Developing new library services using AI (machine learning) : an introduction to the Next Digital Library

Tahee ONUMA, Chief, R&D Office National Diet Library, Japan 18 September 2021



# The R&D Office

### Founded in October 2011

- Conducts research into cutting-edge technologies in collaboration with researchers from outside the NDL
- Releases experimental or newly-developed services on the NDL Lab website(=> next slide), such as:
  - The Next Digital Library \*today's main topic
  - ► NDC Predictor
  - ► Japan Search

. . .

Bibliographic Information Retrieval and Visualization System



# The NDL Lab

Website for conducting field trials of new library services

First launched in May 2013, renovated in March 2020

Recent activities focus on <u>machine learning</u>



#### ピックアップ



#### https://lab.ndl.go.jp/



# Why machine learning?

Machine learning helps improve

searchability based on content Keyword full-text search and illustration search readability of digitized images Bleaching, automatic image processing Functionality like this is implemented and tested on the Next Digital Library.



# The Next Digital Library

- Launched in 2019 by the R&D Office
- A testing ground for experimental functionality developed by the NDL Lab
- Can be used to search the NDL Digital Collection for documents in the public domain
- Two search modes:
  Keyword full-text search
  Illustration search





# Keyword full-text search

30,000 items on industrial and commercial subject matter

(Category 6 of the Nippon Decimal Classification)

Searches full texts generated by OCR

Identify and access individual pages that include keywords



# Illustration search (1)

#### Search for similar images

- 23,000,000 public-domain images from 336,000 digitized books, rare books, and historical materials in the NDL Digital Collection
- Increased potential of information exploration that goes beyond keyword search







# Illustration search (2)

 Image extraction: deep learning method "semantic segmentation" → Will change in the near future
 Feature extraction: trained using ImageNet dataset
 Search engine: vdaas/vald (https://github.com/vdaas/vald)



Blue regions indicate areas extracted to illustrations.



# **Background bleaching**

- Whitens pages discolored due to aging
- Improves readability
- Based on Neural Network model "pix2pix:GAN"

		-
をできてのおいて見上う。	<del>.</del>	
七雷災防止		- North
昭和十五年、日本御梅装護者の中に、常見訪	の非委員商といふものが指来た。日本	•
である。その研究は雷の客を選けること、即ちである。その研究は雷の客を選けること、のようなに関係した事者が数十九も集つて、	「一五式協力して常の研究を始めたの	
そんなごとを聞くと、リフンクリン以来、厳富	「針といふものが出来てゐるのに、今	
頃になって何を研究するのか、と不忌義の思ふ	人があるか苦しれない。しかし、今	27.
くらも読みれてわるといふのが、いなむことの	できない事實ならである。	
利へば、経営針蛇してち、鹿翅に大丈夫で安	みだといふものはまだないのである。	
2010年の日本の日本市場である。その他にも、4月20日にある人々の間でさべる。これがそれを研究してある人々の間でさべる。これが、4月20日にある。4月20日にある。4月20日にある人々の他にあった。	業通信にあける容量の妨害は害と蜜様 か良い、あれが良いといろいろ塗った	2.2
	精が汚ぐなりたと、発信が出来た	
	* せねだといふ場合がたくさん渋る。	-
	る被害の問題なども、非局病論。	
N'CONSTRUCTION OF THE OWNER OF	重要に在の工事でわる。それで明	-
呼撃が崩歩すると、いろいろの自然現象のを空しなければならないことは、あとからぬとか	4件が後々とわかって来る。それで施 からといくらでも撮きないのである。	1 have
Π <del>2</del> 13 4	N. N. N.	





## Automatic image splitting

Automated splitting of pages and removal of margins
 Optimizes browsing on smartphones and tablets



#### **Display on smartphone**





## Further application to other services

### **Research for the Next Digital Library**



Library services from the National Diet Library



# Japan Search

- A national, integrated search platform for digital archives in Japan
- Launched August 2020
- Search 22.8 million metadata from 131 databases (August 2021)
- Similar Image Search function using technology tested on the Next Digital Library



CC BY(表示) 九州国立博物館



https://jpsearch.go.jp/



### GitHub

Provides access to programs (CC BY) and datasets (PD) created by the R&D Office

A hub for the exchange of expertise with engineers from all over the world

https://github.com/ndl-lab





# Ongoing projects

Development of OCR software

- Preparing OCR-generated texts of digitized materials on the NDL Digital Collections
- Providing access to high-quality datasets

