

国立国会図書館のOCRテキスト化事業と 全文テキストデータの提供について

データ編

国立国会図書館 電子情報部 電子情報企画課
青池 亨



(再掲) 2つのOCR事業 (2021年度)

1. デジタル化資料のOCRテキスト化

- ・ 2020年時点でデジタル化済のほぼ全ての図書・雑誌等 **約247万点**

2. 日本語OCR処理プログラム (NDLOCR) の研究開発

- ・ オープンソースで利用可能なOCR処理プログラムの開発
- ・ 当館が今後デジタル化する資料ものは、これを使ってテキスト化を実施予定



国立国会図書館デジタルコレクション



(※)各事業の詳細はNDLラボを参照

https://lab.ndl.go.jp/data_set/ocr/

成果物

1. 大量の全文テキストデータ + 学習用データセット
2. 日本語OCR処理プログラム (NDLOCR) + 学習用データセット

提供データセット活用のご案内

- 当館の資料を活用した研究に

著作権保護期間満了資料のみを対象にウェブブラウザからさっくりお手軽に活用

- ・ 全文テキストダウンロード機能（次世代デジタルライブラリー）
- ・ 全文中のキーワード出現頻度情報のダウンロード機能（NDL Ngram Viewer）

利用申請手続きを行って、当館と協議のうえで著作権法上認められた範囲内で大規模に活用

- ・ テキストデータの機械学習目的利用
- お手元のデジタル化資料のテキスト化に
 - ・ NDLOCR（オープンソース）の活用
- 新しいOCRソフトウェアの開発に（※本日は説明しません）
 - ・ OCR学習用データセット

次世代デジタルライブラリー 全文テキストダウンロード機能

The screenshot displays the NDL Lab digital library interface. On the left, a sidebar titled '帝国図書館一覧' (Imperial Library Overview) contains metadata for a book: '責任表示' (Responsibility Statement), '出版年' (Publication Year) 1912, '出版者' (Publisher) 帝国図書館 (Imperial Library), '請求記号' (Call Number) UL214-E10, and a link to 'デジタルコレクションで見る' (View in Digital Collection). The main area shows a two-page spread of a document titled '帝國図書館一覧' (Imperial Library Overview). The right page contains a table of contents. At the bottom, a navigation bar includes controls for page number (7/25), zoom, and a download menu. The download menu is highlighted with a red box and a red arrow, showing two options: 'この資料の全文テキストデータ' (Full-text text data of this material) and 'この資料の画像データ (IIIF API経由)' (Image data of this material (via IIIF API)).

ダウンロードボタン
(画像・テキストデータ)

<https://lab.ndl.go.jp/dl/book/1907912?page=7>

次世代デジタルライブラリー 全文テキストダウンロード機能

ダウンロード可能なテキストデータの形式

改行なしtxt形式、画像上の矩形座標情報付きjson形式

それぞれコマごとのファイルと、1行1コマで連結したファイル
がzipファイル内に含まれている

json形式の例

```
[{"id":3,"contenttext":"帝國圖書館官制・  
","xmin":4071.0,"ymin":762.0,"xmax":4147.0,"ymax":1259.0},  
{"id":4,"contenttext":"帝國圖書館職員・  
","xmin":3955.0,"ymin":764.0,"xmax":4030.0,"ymax":1259.0},...
```

NDL Ngram Viewer キーワード出現頻度情報のダウンロード機能

<https://lab.ndl.go.jp/ngramviewer/?keyword=.{2,3}図書館&size=100&from=0>

「○○図書館」または「○○○図書館」という条件に合致するキーワードを、ヒット件数の多い順に列挙



ダウンロードボタン

キーワード	総出現頻度	
付属図書館	20310	次世代デジタルライブラリーで検索
帝国図書館	20227	次世代デジタルライブラリーで検索
国会図書館	19334	次世代デジタルライブラリーで検索
学校及図書館	19253	次世代デジタルライブラリーで検索
県立図書館	17409	次世代デジタルライブラリーで検索
学校図書館	15953	次世代デジタルライブラリーで検索
市立図書館	13778	次世代デジタルライブラリーで検索
東京図書館	13265	次世代デジタルライブラリーで検索
上野図書館	12241	次世代デジタルライブラリーで検索

NDL Ngram Viewer

キーワード出現頻度情報のダウンロード機能

検索結果について、
タブ区切りテキストになった出版年代ごとの頻度情報が取れます

Keyword	Total	Frequency	1801	1802	1803	1804	1844	1845	1846	1847	1848	1849	1850
付属図書館	20310	0	0	0	0	0	0	0	1	0	7	8	1
帝国図書館	20227	0	0	0	0	0	1	0	6	0	23	53	99
国会図書館	19334	0	0	0	0	0	248	322	273	283	321	249	331
学校及図書館	19253	0	0	0	0	0	0	0	53	312	148	155	358
県立図書館	17409	0	0	0	0	0	0	0	1	0	0	0	0
学校図書館	15953	0	0	0	0	0	0	6	6	31	17	8	61
市立図書館	13778	0	0	0	0	0	0	1	1	1	0	1	0
東京図書館	13265	0	0	0	0	0	69	22	74	161	287	189	338
大学図書館	9494	0	0	0	0	0	1	3	2	13	11	57	7
公立図書館	7906	0	0	0	0	0	1	1	3	1	3	0	0
幼稚園図書館	7559	0	0	0	0	0	0	0	5	27	93	364	142
私立図書館	7085	0	0	0	0	0	0	0	3	0	3	2	3

(中略)

検索結果をお手元で分析することができます

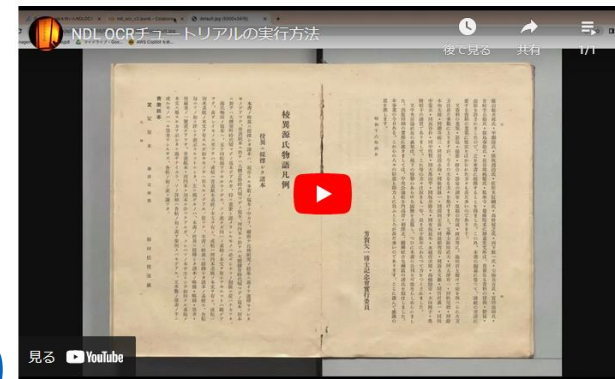
NDLOCR

- 当館が今後デジタル化する資料をテキスト化するために開発
- 明治期のような古い活字資料もテキスト化可能なOCR
- オープンソースとして館外にも公開・提供

https://github.com/ndl-lab/ndlocr_cli

使い方のチュートリアル動画（東京大学 中村寛先生作）

<https://zenn.dev/nakamura196/articles/af12c5fc18ab90>



開発の知見を生かしつつ、古典籍資料をテキスト化するOCR（NDL古典籍OCR）を内製開発。
年内にオープンソースとしての公開を予定

データシート（技術面に興味のある方向け）

資料種別/年代ごとのOCRの性能や、事業における性能改善過程の情報をまとめて公開しています

次世代デジタルライブラリーやNDL Ngram Viewerに利用しているOCRテキストデータの情報

(https://lab.ndl.go.jp/data_set/ocr/r3_line/)

NDLOCRの開発に関する情報

(https://lab.ndl.go.jp/data_set/ocr/r3_morpho/)