# Importance of Textlines
# in Historical Document Classification

**Martin Kišš**, Jan Kohút, Karel Beneš, Michal Hradiš

BRNO FACULTY
UNIVERSITY OF INFORMATION
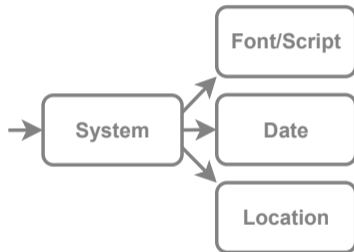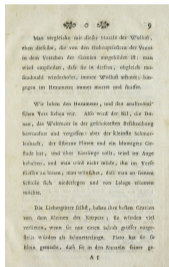OF TECHNOLOGY TECHNOLOGY

Page-level label(s) without any further details

Font & script classification
Up to $N$ possible labels per page
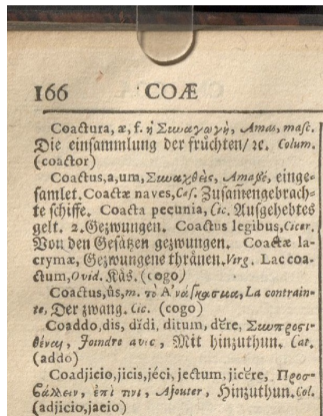
Localization
One label per page
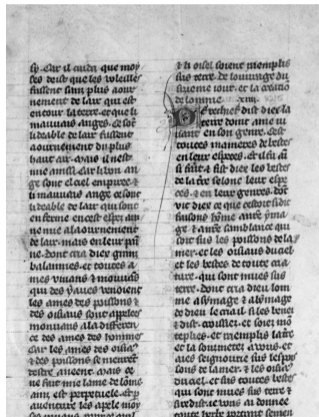
Dating
*(not-before; not-after)* intervals
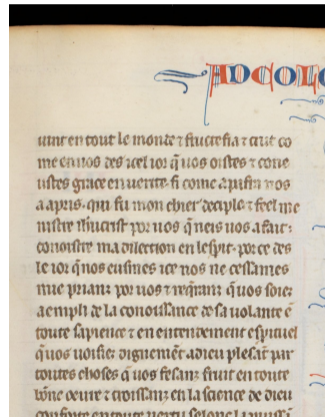
## Font/script

Fraktur, greek, italic, antiqua



## Localization
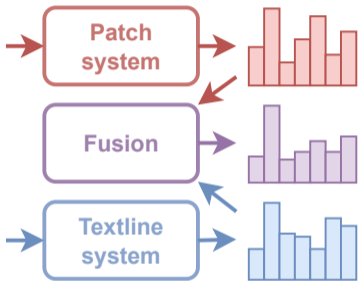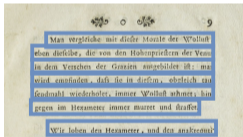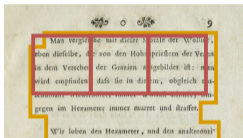
Paris



## Dating

1260-1269

- custom training-validation dataset splits
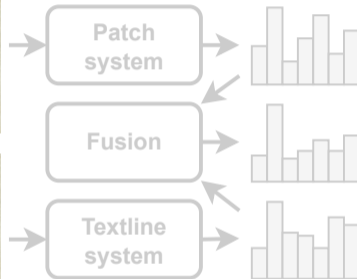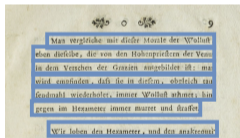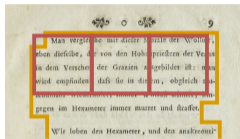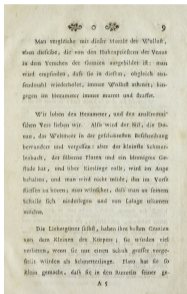- the aim to create class-balanced validation set
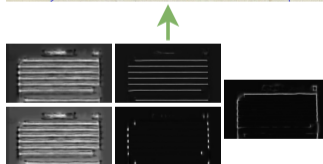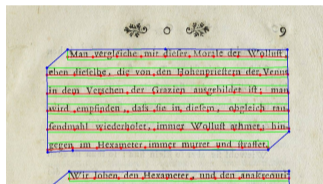- publicly available at `pero.fit.vutbr.cz/hdc_dataset`

| Task | Training | Validation | Testing |
|---|---|---|---|
| Font | 35 382 | 239 | 5 506 |
| Script | 7 594 | 419 | 1 256 |
| Location | 5 397 | 65 | 325 |
| Date | 10 294 | 1 000 | 2 516 |

# Our approach

# Our approach - Layout analysis



Layout analysis

Patch system

Fusion

Textline system

# Our approach - Layout analysis

- ParseNet (U-net-based) neural network[1]
- Text lines and regions detection
- Trained on PERO Layout dataset



**Layout analysis**

---

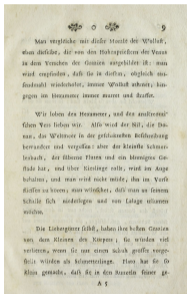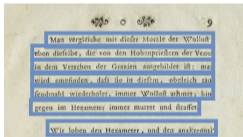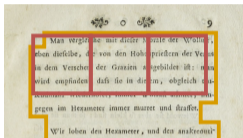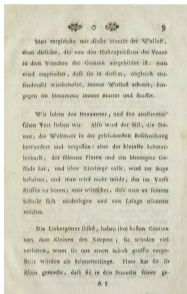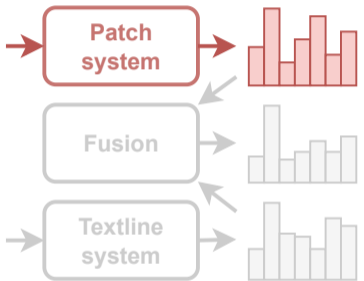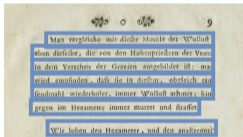[1]Kodym, Oldřich, and Michal Hradiš. "Page layout analysis system for unconstrained historic documents."

# Our approach - Patch system



Layout analysis → Patch system → Fusion → Textline system

# Our approach - Patch system

- ResNeXt-50-based neural network
- Training on four scales of non-overlapping square patches with three training strategies
- The page-level output is calculated as the mean of the local outputs



Page-level output

Patch system

Local outputs

# Our approach - Patch system

# Our approach - Textline system



Layout analysis → Patch system, Textline system → Fusion → Textline system

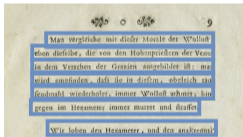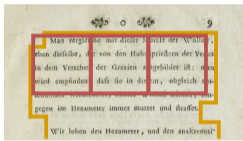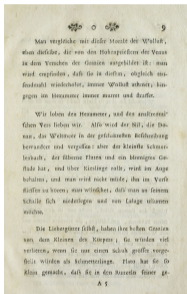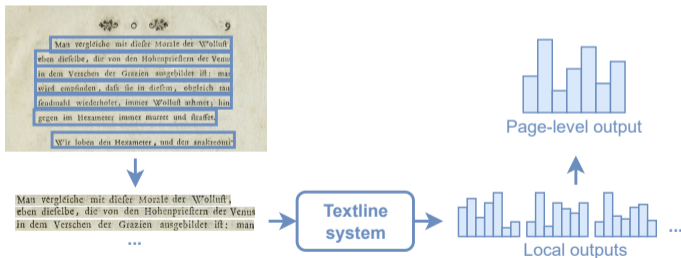# Our approach - Textline system

- VGG-based neural network + global average pooling
- Arbitrarily long height-normalized textline image as an input
- The page-level output is calculated as the mean (classification tasks) or median (dating task) of the local outputs

- Linear (left image) and log-linear (right image) fusions
- Page-level outputs from the patch and the textline systems as the input
- Optimized on the validation sets using 10-fold cross validation

# Our approach - Fusion

- Up to *N possible fonts/scripts* per page
  - e.g. $\mathcal{T} \in \{$*fraktur*, *greek*, *italic*, *antiqua*$\}$
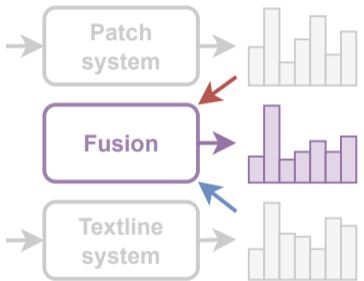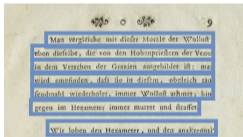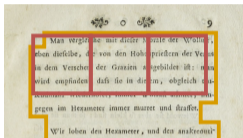- Page-level annotations without location information
- $L_{hard}$ selects the most probable page-level label based on network output
- $L_{soft}$ considers all page-level labels weighted by the network output probability

$$L_{hard} = \min_{i \in \mathcal{T}} \left[ -\log(f(x)_i) \right] \tag{1}$$

$$L_{soft} = \sum_{i \in \mathcal{T}} -\log(f(x)_i) \cdot f(x)_i. \tag{2}$$

- *(not-before; not-after)* interval instead of single point annotations
  - *e.g. (1260; 1269)*
- Modified Huber loss for intervals
- Pulls the output of the network more towards the middle of the interval

| | Font Acc. ↑ | Script Acc. ↑ | Location Acc. ↑ | Date MAE ↓ |
|---|---|---|---|---|
| Textline system | 98.42 % | 88.54 % | 69.85 % | **21.91** years |
| Patch system | 95.68 % | 80.26 % | 75.08 % | 32.45 years |
| Linear fusion | 98.27 % | **88.84** % | 70.77 % | 21.99 years |
| Log-linear fusion | **98.48** % | 88.60 % | **79.69** % | — |
| Baseline | — | 55.22 % | 62.46 % | — |
| The North LTU | 82.80 % | 74.12 % | 43.69 % | 79.43 years |
| CLUZH | 95.66 % | 35.25 % | — | — |
| NAVER Papago | 97.17 % | — | — | — |