

KU-ORCAS国際シンポジウム「デジタルヒューマニティーズ推進のための環境構築とその課題」

歴史ビッグデータ研究基盤のための デジタルツールと相互運用性



北本 朝展（ROIS-DS人文学オープンデータ共同利用センター／国立情報学研究所）

<http://codh.rois.ac.jp/>

はじめに

ROIS-DS人文学 オープンデータ 共同利用セン ター (CODH)



2016年4月～

深化

研究者

巨大化

機械



市民

メンバー
国立情報学研究所
統計数理研究所
センター長＋
特任助教4名

多様化

「オープン」の概念を核として三者
を接続し、知識の深化、巨大化、多
様化を目指す

歴史ビッグデータの統合解析

過去のビッグデータを統合解析するための基盤技術を研究



歴史的資料
(史料)

自然科学的
データ

人文
社会的
データ

気候

地震

噴火

疫病

経済

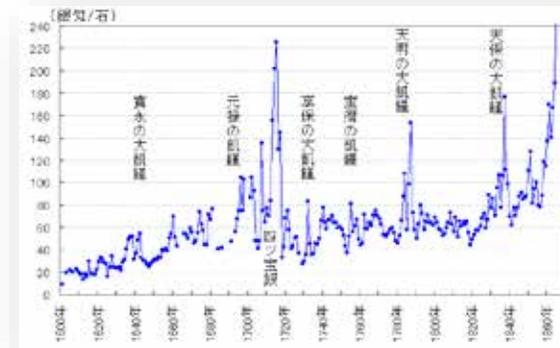
人口

政治

文化

データ
構造化
ワーク
フロー

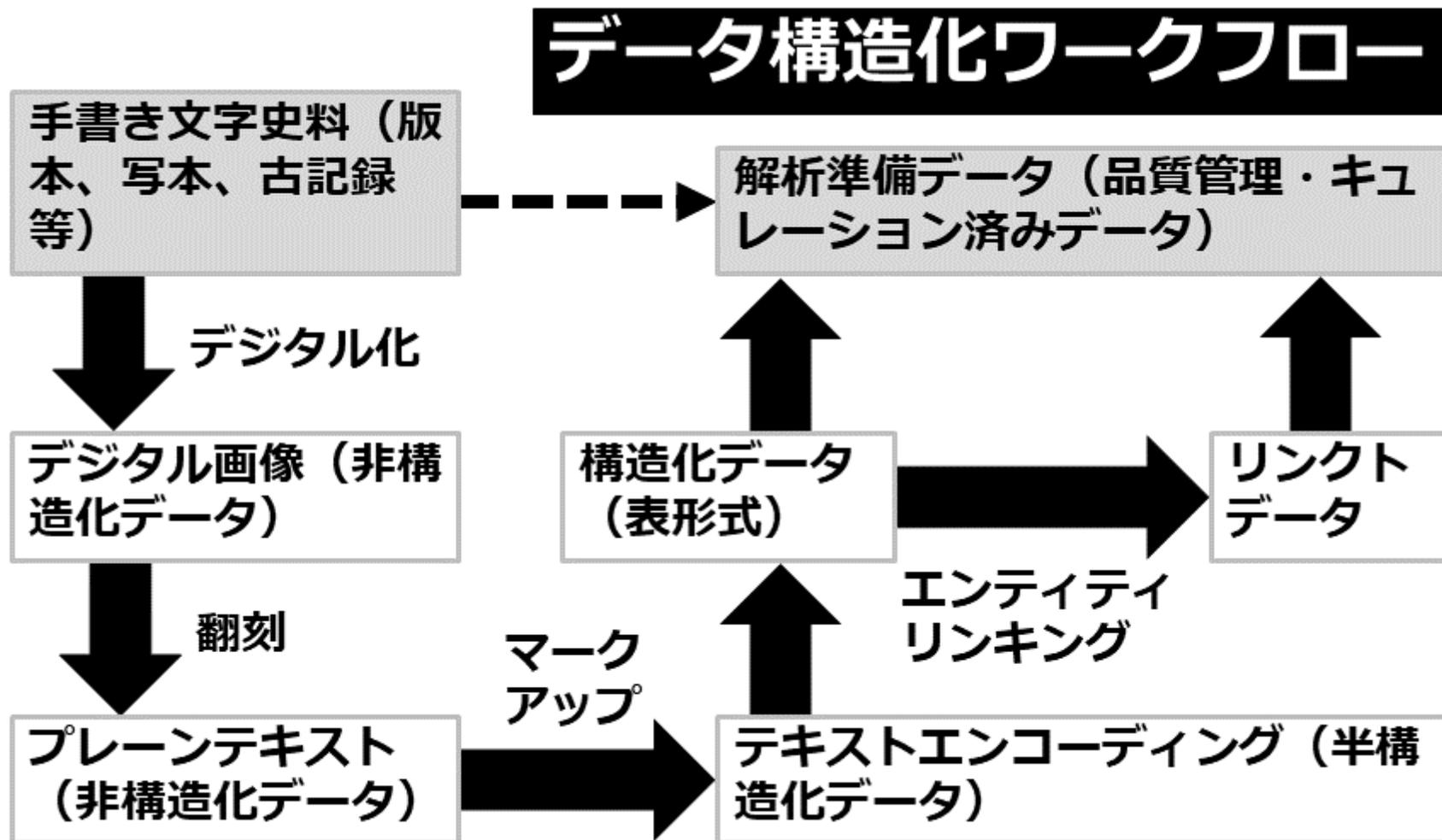
歴史ビッグ
データ研究基
盤 (機械可読
データ)



歴史ビッグデータとは？

1. 現代のデータを想定して開発されたビッグデータの方法論を、過去の記録に延長すること。
2. **生データの構造化**：機械可読なデータを生成し、様々な分野における解析などに利用する。
3. **多様な (Variety) データの統合解析**：単体のデータだけでは見えてこない新しい知見を得る。
4. **人文学研究のデジタルトランスフォーメーション**：研究のデジタル化には、技術要素の深化だけでなく、研究体制や (評価) 制度の変革も関わる。

データ構造化ワークフロー



ビッグデータを巡る言葉の整理

データ駆動型人文学

- データを基盤とする新しい手法により、**人文学の研究手法を革新**する。
- データの共有やエビデンスに基づく検証などの**新しい文化**を人文学に**導入**する。
- **人文学的に新しい知識**を得ることを目標とする。

人文学ビッグデータ

- 人文学から生み出される大規模データにより、**(非)人文学の研究手法を革新**する。
- 人文学データに特有の構造化や解釈などの**技術を分野を超えて共有**する。
- **定量的またはデータに基づく知識**を得ることを目標とする。

AIくずし字認識

NIJL-NWプロジェクト

@国文学研究資料館

<http://www.nijl.ac.jp/pages/cijproject/>



**300,000点の日本古
典籍（1868年以
前）を、デジタル
化し、オープン
データとして公開。**

日本文化のビッグ
データをどのよう
に活用するか？

過去の文化遺産と「くずし字」の問題

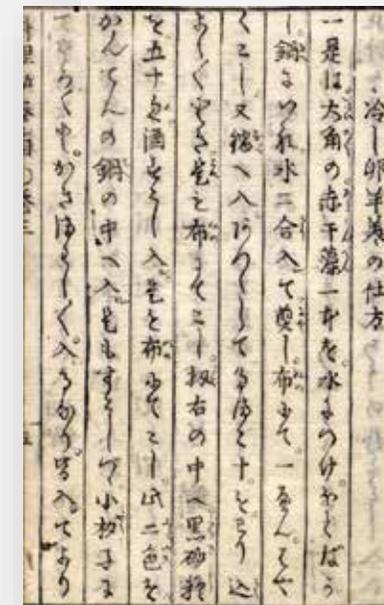


文書
数億点

日本に残された古典
籍や古文書の点数
(冊数) の推定

読者
数千人

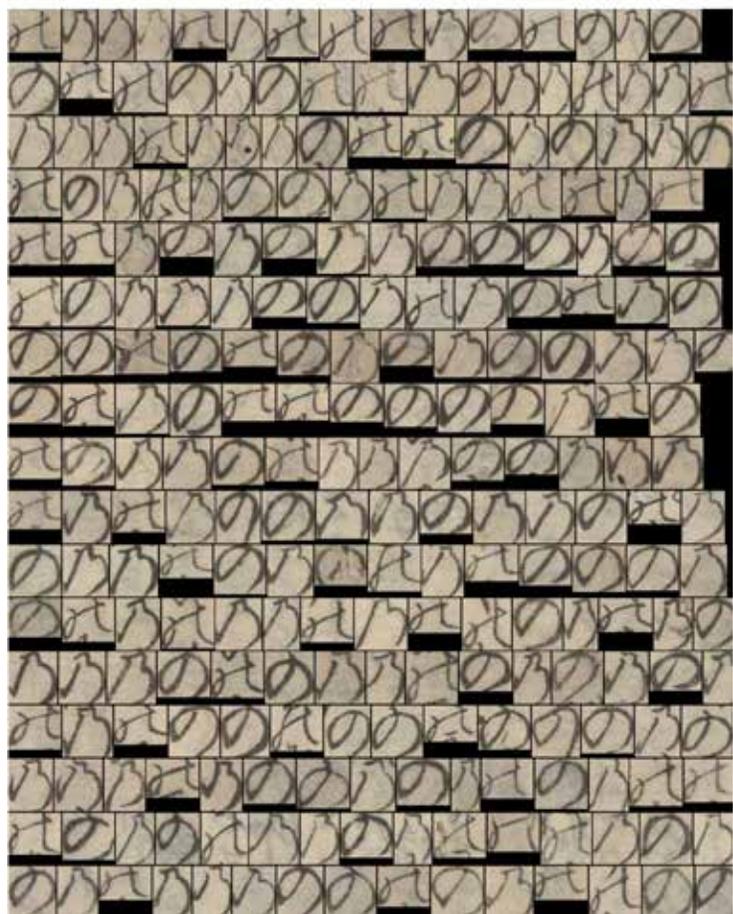
くずし字をきちんと
読める人数の推定
(全人口の0.01%)



くずし字データセット

<http://codh.rois.ac.jp/char-shape/>

雨月物語 (1890)



- 国文学研究資料館が作成、CODHが整理して公開するオープンデータ。
- 2020年7月現在
 - 文字種 = 4,328
 - 文字数 = 1,086,326
- ZIPファイルダウンロード→くずし字認識の訓練データとして利用。
- データセットの公開が、AIによるくずし字認識研究を活性化した。

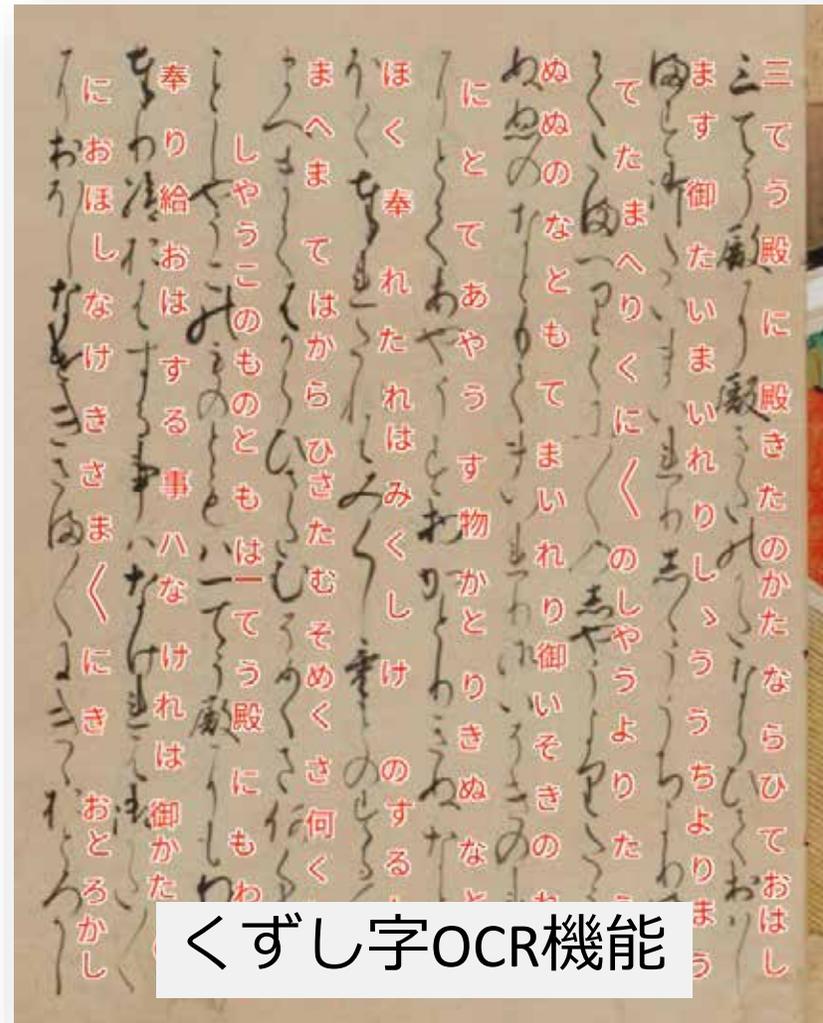
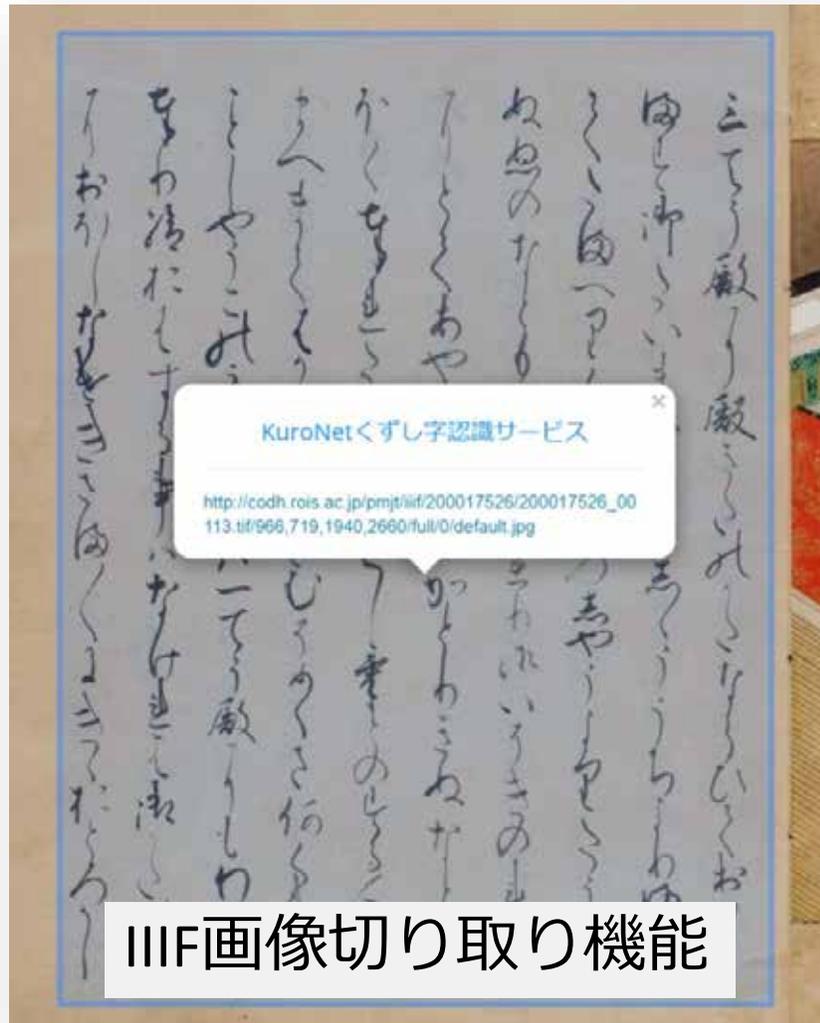
ディープラーニング
によるくずし字認識
手法KuroNetを開発
1枚の画像の認識に約
1秒
条件が良ければ精度
95%



三てう殿に殿きたのかたならひておはし
ます御たいまいれりしううちよりまう
てたまへりくにくのしやうよりたうき
ぬぬのなともてまいれり御いそきのれう
にとてあやうす物かとりきぬなとお
ほく奉れたれはみくしけのする人御
まへまてはからひきたむそめくさ何くれの
としやうこのものともはてう殿にもわかち
奉り給おはする事はなけれは御かたく
におほしなけきさまくにきおとろかし

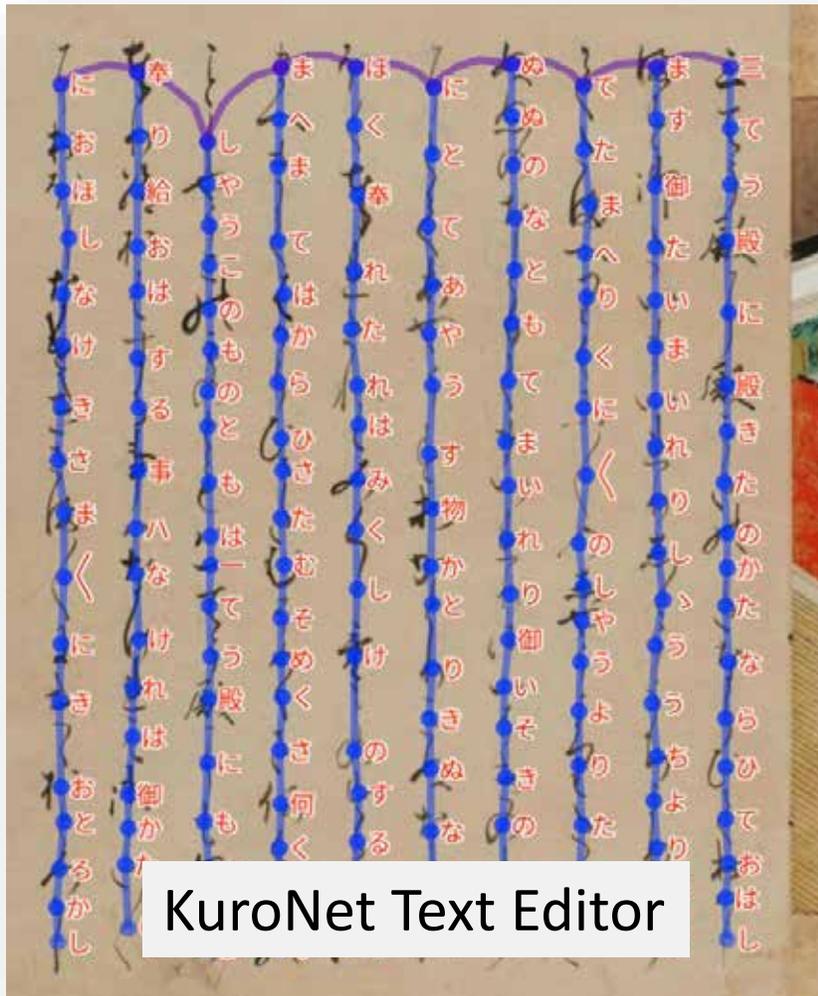
KuroNetくずし字認識サービス

<http://codh.rois.ac.jp/kuronet/>



KuroNetくずし字認識サービス

<http://codh.rois.ac.jp/kuronet/>



KuroNet Text Editor

テキスト化

三てう殿に殿きたのかたならひておはし
ます御たいまいれりしううちよりまう
てたまへりくにくのしやうよりたうき
ぬぬのなともてまいれり御いそきのれう
にとてあやうす物かとりきぬなとお
ほく奉れたれはみくしけのする人御
まへまてはからひさたむそめくさ何くれの
しやうこのものともは一てう殿にもわかち
奉り給おはする事ハなけれは御かたく
におほしなけきさまくにきおとろかし

テキスト化サービス

kaggle くずし字認識コンペ

<http://codh.rois.ac.jp/competition/kaggle/>



機械学習エンジニアが300万人以上登録する、**世界最大のデータサイエンスプラットフォーム Kaggle**にて、くずし字認識コンペを開催。

- **期間**：2019年7月19日～10月14日
- **参加チーム数**：293
- **参加者数**：338
- **結果提出回数**：2652

kaggle コンペの結果

上位の精度
は約95%

1. **さまざまな国の機械学習エンジニア、研究者が上位入賞。**
2. **日本語やくずし字の知識がなくても開発可能。**
3. **コンペの準備を滞りなく進めるには、**情報学者と人文学者の協働が不可欠。****

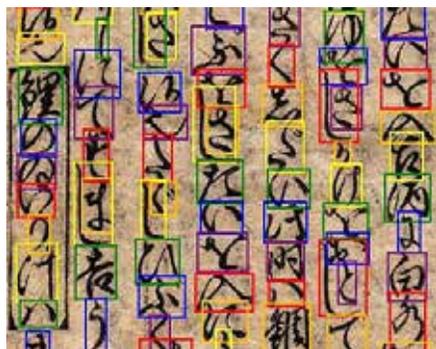
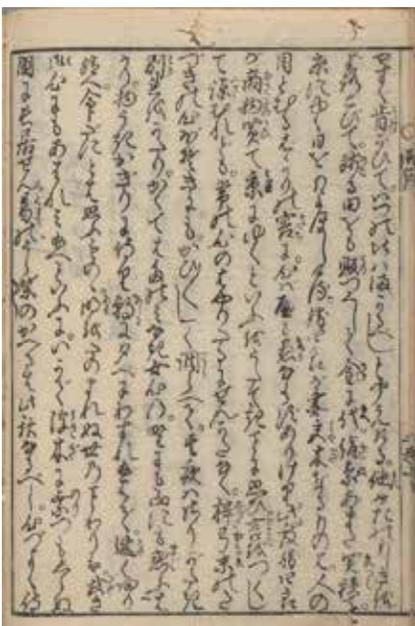
#	Δpub	Team Name	Notebook	Team Members	Score 🏆	Entries
1	—	tascj			0.950	13
2	—	Konstantin Lopuhin			0.950	60
3	—	Kenji			0.944	161
4	▲1	YoudaoOCR			0.942	49
5	▼1	See--			0.940	42
6	—	abc			0.939	15
7	—	K_mat			0.934	20
8	—	t-hanya			0.920	21
9	—	Ollie, Nanashi, and Tom			0.910	35
10	—	Zenkei_R&D			0.903	144
11	—	masayai			0.903	12
12	▲5	Kirill Brodt (shad nsk)			0.901	4
13	▲1	James Day			0.901	33
14	▼1	NEU			0.900	54
15	▼3	s tatsuya			0.900	29

AIくずし字認識

<http://codh.rois.ac.jp/char-shape/>

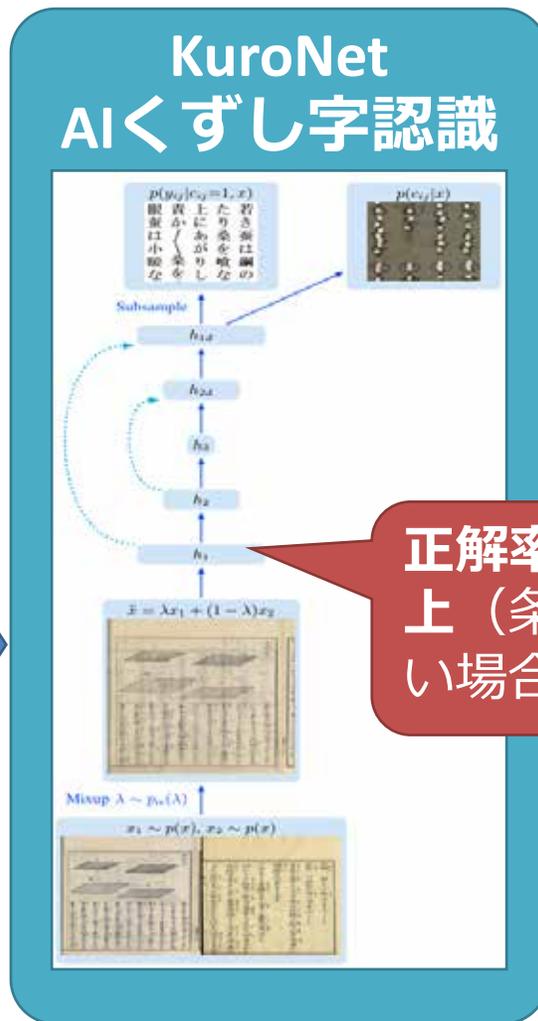
日本古典籍
データセット
(国文研蔵)

くずし字データ
セット (国文研・
CODH作成)



file	char	x	y
200003803_00024_2.jpg	U+3067	416	114
200003804_00024_2.jpg	U+3055	232	115
200003805_00024_2.jpg	U+304A	327	115
200003806_00024_2.jpg	U+3068	145	116
200003807_00024_2.jpg	U+3046	369	116
200003808_00024_2.jpg	U+305F	457	116
200003809_00024_2.jpg	U+5FA1	104	117
200003810_00024_2.jpg	U+3072	191	118
200003811_00024_2.jpg	U+540D	279	120
200003812_00024_2.jpg	U+3061	501	120

CODH カラーヌワット・タリンほか



くずし字認識サービス

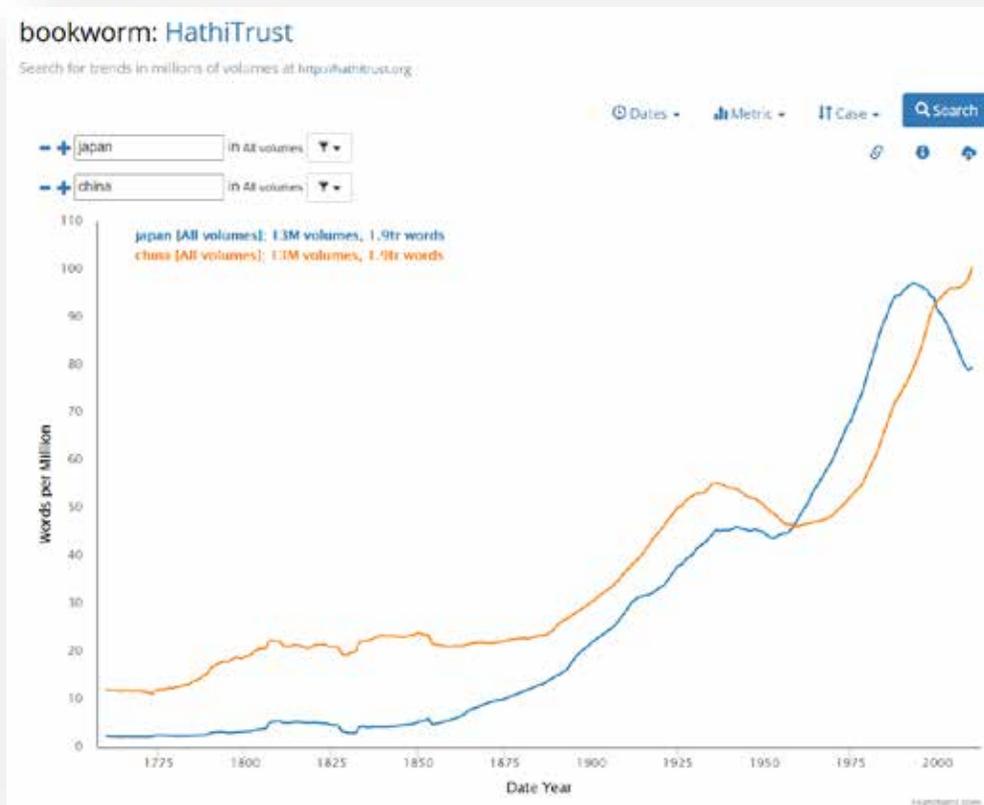
kaggle

くずし字認識コンペ



くずし字認識モバイル
アプリ

人文学ビッグデータと検索



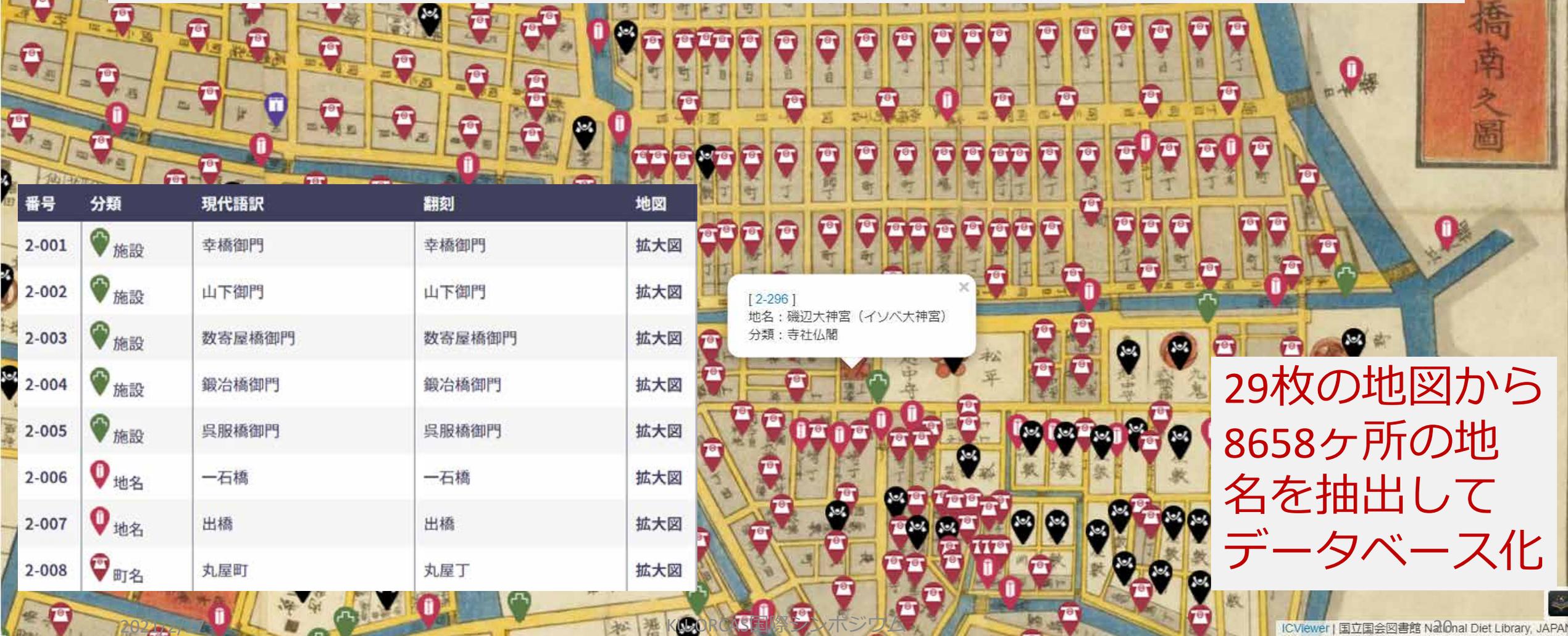
HathiTrust bookworm,
<https://bookworm.htrc.illinois.edu/develop/>

1. 書籍中の単語（N-gram）の出現頻度の推移に基づき、社会の変化を探る。
2. **Google Books Ngram Viewer** や **HathiTrust bookworm** など。
3. **検索機能が人文学研究のスピードとスタイルを変える。**
4. 日本語で同様の検索が困難なのは、OCRが一つの原因。

江戸ビッグデータ

江戸マップ - 江戸切絵図の地名リソース化

<http://codh.rois.ac.jp/edo-maps/>



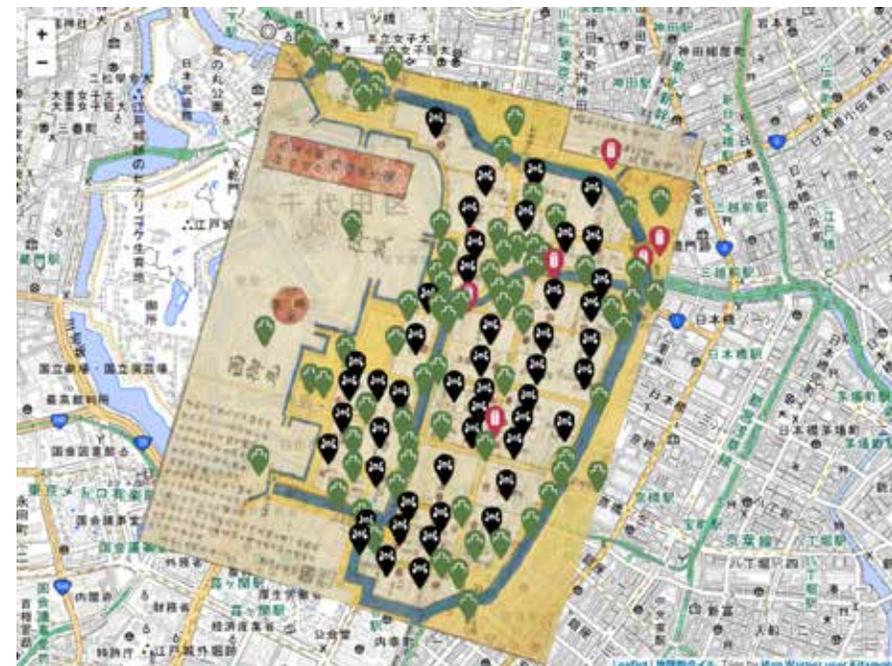
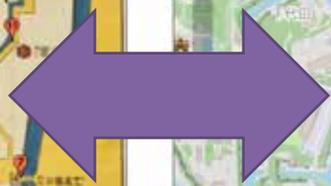
江戸切絵図の構造化

1. 尾張屋版（1849-1862） = 江戸切絵図の中で最も大量に売れて普及した地図。うち29枚を対象とする。
2. 国立国会図書館「デジタルコレクション」で公開されるIIF画像を、IIF Curation Viewerで読み込み。
3. 地名を矩形で囲んで画像座標を記録し、矩形中の文字を翻刻・現代語訳し、データベース化。
4. 「施設」「屋敷地」「寺社仏閣」「店名」「地名」「町名」「海川池」「観光地」「その他」の9分類を付与。ただし旗本屋敷などは数が多いため除く。

古地図のジオレファレンス



GCPの対応付け



国立国会図書館
『江戸切絵図』

立命館大学
日本版Map Warper

江戸マップβ版 + 日本版Map
Warperタイル配信サービス



千代田区

© 2020 ZENRIN

Google Earth

KU-ORCA 国際シンポジウム

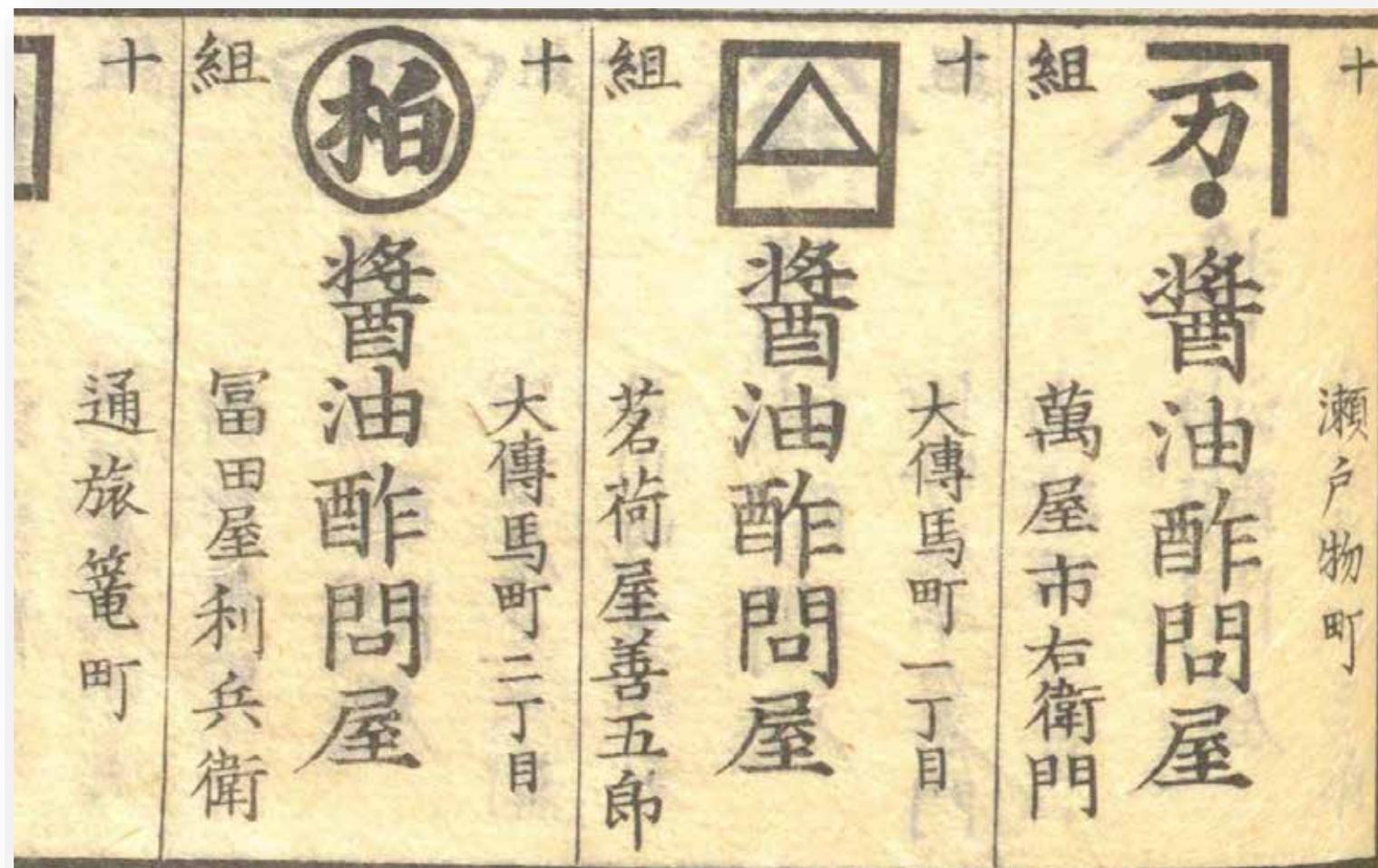
35° 40'52.26" N 139° 45'26.42" E 標高 2.00 km

2021/2/27

1997

江戸買物案内 – 商業ビッグデータ

<http://codh.rois.ac.jp/edo-shops/>



- 『江戸買物独案内』（1824）から、商人名や居所などを抽出しデータベース化。
- 歴博の「江戸商人・職人データベース」は現在の区レベル。
- 江戸買物案内は、江戸マップの町レベルでリンク。

江戸観光案内 – 観光ビッグデータ

<http://codh.rois.ac.jp/edo-spots/>



- 江戸時代の名所記や名所案内を各世紀2点ずつ選択。
- IIF Curation Viewerを用いて挿絵を切り抜き、地名を翻刻し、メタデータを付与。
- 江戸マップ、歴史地名データとリンク。

CODH 鈴木 親彦 ほか

史料とエンティティのリンク



宇田川丁、三島丁・
神明丁、此分潰家多、
土蔵残所なし

固有表現認識

宇田川丁、三島丁・
神明丁、此分潰家多、
土蔵残所なし

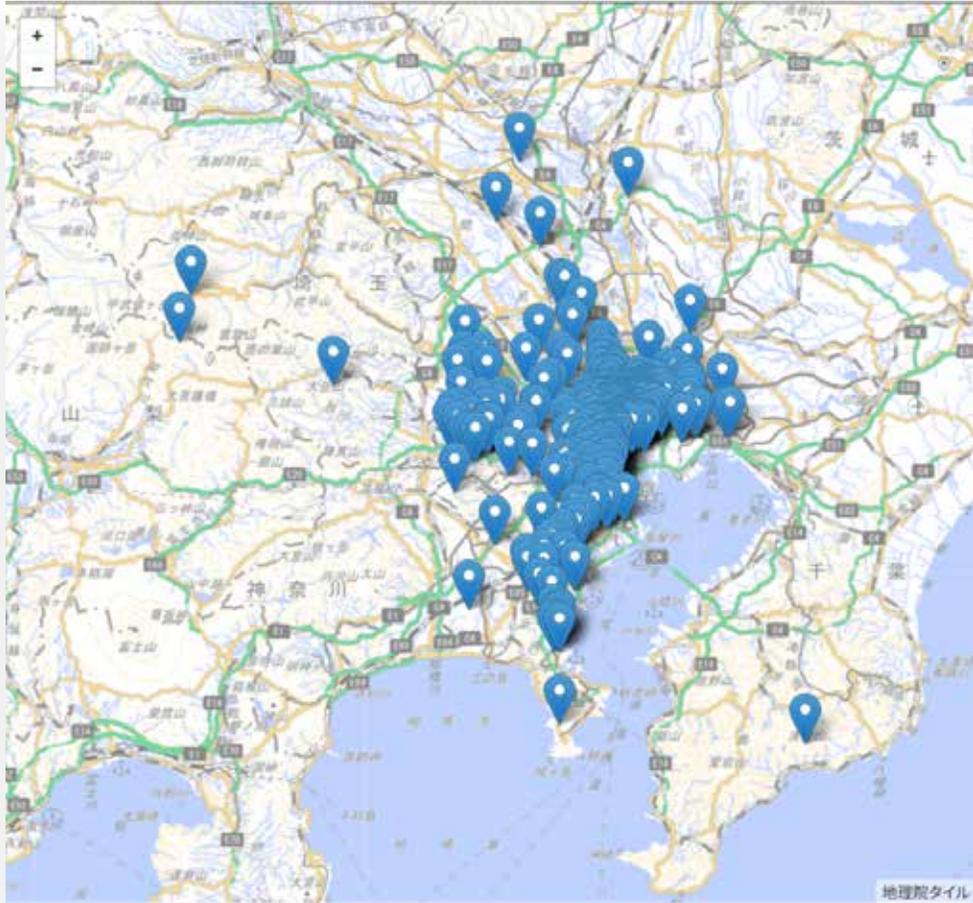
曖昧性解消

御江戸大地震大破并出火類焼場等書上之写（みんなで翻刻）

1. 江戸マップに出現する地名を、地名リソース（エンティティのデータベース）として整備する。
2. 史料の文字列から固有表現（地名）を抽出する。
3. 特定のエンティティとリンクすることで、実世界と紐づける。

原資料表記	江戸マップID	江戸マップ表記
宇田川丁	4-358	宇田川町
三島丁	4-290	三島丁
神明丁	4-294	神明町

江戸観光案内と地名リンク



1. **1255件**の挿絵を収集。その中で地名が明確な挿絵は**1219件**。
2. **歴史地名データ**とのリンクは**765件（全体の6割強）**で可能。
3. 歴史地名データは明治時代の地名が多いが、**観光名所の地名は江戸と明治で大きな変化はなかった可能性**がある。

歴史的行動記録と地名リンク



1. 清河八郎が安政2年（1855）に江戸を訪問した日記『西遊草』を分析。
2. 日記から地名を抽出し、江戸マップや江戸観光案内とリンク。
3. 164件の地名中、江戸マップと132件、江戸観光案内と119件がリンク可能。

歴史ビッグデータの拡大

歴史的記録

1. 歴史的状況記録

- 世界の状況を観察した記録（空の状況を視覚的に観察した記録、地震の揺れを聴覚・触覚で感じた記録、地震の被害状況を視覚で観察した記録など）

2. 歴史的行動記録

- 人間の行動に関する記録（人間がある場所から別の場所に移動した記録、店で何かを買ったり食べたりした記録、観光した記録など）

3. 歴史的状態記録

- 人間やモノの状態に関する記録（人間やモノの属性に関する記録、市場取引の価格に関する記録、モノの輸送に関する記録など）

れきすけ

<https://rksk.ex.nii.ac.jp/>



1. 史資料の所在情報を共有するサービス。
2. 書誌情報に加え、事項ごとの歴史的記録の所在に関する情報を共有できる。
3. 「地震のことが書いてある日記を検索したい」を実現する。

CODH 市野 美夏ほか

れきすけのカード

著作

山脇弁治日記_旧5 (5)

作成日:2020-12-04 19:14:19
更新日:2020-12-21 10:51:02

[編集](#) [関連リスト作成](#) [状況](#) [共有](#)

カード著者 : 平野淳平
作者 :

記録期間

入力 : 安政4年 ~ 明治6年
西暦 : 1857-01-26 ~ 1873-12-31
和暦 : 安政4年1月1日 ~ 明治6年12月31日
検索 : 1857-01-26 ~ 1873-12-31

記述

観測地の変遷:不明
記録期間は、
秋田県公文書館所蔵 郷土資料
https://www.pref.akita.lg.jp/uploads/public/archive_0000003962_00/mokuroku025.xls
に基づいた。

カード情報
カード所有者: Kitamoto

事項

山脇弁治日記_旧5 (388)

作成日:2020-12-04 19:14:19
更新日:2020-12-05 17:44:20

[編集](#) [関連リスト作成](#) [状況](#) [共有](#)

カード著者 : 平野淳平
事項タイプ : 地震,天気,水害
事項タグ : 記述の連続性,天気の詳細度

記録期間

入力 : 安政4年 ~ 明治6年
西暦 : 1857-01-26 ~ 1873-12-31
和暦 : 安政4年1月1日 ~ 明治6年12月31日
検索 : 1857-01-26 ~ 1873-12-31

記述

記述の連続性:9割程度
天気の詳細度:天候の時間変化、風の強さ、雨・雪の降り方などを詳細に記述
記録期間は、
秋田県公文書館所蔵 郷土資料
https://www.pref.akita.lg.jp/uploads/public/archive_0000003962_00/mokuroku025.xls
に基づいた。

カード情報
カード所有者: 橋本幸恵

資料

山脇弁治日記_みんなて割程度 (158)

作成日:2020-12-04 19:14:19
更新日:2020-12-15 18:29:39

[編集](#) [関連リスト作成](#) [状況](#) [共有](#)

カード著者 : 平野淳平, 市野美菜
別名 : 山脇弁治日記
資料タイプ : 割程度テキスト
資料タグ : みんなて割程度のフォーマット
資料作者 : みんなて割程度の発行者

記録期間

入力 : 安政4年 ~ 明治6年
西暦 : 1857-01-26 ~ 1873-12-31
和暦 : 安政4年1月1日 ~ 明治6年12月31日
検索 : 1857-01-26 ~ 1873-12-31

記述

その他:画像は個人撮影のもののみ
画像をくずし字割程度プロジェクト「みんなて割程度」の
制作作業を実施中
書影の記録期間:1863年~1868年
書影格納ソース:秋田県立公文書館
資料タイプ:画像(電子媒体)
記録期間は、
秋田県公文書館所蔵 郷土資料
https://www.pref.akita.lg.jp/uploads/public/3962_00/mokuroku025.xls
に基づいた。
みんなて割程度の山脇日記は田バージョンのみんな
ての下記URLで割程度完了済みの資料も表示する。を
示される
<https://vl.hankoku.org/app/#/idashboard/en>

カード情報
カード所有者: 市野美菜

資料

山脇弁治日記_旧5 (157)

作成日:2020-12-04 19:14:19
更新日:2020-12-05 11:43:17

[編集](#) [関連リスト作成](#) [状況](#) [共有](#)

カード著者 : 平野淳平
別名 : 山脇弁治日記
資料タイプ : 原資料(写本)
資料タグ :
資料作者 :

記録期間

入力 : 安政4年 ~ 明治6年
西暦 : 1857-01-26 ~ 1868年12月31日
和暦 : 安政4年1月1日 ~ 明治6年12月31日
検索 : 1857-01-26 ~ 1868年12月31日

記述

著作種?:
その他:原資料は秋田県
地震については、日本の
月7日・1857年4月1日
の書影格納ソース:秋田県
資料タイプ:原資料
記録期間は、
秋田県公文書館所蔵 郷
土資料
https://www.pref.akita.lg.jp/uploads/public/3962_00/mokuroku025.xls
に基づいた。

カード情報
カード所有者: 橋本幸恵

地名

秋田県秋田市 (733)

作成日:2020-12-04 19:14:19
更新日:

[編集](#) [関連リスト作成](#) [状況](#) [共有](#)

カード著者 : Yukie Hashimoto

地名情報

GeoLOD ID :
緯度・経度 :
都道府県 :

記述

秋田県秋田市

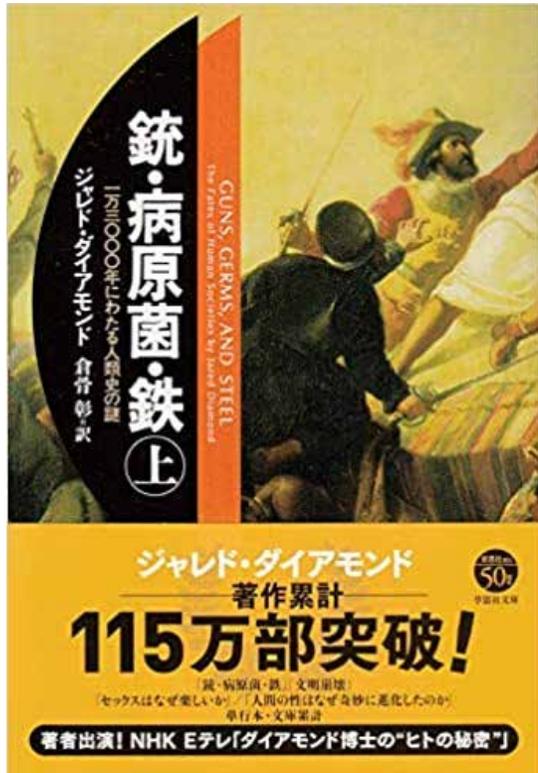
カード情報
カード所有者:

カードの種類一覧

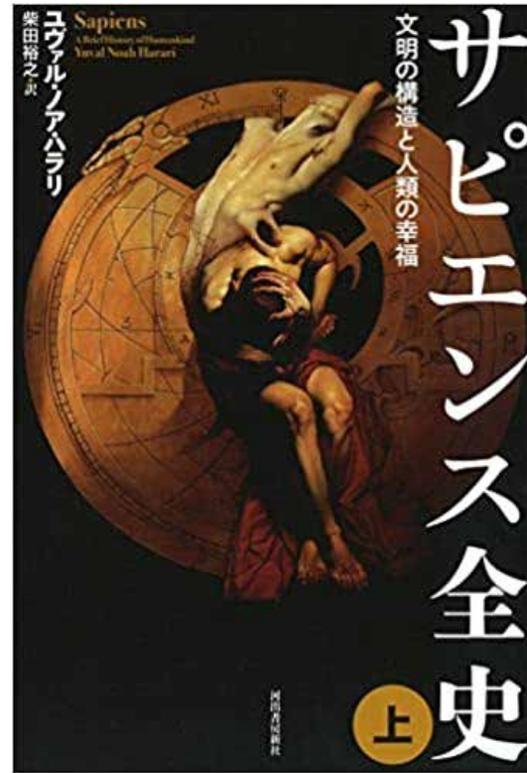
カードの種類	カードの情報
著作カード	作品そのものについての情報
資料カード	それぞれの実体についての情報
事項カード	歴史的状況記録に関する情報や、さまざまな知識や経験から得た情報
所蔵カード	史料の所蔵者である機関、団体、個人の情報
地名カード	空間情報
参考カード	史料を利用した研究成果、カードに登録した情報の典拠

グローバル／ビッグヒストリー

「歴史」の
対象が空間的・社会的
に拡大し、
より分野横
断的になっ
てきた



ジャレド・ダイヤモンド、銃・病原菌・鉄

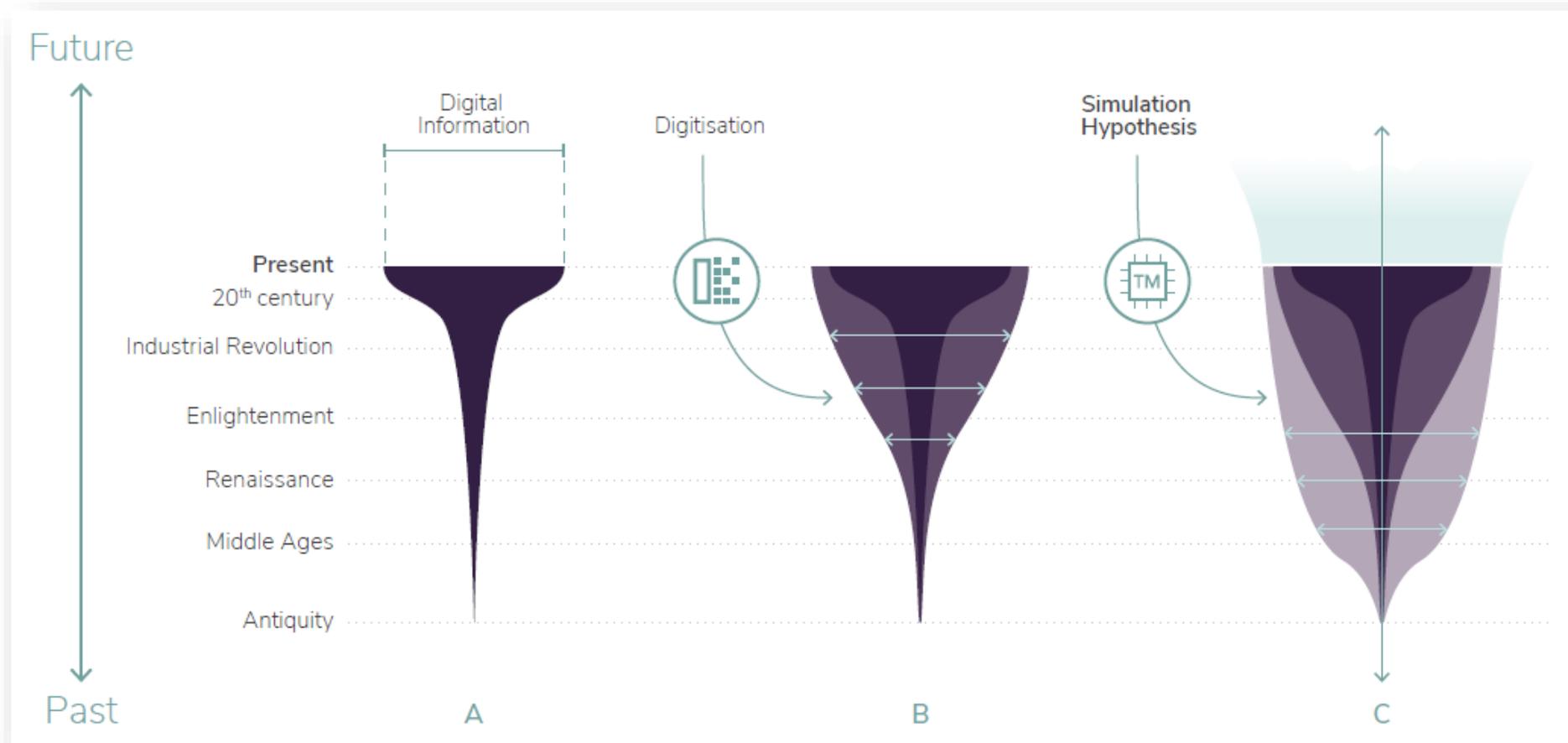


ユヴァル・ノア・ハラリ、サピエンス全史



デヴィッド・クリスチャン、ビッグヒストリー われわれはどこから来て、どこへ行くのか——宇宙開闢から138億年の「人間」史

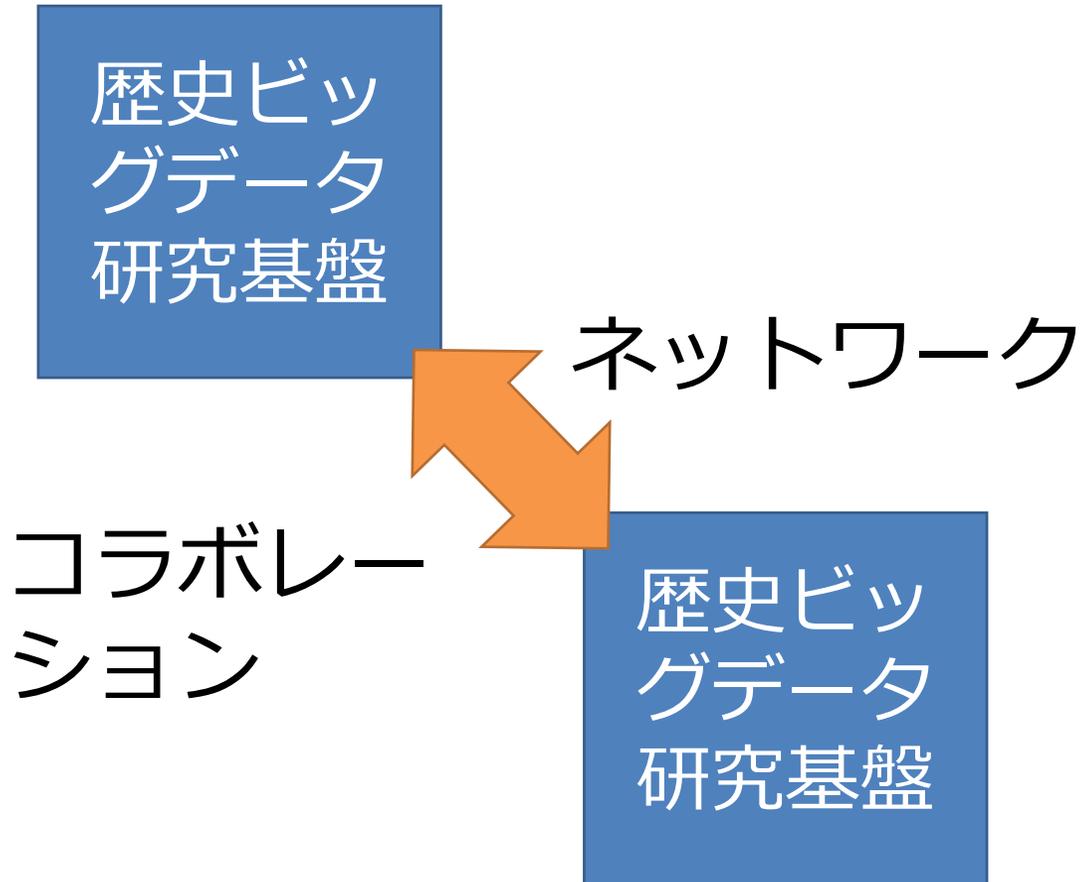
4D Mirror World



Time Machine Manifesto, <https://www.timemachine.eu/wp-content/uploads/2019/06/Time-Machine-Manifesto.pdf>

歴史ビッグデータ研究基盤

歴史ビッグデータ研究基盤



1. 様々な機関が「**自分の得意なこと**」に力を注ぐ。
2. 成果を相互利用する基盤となる「**オペレーティングシステム**」を設計する。
3. **ネットワークの形成とコラボレーションの活発化**により、長期的な視点で研究が高速化する。

UNIX哲学

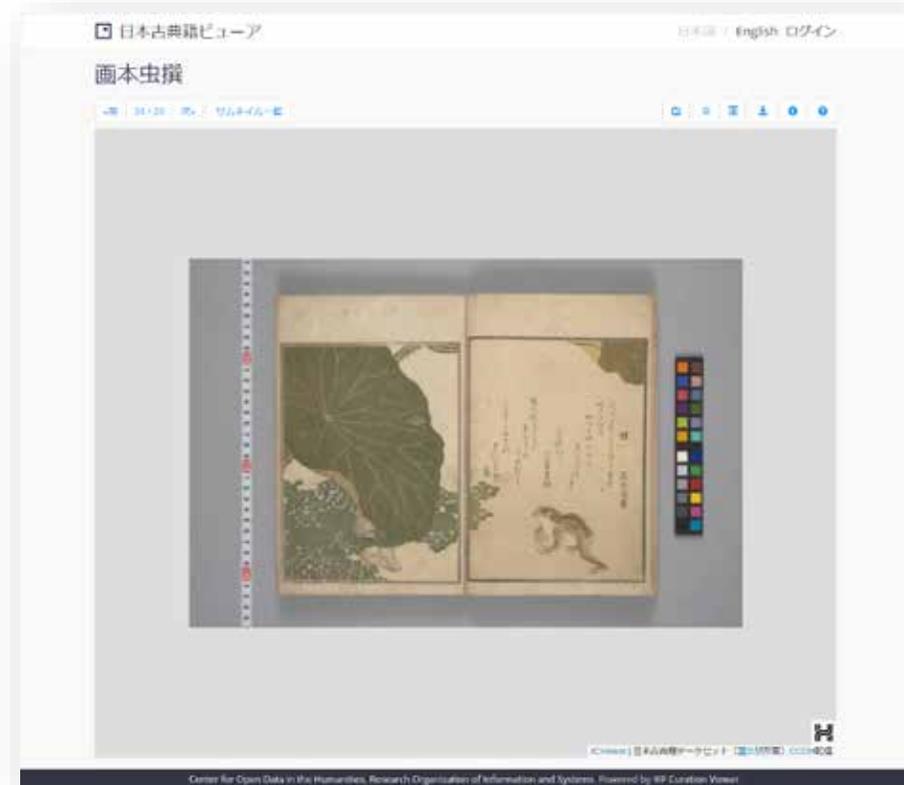
<https://ja.wikipedia.org/wiki/UNIX%E5%93%B2%E5%AD%A6>

1. **一つのこと**を行い、
またそれをうまくやるプログラムを書け。
2. **協調して動く**プログラムを書け。
3. **Keep it Simple, Stupid**
(KISS原則、「シンプルでつまらないものに保て」)



<https://twitter.com/cdixon/status/505118160811728896>

単一の機能を果たすツール



IIF Curation Viewer
画像の閲覧と切り取り



みんなて翻刻
画像のテキスト化（翻刻）

APIの相互運用性

時間 = **HuTime**
By 関野樹氏 (日文研)

API
呼び出し

地名 = **GeoLOD**
By CODH

記録期間
入力ボックスには、西暦(年、年-月、年-月-日)だけでなく、和暦(年号、年、年月、年月日)も入力できます。また和暦と西暦は自動変換します。

開始 ~ 終了

和暦: 西暦:

◆ 暦の変換にはHuTimeの「暦に関する WebAPI」を利用しています。◆

れきすけでの時空間情報入力

東京

検索 + 地名を新規入力する

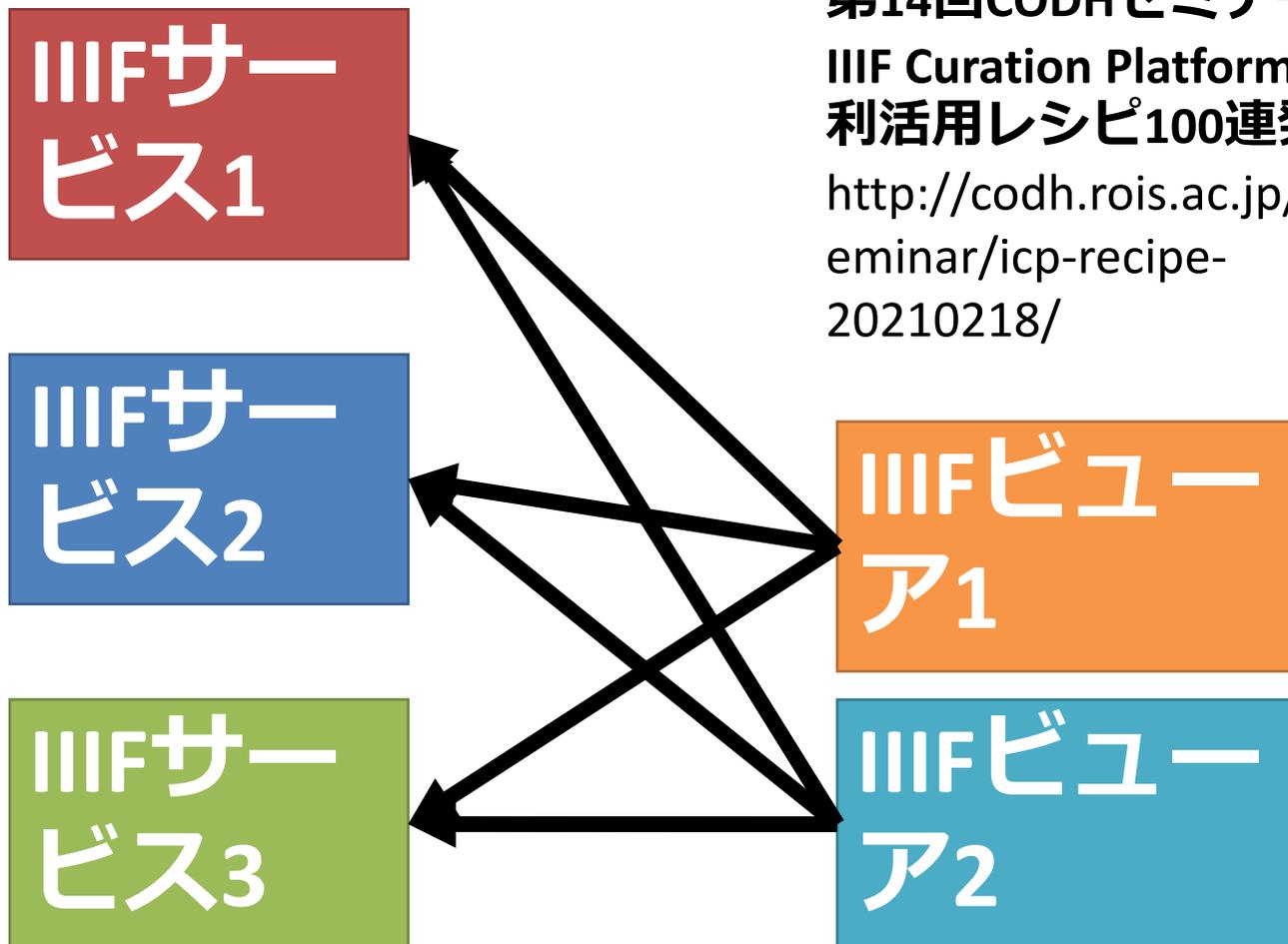
西東京
東京
東京海湾
東京博物館
東京
東京
東京上水
東京城
東京大学

IIIF (トリプルアイエフ) とは？

International Image
Interoperability
Framework = 国際的な
画像配信方式



Web : HTML
画像 : IIIF



第14回CODHセミナー
IIIF Curation Platform
利活用レシピ100連発
<http://codh.rois.ac.jp/seminar/icp-recipe-20210218/>

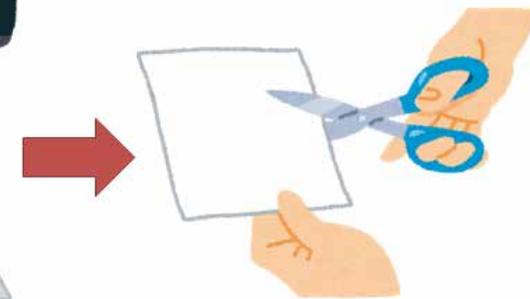
IIF Curation Platform

キュレーションとは、もともと、ミュージアムにおいて、資料の収集や作品の展示などの活動を指すことばである。

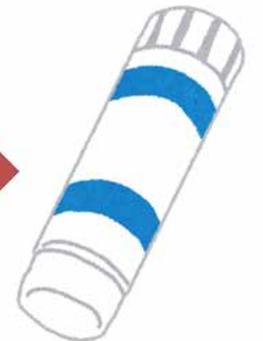
1. あるテーマに沿って**コンテンツ**を集める。
2. 適切な**順番（配置）**に並べる。
3. **新たなコンテンツ**として提示・共有する。



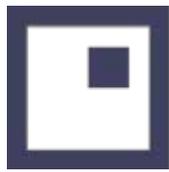
1 ハサミ



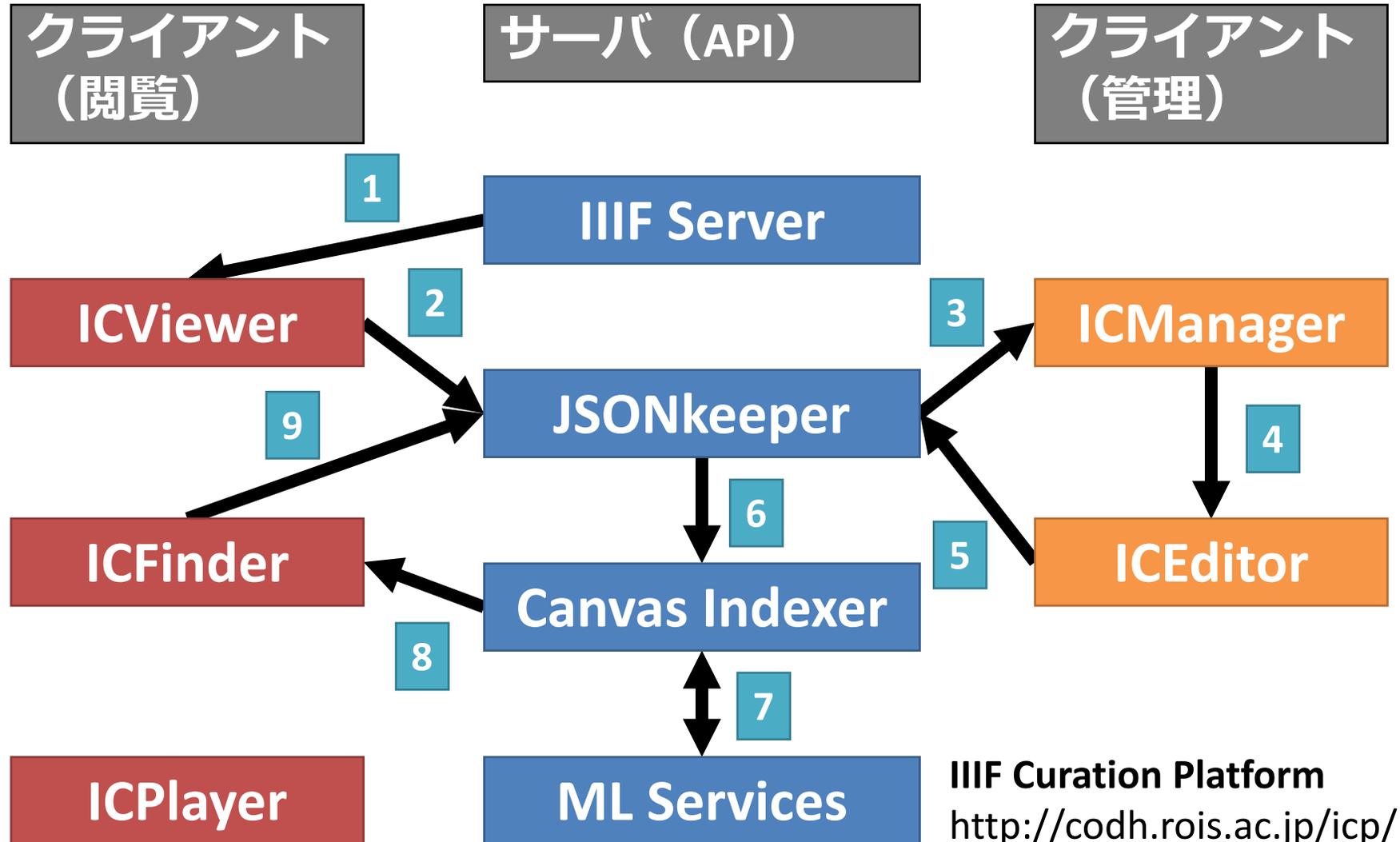
2 ノリ



出典：いらすとや、<http://www.irasutoya.com/>



IIIF Curation Platform

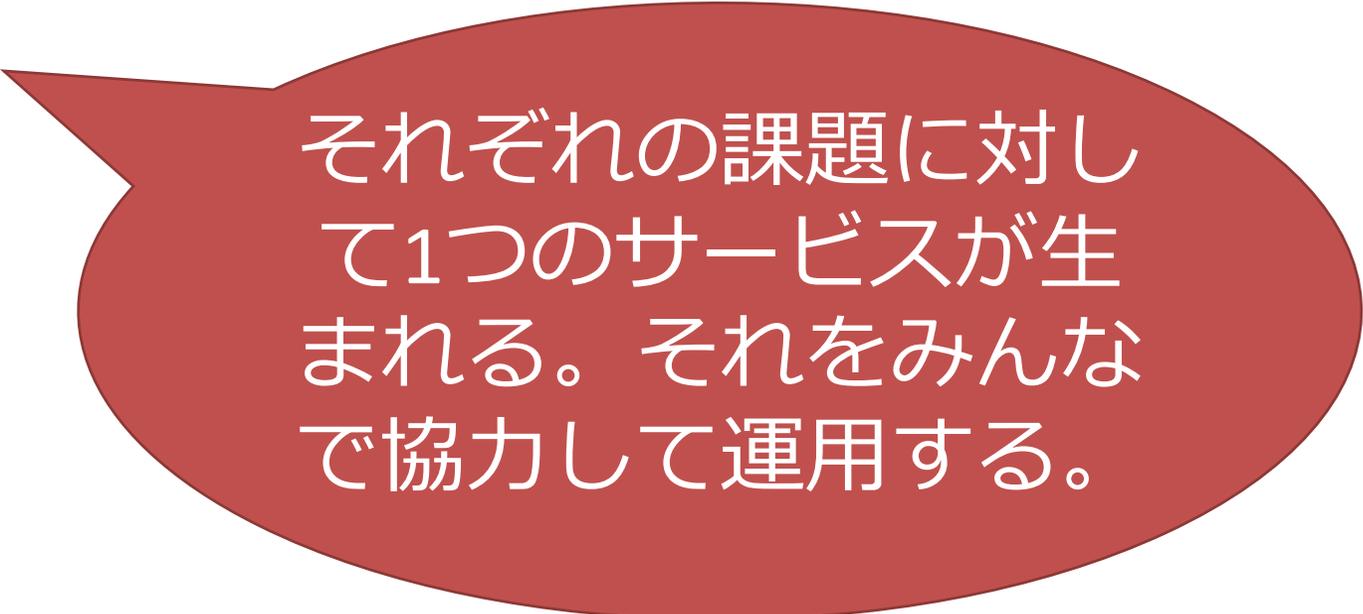


歴史的記録 = 5W1Hの構造化？

1. ニュースの基本情報は**5W1H (When、Where、Who、What、Why、How)** と言われる。
2. ニュースを過去に延長すると、**歴史的記録の基本情報も、5W1Hの形式でまとめられるのではないか？**
3. 5W1Hで表現された「事実」の集積が、「**歴史ビッグデータ**」の**基本データ**となるのではないか？
4. **基本データをどのように記述すべきか、また収集を支援するシステムをどう設計すべきか？**

1つのことをうまくやる

1. When → HuTime
2. Where → GeoLOD
3. Who → ?
4. What → ?
5. Why → ?
6. How → ?



それぞれの課題に対して1つのサービスが生まれる。それをみんなが協力して運用する。

歴史のビッグデータ化

1. 過去の歴史の「**記述の断片（例：5W1H） = 記録**」を**構造化**し、**機械可読化**し、**データベース化**する。
2. 歴史的記録を**数値化**し、様々な観点から**可視化**することで、データの傾向を大まかに把握する。
3. 歴史的記録を**ボトムアップに組織化**し、人間が理解しやすい**ストーリー（物語）化**する。
4. 歴史を**モデル化**し、**シミュレーション**や**シナリオ予測**することで、**現在や未来に資する知識**を生み出す。

参考文献

1. 北本 朝展, 鈴木 親彦, 寺尾 承子, 堀井 美里, 堀井 洋, "地理的史料を対象とした歴史地名の構造化と統合に基づく江戸ビッグデータの構築", 人文科学とコンピュータシンポジウム じんもんこん2020論文集, pp. 171-178, 2020年
2. 鈴木 親彦, 高岸 輝, 本間 淳, Alexis Mermet, 北本 朝展, "日本中世絵巻における性差の描き分け - IIF Curation Platform を活用した GM 法による『遊行上人縁起絵巻』の様式分析", 人文科学とコンピュータシンポジウム じんもんこん2020論文集, pp. 67-74, 2020年
3. 市野 美夏, 増田 耕一, 北本 朝展, "れきすけ : 歴史ビッグデータで知識と経験を共有する異分野間協働プラットフォーム", 人文科学とコンピュータシンポジウム じんもんこん2020論文集, pp. 31-38, 2020年
4. カラーヌワット タリン, 北本 朝展, "くずし字認識の進化とサービス化の展開", 人文科学とコンピュータシンポジウム じんもんこん2020論文集, pp. 3-10, 2020年

謝辞

- 本発表には、CODHのカラーヌワット タリン、鈴木 親彦、市野 美夏特任助教らによる成果を含みます。
- 本研究は以下の支援を受けました。
 - 歴史ビッグデータ研究基盤による過去世界のデータ駆動型復元と統合解析, 科学研究費補助金 基盤研究(A), 日本学術振興会 [No. 19H01141] / 研究代表者: 北本 朝展 / 期間: 2019-2021
 - 人文学ビッグデータにおける構造化ギャップの克服と分野横断的利用の検証, 機構間連携・文理融合プロジェクト, 情報・システム研究機構 / 研究代表者: 北本 朝展 / 期間: 2018-2020
 - 検索機能の高度化に係る総合的研究, 日本語の歴史的典籍の国際共同研究ネットワーク構築事業 研究開発系共同研究, 国文学研究資料館 / 研究代表者: 北本 朝展 / 期間: 2015-2020



<http://codh.rois.ac.jp/>