

Between automatic and manual encoding

Towards a generic TEI model for historical prints and manuscripts

Ariane Pinche ¹, Kelly Christensen ², Simon Gabay ³

¹École nationale des chartes — PSL ²Inria

³Université de Genève

15 September 2022

Between
automatic and
manual encoding

A. Pinche et al.

The *Gallic(orpor)a*
project

Extracting
information

TEI document
modelling

General principles

<teiHeader>

<sourceDoc>

<body>

Conclusion

The Gallic(orpor)a project

• Who and When

- Born from a collaboration between the INRIA, the École nationale des chartes and Geneva University
- Started at the end of 2021
- Funded by the BnF dataLab

• Objectives

- Valorisation of digitized collection from Gallica
- Production of automatically generated corpus from the 15th to the 18th century
- Ensure data compatibility with a wide range of use cases: between strict structured data and flexibility

Data and scripts, created during the project, are available online:

<https://github.com/Gallicorpora>

The *Gallic(orpor)a*
project

Extracting
information

TEI document
modelling

General principles

<teiHeader>

<sourceDoc>

<body>

Conclusion

Extracting information from digital facsimiles

Between
automatic and
manual encoding

A. Pinche et al.

The *Gallic(orpor)a*
project

Extracting
information

TEI document
modelling

General principles

<teiHeader>

<sourceDoc>

<body>

Conclusion

Designing common practices for documents annotation

Between
automatic and
manual encoding

A. Pinche et al.

Design of common practices to enable pipeline automation and data consistency

① Harmonizing layout description

- Using the same codicological vocabulary for description;
- Managing the differences between different layout;
- Results: models for page segmentation

② Harmonizing Transcription

- Using the same transcription rules and the same unicode characters;
- Managing the differences between manuscript vs printed characters, gothic vs antiqua;
- Results : two models, one for Manuscripts and incunabula : Cortado and another for incunabula and prints : **Gallicorpora+**

Models are available here :

<https://github.com/Gallicorpora/Segmentation-and-HTR-Models>

The *Gallicorpora*
project

Extracting
information

TEI document
modelling

General principles

<teiHeader>

<sourceDoc>

<body>

Conclusion

Designing common practices: Segmentation

The *SegmOnto* project consists of a generalist controlled vocabulary and a basic syntax to describe the layout of historical sources.

<i>Zones</i>	<i>Lines</i>
DropCapitalZone	DefaultLine
GraphicZone	DropCapitalLine
MainZone	HeadingLine
MarginTextZone	InterlinearLine
MusicZone	MusicLine
NumberingZone	
QuireMarksZone	
RunningTitleZone	
TableZone	
TitlePageZone	

Example: Zone:subZone e.g. MarginTextZone:footnote

Between
automatic and
manual encoding

A. Pinche et al.

The *Gallic(orpor)a*
project

Extracting
information

TEI document
modelling

General principles

<teiHeader>

<sourceDoc>

<body>

Conclusion

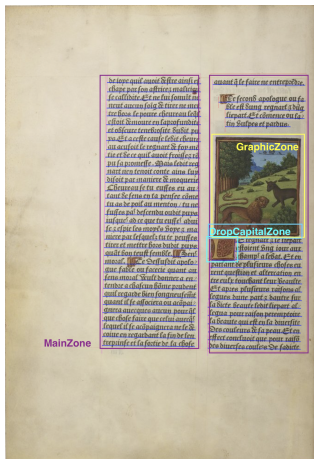


Figure: BnF, Réserve des livres rares, vélin, 611, 15^e s.

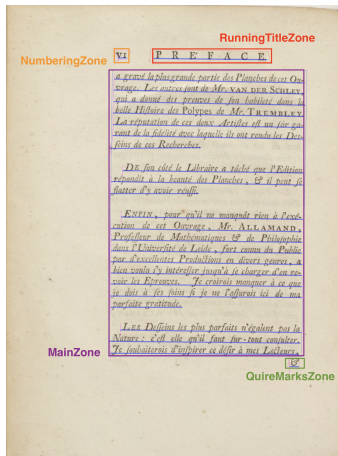


Figure: BnF, Arsenal, 4-S-1534, 18^e s.

Between automatic and manual encoding

A. Pinche et al.

The Gallic(or)por) project

Extracting information

TEI document modelling

General principles

<teiHeader>

<sourceDoc>

<body>

Conclusion

Segmentation and first failures

Between
automatic and
manual encoding

A. Pinche et al.

Training a segmentation model with *Kraken* was not successful. Is the problem due to the heterogeneity of the layouts? The granularity of the description? Or a methodological problem?

The *Gallic(or)por*
project

Extracting
information

TEI document
modelling

General principles
<teiHeader>
<sourceDoc>
<body>

Conclusion



Figure: BnF,
Réserve des livres rares,
vélin 611, 15^e s.

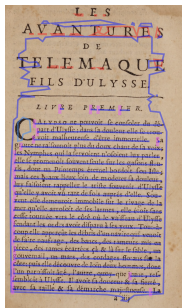


Figure: BnF,
Réserve des livres rares,
RES-Z-2442, 16^e s.

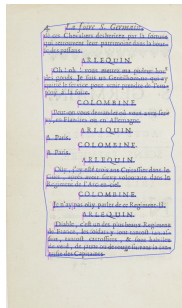


Figure: BnF, Arts du
spectacle, Réserve
8-10-1702, 17^e s.

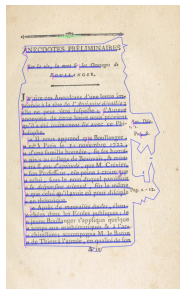


Figure: BnF,
département Droit,
économie, politique,
2012-39571, 18^e s.

Since segmentation is crucial to our pipeline, we needed a solution:

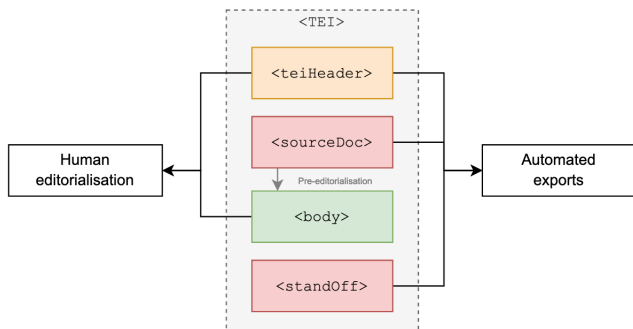
- A new approach : **YALTAi**, see T. Clerice, *You Actually Look Twice At it (YALTAi): Using an Object Detection Approach Instead of Region Segmentation Within the Kraken engine*
 - Different from Kraken pixel categorization
 - Using object detection, YOLOv5 solution
- New scores (average precision in percent) :

Zone	Main	Graphic	DropCapital	MarginText	Numbering	RunningTitle
Kraken	43.5	16.1	23.3	0	0	0
Yolo V5	91.7	48.4	69.2	48.3	75.8	45.6

TEI Document Modelling

Main structure of the TEI model

- Aim : to produce a generic model able to receive a huge variety of documents
- TEI file generated by a python script written by K. Christensen
- Script and documentation available here : <https://github.com/kat-kel/alto2tei>



Between
automatic and
manual encoding

A. Pinche et al.

The *Gallic(orpor)a*
project

Extracting
information

TEI document
modelling

General principles

`<teiHeader>`

`<sourceDoc>`

`<body>`

Conclusion

Automatic generation of the XML Tree

Between
automatic and
manual encoding

A. Pinche et al.

The *Gallic(orpor)a*
project

Extracting
information

TEI document
modelling

General principles

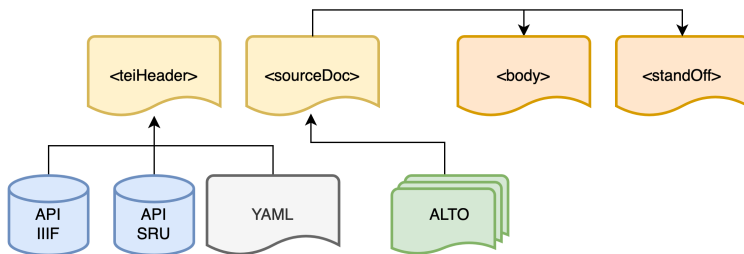
<teiHeader>

<sourceDoc>

<body>

Conclusion

Generation of the XML Tree from external information



<teiHeader>

<teiHeader>

Between
automatic and
manual encoding

A. Pinche et al.

The *Gallica(orpor)a*
project

Extracting
information

TEI document
modelling

General principles

<teiHeader>

<sourceDoc>

<body>

Conclusion

<teiHeader>

```
<teiHeader>
  <fileDesc>
    <titleStmt>
      <title> <!-- ... --> </title>
      <author> <!-- ... --> </author>
      <respStmt> <!-- team --> </respStmt>
    </titleStmt>
    <extent><!-- count of pages --></extent>
    <publicationStmt> <!-- project -->
    </publicationStmt>
    <sourceDesc> <!-- source--> </sourceDesc>
  </fileDesc>
  <profileDesc> <!-- language --> </profileDesc>
  <encodingDesc> <!-- models, HTR engine and
  Ontology -->
  </encodingDesc>
</teiHeader>
```

Between
automatic and
manual encoding

A. Pinche et al.

The *Gallic(or)pora*
project

Extracting
information

TEI document
modelling

General principles

<teiHeader>

<sourceDoc>

<body>

Conclusion

- APIs (IIIF API and SRU API)
 - **Title** - Le Romant comique [1re partie]
 - **Responsibility** - Paul Scarron
 - **Publication** - 1655, Leide, J. Sambix (ed.)
 - **Description of the object** - français, texte imprimé
 - **Conservation** - Bibliothèque nationale de France (8-Y2-55998)
- Config file (YAML)
 - **Responsibility** : Kelly Christensen
 - **Publication** : BnF DataLab
 - **Object Description** : 20 pages

Example of teiHeader automatically generated

```
1 <?xml version="1.0" encoding="UTF-8"?> 43
2 <TEI xmlns="http://www.tei-c.org/ns/1.0" xml:id="ark:12148_bpt6k6424218b"> 44
3 <teiHeader> 45
4 <fileDesc> 46
5 <titleStmt> 47
6 <title>Le Romant comique [Ire partie], par Mr Scarron</title> 48
7 <author xml:id="Sci1"> 49
8 <persName> 50
9 <forename>Paul</forename> 51
10 <surname>Scarron</surname> 52
11 <ptr type="isni" target="0000000120990126"/> 53
12 </persName> 54
13 </author> 55
14 <respStmt> 56
15 <resp>Transformation from ALTO4 to TEI by</resp> 57
16 <persName> 58
17 <forename>Kelly</forename> 59
18 <surname>Christensen</surname> 60
19 <ptr type="orcid" target="000000027236874X"/> 61
20 </persName> 62
21 <persName> 63
22 <forename>Simon</forename> 64
23 <surname>Gaboy</surname> 65
24 <ptr type="orcid" target="0000000190944475"/> 66
25 </persName> 67
26 <persName> 68
27 <forename>Arlane</forename> 69
28 <surname>Pinche</surname> 70
29 <ptr type="orcid" target="0000000278435050"/> 71
30 </persName> 72
31 </respStmt> 73
32 </titleStmt> 74
33 <extent> 75
34 <measure unit="images" n="20"/> 76
35 </extent> 77
36 <publicationStmt> 78
37 <publisher>Gallic(orpor)</publisher> 79
38 <authority>BnF DATALab</authority> 80
39 <availability status="restricted" n="cc-by"> 81
40 <license target="https://creativecommons.org/licenses/by/4.0/"> 82
41 </availability> 83
42 <date when="2022-06-10"/> 84
43 </publicationStmt> 85
44 <sourceDesc> 86
45 <bibl> 87
46 </bibl> 88
47 </sourceDesc> 89
48 </publicationStmt> 89
49 </TEI> 90
```

The *Gallic(orpor)*
project

Extracting
information

TEI document
modelling

General principles

<teiHeader>

<sourceDoc>

<body>

Conclusion

<sourceDoc>

Between
automatic and
manual encoding

A. Pinche et al.

The *Gallic(orpor)a*
project

Extracting
information

TEI document
modelling

General principles

<teiHeader>

<sourceDoc>

<body>

Conclusion

<sourceDoc>

From XML Alto to the <sourceDoc>



Between automatic and manual encoding

A. Pinche et al.

The Gallic(orpor)a project

Extracting information

TEI document modelling

General principles

<teiHeader>

<sourceDoc>

<body>

Conclusion

ALTO → <sourceDoc>

ALTO

```
<TextLine ID="line_3" TAGREFS="LT832"  
  BASELINE="277 985 734 990" HPOS...>  
<Shape>  
  <Polygon POINTS="277 985 275 940..." /> </Shape>  
<String CONTENT="CHAPITRE I." HPOS="275"  
  VPOS="929" WIDTH="460" HEIGHT="70" ></String>...
```

TEI

```
<zone xml:id="f15_z1_l1" type="HeadingLine"  
  subtype="none" n="1"  
  points="277,985 275,940..."  
  source="https://f15/275,929,460,70...jpg">  
<path xml:id="f15_z1_l1_p"  
  points="277,985 734,990"/>  
<line xml:id="f15_z1_l1_t">CHAPITRE I.</line> ...
```

Between
automatic and
manual encoding

A. Pinche et al.

The *Gallic(or)pora*
project

Extracting
information

TEI document
modelling

General principles

<teiHeader>

<sourceDoc>

<body>

Conclusion

From ALTO <TextLine> to TEI <zone>

Between
automatic and
manual encoding

A. Pinche et al.

The *Gallic(or)pora*
project

Extracting
information

TEI document
modelling

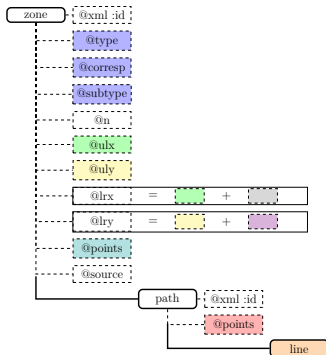
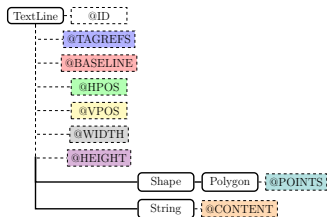
General principles

<teiHeader>

<sourceDoc>

<body>

Conclusion



<sourceDoc>

```
<sourceDoc>
  <surface><!--Page-->
    <zone type="SegmOntoZone"><!--Text Block-->
      <zone type="SegmOntoLine"><!-- TextLine
        (ex. "Text here.") -->
        <zone type="String"><!-- Segment
          (ex. "Text") -->
          <zone type="Glyph">
            <!-- Character (ex. "T") -->
            <c>T</c>
          </zone>
          <zone type="Space"/>
        <line>Texte here.</line>
      </zone>
    </zone>
  </surface>
</sourceDoc>
```

Between
automatic and
manual encoding

A. Pinche et al.

The *Gallic(orpor)a*
project

Extracting
information

TEI document
modelling

General principles

<teiHeader>

<sourceDoc>

<body>

Conclusion

<body>

Between
automatic and
manual encoding

A. Pinche et al.

The *Gallic(orpor)a*
project

Extracting
information

TEI document
modelling

General principles

<teiHeader>

<sourceDoc>

<body>

Conclusion

<body>

<body> and pre-editorialised text

Between
automatic and
manual encoding

A. Pinche et al.

The *Gallic(orpor)a*
project

Extracting
information

TEI document
modelling

General principles

<teiHeader>

<sourceDoc>

<body>

Conclusion

<body> and pre-editorialized text

- Created from <sourceDoc> content
- Uses SegmOnto and layout analysis to organise the information
- The link between the image and the text is preserved
- Gives a first draft of an edition
- Made to be customized according to the problematic of each project
- Limits : As the text is produced automatically, there are always errors.

From XML ALTO to the pre-editorialisation of the <body>

Between
automatic and
manual encoding

A. Pinche et al.

The *Gallic(or)pa*
project

Extracting
information

TEI document
modelling

General principles

<teiHeader>

<sourceDoc>

<body>

Conclusion

Segmentation standardized description allows a basic structuring of the <body> element.

SegmOnto	TEI
NumberingZone	<fw type="NumberingZone" >
QuireMarksZone	<fw type="QuireMarksZone" >
RunningTitleZone	<fw type="RunningTitleZone" >
MarginTextZone	<note type="MarginTextZone" >
MainZone	<ab type="MainZone" >
DefaultLine	<lb>
HeadingLine	<hi rend="HeadingLine" >
DropCapitalLine	<hi rend="DropCapitalLine" >

Mapping <sourceDoc> → <body>

<sourceDoc>

```
<zone xml:id="f15_z1_l1" type="HeadingLine"
      subtype="none" n="1"
      points="277,985 275,940...>
  <path xml:id="f15_z1_l1_p"
        points="277,985 734,990"/>
  <line xml:id="f15_z1_l1_t">CHAPITRE I.</line>
</zone>
```

<body>

```
<pb corresp="#f15"/>
<ab corresp="#f15_z1" type="MainZone">
  <hi rend="HeadingLine">
    <lb corresp="#f15_z1_l1"/>CHAPITRE I.
  </hi>
```

...

Between
automatic and
manual encoding

A. Pinche et al.

The *Gallic(or)pora*
project

Extracting
information

TEI document
modelling

General principles

<teiHeader>

<sourceDoc>

<body>

Conclusion

Example of body encoding

```
<pb corresp="#f23"/>
<fw corresp="#f23-eSc_textblock_a44af781-blockCount1" type="NumberingZone"><lb
corresp="#f23-eSc_textblock_a44af781-eSc_line_543c5953-lineCount1"/>10</fw>
<ab corresp="#f23-eSc_textblock_cb6c5a03-blockCount2" type="MainZone"><hi rend="HeadingLine"
><lb corresp="#f23-eSc_textblock_cb6c5a03-eSc_line_66797451-lineCount2"
/>BRADAMANTE,<lb corresp="#f23-eSc_textblock_cb6c5a03-eSc_line_637e1bb9-lineCount3"
/>TRAGECOMEDIE.<lb corresp="#f23-eSc_textblock_cb6c5a03-eSc_line_fb9ce0c7-lineCount4"
/>ACTE I. SCENE I.<lb corresp="#f23-eSc_textblock_cb6c5a03-eSc_line_d87052a6-lineCount5"
/>Charlemagne.</hi><hi rend="DropCapitalLine"><lb
corresp="#f23-eSc_textblock_cb6c5a03-eSc_line_e55d1a40-lineCount6"/>L</hi><lb
corresp="#f23-eSc_textblock_cb6c5a03-eSc_line_5598aaae-lineCount7"/>Es ceptres des
grands Rois vien-<lb corresp="#f23-eSc_textblock_cb6c5a03-eSc_line_1e7bb869-lineCount8"
/>nent du Dieu suprême,<lb
corresp="#f23-eSc_textblock_cb6c5a03-eSc_line_26556dbe-lineCount9"/>C'eft luy qui ceint
nos chefs d'vn<lb corresp="#f23-eSc_textblock_cb6c5a03-eSc_line_eda21fb8-lineCount10"
/>royal diadème,<lb corresp="#f23-eSc_textblock_cb6c5a03-eSc_line_d59a4e1f-lineCount11"
/>Qui nous fait quand il veut re-<lb
corresp="#f23-eSc_textblock_cb6c5a03-eSc_line_db058f24-lineCount12"/>gner fur
l'Vniuers,<lb corresp="#f23-eSc_textblock_cb6c5a03-eSc_line_6d2c8709-lineCount13"/>Et
quand il veut fait cheoir no-<lb
corresp="#f23-eSc_textblock_cb6c5a03-eSc_line_d5bf7dce-lineCount14"/>ftré empire à
l'euers.<lb corresp="#f23-eSc_textblock_cb6c5a03-eSc_line_a94ed2a3-lineCount15"/>Tout
depend de fa main, tout de fa main procede,<lb
corresp="#f23-eSc_textblock_cb6c5a03-eSc_line_2e0da9de-lineCount16"/>Nous n'auons rien
de nous, c'eft luy qui tout poffede,<lb
corresp="#f23-eSc_textblock_cb6c5a03-eSc_line_ea8ddb95-lineCount17"/>Monarque vniuerfel,
&amp; fes commandemens<lb
corresp="#f23-eSc_textblock_cb6c5a03-eSc_line_b7482d3c-lineCount18"/>Font les spherés
mouuoir &amp; tous les elemens.<lb
corresp="#f23-eSc_textblock_cb6c5a03-eSc_line_587530f5-lineCount19"/>Il a mis fur mon
chef la Françoisé couronne,<lb
corresp="#f23-eSc_textblock_cb6c5a03-eSc_line_8fbc0e66-lineCount20"/>Il a fait que ma
voix toute la terre eftonne,<lb
corresp="#f23-eSc_textblock_cb6c5a03-eSc_line_88955c54-lineCount21"/>Et que l'Aigle
Romain perche en mes eftendars,<lb
corresp="#f23-eSc_textblock_cb6c5a03-eSc_line_dd0cbcd2-lineCount22"/>Guide des efcadrons
de mes vaillans foudars.<lb
corresp="#f23-eSc_textblock_cb6c5a03-eSc_line_9e02c03d-lineCount23"/>L'Itale m'obeit, la
fuperbe Alemagne,<lb corresp="#f23-eSc_textblock_cb6c5a03-eSc_line_242d0072-lineCount24"
/>Et les Rois reculez de l'ondeuf Bretagne.</ab>
```

Between
automatic and
manual encoding

A. Pinche et al.

The *Gallic(or)por*a
project

Extracting
information

TEI document
modelling

General principles

<teiHeader>

<sourceDoc>

<body>

Conclusion

Conclusion

- The experiment is conclusive: we can automatically produce TEI XML files from the automatic acquisition of text and metadata.
- Further developments
 - To add linguistic information in our pipeline and optimise its usefulness by linking the body and the <standOff> element.
 - To provide a simple interface for viewing and annotating through `teiPublisher`