Horae project
oooooo

Pages selection process
oooo

Annotation results
oo

Document layout analysis
ooooo

# HORAE: an annotated dataset of books of hours

Mélodie Boillet, Marie-Laurence Bonhomme, Dominique
Stutzmann, Christopher Kermorvant

Teklia SAS, Paris, France
LITIS, Rouen-Normandie University, France
IRHT-CNRS, Paris, France
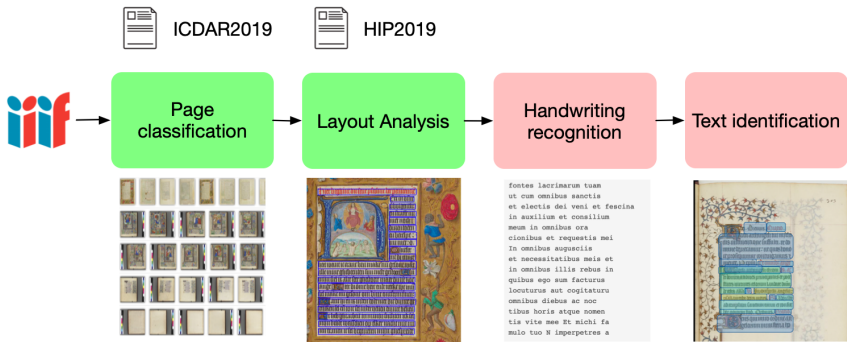
HIP 2019, 20th September 2019

## Horae project

- Book of hours, the medieval *best-seller*: more than 10,000 witnesses
- Personal prayer books, owned by rich laypersons
- Content:
    - perpetual calendar of the Church feasts
    - texts for each of the eight canonical hours (payer times) of the day
    - rich illustrations
- 300 pages, complex organization
- Surprisingly, no complete transcriptions of books of hours
- HORAE Project: automatic text recognition and structuration of book of hours
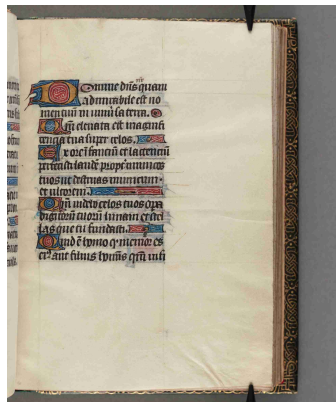
# Les Très Riches Heures du duc de Berry

# Project overview



ICDAR2019     HIP2019

Page classification → Layout Analysis → Handwriting recognition → Text identification

## Manuscripts collection

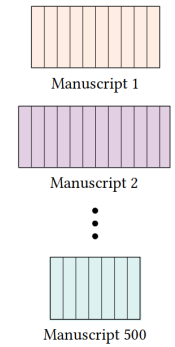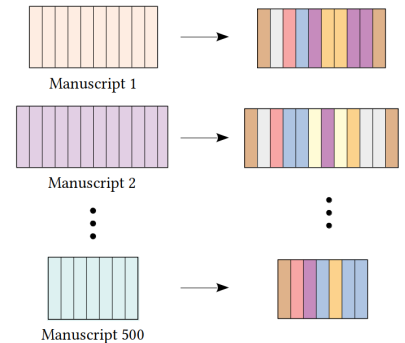| Provider | City | Manuscripts |
|---|---|---|
| UGent | Gent | 1 |
| BVMM | $\leq 10$ | 124 |
| | Angers | 21 |
| | Autun | 12 |
| | Beaune | 15 |
| | Chantilly | 30 |
| | Nantes | 18 |
| | Paris | 17 |
| | Rennes | 23 |
| | Toulouse | 15 |
| Gallica | Paris | 183 |
| Harvard | Cambridge | 32 |
| UBC | Vancouver | 1 |
| Stanford University | Stanford | 6 |
| WDL | Baltimore | 2 |
| **Total** | | **500** |

# Layout examples I

# Layout examples II

How to select the most representative set of pages ?

✗ Randomly : overrepresentation of the text pages and the large
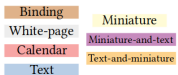  manuscripts;
✓ Selection process.

Horae project
oooooo

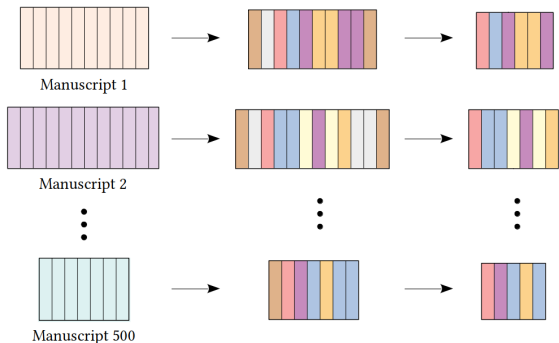Pages selection process
o●oo

Annotation results
oo

Document layout analysis
ooooo

# Selection process schema

Horae project
oooooo

Pages selection process
o●oo

Annotation results
oo

Document layout analysis
ooooo

# Selection process schema

Horae project
oooooo

Pages selection process
o●oo

Annotation results
oo

Document layout analysis
ooooo

# Selection process schema

Horae project
oooooo

**Pages selection process**
o●oo

Annotation results
oo

Document layout analysis
ooooo

# Selection process schema



| Classification | | Filtering | Clustering |
|---|---|---|---|
| Binding | Miniature | Binding | HDBSCAN |
| White-page | Miniature-and-text | White-page | • 141 clusters |
| Calendar | Text-and-miniature | • 2 max by class | • 2 200 outliers |
| Text | | | |

Horae project
○○○○○○

Pages selection process
○●○○

Annotation results
○○

Document layout analysis
○○○○○

# Selection process schema



| Classification | | | Filtering | Clustering HDBSCAN | Selection |
|---|---|---|---|---|---|

Classification
- Binding
- White-page
- Calendar
- Text
- Miniature
- Miniature-and-text
- Text-and-miniature

Filtering
- Binding White-page
- 2 max by class

Clustering
HDBSCAN
- 141 clusters
- 2 200 outliers

Selection
- 141 centroids
- 459 outliers

Horae project
000000

Pages selection process
OO●O

Annotation results
OO

Document layout analysis
00000

## Random selection



Mostly text pages

Horae project
000000

Pages selection process
000●

Annotation results
00

Document layout analysis
00000

Our selection



More illustrations

Horae project
oooooo

Pages selection process
oooo

Annotation results
●o

Document layout analysis
ooooo

# Distribution of the annotated elements using Transkribus

Horae project
oooooo

Pages selection process
oooo

Annotation results
o●

Document layout analysis
ooooo

# Annotation examples



illustrated border

historiated initial

miniature

simple initial

# Annotation examples



decorated border

text line

line filler

Horae project
000000

Pages selection process
0000

Annotation results
00

Document layout analysis
●0000

## How many documents to annotate ?

Line and region detection with dhSegment

| Training size | Task | IoU with post-processing |
|---|---|---|
| 220 | Line detection | 0.88 |
| | Layout analysis | 0.71 |

Horae project
000000

Pages selection process
0000

Annotation results
00

Document layout analysis
0●000

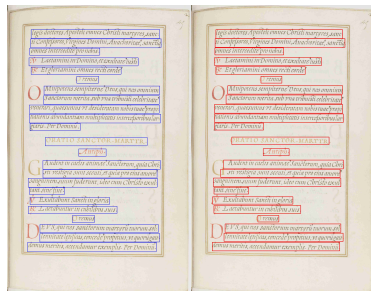## How many documents to annotate ?

Line and region detection with dhSegment

| Training size | Task | IoU with post-processing |
|---|---|---|
| 220 | Line detection | 0.88 |
| | Layout analysis | 0.71 |
| 510 | Line detection | 0.88 |
| | Layout analysis | 0.72 |

More data not needed with dhSegment model

# Visualization of the predictions I

## Visualization of the predictions II

Horae project
000000

Pages selection process
0000

Annotation results
00

Document layout analysis
0000●

## Conclusion and future work

- Introduction of a new dataset Horae including a large variety of types of pages;

- First reference results for line segmentation and layout analysis;

- Satisfactory results that can be improved using more complex neural networks.

- Classification for double-pages $\rightarrow$ only one class assigned;

- Ambiguity considering the initials $\rightarrow$ Inside or outside the text lines;

- Confusions between the initials;

- Problem with the post-processing step $\rightarrow$ Only rectangles are created for now.

# Freely available

https://github.com/oriflamms/HORAE

# Bibliography

▶ Dominique Stutzmann et al. "Integrated DH. Rationale of the HORAE Research Project". In: *Digital Humanities*. July 9, 2019. published.

▶ Emanuela Boros et al. "Automatic page classification in a large collection of manuscripts based on the International Image Interoperability Framework". In: *International Conference on Document Analysis and Recognition*. Sept. 1, 2019. published.

▶ Leland McInnes, John Healy, and Steve Astels. "HDBSCAN: Hierarchical density based clustering". In: *The Journal of Open Source Software* 2.11 (2017). DOI: 10.21105/joss.00205. URL: https://doi.org/10.21105%2Fjoss.00205.

▶ Sofia Ares Oliveira, Benoit Seguin, and Frederic Kaplan. "dhSegment: A generic deep-learning approach for document segmentation". In: *Frontiers in Handwriting Recognition (ICFHR), 2018 16th International Conference on*. IEEE. 2018, pp. 7–12.