# Open Data
## at the BNL

Bibliothèque nationale du Luxembourg

Nationalbibliothéik

@bnluxembourg

**Ralph Marschall**

Project Manager Digitisation

ONB LABS SYMPOSIUM 2021
24.11.2021

# Quick overview

- Digitization @ National Library of Luxembourg (BnL)

- Open Data & Tools

- A.I. / Machine Learning Project

# Digitization at National Library of Luxembourg (BnL)

# Digitization @ BnL

- **Preserve** old documents

  - **Restrict** access to physical volumes (scan once, minimize physical usage)

  - Some documents are **unique** (1 copy)

  - Replace microfilms with **improved digital access**

- **Search** Full Text – OCR (optical text recognition)

- Decompose images into **articles** (semantic units) – OLR (optical layout recognition)

- **Present** digital results online: eluxemburgensia.lu, a-z.lu

# Digitization timeline

**2003 to 2006**
- Image only
- Newspapers, postcards, books, manuscripts

**2006+**
- Luxembourg newspapers
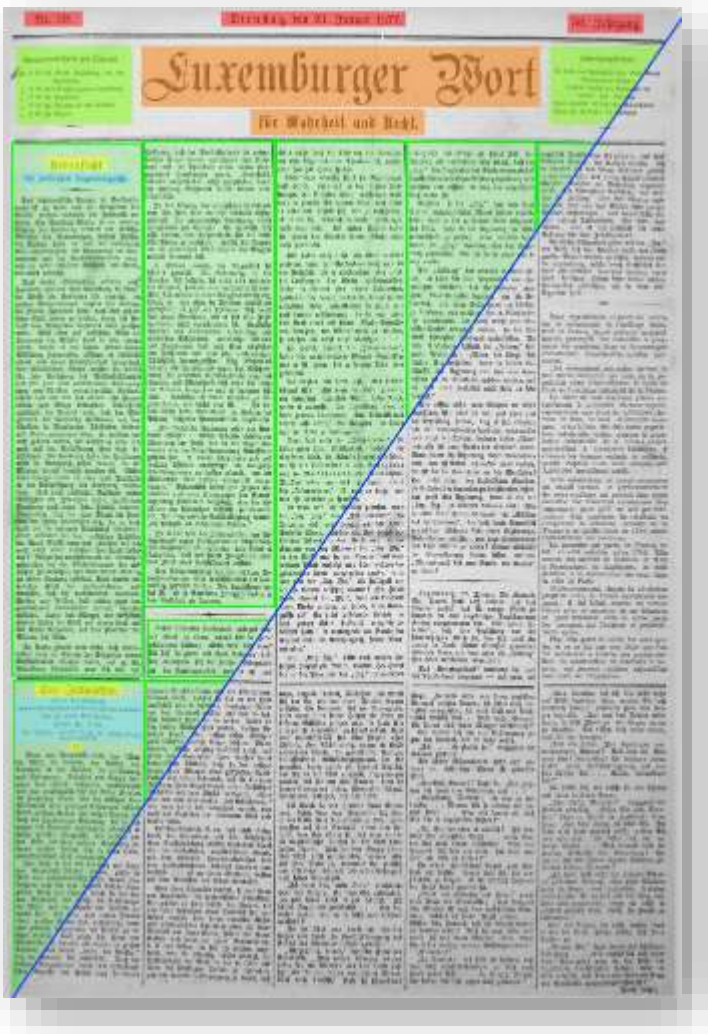- Image & rich metadata (METS/ALTO)

**2015**
- 15 medieval manuscripts with metadata
- 441 books with METS/ALTO

**2016+**
- Newspapers, books, postcards etc.

# Standards



- **300 DPI** Images
- **Uncompressed**
- 8 bit
- 256 levels of grey or full color

- **PDF/A** for online download
- With **Full Text** and **TOC**

- **METS:** Container Format
- **ALTO:** OCR Format

**For each issue / book:**
- 1 METS file
- 1 PDF file

**For each page:**
- 1 TIFF
- 1 PNG (B&W)
- 1 PDF
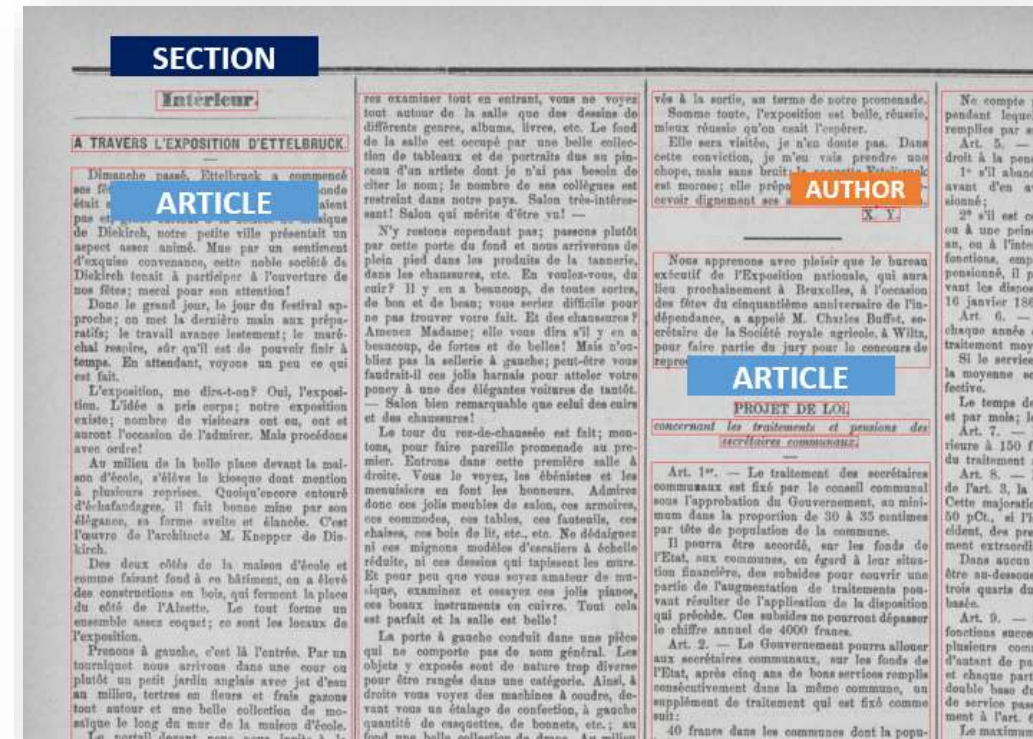- 1 ALTO

# METS / ALTO

- **METS:**
  - 1 METS per document
  - Describes the logical and physical structure
  - Contains all the metadata (technical, descriptive, access rights etc.)
  - Links to ALTO blocks

- **ALTO:**
  - 1 ALTO file per page
  - OCR results with text coordinates
  - Organized in Text Blocks, Composed Blocks, Lines and Strings elements
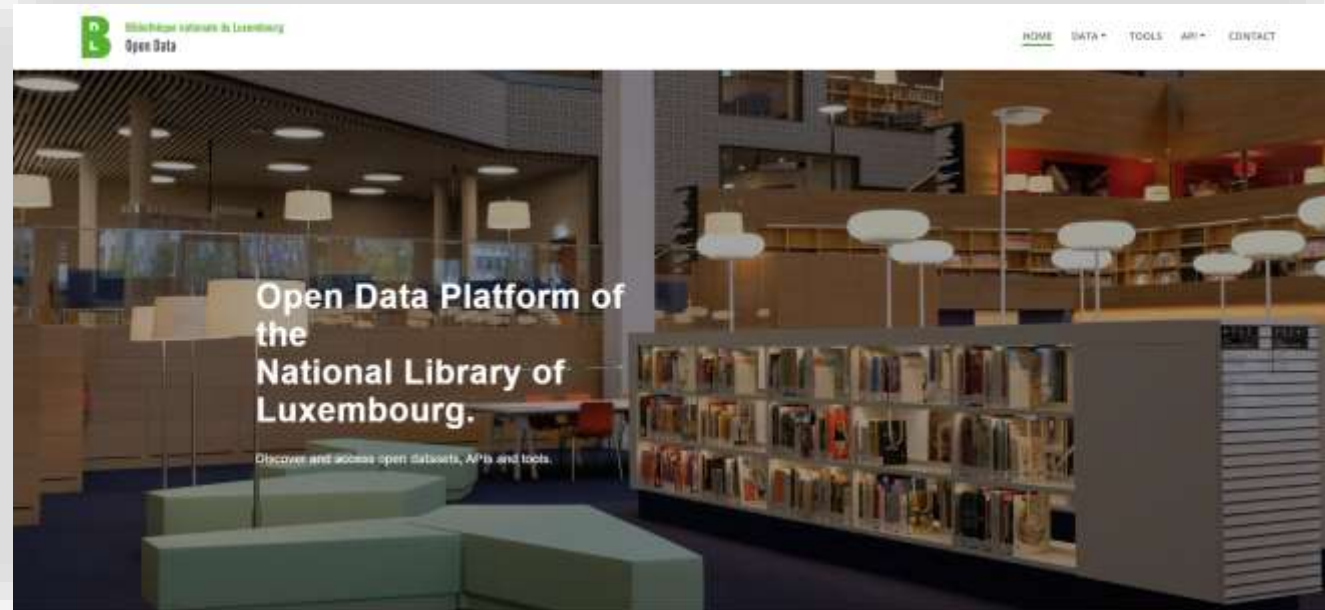
# Open Data & Open Tools

# 2019: data.bnl.lu



**Data**          **Tools**          **APIs**

# data.bnl.lu

- Organised in Packs

- Download what you need

- Raw or Processed Data

- **Technical Specifications**

# Hackathon – Open Data Challenge

- Create novel interfaces / tools ideally using machine learning techniques

- METS/ALTO is not user-friendly when you have 2 days to learn it & use it

# Quality Assurance Tool

## Namalysator

➢ Checks METS/ALTO files

➢ Measures accuracy on the "titles to check manually"

➢ Generates Report (custom error codes)

➢ Open Source

https://github.com/natliblux/Namalysator

# METS / ALTO Viewer

# METS / ALTO Viewer



**Article Segmentation**

**Full Text**

**Block-level Links**

**Word-Highlight**

**Word-To-Image Link**

# METS / ALTO Viewer



Metadata

Table of Contents

Calendar

Search

IIIF Images

Redaction

# A.I. / Machine Learning Projects

# eLmA = eLuxemburgensia meets AI

**Improving OCR**

- Ground Truth
- Engine Training
- Re-OCR

**Article Recommendation**

- Visual Exploration
- Named Entities

# OCR

Search

Police correctionnelle de Luxembrg.

Page 4

⊙ Police correctionnelle de Luxembrg.

$). ©,*©p., ^anbelèinmm unb gctïber tnt ©vttttb, Sovftabt Suiemburg , tft tn ber te|ten ©t'Çuitg beô Budjtpcltjct'gcricljté ju etncrQefbbufie non 100 gran* ïen mib in bte tïoftcn ocvurt^cilt «jorben, afè liber* füf;vt, ben Servit ÿ. t?., ©erber, tn ber 9ïad;t, aitf offenev ©trafic unb cljne ba$u aufgeretjt toorben ju fetn, gcfélagett unb mtfifmnbcft 311 ^abeit.

Polîce correctionnelle de Luxembrg.

P. G.=Sp., Handelsmann und Färber im Grund, Vorstadt Luxemburg, ist in der letzten Sitzung des Zuchtpolizeigerichts zu einer Geldbuße von 100 Franken und in die Kosten verurtheilt worden, als überführt, den Herrn P. K., Gerber, in der Nacht, auf offener Straße und ohne dazu aufgereizt worden zu sein, geschlagen und mißhandelt zu haben.

# Varying quality for OCR

La machine parlementaire a recom-
mencé à fonctionner. Sauf quelques in-
terruptions, elle va marcher pendant dix
mois, sur lesquels les députés en per-
dront largement quatre ou cinq à des
choses qui leur sembleront, à eux et au
public, très importantes, et qui seront
parfaitement inutiles.

La machine parlementaire a recom-
mencé à fonctionner. Sauf quelques in-
terruptions, elle va marcher pendant dix
mois, sur lesquels les députés en per-
dront largement quatre ou cinq à des
choses qui leur sembleront, à eux et au
public, très importantes, et qui seront
parfaitement inutiles.

Hierauf werden Gerüchte zurückgeführt, wonach
die Prinzessin infolge der Aufregung über den Erlaß
des Königs Georg von Sachsen, einen Selbstmord
versuch begangen habe.

Hierauf werden Gerüchte zim'cl.elührt, wonach
bie P,i zeffüt infolge ber AufregU'g üdtc d',< E,Il>^
des Königs Georg von Hachsen, einen Selb fl!u»r b»
uerfuch begangen habe.

**Ground Truth**
≈100 thousand text lines

# Newspaper Corpus
≈200 million text lines

Ground Truth

unsupervised quality evaluation • kNN algorithm

🟢 score > 0.95

🔵 else

score = 1 - ( editDistance / |chars| )

**85% Accuracy on test set**

Original Text Block



Binarized Text Block

cleaning • dilation • padding • inversion • white on black detection

P. G.=Sp., Handelsmann und Färber im Grund,
Vorstadt Luxemburg, ist in der letzten Sitzung des
Zuchtpolizeigerichts zu einer Geldbuße von 100 Fran=
ken und in die Kosten verurtheilt worden, als über=
führt, den Herrn P. K., Gerber, in der Nacht, auf
offener Straße und ohne dazu aufgereizt worden zu
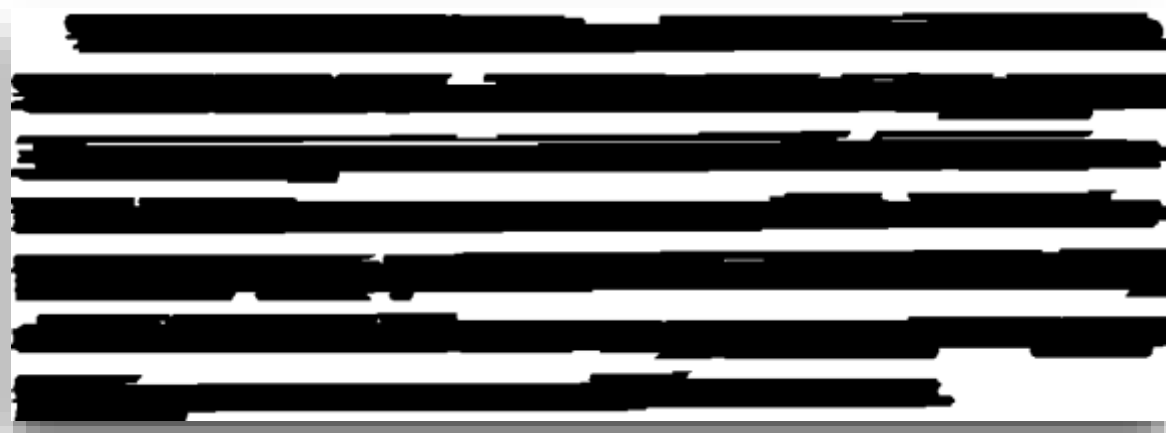sein, geschlagen und mißhandelt zu haben.

**99% Accuracy**
on test set

morphology • connected components • horizontal histogram projection

Convolutional Neural Network

≈50 thousand chars

*Fraktur*

Regular

**99% Accuracy**
on test set

font recognition • binary classifier

**kraken**                    **Tesseract**                    **Calamari**

**96% Accuracy**
on test set

open source software • custom model training
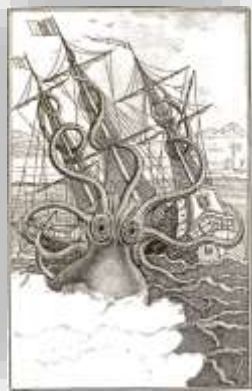
```
18      <TextLine HPOS="0" VPOS="41" WIDTH="761" HEIGHT="41">
19          <String CONTENT="Vorstadt" CC="00000000" HPOS="0" VPOS="41" WIDTH="135" HEIGHT="41"/>
20          <SP/>
21          <String CONTENT="Luxemburg," CC="0000000000" HPOS="135" VPOS="41" WIDTH="189" HEIGHT="41"/>
22          <SP/>
23          <String CONTENT="ist" CC="000" HPOS="324" VPOS="41" WIDTH="49" HEIGHT="41"/>
24          <SP/>
25          <String CONTENT="in" CC="00" HPOS="373" VPOS="41" WIDTH="42" HEIGHT="41"/>
26          <SP/>
27          <String CONTENT="der" CC="000" HPOS="415" VPOS="41" WIDTH="63" HEIGHT="41"/>
28          <SP/>
29          <String CONTENT="letzten" CC="0000100" HPOS="478" VPOS="41" WIDTH="97" HEIGHT="41"/>
30          <SP/>
31          <String CONTENT="Sitzung" CC="1000000" HPOS="575" VPOS="41" WIDTH="118" HEIGHT="41"/>
32          <SP/>
33          <String CONTENT="des" CC="000" HPOS="693" VPOS="41" WIDTH="68" HEIGHT="41"/>
34      </TextLine>
```

word bounding boxes • ALTO XML schema

spaCy

Language Models

P.G.≶Sp. , Handelsmann und Färber im Grund,
Vorstadt Luxemburg, ist in der letzten Sitzung des
Zuchtpolizeigerichts zu einer Geldbuße von 100 Frank≶
ken und in die Kosten verurtheilt worden, als über≶
führt, den Herrn P. K., Gerber, in der Nacht, auf
offener Straße und ohne dazu aufgereizt worden zu
sein, geschlagen und mißhandelt zu haben.

≈35 thousand entities

Improved
OCR

Location

Person

Organisation

Date

ground truth generation • named entity recognition (NER)

Bibliothèque nationale du Luxembourg
Nationalbibliothéik

# Results

## 1) Improve Optical Character Recognition

**Accuracy**
Improvement of ~30% text blocks

**Production**
OCR application on entire corpus

## 2) Improve Newspaper Exploration

**Named Entity Recognition**
Finished ground truth - observing first results

**New User Interface**
Proof of concept

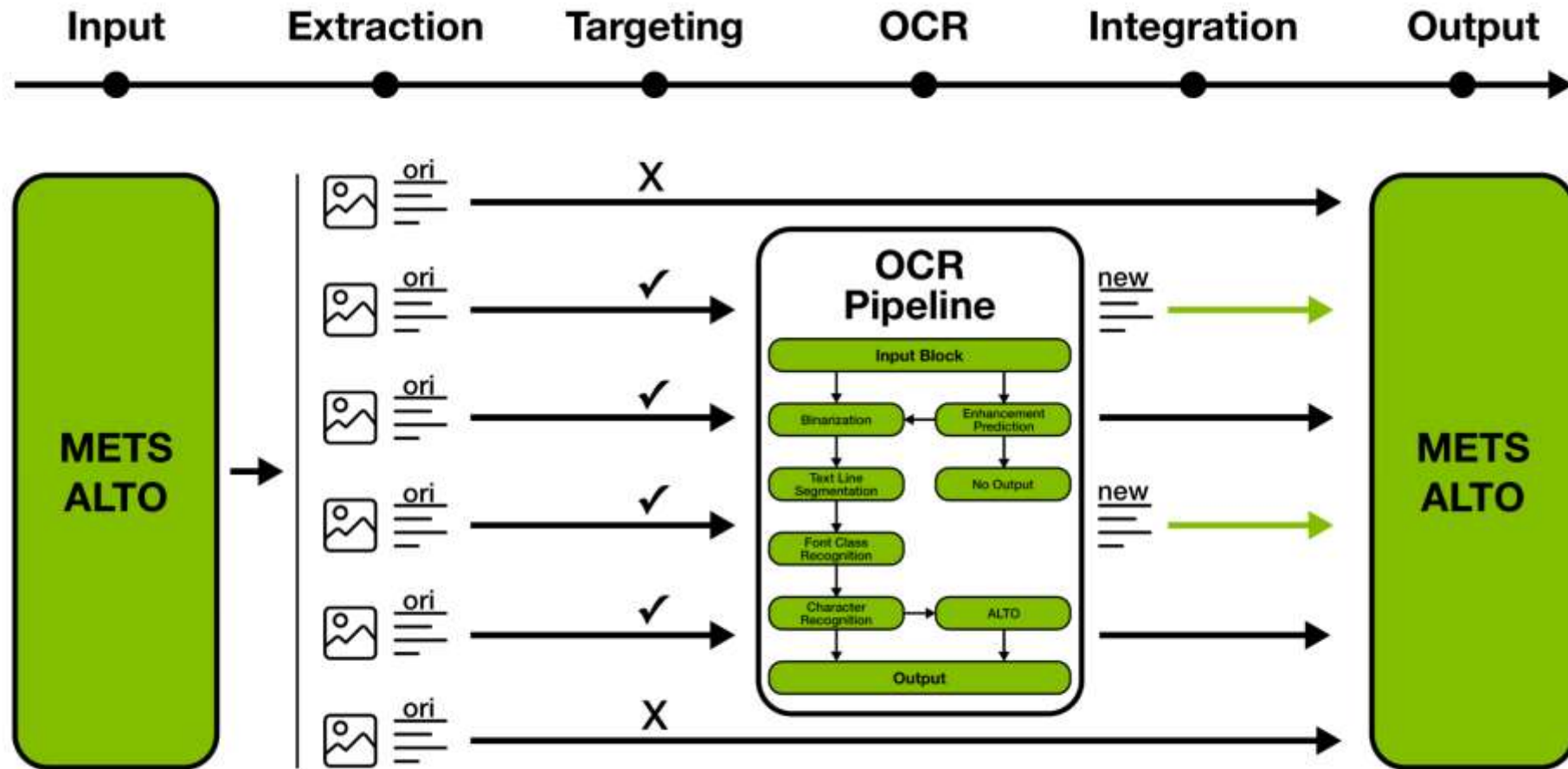Bibliothèque nationale du Luxembourg
Open Data

data.public.lu

Bibliothèque nationale du Luxembourg
Nationalbibliothéik

# POC



- Named Entities
- Entity Relations
- Timeline
- Maps
- Wikidata

# Nautilus-OCR



https://github.com/natliblux/nautilusocr

# data.bnl.lu

## OCR Datasets

As part of BnL's AI strategy, we provide the ground truth data that falls into the public domain (CC0, see copyright notice). Available in two variations, the datasets cover historical newspapers published before 1878. The data is generally in German, French or Luxembourgish and has been manually corrected for a minimum accuracy of 99.95%.

### GROUND TRUTH PACK

## 33.000

transcribed text lines

- Text line based OCR
- 19.000 text lines in Antiqua
- 14.000 text lines in Fraktur
- Transcribed using double-keying (99.95% accuracy)
- Public Domain, CC0 (See copyright notice)
- Best for training an OCR engine

**⬇ DOWNLOAD (ZIP)**

### RAW GROUND TRUTH PACK

## 1.700

text blocks

- Raw uncropped text blocks
- Pairs consist of block image and ALTO XML
- Public Domain, CC0 (See copyright notice)
- Best for testing/training text line segmentation

**⬇ DOWNLOAD (ZIP)**

Bibliothèque nationale du Luxembourg
Nationalbibliothéik

# Thank you! Questions?

**Ralph Marschall**

Ralph.Marschall@bnl.etat.lu

**@RalphMarschall**

Bibliothèque nationale du Luxembourg
Nationalbibliothéik

**@bnluxembourg**