

CIHAM (UMR 5648), CNRS

Automatic Text Recognition and Historical Documents

Ariane Pinche

ariane.pinche@cnrs.fr

EnExDi, Lyon, 3 juin 2025

1 Automatic Text Recognition

- 1.1 Definition
- 1.2 The Steps of ATR
- 1.3 How to train an ATR model
- 1.4 Vocabulary
- 1.5 Evaluating a Model

2 Towards Training Generic Models

- 2.1 ATR and Specific Challenges of Historical Documents
- 2.2 Preparing Training Data
- 2.3 Sharing Your Data

What is ATR?



Figure: ATR Prediction

- Prediction of textual content
- from an image of the source by
- an artificial intelligence trained by a human
- through a process alternating
 - phases of human intervention
 - and computation phases

- **OCR:** Optical Character Recognition
- **HTR:** Handwritten Text Recognition
- **ATR:** Automatic Text Recognition

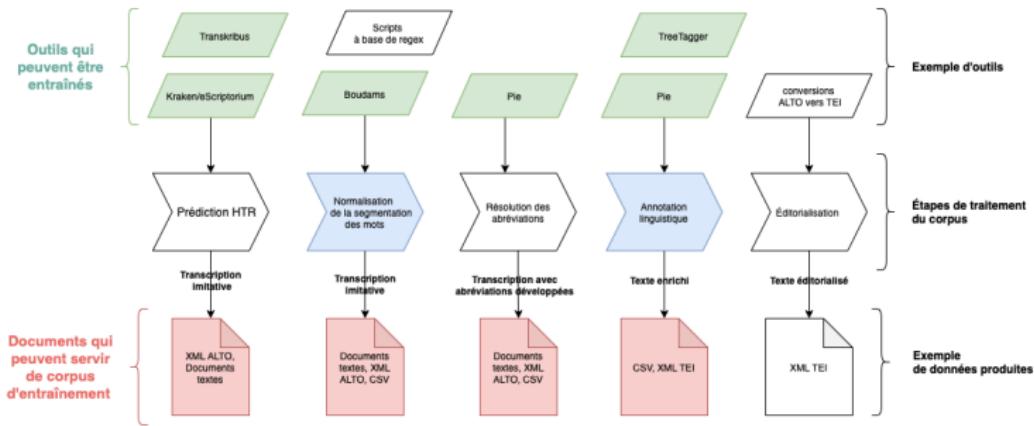
Some Tools

- eScriptorium [B. Kiessling, Tissot, Stokes, and Ezra 2019 and Benjamin Kiessling 2019]
- Transkribus [Kahle, Colutto, Hackl, and Mühlberger 2017]
- Calfa [Vidal-Gorène et al. 2021]

- 2000–2010: Development of HTR usage (experimental phase)
- Early research and pioneering projects:
 - Alex Graves and Jürgen Schmidhuber, “Offline Handwriting Recognition with Multidimensional Recurrent Neural Networks,” *Advances in Neural Information Processing Systems*, 2008.
 - Andreas Fischer, Emanuel Indermühle, Horst Bunke, [et al.], “Ground truth creation for handwriting recognition in historical documents,” 2010.
 - Pioneering projects: Himanis project (2015), ANR Horae project (2017), both led by Dominique Stutzmann.
- Late 2010s: Growing use of HTR in research projects, driven by the development of two main tools: Transkribus (using Pylaia and HTR+ engines) and eScriptorium (using the Kraken engine)
- Organization of international conferences:
 - ICDAR: International Conference on Document Analysis and Recognition
 - HIP: Historical Document Imaging and Processing Workshop

- ATR is a well-mastered task from a computer science point of view
 - Nowdays, with model that can reach a CER (character error rate) between 8% and 2% for manuscripts, “from a computer science point of view, the recognition of handwriting seems to be a resolved task. The latest recognition engines allow for the successful recognition of specifically trained hands producing a text as reusable data” , *Hodel, Tobias, David Schoch, Christa Schneider, and Jake Purcell. 2021. “General Models for Handwritten Text Recognition: Feasibility and State-of-the Art. German Kurrent as an Example.” Journal of Open Humanities Data 7(0):13. doi: 10.5334/johd.46.*
- L'ATR est devenu une étape courante dans les pipelines d'acquisition textuelle qui intéresse aussi bien les institutions patrimoniales que les projets de recherche.

Example of textual acquisition workflow

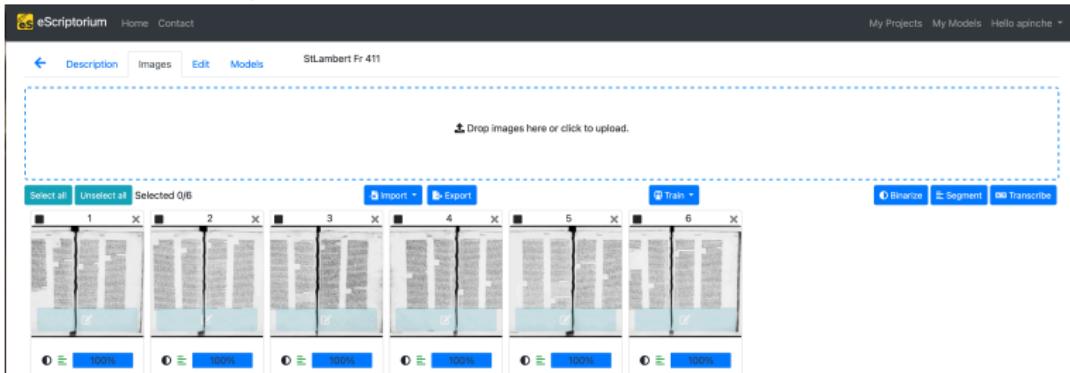


Why Use ATR Today?

- To accelerate the text acquisition phase. The prediction can serve:
 - as a basis for an edition: high level of accuracy, above 95%
 - for providing raw text: medium accuracy, between 90% and 95%
 - as a basis for quantitative analyses: lower accuracy, above 80%
(see EDER, Maciej, “Mind your corpus: systematic errors in authorship attribution,” *Literary and Linguistic Computing*, vol. 28/4, December 2013, pp. 603–614.)

The Steps of ATR

- Image loading
 - Load a collection of JPG or TIF images locally
 - Load from an IIIF manifest (e.g., collections from Gallica or e-Codices)
- Image preprocessing (optional)
 - 300 dpi resolution
 - Colour or grayscale
 - Optional linearisation to reduce noise
 - Multispectral imaging (in the case of heavily damaged documents)



- Segmentation of the zones



Figure: Bnf, fr. 412, fol.10r

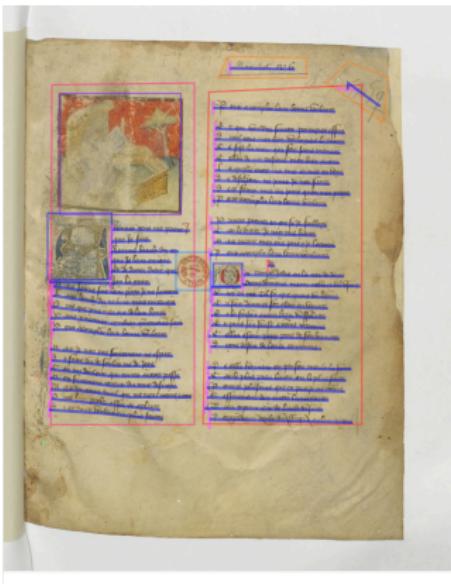
The Steps of ATR

- Segmentation of the Text Lines



Figure: Bnf, fr. 1728, fol.8v

- Text prediction



1 ucomes gens me prient ¶.
2 que le face
3 Aucuns beaux ditz et
4 que le leur enuoye
5 Et de ditz dient que
6 A
7 Mais sauve soit leur paix le ne sauroye
8 lay la gracie
9 Faire beaux ditz ne bons, mais touseuoye
10 Puis que pris men ont de leur bonte
11 Peine y mettray combien qu'ignouli soyte
12 Pour accomplir leur bonne oulemente
13 Mais le n ay pas sentement ne espere
14 De faire ditz, de soules ne de loye
15 Car ma douleur qui toutes autres passe
16 Mon sentement loyeux du tout desuoye
17 Mais du grant dueil qui me tient morne n'coye
18 Puis bien parler assez et aplante
19 Si en diray oulementiers plus feroye
20 6259
21 Pour accomplir leur bonne oulemente
22 Et qui voudra sauoir pourquoi efface
23 Dueil tout mon bles oulementiers le diroye
24 Ce fist la mort qui ferri sans mercie
25 Celli de qui trestout mon bien ausye
26 Laquelle mort ma mis et mat en uoye
27 De dessoپoir ne puis le nos sante
28 De ce feray mes ditz puis quon men proye
29 Pour accomplir leur bonne oulemente
30 Princes premes en gre se le fallroye
31 Car le bister le nay mie hante
32 Mais mainten ont prie ¶ le lotroye
33 Pour accomplir leur bonne oulemente
34 u temps ladiis en le cite de Rôme
35 .Ji.
36 O
37 ung en yst. Tel fu que quant un hôme
38 Orient Rômaine maint noble ¶ bel usage

Figure: Bnf, fr. 12779, fol.9r

The Steps of ATR

- Data export (txt, alto, page)

```
<Layout>
  <Page WIDTH="4648" HEIGHT="3407" PHYSICAL_IMG_NR="8" ID="eSc_dummypage_>
    <PrintSpace HPOS="0" VPOS="0" WIDTH="4648" HEIGHT="3407">

      <TextBlock HPOS="693" VPOS="321" WIDTH="1701" HEIGHT="2451"
                 ID="eSc_textblock_08b9f915" TAGREFS="BT3852">
        <Shape>
          <Polygon
            POINTS="693 413 693 2772 2394 2772 2254 321"/>
        </Shape>

        <TextLine ID="eSc_line_d939596f" TAGREFS="LT1299"
                  BASELINE="746 476 2143 428" HPOS="743" VPOS="352"
                  WIDTH="1400" HEIGHT="156">
          <Shape>
            <Polygon
              POINTS="2078 388 2050 388 2021 386 1993 383 1964 383 1936 380 1908 377 1876 374 1848 374 1820 371 1811
              />
          </Shape>
          <String
            CONTENT="fors de la ville. Tant fut l'assault merveilleux et"
            HPOS="743" VPOS="352" WIDTH="1400" HEIGHT="156"/>
        </TextLine>
    </PrintSpace>
  </Page>
</Layout>
```

Figure: Example of ALTO xml

- **Training corpus:** A dataset composed of images and their corresponding line-by-line transcriptions, used to teach the model how to recognise text.
- **Model:** A file generated at the end of the training process, containing all optimised parameters. It allows the Automatic text recognition tool to predict text from previously unseen images.
- **Prediction:** The process by which the model uses images to predict the text, based on the parameters learned during training.
- **Training:** A sequence of cycles during which the model adjusts its parameters to improve its ability to correctly recognise and predict text from images.

Training an ATR Model

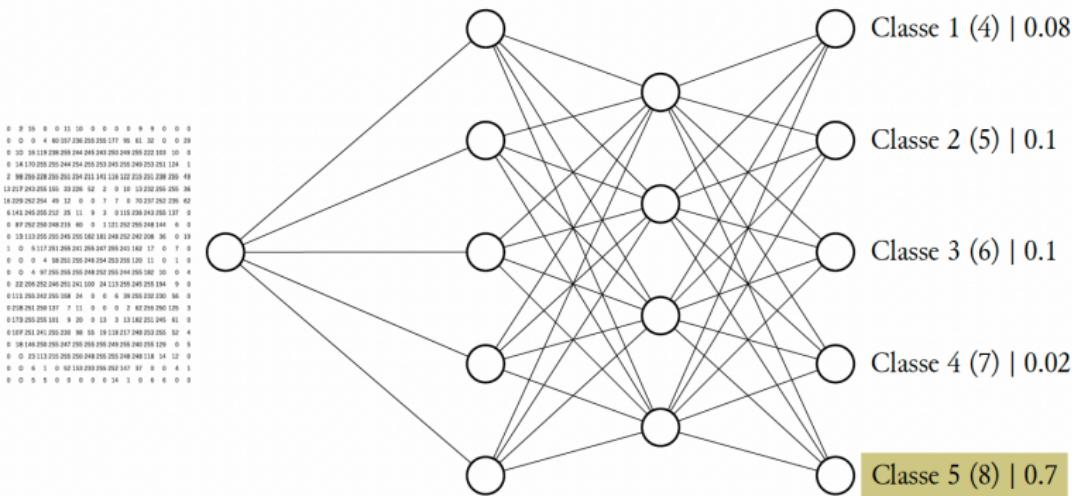


Figure: Simplified representation of a neural network

Training an ATR Model

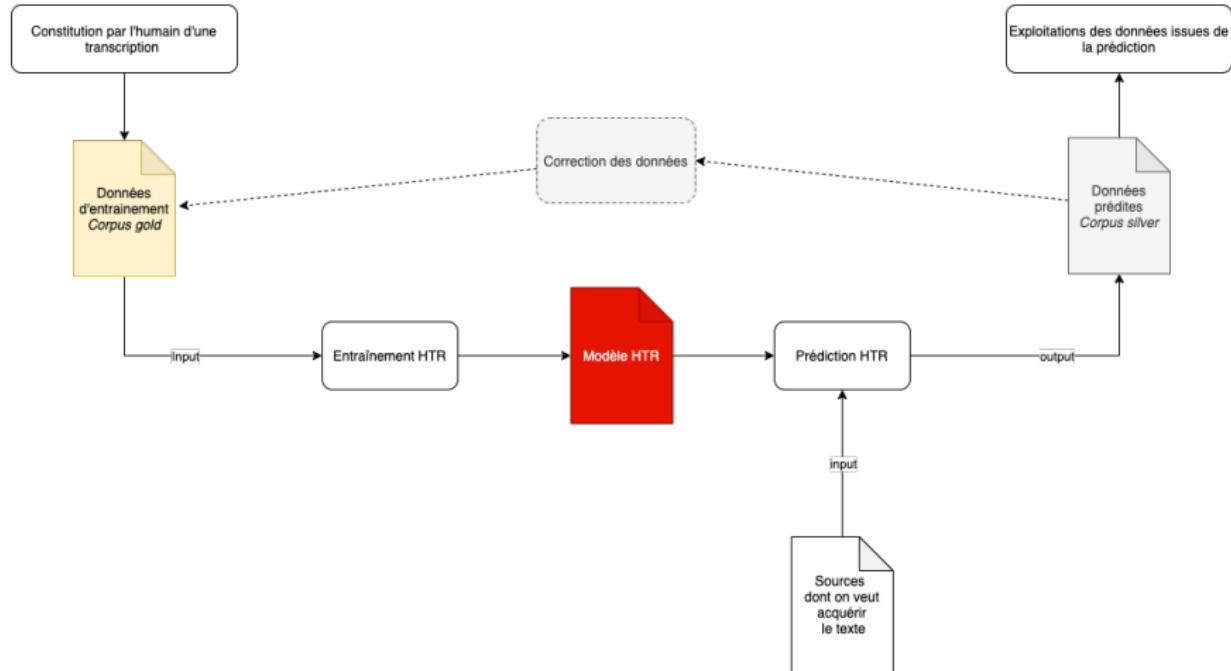


Figure: Representation of a training cycle

Training an ATR Model

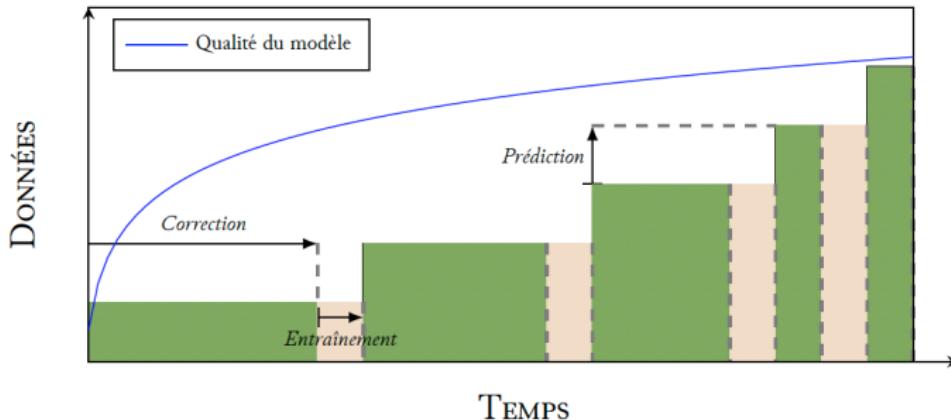
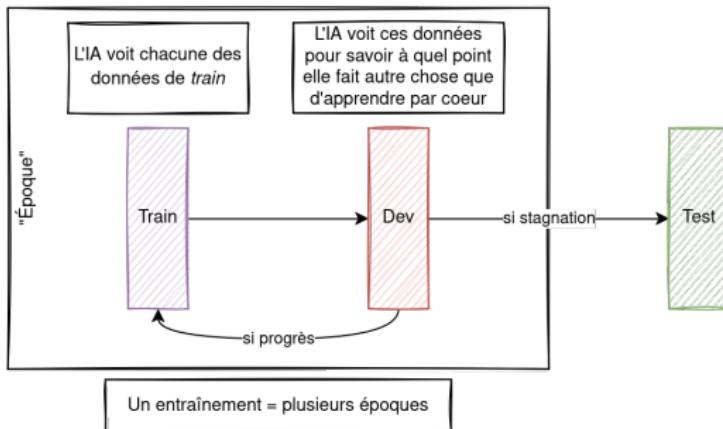
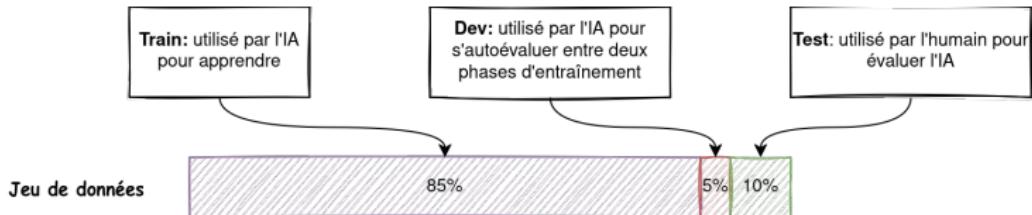
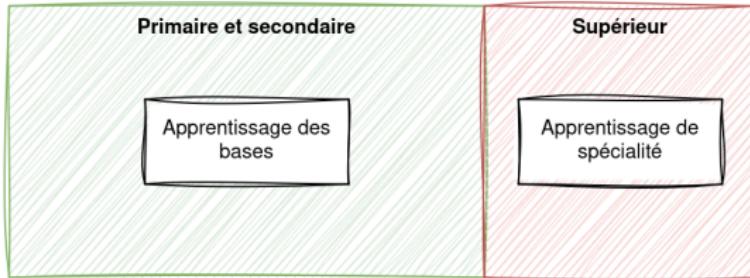


Figure: Evolution of data correction time depending on the model quality,
figure from the work of M. Gilles Levenson

Train? Dev? Test? The splits.



Fine-tuning?



Vous rappelez vous de tout ce que vous avez appris au lycée ?

En particulier ce qui n'a aucun rapport avec votre spécialité ?



Oubli d'une partie de (1) mais amélioration sur (2)

To evaluate an ATR model:

- prepare the dataset
 - train set (80%) – training
 - dev set (10%) – evaluation during training cycles
 - test set (10%) – data never seen during training
- compare:
 - a ground truth (GT) produced by a human (test set)
 - with the model's prediction of the same lines
 - to calculate a score which can be either:
 - a CER (Character Error Rate)
 - or an Accuracy (the percentage of correct predictions by the model)

STEAM

STEAM

STEAM

STEAL

TEAM

STREAM

 Substitution Deletion Insertion

$$CER = \frac{S + D + I}{N}$$

Where to Find Models?

- A performant model already exists. Where to find it? Where to find data to train your own?

Data HTR-united

Models Zenodo Community “OCR/HTR Models”

Table of Contents

1 Automatic Text Recognition

- 1.1 Definition
- 1.2 The Steps of ATR
- 1.3 How to train an ATR model
- 1.4 Vocabulary
- 1.5 Evaluating a Model

2 Towards Training Generic Models

- 2.1 ATR and Specific Challenges of Historical Documents
- 2.2 Preparing Training Data
- 2.3 Sharing Your Data

- Deciphering handwritten scripts in historical documents presents unique challenges:
 - Non-standardised layouts
 - Degraded supports
 - Irregular handwriting(s)
 - Graphic and/or dialectal variations

<i>Caroline</i> 8 th -13 th	eute ^x t ^{er} iticus qui iner ^{pre} tat ^x for ^x tunat ^x	potūt̄ hēri i bona ḡtitate f; q̄ tibi videb̄t̄	<i>Cursiva</i> 14 th -15 ^{tg}
<i>Praegothica</i> 12 th -13 th	eutex t euticus qui iner ^{pre} tat ^x fortunat ^x	potūt̄ hēri i bona ḡtitate f; q̄ tibi videb̄t̄	
<i>Gothica Textualis</i> 13 th -16 th	iora qb; ifort̄ ymago d̄ ppe sic pictā	domos .s. uii milia ho ^m quos 7 ipē	<i>Hybrida</i> 15 th -16 th
<i>Semitextualis</i> 13 th -16 th	e q̄re pauonē q̄re anserē gallina refugi	inandria sic aut inopia & cognatorum negligē	<i>Humanistic</i> 16 th
	C omo que creo no fuessēn	De Sperchio fluuio.xxiii.oceāi filio	<i>Incunabulum</i> 15 th
	C omo que creo no fuessen	De Sperchio fluuio.xxiii.oceāi filio	

Figure: Sample of medieval scripts

Non-unified Sources



Figure: BnF, Latin, 8001,
13th century



Figure: Strasbourg, ms.
1.916, 13th century

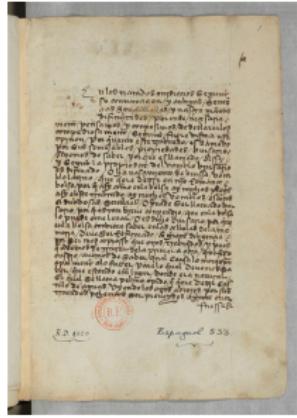


Figure: BnF, Spanish, 533,
15th century

How to Transcribe Documents?

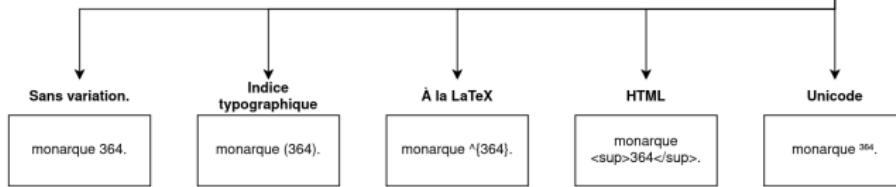
- Transcribing means describing a source
 - Transcribing means translating the source for non-palaeographers
 - Transcribing means interpreting the text
 - Transcribing means making choices
- “Well-prepared material is key to producing general recognition models. It is unthinkable that single scholars and small project teams could provide enough training material to train a general model independently.”
- Tobias Hodel, David Schoch, Christa Schneider, [et al.], “General Models for Handwritten Text Recognition: Feasibility and State-of-the-Art. German Kurrent as an Example,” 2021.

Making Choices and Documenting Them

Par ailleurs, comme le fait remarquer Francisco Tomás y Valiente, cette instance humaine de pouvoir, non sacré, favorise l'expression de critiques et de mécontentements, difficile à envisager lorsqu'il s'agit du monarque³⁶⁴.

Comment le transcrire ?

Villamediana est l'opposant le plus célèbre et le plus virulent aux favoris de Philippe III, Lerma et Uceda. Mais il convient de se demander si les satires politiques de Villamediana visent à abattre les personnes ou si, de manière plus idéologique, elles condamnent le système de gouvernement mixte. Le point de vue de Francisco Tomás y Valiente est tranché :



How to Transcribe Documents?

- How to transcribe documents for the machine?
- How to transcribe consistently within a project?
- How to transcribe so that my data can be reused?



jmfrajejas.bsky.social José Manuel Frajeda
@JMFraeRue

...

Examination of the output of the previous ms shows one of the problems with this model. Being a snowball model, it mixes transcription criteria. Lines 24 and 25 show that there're models that don't develop abbreviations (q); line 27 tells that some use the HSMS system (q<ue>) ->

[Traduire le post](#)

- 1-23 fechura no deue paran mjero
- 1-24 ala color \$da q qere\$a los
- 1-25 fralcons q soy cntrados o
- 1-26 f Fuara a manallos.
- 1-27 oq<ue> torna contra umneio prriua

12:46 PM · 30 nov. 2023 · 134 vues

How to Transcribe Documents?

- Define transcription methods adapted to the research question and machine learning.
- Define the desired level of precision in transcription
- Use a predefined character set and document your choices.
- Ensure compatibility of transcription data.



Figure: Turin Manuscript,
Ségurant the Knight and the
Dragon, 15th century



Figure: BnF, Arsenal, 3516,
12th century

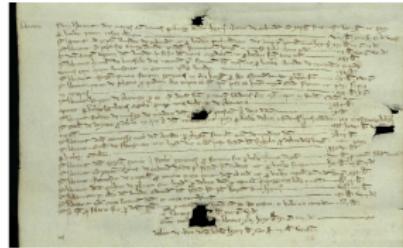


Figure: Departmental Archives of Côte
d'Or, B6739, 13th century

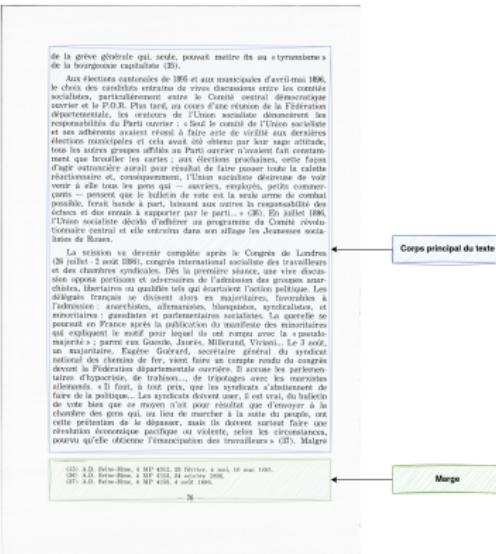
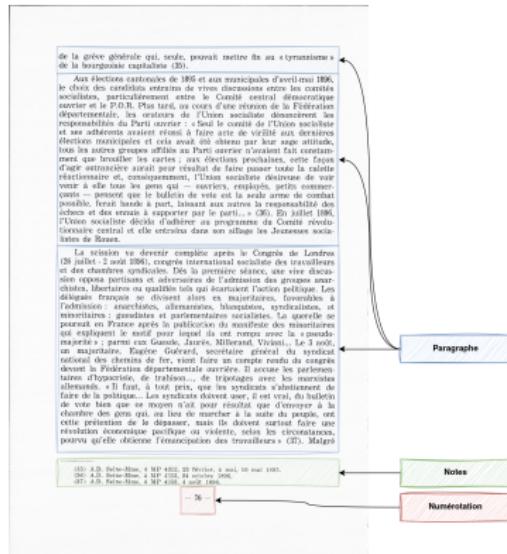
Layout Analysis: Segmentation

- Identification of the different zones in the document: use controlled vocabulary such as SegmOnto.

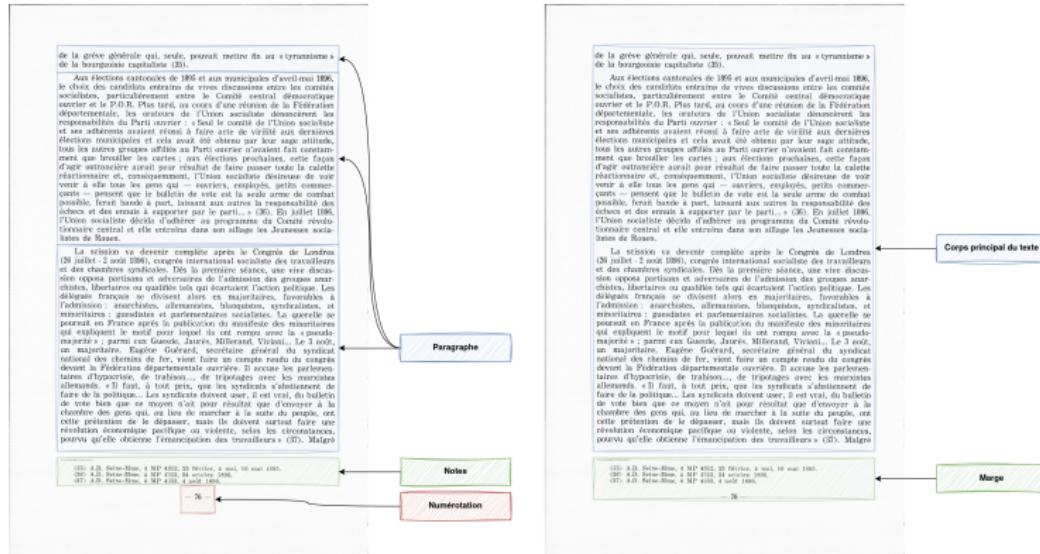


Figure: BnF, fr. 412, fol. 10r

Options?



Options ?



<https://segmonto.github.io>

Page

DamageZone
DropCapitalZone
FigureZone
MainZone
MarginZone
MusicNotationZone
NumberingZone
QuireMarksZone
RunningTitleZone

Line

DefaultLine
Interlinear
HeadingLine

To go further...

See extended dataset LADaS :

Thibault Clérice et al. *Layout Analysis Dataset with*

SegmOnto (LADaS). URL:

<https://github.com/DEFI-COLaF/LADaS>

Today, segmentation is the main source of errors in Handwritten Text Recognition (HTR):

- Incorrect identification of zones
- Incorrect labelling of zone
- Merging of nearby zones
- Lines not linked to the correct zone
- Since 2011. This task has been the focus of competitions held at ICDAR and HIP conferences.

Harmonising data allows the exchange of data and HTR models. How to achieve this?

- Define transcription methods adapted to your research questions and machine learning.
- Use a predefined character set and document your choices.
 - See the Medieval Unicode Font Initiative (MUFI)
 - See the transcription guidelines proposed by CREMMALab for medieval texts, and the CREMMA guidelines for modern transcriptions.
- Use controlled vocabulary to describe the layout and document it.

- Deposit your data on an accessible online repository:
 - Github
 - GitLab
- Document your data:
 - Data format
 - Number of transcribed lines
 - Segmentation tools
 - HTR engine
 - Corpus language
 - Date
 - Document and handwriting type
 - Transcription method
- Increase data visibility: integrate into a catalogue, see [HTR-united](#).