# Book Barcoding for Differential Reading -Application to Woodblock-printed Books in the Bukan Complete Collection-

**Asanobu KITAMOTO**

ROIS-DS Center for Open Data in the Humanities (CODH)
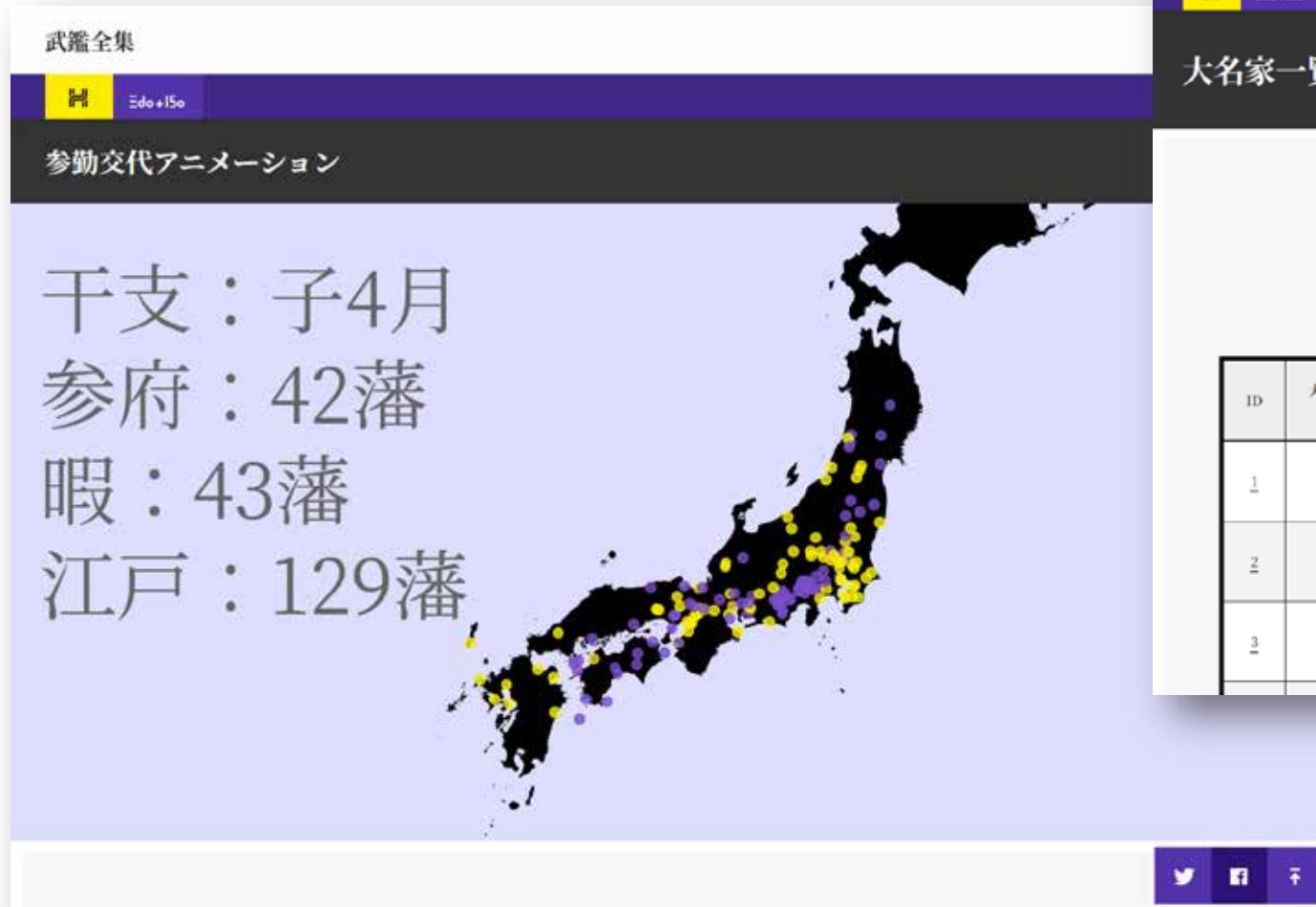
National Institute of Informatics

https://researchmap.jp/kitamoto/

kitamoto@nii.ac.jp

# What is Bukan 「武鑑」?



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.

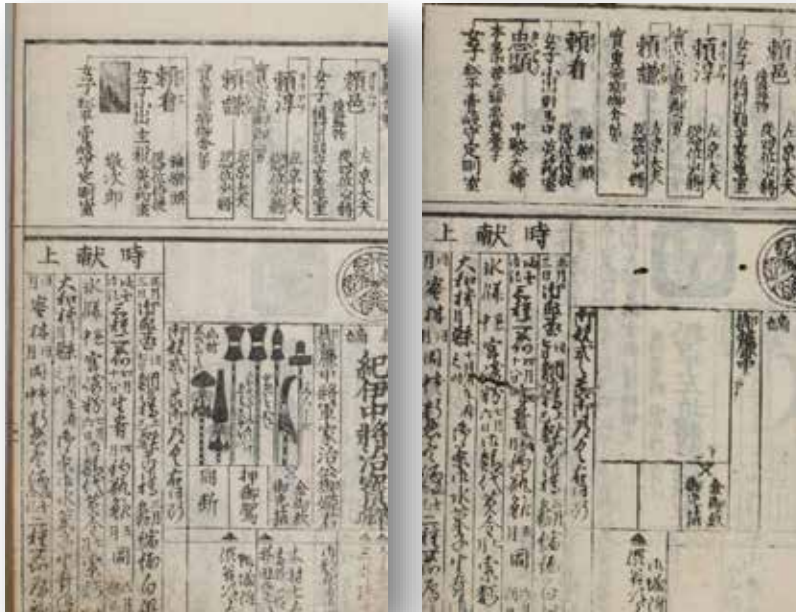Reference: Kumiko Fujizane, 2008

# Bukan Complete Collection

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



List of Daimyos

Sankin Kotai Dynamic Map

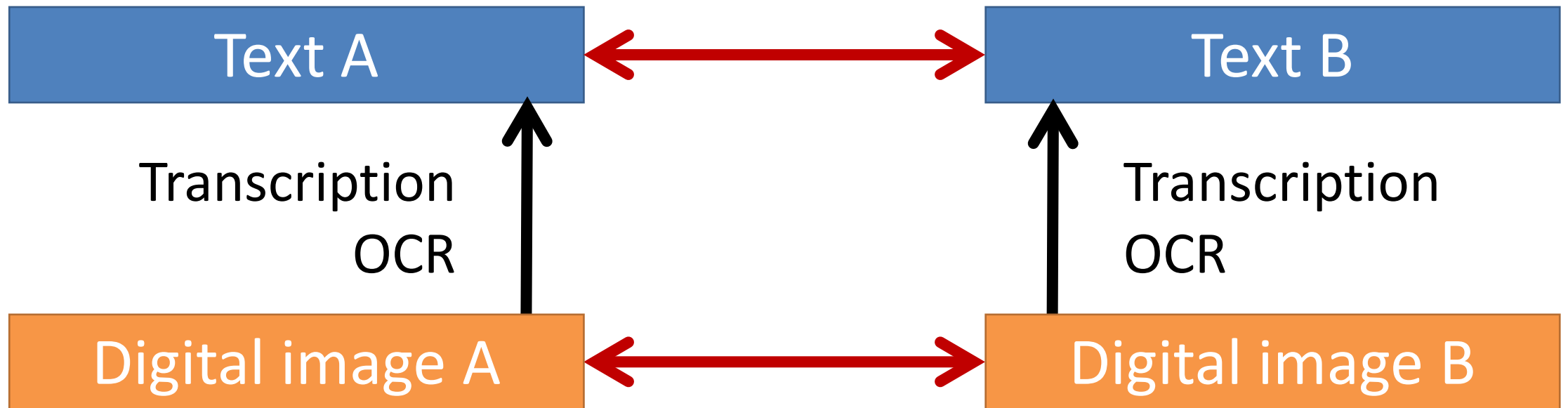# Woodblock-printed Books and Editions



Left: Kansei Bukan (1789)
Right: Kansei Bukan (1791)

Publications in the Edo Period: mainstream was woodblock printing (not movable type printing).

1. Publication: woodblock is completely recreated (= **major version**).
2. Print: multiple prints are produced from the same woodblock (= **instance**).
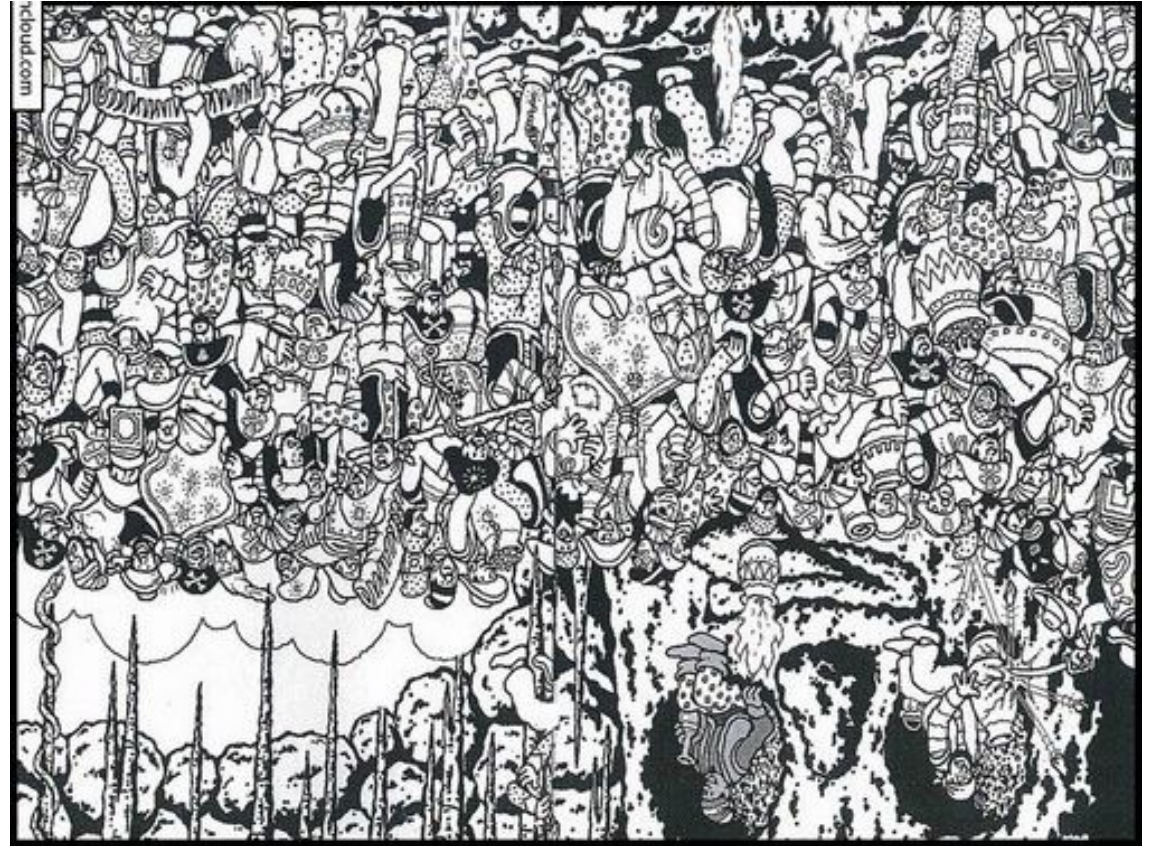3. **Correction**: woodblock is carved or patched by a small plate  (= **minor version**).
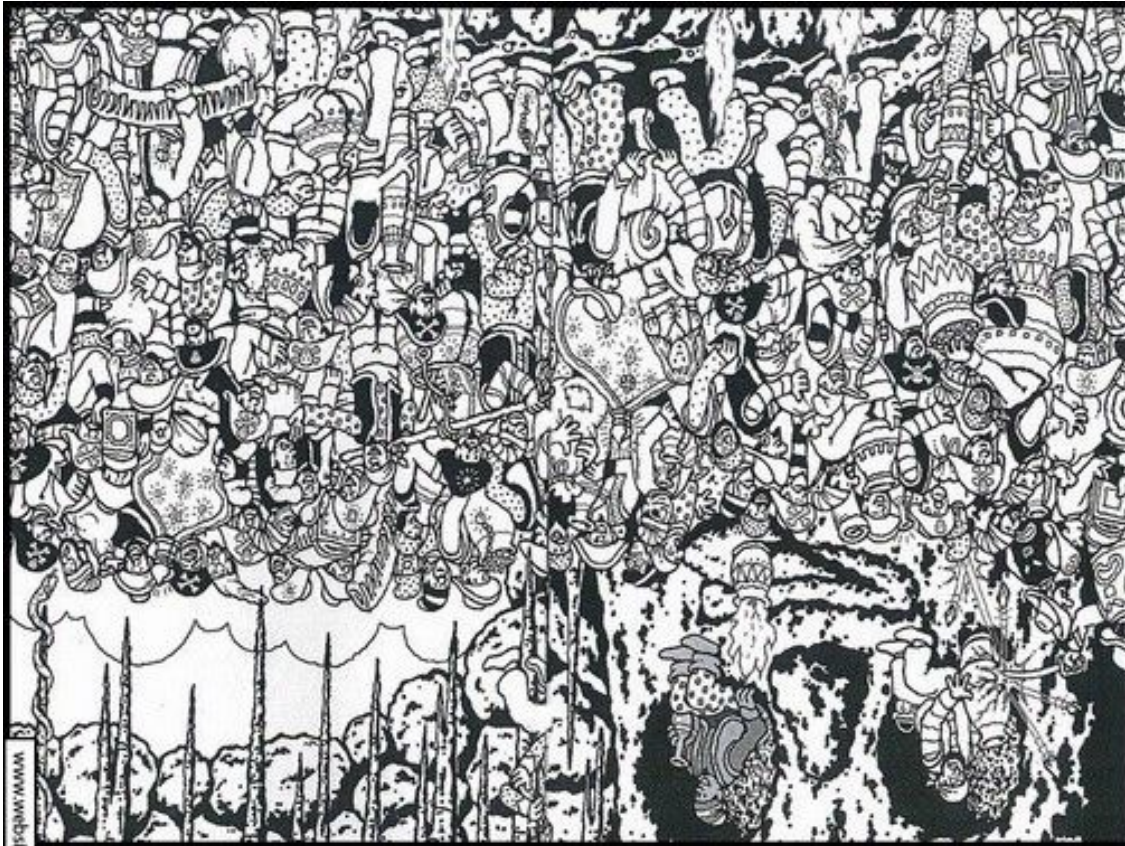
# Text-based and Image-Based Comparison

**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

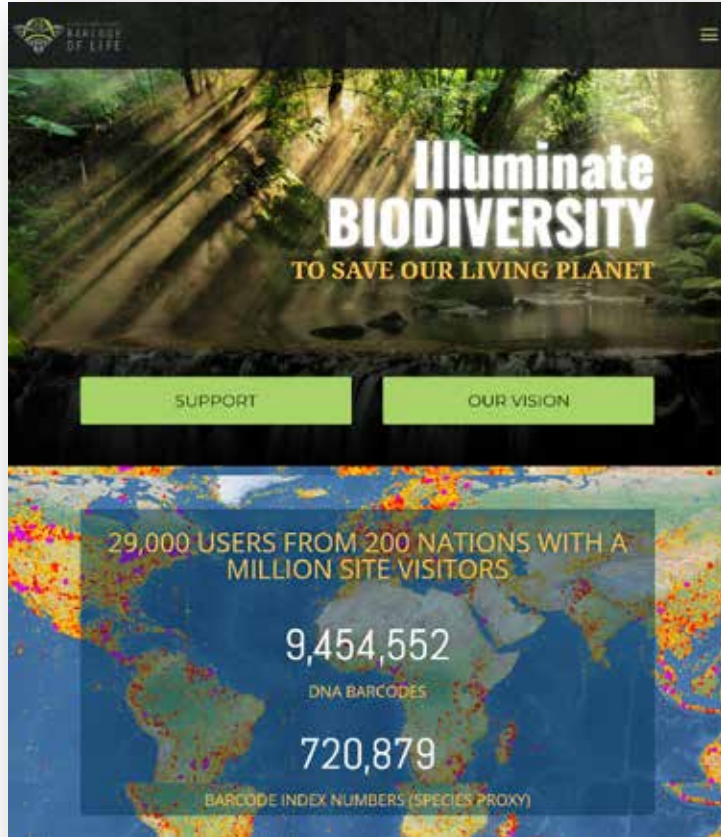http://codh.rois.ac.jp/software/vdiffjs/demo/local.html

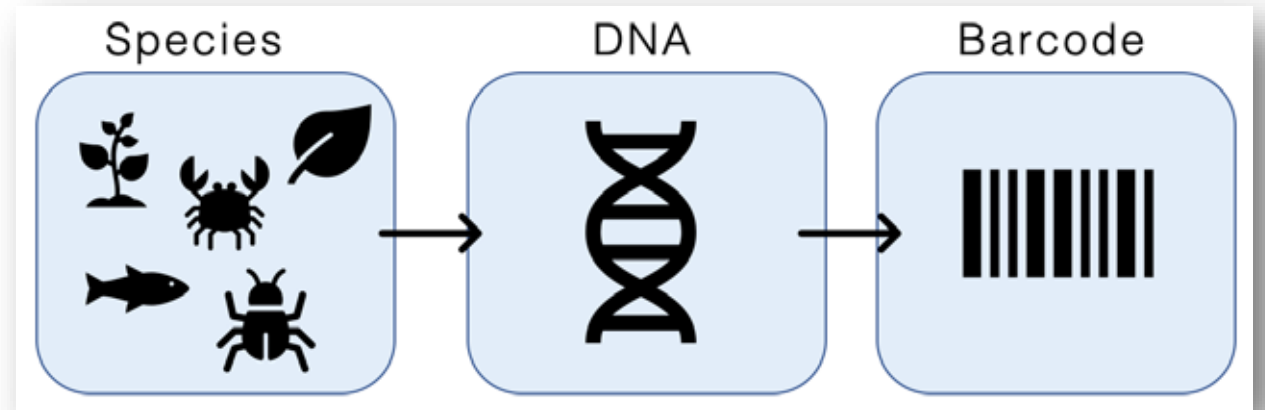Red and blue colors were used to emphasize the difference.

# 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.  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).

# DNA Barcoding



https://en.wikipedia.org/wiki/DNA_barcoding



Illuminate BIODIVERSITY
TO SAVE OUR LIVING PLANET

SUPPORT    OUR VISION

29,000 USERS FROM 200 NATIONS WITH A MILLION SITE VISITORS

9,454,552
DNA BARCODES

720,879
BARCODE INDEX NUMBERS (SPECIES PROXY)

International Barcode of Life
https://ibol.org/

1. DNA barcodes are unique DNA sequences to assign identities to sequences of unknown origin.

2. Barcode of Life Data Systems (BOLD) database is an online workbench that includes a reference library of DNA barcodes.

# Book Barcoding

Species → DNA → Barcode

Digital Image
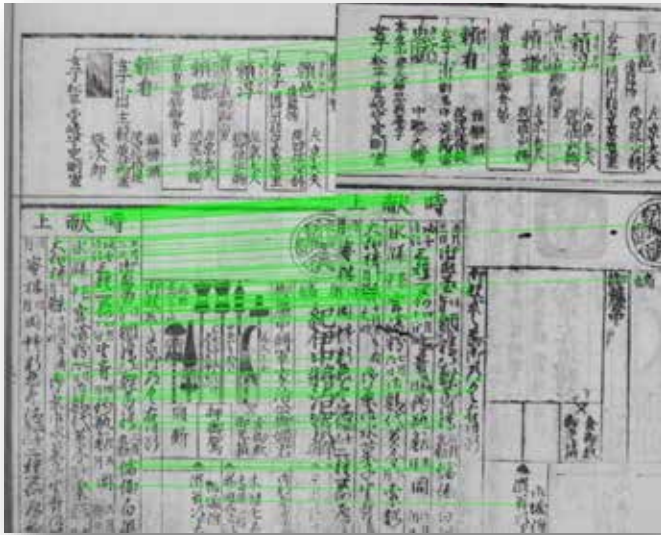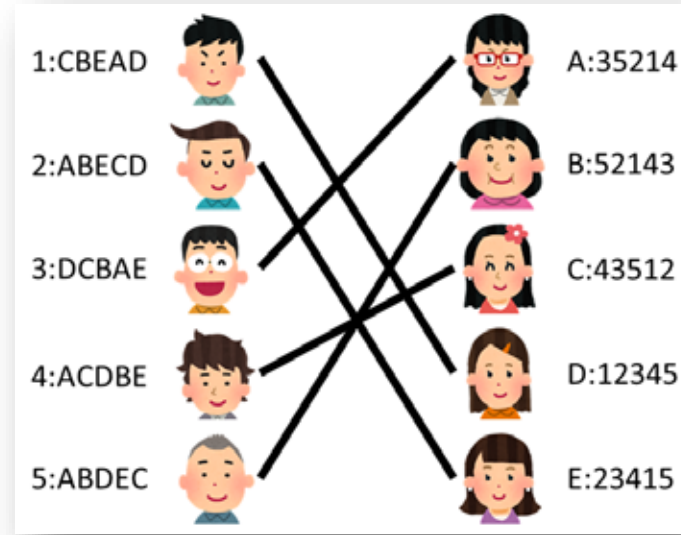
Keypoints

Barcode

# Book Barcoding Method



**1.** **Page-by-page collation**: Keypoint detection and matching



**2.** **Book-by-book collation**: stable marriage algorithm using inlier keypoints
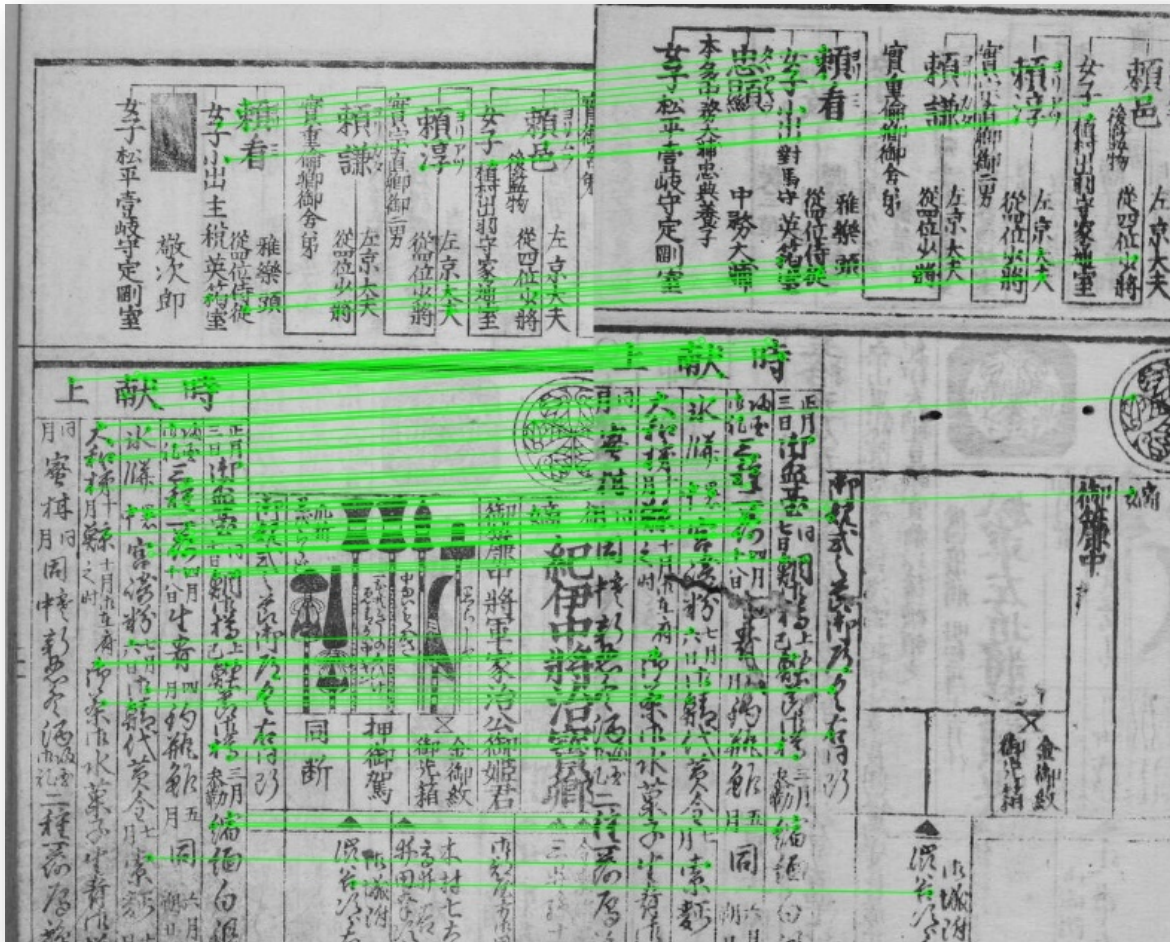


**3.** **Bukan differential reading platform**: Visualization and navigation of the results
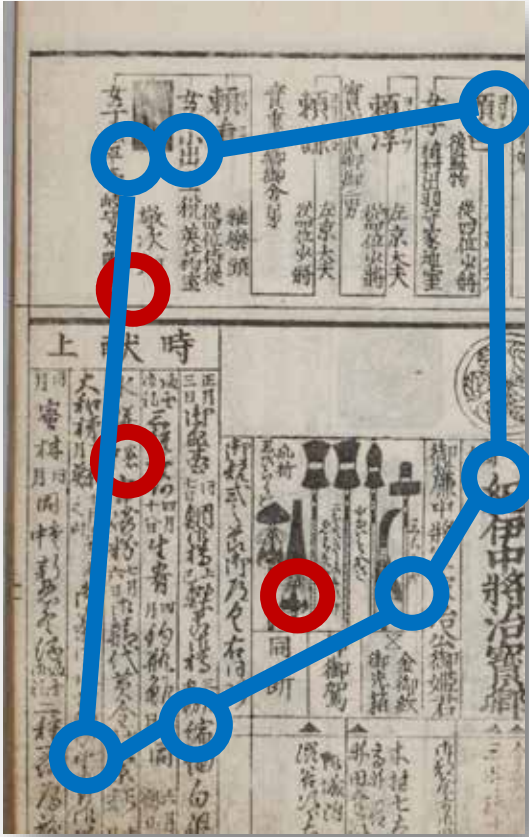
# Page-by-Page Collation - Detection



1. **Keypoint**: locations of the most distinctive features on each image.

2. **Local descriptor**: a vector of numbers that describes the visual appearance of the keypoint.

3. Among several keypoint detectors in OpenCV, we used the **AKAZE detector**.

# Page-by-Page Collation - Matching



1. Keypoints are compared by local descriptors with the Hamming distance metric.

2. **A projective transformation matrix** (used in **vdiff.js**) is estimated using the RANSAC algorithm.

3. **The number of inlier keypoints** (used in **book-by-book collation**) is counted for the goodness of matching.
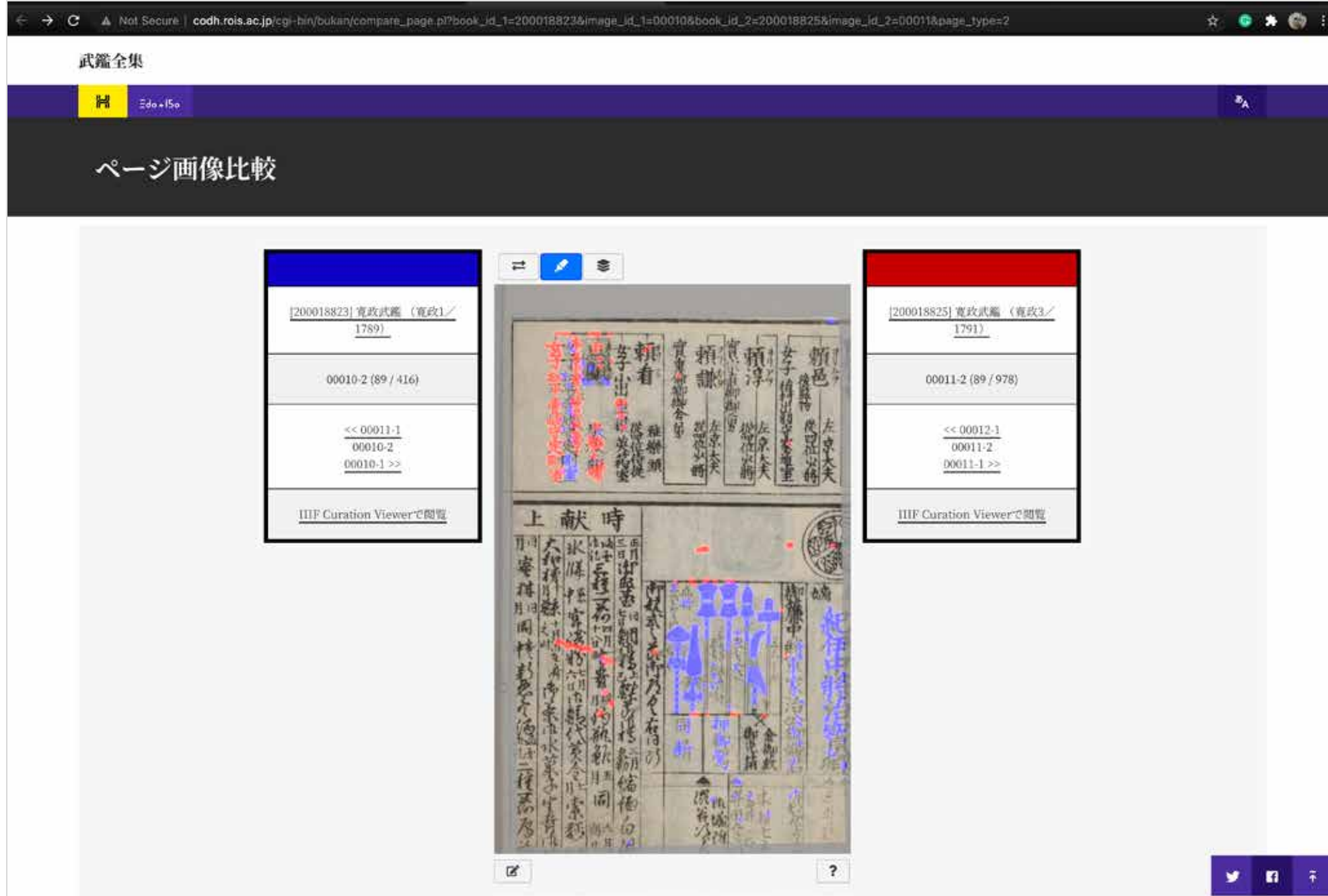
# Page-by-Page Collation - Projection



Compute the convex hull
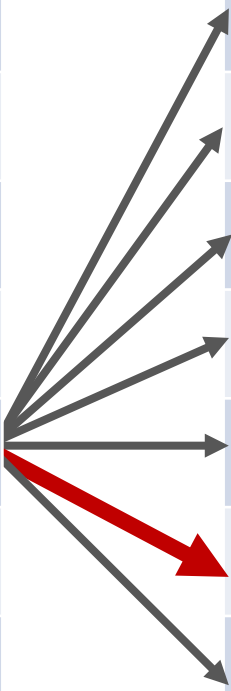
Select four corresponding points

1. The image comparison tool vdiff.js accepts four corresponding points, not a projective transformation matrix.

2. Choose points using the convex hull algorithm.

3. Choose four points which maximize distances.

# Page-by-Page Collation – Visualization by vdiff.js

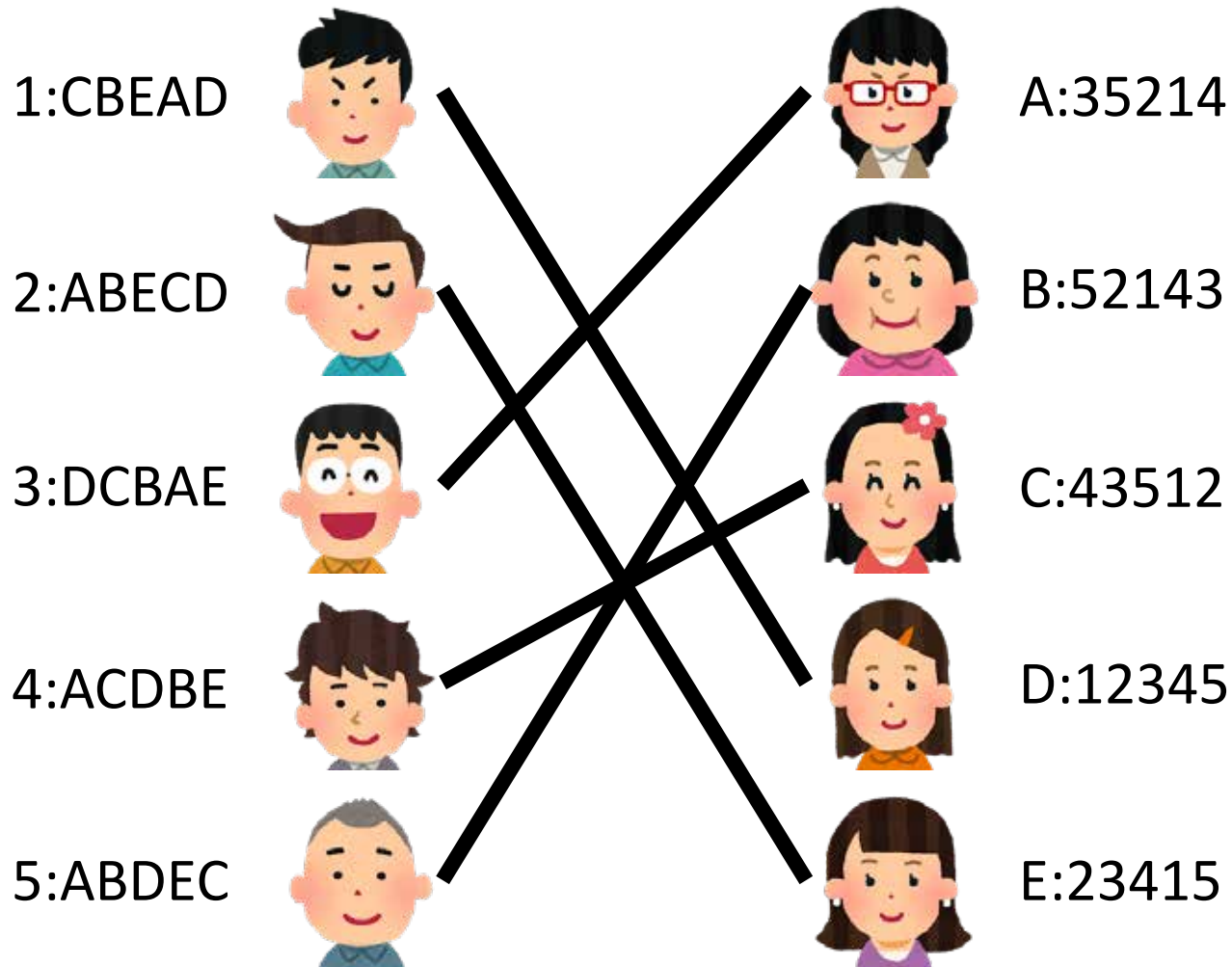http://codh.rois.ac.jp/software/vdiffjs/

# Book-by-Book Collation – Matching Score

| 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. **The number of inliner keypoints** is used as the matching score.

2. Seeing from Book A - Page 5, the order of preference for Book B is Page 6, Page 3, …

3. **What is the best way to choose the best page pairs between two books?**

# Book-by-Book Collation - Stable Marriage Problem

1:CBEAD

2:ABECD

3:DCBAE

4:ACDBE

5:ABDEC

A:35214

B:52143

C:43512

D:12345

E:23415

1. **A matching is stable** when there does not exist any match (*A*, *B*) which both prefer each other to their current partner under the matching.

2. **Gale-Shapley algorithm** is a classic solution to this problem.

Image source: Irasutoya

# Statistics

| Item | Number |
|------|--------|
| Books | **336** |
| Images | **143,616** (111,114 portrait 32,502 landscape) |
| Keypoints | **67,071,993** (467 keypoints per image) |
| Tested book pairs | **3,323** |
| Tested page pairs | **27,821,763** (8,375 page pairs per book) |
| Married page pairs | **419,118** (about 1.5% of tested page pairs) |
| Final page pairs | **418,651** |

# Bukan Differential Reading Platform

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

寛政武鑑 （寛政1／1789） [200018823]

寛政武鑑 （寛政3／1791） [200018825]

1. Choose the base edition of Bukan from the whole list.

2. Choose the target edition of Bukan from the suggested list.

3. Choose the page pair from the list of page pairs (left). Color suggests reliability of the collation (red: low, blue: high, gray: none).

http://codh.rois.ac.jp/cgi-bin/bukan/select_page.pl?book_id_1=200018823&book_id_2=200018825

# Differential Transcription

Machines can help humans



Bukan editions in time-series

Comparison with the base edition

Difference detection

Reduce the cost of transcription

Complete transcription ⟷ Differential transcription

- Base transcription: all the pages are transcribed.
- Differential transcription: only changes are transcribed.
- Recreation detection: refresh the basic transcription.
- Advantage: transcription cost is reduced to the percentage of change.

# New Research Questions

1. **Complete ordering**: Given two books, which is the newer edition? What evidence? What is the lineage of books?

2. **Publishing industry**: How fast the error was fixed? How long the woodblock had been used?

3. **Career path**: How many and how often people were promoted or disappeared?

4. **Historical big data**: How economic situation affected the management of human resources?

# Acknowledgment

- We thank **Mr. Jun Homma** in FLX Style for the development of vdiff.js as the core contributor.

- We thank **Prof. Kumiko Fujizane** and **Prof. Kazuaki Yamamoto** in the National Institute of Japanese Literature for helpful discussion.

- A part of the research is based on the work of **Mr. Thomas Leyh** who contributed to this project while he was an NII internship student.

- This work is partially supported by JSPS KAKENHI Grant Number JP19H01141.

  - Bukan Complete Collection http://codh.rois.ac.jp/bukan/
  - Vdiff.js http://codh.rois.ac.jp/software/vdiffjs/