

Mary Shelley HTR: a Handwritten Text Recognition campaign with Transkribus

Project for the Semantic Digital Libraries course, a.a. 2023–2024

Gaia Ortona, Chiara Parravicini, Alice Picco

Introduction

The primary objective of our project is to develop a model capable of automatically transcribing the handwritten texts of Mary Shelley with a reasonably low error rate. We have utilized the online platform Transkribus¹ for training two distinct models, using approximately 70 pages (around 10000 words) from Shelley's renowned novels, *Frankenstein* and *Mathilda*.

The context

The Shelley-Godwin Archive² is a digital repository that provides access to the digitized manuscripts of prominent literary figures such as Percy Bysshe Shelley, Mary Wollstonecraft Shelley, William Godwin, and Mary Wollstonecraft. It serves as a platform to collect, for the first time online, the dispersed handwritten legacy of this family of writers. Through partnerships between institutions like the New York Public Library³, the Maryland Institute for Technology in the Humanities⁴, Oxford's Bodleian Library⁵, the Huntington Library⁶, the British Library⁷, the Houghton Library⁸, and the Victoria and Albert Museum⁹, the archive contains over 90% of all known relevant manuscripts.

Mary Wollstonecraft, William Godwin, Mary Wollstonecraft Godwin Shelley, and Percy Bysshe Shelley collectively form "England's first family of writers." Their literary and philosophical contributions are remarkable, as they include crucial works for English literature and political philosophy. Of particular note is Mary Shelley's *Frankenstein, or the Modern Prometheus*, published in 1818, which profoundly influenced Western literature over the years.

¹ <https://www.transkribus.org/>

² <http://shelleygodwinarchive.org/>

³ <https://www.nypl.org/>

⁴ <https://mith.umd.edu/>

⁵ <https://www.bodleian.ox.ac.uk/home>

⁶ <https://huntington.org/>

⁷ <https://www.bl.uk/>

⁸ <https://library.harvard.edu/libraries/houghton>

⁹ <https://www.vam.ac.uk/>

The technological infrastructure of the Shelley-Godwin Archive is built on linked data principles and adheres to emerging standards such as the Shared Canvas data model¹⁰ and the Text Encoding Initiative (TEI)¹¹. These frameworks support a participatory platform where scholars, students, and the public can engage in the curation and annotation of the archive's contents. Additionally, the archive's transcriptions and associated software applications are openly available on GitHub¹² under open licenses ([Creative Commons Attribution-NonCommercial-ShareAlike 2.0 Generic \(CC BY-NC-SA 2.0\)](#)), fostering collaboration and innovation within the academic community.

Our project fits well in the context of the Shelley-Godwin Archive because it aims to speed up the process of adding new material to the archive. While we're initially concentrating on transcribing Mary Shelley's works, our models could easily be adjusted to transcribe texts from other authors featured in the archive. This flexibility could improve significantly to the archive's expansion and accessibility, offering a wider array of literary materials for academic study. By automating transcription, our goal is to simplify archival tasks, encouraging more scholars and enthusiasts to delve into literary history.

Workflow

As mentioned in the section above, in order to develop our models we had to select the works in the archive that were written exclusively by Mary Shelley. Therefore our overall documents set ended up consisting of 36 pages of *Frankenstein* (more specifically the chapters 5,6 and 7) and 28 pages of *Mathilda*, Shelley's second novel. Out of these 64 pages, a percentage of 10% of them was randomly selected by Transkribus as part of the validation set before the training to assess the model performance.

The number of pages was set based on what is specified in the Transkribus guidelines (it's recommended to have at least 10k words to train a HTR model¹³).

After the selection process, we leveraged the pre-trained public model available on Transkribus (the Universal Lines model) for the setting of the layout. Regarding the textual part, we relied on the transcription provided by the archive, considering it a reliable ground truth.

¹⁰ <https://iiif.io/api/model/shared-canvas/1.0/>

¹¹ <https://tei-c.org/>

¹² <https://github.com/umd-mith/sqa>

¹³ <https://help.transkribus.org/data-preparation>

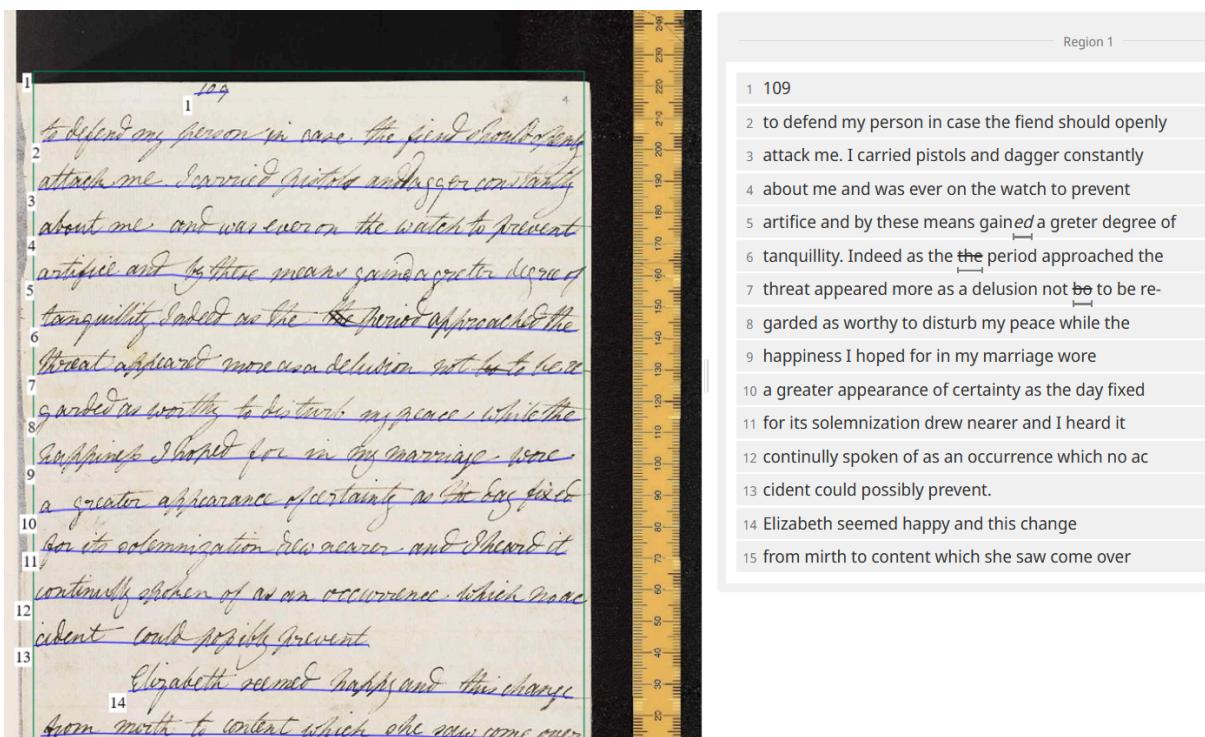


Figure 1 - "Frankenstein" Mary Shelley, page 109

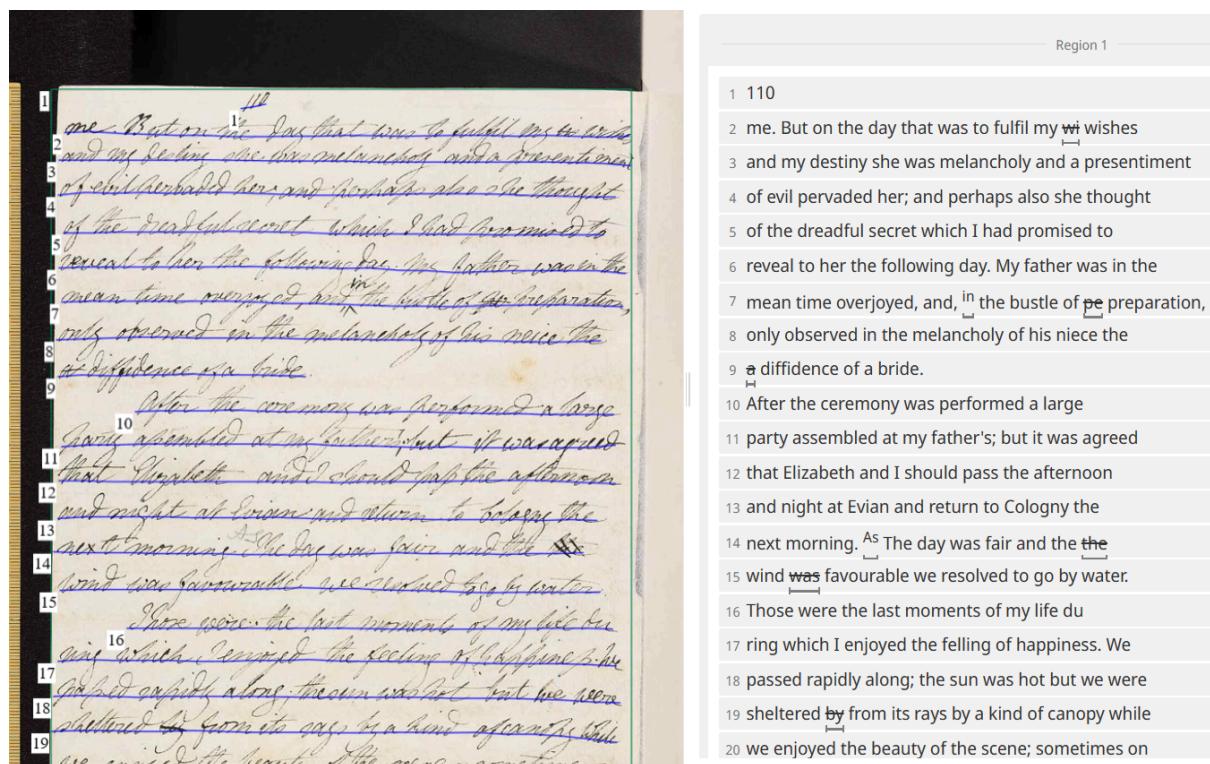


Figure 2 - "Frankenstein" Mary Shelley, page 110

As it shows in *Figure 1* and *Figure 2*, we opted for a one-region-layout and we used two tags to underline the presence in the text of deletions and additions.

Training the models with Transkribus

In the subsequent phase following transcription, our project shifted focus to model training. We devised two private models: one fine-tuned using an existing model (specifically, the English model M4), and the other without any pre-existing tuning. The latter approach allowed us to systematically assess potential disparities and inconsistencies.

The main parameters we focused on were the batch size and the number of epochs: we opted for a small batch size number (10) in order to avoid overfitting and improving the computational efficiency; regarding the number of epochs, we opted for the default setting (100) because we thought it was an appropriate number for data processing and evaluation since the dataset is not large.

Results

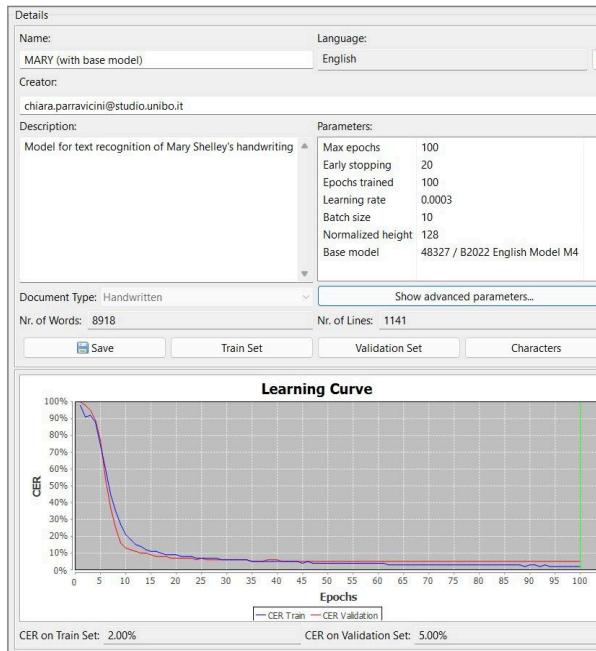


Figure 3 - HTR MaryModel (with base model)

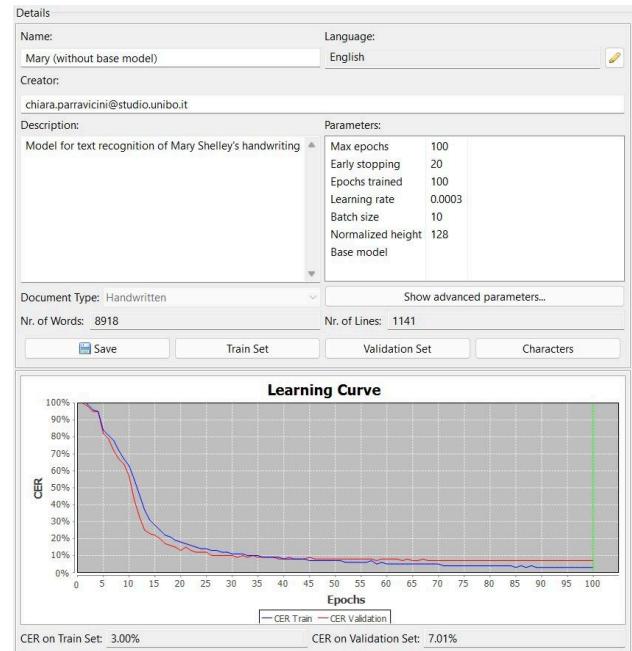


Figure 4 - HTR MaryModel (no base model)

As it shows in *Figure 3* and *Figure 4*, the graphs showcase the learning curves of both the training and validation sets, which represent how much the Character Error Rate is diminishing over the epochs. Generally speaking if CER is over 10% the transcriptions are not so accurate because a lot of manual corrections are needed. In our case, both models performed quite well, even though it's pretty clear that the one without fine-tuning has a higher CER percentage both on the training set and the validation set.

Testing the model on an unseen page of *Mathilda*

The last step of the project was to check the functionality of the model by testing it on an unseen page written by Mary Shelley. We chose to do it on a page taken from *Mathilda* that wasn't included in the initial dataset.

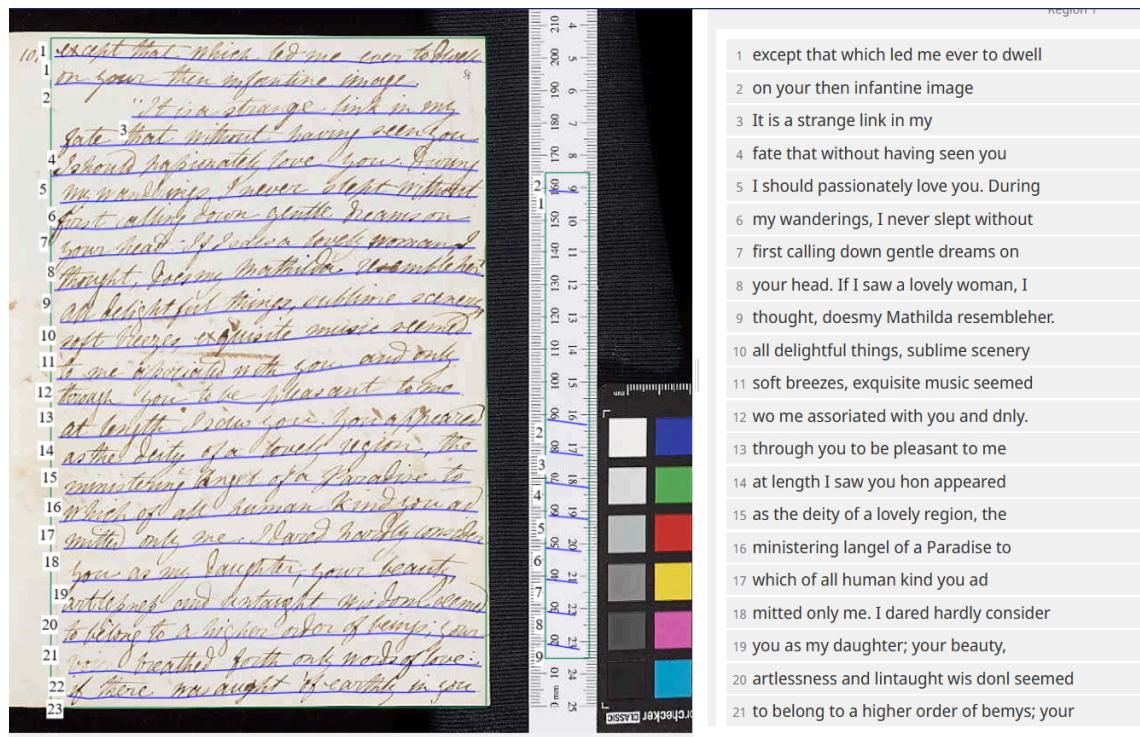


Figure 5 - Test HTR MaryModel with base model (from “Mathilda”, pag. 101)

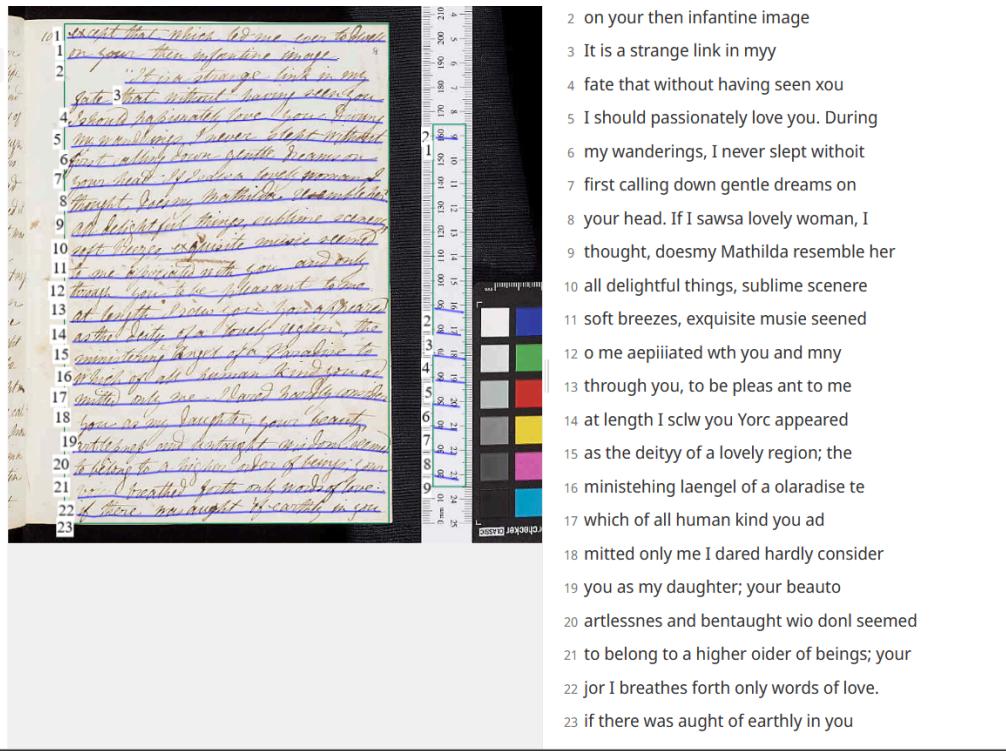


Figure 6 - Test HTR MaryModel without base model (from "Mathilda", pag. 101)

As showcased above, the results of the test are clearly more efficient and reliable in the case of the model that was fine-tuned. This test is completely in line with what was discussed in the previous paragraph.

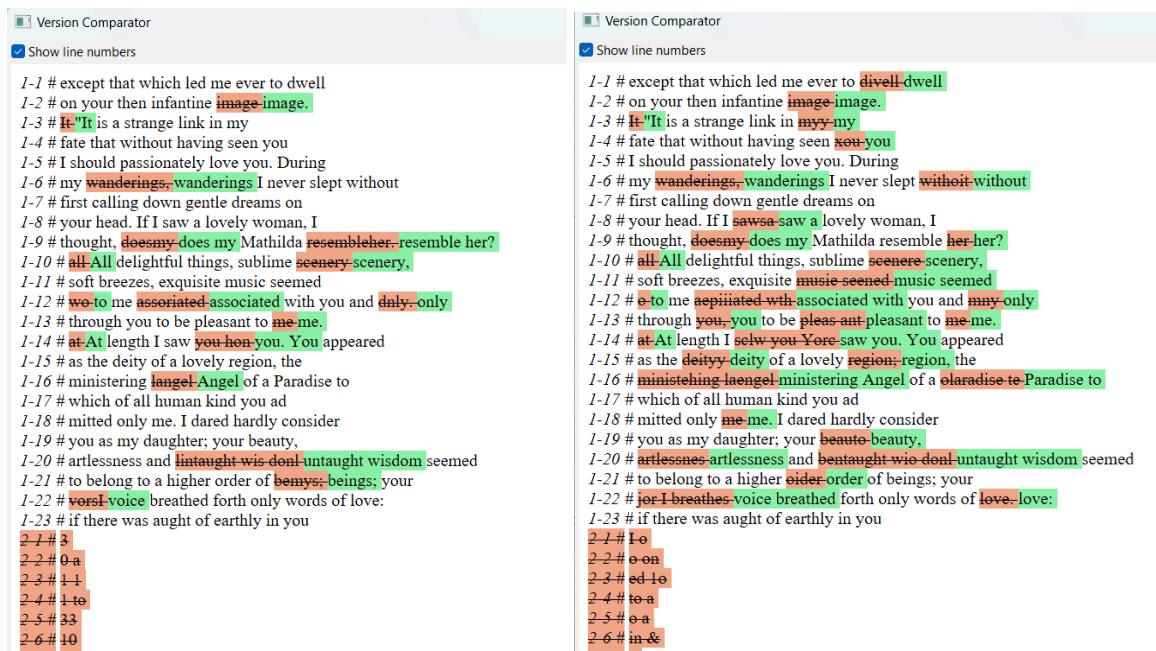


Figure 7 - Text Versions Comparison between the ground truth and the prediction of Mary Model with base model (on the left) and without base model (on the right)

Created	Status	Queries	Duration	Scope	Type	Results
14.05.24 10:05:28	Completed	Page(s) : 1 Option : Quick Comp...	0.46 sec.	Document ...	HTR	CER/WER: 9.55%/21.64%

Created	Status	Queries	Duration	Scope	Type	Results
14.05.24 09:52:08	Completed	Page(s) : 1 Option : Quick Comp...	0.35 sec.	Document ...	HTR	CER/WER: 14.09%/38.01%

Figure 8 - CER/WER percentages of the comparison between the ground truth and the prediction of Mary Model with base model (above) and without base model (below)

One issue we wanted to focus on is the fact that our model was not well-trained for the layout recognition: in fact the last few lines in *Figure 7* refer to the numbers of a ruler that was present in the image but was not supposed to be taken into consideration in the transcription. Of course we are aware of the fact that this issue could be solved at the beginning by cropping the image, so it's something we'll consider as a possible future development of the project.

Conclusions and critics

In conclusion, we are satisfied with the results obtained, and we believe that these models could serve as a useful starting point for future automation of the transcription and digitization process for the texts in the Shelley-Godwin archive. However, we want to highlight some critical issues encountered during the project. After completing the transcription process, we trained the model on Transkribus and had to wait three days for our request to be processed by the platform's server. Due to previous unsuccessful attempts, we didn't have time to recalibrate the model and make it more efficient. In our opinion, this limitation makes Transkribus less reliable in its open-access version.

Bibliography

Mercer, A. (2016). Beyond *Frankenstein*: The Collaborative Literary Relationship of Percy Bysshe and Mary Shelley. *The Keats-Shelley Review*, 30(1), 80–85.
<https://doi.org/10.1080/09524142.2016.1145937>

Muñoz, T., & Viglianti, R. (2014-2015). Texts and Documents: New Challenges for TEI Interchange and Lessons from the Shelley-Godwin Archive. *Journal of the Text Encoding Initiative*, Issue 8. Retrieved May 10, 2024, from
<journals.openedition.org/itei/1270>. <https://doi.org/10.30682/aldo2201c>

Sitography

The Shelley-Godwin Archive: <http://shelleygodwinarchive.org/>

Transkribus How To Guides:

<https://readcoop.eu/transkribus/resources/how-to-guides/>

- How To Train and Apply Handwritten Text Recognition Models in Transkribus:
<https://readcoop.eu/transkribus/howto/how-to-train-a-handwritten-text-recognition-model-in-transkribus/>