国立国会図書館
National Diet Library, Japan

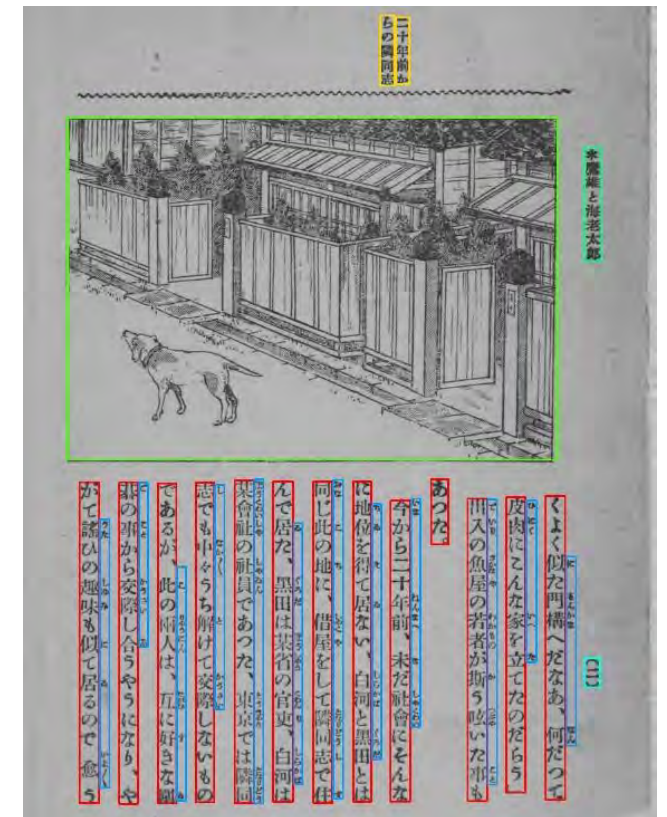# New Features of Digital Services at the NDL

**Onuma, Tahee**

**Assistant Director, Digital Information Planning Division**

**National Diet Library, Japan**

lab@ndl.go.jp

# Impact of Machine Learning Technology

- The development of machine learning/AI technology has greatly expanded possibilities for information exploration.
  - Data mining
  - Natural language processing
  - Computer vision

  ...

- Significant recent advances at the NDL:
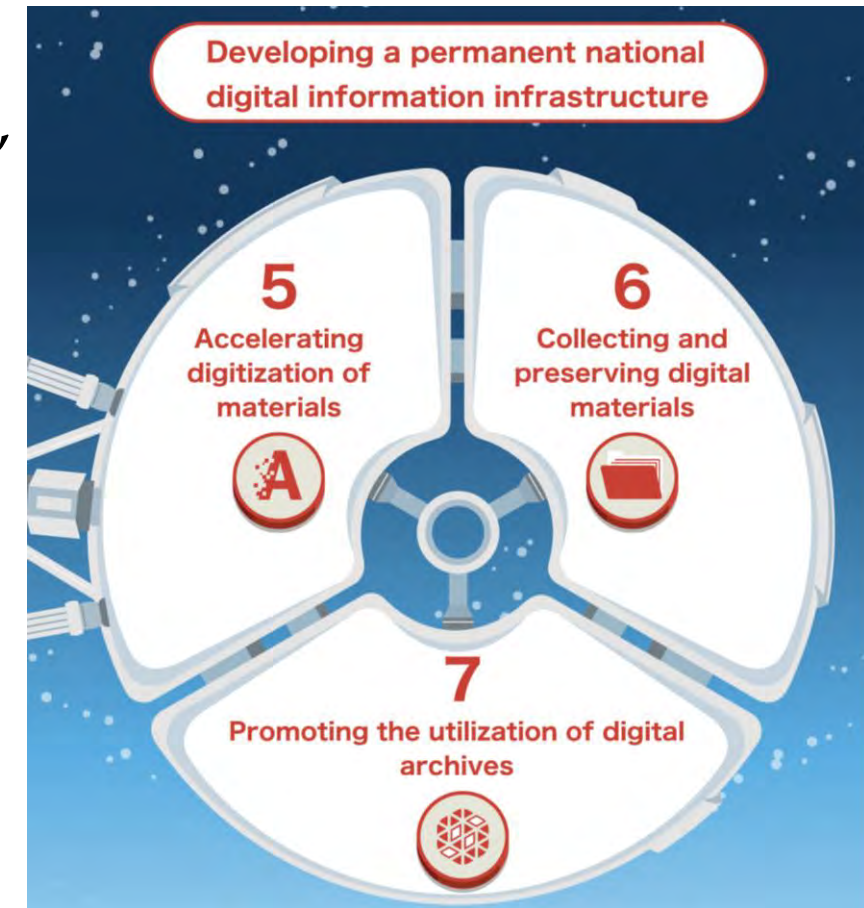  - Services using **full-text data**
  - **Image search** functionality

# 1. Use of Full-Text Data

- *Vision 2021-2025*

  Initiative 5: "**Accelerating digitization of materials**," which also includes <u>full-text conversion</u>

- Benefits of full-text data

  - For search purposes

  - Datasets for machine learning



Developing a permanent national digital information infrastructure

5 Accelerating digitization of materials

6 Collecting and preserving digital materials

7 Promoting the utilization of digital archives

# Challenges in OCR for Modern Japanese Materials

1. **Complexity of writing system**
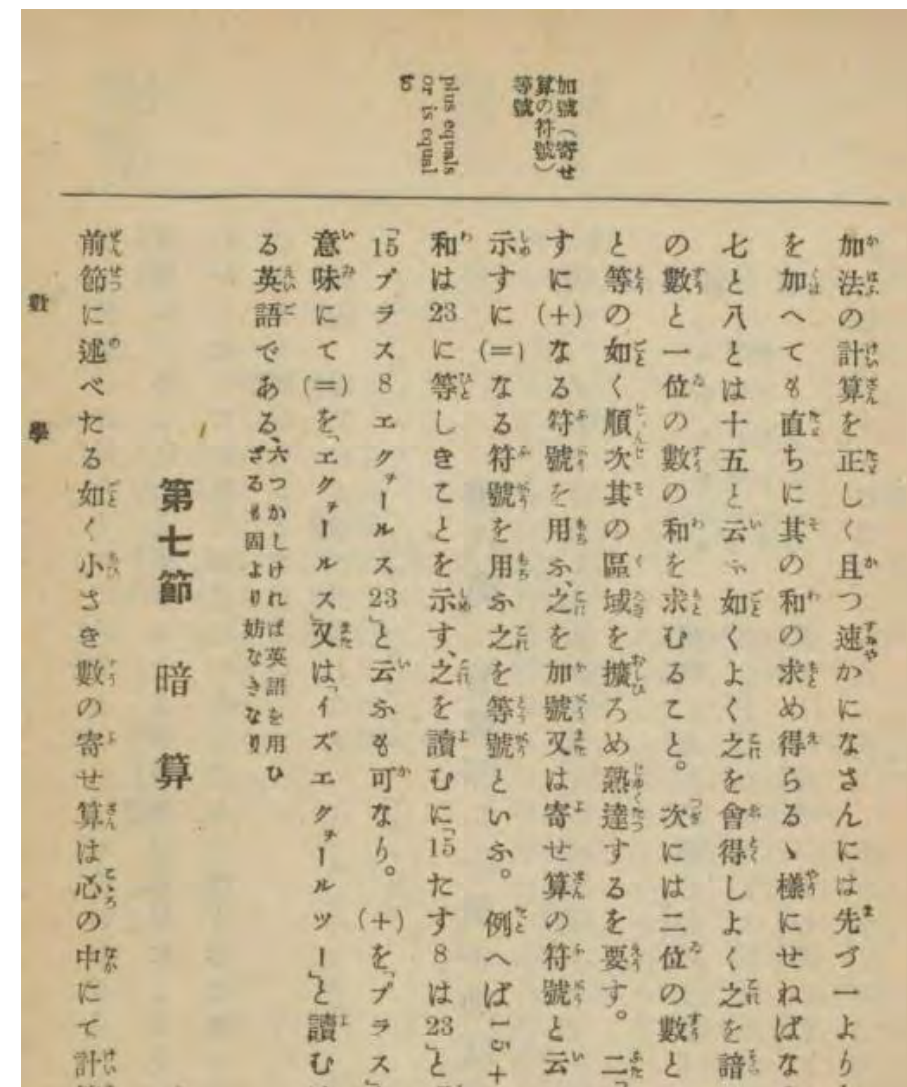
   - **Kanji** (漢字), **hiragana** (ひらがな), **katakana** (カタカナ)

   - Latin alphabet and other signs (abc, 123, #$%&...)

2. **Mixture of reading direction**

   - Vertical or horizontal, sometimes mixed (!)

3. **Some typographic features**

   - Ruby characters (ルビ)

   - Headnotes (頭注), split annotations (割注)

   …

4

# OCR-related Projects in FY2021: Overview

➢ **Project 1: Mass conversion from digitized images into text data**

- Target: **2.47 million items** (223 million images)

  = almost **all** documents that had been digitized as of 2020

➢ **Project 2: Development of AI-OCR software for Japanese materials**

- To be used for the text conversion of materials digitized after 2021

- To be made freely available as open-source software (**OSS**)

# OCR-related Projects in FY2021: Challenges

- How to generate full-text data for massive amounts of digitized images (*cost, time*)

- How to deal with characteristics of and variety in pre-WW II (1868-1945) Japanese-language materials, such as outdated character forms or fonts (*quality*)

- Development of an AI-OCR model optimized for digitized materials of the NDL

# OCR-related Projects in FY2021: Output

1. **A huge amount of full-text data**

   - Books, periodicals and other digitized collections (uncorrected)

   - Integrated into **NDL Digital Collections** for full-text search

2. **AI-OCR software NDLOCR**

   - Customizable and trainable OSS

   - Freely available on GitHub (CC-BY)

   For more details of these OCR projects, please see:
   https://lab.ndl.go.jp/data_set/ocr_en/

# Summary: Digitized and OCRed Materials at the NDL

To be OCRed from FY2023 by NDLOCR (0.58 million items)

1800    1868    1900    1969    1987    1995    2000

Pre-modern materials (80,000 items)

**Books (0.97 million items)**

OCRed in FY2022 experimental project (mentioned later)

OCRed in FY2021

**Periodicals (1.33 million items)**

: Digitized    : Prepared for digitization from FY2023

# NDL Digital Collections: Overview

- Platform for digital materials collected by the NDL
  - Digitized documents (3.33 million items)
  - Online publications (1.46 million items)
  - Doctoral dissertations in digital format (0.1 million items)

- 3 levels of access:
  - Available online (mostly PD)
  - Available through the Digitized Contents Transmission Service (see next slide)
  - On-premises only

https://dl.ndl.go.jp/

# Digitized Contents Transmission Service

- Provides online access to digitized out-of-print or otherwise difficult-to-obtain materials (including those under copyright protection)

i. **For libraries** (January 2014-)

- In partner libraries, patrons can access materials at terminals on premises

ii. **For individuals** (May 2022-)

- Official registered users (<u>Japanese residents only</u>) can access them via internet



Digitized materials: 3.33 million items

1.84 million items

i. For libraries

ii. For individuals

Out-of-print / Difficult-to-obtain materials

10

# Digitized Contents Transmission Service (for Libraries)

- **List of overseas partner libraries** (as of March 1, 2023)

| Country | Libraries |
|---|---|
| United States of America | C.V. Starr East Asian Library, University of California, Berkeley ↗ 【View Only】<br>University of Iowa Libraries ↗ 【View Only】 |
| Belgium | Library of Arts and Philosophy - Ghent University ↗ 【View Only】 |
| Ireland | University College Cork Library ↗ 【View Only】 |
| Italy | Istituto Giapponese di Cultura Biblioteca (ローマ日本文化会館図書館) ↗ 【View Only】 |
| Spain | Biblioteca de Humanidades, Universidad Autónoma de Madrid ↗ 【View Only】 |
| Israel | Bloomfield Library for the Humanities and Social Sciences, the Hebrew University of Jerusalem ↗ 【View Only】 |

https://dl.ndl.go.jp/en/soshin_librarylist

# NDL Digital Collections: Recent Progress

- Acceleration of digitization (as of February 2023)
  - Available online: 0.58 million items
  - Available through transmission service: 1.84 million items
  - On-premises: 0.91 million items

- Totally renewed in December 2022 with new functionalities:
  - **Full-text search**
  - **Image search**

cf.
"[Digitizing Library Materials at the NDL (part 1)](#)", *National Diet Library Newsletter*, (249).
"[Renewal of the National Diet Library Digital Collections](#)", ibid.

# NDL Digital Collections: Full-Text Search

- Text data produced in FY 2021 has been integrated. (2.47 million items)

- In the search results page, query keywords and their surrounding text are displayed as snippets.

- This functionality has led to a significant improvement in possibilities for information seeking. (e.g. first appearance of a specific term in the corpus)

# NDL Digital Collections: Full-Text Search

- Full-text search within a document is also available.
- Snippets are displayed for PD documents (as shown below).



https://doi.org/10.11501/1907107

14

# NDL Ngram Viewer

- Experimental analytic tool for the corpus of digitized **books** and **periodicals** (2.3 million items)

- Visualizes frequency of query keywords by publication year

- Useful features:
  - Supports **regular expression** search
  - Target corpus can be selected (only books, periodicals, or PD books…)



https://lab.ndl.go.jp/ngramviewer/

For more information, see:  https://lab.ndl.go.jp/service/ngramviewer/en/

15

# NDL Ngram Viewer

- Example: Comparison of two variants for "Harvard University"
  in Japanese writing ("ハーバード大学" or "ハーヴァード大学")

# 2. Use of Image Data

- Image-based search services offered by the NDL:

  - **NDL Digital Collections**
  - **Japan Search**

# NDL Digital Collections: Image Search

- Image search functionality: released in December 2022

- Queries can be launched by:
  - Uploading an image file
  - Referring an image by URI

# Japan Search: Overview

- Official release in August 2020
- National platform for promoting the use of content from Japan's digital archives (such as GLAMs)
- Includes various types of digital materials (images, videos, sounds, 3D data…)
- 3 main functions: *Search, Galleries, API* (SPARQL)



https://jpsearch.go.jp/

19

# Japan Search: Image Search

- Image search can be launched from any thumbnails shown

- Feature-based, not metadata-based

# Next Digital Library: Overview

- Experimental search/retrieval system for field trials of new functionalities developed by the R&D Office (March 2019-)

- Contains **books** and **<u>pre-modern materials</u>** in PD (350,000 in total)

21

# Challenges in OCR for Pre-modern Japanese Materials

**Cursive characters (崩し字)**

**Variants of hiragana (変体仮名)**



＝ つれづれなるままに

22

# OCR Experiment for Pre-modern Materials

- In FY2022, the R&D Office developed in-house an AI-OCR software specialized in pre-modern materials (mostly before 1868) and generated full-text data.

- The full-text search for all digitized pre-modern materials (80,000 items) is available.



https://doi.org/10.11501/2558997

# Next Digital Library: Full-Text Search

- In search results, snippets are displayed with query keywords highlighted.



**Search for "弥次郎兵衛"**
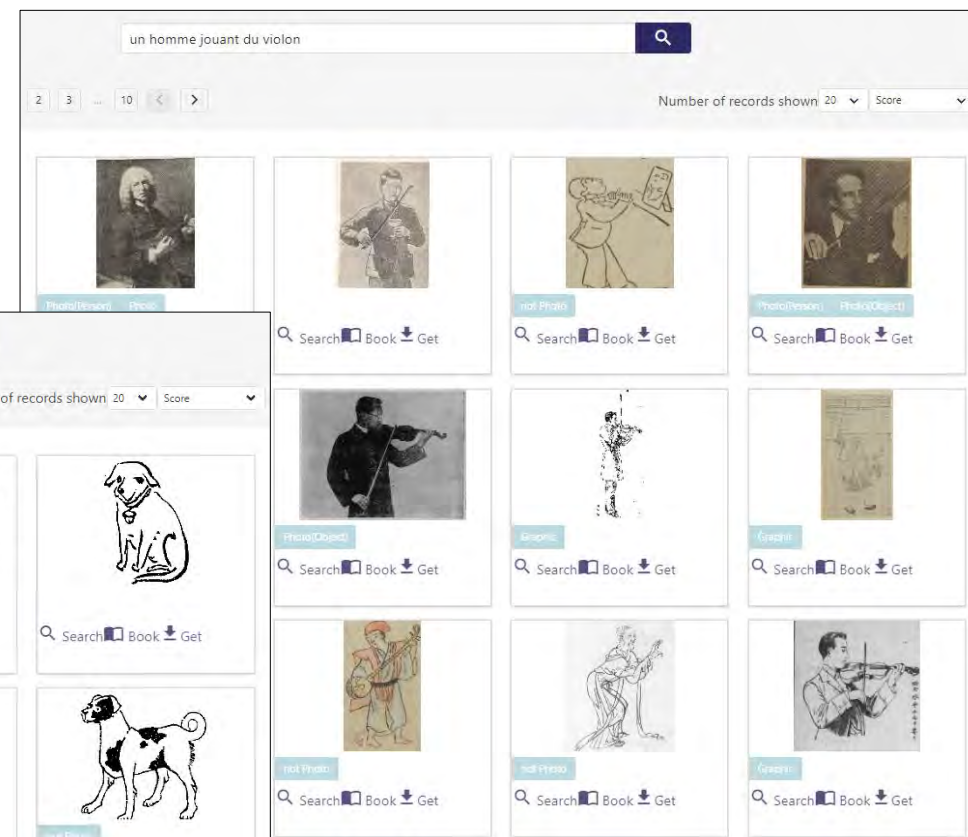
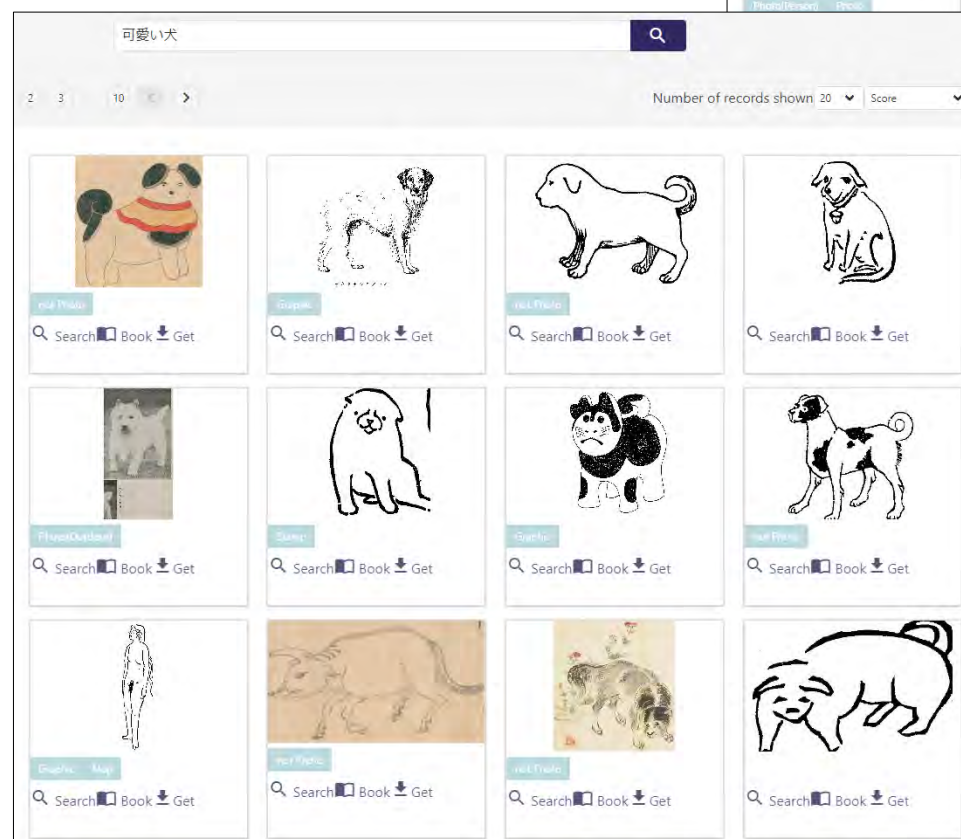# Next Digital Library: "Text-to-Image" Search

- Image search function by free text

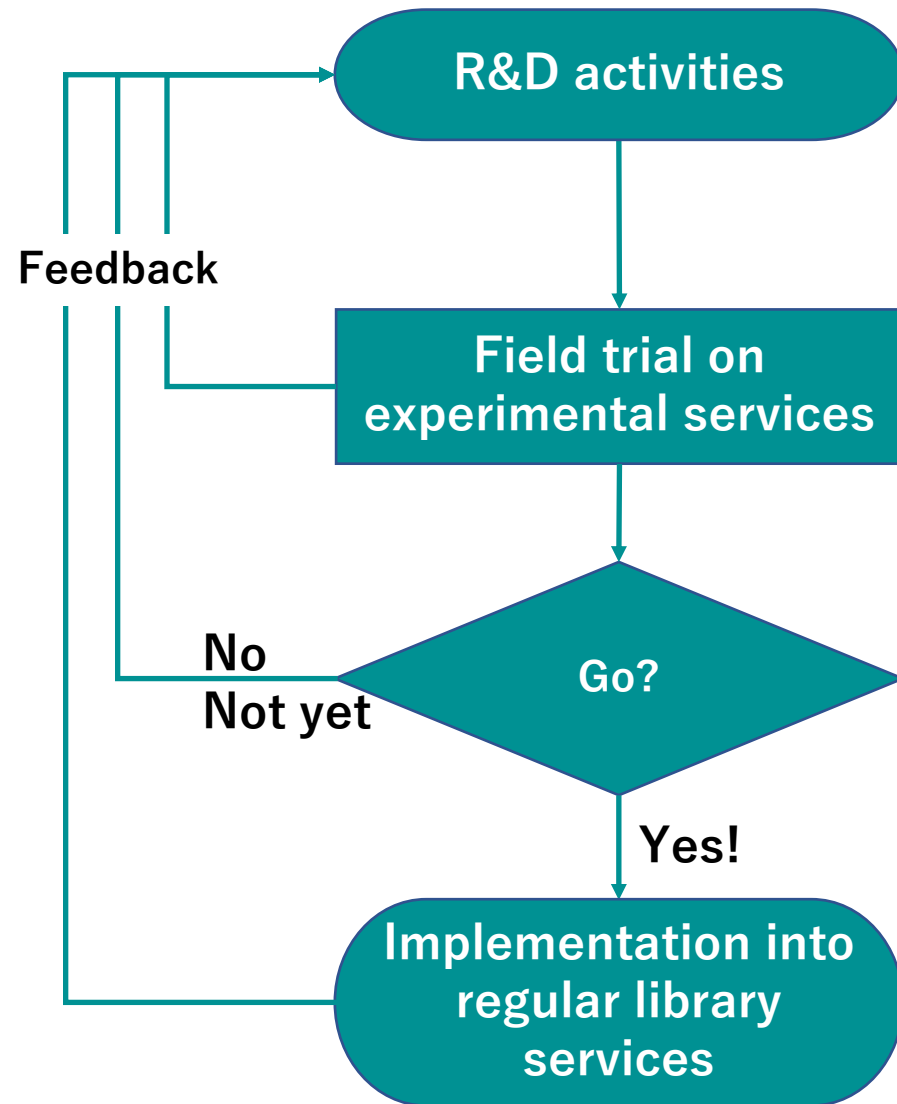- Supports multilingual queries by using machine translation
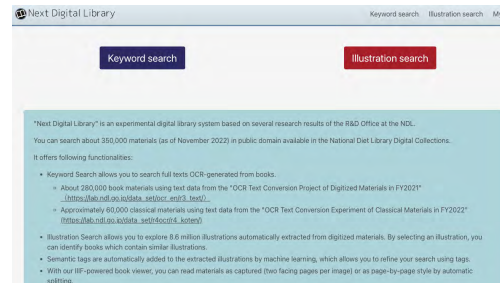
Examples of:

"可愛い犬" (cute dog(s))

and

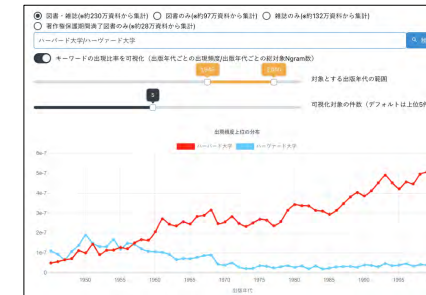"un homme jouant du violon" (a man playing the violin)

# Ecosystem Towards New Library Services



R&D activities

Feedback

Field trial on experimental services

Go?

No
Not yet

Yes!

Implementation into regular library services

**Next Digital Library**

**NDL Ngram Viewer**

...

**NDL Digital Collections**

**Japan Search**

...