



**HAL**  
open science

# Which TEI representation for the output of automatic transcriptions and their metadata? An illustrated proposition

Hugo Scheithauer, Alix Chagué, Laurent Romary

## ► To cite this version:

Hugo Scheithauer, Alix Chagué, Laurent Romary. Which TEI representation for the output of automatic transcriptions and their metadata? An illustrated proposition. 2022. hal-04001303v2

**HAL Id: hal-04001303**

**<https://inria.hal.science/hal-04001303v2>**

Preprint submitted on 18 Sep 2023 (v2), last revised 30 May 2024 (v3)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

# Which TEI representation for the output of automatic transcriptions and their metadata? An illustrated proposition

Hugo Scheithauer<sup>1</sup>, Alix Chagué<sup>1,2,3</sup>, Laurent Romary<sup>4</sup>

<sup>1</sup> ALMAAnaCH - Automatic Language Modelling and ANALysis & Computational Humanities, Inria Paris

<sup>2</sup> UdeM - Université de Montréal

<sup>3</sup> EPHE - École pratique des hautes études

<sup>4</sup> Inria, Directorate for Scientific Information and Culture

## Abstract

The recent and fast development of automatic transcription software is accompanied by a growing heterogeneity of formats to save the output of such a task. TEI P5 can be helpful to simplify workflows and bring in more coherence in digitization pipelines. We present a twofold modelization in TEI which brings together essential information resulting from the transcription phase with the editorial layers. The usefulness of this modelization is illustrated with several examples showing how such an approach can be leveraged at different stages of a digitization pipeline.

Keywords: HTR ; OCR ; digital edition ; metadata ; layout ; eScriptorium ; TEI Publisher.

## Introduction

The recent growth of computer power available for processing data with machine learning techniques – particularly for Deep Learning contexts – have led to the faster development of technologies such as Optical Character Recognition (OCR), and recently Handwritten Text Recognition (HTR). It is now possible to envision reading manuscripts from the 20th century (Massot, Sforzini, and Ventresque 2019) as easily as those from the 13th century.<sup>1</sup> The main challenge for automatic transcription is the variation of letter shapes through time and space, even within the same alphabetical system. However, thanks to the development of the appropriate software and tools with rather accessible interfaces, the technology has now moved from Computer

---

<sup>1</sup> See for example the CREMMA Lab post-doctorate of Ariane Pinche (ENC - CJM - INRIA): <https://cremmalab.hypotheses.org/> (viewed on 02/12/2022)

Vision labs to those in Social Sciences and the Humanities as well as to cultural institutions. Such software accelerates our capacity to create texts which can then be used for research projects or for the digital editions of large patrimonial collections. These evolutions are highlighted by the emergence of new workflows for digitizing and editing texts. It has now become crucial to address the question of knowing how we can keep and organize all the information deriving from the automatic transcription process in a reliable and sustainable manner.

## Expanding the role of the TEI guidelines in the context of HTR

In the field of automatic text recognition, whether applied to printed documents or to manuscripts, two *de facto* standards are currently used to record the output of the recognition process: ALTO XML<sup>2</sup> and PAGE XML (Pletschacher & Antonacopoulos, 2010).<sup>3</sup> These two XML formats are mainly intended to render the details of the layout of the original documents in combination with the textual elements that have been recognized. They thus follow a representation paradigm which is far from sufficient when the underlying objective is to produce complex digital editions. On the contrary, TEI XML is specifically appropriate for digital editions with the specificity that it enables one to simultaneously handle the intellectual organization of the document but also its physical incarnation. ALTO XML and PAGE XML also work at the image level, whereas TEI allows to handle whole documents into a single XML file. For these reasons, it is common that once the transcription process is over, and the edition begins, we move on to TEI XML and keep only the information necessary for the edition. This usually forces projects to manage an aggregate of heterogeneous formats in order to a) keep the precise information resulting from the HTR process and b) enrich the content in the form of a structured edition, while avoiding to lose the connection between both.

We would like to advocate here for a change of paradigm where the TEI guidelines would be given a stronger and earlier role in these workflows, with the advantage of better organizing the articulation between the various transcription, edition and enrichment levels. Thus, our question is the following: how far can we map the results of an automatic transcription process to TEI and by which representational means? We also argue that both ALTO XML or PAGE XML can be used as intermediary representations to generate a comprehensive TEI file containing all data produced by automatic text recognition models.

Initially arising from reflections conducted within the framework of the LECTAUREP project (which we will present later), we maintained an effort to generalize our

---

<sup>2</sup> See the Analyzed Layout and Text Object (ALTO) 4.2 schema specifications at <https://www.loc.gov/standards/alto/news.html#4-2-released> (viewed on 03/23/2022)

<sup>3</sup> See also the PAGE XML Github repository: <https://github.com/PRImA-Research-Lab/PAGE-XML> (viewed on 03/24/2022)

proposal to as many use cases as possible. Our experiment is based on the use of eScriptorium, a virtual research environment which combines a web interface and the OCR/HTR engine Kraken (Stokes et al., 2021).<sup>4</sup> Another notable HTR engine used by the Digital Humanities community is Transkribus (Mühlberger et al., 2019). They both offer the possibility to export the output of the automatic transcription in formats such as plain text, ALTO XML and PAGE XML.<sup>5</sup> Transkribus has also implemented the export to TEI XML<sup>6</sup> and other formats like Docx. Apart from ALTO and PAGE, all these formats generate representations that cannot be reinjected into the software since they lose a great deal of the segmentation and layout information. It is worth noting that ALTO and PAGE suffer from a lack of normalization and compatibility. For instance, in 2020, ALTO files exported with Transkribus were not compatible with eScriptorium import module because the former used ALTO 2 and the latter ALTO 4.<sup>7</sup> As of now, Transkribus switched to ALTO 4.0, and eScriptorium to ALTO 4.2, but this incompatibility still remains. Users also have the possibility to opt for the PAGE format, but this choice amplifies the potential for numerous incompatibilities across different formats and their respective versions. It also goes without saying that PAGE and ALTO are not compatible with each other, and without a clear choice from the community, archives in future will have to deal with two different formats that need to be converted back and forth, if the conversion is made possible at one point. PAGE also appears to receive infrequent updates, with its most recent update dating back to 2019. In contrast, ALTO demonstrates ongoing evolution, with its latest update being as recent as April 2023. Using TEI guidelines would allow final end users to rely on one single format to standardize automatic text recognition output in the context of textual edition, and would then ease its archival. On the contrary, one of our objectives is to improve the reusability of contents generated through HTR and OCR at various stages of the workflow leading to a digital edition. Using TEI XML to encode pivot files would mean aggregating the raw transcription with its edited version into one single file. Ideally, this would allow users to go back and forth between these two states of the text while avoiding complex manipulations of the data. Such a pivot file entails making the TEI encoding compatible with the requirements of automatic transcription software and their traditional formats as all information is encoded into TEI and can then be easily

---

<sup>4</sup> See also the Kraken repository on Github: <https://github.com/mittagessen/kraken> (viewed on 03/23/2022)

<sup>5</sup> The implementation of these two export formats does not go without raising issues, namely due to the absence of a normalization strategy which results in tedious transitions from one interpretation of the format to another (Transkribus' ALTO and eScriptorium's ALTO are not immediately compatible) and from one version to the next (Transkribus' shift from ALTO 2 to ALTO 4 in March 2021 broke many pipelines based on ALTO 2).

<sup>6</sup> Said XSLT can be found on the following Github repository: <https://github.com/dariok/page2tei> (viewed on 03/23/2022)

<sup>7</sup> See Aspyre-GT, a python tool developed in 2020-2021 to make Transkribus' ALTO 2 files compatible with eScriptorium. The transformation scenario was later extended to other sorts of ALTO files produced by software like LIMB or PDFALTO: <https://github.com/alix-tz/aspyre-gt> (viewed on 07/08/2023)

converted back,<sup>8</sup> but it could also improve the propagation of the metadata throughout the pipeline.

New issues are also currently arising following the creation of large corpora produced with automatic transcription technologies, such as the need to reconstruct the logical structure of digitized text documents (Clérice 2023). Having access to layout information allows semi-automatic and automatic processing in order to do so. Layout analysis is at the interface of computer vision and information extraction, and can be done manually with segmentation software, or with machine learning models (Clérice 2023; Oliveira et al. 2018). Converting annotations directly to TEI is not the scientific norm. To the best of our knowledge, only GROBID takes an image with a text layer as input, and outputs a TEI based on the inference of different CRF models using layout token features (Khemakhem 2020).<sup>9</sup> Otherwise, layout encoding is usually done manually or semi-automatically using scripts that convert ALTO XML and PAGE XML to TEI, based on text line annotations and content.

Our approach has also been inspired by the vision adopted in the “TEI in libraries” initiative, which also had in mind the idea of bringing together automatic transcription and more structured TEI-based representations. Born in 1999, the “TEI in libraries” initiative was grounded as a TEI workgroup<sup>10</sup> and released its last version in 2018. It listed best practices for using TEI in the context of library text digitization projects.<sup>11</sup> A notable aspect is that it proposed to keep OCR output in a TEI file, and established key concepts which echoes the direction we decided to take: TEI can “be suited to the goals of a preservation unit or mass digitization initiative,” and it allows texts “to be a faithful representation of the appearance of the source document derived from OCR.”<sup>12</sup> In this paper, we want to go a step further by considering text as data, as argued in (Burnard et al. 2021), and specifically address the challenges now posed by a greater use of OCR and HTR in a variety of new contexts, with a specific emphasis on modeling layout information in TEI documents.

## OCR and HTR terminology

Before entering into the details of our TEI modelization for the encoding of automatic transcription output, we would like to provide the reader with a few elements about the OCR and HTR terminology, which we will then map onto TEI concepts later.

---

<sup>8</sup> With an XSLT, or a script, for instance.

<sup>9</sup> See for example the machine learning library GROBID:

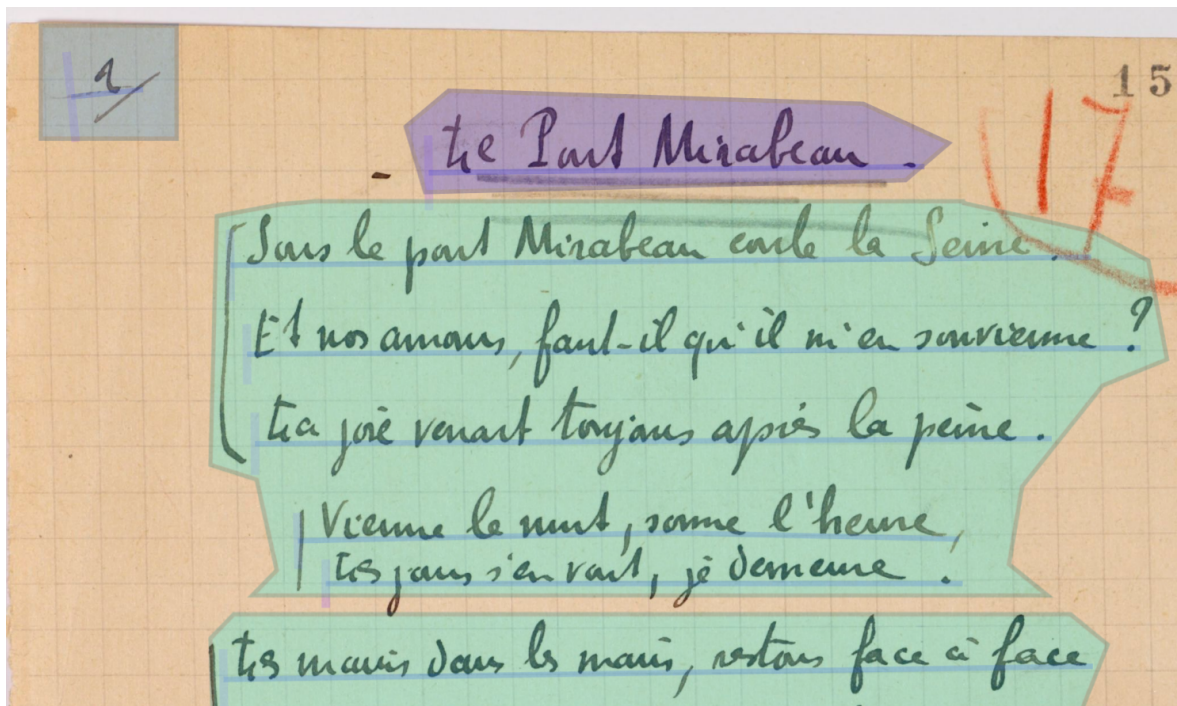
<https://grobid.readthedocs.io/en/latest/Principles/#layout-tokens-not-text> (viewed on 03/23/2022)

<sup>10</sup> See [https://wiki.tei-c.org/index.php/Workgroup\\_to\\_revise\\_the\\_Best\\_Practices\\_for\\_TEI\\_in\\_Libraries](https://wiki.tei-c.org/index.php/Workgroup_to_revise_the_Best_Practices_for_TEI_in_Libraries) and the workgroup’s Github repository: <https://github.com/kshawkin/Best-Practices-for-TEI-in-Libraries> (viewed on 03/25/2022)

<sup>11</sup> See <https://tei-c.org/extra/teiinlibraries/> (viewed on 03/23/2022)

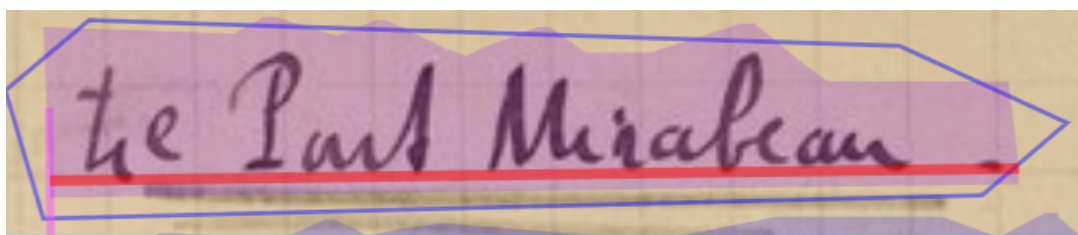
<sup>12</sup> See TEI in Libraries and especially the 4.0.0 version which can be found here: <https://tei-c.org/extra/teiinlibraries/4.0.0/bptl-driver.html> (viewed on 03/23/2022)

In most automatic transcription software components, an image can be divided into one or more **text regions** (also commonly called zones, or text zones). They are used to mark zones in a page that usually bear a semantic significance. Figure 1 represents a poem, “Le Pont Mirabeau,” written by Guillaume Apollinaire, circa 1912, kept at the French National Library (BnF). The layout is rather simple, and consists of a page number, a title, and stanzas. All of them are directly identified on the image as text regions. We should mention that there exist initiatives that intend to build ontologies to normalize the naming of the zones and their use, such as SegmOnto (Gabay et al. 2021).



[fig. 1] Layout annotation of a poem written by Guillaume Apollinaire, *Le Pont Mirabeau*, ca. 1912, Bibliothèque nationale de France (BnF). <https://gallica.bnf.fr/ark:/12148/btv1b525056707/f33/>.

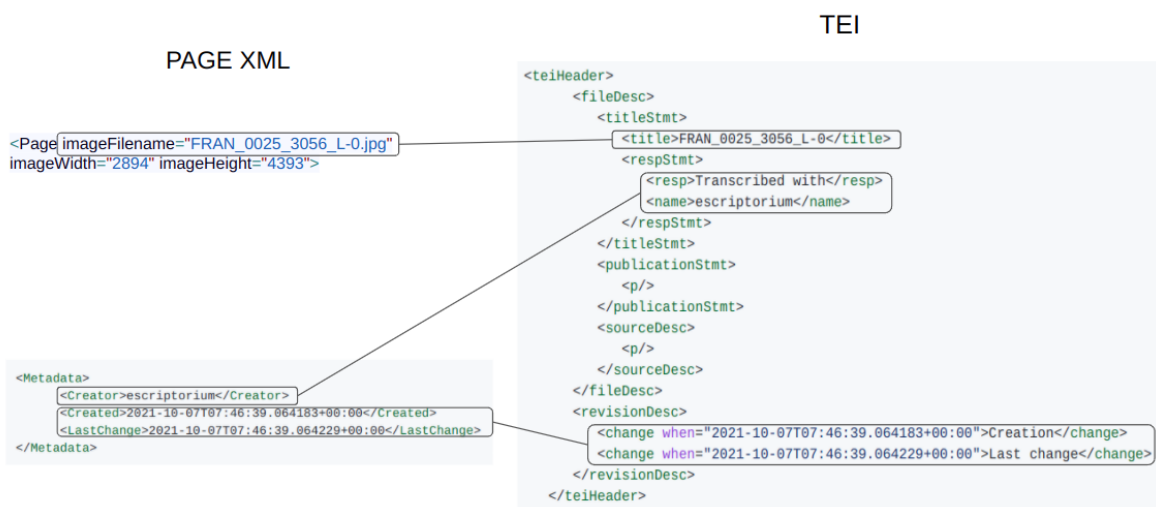
A zone normally contains lines of text. They are themselves composed with three elements (see Figure 2): a **baseline/topline** defines a virtual line, passing through at least two points, on which the text is written or from which it is hanging; a **mask** is a polygon, defined by at least three points, which delimits the area of pixels containing the text of the line; and lastly, the **text** itself.



[fig. 2] A text line corresponding to the poem's title, its baseline (red), and mask (purple). Guillaume Apollinaire, *Le Pont Mirabeau*, ca. 1912, Bibliothèque nationale de France (BnF).  
<https://gallica.bnf.fr/ark:/12148/btv1b525056707/f33/>

## Encoding automatic text transcription in TEI

As of March 2022, the ALTO and PAGE XML formats used in eScriptorium to export the recognition output are quite rudimentary regarding metadata elements about the transcription itself. In the resulting PAGE XML files, for instance, it includes the identification of the creator (i.e. the application used: eScriptorium) and changes made to the transcription on the platform, such as the dates of creation and of the last modification.<sup>13</sup> Such information is quite straightforward to map onto **<teiHeader>** components: in **<respStmt>** for information about the users and the software, and **<revisionDesc>** for temporal information. This can be complemented with another mapping, for information such as the title of the document or any identifier linked to the image, within **<titleStmt>** (see Figure 3).



[fig. 3] Documenting a transcription: metadata representation from ALTO XML to TEI. Created by the authors.

As of now, eScriptorium does not provide enough metadata for ensuring FAIRness. However, with the continuous development of automatic transcription software, more complex and detailed metadata will become available in the future. The **<teiHeader>** already offers the needed components to fully document most of those we can already foresee. For instance, the metadata attached to the transcription models – such as the name of the model, who trained it, the OCR/HTR engine used and its version, or the list of codecs known by the model –, could be given in the

<sup>13</sup> Note that, as of March 2022, the current version of eScriptorium considers that date of creation equals the date of last modification.

**<respStmt>**. The name of the software and of the individuals or organizations who modify the transcription afterwards could be indicated there as well.

Details on the transcription itself, such as the transcription guidelines, could be supplied in the **<editorialDecl>** (see Figure 4). This element already offers a set of sub-elements which can be used to describe regular decisions often made before training and applying a transcription model, such as hyphenation (**<hyphenation>**), normalization (**<normalization>**), i.e. how diacritics and other typographic elements were transcribed; or else how punctuation (**<punctuation>**) is handled and transposed into unicode characters. The output of OCR and HTR often needs to undergo manual or automatic correction to remove the remaining errors. Post-transcription corrections, either done automatically, semi-automatically or by hand, can thus be detailed in **<correction>** (Rigaud et al., 2019; Nguyen et al., 2019).<sup>14</sup>

```
<editorialDecl>
  <normalization>
    <p>In the training corpus, abbreviations were transcribed using "^".</p>
  </normalization>
  <correction>
    <p>No post-transcription corrections.</p>
  </correction>
  <hyphenation eol="all">
    <p>All hyphenations were kept. They were transcribed with "-".</p>
  </hyphenation>
</editorialDecl>
```

[fig. 4] An example of **<editorialDecl>**. Created by the authors.

As for the representation of the raw transcription itself, we chose to keep a clear separation between it and the actual edited text. While the latter would go in the **<body>** element, all the raw output of the transcription is defined in a **<sourceDoc>** element. As stated in the TEI P5 Guidelines, **<sourceDoc>** contains the transcription or other representations of a single source document.<sup>15</sup> Several child elements are available, namely: **<graphic>**, **<surface>** and **<zone>**.

Our implementation of **<sourceDoc>** follows two key principles. First, **<sourceDoc>** must be the strict transposition of any output information resulting from the HTR or OCR process. For instance, since ALTO and PAGE XML exports include masks' coordinates, we need to reflect these in the TEI representation. Secondly, we want to

---

<sup>14</sup> See the TEI guidelines for more documentation about said elements: respectively, <https://www.tei-c.org/release/doc/tei-p5-doc/en/html/ref-normalization.html>, <https://www.tei-c.org/release/doc/tei-p5-doc/en/html/ref-punctuation.html>, <https://www.tei-c.org/release/doc/tei-p5-doc/en/html/ref-correction.html> (viewed on 03/23/2022)

<sup>15</sup> See the TEI Guidelines for the **<sourceDoc>** element: <https://www.tei-c.org/release/doc/tei-p5-doc/fr/html/ref-sourceDoc.html> (viewed on 03/24/2022)



keep **<body>** free from any HTR or OCR information, meaning that what we find there will be the sole responsibility of the editor. As a result, any interpretation of the output of the transcription will be contained in **<body>** and never inside **<sourceDoc>**. Distinguishing the raw transcription from the edition, understood as a stable transcription with information on its layout, guarantees a continuum between the original document and its final publication. We deliberately decided to double the textual content in order to have two distinct blocks: the initial output of the transcription, and its edition (see Figure 5).

In conjunction with our specification work to map the OCR/HTR output to TEI, we have implemented a full-fledged XSLT transform from PAGE XML to TEI which generates the content of **<sourceDoc>** as presented here.<sup>16</sup> Note that a similar transformation scenario could be created to go from ALTO to TEI. We managed to ensure that all elements available in a PAGE XML document could be retrieved during the transformation and transposed into the resulting TEI document.

In a PAGE XML document, a **<Page>** element represents the transcribed image. Basic metadata can be assigned to various attributes, for instance **@imageFilename** will store the name given to the source image, and **@imageWidth** and **@imageHeight** give a set of x,y values in pixels, defining a two dimensional space bearing text lines. With TEI, and nested inside **<sourceDoc>**, we use **<graphic>** for documenting such information, with the **@url**, **@width** and **@height** attributes. In addition, we give the image an identifier with **@xml:id** for linking it later to its transcription, allowing one to go back and forth between its raw state and its edition.

We then use **<surfaceGrp>** to represent the sum of all text regions and their associated text lines for a given image. **<surface>** and **<surfaceGrp>** are linked using an **@xml:id** attribute on the first element, and **@facs** for the second. For each text region, a **<surface>** element is created and nested within this parent element.

PAGE XML uses **<textRegion>** and **<Coords>** elements to document text regions. The former gives an identifier, and the latter gives its coordinates as pairs of points, which can be located on the image. In TEI, we distribute this information with three attributes inside **<surface>**: **@xml:id** gives an identifier; **@type** will be associated with values defined with an ontology in the transcription software to qualify the different text regions (for example “numbering”, “title” and “line\_group”, as seen in our example in Figure 1). Finally, **@points** indicates the coordinates of the text region. This last attribute will be found in any other element bearing layout information and defining an area or a line.

Nested in **<surface>**, each text line associated with a text region is represented with **<zone>** and its children elements: **<path>** and **<line>**. **<zone>** corresponds to the

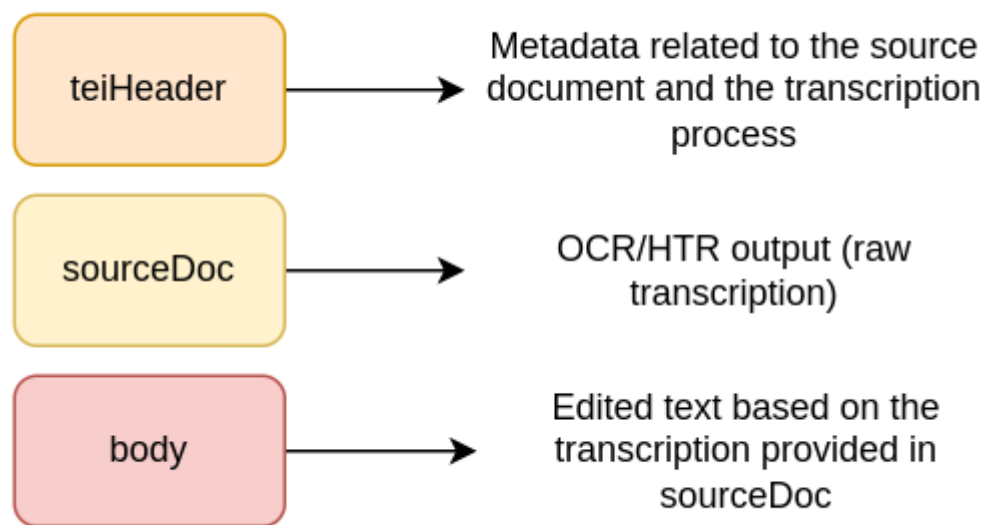
---

<sup>16</sup> The XSLT is available on Github repository: <https://github.com/TEI4HTR/page2tei> (viewed on 03/24/2022)

mask of the text line. Then, **<path>** and **<line>** are respectively equivalent to the baseline/topline and the text node (see Figure 6).

The **<body>** can then be encoded based on the content of the transcription, text lines and ideally layout information, stored in the **<sourceDoc>**, as found in ALTO XML or PAGE XML.

Fully encoded TEI examples can be found on Github.<sup>17</sup>



[fig. 5] A sourceDoc-based encoding components. Created by the authors.

```
<sourceDoc>
<graphic xml:id="f33" url="ark:/12148/btv1b525056707/f33/" width="2312px"
height="3469px"/>
  <surfaceGrp facs="#f33">
    <!-- ... -->
    <surface xml:id="eSc_textblock_2ed880a5" type="title" points="1594,196
1744,265 1668,330 827,351 793,236 868,173">
      <zone xml:id="eSc_line_08a9f120" type="mask" points="833,317 827,179
998,185 1024,170 1027,170 1027,170 1027,170 1030,170 1030,170 1033,170
1033,170 1033,170 1035,170 1090,187 1125,173 1125,173 1128,173 1128,173
1128,173 1131,173 1131,173 1131,173 1134,173 1134,173 1189,199 1206,205
1218,199 1270,179 1270,179 1273,179 1273,179 1273,179 1276,179 1276,179
```

<sup>17</sup> Two fully encoded TEI examples can be found at [https://github.com/TEI4HTR/page2tei/blob/main/ressources/FRAN\\_0025\\_0227\\_L-0.tei.xml](https://github.com/TEI4HTR/page2tei/blob/main/ressources/FRAN_0025_0227_L-0.tei.xml) and [https://github.com/TEI4HTR/page2tei/blob/main/tei/32\\_c42c1\\_default.tei.xml](https://github.com/TEI4HTR/page2tei/blob/main/tei/32_c42c1_default.tei.xml). The first file results from an automatic transcription, and the second was manually transcribed (viewed on 03/29/2022)

```
1278,179 1278,179 1278,179 1281,179 1313,202 1371,202 1452,176 1455,176
1455,176 1455,176 1458,176 1458,176 1461,176 1461,176 1461,176 1464,176
1464,176 1513,216 1530,228 1692,228 1701,306 1701,320 833,332">
  <path type="baseline" points="833,317 1701,306"/>
  <line>Le Pont Mirabeau.</line>
</zone>
</surface>
<!-- ... -->
</surfaceGrp>
</sourceDoc>
```

[fig. 6] An example of the <sourceDoc>. Created by the authors.

## The example of LEPIDEMO: TEI as a pivot file

As stated earlier, we do not envision **<sourceDoc>** as the place where the transcription is interpreted, nor where the edition takes place. However, the text provided in **<sourceDoc>** will often be the foundation for later editions and interpretations which will be rendered in **<body>**. The LECTAUREP project provided a use case to demonstrate how a TEI file built according to our modelization can usefully serve both the raw transcription and its edition.

LECTAUREP is a project jointly led by Inria (ALMAAnaCH team) and the Archives nationales de France (DMC team) between 2018 and 2021 (Chagué and Rostaing, 2021). It aimed at facilitating the exploration of directories listing minutes and deeds redacted by Parisian notaries between the beginning of the 19th century and the mid-20th century. Once the transcriptions were acquired with eScriptorium, we wanted to detect named entities before publishing the final enriched transcriptions. However, due to the complex layout used in the documents, we realized that it would be difficult to automatically parse and analyze hundreds of pages without structuring the output of the transcription. Such structuration would help us target specific text regions where named entities can be found. As illustrated in Figure 7, each page contains a table consisting of a header and several columns (usually seven) dividing the information in different categories.

*Seizième dix neuvième feuille*  
79

N <sup>os</sup> DU RÉPERTOIRE	DATES DES ACTES	NATURE ET ESPÈCE DES ACTES :		NOMS, PRÉNOMS ET DOMICILES DES PARTIES  INDICATIONS, SITUATIONS ET PRIX DES BIENS	RELATION de l'Enregistrement.		
		EN BREVETS	EN MINUTES		DATES	DROITS	
<b>An 1937, mois d'Avril</b>							
296	23		Donation	Morisset (par Mari) à son mari sus-nommé Pierre-Propriété			
297	24	Certificat de Propriété		Perisard d'une action Bon Marché au nom de Joseph 64 Rue Linc - à Paris, décedé	26 Av.	11 50	
298	24	Certificat de Propriété		Dumont (de 1515) le comte Crisse d'Espagne au nom de Arsène Mari) 2 Rue Hoche à Paris décedé.			
299	24		Prêt	Lucas (par Et Fournier de France 19 Rue des Capucines Paris à Marie Auguste) 5 Rue Germain Pétay à Paris 155.000 <sup>f</sup>	27 Av.	45 -	
300	26		Restation de servitude	20 ans de 6 <sup>h</sup> 30 <sup>0</sup> Dubonnet (par Robert) 19 Av. de Versailles à Paris, servitude du 15-12-36 devant Collet gmt. à Nantes	d <sup>e</sup>	22 50	
301	26 (1/2) (mars et)		Bail	Normandier (par Germain) Machine à vapeur de Jean Chiron (Joseph) 1 Boulevard Bineau à Wallons Terret à Robert Lemaire de Gall, 111. Rue de la Garonne à Courbevoie, louée à Wallons 1450 Bineau 6 ans 1900 <sup>f</sup> + 500 <sup>f</sup>	28 Av.	99 10	
302	26		Vente	Henrat (par Louis Auguste) et Adèle Marie Amélie Susbrace sa de 154 Rue de Grenelle Paris à Jean Louis Gabriel Helvie Roche	d <sup>e</sup>	4.961 -	

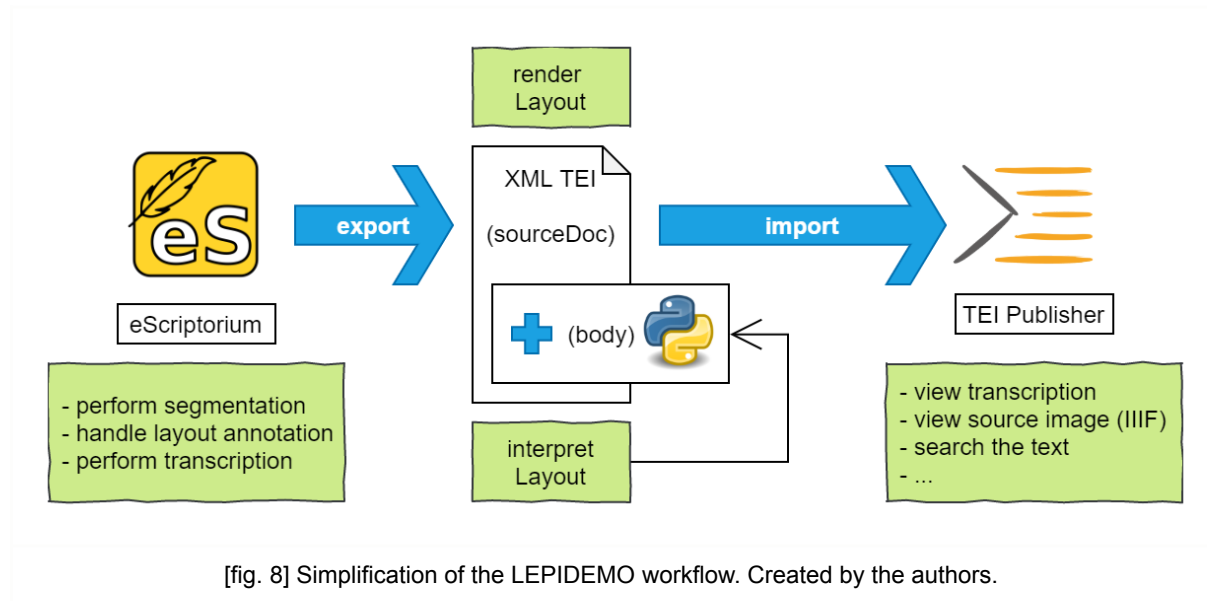
[fig. 7] A notary directory dated from April, 1937 (inventory number: FRAN\_0025\_4648\_L-1). From left to right, the header indicates: repertory number; date; the type of deeds type; the persons who signed it and a summary; when it was paid and the sum paid to the notary. With a fully encoded table in TEI, it would theoretically be possible to aim for a specific column and then tag elements according to its nature. For example, we could tag every type of act in the third and fourth column with simple rules because the information is rather simple, consisting only of tokens specifically describing the type of act. In our use case, named entity recognition would then be performed semi-automatically, without the need to train a model for every column, and automatically when information is syntactically more complex (e.g. the fifth column).

We built a demonstration pipeline called LEPIDEMO (for LECTAUREP Pipeline Demonstration, see Figure 8) during which images are first loaded on eScriptorium, where the layout is annotated and the documents transcribed before being exported as PAGE XML files.<sup>18</sup> The Page XML files were transposed into TEI with an XSLT.<sup>19</sup> All files could then be aggregated to have a single file representing the whole document. As mentioned before, the extracted information populated **<teiHeader>** and **<sourceDoc>** only. Then, we processed coordinates available in **<sourceDoc>** thanks to a Python script in order to reconstruct each page's logical structure. This helped us create the content of **<body>**: mainly a **<table>** element with child elements such as **<row>** and **<cell>** encoding its content. We also created relations between the **<sourceDoc>** and the **<body>**, linking the lines of texts by means of a pair of **@xml:id** in the former, and **@facs** in the latter. Finally, we were able to normalize dates given that one line specifies the year and month applicable to the whole page while a column gives only the date of the day applicable to the minute (see Figure 8). We are also able to narrow down the scope of the search for named

<sup>18</sup> See the pipeline at <https://github.com/lectaurep/lepidemo> (viewed on 03/28/2022). Restructured files could be accessed at: <https://github.com/lectaurep/lepidemo/tree/master/data/output> (viewed on 03/29/2022)

<sup>19</sup> We also developed a python script that allows users to download TEI files from eScriptorium. See <https://github.com/lectaurep/TEI-From-eScriptorium> (viewed on 03/28/2022). We hope to implement a TEI export directly in eScriptorium in the future.

entities in the text. Corrections and the annotation of named entities would be appended to the content of **<body>**.



With the example of LEPIDEMO, we show the importance of distinguishing the raw output of OCR/HTR from the edition phase where the coordinates are interpreted and the text annotated. Without this interpretation of **<sourceDoc>**, we would not be able to make sense of the different resulting text lines. However, the content of **<body>** is irrelevant to an OCR/HTR software because in order to become readable the information no longer follows the physical structure rendered on the image. Keeping both **<sourceDoc>** and **<body>** in the same document enables one to show a physical structure compatible with an OCR/HTR software (to train a new transcription or segmentation model, for example) while still being able to move on to more detailed edition steps.

*Sommaire des neuvième-feuilles*

N <sup>os</sup> DE RÉPERTOIRE	DATES DES ACTES	NATURE ET ESPÈCE DES ACTES :		NOMS, PRÉNOMS ET DOMICILES DES PARTIES	RELATION de l'Enregistrement.	
		EN-BREVETS	EN-MINUTES		DATES	DROITS
				An 1937, mois d'avril		
286	23		Donation	Morlesan (épouse) à son mari sur nommé Pierre-François		
297	24	Certificat de Propriété		Penissard (d'une action Bay March au nom de Joseph)	26 Ar.	11 50
298	24	Certificat de Propriété		Beuret (de 51512 d'arr. de Crissac d'Espagne au nom de Arsène		
299	24		Prêt	Mari) à Rue Hoche à Paris, décès	27 Ar.	45 -
300	26		Restation de somme	à Paris (Paris) Foncin de France 19 Rue des Capucines Paris		
301	26 (1/2 Mars et)		Bail	à Paris (Paris) 5 Rue Germain Pilon à Paris 15.000 <sup>fr</sup>		
				20 ans à 6 <sup>fr</sup> 50		
				Dubonnet (par Robert) 19 Av. de Versailles à Paris, inventaire	dé	11 50
				du 15-12-36 devant Collot gmt. de Nantes		
				Normandin (par Germain) Bruchon et de Prêt venant de Jean Chiquet	28 Ar.	39 10
				(Joseph) à Paris (Paris) Bureau à Lulliers Verret à Robert Lince		

when="1937-04-24"

[fig. 9] Building an ISO-compliant representation of the date based on the elements present on the page (excerpt from FRAN\_0025\_4648\_L-1.jpg). Created by the authors.

## Using layout information when publishing TEI

Having access to layout information inside the final encoded text also allows for further usage during the publication phase.

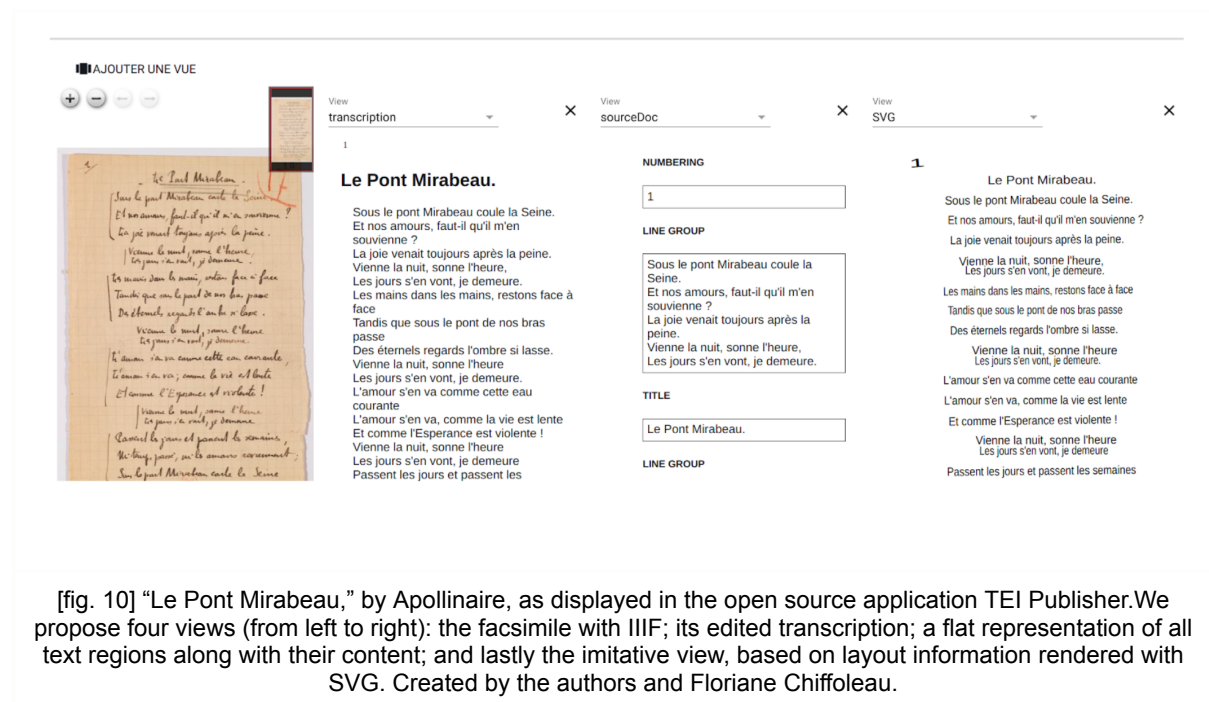
With TEI Publisher (e-editions 2021), we propose a solution to display the resulting TEI file where the transcription output is associated with an image facsimile. TEI Publisher is an open source publication application powered by eXist-db. It offers a fully customizable interface, and uses ODDs and web templates to visualize a corpus of files encoded with TEI. Furthermore, it allows faceted search and corpora exploration.<sup>20</sup>

We used the modelization presented above to propose a visualization of the resulting TEI files hinged on four views, as illustrated in Figure 10. The first one is a IIIF viewer displaying a facsimile: a IIIF identifier referring to the image is stored with a `@url` inside `<graphic>` in `<sourceDoc>`. Secondly, a diplomatic representation of the document is based on the content of the `<body>` element. Then thirdly, a flat representation of the text provided in sourceDoc is displayed region by region, allowing users to see the original image's segmentation. Lastly, an imitative transcription of the text, based on the information contained in `<sourceDoc>`, including the canvas' size, the segments' coordinates and their associated text

<sup>20</sup> See TEI Publisher's website: <https://teipublisher.com/index.html> (viewed on 03/28/2022)

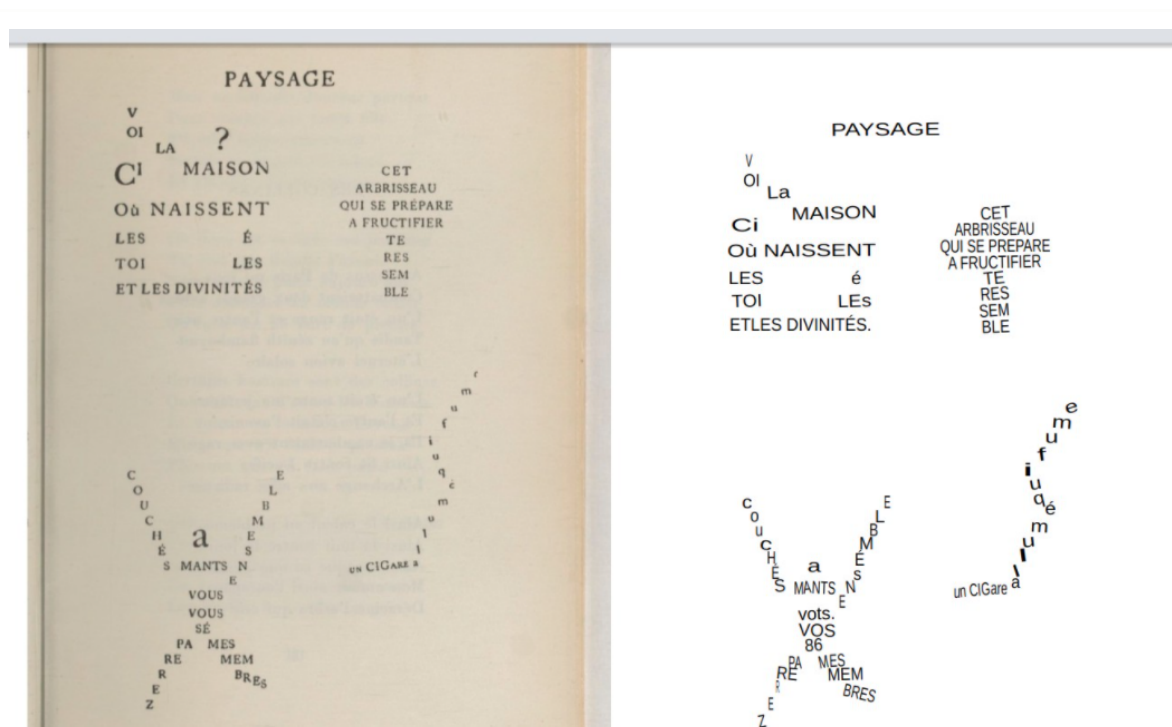
nodes are rendered thanks to SVG, an XML-based vector image format for two-dimensional graphics.<sup>21</sup>

Whole documents would then be displayed on an interface such as TEI Publisher with its layout thanks to the information stored in the `<sourceDoc>`. However, as TEI files become bigger, hardware issues are at stake. Splitting such a TEI file on the `@xml:id` specified in `<graphic>` and nested in the `<sourceDoc>` could allow bypassing hardware limitations.



The SVG visualization provides a solution to confront a transcription with its original layout. This gives the community an opportunity to publish text with a complex layout, as seen in Figure 11, where the meaning of the text is extremely tied to the layout.

<sup>21</sup> See the Scalable Vector Graphics (SVG) 2 recommendations at <https://www.w3.org/TR/SVG2/> (viewed on 03/28/2022)



[fig. 11] A visualization in TEI Publisher of a facsimile of calligrams written by Guillaume Apollinaire. Guillaume Apollinaire, *Calligrammes*, ed. Roger de La Fresnaye, Lausanne, 1952, <https://gallica.bnf.fr/ark:/12148/bpt6k9775732c/f31>. Created by the authors. Note: The transcription was acquired with an HTR model and was not corrected, hence some characters appearing as lowercase although they should be uppercase.

## Conclusion

We argued that all elements resulting from an OCR or an HTR process can be mapped to a **<sourceDoc>** element as defined by the TEI Guidelines. With a bipartite organization revolving around the articulation of **<sourceDoc>** and **<body>**, we can leave the output of the transcription intact and render any interpretation or correction of the text (otherwise incompatible with a transcription software) inside **<body>**, with the appropriate elements. By adopting this architecture, we propose that digitization pipelines switch earlier to TEI and moving away from relying on ALTO or PAGE as an archival format, with the objective of simplifying the workflow and better documenting the transcription phase and its later edition, either for keeping track of this process, or for further usage, such as described in our examples. This is particularly crucial for archival purposes, as it allows to rely on a single well-maintained format. OCR, and HTR in particular, are rapidly evolving fields, this concerns formats as much as software. Our modelization is intended as a first step towards a better stabilization of the landscape thanks to TEI. New use cases and approaches will certainly arise, we therefore welcome any feedback on that matter and are open to collaborations.



# Acknowledgments

We would like to mention the essential contribution of our colleague Floriane Chiffolleau on the elaboration of the visualization prototype with TEI Publisher. We also wish to thank Simon Gabay, Juliette Janes, Claire Jahan and Ariane Pinche for their insightful feedback.

# Bibliography

- Burnard, Lou, Christof Schöch, and Carolin Odebrecht. 2021. "In search of comity: TEI for distant reading." *Journal of the Text Encoding Initiative* (Issue 14). doi: 10.4000/jtei.3500.
- Clérice, Thibault. 2023. "You Actually Look Twice At It (YALTAi): Using an Object Detection Approach Instead of Region Segmentation within the Kraken Engine."
- Chagué, Alix, and Rostaing Aurélia. "LECTAUREP: Paris Notary Record Books Automated Reading." In *Fantastic Futures 2021 / Futures Fantastiques 2021*. Paris, France: AI4LAM and BnF and Université Paris Saclay, 2021. <https://hal.inria.fr/hal-03479258>.
- e-editions. "TEI Publisher: Documentation." TEI Publisher, August 2021. <https://teipublisher.com/exist/apps/tei-publisher/doc/documentation.xml?odd=docbook.odd>.
- Gabay, Simon, Jean-Baptiste Camps, Ariane Pinche, and Claire Jahan. "SegmOnto: Common Vocabulary and Practices for Analysing the Layout of Manuscripts (and More)." In *1st International Workshop on Computational Paleography (IWCP@ICDAR 2021)*. Lausanne, Switzerland, 2021. <https://hal.archives-ouvertes.fr/hal-03336528>.
- Khemakhem, Mohamed. 2020. "Standard-Based Lexical Models for Automatically Structured Dictionaries." Thèse, Université Paris Cité.
- Massot, Marie-Laure, Arianna Sforzini, and Vincent Ventresque. "Transcribing Foucault's Handwriting with Transkribus." *Journal of Data Mining and Digital Humanities* Atelier Digit\_Hum (March 2019). <https://doi.org/10.46298/jdmdh.5043>.
- Mühlberger, Guenter, Louise Seaward, Melissa Terras, Oliveira Sofia Ares, Vicente Bosch, Maximilian Bryan, Sebastian Colutto, et al. "Transforming Scholarship in the Archives through Handwritten Text Recognition: Transkribus as a Case Study."

*Journal of Documentation* 75, no. 5 (2019): 954–76.

<https://doi.org/10.1108/JD-07-2018-0114>.

Nguyen, Thi-Tuyet-Hai, Adam Jatowt, Mickaël Coustaty, Nhu-Van Nguyen, and Antoine Doucet. “Deep Statistical Analysis of OCR Errors for Effective Post-OCR Processing.” In *2019 ACM/IEEE Joint Conference on Digital Libraries (JCDL)*, 29–38. Champaign, France: IEEE, 2019. <https://doi.org/10.1109/jcdl.2019.00015>.

Oliveira, Sofia Ares, Benoit Seguin, and Frederic Kaplan. 2018. “DhSegment: A Generic Deep-Learning Approach for Document Segmentation.” Pp. 7–12 in *2018 16th International Conference on Frontiers in Handwriting Recognition (ICFHR)*.

Pletschacher, Stefan, and Apostolos Antonacopoulos. “The PAGE (Page Analysis and Ground-Truth Elements) Format Framework.” In *2010 20th International Conference on Pattern Recognition*, 257–60, 2010. <https://doi.org/10.1109/ICPR.2010.72>.

Rigaud, Christophe, Antoine Doucet, Mickaël Coustaty, and Jean-Philippe Moreux. “ICDAR 2019 Competition on Post-OCR Text Correction.” In *15th International Conference on Document Analysis and Recognition*, 1588–93. Sydney, Australia, 2019. <https://hal.archives-ouvertes.fr/hal-02304334>.

Stokes, Peter A., Benjamin Kiessling, Daniel Stökl Ben Ezra, Robin Tissot, and El Hassane Gargem. “The EScriptorium VRE for Manuscript Cultures.” *Classics@Journal*. 2021. <https://classics-at.chs.harvard.edu/classics18-stokes-kiessling-stokl-ben-ezra-tissot-gargem/>.

## Biographies

**Hugo Scheithauer**, Research & Development Engineer, ALMAnaCH (Inria), Paris: [hugo.scheithauer@inria.fr](mailto:hugo.scheithauer@inria.fr).

Research and Development Engineer in the Inria ALMAnaCH team, **Hugo Scheithauer** holds a master's degree in art history and in "Technologies numériques appliquées à l'histoire" at the École nationale des chartes. He works on the automatic segmentation of sale catalogs for the DataCatalogue project, jointly led by Inria, the National Library of France (BnF) and the National Institute for Art History (INHA).

**Alix Chagué**, PhD Candidate in Digital Humanities, ALMAAnaCH (Inria), Paris and CRIHN (Université de Montréal), Montreal: [alix.chague@inria.fr](mailto:alix.chague@inria.fr).

**Alix Chagué** is a PhD candidate in Digital Humanities, focusing on the resources and methods to apply HTR technologies in the various fields of the Humanities. As a co-founder of the HTR-United initiative (<https://htr-united.github.io>), she defends the construction by the community of homogenous practices within the Open Science framework. She is also a coordinator for the Master Degree in Documentation and Digital Humanities at the Ecole du Louvre in Paris.

**Laurent Romary**, Director for Scientific Information and Culture, Inria, Paris: [laurent.romary@inria.fr](mailto:laurent.romary@inria.fr)

Laurent Romary is a Research Director and the Director for Scientific Information and Culture at Inria, Paris. He conducted multiple research projects in computer science, linguistics, and digital humanities. He developed a strong interest in standardization and open data. He also identified the need to implement shared infrastructures for open science.