# AI and Data Science - Machine Learning as Digital Catalyst for Data Curation

KITAMOTO Asanobu

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

National Institute of Informatics

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

# Humanities Data

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

Our team consists of 1 professor, 4 post-docs, and 5 appointed professors.

**2016** Pre-center started.

**2017** Officially launched.

**Member**: One director and four project researchers (NII and ISM).

**Direction 1**: **Innovate humanities research** by computer science and statistical technologies and tools.

**Direction 2**: **Innovate non-humanities research** by data and questions from humanities.

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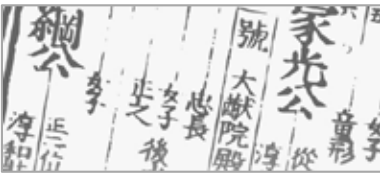
http://codh.rois.ac.jp/dataset/

Dataset of Pre-Modern Japanese Text

Kuzushiji Dataset

Dataset of Edo Cooking Recipes

Bukan Complete Collection

Collection of Facial Expressions

Dataset of Historical Administrative Boundaries

# How to Access Humanities Data?

Humanities data are mainly textual data, but visual and spatial data requires metadata and annotation to enable **deep access** to content.
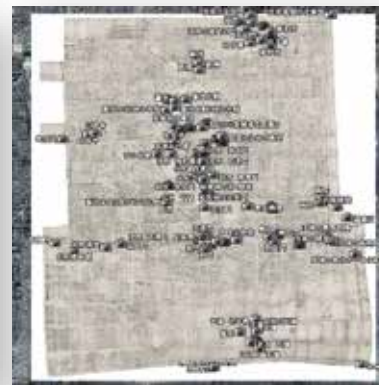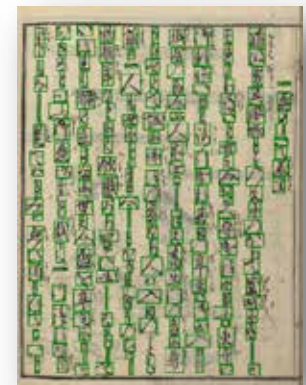
**Images**          **Photographs**          **Maps**          **Characters**

# Knowledge Representation

**Google Dataset Search – Schema.org**

https://toolbox.google.com/datasetsearch

**DataCite - DOI**

https://search.datacite.org/

Interoperable metadata and Semantic Web can increase findability.

# Manual Image Annotation

## https://tropy.org/

# Machine Learning (ML)



Credit: David Stanley,
https://www.flickr.com/photos/davidstanleytravel/

Tags ?

Coconut    Grove    Beach

Resort    Elmina    Ghana

landscape    shore

seaside    outdoor    coast

sand

Gray: Human annotated tags, White: Machine annotated tags.

# Machine Learning for Photographic Database

# Best Practices for
# AI-assisted Data Curation

1. **What could be done by AI, and not by AI?** Hype and criticism should be corrected.

2. **Machine learning**: especially effective for learning patterns from image data.

3. **Images, especially photographs**: selected as the initial target of the work.

4. **General numerical datasets**: content-based access is still a challenge.

# Open Images Dataset V2

https://github.com/openimages/dataset/blob/master/READMEV2.md



Annotated images from the Open Images dataset. Left: FAMILY MAKING A SNOWMAN by mwvchamber. Right: STANZA STUDENTI.S.S. ANNUNZIATA by ersupalermo. Both images used under CC BY 2.0 license.

# Deep Learning Model



https://medium.com/@siddharthdas_32104/cnns-architectures-lenet-alexnet-vgg-googlenet-resnet-and-more-666091488df5

1. ResNet 101 classifier learns 5000 tags from 9 million images (Open Images Dataset V2).

2. We used the model already trained on general photographs, not on our dataset.

# Case 1: Ethnology Field Work



1. Field work in Ghana, in August 2017.

2. About 3,700 photographs, yet to be released to the public.

3. Collaboration with National Museum of Ethnology (Prof. Yoshida, Prof. Iida and others).

# Tag: Person

# Tag: Food

# Tag: Beer

# Tag: Art

# Case 2: Archaeology Field Work
http://dsr.nii.ac.jp/photograph/



1. Photographs of the Silk Road, mainly about old ruins.

2. 6,129 photographs across long time span.

3. Many photographs were taken by Dr. Nishimura in Toyo University.

# Tag: Archaeological Site

# Tag: Snow

# Tag: Wood

# Different Responses from Users

| Ethnology Field Work | Archaeology Field Work |
|---|---|
| Image tagging has great potential for grouping photographs by theme. | Grouping by machine-generated tags is less useful than grouping by entity names. |
| Even if the tag is not correct, it gives some hints about the content. | Some tags are simply wrong due to different training images and domains. |

# Why Different Reponses?

| General Noun Metadata | Proper Noun Metadata |
|---|---|
| Ethnology photographs are so diverse that tags of general nouns are effective for grouping. | Archaeology photographs are usually taken with intentions. |
| It motivates experts to describe deeper metadata. | Entity names are difficult to identify by machine learning. |

# Case 3: Post-Disaster Survey



1. Photographs of East Japan Earthquake 2011 and Kumamoto Earthquake 2016.

2. More than 10,000 photographs, yet to be released to the public.

3. Collaboration with National Research Institute for Earth Science and Disaster Resilience (NIED).

# Serendipity
# Tag: Stadium

# Semantic Grouping of Low-Level Tags

Vehicle 55080  Transport 54615  Road 36838
Residential area 31354  Infrastructure 30770  Wall 28830
Tree 28183  Asphalt 27873  Lane 23916  Soil 23566
Waterway 21086  Rural area 20360  House 19245  Highway 18228
Sea 17242  Town 16648  Coast 16262  Urban area 16186
Building 15509  Car 15174  Neighbourhood 14971  Property 13790
Home 12801  Road surface 12307  Facade 11968  Art 11041
Walkway 11012  Plant 10236  Sport venue 10224  Track 10124
Mountain 9916  Hill 9788  Land vehicle 9671  Street 9633
Village 9491  Automotive exterior 9387  Shore 8986  Driving 8904

Domain experts need higher-level semantic grouping of low-level tags.

ArtificialObj 107508    HumanAct 93060    Natural 38392

# Case 4: Historical Photographs

http://codh.rois.ac.jp/north-china-railway/



1. Photographs of North China Railway, a company existed around 1940.

2. More than 35,000 photographs will be released in Feb. 2019.

3. Collaboration with Kyoto University.

# Image Tagging

Road, Street, Black-and-white, Monochrome photography, Monochrome, Infrastructure, Transport, Lane, Vehicle, Photograph

# Image Colorization

# Image Tagging after Colorization



Road, Street, Infrastructure, Town, Transport, Photograph, Urban area, Vehicle, Lane, Pedestrian

# Lessons from Two Collections

1. Two photographic collections are too large for humans to annotate one by one.

2. Automatic tagging may be useful as the initial step for improving findability.

3. Statistical research questions, such as thematic distribution may be answered.

4. Other methods can improve findability, such as colorization and object detection.

# The Value of Data and FAIR Principle

# The Value of Data

**1. Intrinsic Value**
Raw data
**scientists / scholars**

**2. Basic Value**
Organized data
**(data) librarians**

**3. Added Value**
Integrated data
**(data) curators**

**4. Persistent Value**
Preserved data
**(data) archivists**

# Machine Learning for Increasing the Value of Data

1. Basic value and added value need high quality metadata for higher value.

2. FAIR (Findable, Accessible, Interoperable, Reusable) principle asks for good metadata.

3. Humans procrastinate in adding metadata, hence the workflow does not start.

4. Use machines to quickly reach **a state which is better than nothing.**

# Digital Catalyst



https://commons.wikimedia.org/wiki/File:CatalysisScheme.png

To reach a state of curated data, we need to go beyond the high energy barrier.

**Machine learning as digital catalyst** reduces the barrier, requiring less human motivation to pass the barrier.

# Human-Machine Collaborative Workflow

1. **Machines** can automatically add general noun tags for coarse grouping.

2. **Humans** can manually add proper noun tags for fine meaning as metadata.

3. **Domain experts** can add high-level metadata and semantic grouping.

4. **ML models** can use added metadata as new training data to improve performance.

# Conclusion

1.  Machine learning, e.g. image tagging, is <span style="color:red">beneficial for improving findability</span>.

2.  <span style="color:red">General nouns</span> are useful for some apps; other apps require higher level metadata.

3.  Better findability (curation) increases the <span style="color:red">basic value</span> and <span style="color:red">added value</span> of data.

4.  <span style="color:red">Digital catalyst</span> is a concept of machine-assisted data curation to motivate humans.

# Acknowledgment and Links

Photograph collections were provided from the following collaborators:

Dr. Taku Iida in National Museum of Ethnology

Dr. Yoko Nishimura in Toyo University

Ms. Hinako Suzuki in National Research Institute for Earth Science and Disaster Resilience

Dr. Toshihiko Kishi and his colleagues in Kyoto University.

A part of the machine learning workflow was developed by:

Hoàng Văn, Hà (Vietnam National University, HCMC) during NII internship.

- Center for Open Data in the Humanities
  - http://codh.rois.ac.jp/
- Open Science
  - http://agora.ex.nii.ac.jp/~kitamoto/research/open-science/
- Researchmap
  - http://researchmap.jp/kitamoto/