

# “I’m here to fight for ground truth”: HTR-United, a solution towards a common for HTR training data

**Chagué, Alix**

alix.chague@inria.fr

ALMAnaCH, Inria, France; Université de Montréal, Montréal, Canada

**Clérice, Thibault**

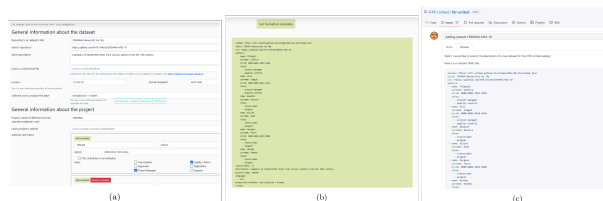
clerice.thibault@algorithme.net

Centre Jean Mabillon, PSL-Ecole nationales des chartes; ALMAnaCH, Inria, France

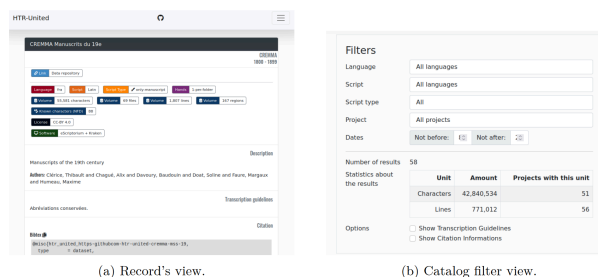
The growth of computation power and rise of artificial intelligence (in particular Deep Learning) allowed for the development of automatic text recognition, both on printed texts (OCR) and handwritten ones (HTR). Such technologies can now make millions of images of texts from various periods of time, held in patrimonial institutions, available for further search and processing.

HTR became more accessible when user friendly interfaces started to be developed: namely Transkribus from 2015 (Muehlberger et al. 2019) and eScriptorium from 2019 (Kiessling et al. 2019). In the case of HTR and old prints though, one of the hurdles remaining to be overcome is the access to robust models capable of recognizing coherent texts despite the multiple variations in handwriting or fonts. Such models usually necessitate users to produce large amounts of manual transcriptions considered as perfect-called ground truth-, which is a costly task. It requires a good understanding of the way deep learning functions, skills in paleography, and time. An easy way to reduce the costs of creating the training data to obtain a model is to rely on the data produced by other projects. Unfortunately, they are hard to find and not always published, because there is no incentive to put in this extra effort, neither for their publication nor for their documentation.

HTR-United is a collaborative initiative whose main purpose is to improve the findability of these open datasets, covering as many periods, scripts and languages as possible. Through this initiative, we support the creation a public catalog of dataset descriptions, contributed by individuals volunteering their own datasets. In general, descriptions are submitted as a YAML file filled with the help of a form available on HTR-United website (Figure 1)<sup>1</sup>. Raising awareness on the necessity to correctly document such shared datasets, HTR-United favors the implementation of the FAIR principles in the specific case of text recognition training datasets (Chagué / Clérice 2022b). The catalog<sup>2</sup> can be browsed using filters (script, language, type of font, period, etc.) and offer means to easily cite a dataset (Figure 2).



**Figure 1:** (a) Excerpt of the form to record the description of a new dataset; (b) YAML content generated by the form; (c) YAML description of a dataset submitted to HTR-United with Github.



**Figure 2:** View of records in the catalog: records can be seen in their own page (a) or browsed in the catalog, including after using filters (b).

The initiative is set up as an ecosystem of public Github repositories<sup>3</sup>, which guarantees the existence of precious versioning features for an ever-evolving catalog, transparency from all the parties as well as the possibility for us to rely on minimalistic developments. For example, anytime a dataset description is validated by our team, a Github Action processes all the existing descriptions in order to generate a new version of the catalog in the form of a pivot YAML file<sup>4</sup>: the catalog is never directly edited manually which reduces the risks of introducing errors. While a repository is dedicated to gathering all the descriptions feeding the catalog, another one hosts the specifications of the schema used to control the conformity of the descriptions<sup>5</sup>. Anyone can open a discussion to suggest the addition of new features in the specifications, or access the details of the arguments having led to the modification of the schema. Additionally, we aim to provide and maintain a suite of tools, available locally or through Github Actions and continuous integration, which help control, document and manage dataset on the short and long term, specifically in heavily collaborative contexts<sup>6</sup>.

During the DH2023 conference, we would like to introduce HTR-United to the international DH community by presenting how the ecosystem is organized, how contributors can submit new entries to the catalog as well as the stakes of contributing to such an initiative. HTR-United can be useful for the entire community of users of HTR technologies as the datasets listed in the catalog cover more and more languages or writing systems.

We would like to present some of the most interesting outcomes of such a collaborative catalog. First, various generic models for HTR were trained thanks to having access to a great variety of ground truth datasets<sup>7</sup>, which are of tremendous importance for the successful development of HTR for the humanities and the cultural institutions. The existence of such models allows smaller institutions or groups of researchers to quickly train robust models by simply fine-tuning generic models instead of starting from scratch (Chagué et al. 2021). Secondly, one of the most exciting aspects of possessing such a space for exchanging information

about ground truth datasets is the fact that it creates opportunities to pave the way towards a (international) standardization of transcription practices in the context of ground truth creation.

## Notes

1. See <https://htr-united.github.io/document-your-data.html> (27/04/2023).
2. See <https://htr-united.github.io/catalog.html> (27/04/2023).
3. See <https://github.com/HTR-United> (27/04/2023).
4. See in particular <https://github.com/HTR-United/htr-united/blob/master/htr-united.yml> (27/04/2023).
5. See <https://github.com/HTR-United/schema> (27/04/2023).
6. See <https://htr-united.github.io/actions.html> (27/04/2023).
7. See the CREMMA Medieval model (Pinche 2023), the Manus McFrench models (Chagué / Clérice 2022a) and other experiments on large training datasets such as Hodel et al. (Hodel et al. 2021).

## Bibliography

**Chagué, Alix / Clérice, Thibault** (2022a). *HTR-United—Manus McFrench VI (Manuscripts of Modern and Contemporaneous French)*. <https://doi.org/10.5281/zenodo.6657809> (27/04/2023)

**Chagué, Alix / Clérice, Thibault** (2022b, June 23). *Sharing HTR datasets with standardized metadata: The HTR-United initiative*. Documents anciens et reconnaissance automatique des écritures manuscrites. <https://inria.hal.science/hal-03703989> (27/04/2023)

**Chagué, Alix / Clérice, Thibault / Romary, Laurent** (2021, November 15). *HTR-United: Mutualisons la vérité de terrain ! DHNord2021 - Publier, partager, réutiliser les données de la recherche# : les data papers et leurs enjeux*. <https://hal.science/hal-03398740> (27/04/2023)

**Hodel, Tobias / Schoch, David / Schneider, Christa / Purcell, Jake** (2021). *General Models for Handwritten Text Recognition: Feasibility and State-of-the Art. German Kurrent as an Example* (No. 0). 7(0), Article 0. <https://doi.org/10.5334/johd.46> (27/04/2023)

**Kiessling, Benjamin / Tissot, Robin / Stokes, Peter / Stökl Ben Ezra, Daniel** (2019). eScriptorium: An Open Source Platform for Historical Document Analysis. *2019 International Conference on Document Analysis and Recognition Workshops (ICDARW)*, 2, 19–19. <https://doi.org/10.1109/ICDARW.2019.10032> (27/04/2023)

**Muehlberger, Guenter / et al.** (2019). Transforming scholarship in the archives through handwritten text recognition: Transkribus as a case study. *Journal of Documentation*, 75(5), 954–976. <https://doi.org/10.1108/JD-07-2018-0114> (27/04/2023)

**Pinche, Ariane** (2023). *Generic HTR Models for Medieval Manuscripts. The CREMMALab Project*. <https://hal.science/hal-03837519> (27/04/2023)