# How to read a million books

Clemens Neudecker
Staatsbibliothek zu Berlin
17 May 2017



Staatsbibliothek
zu Berlin
Preußischer Kulturbesitz

# How to read a million books?

Clemens Neudecker
Staatsbibliothek zu Berlin
17 May 2017

# How to read a million ~~books~~ newspapers?

Clemens Neudecker
Staatsbibliothek zu Berlin
17 May 2017

# How to read* a million newspapers?

## (* hint: <u>do</u> try this at home)

Clemens Neudecker
Staatsbibliothek zu Berlin
17 May 2017

# Agenda

- Digital Collections

- Data, Tools, Formats

- Europeana Newspapers

- NLP challenges

- Experiments & use cases
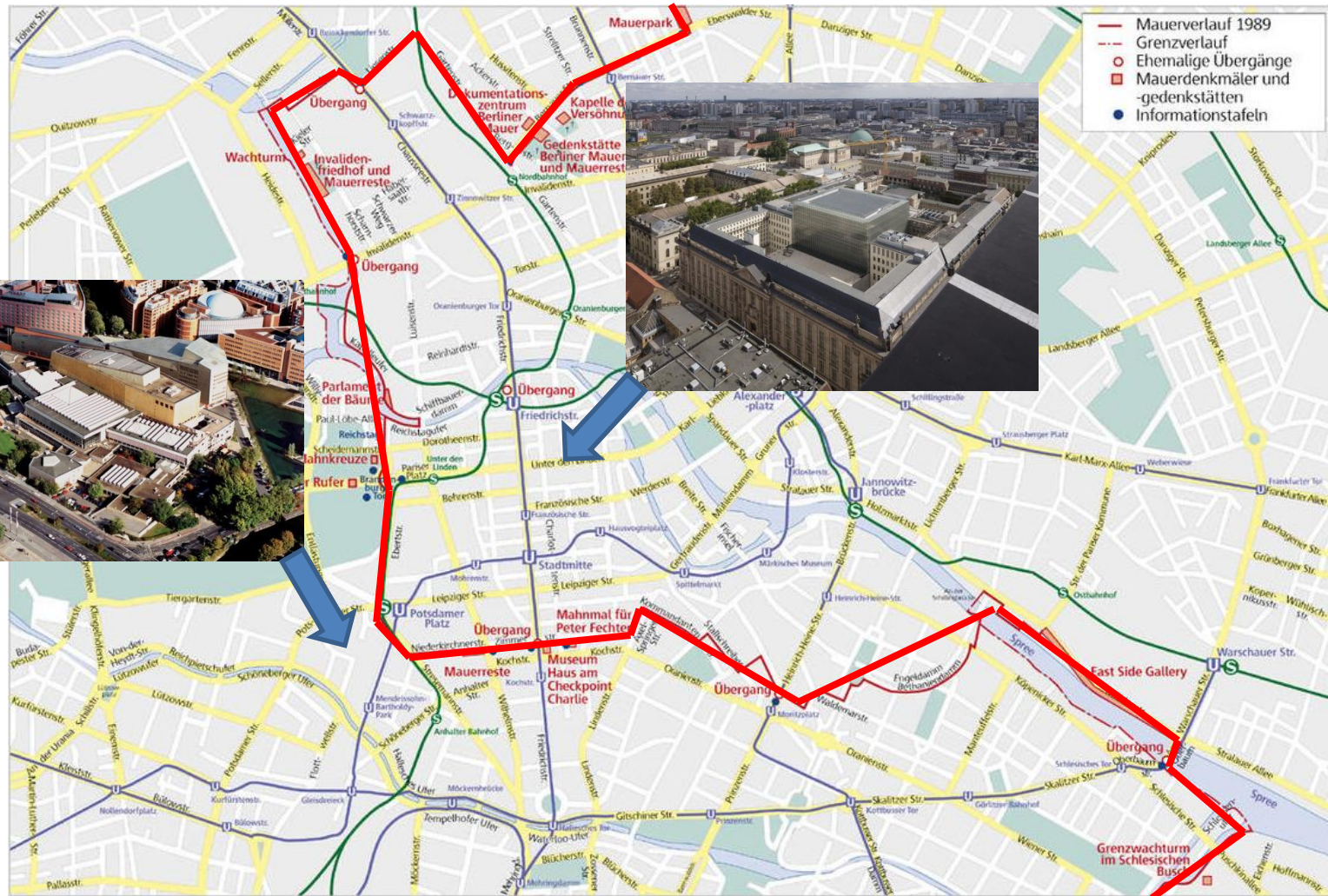


https://xkcd.com/1838

# About me

- Research Coordinator @ Berlin State Library

- M.A. Philosophy, Computer Science, Political Science

- Mostly curious about
  - Optical Character Recognition, Document Analysis
  - Natural Language Processing
  - Digital Humanities

- More: @cneudecker, cneud.net

# Staatsbibliothek zu Berlin

- Established 1661 as the Library of the King of Prussia

- Today largest research library in Germany,
  with approx. 11.5m volumes (23m objects)

- Part of the „Stiftung Preußischer Kulturbesitz",
  a unique union of museums, archives, libraries
  and research institutes from Berlin

- http://staatsbibliothek-berlin.de/

# Digitisation 2.0



SCANROBOT® 2.0 MDS
MASS DIGITIZATION SYSTEM

**High speed scanning**
> 3,000 pages / hour

# Europeana

- http://www.europeana.eu/portal/

- Europe's Digital Library



- > 53m objects incl.
  art, sound, fashion,...

- API: http://labs.europeana.eu/api

# DPLA

- [https://dp.la/](https://dp.la/)

- Digital Public Library of America

- > 16m objects

- API: [https://dp.la/info/developers/codex/](https://dp.la/info/developers/codex/)

# Hathi Trust

- https://www.hathitrust.org/

- Public copy of Google Books

- > 15m volumes

- API: https://www.hathitrust.org/data

# DDB

- http://ddb.de/

- Germany's federal Digital Library

- \> 9m objects

- API: https://api.deutsche-digitale-bibliothek.de/

# Trove

- http://trove.nla.gov.au/

- Digital Library
  of Australia

- > 540m objects

- API:
  http://help.nla.gov.au/trove/building-with-trove/api

# Formats & Standards

- What data is available?

- Typically, a digital object is composed of:
  - Scanned Images in TIFF, JP2 or JPEG
  - Descriptive metadata in DublinCore
  - Structural metadata in METS
  - Text content in ALTO or TEI
  - Europeana in EDM
  - Linked Data in RDF or JSON-LD

# 🐍 Tools (1/3)

- OAI-PMH
  - [https://pypi.python.org/pypi/Sickle](https://pypi.python.org/pypi/Sickle)
  - [https://pypi.python.org/pypi/pyoai](https://pypi.python.org/pypi/pyoai)
  - [https://pypi.python.org/pypi/oaiharvest](https://pypi.python.org/pypi/oaiharvest)

- METS
  - [https://pypi.python.org/pypi/metsrw](https://pypi.python.org/pypi/metsrw)
  - [https://pypi.python.org/pypi/pymets](https://pypi.python.org/pypi/pymets)

# 🐍 Tools (2/3)

- DublinCore
  - https://pypi.python.org/pypi/pydc
  - https://pypi.python.org/pypi/dcxml
- Europeana
  - https://pypi.python.org/pypi/europeana-search
  - https://pypi.python.org/pypi/django-europeana

# 🐍 Tools (3/3)

- IIIF
  - https://pypi.python.org/pypi/iiif/
  - https://pypi.python.org/pypi/Flask-IIIF/
- KB NL
  - https://pypi.python.org/pypi/kb
  - https://github.com/KBNLresearch/intro-kb-apis
  - http://lab.kb.nl/

# Europeana Newspapers

- EU-project to make Europe's historical newspapers searchable & accessible

EXPLORE
Historic Newspapers on www.europeana.eu

SEARCH
Our Historic Newspapers
Try our prototype newspaper browser!

USE
Our Open-Source Tools
Designed to help libraries and others refine and improve the quality of digital historic newspapers.

- http://www.europeana-newspapers.eu/

# Europeana Newspapers Collection

- 12 million historic newspaper pages text (> 10.000.000.000 tokens)

- 40 languages, 4 alphabets

- 400 years (1618 – 2016)

- http://www.theeuropeanlibrary.org/tel4/newspapers

# OCR / OLR



(U.lag nul «chestttetrung- ■geeinoel II, Setch«it,zen I—Ig Ufr sterntpeechee g» U II. für ftrene-geingelpilche: 13 01191 nnd 13 03 11 io"gl f l««lt-beOeu; OetHn *1, blnftraße IS IZeinsptechee; H I Sanemeinummet gurfilrft 8MB); ««de«: gdn.o(tio||e III (ZemlpreAei 284.3»). lie.gonlen nur nnier heimnnn » Erden bei der veutlchen Bonl »n« Vtdennld-Getelltchoil gttloto bumduig. Commerz- nndprinoldonlN voINchrSomI bomduig u 189 ß>, .»ontbeegee Kochelchlen- eitchelne» 12 mal wSchenNIch. täglich zweimal — morgen« nnd ndendn —, Sonntage nnr morgen». Toonlnge nur abend» Zn den Kochdorerlen wird die Ndend-Nuegode noch am üben!

Dieser Entwurf ist. wie Bürgermeister Roß mitteilte, den Fraktionen zur Stellungnahme vorgelegt worden. Zum Donnerstag war eine zweite Sitzung der Fraktionsfübrer vom Vertreter des Senats angeordnet worden, zu der ober zwei Fraktionen, die Teutschnationalen und die Nationalsozialisten, nicht erschienen waren. Von den Nationalsozialisten ist kurz vor Beginn der Sitzung eine telephonische Erklärung abgegeben worden, etwa des Inhalts, daß die Fraktion sich den sachlichen Verhandlungen entziehen müsse, solange nicht gewisse Vorbedingungen erfüllt sein.

http://www.theeuropeanlibrary.org/tel4/newspapers/issue/Hamburger_Nachrichten/1932/12/31
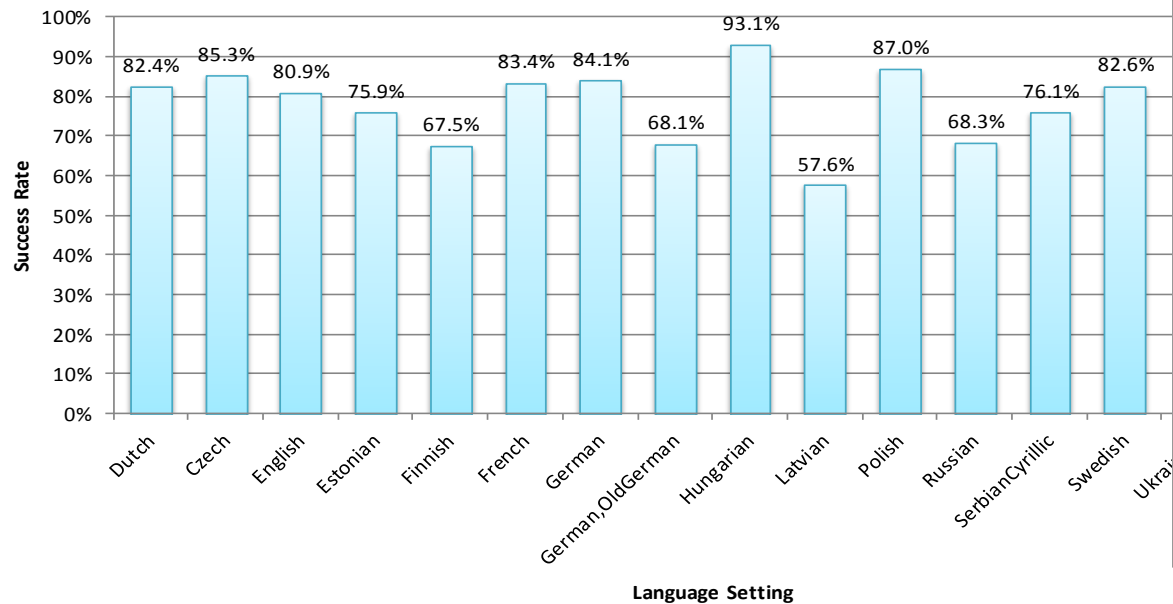
```xml
<Page ID="ID1" HEIGHT="25941" WIDTH="20463" PHYSICAL_IMG_NR="1" PRINTED_IMG_NR="[1]" PROCESSING="ID0">
    <PrintSpace HEIGHT="24453.0" WIDTH="19500.0" HPOS="294.0" VPOS="954.0" PC="0.93230003">
        <TextBlock ID="ID2" HEIGHT="1701" WIDTH="16026" HPOS="294" VPOS="972" ns1:type="simple" language="en">
            <TextLine HEIGHT="90.0" WIDTH="3570.0" HPOS="3960.0" VPOS="972.0">
                <String HEIGHT="81.0" WIDTH="36.0" HPOS="3960.0" VPOS="981.0" CONTENT="J" WC="0.8095238"/>
                <SP WIDTH="36.0" HPOS="3996.0" VPOS="981.0"/>
                <String HEIGHT="66.0" WIDTH="42.0" HPOS="4032.0" VPOS="993.0" CONTENT="a" WC="0.8095238"/>
                <SP WIDTH="68.0" HPOS="4074.0" VPOS="993.0"/>
                <String HEIGHT="81.0" WIDTH="228.0" HPOS="4143.0" VPOS="972.0" CONTENT="Ira" WC="0.95238096"/>
                <SP WIDTH="447.0" HPOS="4371.0" VPOS="972.0"/>
                <String HEIGHT="84.0" WIDTH="81.0" HPOS="4818.0" VPOS="975.0" CONTENT="mj" WC="0.8095238"/>
                <SP WIDTH="687.0" HPOS="4899.0" VPOS="975.0"/>
                <String HEIGHT="63.0" WIDTH="105.0" HPOS="5586.0" VPOS="984.0" CONTENT="iI" WC="0.8095238"/>
                <SP WIDTH="416.0" HPOS="5691.0" VPOS="984.0"/>
                <String HEIGHT="90.0" WIDTH="810.0" HPOS="6108.0" VPOS="972.0" CONTENT="tE1r" WC="0.8095238"/>
                <SP WIDTH="585.0" HPOS="6918.0" VPOS="972.0"/>
                <String HEIGHT="51.0" WIDTH="27.0" HPOS="7503.0" VPOS="1008.0" CONTENT="3" WC="0.8095238"/>
            </TextLine>
            <TextLine HEIGHT="39.0" WIDTH="24.0" HPOS="6792.0" VPOS="1065.0">
                <String HEIGHT="39.0" WIDTH="24.0" HPOS="6792.0" VPOS="1065.0" CONTENT="i" WC="0.8095238"/>
            </TextLine>
            <TextLine HEIGHT="42.0" WIDTH="30.0" HPOS="2517.0" VPOS="1125.0">
                <String STYLEREFS="ID3" HEIGHT="42.0" WIDTH="30.0" HPOS="2517.0" VPOS="1125.0" CONTENT="c" WC="0.8095238"/>
            </TextLine>
            <TextLine HEIGHT="1269.0" WIDTH="16026.0" HPOS="294.0" VPOS="1404.0">
                <String HEIGHT="1269.0" WIDTH="2052.0" HPOS="294.0" VPOS="1404.0" CONTENT="JiLas" WC="0.8095238">
                    <ALTERNATIVE>Jails</ALTERNATIVE>
                </String>
                <SP WIDTH="272.0" HPOS="2346.0" VPOS="1404.0"/>
                <String STYLEREFS="ID4" HEIGHT="219.0" WIDTH="1023.0" HPOS="2619.0" VPOS="2184.0" CONTENT="Edition" WC="0.95238096"/>
                <SP WIDTH="677.0" HPOS="3642.0" VPOS="2184.0"/>
                <String STYLEREFS="ID5" HEIGHT="771.0" WIDTH="1920.0" HPOS="4320.0" VPOS="1839.0" CONTENT="THE" WC="1.0">
                    <ALTERNATIVE>TIIE</ALTERNATIVE>
                    <ALTERNATIVE>TIE</ALTERNATIVE>
                </String>
                <SP WIDTH="503.0" HPOS="6240.0" VPOS="1839.0"/>
                <String STYLEREFS="ID5" HEIGHT="792.0" WIDTH="6312.0" HPOS="6744.0" VPOS="1833.0" CONTENT="WINCHESTER" WC="1.0"/>
                <SP WIDTH="527.0" HPOS="13056.0" VPOS="1833.0"/>
                <String STYLEREFS="ID5" HEIGHT="771.0" WIDTH="2736.0" HPOS="13584.0" VPOS="1824.0" CONTENT="NEWS" WC="1.0"/>
            </TextLine>
        </TextBlock>
    </Printspace>
</Page>
```

https://github.com/cneud/alto-tools

# Performance



**Bag of Words OCR Evaluation**
Per Language

**Layout Analysis Performance**
Per evaluation profile

# Experimental(!) Downloads

- http://data.theeuropeanlibrary.org/download/newspapers-by-country/README.html
- http://research.europeana.eu/itemtype/newspapers
- http://test-solr-mongo.eanadev.org/europeana-research-newspapers-dump/sample-2017-04-26/Staatsbibliothek_zu_Berlin_Preu%253Fischer_Kulturbesitz/titles.html

# OCR

- EU: IMPACT project (2008-2012)
- US: eMOP project (2013-2015)
- DE: OCR-D project (2016-2018)

- Google:
  - Tesseract
  - ocropy (fka OCRopus)
  - Aksara

# Named Entity Recognition

- 3 Categories:
  - PERSON; LOCATION; ORGANIZATION
- 3 Languages:
  - Dutch; French; German

- Powered by Stanford CoreNLP - CRF-NER

The Stanford Natural Language Processing Group

# Annotations

| Language | # tokens | # PER | # LOC | # ORG |
|---|---|---|---|---|
| French | 207,000 | 5,672 | 5,614 | 2,574 |
| Dutch | 182,483 | 4,492 | 4,448 | 1,160 |
| German | 96,735 | 7,914 | 6,143 | 2,784 |

| Language | # tokens | # PER | # LOC | # ORG |
|---|---|---|---|---|
| French | 100% | 2,75% | 2,71% | 1,24% |
| Dutch | 100% | 2,46% | 2,44% | 0,64% |
| German | 100% | 8,18% | 6,35% | 2,88% |

| Language | Word-Error-Rate (Bag of Words) | Reading Order Success Rate |
|---|---|---|
| French | 16,6% | 19,9% |
| Dutch | 17,6% | 23,2% |
| German | 15,9% / 21,9% | 13,6% |

# Evaluation



Dutch

French

# Challenges



- [https://github.com/EuropeanaNewspapers/ner-corpora/wiki/Corpus-cleanup](https://github.com/EuropeanaNewspapers/ner-corpora/wiki/Corpus-cleanup)

# Lack of metadata

**Issue**

There is no associated metadata for the annotated text (newspaper title, date, etc.)

**Solution**

Automatically match lines with newspaper pages through keyword search

# OCR errors vs. historical spelling

**Issue**

Text contains OCR errors but also valid(!) historical spelling variants

$$teil \xrightarrow{t \to th, ei \to ey} theyl$$

**Solution**

Document language profiling to distinguish OCR errors and spelling variants

# Sentence splits

**Issue**

　During data pre-processing, (parts) of sentences have been erroneously cut

**Solution**

　Reconstruct sentences through keyword search and matching procedure

# Hyphenation

**Issue**

Text contains hyphenation to be removed but hyphens do also occur in regular text

**Solution**

Use a tokenizer to determine hyphens to be removed

# Missing tags

**Issue**

Human operators forgot to tag some entities or tagged them with the wrong category
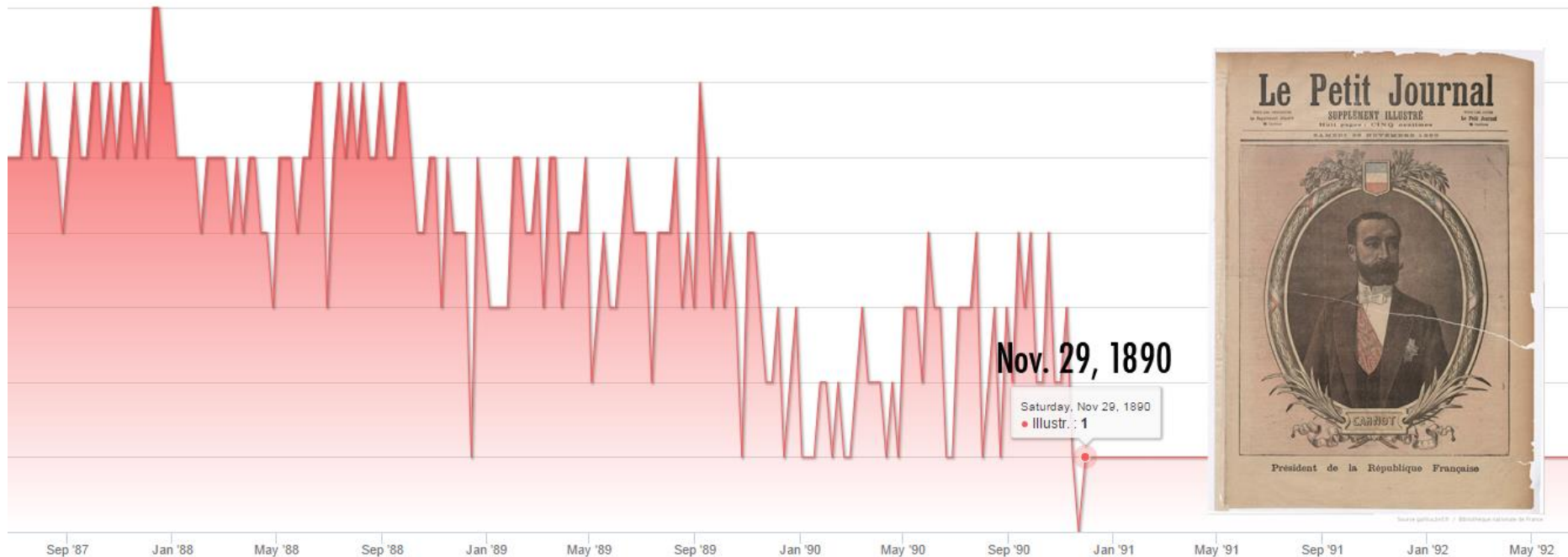
**Solution**

???

# Punctuation

**Issue**

According to CONLL, punctuation should be in a separate line from the token - but abbreviations…

**Solution**

???

https://altomator.github.io/EN-data_mining/

altomator / **EN-data_mining**

Watch 2    Star 7    Fork 1

<> Code    ① Issues 0    Pull requests 2    Projects 0    Pulse    Graphs

Data Mining Historical Newspaper Metadata (METS/ALTO formats)    http://altomator.github.io/EN-data_mi...

45 commits    2 branches    0 releases    1 contributor

Branch: master    New pull request    Find file    Clone or download

altomator committed on GitHub Delete .DS_Store    Latest commit 72127df on 3 Mar

| | | |
|---|---|---|
| API XQuery | XQUERY | 3 months ago |
| Charts | XQUERY | 3 months ago |
| Datasets | XQUERY | 3 months ago |
| Scripts | XQUERY | 3 months ago |
| README.md | Update README.md | 5 months ago |

https://github.com/altomator/EN-data_mining

# Dictionary viewer



http://www.kbresearch.nl/dictionary/

http://www.kbresearch.nl/telraam/

```python
 1    #!/usr/bin/env python3
 2
 3    #
 4    #    Telraam historische kranten v0.03
 5    #    Author: WillemJan Faber -- http://www.kb.nl/lab/
 6    #
 7    #    This program is free software: you can redistribute it and/or modify
 8    #    it under the terms of the GNU General Public License as published by
 9    #    the Free Software Foundation, either version 3 of the License, or
10    #    (at your option) any later version.
11    #
12    #    This program is distributed in the hope that it will be useful,
13    #    but WITHOUT ANY WARRANTY; without even the implied warranty of
14    #    MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE.  See the
15    #    GNU General Public License for more details.
16    #
17    #    You should have received a copy of the GNU General Public License
18    #    along with this program.  If not, see <http://www.gnu.org/licenses/>.
19    #
20
21    import requests
22    import urllib
23    import xml.etree.ElementTree as ET
24
25    SRU_DDD = "http://jsru.kb.nl/sru/sru"
26    SRU_DDD += "?version=1.2&operation=searchRetrieve"
27    SRU_DDD += "&x-collection=DDD_krantnr&recordSchema=ddd"
28    SRU_DDD += "&startRecord=%i&maximumRecords=1000000&query=(%s)"
```
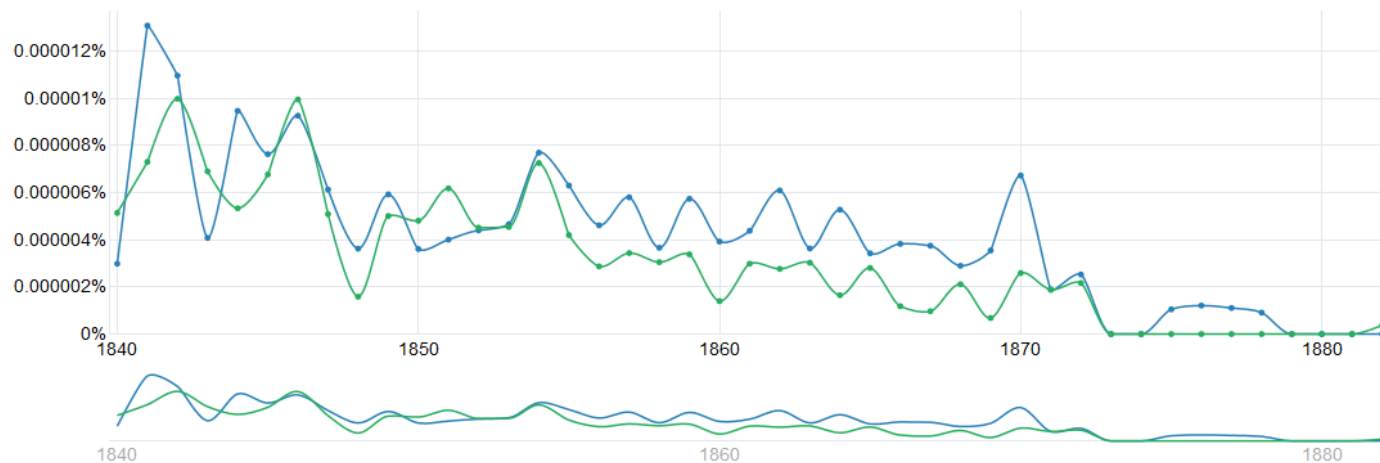
https://gist.github.com/WillemJan/6ab02c48af576ba47b68

http://ngramviewer.kbresearch.nl/

ILPS / Political Mashup / pm-ngramviewers-kbkranten

## Overview

| | | | |
|---|---|---|---|
| Last updated | 2016-02-17 | **0** Open PRs | **2** Watchers |
| Language | Python | **1** Branch | **0** Forks |
| Access level | Read | | |

### KBKranten Ngramviewer

The KBKranten Ngramviewer analyzes and indexes the <u>Historische Kranten</u> corpus of the <u>Royal Library</u> up to 5-grams in Elasticsearch, and provides an API that allows for `suggest/<term>` and `search/<term>` on that corpus. The `search/<term>` endpoint filters the corpus on the `<term>` and returns a histogram of the frequency and the *relative* frequency (frequency of that ngram in a year / total number of ngrams with that n in that year).

Also, there is a <u>webinterface</u> on top of the API that visualizes the `<term>` histograms, and links back to the <u>Historische kranten</u> website to provide the user with some context.

**https://bitbucket.org/ilps/pm-ngramviewers-kbkranten**

http://www.digitalvictorianist.com/
https://twitter.com/VictorianHumour

BL-Labs / **embellishments**

Watch ▾ 4    ★ Unstar 13    Fork 1

<> Code    ⓘ Issues 0    Pull requests 0    Projects 0    Wiki    Pulse    Graphs

No description, website, or topics provided.

🕐 19 commits    1 branch    0 releases    1 contributor    MIT

Branch: master ▾    New pull request    Create new file    Upload files    Find file    Clone or download ▾

benosteen Adding flickr upload queue loader    Latest commit 6bdff60 on 22 Nov 2013

| | | |
|---|---|---|
| 📁 haarcascades | Added current state of code | 4 years ago |
| 📄 ALTO_Example_000411385_000100.xml | Added some example xml files | 4 years ago |
| 📄 LICENSE | Initial commit | 4 years ago |
| 📄 Metadata_Example.xml | Added some example xml files | 4 years ago |
| 📄 README.md | Update README.md | 4 years ago |
| 📄 auth_clean.py | Added cleaned versions of the tumblr oauth code | 4 years ago |

https://github.com/BL-Labs/embellishments

http://networks.viraltexts.org/1836to1899/index.html

https://github.com/dasmiq/passim

EUROPEANA TRANSCRIBATHON CAMPUS BERLIN 2017

22-23 June 2017

Berlin State Library

http://pro.europeana.eu/event/europeana-transcribathon-campus-2017

---

europeana

**TRANSCRIBATHON CAMPUS BERLIN 2017**

Berliner Tageblatt
und Handels-Zeitung

Die ruffifche Front bei Stryj durchbrochen.

READ THE FULL STORY!

Abb.: Europeana 1914-1918,
Rolf Kranz, CC-BY-SA.
Europeana Newspapers,
Berliner Tageblatt,
28.05.1915, Public Domain.

**STAATSBIBLIOTHEK ZU BERLIN**

22. Juni 2017, 14.00 Uhr bis
23. Juni 2017, 21.00 Uhr

**Potsdamer Straße 33
10785 Berlin**

Entziffern Sie alte Handschriften und verknüpfen Sie Dokumente mit historischen Zeitungen von **Europeana 1914-1918** und **Europeana Newspapers**. Machen Sie mit bei einem internationalen Transkriptionswettbewerb und Data-Enrichment Workshops.

**www.transcribathon.eu/berlin2017**

Staatsbibliothek zu Berlin Preußischer Kulturbesitz    WIKIMEDIA FOUNDATION    Facts & Files    Co-financed by the European Union Connecting Europe Facility

# Thank you for your attention!

Clemens Neudecker
Staatsbibliothek zu Berlin
17 May 2017