# Building a national Infrastructure to enable research with Dutch Digital Heritage Collections

**Dr. Martijn Kleppe**
Board Member Research & Discovery
KB, national library of the Netherlands

martijn.kleppe@kb.nl
@martijnkleppe
https://www.kb.nl/over-ons/experts/martijn-kleppe

KB 〉 nationale bibliotheek

# Today

➢ What is the KB and what do we do with data for research?

➢ Case 1: FAIR Cultural Heritage data for research

➢ Case 2: Secure Data Analysis Environment

➢ Wrap up & some lessons

**7 million items
120 kilometres publications**

Full text (OCR) access to:

✓ 1.000.000 books (1486 – 2013)

✓ 1.700.000 newspapers (1618 – 1995)

✓ 10.000.000 periodical pages (1840 – 1940)

✓ 1.500.000 million ANP radio items (1937 – 1984)

https://www.delpher.nl/

www.delpher.nl



www.kb.nl/dataservices

# Dataservices

= result of

more than 200 years of collecting

over 30 years of digitisation

10+ years of collecting born-digital publications

= machine readable, mostly textual

= structured or semi-structured

= legally as open as possible



## Datasets of the KB, National Library of the Netherlands

The collections of the KB National Library of the Netherlands are being digitised at a large scale. The KB publishes books, periodicals, newspapers and other textual heritage freely accessible on the Web.

A significant amount is also made available in bulk for research purposes. Datasets, consisting of digital texts, images and metadata, can be accessed through www.kb.nl/dataservices.

### Characteristics
- extensive Dutch text corpora
- optical character recognition with word coordinates
- machine readable access
- documents in PDF and/or JPEG
- full text as XML
- metadata in Dublin Core and MPEG21 DIDL
- access via SRU and OAI-PMH

### Research projects
HiTiMe, BILAND, CLARIN, PoliticalMashup, PoliMedia, CATCH, SEALINC, Radicale Politieke Verbeelding, WHASP, and more...

### Datasets

**Medieval Illuminated Manuscripts**
11,000 images from 400 medieval manuscripts (manuscripts.kb.nl)

**Historical Newspapers**
8.5 million newspaper pages from the Netherlands and its former colonies, 1618-1995 (kranten.kb.nl)

**Early Dutch Books Online**
10,000 books from the Dutch-speaking region, 1781-1800 (www.earlydutchbooksonline.nl)

**Staten-Generaal Digitaal**
450,000 Dutch parliamentary papers, 1814-1995 (www.statengeneraaldigitaal.nl)

**Periodicals**
1.5 million pages from Dutch periodicals, 1850-1940 (tijdschriften.kb.nl)

**ANP Radiobulletins**
1.5 million typoscripts of radio bulletins, 1937-1984 (anp.kb.nl)

And more...

**KB Koninklijke Bibliotheek**
National Library of the Netherlands
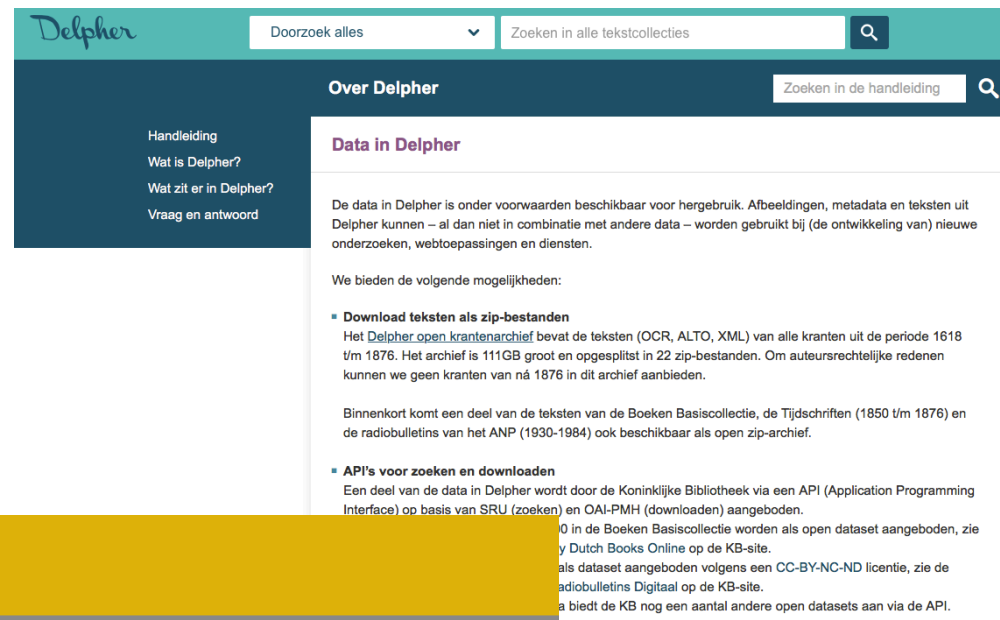
### Contact
dataservices@kb.nl
twitter: @sclaeyssens
www.kb.nl/dataservices

KB national library of the netherlands

We can provide data in several ways:

1. There are two APIs: a metadata harvest API on the basis of OAI-PMH, and a search API on the basis of SRU. Manuals for these APIs can be supplied once legal access has been granted via dataservices@kb.nl. Please note: users must have some experience of programming.
2. The Delpher newspapers comprise the texts (OCR, ALTO, XML) from all newspapers dating from 1618 to 1879. The archive is 111 GB and split into 23 ZIP files.

We can sometimes provide customised support. Ask your question via dataservices@kb.nl.

Delpher | Doorzoek alles | Zoeken in alle tekstcollecties | 🔍

**Over Delpher**

Zoeken in de handleiding 🔍

Handleiding
Wat is Delpher?
Wat zit er in Delpher?
Vraag en antwoord

**Data in Delpher**

De data in Delpher is onder voorwaarden beschikbaar voor hergebruik. Afbeeldingen, metadata en teksten uit Delpher kunnen – al dan niet in combinatie met andere data – worden gebruikt bij (de ontwikkeling van) nieuwe onderzoeken, webtoepassingen en diensten.

We bieden de volgende mogelijkheden:

- **Download teksten als zip-bestanden**
  Het Delpher open krantenarchief bevat de teksten (OCR, ALTO, XML) van alle kranten uit de periode 1618 t/m 1876. Het archief is 111GB groot en opgesplitst in 22 zip-bestanden. Om auteursrechtelijke redenen kunnen we geen kranten van ná 1876 in dit archief aanbieden.

  Binnenkort komt een deel van de teksten van de Boeken Basiscollectie, de Tijdschriften (1850 t/m 1876) en de radiobulletins van het ANP (1930-1984) ook beschikbaar als open zip-archief.

- **API's voor zoeken en downloaden**
  Een deel van de data in Delpher wordt door de Koninklijke Bibliotheek via een API (Application Programming Interface) op basis van SRU (zoeken) en OAI-PMH (downloaden) aangeboden.

... 0 in de Boeken Basiscollectie worden als open dataset aangeboden, zie ... y Dutch Books Online op de KB-site.

... als dataset aangeboden volgens een CC-BY-NC-ND licentie, zie de ... adiobulletins Digitaal op de KB-site.

... a biedt de KB nog een aantal andere open datasets aan via de API.

**http://data.bibliotheken.nl/**

De toegang tot alle Linked Open Data (LOD) zoals beschikbaar gesteld door de Koninklijke Bibliotheek. Alle data is beschikbaar onder de CC0-licentie. Zie de hulptekst voor tips en voorbeelden van het gebruik van deze data. Deze dienst is een bètaversie. Kijk bij Dataservices en API's voor meer datasets en -services.

**KB LAB** | Datasets | Tools | News and events | Blogs | About us | NL | Light ⬤ Dark

| Dataset | URI |
| --- | --- |
| lba amicorum van de Koninklijke Bibliotheek"@nl | http://data.bibliotheken.nl/id/dataset/rise-alba |
| rinkman trefwoordenthesaurus"@nl | http://data.bibliotheken.nl/id/dataset/brinkman |
| entsprenten"@nl | http://data.bibliotheken.nl/id/dataset/rise-centsprenten |
| emeenschappelijke Trefwoordenthesaurus (GTT)"@nl | http://data.bibliotheken.nl/id/dataset/gtt |
| ederlandse Bibliografie Totaal (NBT)"@nl | http://data.bibliotheken.nl/id/dataset/nbt |
| rganisaties uit de corporatiethesaurus van de Koninklijke Bibliotheek"@nl | http://data.bibliotheken.nl/id/dataset/corps |
| ersonen uit de Nederlandse Thesaurus van Auteursnamen (NTA)"@nl | http://data.bibliotheken.nl/id/dataset/persons |
| hort-Title Catalogue Netherlands (STCN)"@nl | http://data.bibliotheken.nl/id/dataset/stcn |
| hesaurus Auteurs DBNL"@nl | http://data.bibliotheken.nl/id/dataset/dbnla |
| hesaurus KBcode"@nl | http://data.bibliotheken.nl/id/dataset/kbcode |
| itels DBNL"@nl | http://data.bibliotheken.nl/id/dataset/dbnlt |

## Dataset

**Dutch Novels 1800-2000**

Dataset that contains a corpus of 1346 novels from DBNL.

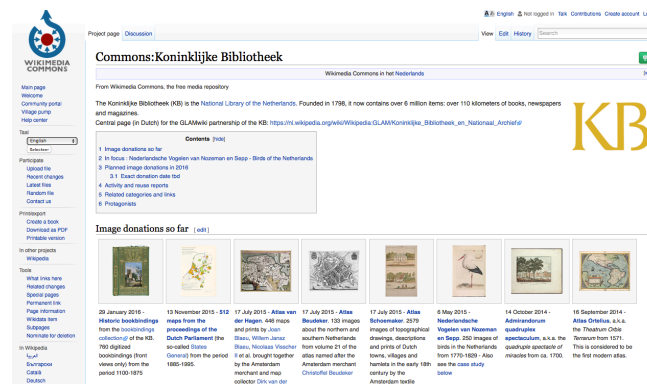**Accessible e-books and audiobooks**

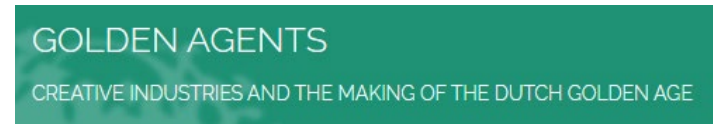We asked specialists to convert several public domain publications into accessible versions.

https://www.kb.nl/en/research-find/for-researchers/data-services-apis-and-downloads

**Commons:Koninklijke Bibliotheek**

From Wikimedia Commons, the free media repository

The Koninklijke Bibliotheek (KB) is the National Library of the Netherlands. Founded in 1798, it now contains over 6 million items: over 110 kilometers of books, newspapers and magazines.

Central page (in Dutch) for the GLAMWiki partnership of the KB.

**Image donations so far**

humanities and social sciences
**Mining Shifting Concepts Through Time (ShiCo)**

PoliMedia

BEELD EN GELUID AVResearcher XL

GOLDEN AGENTS
CREATIVE INDUSTRIES AND THE MAKING OF THE DUTCH GOLDEN AGE

**DIVE+**
*Dynamically Linking Collections on the Basis of Events*

**Time Capsule**
Making Pharmaceutical and Botanical Digital Heritage Accessible and Usable

**HiTiME**
Historical Timeline Mining and Extraction

**Translantis**
Digital Humanities Approaches to Reference Cultures: The Emergence of the United States in Public Discourse in the Netherlands, 1890-1990

**PoliticalMashup**

KB ) national library / of the netherlands

http://ngramviewer.kbresearch.nl/

https://www.nederlab.nl/

https://mediasuite.clariah.nl/

How can we better fit the needs of Digital Humanities scholars

by co-designing the solutions together with the researchers?

**dialogic**
*innovatie · interactie*

## USER NEEDS FOR A TEXT SUITE FOR ADVANCED DIGITAL RESEARCH

### Recommendations

**We conclude there are sufficient opportunities for developing a text suite.** For this development, we provide the following recommendations:

1. **Position a text suite as a corpus selection tool and support the discovery and selection research phases.** A text suite hereby functions as a user-friendly front end to Dataservices with more advanced features that do not fit within Delpher and DBNL. Users can then make their own selection of sources and export them (possibly after approval by a KB employee).

dr. Max Kemman
ir. Nick Jelicic
Guido de Moor MSc
Marenne Massop MSc
ir. Tommy van der Vorst

**COMMISSIONED BY**
KB

**PUBLICATION NUMBER**
2022.024-2212

**DATE**
Utrecht, March 31st 2022

https://zenodo.org/record/6591572

KB national library of the netherlands

# Case 1
# FAIR Cultural Heritage Data

And not only KB data but data from different heritage institutes

KB 〉 nationale bibliotheek

# Network Digital Heritage

The Dutch Digital Heritage Network (NDE) aims **at increasing the social value of the cultural heritage information** maintained by libraries, archives, museums and other cultural institutions.

The NDE strategy starts from the **end user perspective** and encourages institutions to provide digital heritage information that is more **visible, usable and sustainable**.

The NDE program is about building strong **cross sector networks** on the level of **expertise and information**. **Linked Data** is regarded as one of the enabling technologies.



National Digital Heritage Strategy

March 2021

This is a publication of the Dutch Digital Heritage Network and the Ministry of Education, Culture and Science

https://netwerkdigitaalerfgoed.nl/en/

KB ) national library of the netherlands

# Roadmap for the NDE discovery infrastructure

https://datasetregister.netwerkdigitaalerfgoed.nl/?lang=en

# CLARIAH

CLARIAH develops, facilitates and stimulates the use of Digital Humanities resources and infrastructures. We offer these resources to researchers and other professionals in an insightful and user-friendly way.



https://www.clariah.nl/

# Data

= result of

   more than 200 years of collecting

   over 30 years of digitisation

   10+ years of collecting born-digital publications

= machine readable, mostly textual

= structured or semi-structured

= legally as open as possible

# Data

= result of

   more than 200 years of collecting

   over 30 years of digitisation

   10+ years of collecting born-digital publications

= machine readable, mostly textual

= structured or semi-structured

= **legally as open as possible**

# Current workflow

| Researcher signs contract with KB | → | Researcher harvests data | → | Researchers does what he/she wants |
|---|---|---|---|---|

Challenges for both researcher as well as KB:

- More data

- Better & richer data

- More recent data (legal issues)

https://www.clariah.nl/projects/tools-to-data

# Future workflow



Sandbox

1. Data provider authorises researcher → 2. Researcher selects data → 3. Researcher develops/chooses algorithm → 4. Researcher uploads algorithm

5. Data provider approves access → 6. Algorithm runs in secure environment → 7. Results checked by data provider → 8. Results available to researcher

WORK IN PROGRESS

https://zenodo.org/record/7254517

# Wrap up

How can we better fit the needs of Digital Humanities scholars

by co-designing the solutions together with the researchers?

KB 〉 nationale bibliotheek

# Some lessons
## COLLABORATE

- ➢ Collaborate with researchers

- ➢ Collaborate with other heritage institutes

- ➢ Collaborate with research IT partners

- ➢ Build generic solutions. There is no 1 size fits all

- ➢ Make access to data as simple and secure as possible

# Questions?

## Building a national Infrastructure to enable research with Dutch Digital Heritage Collections

*Thanks fantastic colleagues*

Marian Hellema

**Dr. Martijn Kleppe**
Board Member Research & Discovery
KB, national library of the Netherlands

martijn.kleppe@kb.nl
@martijnkleppe
https://www.kb.nl/over-ons/experts/martijn-kleppe

Enno Meijers

Steven Claeyssens

KB ⟩ nationale bibliotheek