# Reading Edo: Data-driven Approaches for Japan Studies

Asanobu KITAMOTO

Director, ROIS-DS Center for Open Data in the Humanities (CODH) and

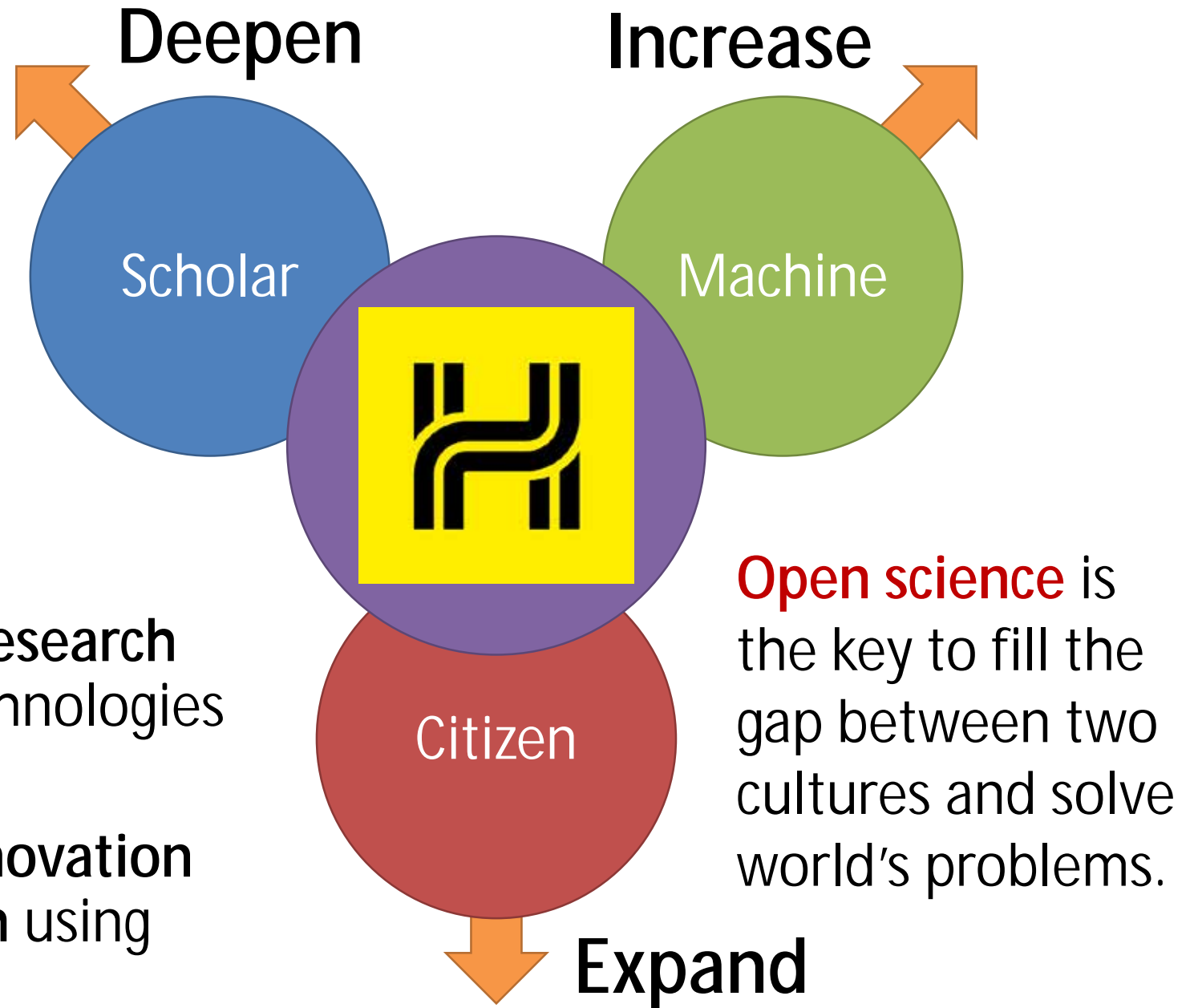Professor, National Institute of Informatics

http://codh.rois.ac.jp/ @rois_codh

# ROIS-DS Center for Open Data in the Humanities (CODH)

http://codh.rois.ac.jp/

**1. Data-driven Humanities**: **Innovation in humanities research** using computer science technologies and tools.

**2. Humanities Big Data**: **Innovation in non-humanities research** using humanities data.

**Deepen**

**Increase**

Scholar

Machine

Citizen

**Open science** is the key to fill the gap between two cultures and solve world's problems.

**Expand**

# Open Datasets
http://codh.rois.ac.jp/dataset/



Pre-modern Japanese Text (3126 books, 609631 pages)



Edo Cooking Recipes (103 egg dishes)

# NIJI-NW Project

http://www.nijl.ac.jp/pages/cijproject/index_e.html

**300,000 Pre-modern Japanese Books** (before 1868) are being digitized and released as open data.

Japanese culture finally entered into the big data era...
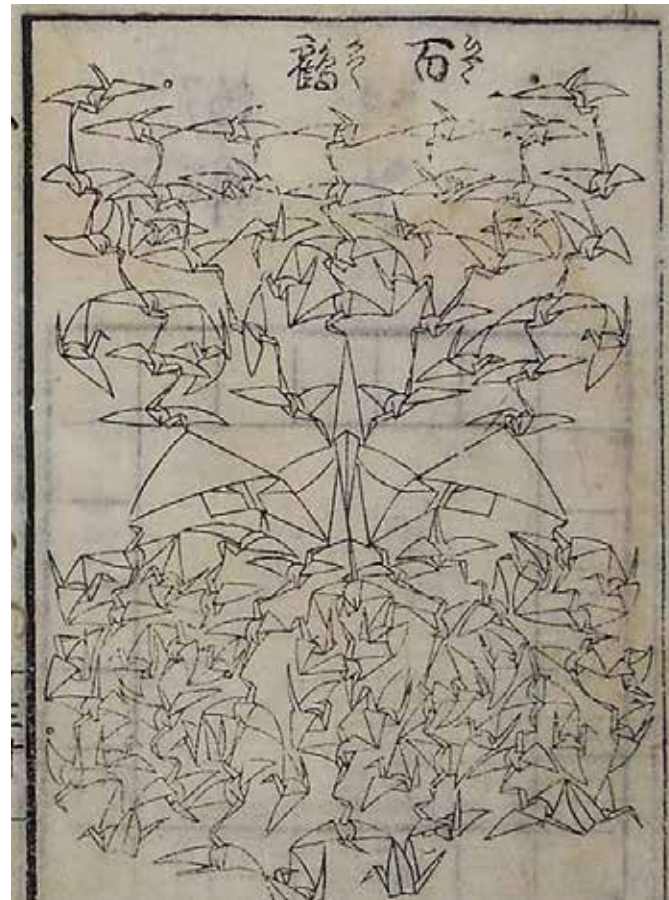
# AI Kuzushiji Recognition from the Dataset to the App

Collaborator: Tarin Clanuwat (Google Brain, formerly CODH)

# Japanese Knowledge over 1000 Years



How to wear makeup



How to fold 100 cranes using one piece of paper



How to build automata

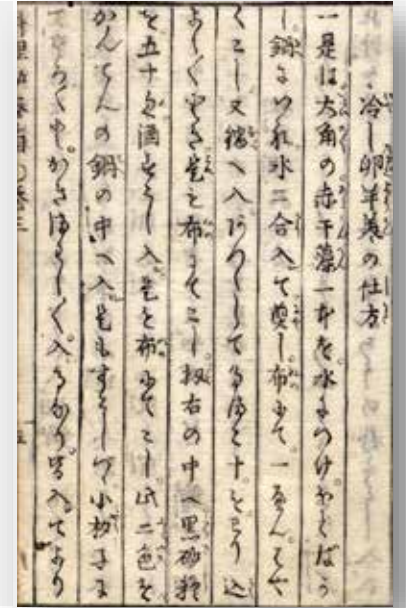# Massive Documents vs. Few Readers

**1 billion documents**

Estimated number of old books and documents in Japan
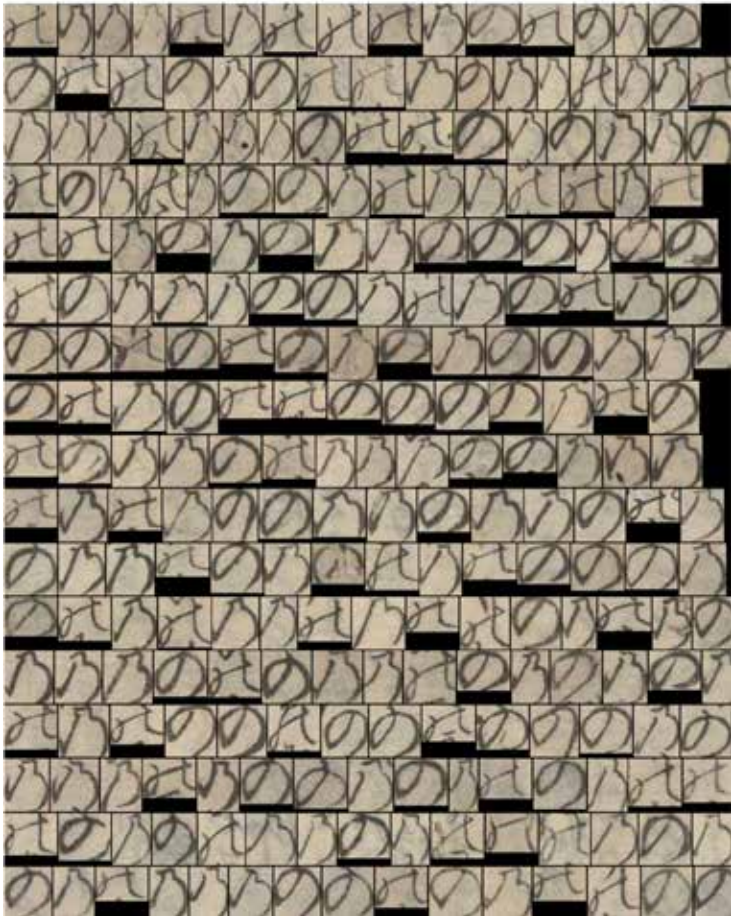
**10000 readers**

Estimated number of people with fluency in reading Kuzushiji
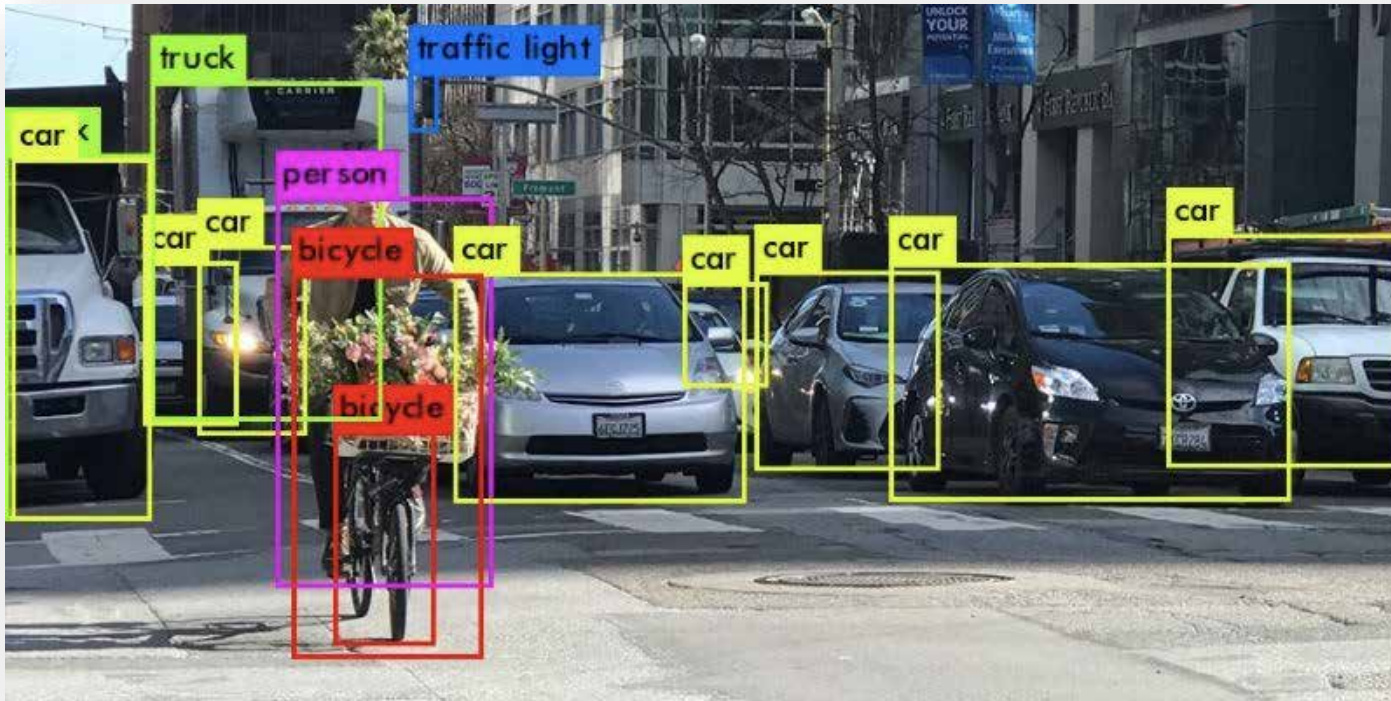
# Kuzushiji Dataset
[http://codh.rois.ac.jp/char-shape/](http://codh.rois.ac.jp/char-shape/)

雨月物語 (1890)



- **National Institute of Japanese Literature** created and **CODH** curated.
- The open data consists of
  - **Character types: 4,328**
  - **Character shapes: 1,086,326**
- Download the Zip file and use it as training data for machine learning.
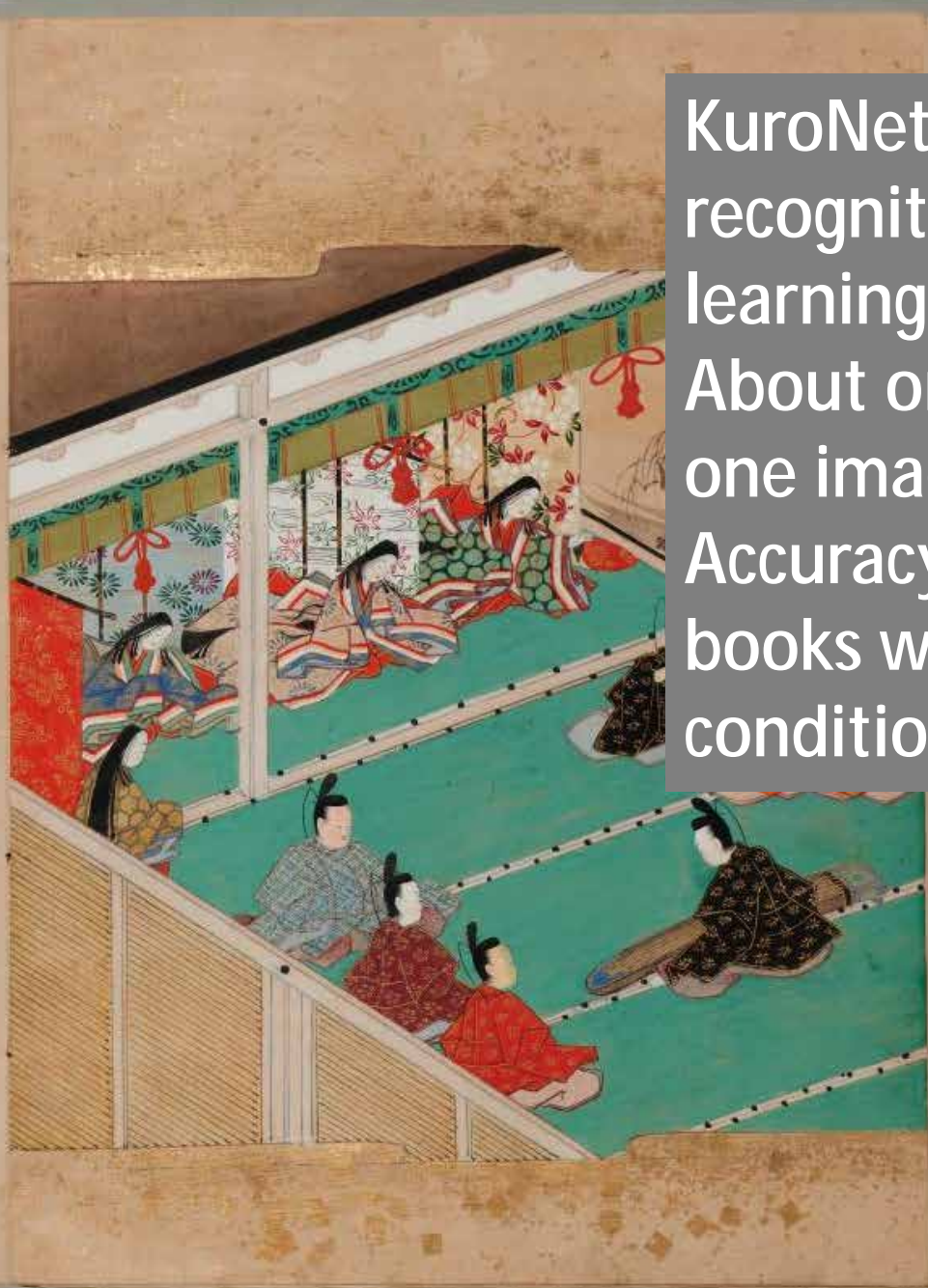- The release of dataset stimulated research on **AI kuzushiji recognition.**

# Computer Vision-based Object Detection



1. Object detection is a vibrant research area with industrial value such as autonomous driving.

2. **Can we apply this technology for kuzushiji?** A simple idea, but it was not possible before.

KuroNet: Kuzushiji recognition using deep learning.
About one second for one image.
Accuracy is 95% for books with the best condition.

# KuroNet Kuzushiji Recognition

[http://codh.rois.ac.jp/kuronet/](http://codh.rois.ac.jp/kuronet/)

# kaggle Kuzushiji Recognition

http://codh.rois.ac.jp/competition/kaggle/



**Kaggle** is the largest **AI competition** platform.
Our competition was the first in the humanities domain.

- Period: July 19 to October 14, 2019
- Teams: 293
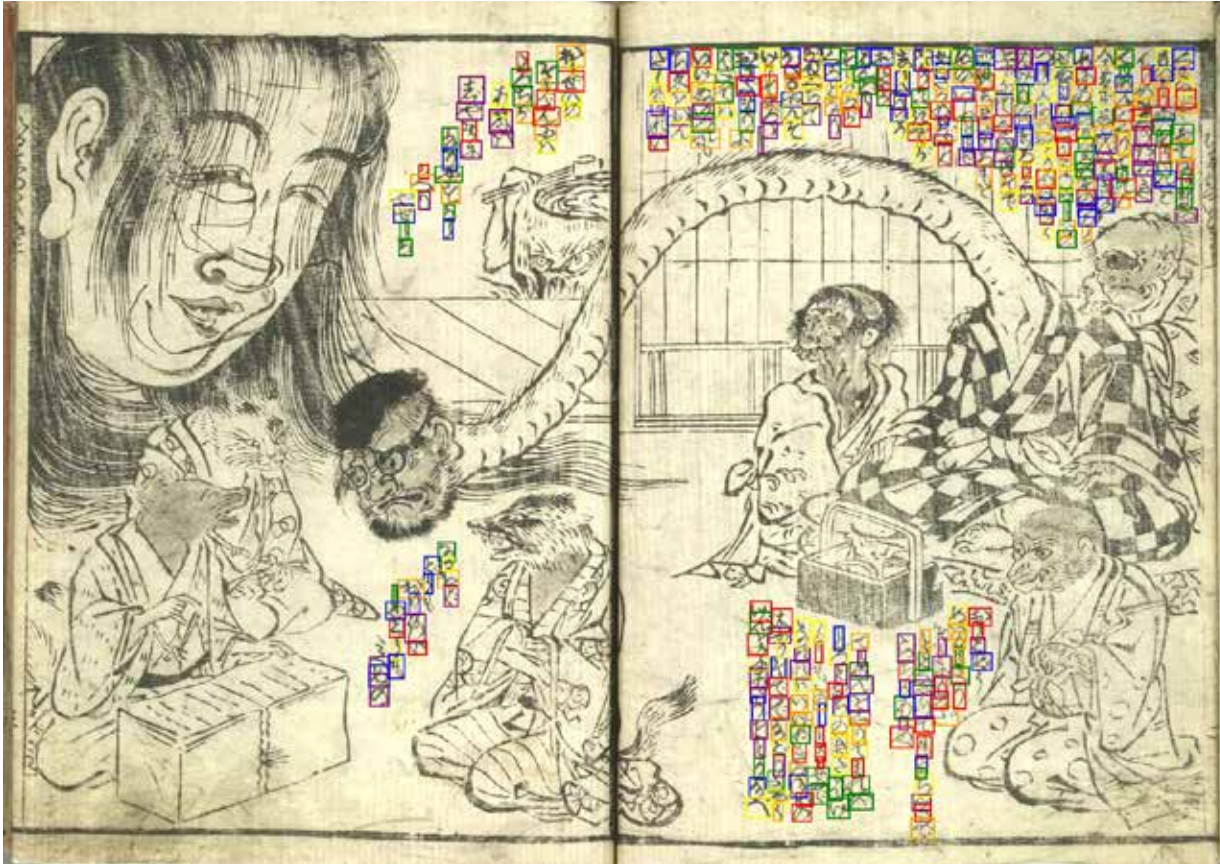- Members: 338
- Submissions: 2652

# Result of the Competition



The top accuracy was 95%

The winner model (tascj team) was applied to an image from Waseda University Kotenseki Database

1. All winners do not read kuzushiji, but have developed good machine learning models.
2. This is because the domain knowledge was embedded in the dataset.

# miwo – App for AI Kuzushiji Recognition

http://codh.rois.ac.jp/miwo/
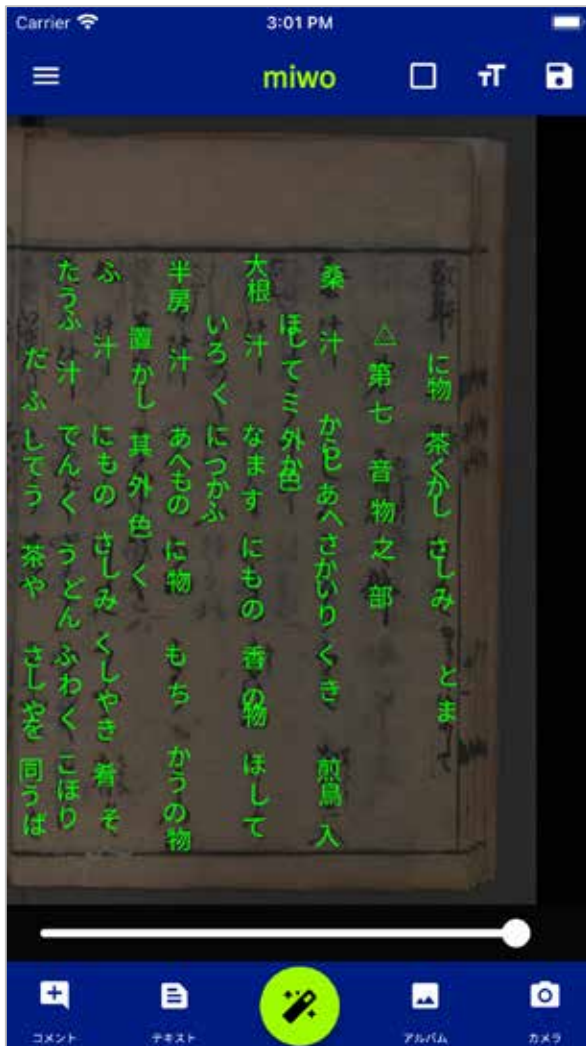
1. The name comes from the 14th chapter of The Tale of Genji "miwotsukushi," referring to waterway signs.

2. Just as the miwotsukushi is a guide for boats in the sea, we aim to make our "miwo" app as a guide for traveling the ocean of historical documents.

miwo app prototype version at the KeMCo Museum (April 2021)

Show a recognition result in characters

Show a recognition result with bounding boxes

Modify the error with reference to root characters.

Generate the text output from the recognition result

# Impact and Future of Kuzushiji Recognition

1. The miwo app was downloaded more than **42,000** times, and more than **337,000** images were recognized.

2. The daily uploaded images is constantly above **2,000** images, which indicates steady demand from the public.

3. Future of kuzushiji recognition is the full-text search engine of historical documents (we'll call it "**tsukushi**").

4. The full-text search engine will be the driver of digital transformation in the humanities research.

# Bukan Complete Collection

Collaborator: Kumiko Fujizane (National Institute of Japanese Literature)

# Textual and Non-textual Digital Humanities

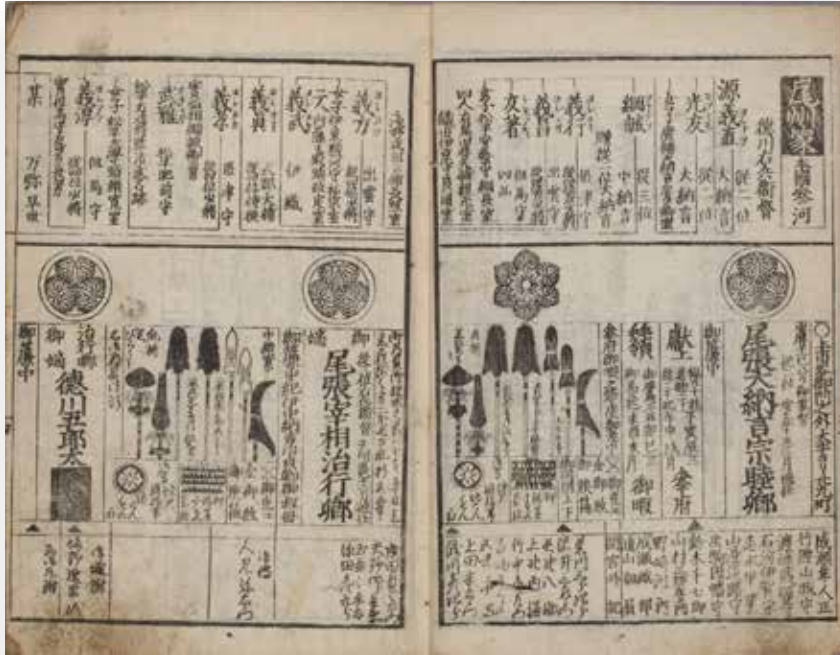**Images**  **Photographs**  **Maps**  **Characters**



**Digital humanities is not only about text.**

Structured and unstructured data (visual and spatial sources) requires its own analysis and interpretation framework.

# What is Bukan　　　　　？
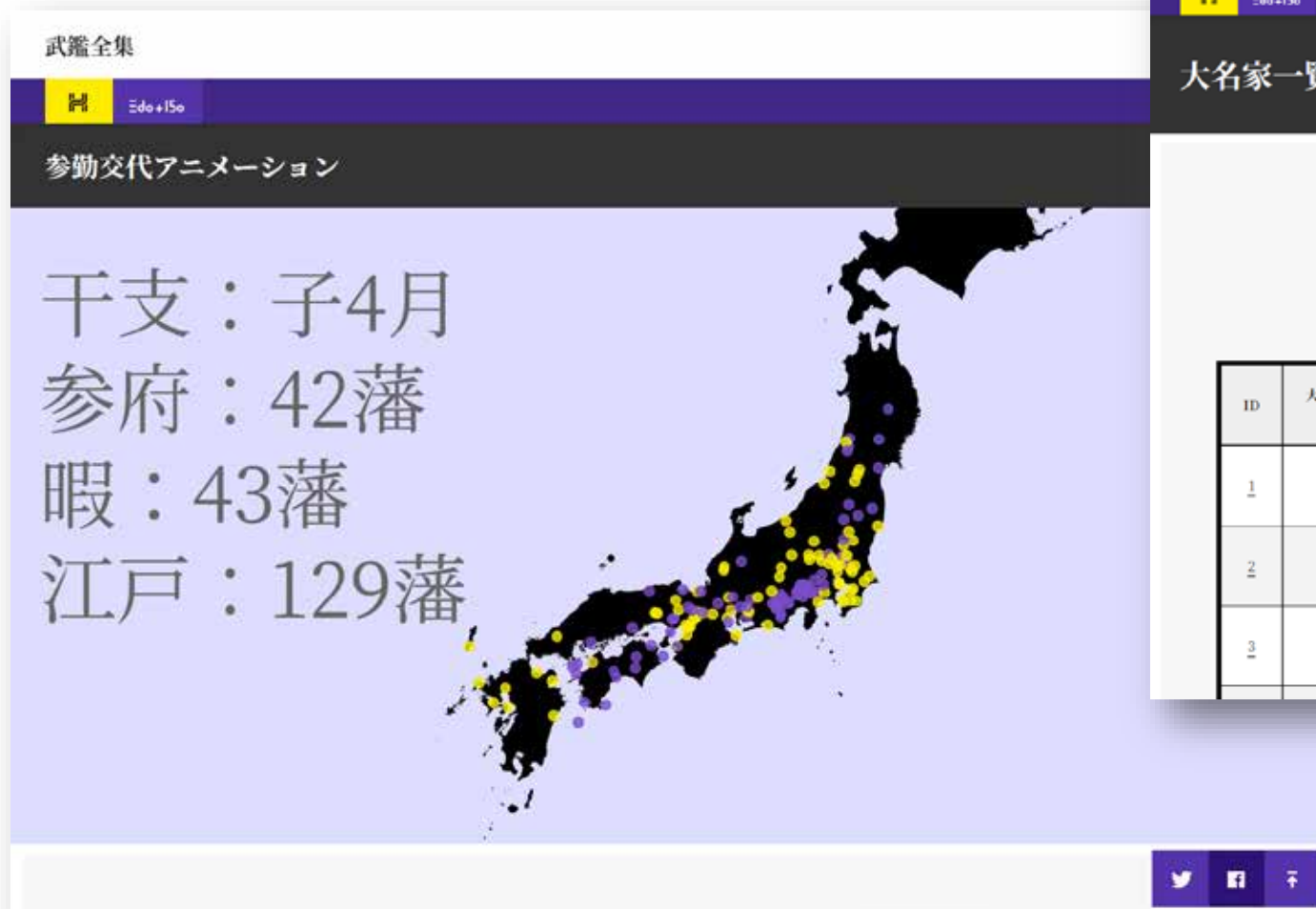


Kansei Bukan (1789), Dataset of
Premodern Japanese Text (NIJL)
http://codh.rois.ac.jp/pmjt/book/200018823/

1. Bukan is a "data book" of Daimyo and personnel in the Edo Bakufu compiled in a structured format.

2. Published for 200+ years before 1867, until the end of the Edo Period.

3. Long-seller books with practical usage.

4. The frequency of updates had increased to a few times a month at the peak.

Reference: Kumiko Fujizane, 2008

# Bukan Complete Collection

http://codh.rois.ac.jp/bukan/



武鑑全集

参勤交代アニメーション

干支：子4月
参府：42藩
暇：43藩
江戸：129藩



武鑑全集

大名家一覧

日本古典籍データセットで公開する寛政武鑑（1789）の大名家一覧です。IDは寛政武鑑（1789）での出現順に付与しています。

| ID | 大名当主名（現代通称） | 大名当主名（武鑑表記） | 藩名（現代通称） | 居城地（武鑑表記） | 領知高（単位：石） | 参勤交代年月（参府） | 参勤交代年月（暇） |
|---|---|---|---|---|---|---|---|
| 1 | 徳川宗睦 | 尾張大納言宗睦 | 尾張 | 尾州愛知郡名古屋 | 619,500 | 子寅辰午申戌 3月 | 丑卯巳未酉亥 3月 |
| 2 | 松平義裕 | 松平摂津守義裕 | 高須 | 濃州石津郡高須 | 30,000 | 子寅辰午申戌 4月 | 丑卯巳未酉亥 4月 |
| 3 | 徳川治貞 | 紀伊中納言治貞 | 紀州 | 紀州名草郡和歌山 | 555,000 | 丑卯巳未酉亥 3月 | 子寅辰午申戌 |

## List of Daimyos

## Sankin Kotai Dynamic Map

# Differential Transcription

Question: how can we transcribe books over 200+ years?

Solution: detect and transcribe the difference to create time-series data.

# Text-based and Image-Based Difference

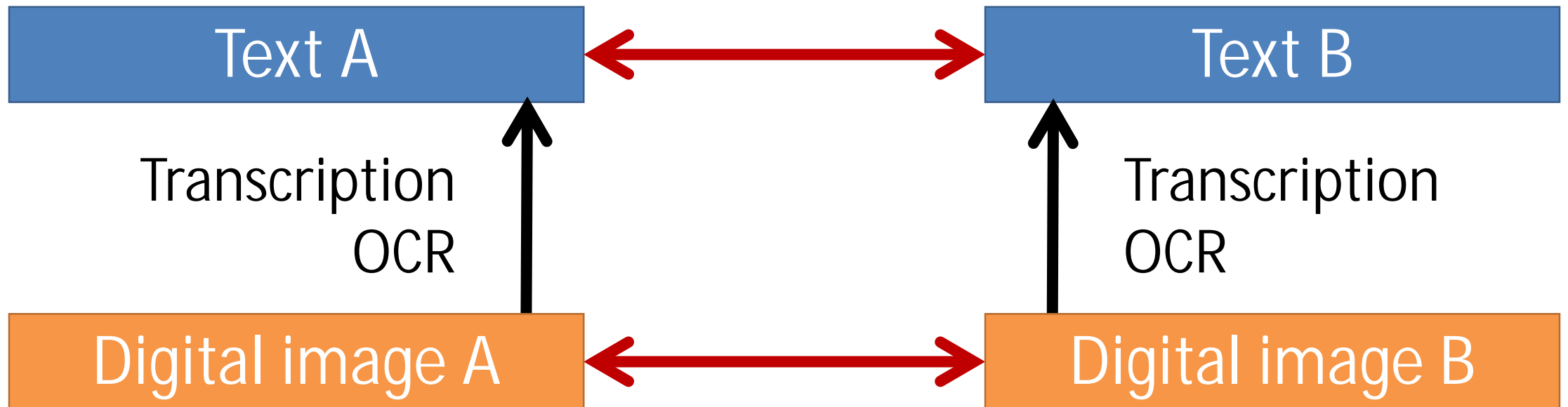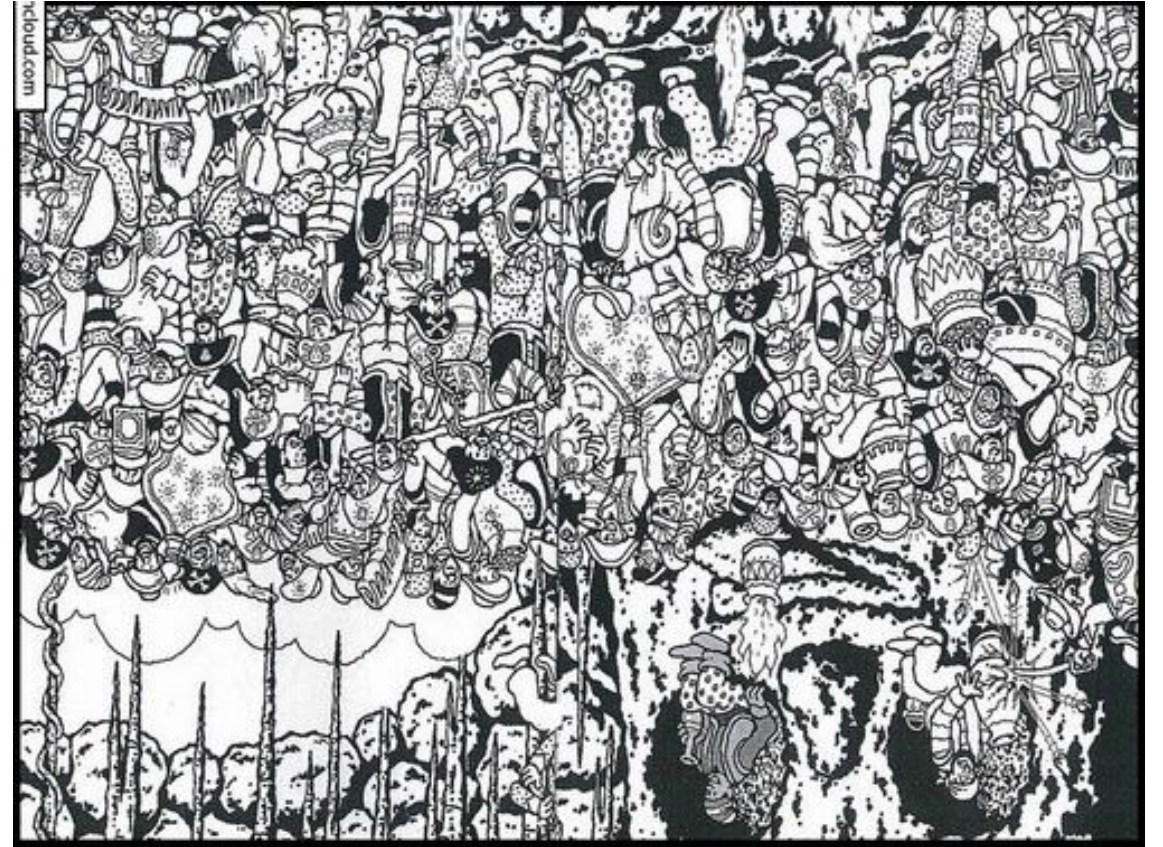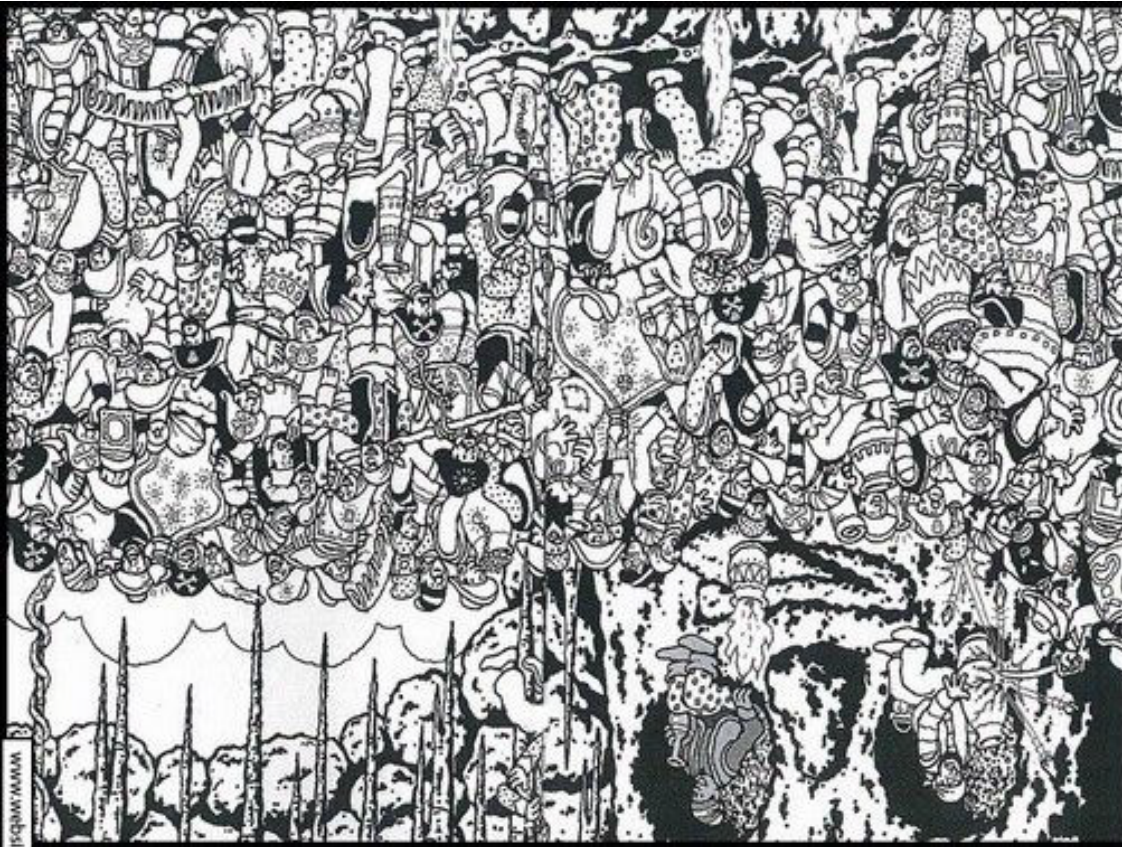**Text-based difference = <span style="color:darkred">many tools available</span>**

| Text A | ←——→ | Text B |
|---|---|---|

↑ Transcription OCR    ↑ Transcription OCR

| Digital image A | ←——→ | Digital image B |
|---|---|---|

**Image-based (non-textual) difference =
<span style="color:darkred">no standard tools available (side-by-side comparison)</span>**

# Visual Comparison = Find the Difference!



https://www.activities.websincloud.com/finddifferences/whereswally/21.html

# Answer

http://codh.rois.ac.jp/differential-reading/file/

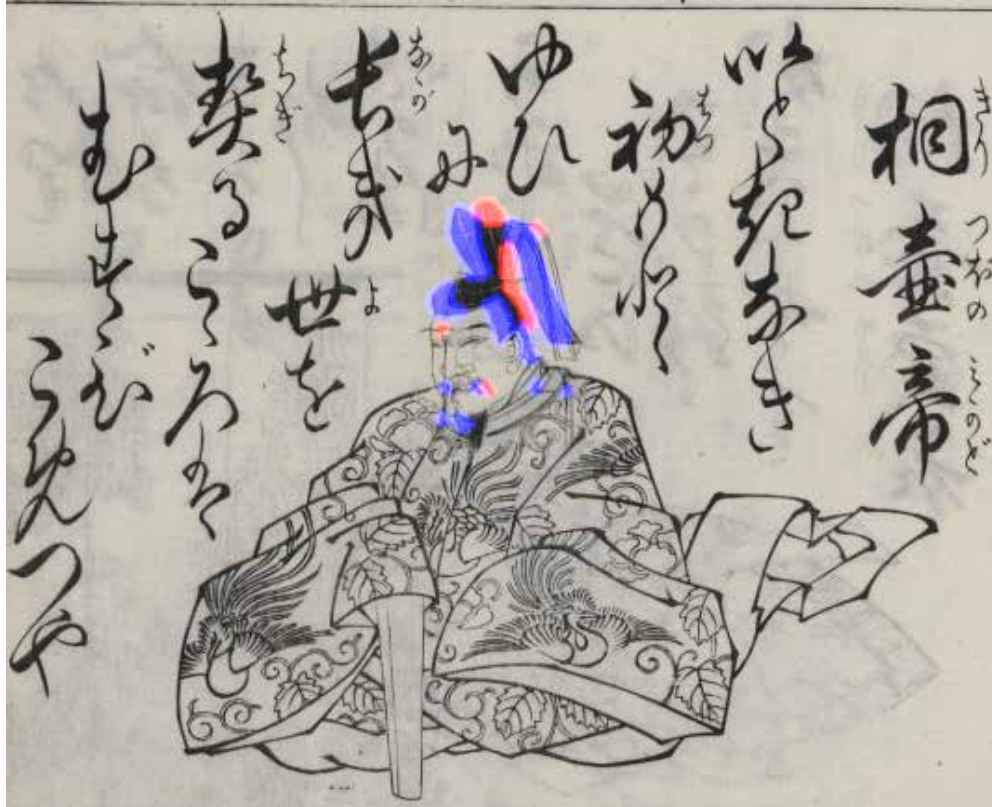Red and blue colors were used to emphasize the difference.

# Differential Reading

1.  **For humans**: visual comparison requires an effort comparable to playing games.

2.  **For machines**: visual comparison is an easy game using a computer vision-based image matching algorithm.

3.  Let's turn a difficult task (reading difference) into an easy one with the help of machines.

4.  **Differential reading**: A new mode of reading books focusing on difference between editions (versions).

# Image Collation for Differential Reading
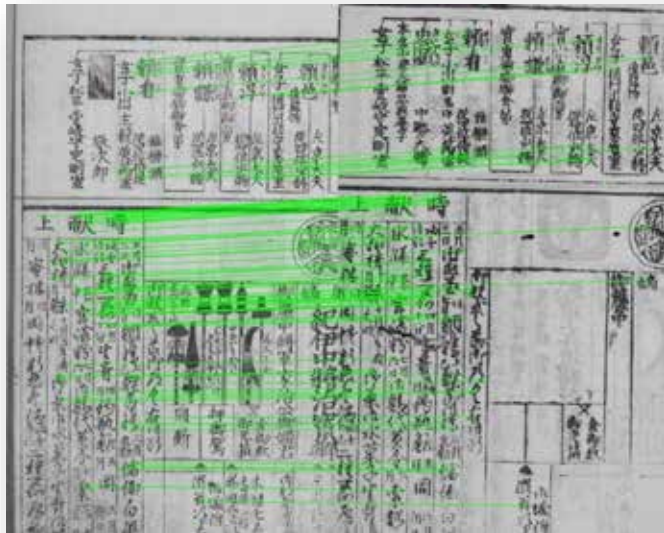
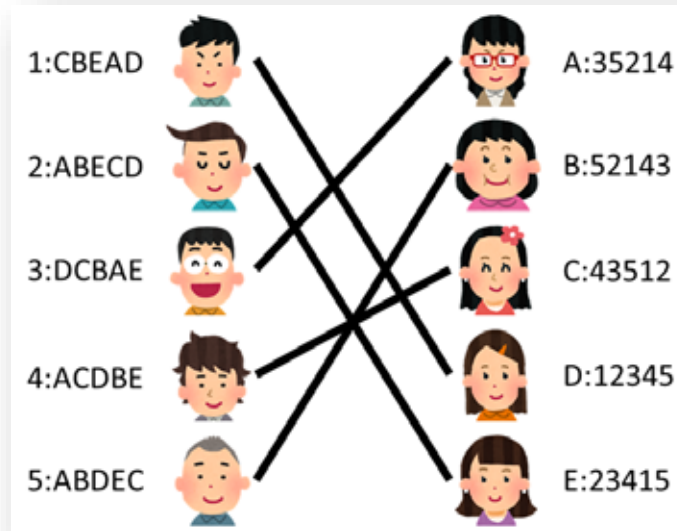http://codh.rois.ac.jp/differential-reading/



Genji Hyakunin Isshu Comparison,
University of Tokyo Library.

1. A JavaScript-based tool "vdiff.js" for comparing images.

2. Anyone can upload two images (or specify URLs).

3. The system can automatically match two images and emphasize the difference.

4. When the system fails, you can manually improve the matching.

# Large-Scale Book Collation



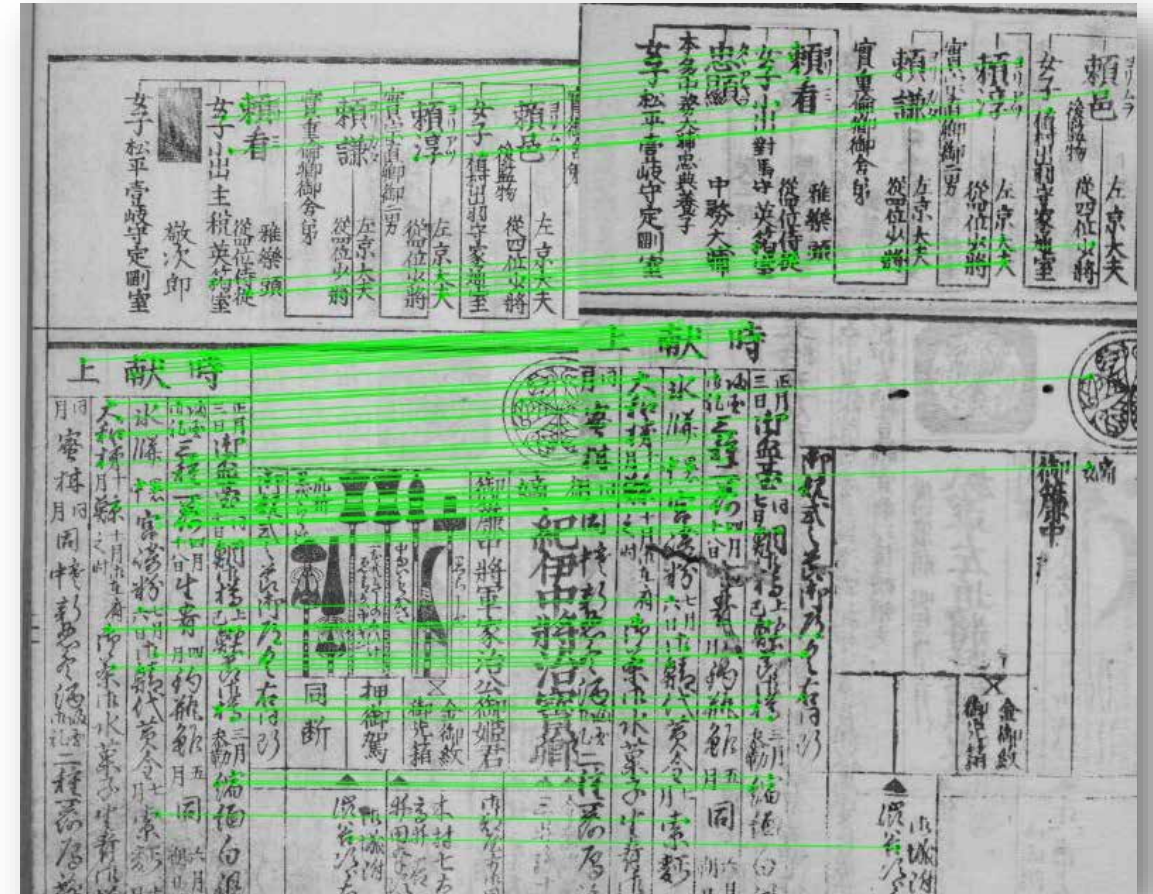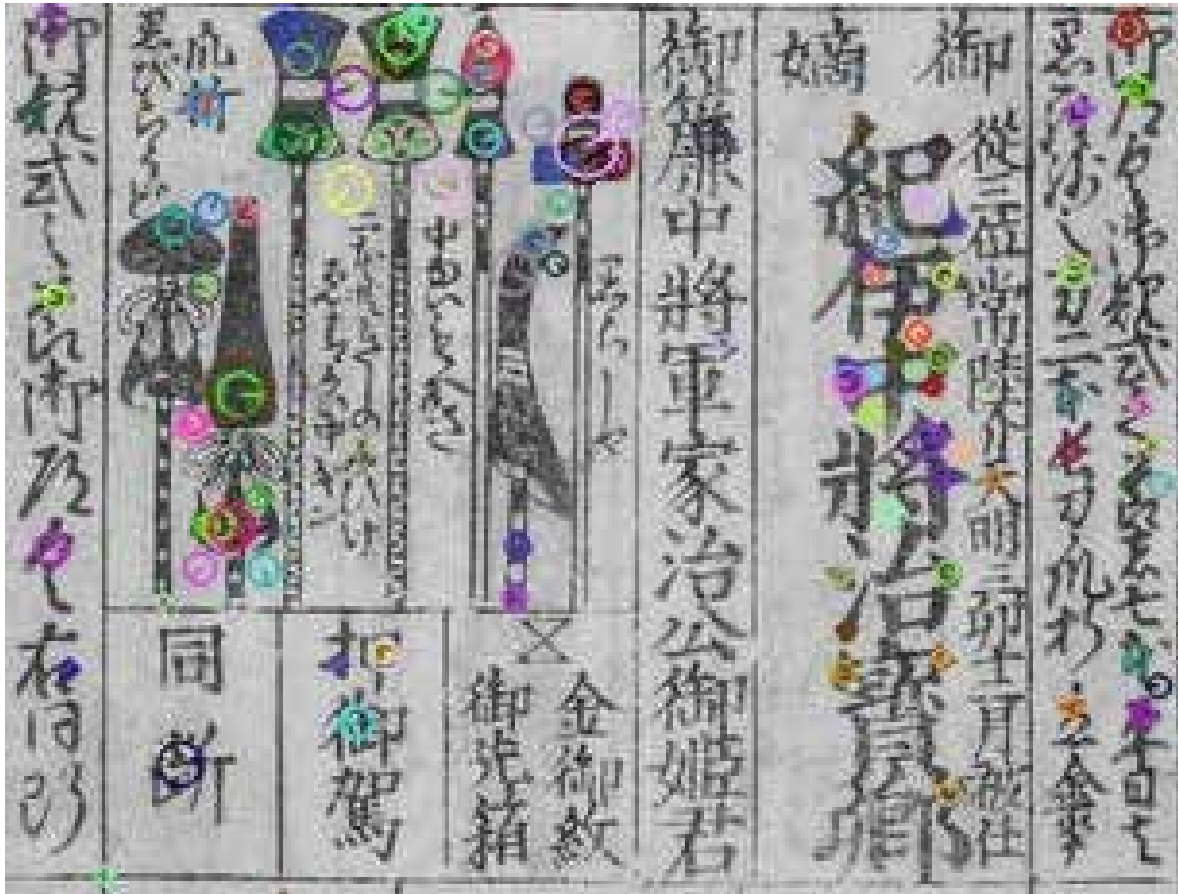**1. Page collation:** image matching using keypoints.

**2. Book collation:** stable marriage algorithm based on page collation.

**3. Woodblock tracking:** The same woodblock is estimated and connected across books.

# Page Collation – Keypoint Matching

# Book Collation – Stable Marriage Algorithm

| Book A | | Book B | Score |
|--------|---|--------|-------|
| 1 | | 1 | 0 |
| 2 | | 2 | 5 |
| 3 | | 3 | 10 |
| 4 | | 4 | 4 |
| 5 | | 5 | 6 |
| 6 | | 6 | **50** |
| 7 | | 7 | 8 |

1:CBEAD    A:35214

2:ABECD    B:52143

3:DCBAE    C:43512

4:ACDBE    D:12345

5:ABDEC    E:23415

Image source: Irasutoya

# Page-by-Page Collation – Visualization by vdiff.js
http://codh.rois.ac.jp/software/vdiffjs/

# Woodblock Tracking

http://codh.rois.ac.jp/bukan/diff/woodblock/
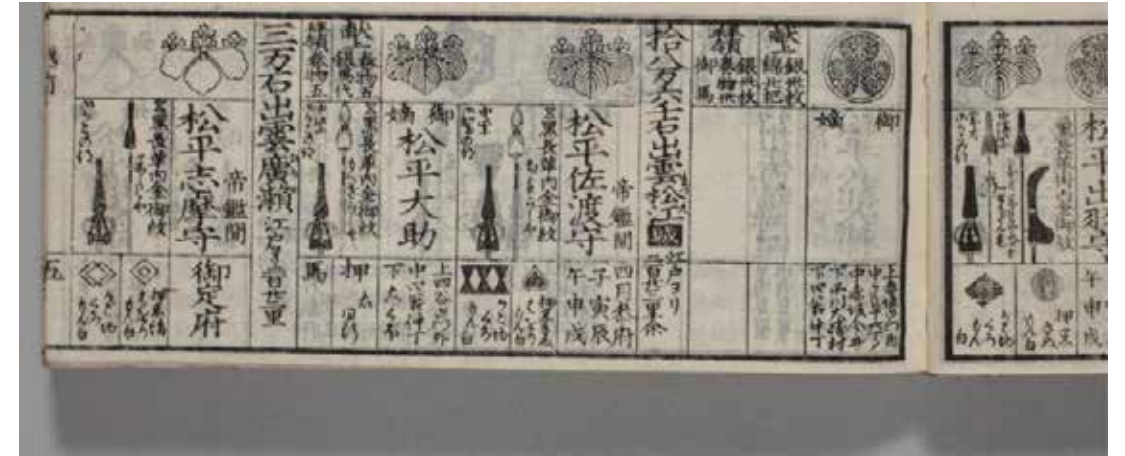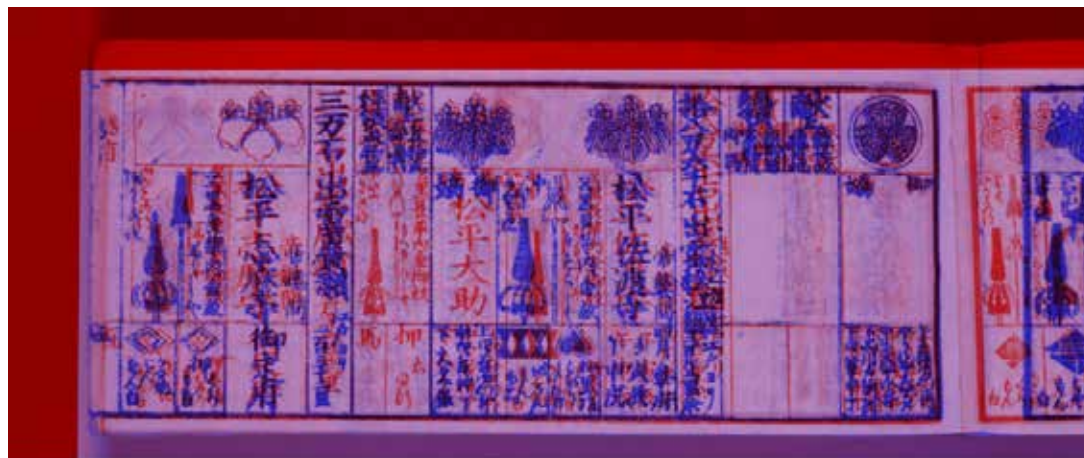


1809 [200019413]



1841 [200019430]



The same woodblock can be tracked to analyze the evolution of information on the woodblock.

# KaoKore and IIIF Curation Platform

Collaborator: Chikahiko Suzuki (CODH), Jun Homma (FLX Style), Yingtao Tian (Google Brain)

# What is IIIF ("triple-I F")?

IIIF = International Image Interoperability Framework
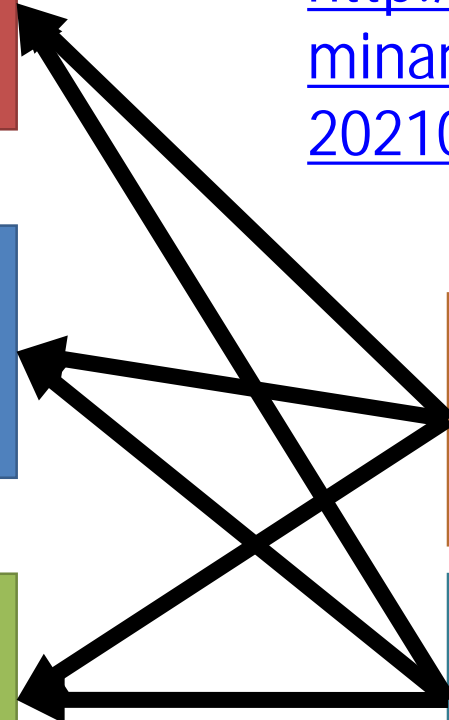
Web: HTML
Images: IIIF

IIIF service 1

IIIF service 2

IIIF service 3

IIIF viewer 1

IIIF viewer 2

14th CODH Seminar - 100 Recipes for **IIIF** Curation Platform
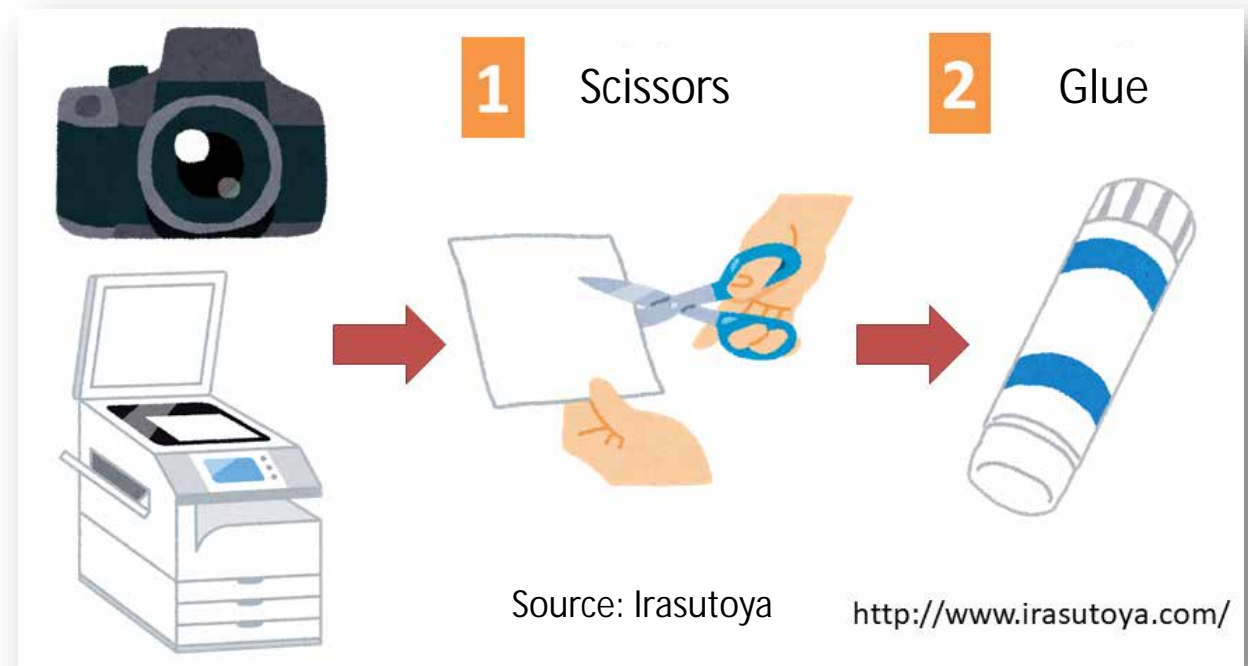http://codh.rois.ac.jp/seminar/icp-recipe-20210218/

# What is Curation?

"Curation" is a word that originally means activities at museums such as collecting materials and exhibiting artworks.

1. Collect materials under a certain theme.
2. Arrange them in an appropriate order (layout).
3. Present or share the result as a new material.

Scissors

Glue

Source: Irasutoya

http://www.irasutoya.com/

# IIIF Curation Viewer

http://codh.rois.ac.jp/software/iiif-curation-viewer/



1. **2** is the "crop" button →
   Selects a rectangular region
2. **1** is the "favorite" button →
   Collects regions you need

# Collection of Facial Expressions (KaoKore)
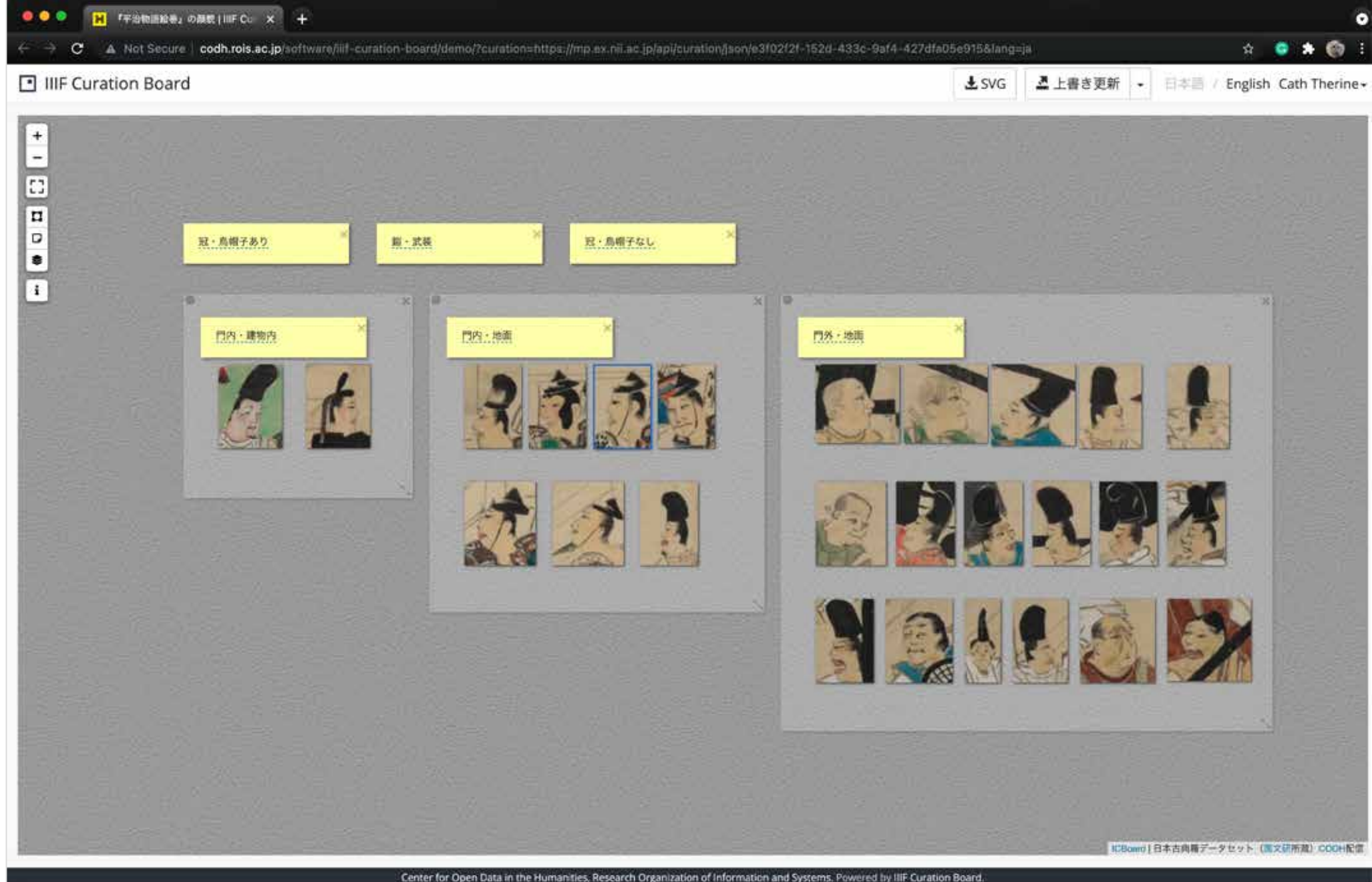
http://codh.rois.ac.jp/face/



1. **IIIF Curation Viewer** for cropping and collecting a part of images.

2. **IIIF Curation Finder** for searching the collection by metadata.

3. **IIIF Curation Board** for analyzing the collection for art history research (**digital humanities**).

# IIIF Curation Board

http://codh.rois.ac.jp/software/iiif-curation-board/

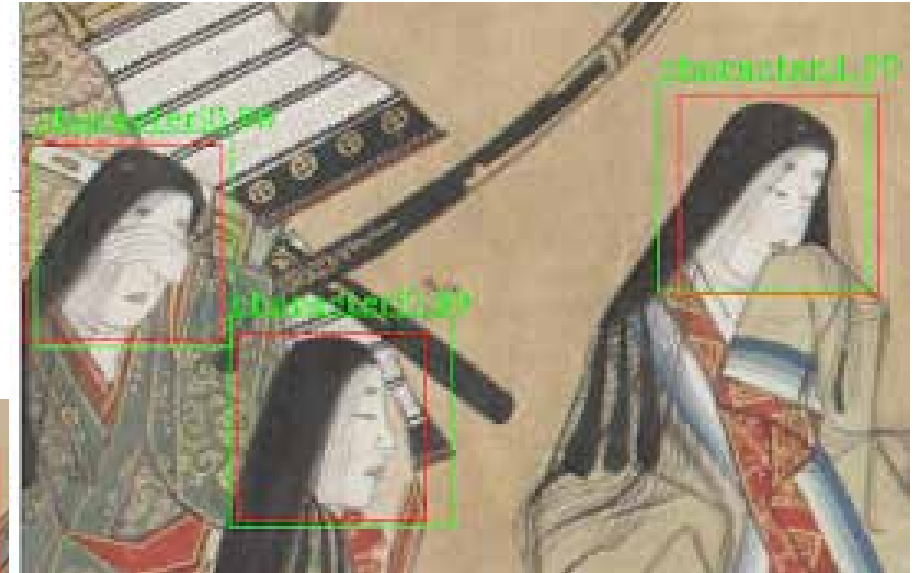# Face Detection by Machine Learning



Alexis Mermet, Asanobu KITAMOTO, Chikahiko SUZUKI, Akira TAKAGISHI, "Face Detection on Pre-modern Japanese Artworks using R-CNN and Image Patching for Semi-Automatic Annotation", Proceedings of the 2nd Workshop on Structuring and Understanding of Multimedia heritAge Contents (SUMAC'20), pp. 23-31, doi:10.1145/3423323.3423412, 2020.

Source: Kaokore dataset

# ML-assisted Annotation

1. Learning from the KaoKore Dataset, **about 80%** of the faces were automatically detected.

2. **About 70%** of the faces were automatically detected when applied to artworks from different time periods.

3. If **two thirds** can be detected by machines, the amount of work by humans is reduced to **one thirds**.

4. Art historians can analyze more data, and more data leads to richer evidence and higher reliability of the results.

# Ukiyo-e Faces Dataset

http://codh.rois.ac.jp/ukiyo-e/face-dataset/



"ARC Ukiyo-e Faces Dataset" (Created by Yingtao Tian, ROIS-DS CODH; Collected from ARC　, https://doi.org/10.20676/00000394

1. **Art Research Center** of Ritsumeikan University has Ukiyo-e Dataset.

2. ML researcher from **Google Brain** found that existing API can crop the faces.

3. A new dataset was released for visual Ukiyo-e research.

# Edo Maps, edomi, and Historical Big Data

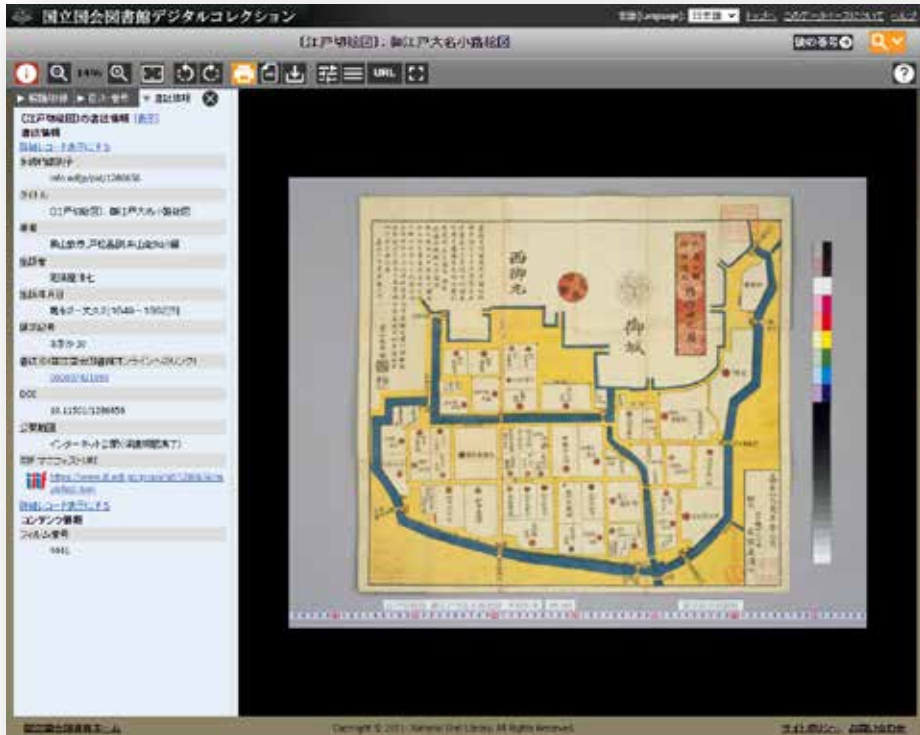Collaborator: Chikahiko Suzuki, Mika Ichino (CODH)

# Edo Maps Beta

http://codh.rois.ac.jp/edo-maps/

| 番号 | 分類 | | 現代語訳 | 翻刻 | 地図 |
|---|---|---|---|---|---|
| 2-001 | | 施設 | 幸橋御門 | 幸橋御門 | 拡大図 |
| 2-002 | | 施設 | 山下御門 | 山下御門 | 拡大図 |
| 2-003 | | 施設 | 数寄屋橋御門 | 数寄屋橋御門 | 拡大図 |
| 2-004 | | 施設 | 鍛冶橋御門 | 鍛冶橋御門 | 拡大図 |
| 2-005 | | 施設 | 呉服橋御門 | 呉服橋御門 | 拡大図 |
| 2-006 | | 地名 | 一石橋 | 一石橋 | 拡大図 |
| 2-007 | | 地名 | 出橋 | 出橋 | 拡大図 |
| 2-008 | | 町名 | 丸屋町 | 丸屋丁 | 拡大図 |

[ 2-296 ]
地名：磯辺大神宮（イソベ大神宮）
分類：寺社仏閣

**From 29 sheets, 8719 place names were extracted.**

# Annotating the Maps



Edo Kiriezu Owariya Version (1849-1862) from the Digital Collection of National Diet Library. doi:10.11501/1286656

With IIIF, you can add value by annotating information without copying original images.



Read the image on the IIIF Curation Viewer, draw a rectangle to record the coordinate, transcribe characters, and save them using the IIIF curation format.
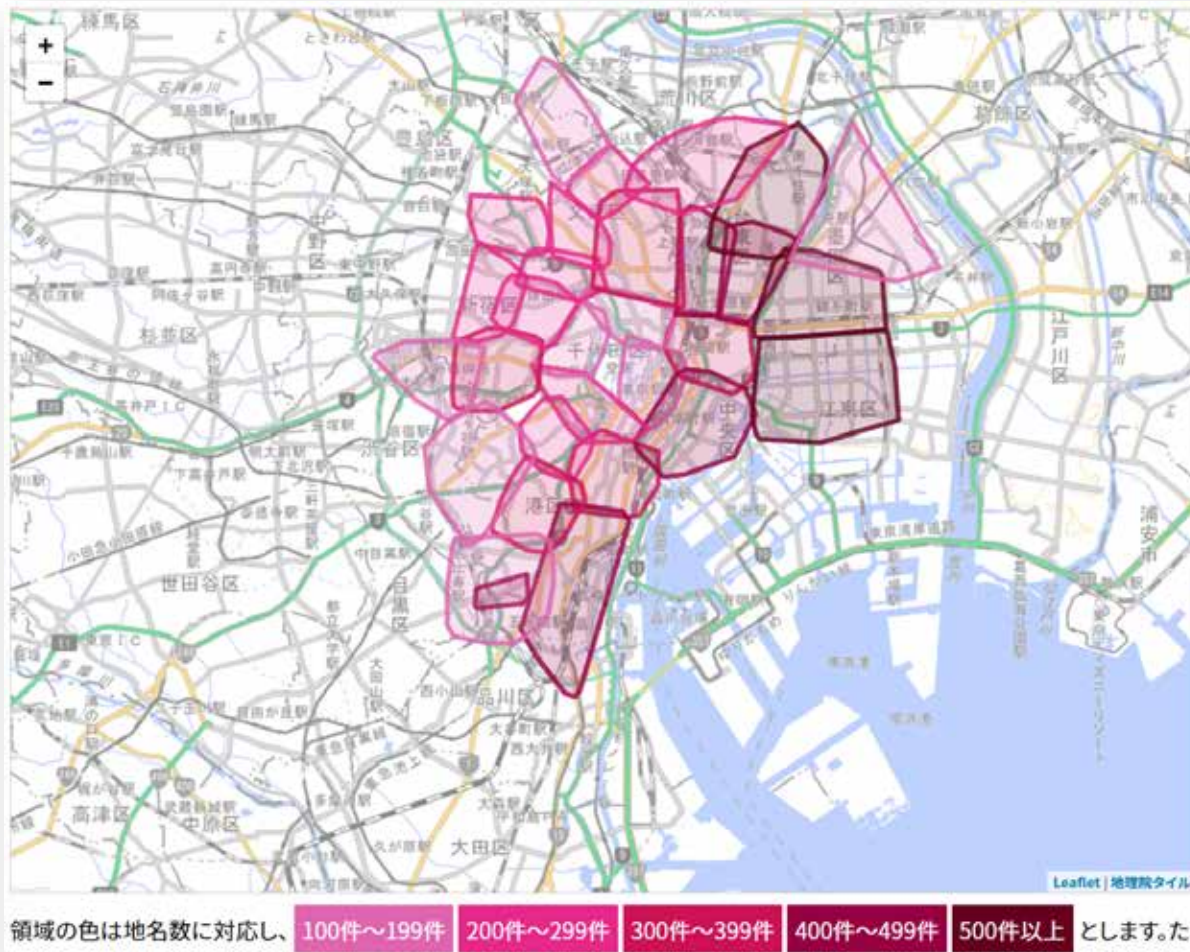
# Georeferencing the Maps



Correspondence of GCP

National Diet Library
"Edo Kiriezu"

Ritsumeikan University
Map Warper for Japanese

Edo Maps + Map Warper tile
service

# Overview of Edo Maps

http://codh.rois.ac.jp/edo-maps/owariya/



| 分類 | | 内容 |
|---|---|---|
| | 施設 | Facilities = 1326　屋敷、門、河岸、馬場、囚猿 |
| | 屋敷地 | Residence = 1647　泰屋敷」、○○組、同心 |
| | 寺社 | Temples and shrines = 1990　]絵図屋 |
| | 商店 | Shopping sites = 56 |
| | 地名 | Place names = 808　、橋、渡、新田、清水、上 |
| | 町村字 | Town names = 2780　「○○町蔵地」、町屋 |
| | 海川池 | Water areas = 52 |
| | 名所 | Sightseeing spots = 27 |
| | その他 | Others = 36 |

領域の色は地名数に対応し、 100件～199件 200件～299件 300件～399件 400件～499件 500件以上 とします。た

千代田区

© 2020 ZENRIN

Google Earth

2022/02/08

# Edo Maps

# GeoLOD

https://geolod.ex.nii.ac.jp/

1. An identifier designed for toponyms (GeoLOD ID).

2. **IIF canvas coordinate** is converted to **(lat, lng)** by georeferencing.

3. Metadata is integrated under an identifier.

Curations are converted to the gazetteer format for GeoLOD.

Name: Isobe Shrine
GeoLOD ID: G8AYsq
Lat: 35.676326
Lng: 139.774755

https://geolod.ex.nii.ac.jp/resource/G8AYsq

# Edo Sightseeing Guide

http://codh.rois.ac.jp/edo-spots/



1. Selected two travel guidebooks for each century.

2. Used IIIF Curation Viewer to crop pictorial parts and added metadata by transcription.

3. Assigned GeoLOD and other identifiers.

# Edo Shopping Guide

http://codh.rois.ac.jp/edo-shops/



1. Extracted the name and address of merchants from the shopping guidebook published in 1824.
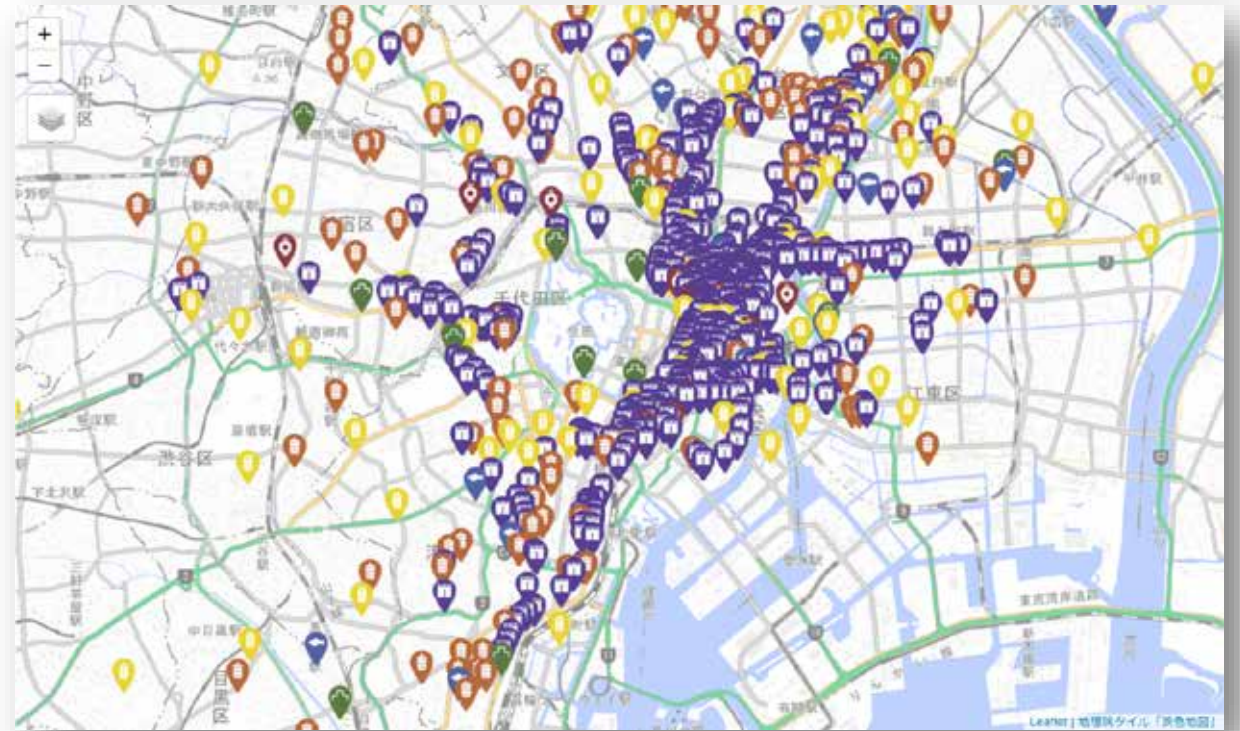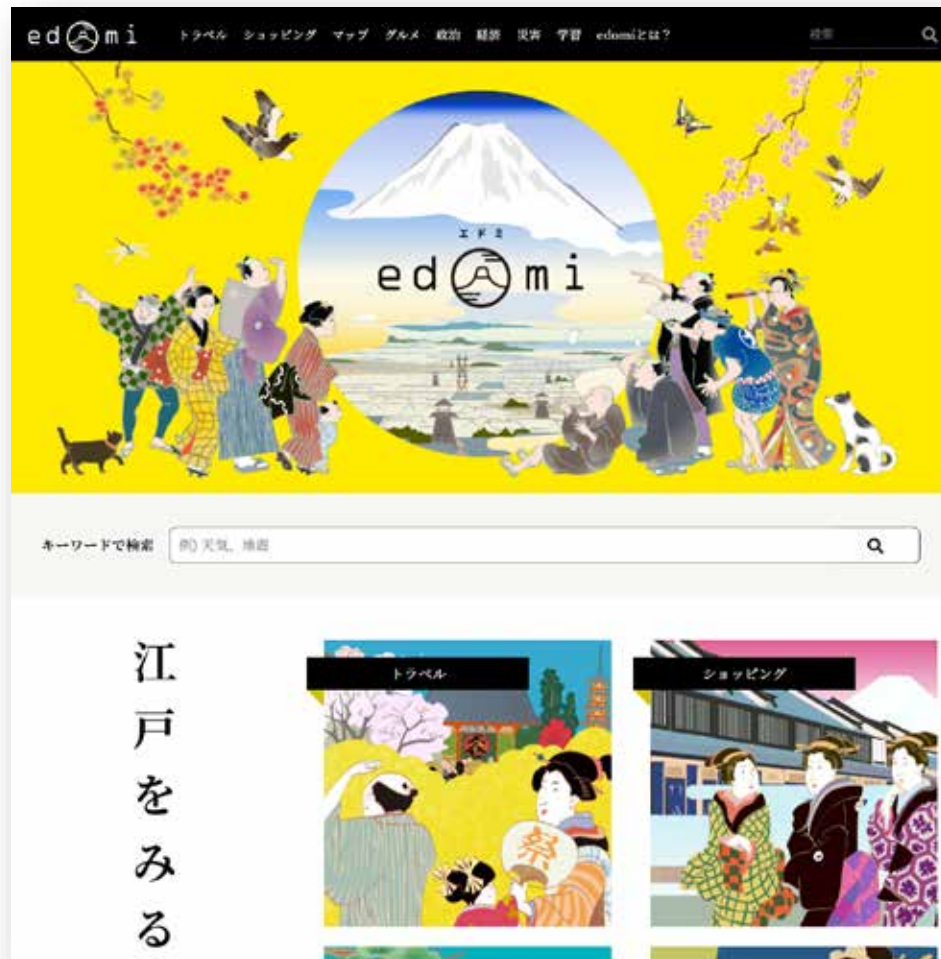
2. Classified the type of merchants' business according to **today's classification system**.

3. Assigned GeoLOD and other identifiers.

# edomi – Data Portal for the Historical Edo
http://codh.rois.ac.jp/edomi/





The distribution of geographic features (e.g. sightseeing spots and commercial stores) in the city of Edo.

# Look Back the History

| Type | Number |
|------|-------:|
| Medicine | 205 |
| Household hardware | 175 |
| Household goods | 130 |
| Confectionary | 116 |
| Restaurant | 102 |
| Cosmetics | 92 |
| Toy | 88 |
| Grain | 83 |
| Medicine (wholesale) | 71 |
| Tobacco | 69 |
| Cloth | 65 |
| Dry foods | 62 |
| Publisher | 60 |
| Shoes | 58 |

A mapping from the type of business in the shopping guide to the Japan Standard Industrial Classification (2013).

Experts try to understand the past as it was.

The general public wants to look back the history from the present.

# Time Machine Europe

https://timemachine.eu/



1. **Big Data of the Past**: create machine-readable data of the past using AI and simulation.

2. Developing new critical reflections on the past and future.

# Historical Big Data

**Historical sources**

Nature data → Weather, Earthquake, Eruption, Disease

Culture data → Economy, Population, Politics, Culture

Data structuring workflow

Platform for the integrated analysis of HBD.

# Data Structuring Workflow

Handwritten characters (published, copied, written)

Analysis-ready data (quality controlled and curated data)

Digitization

Digital image (unstructured data)

Tabular data (structured data)

Linked data

Entity linking

Transcription

Gap between dual spaces

Markup

Plain text (unstructured data)

Encoded or annotated text (semi-structured data)

# Data-driven Approaches for Japan Studies

1. **AI kuzushiji recognition** illustrates how a <span style="color:red">machine learning project</span> can be started and developed into the real world.

2. **Bukan Complete Collection** shows how the idea of <span style="color:red">differential reading</span> can reduce the burden of humans.

3. **Kaokore** demonstrates how <span style="color:red">interoperability such as IIIF</span> plays a critical role in a digital humanities platform.

4. **Edo maps and historical big data** explores new possibilities for <span style="color:red">linking the past, present and future</span>.

# Acknowledgments and More Information

- The presentation includes research results from CODH researchers, Chikahiko Suzuki, Mika Ichino, and Tarin Clanuwat.
- IIIF Curation Platform was mainly developed by Jun Homma and Tarek Saier.
- Some results are based on the work of NII internship students.
- Many data are from National Institute of Japanese Literature.

Visit our Website http://codh.rois.ac.jp/
Collaboration is welcome.