# Making Historic Newspapers Available Online: Why, Where and How

IFLA Newspaper Pre-Conference

14 August 2014, Geneva

*Hans-Jörg Lieder, Staatsbibliothek zu Berlin – Preußischer Kulturbesitz | Berlin State Library*
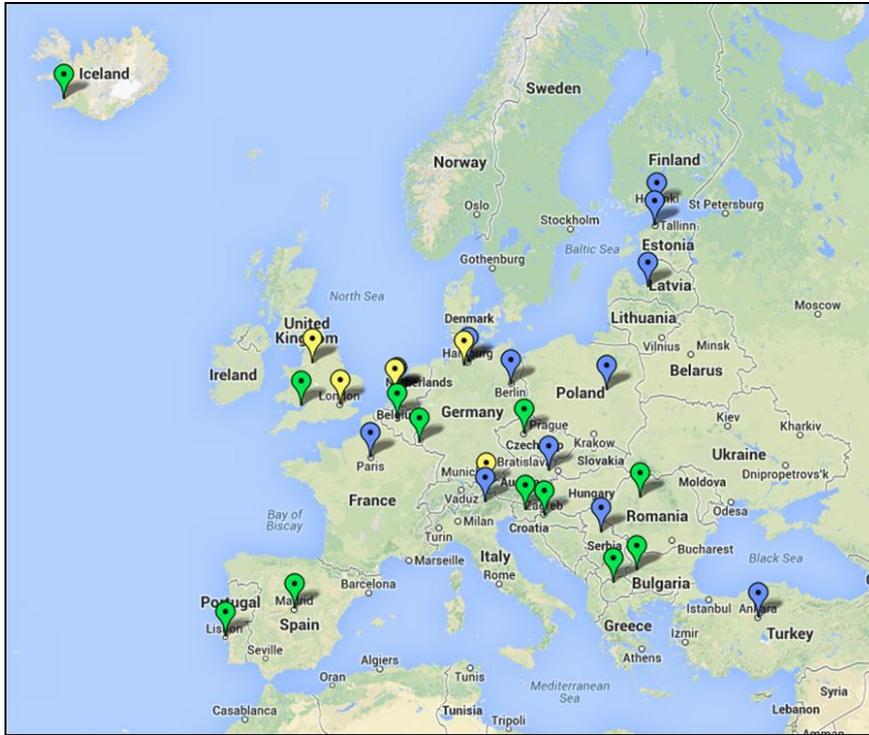
# Why Newspapers?

Cons:

- Originals are cumbersome objects
  - Prone to damage and destruction due to paper quality
  - Missing issues and pages
  - Difficult to deal with from a catalogueing point of view
  - Poor bindings
  - Funny fonts and fading ink
- Microforms may also be cumbersome objects
  - Skewed images, text loss
  - More missing issues and pages, plus duplicate pages

# That's Why!

Pros:

- "Newspapers are the second hand of history"
  - Provide insights into history's microstructure
  - Unlimited thematic scope
  - Interesting for all fields of scholarship, but also for the layman
  - Massive digital newspaper text corpora allow for new ways of research
  - A European perspective: significant contribution to the shaping of identities of peoples and individuals

# The Europeana Newspaper Project – Who?



Blue – Content Providers

Yellow – Service Providers

Green – Associated Partners

# The Europeana Newspaper Project – What?

| Country | Partner | Titel | Startdate (overall) | Enddate (overall) |
|---|---|---|---|---|
| France | BnF | 80 | 1814 | 1944 |
| Germany | SBB | 6 | 1872 | 1940 |
| | SUB-HH | 16 | 1721 | 1945 |
| Netherlands | KB | 203 | 1618 | 1900 |
| Italy | LFT | 15 | 1813 | 1949 |
| Estonia | NLE | 43 | 1852 | 1944 |
| Finland | NLF | 11 | 1900 | 1910 |
| Latvia | NLL | 117 | 1868 | 1955 |
| Poland | NLP | 118 | 1914 | 1939 |
| Turkey | NLT | 22 | 1818 | 1928 |
| Austria | ONB | 275 | 1686 | 1945 |
| Serbia | UB | 45 | 1830 | 1944 |

20 languages

ca. 950 titles

ca. 10m pages refined
- 8m OCR
- 2m OLR
- 2m NER

europeana
newspapers

ICTPSP
ICT POLICY SUPPORT PROGRAMME

# The Europeana Newspaper Project – What else?

- Tools for informed selection of newspapers for digitisation
- Specifications and tools for the creation and validation of OCR-ready images
- Large-scale, highly automated workflows for refinement (OCR, OLR, NER)
- Metadata best practice recommendations
- Transmission of data to European Portals and the Union Catalogue of Serials
- Presentation of results

# What does it look like … in TEL?

# What does it look like … in Europeana?

# What does it look like … in the Union Catalogue of Serials?

# What about Services?

- Richest service portfolio available at local web pages (if you're lucky)
  - Calendar navigation, search in texts
  - filters to narrow down queries or result sets
  - mark-ups, annotations, links to other information resources, etc.

- Services at TEL
  - Calendar navigation, search in texts
  - Filters for searches: title, date, owning library
  - Filters for results: title, date, owning library, country, language
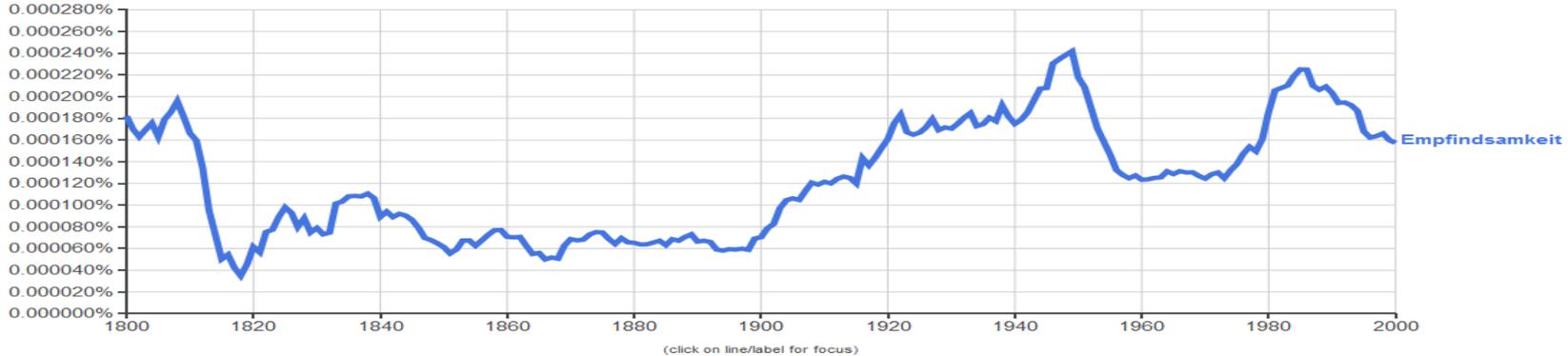
# What about further Services? – An Example
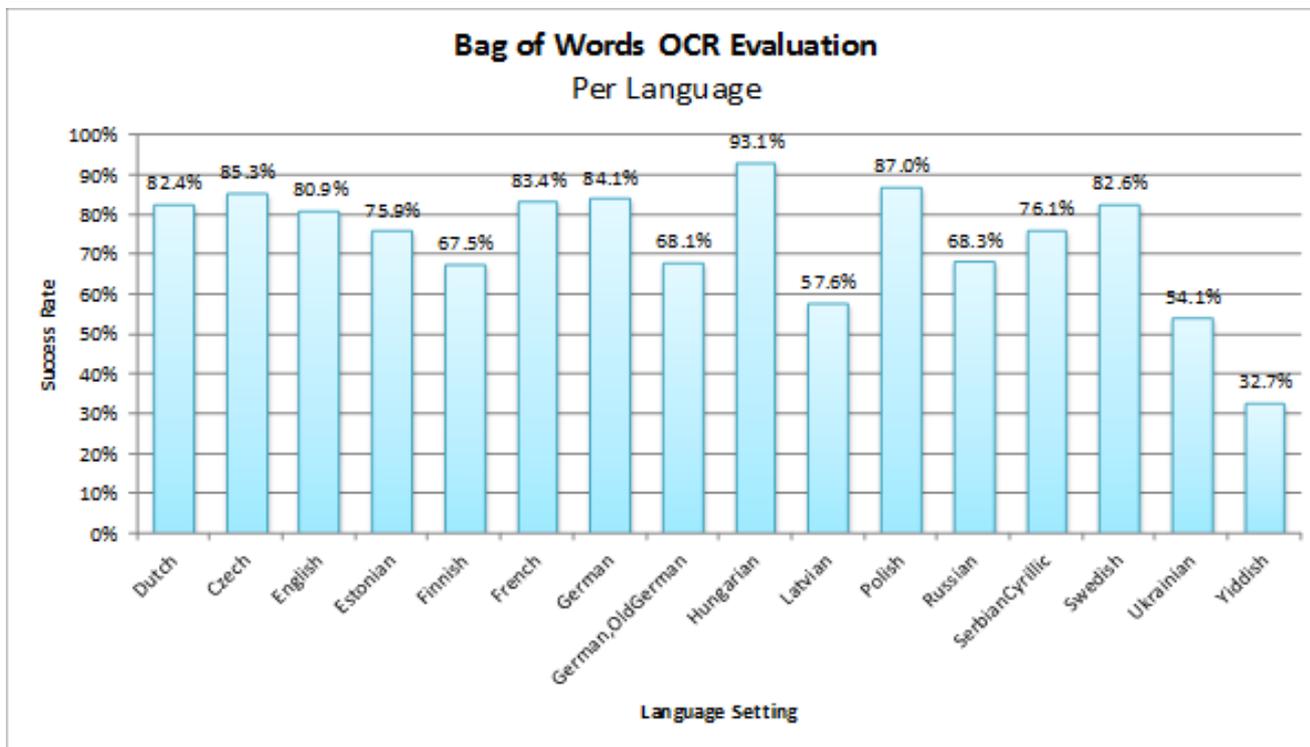


Empfindsamkeit
(ca. 1720-1800)
=
Sentimentalism

# What about further Services?

- Natural language processing
- Text mining
- Visualisations
- Cross-media linking
- Semantic field analysis
- Links to other resources, librarian and non-librarian
- …
- LIBERATE YOUR DATA AND LEARN FROM YOUR USERS!

# Digital Text Corpora: The Inconvenient Truth



Bag of Words OCR Evaluation Per Language

# What About Digital Text Corpora?

- Provide possibilities for corrections where data is presented
- Options for improvement
  - Automated corrections (index and page level)
  - Software aided corrections
  - Crowdsourcing
- Challenges: data synchronisation, update intervals, versioning…

europeana
newspapers

ICT PSP
ICT POLICY SUPPORT PROGRAMME

# Thank you for your attention!

IFLA Newspaper Pre-Conference

14 August 2014, Geneva

*Hans-Jörg Lieder, Staatsbibliothek zu Berlin – Preußischer Kulturbesitz | Berlin State Library*