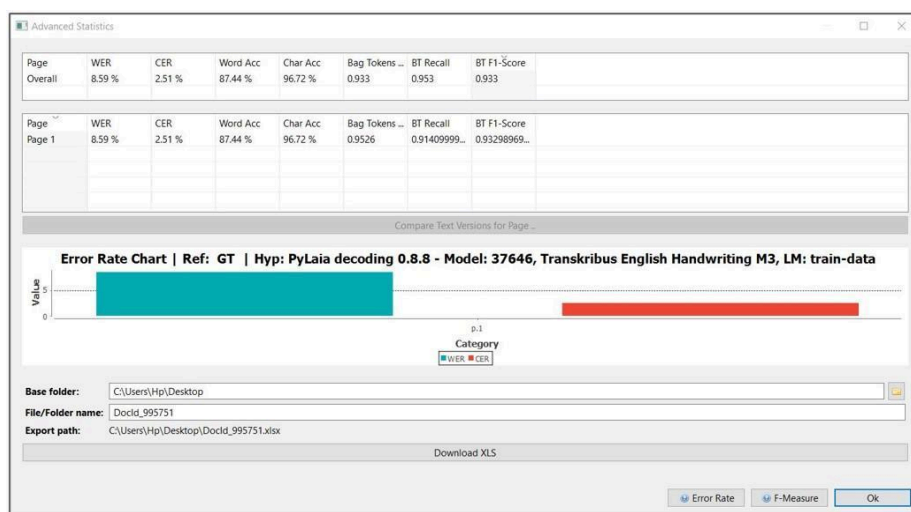
Compare and measure the CER:

This accuracy check enables you to compare the two versions of the same page and measure the Character Error Rate (CER) and Word Error Rate (WER) on that page.

To access it, click on “Compare...”. Firstly, please make sure the right versions have been selected as Reference and Hypothesis in the upper section of the appearing window. Then hit the “Compare” button. The result will be shown in the lower section of the window after a few seconds.



The values are calculated for the page you have currently loaded in the background. In the image above, we have a CER of 2.51% on that page, which means that 97.49% of the characters in the automated transcript are correct. By double-clicking the date and time within the “Created” column of the simple comparison tab, you will automatically arrive at the “Advanced Statistics” window. Here you will get more detailed indications and values, and the results can be exported into an Excel file.





The tables and the chart displays the accuracy of the Text Recognition model on that page in terms of Word Error Rate and Character Error Rate. By double-clicking the line with the page number, you will arrive at the text comparison, where you can check which words or text passages have been challenging.

If you want to check the accuracy for more pages at once, when opening the “Compare” window, you can choose the “Advanced Compare” tab on the right. Add the pages you would like to evaluate (e.g. 1-6), or click on the button with the three points on the far right to choose individual pages.

After starting the accuracy check by clicking on “Compare”, the results will be shown in the table below and by double-clicking the date and time cell in the “Created” column, you will arrive at the “Advanced Statistics” window again.

The overview display shows two tables: one with the “Overall” value, which is the average value of the recognition on all the computed pages in the document. In the table below, you can find the values for the individual pages. This way, you can compare the results on different pages, and by double-clicking the line, you will arrive at the text comparison, where you can check which words or text passages have been challenging. Note: The weighting of pages for the “Overall” value is calculated based on the number of recognised words on a page.

Compare Samples

The “Compare Samples” functionality is useful if you are planning a bigger recognition project and would like to evaluate which model to choose before you run it on the whole document. This comparison feature chooses random lines from the sample document and tests the performance of the model on these lines.

It makes sense to put some pages aside at the beginning in order to use them as a sample document. This is advantageous, as the material the model will be tested on has not been seen before, and therefore the evaluation result will be more reliable.

The “Compare Samples” functionality is situated within the “Tools” tab in the “Compute Accuracy” section. To open it, click on “Compare Samples”, and under



“Create New Samples”, fill out the required information.

The screenshot shows the 'Compare Samples' application window. The 'Create New Sample' tab is selected. The 'Sample Title' field is filled with 'Sample_Siege de Gibraltar'. The 'Nr. of lines for sample' field is filled with '100'. The 'Documents' list on the left contains several entries, with '1012810 - Siege de Gibraltar (26 pages)' selected. The 'Documents added to Sample Set' table on the right shows one entry: ID 101..., Title 'Siege de Gibraltar', and Pages 1-26. The 'Add to Sample Set' button is visible below the document list. The 'Remove selected entries from sample set' button is visible below the table. The 'Create Sample' button is visible at the bottom right of the window. The 'Help' and 'Cancel' buttons are also visible at the bottom right.

At “Nr. of lines for sample”, you can define how many lines you would like to test. 500 is the recommended average. The more lines you put here, the lower the variation will be in the result, and the prognosis will be more precise. For a large project with many pages, it might be reasonable to say 1000 lines; for a very small attempt, maybe 100 lines are already enough. Here too, the best way to go about it is a “trial and error” approach, as it always depends on the individual goal.

With the “Baseline length threshold”, you can control the length of lines, which is practical if you have a lot of short lines in your material, which often happens e.g. with tables and newspapers. The value to enter should be between 0 and 1 (e.g. 0.5 for half the page width). If set, baselines with a length smaller than this fraction of the page width will be discarded. For handwritten material with only one column, this step is probably not necessary.

By clicking the “Keep line text” option, you can literally keep the text you already have in your documents and only need to correct the lines after creating the sample.

From the list on the left side, choose the documents of which the sample should consist via the “Add to Sample Set button”. Then click on “Create sample”.

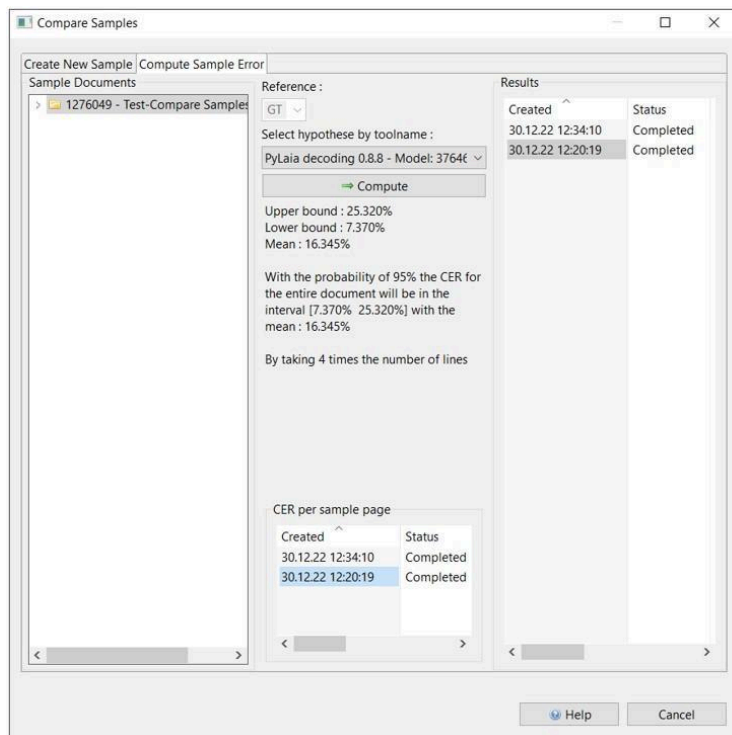


Transkribus will now randomly choose the defined number of lines in the selected documents.

The next step is to open the sample document (you can find it in your collection) and manually transcribe the line snippets (if you haven't kept the text like described above). It will be only one line per page, and therefore the transcribing in most cases will be quick. If you have finished one line, jump to the next page of the sample document to proceed.

When you are finished with transcribing, run the model which you would like to test on the sample document to produce the transcription.

You can then compare the two versions. To do so, open the "Compute Sample Error" tab in the "Compare Samples" window and choose the sample you would like to evaluate. Then click on "Compute" to start the job. As soon as "Completed" appears in the "Status" column, you can double-click the cell in the "Created" column to view the results.



By double-clicking the date and time in the "Created" column in the "CER per sample page" table, you will automatically arrive at the "Advanced Statistics" window and look at the accuracy measurements more in detail.



19. Structural Tags

With structure tags, you can divide up your documents into structural sections like paragraphs, headers or page numbers and also add customized tag categories for your individual needs. Moreover, it is possible to train P2PaLa models to automatically recognise your documents' structure.

It is not possible to save structure tags at the Collection level in Transkribus eXpert.

There is no need to tag every feature of your documents: focus on marking up the sections that you are interested in.

Structure type	Color	Shortcut
paragraph	Light Blue	
heading	Pink	
caption	Green	
header	Blue	
footer	Purple	
page-number	Yellow	
drop-capital	Teal	
credit	Light Green	
floating	Brown	
signature-mark	Dark Blue	
catch-word	Light Green	
marginalia	Brown	

Type	Structure	Text	ID
> TextRegion	header		TextRe...
TextRegion			TextRe...
> TextRegion	title		r
TextRegion			r_580
> TextRegion	sub-title		r_584
TextRegion			r_599
> TextRegion	page-number		r4
TextRegion			r_586
TextRegion	drop-capital		

First, open your document in Transkribus eXpert. The structural tagging interface can be found by clicking the “Metadata” tab and then the “Structural” tab. In the centre of the tab, you can see the different predefined structure types.

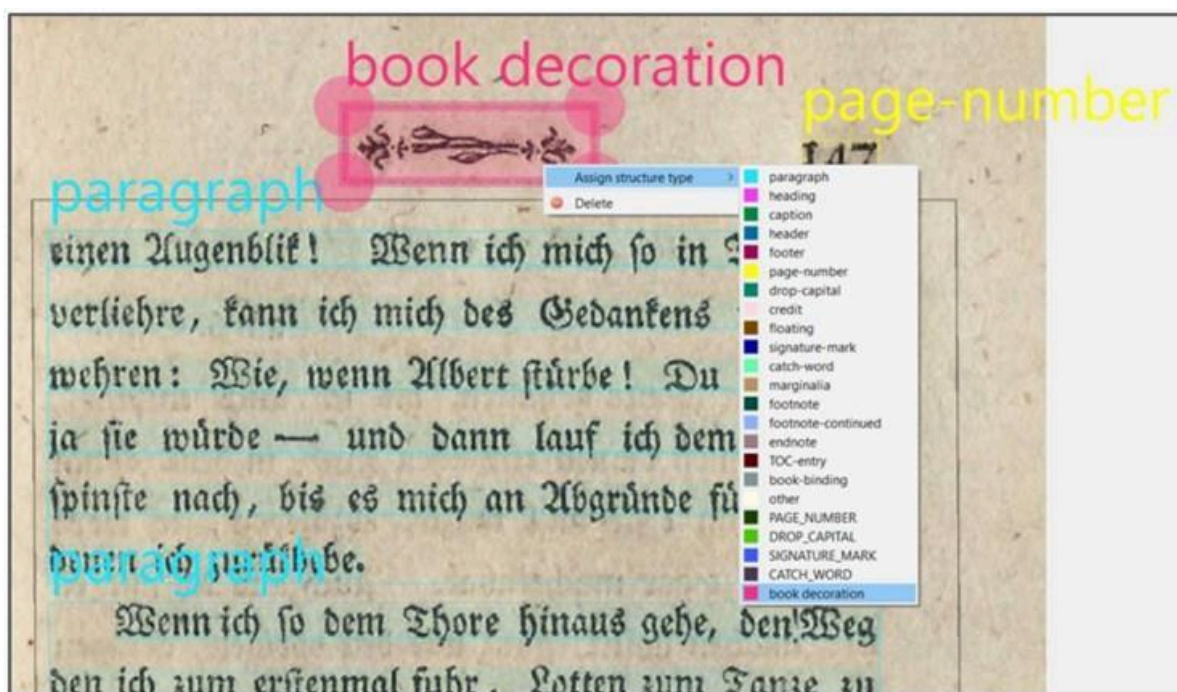
To create your own tag categories, click the “Customize” button. The “Tag configuration” window will open up. In order to create a new tag category, simply type in the name in the blank box at the bottom of the window, then click the green plus button. In this window, you can also customize the tag colours by clicking on the coloured section next to a tag and then choosing your desired colour. The new tags you



created will also be automatically available for all your documents in all your collections.

You can assign tags to text regions and line regions on each page in your document. To place a tag first, click on the “Item visibility” button in the Main menu and make sure that text regions and line regions are visible on your document. Select the text region or line in the image window, right-click the selected shape and then choose the desired tag under “Assign structure type”. Or alternatively, you can add the tag by clicking the green plus button on the right of the desired tag category in the “Structural” tab.

You can select and tag several regions at once by holding down the “CTRL” key on your keyboard and then clicking on your document.

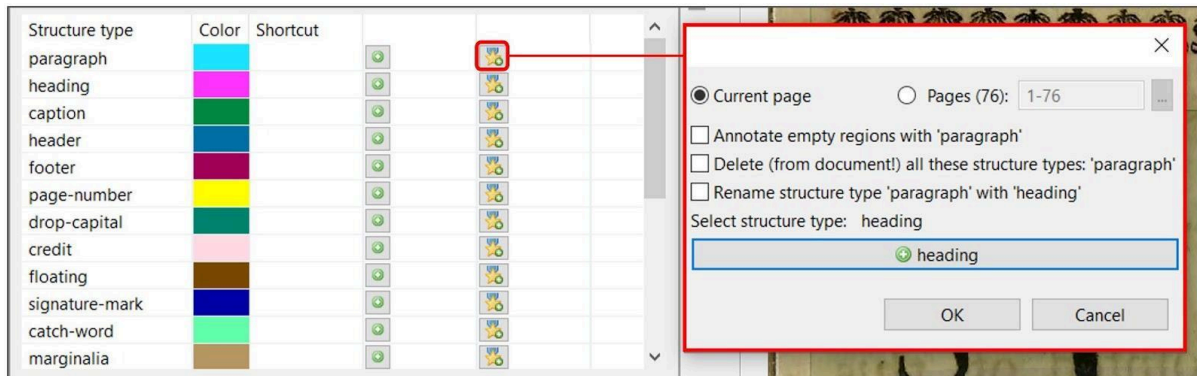


The structural tab also enables you to:

- Assign a “Page type” to each page of your document. Possible options are: Front cover, Back cover, Title, Table-of-contents, Index, Content, Blank, Other. When you have your page open, choose the appropriate definition by clicking the arrow next to the “Page type” options and then choosing the desired type. The page type is not relevant for the P2PaLA training.
- Link two structural tags together with the “Links” buttons, e.g. a link between a line and the footnote connected with that line. The first button is to create such a link, and the second one is to remove it. Please note that for P2PaLA training, the linking of shapes is not relevant.



- Remove a structural tag: select the tagged region and then click the red button;
- Show structural tags names and colours in the image window;
- Click the star button next to each structural tag to access advanced options: here, you can annotate all empty text regions with the structural tag of your choice; delete from all the pages of the document a certain structural tag; rename an assigned structural tag with another name.



- Layout section: here, you find an overview of the structural types in your document and snippets of any transcribed text. You may find it quicker to consult this list rather than search for a particular line or text region in the image. To go to the desired tagged text or line region, double-click the region in the "Layout" section. The image and the Text Editor will automatically jump to this line. The tags you have added will be shown in the "Structure" column. Next to the structure type, there is a small downward arrow. By clicking it, you can quickly change the structure tag; if you click the "delete" (it is the first item of the list), the structural tag will be deleted.



20. Textual Tags

Textual tags enable you to enrich your transcriptions by tagging some words (e.g. abbreviations, places, people...) and adding attributes. You can then search for specific tags and export them in different file formats so that you can go on working with them outside of Transkribus.

To manage the textual tags of a collection, open the collection in question and go to the Metadata-Textual tab in the Managing & Tools sections. Flag the "Collection tags" option to see all the tags available in the collection.

If you want to edit the tags linked to the collection, click "Customize" and select "Collections tags". Here you can add new tags and their attributes to the collection: click on "Create new tag", type its name and add the properties (if necessary) with the "Add properties" green + button.

To delete a collection tag, select it and click the "Remove" red button. Note, however, that the predefined tags and their properties cannot be deleted: they will always be shown in your collection tags list, even if you are not interested in using them.

With the most right green + button, you can take over the selected tags to the User Tags so that you can use them in every collection. You can do the same the other way around: in the User Tags tab, you can select some tags and take them over to the Collections Tags: in this way, all the users who have access to the collection can use the tags.

All the changes done to the Collection Tags are saved only for the collection in question (the one opened in the background). The other collection users can see and use your newly added or edited tags.

The tagging interface can be found by clicking the "Metadata" tab and then the "Textual" tab.



The screenshot shows the 'Metadata' tab of a text editor. At the top, there are tabs for 'Server', 'Overview', 'Layout', 'Metadata', and 'Tools'. Below these are sub-tabs for 'Document', 'Page', 'Structural', 'Textual', and 'Comments'. The main area is titled 'Tags of current Transcript' and contains a table with columns for 'Tag', 'Value', 'Text', and 'Properties'. The table lists 13 entries, each with a tag name (e.g., 'date', 'abbrev', 'person', 'place', 'textStyle', 'unclear', 'gap', 'blackening', 'place'), a value (e.g., 'Augt. 1.st 1914', 'Augt.', 'V.', 'Ballyarthur'), a snippet of text (e.g., 'Augt. 1.st 1914 War declared', 'Augt. 1.st 1914 War'), and properties (e.g., 'year: 1914 day: 1'). Below the table is a 'Tags' section with a list of tag specifications, including 'gap', 'organization', 'person', 'place', 'sic', 'speech', and 'modified'. Each tag has a color swatch and a 'Shortcut' column. At the bottom, there is a 'Props for tag: 'date' - value: 'Augt. 1.st 1914'' section with fields for 'day', 'month', and 'year'. Navigation buttons for 'Previous', 'Next', and 'Apply to selected' are at the bottom.

Tag	Value	Text	Properties
1 date	Augt. 1.st 1914	Augt. 1.st 1914 War declared	year: 1914 day: 1
2 abbrev	Augt.	Augt. 1.st 1914 War	
3 person	V.	paper Afternoon V. & I went	
4 place	Ballyarthur	& I went up to Ballyarthur, ti	
5 person	Mr. Armetrong	the Olivers & Mr. Armetrong	
6 person	Mr. A.	there for tennis. Mr. A. said t	
7 textStyle	channel	said the cross channel stearr	italic: true fontSi
8 unclear	England	could get over to England Iv	
9 gap		she had wished to leave the	
10 textStyle	gun-running	day after the gun-running &	fontSize: 0.0 kern
11 person	Jack	In the evening Jack turned up	
12 blackening	Barrons	turned up said the Barrons h	
13 place	Slandelough	their coming to Slandelough	

To add a textual tag, select the text in the Text Editor and afterwards click on the green + button near the tag you want to apply. Alternatively, after having highlighted the text, right-click with your mouse and choose the suitable tag under "All tags".

In the upper section of the Textual tab, you see the tags present in the transcription of the current page. Clicking on one of them, the image and the Text Editor will automatically jump to the line containing it.

Use the red button here to delete tags: select the tag from the list (press CTRL to select more than one tag at a time) and click the red button. Alternatively, you can delete a tag by highlighting the tagged word or phrase, right-clicking with your mouse and then pressing the "Delete" button. The program will give you two options: "Delete only the highlighted tag" or "Delete all the tags for the current selection".

Below, there is the "Tag" section: here are listed all the textual tags you can use. You can decide if you want to show all the user tags or only the ones linked to the collection. Near each tag, there is a green + button and a star button: the first, as explained before, adds a tag to the highlighted text; the star button gives access to

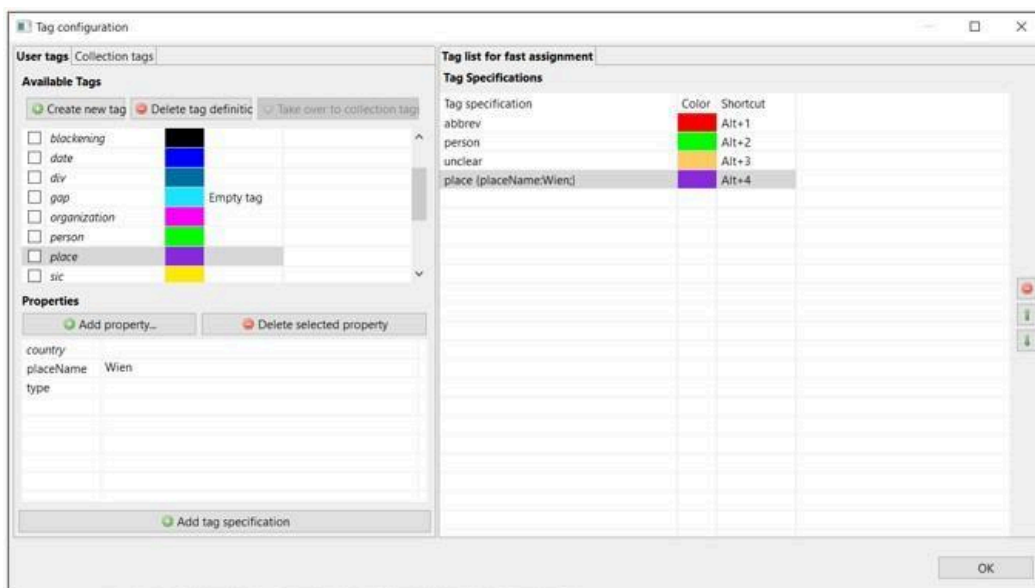


advanced options for tagging. In particular, the advanced options enable you to delete or rename all tags of a certain type.

By clicking on the “Customize button”, the Tag configuration window opens up. Here you can customise the textual tags both at the level of user and collection. In more detail, you can:

- Create new tags;
- Delete tags;
- Change the colour with which the tag is displayed;
- Add/modify properties of new and existing tags;
- Add a tag to the list of “Tag Specification” and assign it a shortcut. Shortcuts allow you to speed up the insertion of frequently used tags: select the text and press the shortcut keys to add the tag.
- Add a shortcut relating to the properties of your tags, e.g. for expanding abbreviations or adding a standardised country name to a place tag. Select the tag and type the property content you would like. Then click “Add tag specification”: now your tag and its property will appear in the “Tag Specification” section of the window, and you can add the shortcut you would like to use. Now you can add the tag and its property by simply highlighting the word or phrase in the Text Editor field and then pressing the shortcut.

In the “Tag Configuration” window, predefined tags are shown in italics; customized ones are shown without italicisation.





Lastly, at the bottom of the Textual tab, there is the Properties section, where you can edit the properties of a tag, if it has any.

Textual tags should be tailored to your specific purposes, but in general, we can say that the following ones are particularly useful when working with historical documents:

- **Abbreviation:**

You can tag the abbreviated word and type the expansion as a property.

There is no right way to deal with abbreviations. Depending on the transcription you want to have in the end, you can adopt one of these approaches:

1. transcribe the expanded version of the abbreviation directly into the text editor (the neural networks are often able to learn to recognise and use expansions, especially if they appear frequently);
2. transcribe the abbreviation using the base characters most similar to the characters written by the writer; you can then tag it and write the expansion as a property;
3. transcribe the abbreviation using the Unicode characters, which are near to the special graphemes of the original document. Remember that the Virtual Keyboard enables you to add special characters. Since it is often hard to decide which Unicode character may be the right one, you may consult the [MUI website](#) to get more information on this issue. As above, if interested, you can tag the abbreviation and write the expansion as a property.

- **Unclear:**

Use this tag when the text can not be transcribed since it is illegible. Highlight the unclear text in the text editor and tag it as “unclear”; you may also add alternatives or suggestions for the illegible word as an attribute of the tag.

- **Gap:**

If the text is impossible to read, add the “gap” tag where the illegible text should appear in the text editor.

- **Blackening:**

Use this tag to redact sensitive information in the export formats. Typically this is used to hide personal data in a document which is made publicly available. The blackening tag is used in conjunction with the “blackening” region, which must be added with the segmentation tools.

To blacken part of your text in the image window, go to the Canvas menu, use the drop-down menu on the “+...” segmentation element button and select



“Blackening”. Use the “Blackening” region to mark the word or section that you want to hide (remember to click the “Item visibility” button on the Main menu and select “Render blackenings” to display the blackened sections on a page). Then, highlight the corresponding word in the Text Editor and select the “Blackening” tag.

In the export of the document, the text will be replaced by asterisks. When you export your document, make sure that “Do blackening” is selected. Please note that in METS and TEI files, the word or phrase is blacked out, but the information behind the blackened section is kept. In other file formats, the text behind the blacked-out section is completely obscured.

Text styles (bold, italic, subscript, superscript, strikethrough, underline...) are added as textual tags. To modify the text style of your text, use the Formatting Bar below the Text Editor.

It is possible to train tags and properties while training an HTR model. The resulting model will both transcribe the document and add the textual tags to the transcription automatically. Read the Model Training page to learn how the training tags feature works.



21. Searching

Fulltext Search

The Fulltext Search is the simplest search within Transkribus: it searches the latest transcription of your documents and returns the pages that contain some or all of the search terms.

To open the Search window, click on the binoculars icon in the Main Tab at the top of the screen. Choose if you want to search any word, all words or the exact phrase and type the term(s) you want to search.

Before launching the search, you can choose to see a word image preview when hovering over the result; restrict the search to the current document (otherwise, the search will be done on all the collections you have access to); perform a Case sensitive search or fuzzy search.

The search results will appear below. For each, you see the result highlighted in yellow within its context, the document name and ID, and the page number where it has been found. By hovering over the highlighted word, the word image preview appears. Double-click on the result to open the corresponding page in the background and be directed to the line containing the result.

To find variations on your search terms, use a question mark (?) for a single-character wildcard search and an asterisk (*) for a multiple-character wildcard search. For example, if you search the term "bo?k", you will get "book", "boek" and "bock" as results. If you search "stat*", you will find "state", "statue", "station", "statutes"...

The lower section of the windows enables you to mark up words with textual tags or replace them easily. Select the results to tag/replace in the search result list (use the green tick to select all of them).

To tag them, select from the drop-down menu the tag and the properties that you have already added to the "Tag Specification" group.

To replace them, type the new word and then click on "Replace selected" to replace only the one selected in the search result list; "Replace all in doc" and "Replace all in coll" to replace the search term all the type that appear, respectively, in the document and in the collection. Reload the page/document to see the replacement in effect.



Fuzzy Search

Fuzzy Search is a search technique that allows you to find similar words in addition to exact matches. It retrieves results that differ by only one or two letters from the search term, and it is useful with misspellings and spelling variations.

To enable it, check the Fuzzy search option and launch the Search. The results will contain both the exact matches and the variants (for example, when searching for "house", fuzzy search finds also "howse, "horse" and "houses").



22. Workspace Management

Managing Collections

Within Transkribus eXpert, your Documents are organized in so-called Collections. These collections can be understood as a folder containing documents.

Collections are typically used on a project basis. For instance, all documents belonging to one project are organized in one collection. One collection can have multiple documents. And documents, on the other hand, can consist of one or more pages.

To see all your collections, go to the Server tab in the Managing & Tool section on the left side of your screen, and click on the collection name next to "Collections". A new window with all the collections that you have access to will open. Here you can:

- Create a new collection with the green + button;
- Delete a selected collection with the red button (only if your role is Owner);
- Modify the collection's metadata (name and description);
- Manage Users;
- Stray documents: after deleting a collection, the documents end up here. You can reassign them to another collection or delete them permanently;
- Credit Manager.

Three elements are managed at the collection level:

- Users
- Textual tags
- Structure tags

Managing Documents

When you open a collection, the list of all its documents appears in the Server tab, in the Managing & Tools section. Double-click on a document to open it.

The toolbar at the top of the documents' list is extremely useful for managing your documents:



Documents		Model Data			
ID	Title	Pages	Uploader	Uploaded	Co
128...	Emily Wynne's World War diary	6	s.mansutti...	Mon Jan 09...	(Tr
127...	Test-Compare Samples	30	s.mansutti...	Fri Dec 30 ...	(Tr
123...	A treatise upon the modes, ...	76	s.mansutti...	Tue Nov 2...	(Ti
111...	MS_Kelvin_B13_1rb.pdf	1	s.mansutti...	Mon Aug 0...	(Tr

- Green + button: to link the document to another collection: it only creates a link to the document in the new collection; it does not duplicate it. The document in the server (and the doc ID) remains one; all saved changes are visible in all collections to which the document is linked, regardless of the collection in which they were made.
In the "Collections" column of the document list, you can see all the collections to which the document is linked.
Be careful when deleting a linked document because it will be deleted in all collections with which it is shared.
- Red - button: to unlinked the document from the opened collection. You will receive an error message if you try to unlink a document that is only present in one collection: link the document to another collection and then unlink it from this (if you want that the images and transcriptions remained within Transkribus for your future use). On the contrary, if you want to eliminate the document from the server definitely, delete it.
- Delete button: use it to delete a document (images, transcriptions, metadata...) from Transkribus. Deleted documents will remain in the Recycle bin for two weeks: after that, they will be discarded, and it will be impossible to restore them.
- Manage users: see the [Managing Users page](#).
- Duplicate button: to duplicate the selected document into another collection. It creates a copy of the document in the server, assigning it a new doc ID. The original document and its copy do not maintain any link, and the changes made in one of them do not appear in the other (differently from what happened when you link a document to another collection). You can decide to copy only the pages with a certain status; to copy all the available transcriptions, i.e. all the saved version, not only the latest. If you select more than one document, you can choose to copy all the selected into one new document.
- Document Manager: it enables to:



- Batch update the document status of all the pages of a document or some selected ones;
- Add new page(s) to a document: some selected pages or all the pages of the document. Select the document and click the green + button, then select the image(s) you want to add. It will be added as the last page, and then you can move it to the right position.
- Add local transcripts (PAGE XML): if you already have the PAGE XML of your documents with the layout coordinates and transcripts, you can import them from here. You need first to prepare a folder with the PAGE XML file for each page, preferably with the same filename of the images. From the Document Manager, select the folder, synchronise the transcriptions and the images (if they have the same filename, no issue should arise here) and click ok. You will now see the images with the text regions, line regions, baselines and text (if present) according to the PAGE XML file.
- Delete one or more pages from a document. This action cannot be undone.
- Set the transcript with a selected status as the latest version.
- Create a Sample: it is useful when you want to have a representative sample of all your documents to train a Text Recognition or Layout model. Select the documents to include in the sample, click "Add to Basic Sample", and then "Create Sample". You can choose to select the sample randomly (x pages from the basic sample), systematically (the x-th page from each document), or to select x pages from each document. A new document will be created within the same collection.
- Recycle bin: it contains the deleted documents in the last two weeks. If you delete a document by mistake, here you can restore it. Note that after two weeks from the deletion, the document will be discharged definitely, and it will not be possible to restore it anymore.

Finally, to edit the metadata of a document, go to the Metadata-Document tab in the Managing & Tools section. Here you can edit the title, author, genre, script type, date, and description; plus, you can add some extended document metadata (authority, hierarchy, backlink and external ID).

At the bottom of the section, there is also the possibility to define and add your own metadata fields, as well as edit the Editorial Declaration.



Since there are always several ways to produce a correct transcript of a text, it is important to be transparent about the way in which the transcription was undertaken, especially if you aim to create a scholarly edition. The “Editorial Declaration” is the space where you can state the editorial choices you applied during the (manual or automatic) transcription process. It offers a set of predefined features and options. Moreover, you are able to create your own descriptions and store them together with your document. The Editorial Declaration can be then exported as an XLSX file.

Managing Models

To view the Text Recognition and Baselines models and manage them, go to the Tools Tab in the Managing & Tools Bar and click the “View models...” button: a window with the public models and all your private models opens.

Here you can browse the public and private models and edit the metadata of your private ones. With the second drop-down menu, you can choose to show the list of Text Recognition models (text) or Baselines models (layout). When you click on a model, its details will appear on the right-side of the window.

In regard to your private models (i.e. the models you trained), you can:

- Share a model: right-click on the model’s name in the list and select “Share model...”: a window will open up, and by clicking on the green plus button, you can add the collection(s) to which you would like to share the model. See the section below for more details on sharing a model.
- Delete a model: right-click on the model’s name in the list and select “Delete model...” If you confirm this action, you cannot undo it.
- Edit model’s metadata: modify and save your changes in the details area of the model, on the right-side of the window.

All the models you train are by default private. If you want to make a model public because it may benefit other Trankribus users, send an email at info@readcoop.eu with a short description of the model and a representative image of the script. You can also decide if, alongside the model, you want to make public the dataset on which the model was trained. If not specified, the model will be made public but the dataset will remain private, i.e. all the Trankribus users can apply the model to their documents without having access to the Ground Truth.

Sharing a Model

Sharing a private model is possible only with collections you have access to (because you have created them or they have been shared with you); it is not possible to share it



directly with individual users. This is done to prevent users from using your private models without your control.

To share a private model with a collection, go to the Tools Tab in the Managing & Tools Bar and click the “View models...” button: a window with all the models will open. In the models’ list, right-click on the private model’s name you want to share and select “Share model...”. In the new window, click on the green plus button to add the collection(s) you would like to share the model with. If you are the creator of the collection, remember then to add the interested users to the collection with the role of Owner, Editor or Transcriber (as explained on the [Managing Users page](#)).

If you want to unshare a model with a collection, follow the steps above but instead of clicking on the green plus button, select the collection name from the list and click the “Remove from current collection” red button.

For more information on all Transkribus products and features please have a look here: <https://help.transkribus.org/>