October 24, 2024

2024 Taiwan Cultural Memory Bank International Forum
- Co-creation of Open Data in Museums

# Practice of Digital Archive Services Applying Machine Learning Technology

Toru AOIKE

Research and Development for Next-Generation Systems Office, National Diet Library

国立国会図書館
National Diet Library, Japan

# INDEX

- **About us**

- **Introduction**

- **Topic 1: Overcoming the Barriers of Language**

    Japan Search features that use machine learning

- **Topic 2: Overcoming the Barriers of Big Data**

    Creating and using large volumes of text data

- **Topic 3: Overcoming the Barriers of Time**

    In-house development of OCR for pre-modern materials

- **Summary & Future Activities**

# About us

Research and Development for Next-Generation Systems Office

This is a relatively new office, established at the National Diet Library (NDL) in 2011.

We are responsible for the research and development of new library services that use advanced information technology.

**Me !**

Office staff
1 office head, 1 chief, 1 staff member,
2 part-time staff members, 3 part-time researchers, and 1 associate member

We are a very small team, but we are tackling very interesting projects!

# Introduction

- Today's key phrase is **"overcoming barriers through technology"**.

- There are some barriers to the use of digital archives.

- We are exploring how machine learning and algorithms can overcome such barriers to make digital archives more usable.

- The three topics I will talk about today are all publicly available services on the Internet.

- I hope you will listen to what I have to say today while using and enjoying these services.
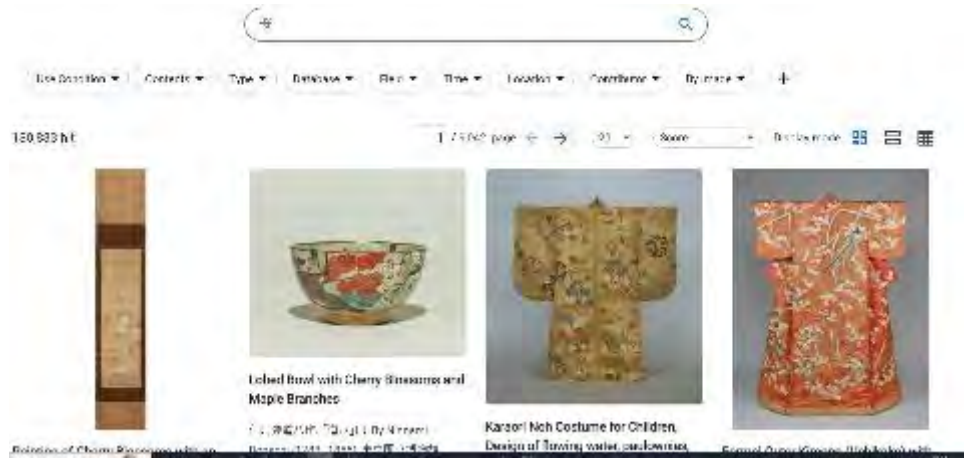
**Topic 1: Overcoming the Barrier of Language**

Japan Search features that use machine learning
Similar image search, Multi-modal search, and Visualization

# The Barrier of Language in Japan Search

Search Example:「桜」「櫻花」「Cherry Blossoms」



Japan Search results
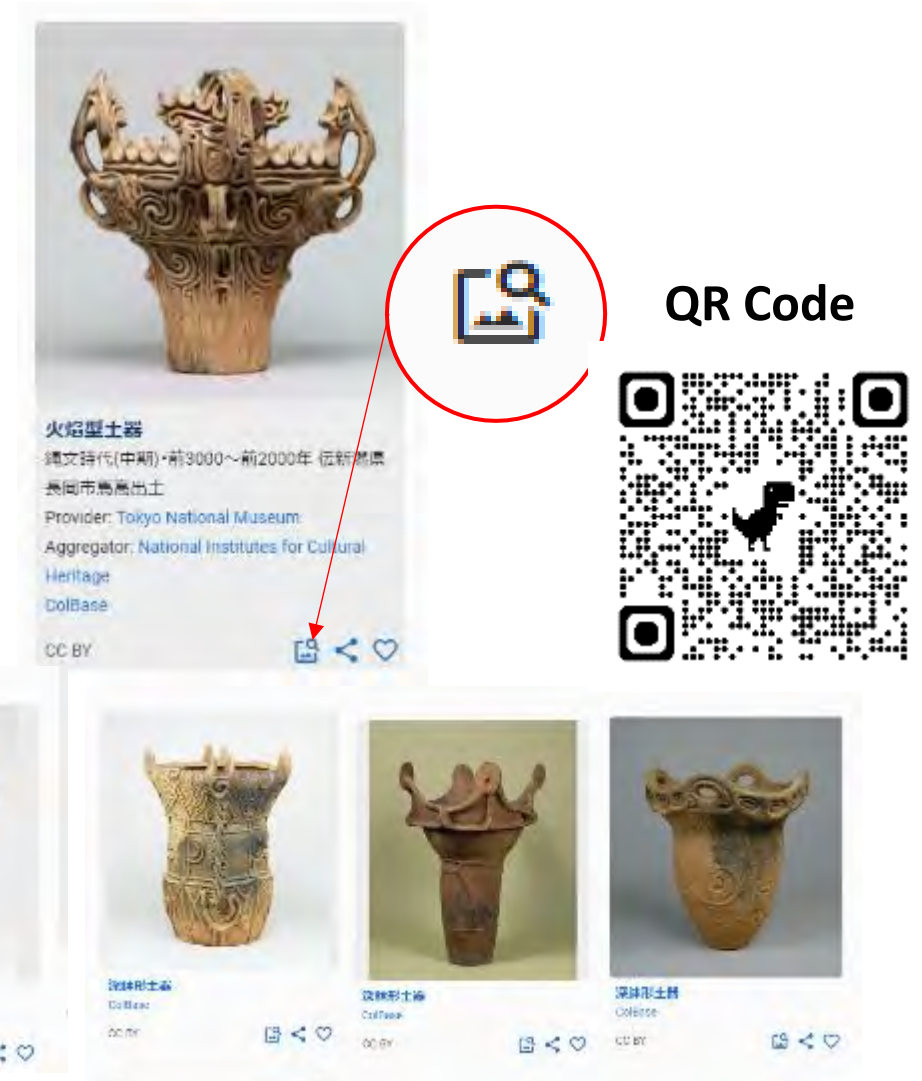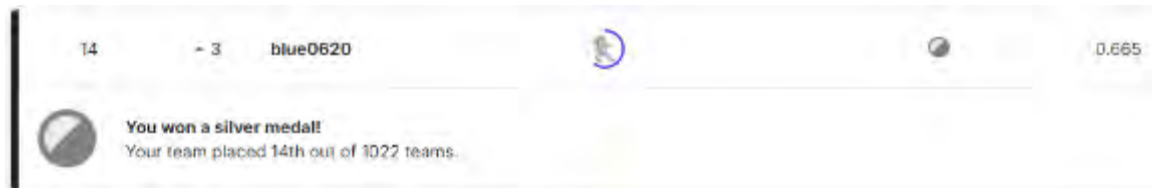桜：180,883 hits
櫻花：5,312 hits
Cherry Blossoms： 12,643 hits

This difference represents a limitation due to the linguistic information in the metadata. 8

# Overcoming the Barrier of Language

## Similar Image Search

This function searches for similar images based on the shape of objects in the image, without using query keywords.

This feature uses AI technology that I developed while participating in an international competition for image search held by Google Inc, in 2022, in which I placed 14th out of 1,022 teams.

**QR Code**

# Overcoming the Barrier of Language

## Multi-modal search

- Using an AI called ViT-CLIP, users can bridge text and image information to search for thumbnail images by keyword.

- Automatic language detection and machine translation enable multilingual search queries.

Multilingual search query:
「馬に乗った男性」
「騎馬的男人」
「A man on horseback」
「Homme à cheval」

QR Code

Nearly identical results, regardless of query language

10

# Overcoming the Barrier of Language

## Multi-modal search

Search for 雞肉飯 (chicken and rice)

Search for 香蕉 (banana)

# Overcoming the Barrier of Language

## Item Visualization Map (Visualization & Multi-modal search)

- In searches without strict keyword matching, such as similar image search and multimodal search, users are interested in the coverage of the search target.
- This service produces a single-screen, bird's-eye view of millions of thumbnail images on Japan Search, thereby providing users with a clear idea of the coverage available from multi-modal searches.
- Based on a modification of the source code for deepscatter (https://github.com/nomic-ai/deepscatter), available under CC BY NC

# Overcoming the Barrier of Language



## An example of visual exploration

## An example of multi-modal search

Search for 煙火（fireworks）

13

# Topic 2: Overcoming the Barrier of Big Data

Creating and using large volumes of text data

Our OCR project and the NDL Ngram Viewer

# Now we can consistently produce high-quality text data with OCR and provide full-text search

But how can we sift through the flood of information available from full-text data?

For example, searching the NDL Digital Collections for 台湾 (Taiwan) gives search results for nearly 1 million materials.

# Topic 3: Overcoming the Barrier of Time

In-house development of OCR for pre-modern materials
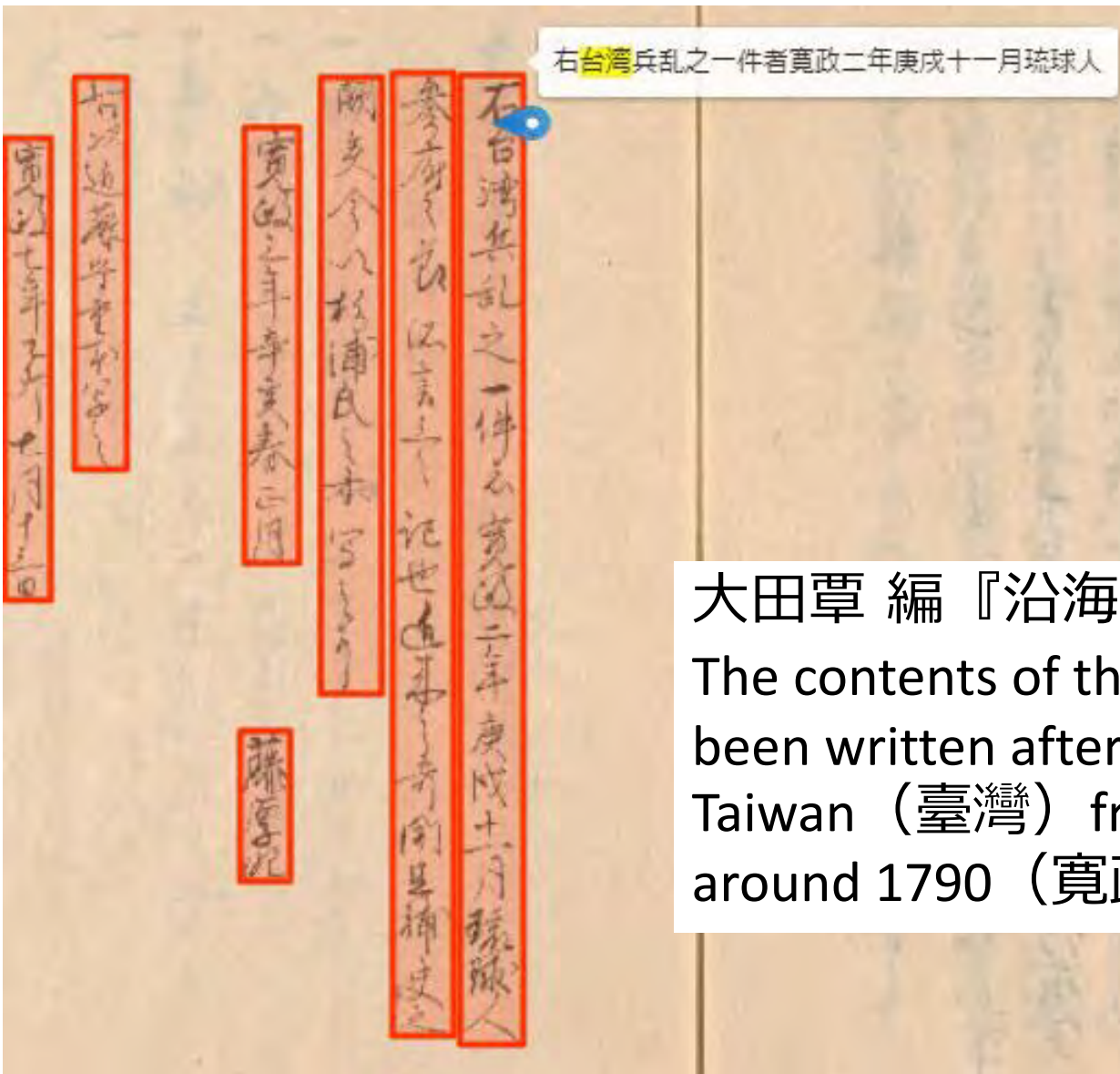NDLkotenOCR and the Next Digital Library

# Let's look for documents about Taiwan（台湾）in the full text data of pre-19th century materials!

https://lab.ndl.go.jp/dl/fulltext?from=0&keyword=台湾&fc-isClassic=true

# Example of search results



森嶋中良 編輯『紅毛雑話 5巻』
A passage describing a voyage through Taiwan （臺灣）
on the way to the equator around 1662（寛文2年）

開帆し。海上凡三百七十里余の波濤を過て。台湾に至り。

# Example of search results



右台湾兵乱之一件者寛政二年庚戌十一月琉球人

大田覃 編 『沿海異聞』[6]
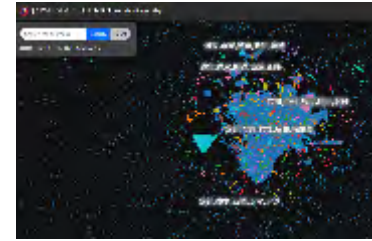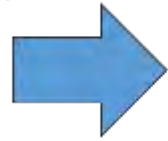The contents of this document are described as having been written after hearing about the domestic situation in Taiwan（臺灣）from the people of Ryukyu（琉球）around 1790（寛政2年）.

# Summary

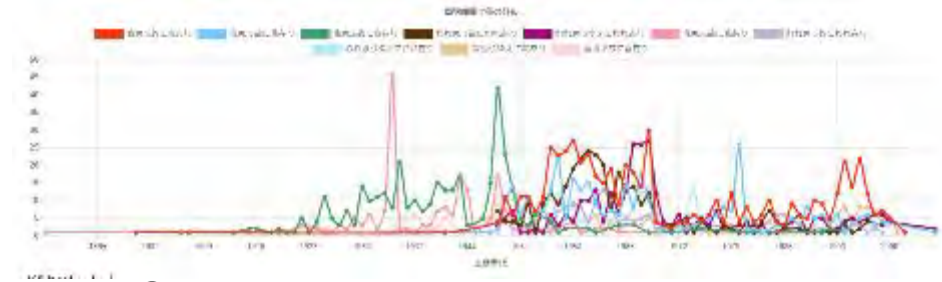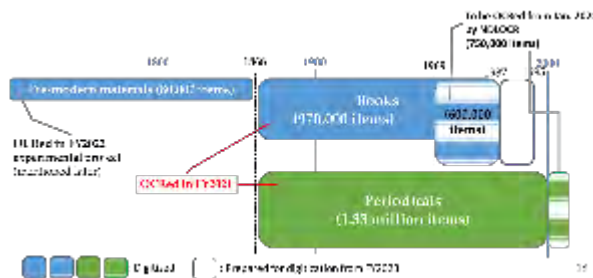- ## Topic 1: Overcoming the Barrier of Language



- ## Topic 2: Overcoming the Barrier of Big Data



- ## Topic 3: Overcoming the Barrier of Time

# Future Activities:

- The topics presented today are in the development stage and have yet to be perfected. It is important to consider better methods and to improve accuracy.

- Here are some new challenges that we are considering.
1. Use of generative AI for reference queries
2. Use of AI technology for video and audio materials