# Book Barcoding: A Framework for the Visual Collation and Woodblock Tracking of Japanese Printed Books

**Asanobu KITAMOTO**

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

National Institute of Informatics

https://researchmap.jp/kitamoto/

Twitter: @kitamotoasanobu, @rois_codh

# Research Question and Goal

**Research Question**

How can we detect differences in multiple Japanese woodblock-printed books with the help of machines?

**Research Goal**

Create a platform for large-scale book collation and differential transcription to create diachronic data efficiently and analyze the genealogy of books.
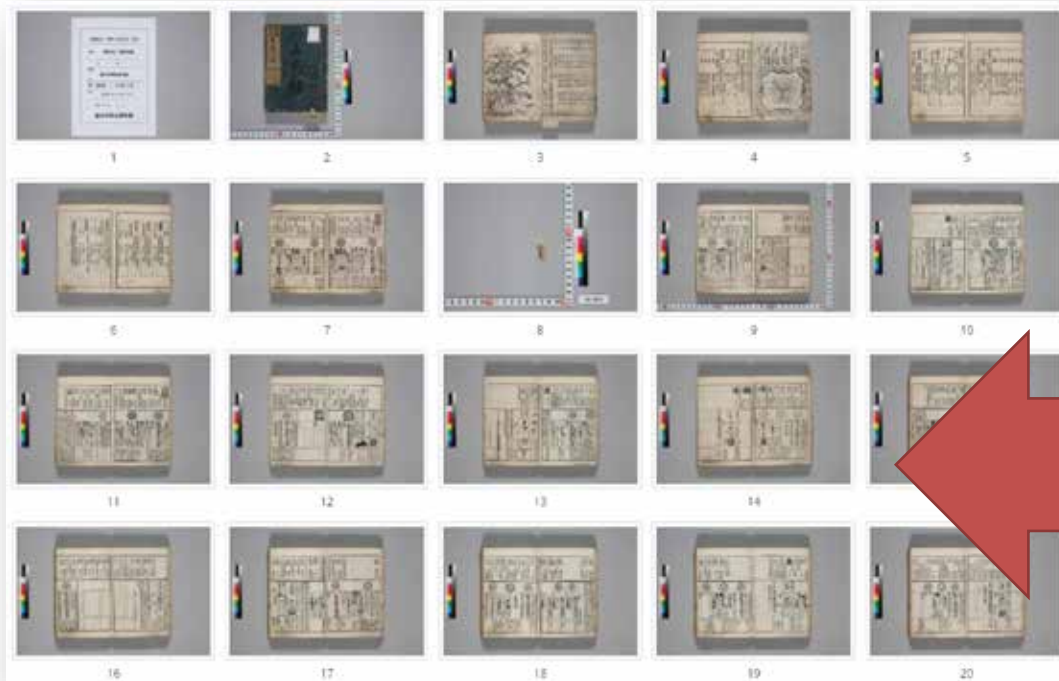
# Image-based Comparison and Differential Reading

Digital Humanities 2022

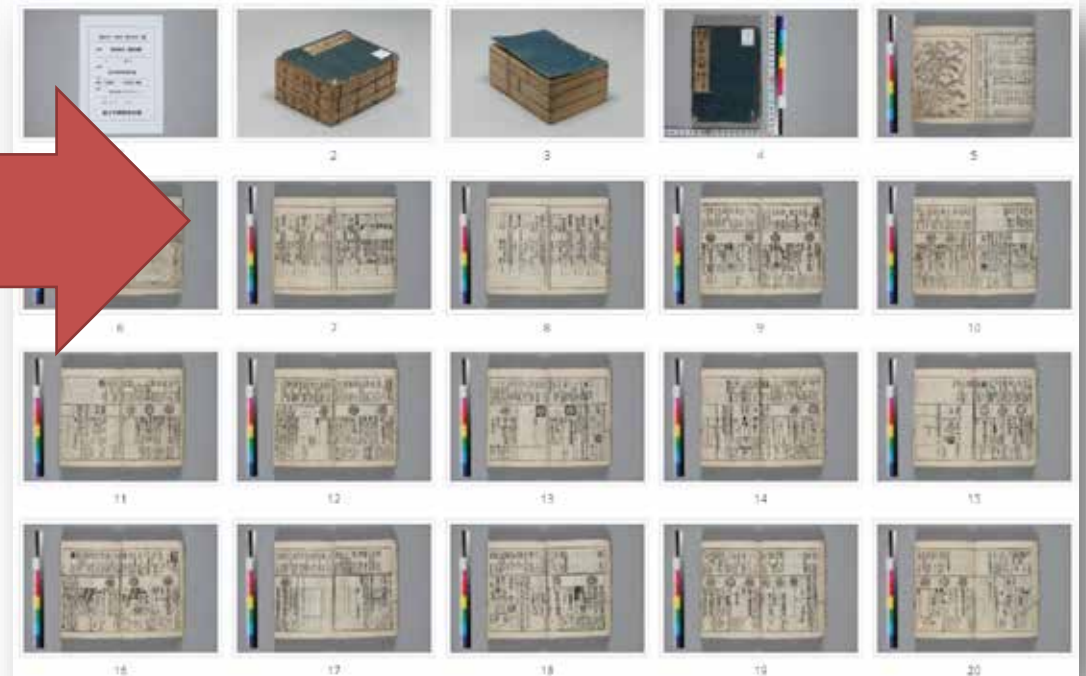# Japanese Woodblock-printed Books

# Searching for a Page Pair Printed by the Same Woodblock
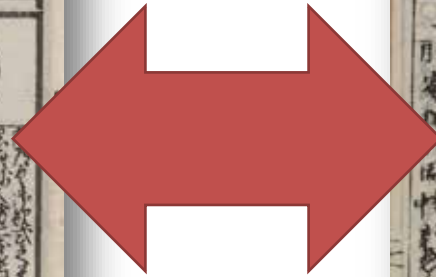


Kansei Bukan 1791

Kansei Bukan 1789

# Comparing Page Images Printed by the Same Woodblock

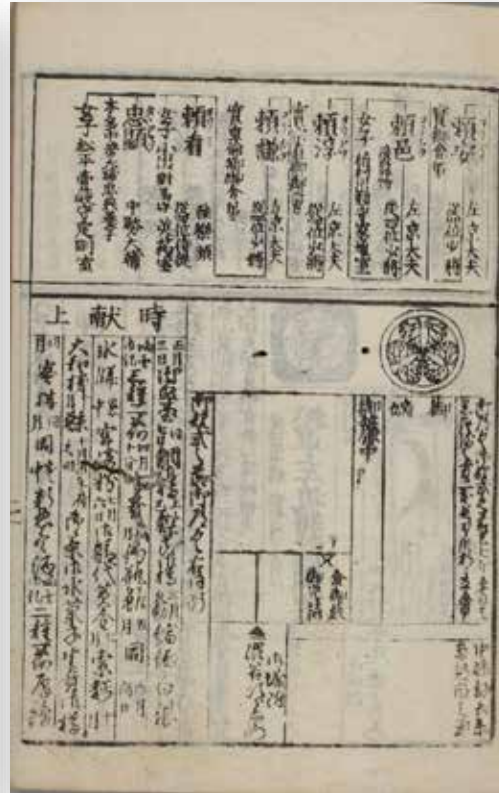Kansei Bukan 1789



Kansei Bukan 1791

# Editions of Woodblock-printed Books

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**).
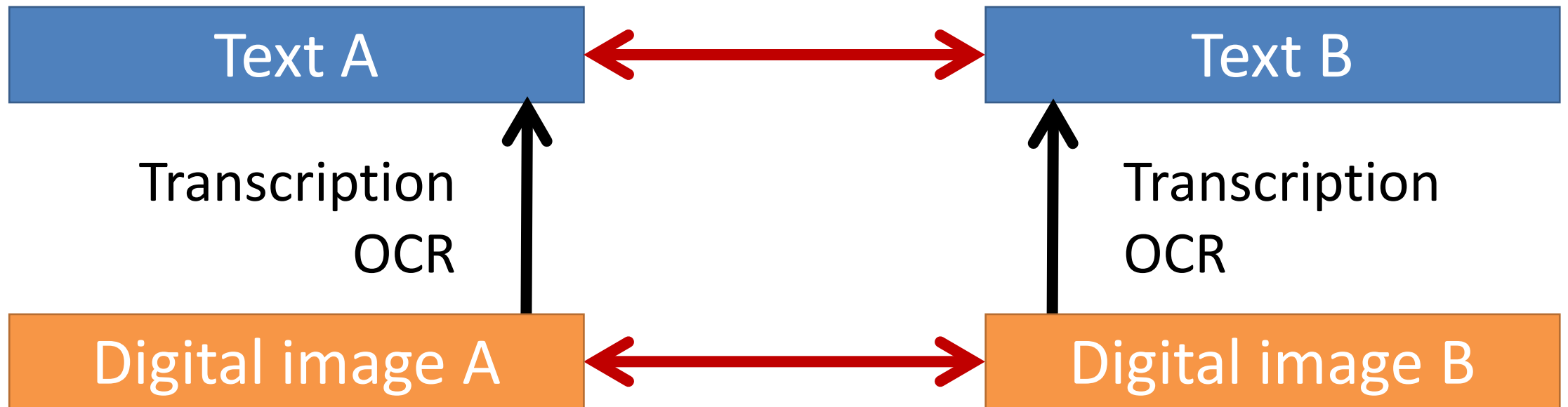
Kansei Bukan 1789

Kansei Bukan 1791

# Differential Reading

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

2. **Difficult for humans**: visual comparison requires an effort comparable to playing games.

3. **Easy for machines**: visual comparison can be done using a computer vision-based image matching algorithm.
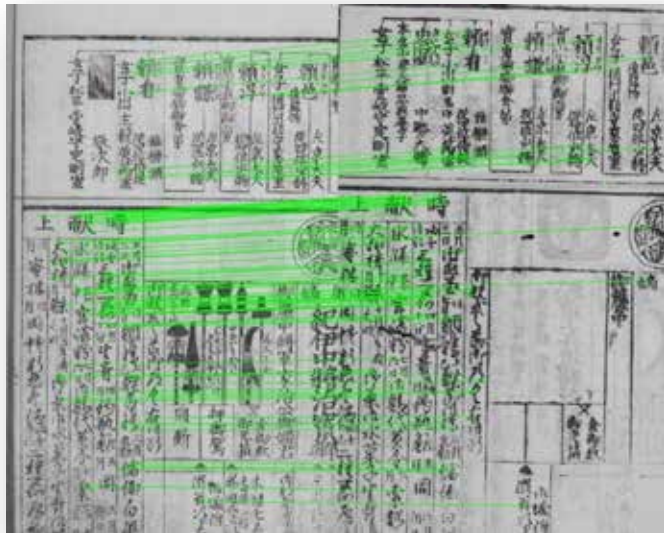
# Text-based and Image-Based Comparison

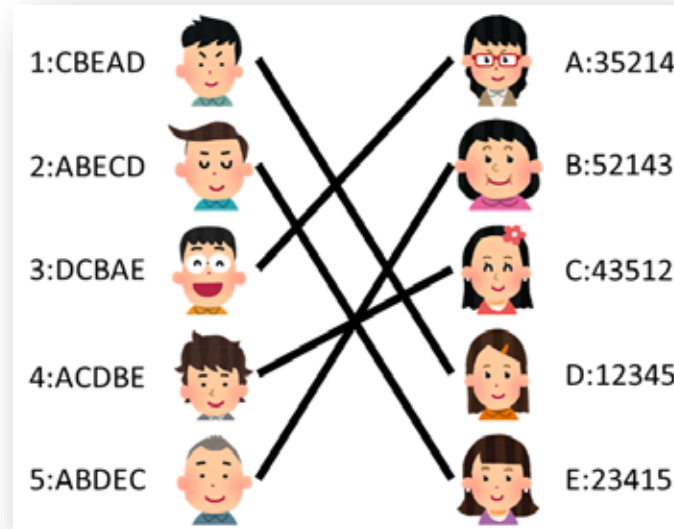**Text-based difference = many tools available**

| Text A | ↔ | Text B |

↑ Transcription OCR          ↑ Transcription OCR

| Digital image A | ↔ | Digital image B |

**Image-based (non-textual) difference =
no standard tools available (side-by-side comparison)**

# Large-Scale Book Collation Algorithm

Digital Humanities 2022

# Large-Scale Book Collation



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



1:CBEAD → 
2:ABECD → 
3:DCBAE → 
4:ACDBE → 
5:ABDEC → 
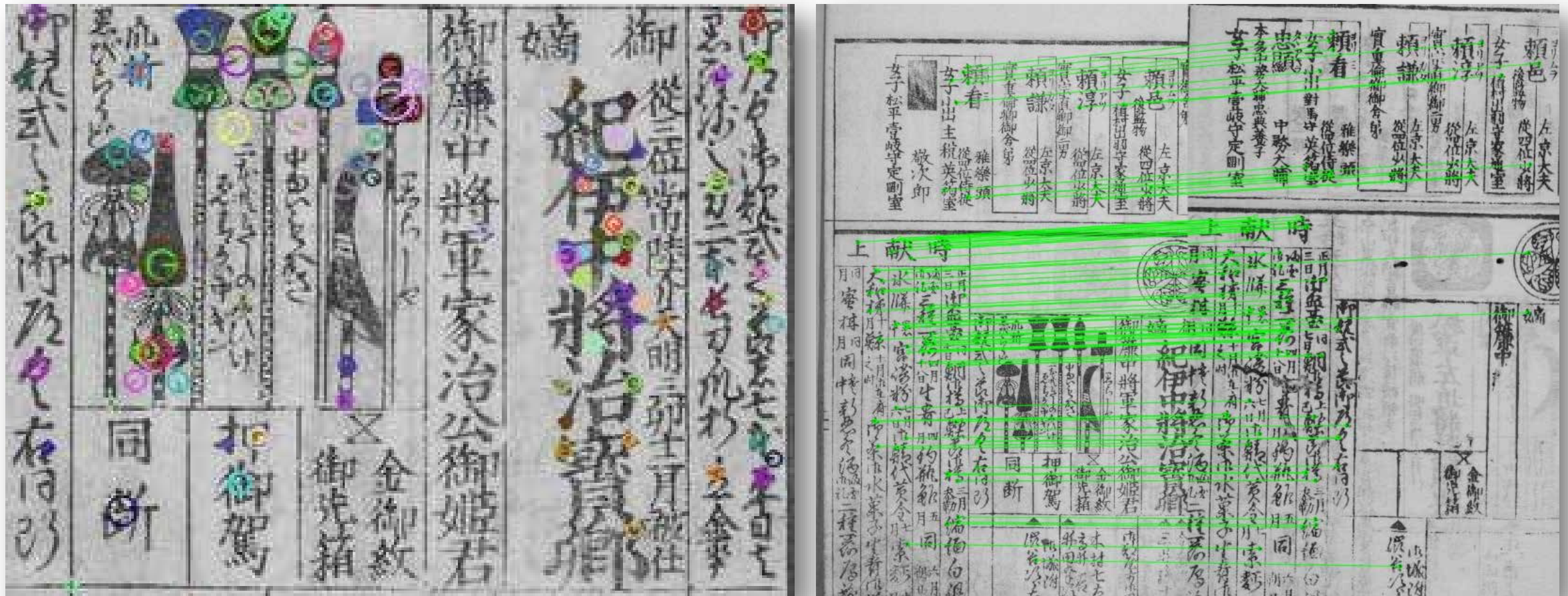
A:35214
B:52143
C:43512
D:12345
E:23415

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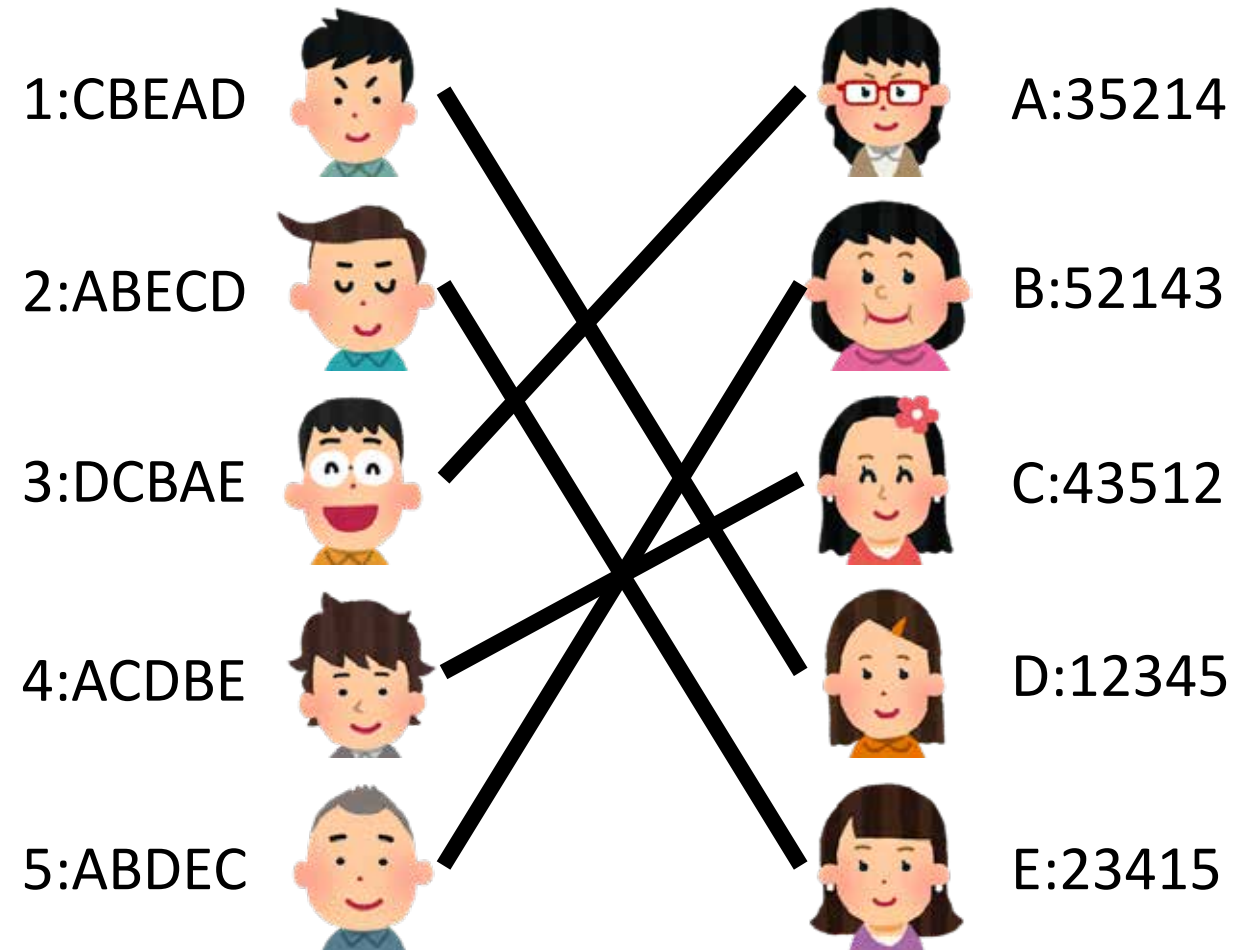


**3. Woodblock tracking:** The same woodblock is estimated and connected across books.
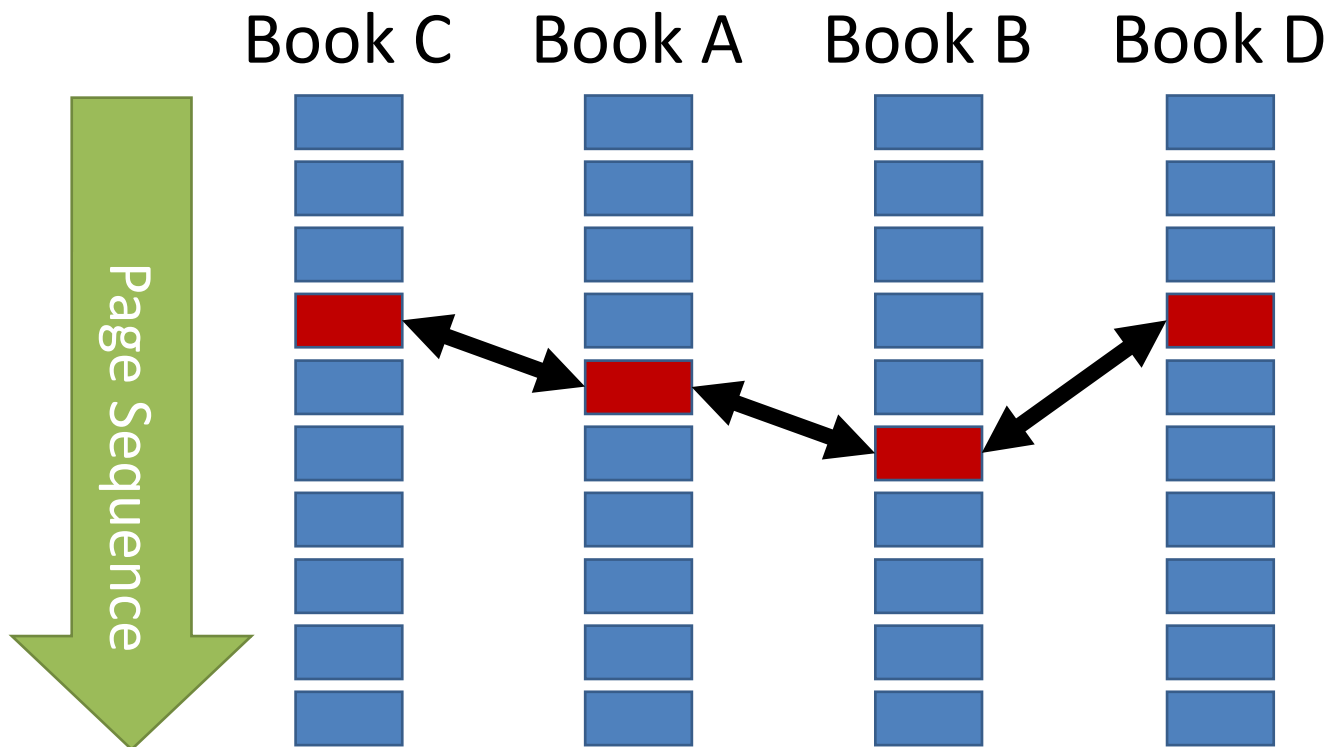
# Page Collation – Keypoint Matching

# Book Collation – Stable Marriage Algorithm



| Book A |
|--------|
| 1 |
| 2 |
| 3 |
| 4 |
| 5 |
| 6 |
| 7 |

| Book B |
|--------|
| 1 |
| 2 |
| 3 |
| 4 |
| 5 |
| 6 |
| 7 |

| Score |
|-------|
| 0 |
| 5 |
| 10 |
| 4 |
| 6 |
| **50** |
| 8 |

1:CBEAD

2:ABECD

3:DCBAE

4:ACDBE

5:ABDEC

A:35214

B:52143

C:43512

D:12345

E:23415

# Woodblock Tracking
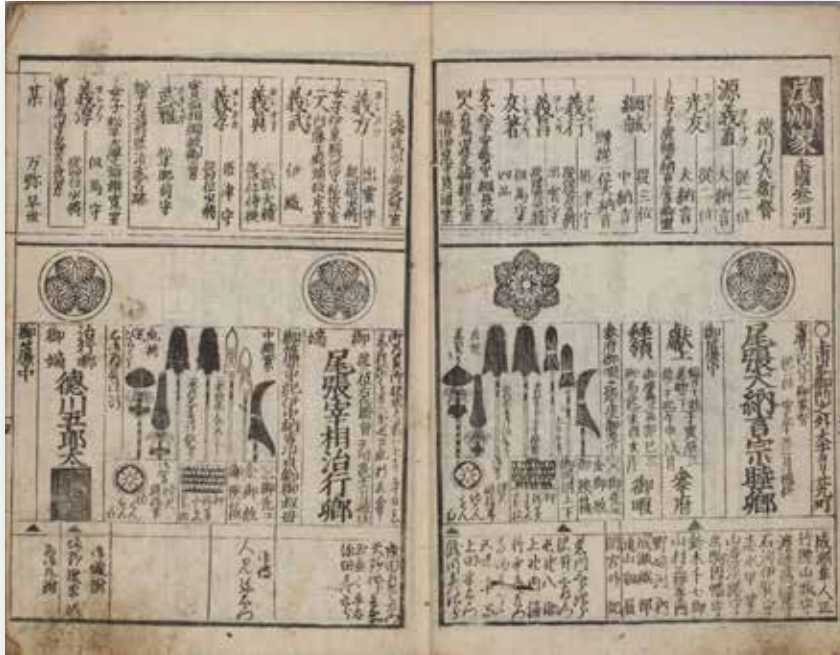
Page Sequence

Book C Book A Book B Book D

1. A page pair suggests the usage of the same woodblock across books.

2. A bottom-up process to **connect page pairs**.

3. **Differential reading** focuses on the **diachronic change of data** occurred on the same woodblock.

# Application to "Bukan Complete Collection"

# Bukan Complete Collection Project

http://codh.rois.ac.jp/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 families and personnel in the Edo Bakufu with structured data and graphical elements.

2. Published for 200+ years before 1867 as long-seller books with practical usage.

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

4. The project creates the database of synchronic and diachronic structured data from **381** different editions of Bukan.

# Bukan Differential Reading Platform

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



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

4. Color suggests the reliability of the collation (red: low, green: middle, blue: high).

# Page Collation

The interface emphasizes difference between page images with **navigation links to move across books**.

**vdiff.js**: JavaScript-based visual differencing tool to compare two images in 4 different modes.

**vdiff-seq.js**: JavaScript-based visual differencing tool for sequential images.

# Statistics

| Item | Number |
|---|---|
| Books | **336** |
| Images | **143,616** (111,114 portrait 32,502 landscape) |
| Keypoints | **70,417,952** (490 keypoints per image) |
| Tested book pairs | **1,326** |
| Tested page pairs | **87,110,176** (65,694 page pairs per book pair) |
| Married page pairs | **545,869** (about 0.63% of tested page pairs) |
| Woodblocks | **40,662** |

# Woodblock Tracking and Differential Reading

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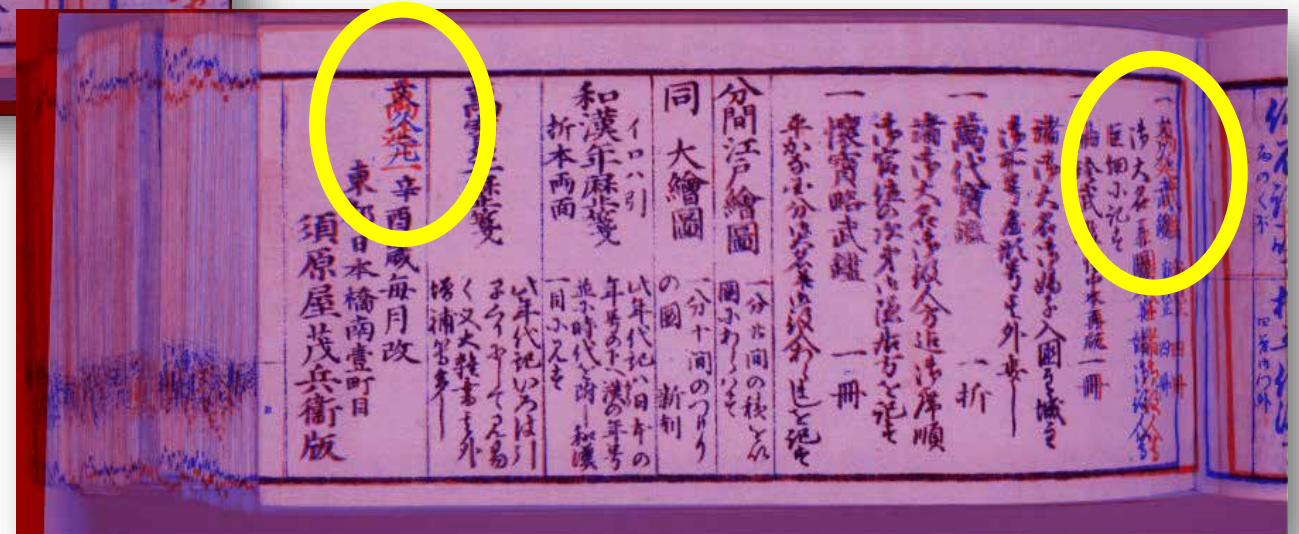


When the **era name** changed, they only **corrected** the era name.

Other parts remain the same, so you don't need to transcribe them again.

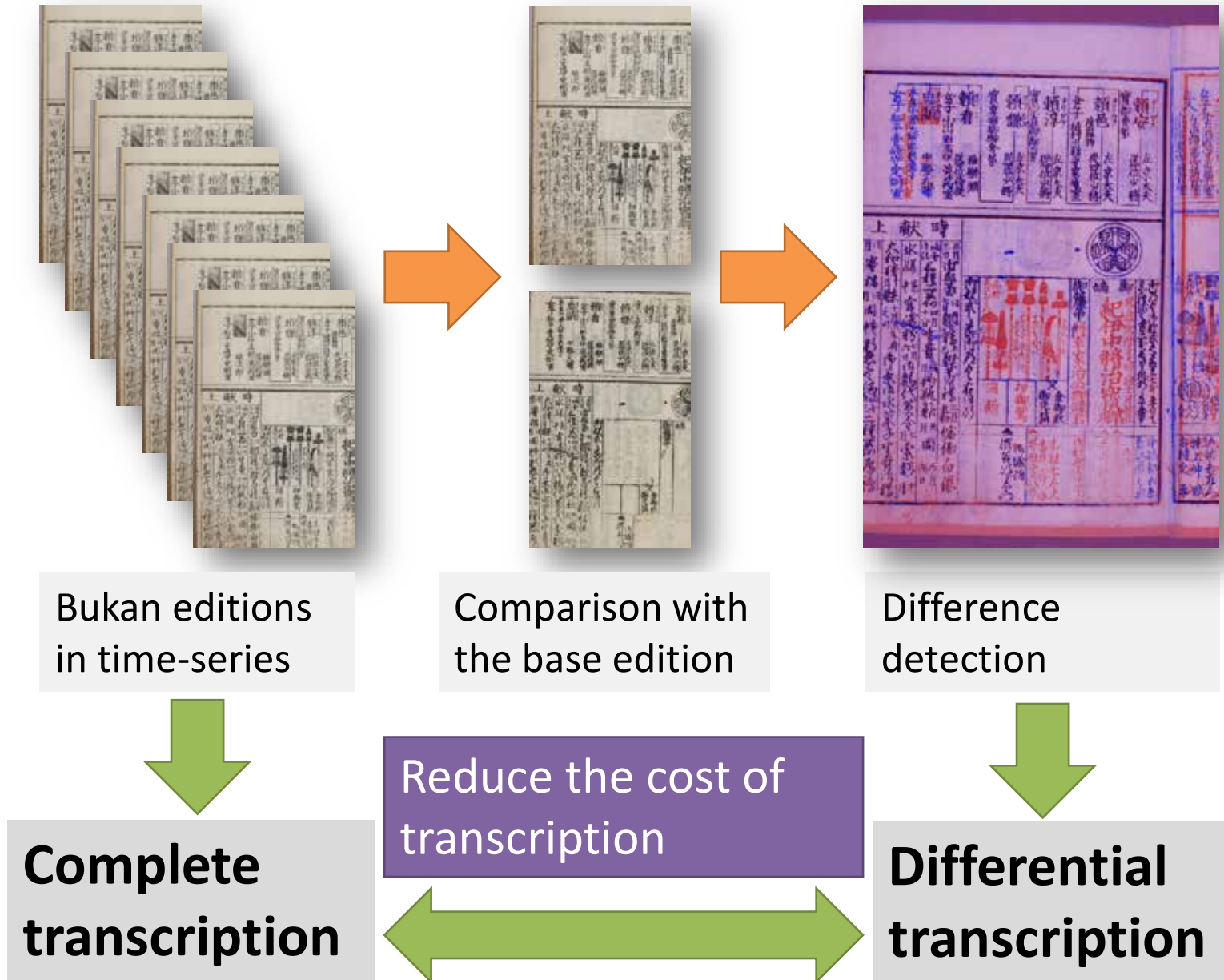# Distribution of the Life of Woodblocks



Page Sequence

Long-lived woodblocks: 48 versions

Publication Years

Digital Humanities 2022

# Differential Transcription

Bukan editions in time-series

Comparison with the base edition

Difference detection

Reduce the cost of transcription

**Complete transcription**
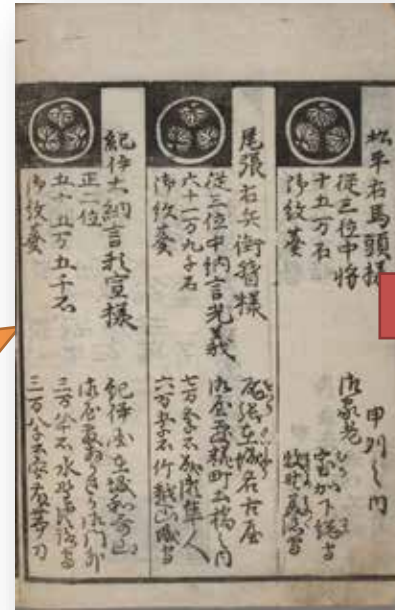
**Differential transcription**

1. **Complete transcription**: all text is transcribed.

2. **Differential transcription**: only changes are updated.

3. **Recreation detection**: woodblock recreation requires a new complete transcription.
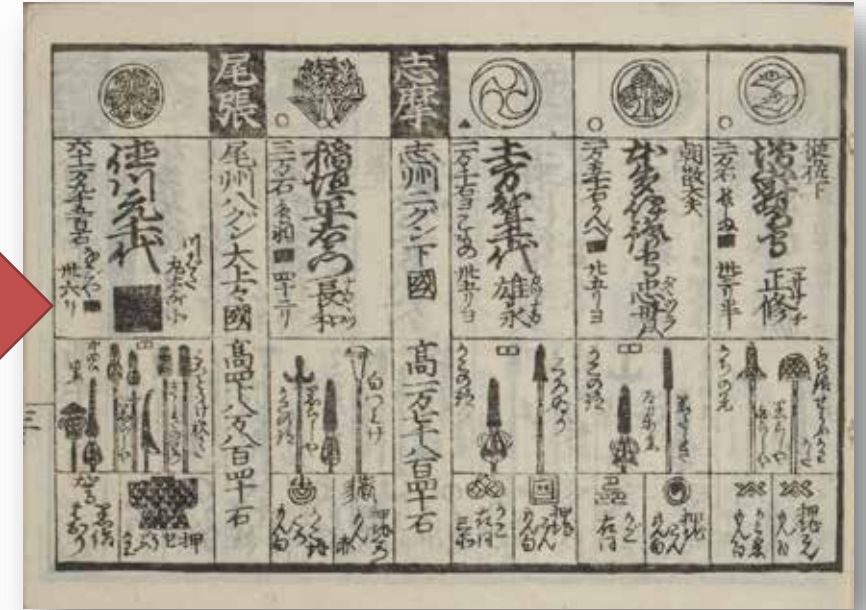
# Differential Transcription for Diachronic Data Structuring

**Research Question**: How **a data element of each daimyo family** has changed over 200 years in the Edo period?

**Differential transcription** was successful for about **85%** of the book pairs, especially in later years, thanks to the usage of the same woodblock.
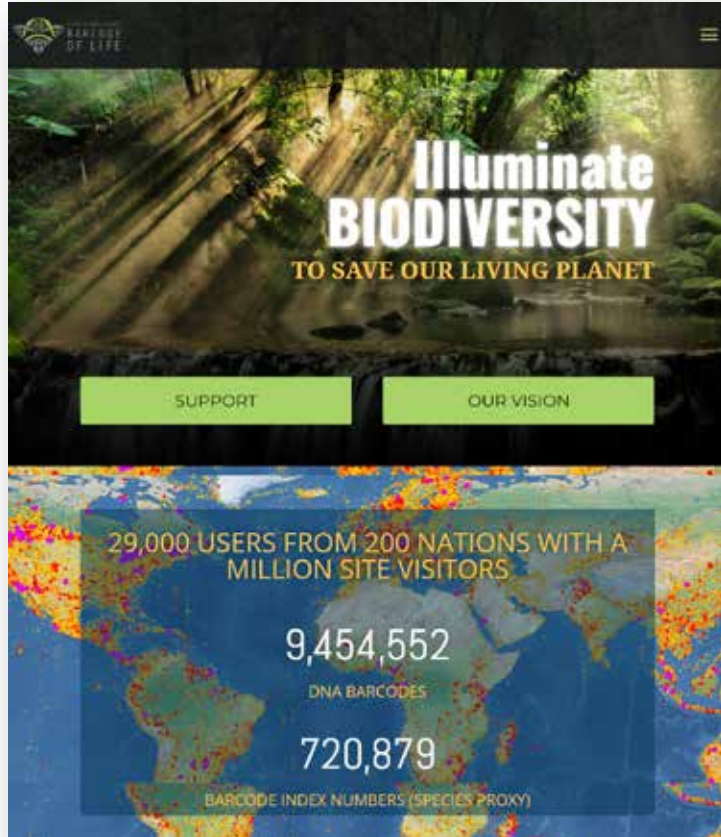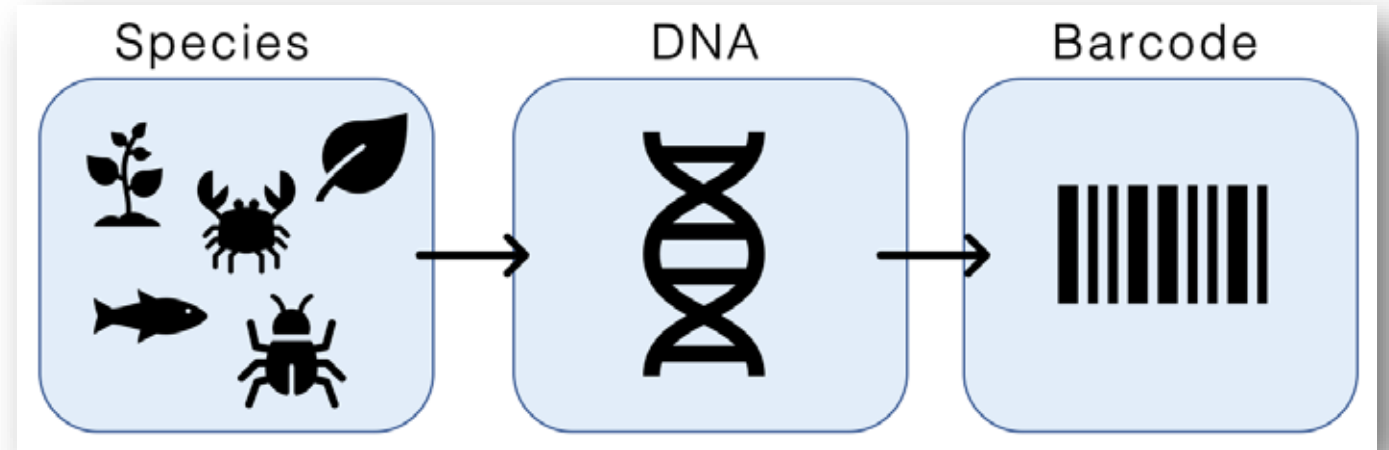


Mondukushi 1658



Okuniwake Bukan 1868

The failure of woodblock tracking is caused by **missing volumes, missing pages, and woodblock recreation**.

# Toward Book Barcoding Platform

# DNA Barcoding



Species DNA Barcode

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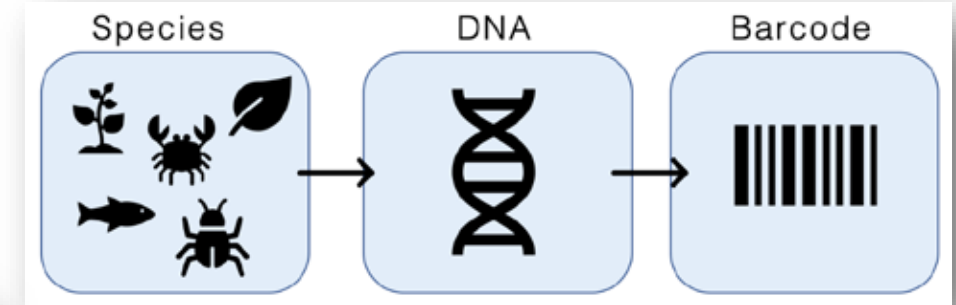
1. DNA barcodes are unique DNA sequences to assign identities to sequences of unknown origin.

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

# Book Barcoding



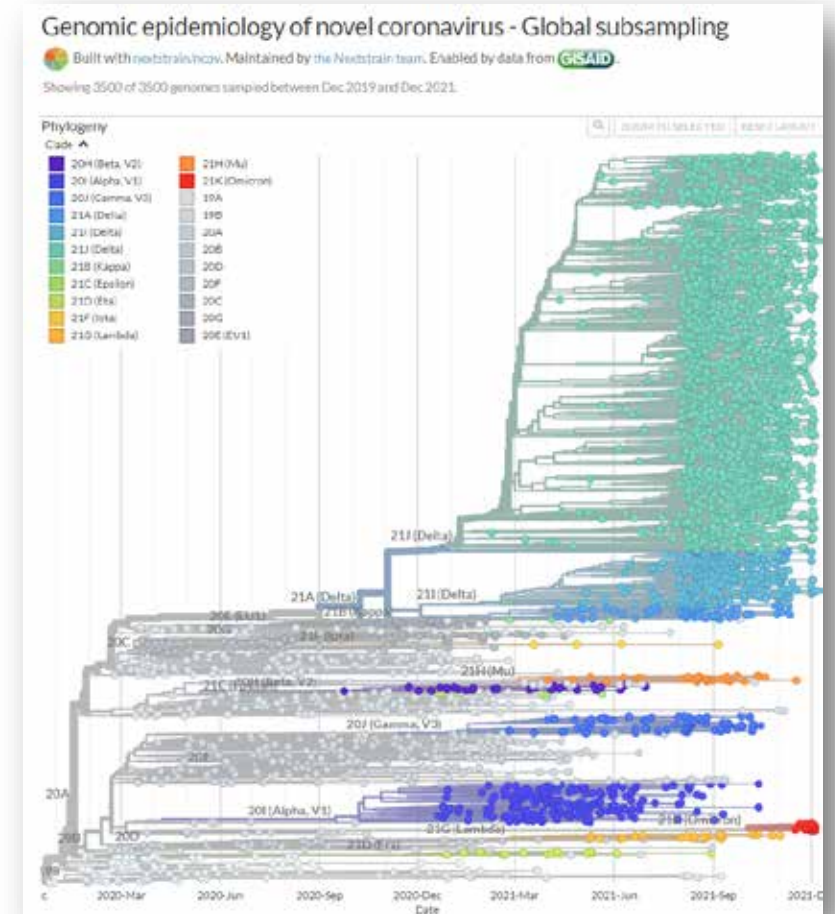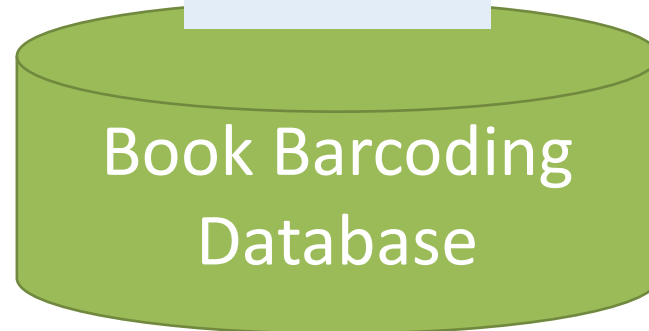Digital Image

Keypoints

Species → DNA → Barcode

Barcode

# Book Barcoding Platform for Japanese Woodblock-printed Books



Compare the barcode of unknown books with entries in the DB.

Book Barcoding Database

Genomic epidemiology of novel coronavirus - Global subsampling

https://nextstrain.org/ncov/gisaid/global

# Summary

1. We proposed image-based comparison and differential reading, suitable for Japanese woodblock-printed books.

2. We built a large-scale book collation algorithm to automate book and page collations.

3. We applied the algorithm to Bukan Complete Collection for the differential transcription of diachronic data.

4. We named the algorithm as book barcoding after the idea of DNA barcoding.

# Acknowledgment

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

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

- We thank **people in AMANE LLC** to transcribe Bukan Complete Collection and give suggestion to the differential reading platform.

- 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/
- **vdiff-seq.js** http://codh.rois.ac.jp/software/vdiffseqjs/