# 人文ＤＸを指向する情報基盤の構築

## Construction of an Information Infrastructure for the Humanities DX

原 正一郎*1、馬場 弘樹*2、森 信介*1、村上 勇介*1、関野 樹*3、
山田 太造*4、後藤 真*5

Shoichiro Hara*1, Hiroki Baba*2, Shinsuke Mori*1, Yusuke Murakami*1, Tatsuki Sekino*3, Taizo Yamada*4, Makoto Goto*5
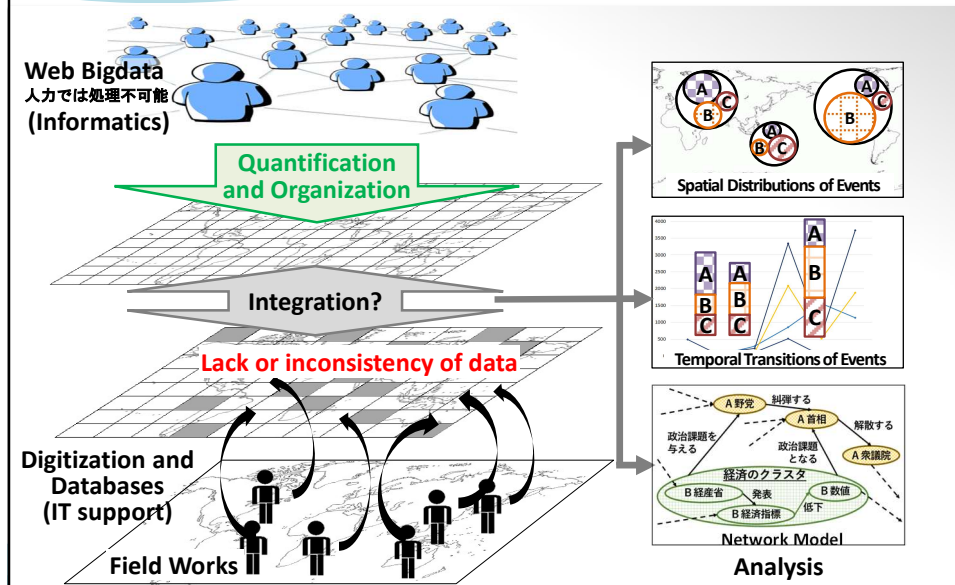
1.京都大学(Kyoto University) 2.一橋大学(Hitotsubashi University)
3.国際日本文化研究センター(International Research Center for Japanese Studies)
4.東京大学(University of Tokyo) 5.国立歴史民俗博物館(National Museum of Japanese History)

AI等の活用を推進する研究データエコシステム構築事業シンポジウム
2023/09/29@学術総合センター
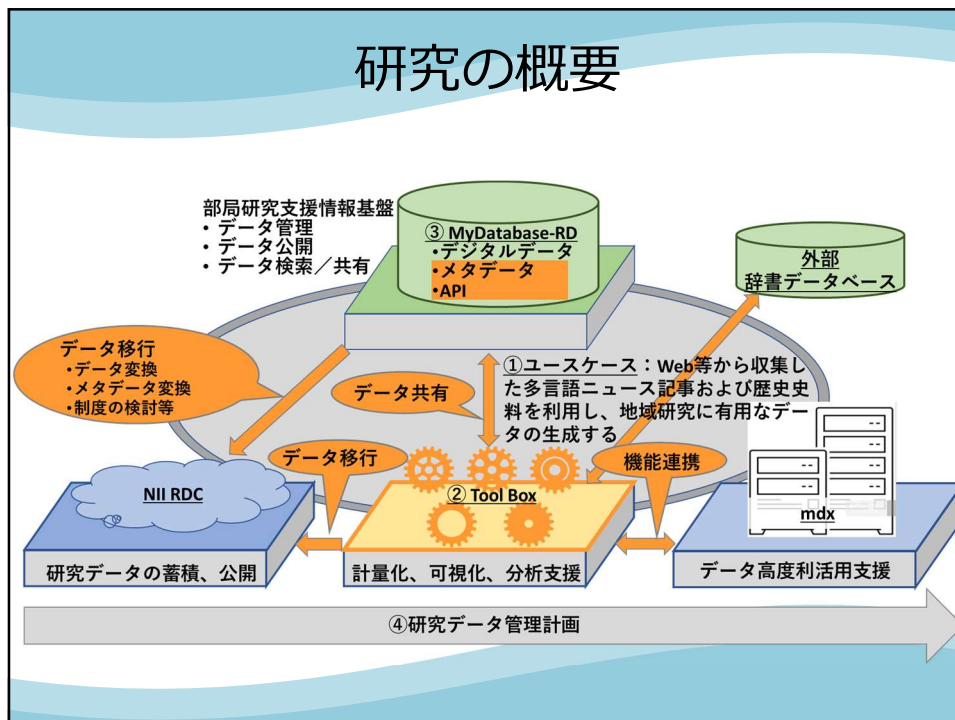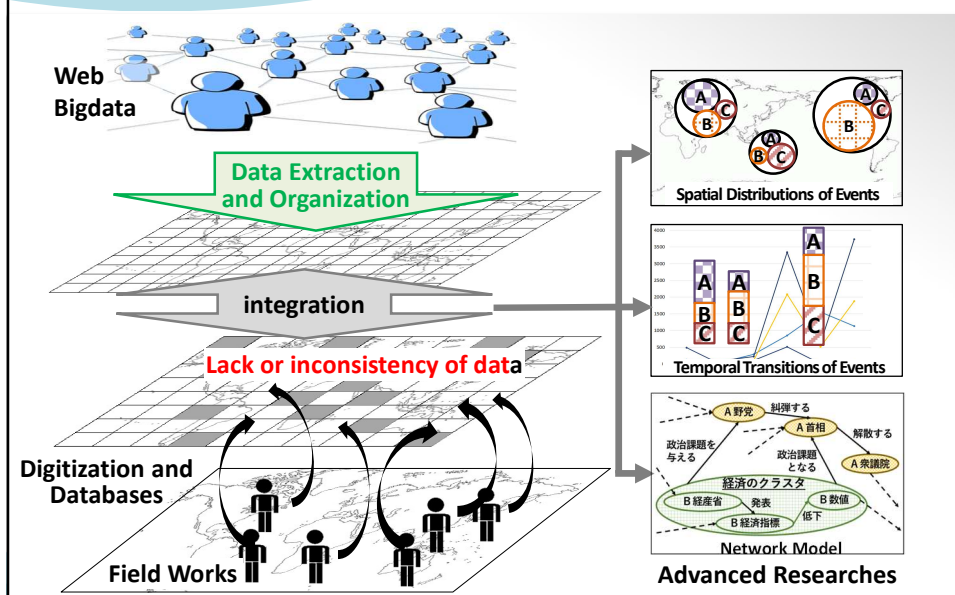
---

# モチベーション:情報学で地域研究を

# 目 的

DX：情報学をドメインとした計量的な地域研究の試みと、データ駆動型
地域研究の基盤整備。そこで本研究では、大量のテキストデータの計
量的処理をユースケースとして以下の研究を行う。

① ユースケース：Webニュースデータから地域変動を検出する計量的手法を
開発する（mdxの利用）

② ツール化と公開：データ分析を支援する使いやすい計算環境を、ツール
ボックスとして構築・公開する（mdxとGakuNin RDMの利用）

③ データ保存と共有：大量・多様な研究データの長期保存・共有・利活をど
のように実現するか？　まずはローカルなデータシステム構築し、次に公
的（？）な外部サーバとの連携を考える（GakuNin RDM の利用）

④ 研究データ管理手法：これまでのデータ管理は自己流。研究データを適切
に管理する知識・スキル・評価の標準的な手法の確立が必要（ Basic Rubric
とGakuNin RDM の利用）

# 研究の概要

ユースケース:Webニュースデータから地域変動を検出する計量的手法の開発

---

# ユースケース:データサンプル

- **収集しているデータ**
  - 毎日新聞: Web 収集および購入CD-ROM)
  - 朝日新聞: Web 収集
  - 読売新聞:Web 収集
  - AFP (English and Spanish news articles via AFP Forum)
- **この発表で利用したデータ**
  - 毎日新聞CD-ROM版
  - 全国版と46地方版
  - 2010-01-01 〜 2019-12-31

| | Articles | Sentences | Words | Characters |
|---|---|---|---|---|
| National | 607,671 | 7,791,338 | 219,677,468 | 321,075,103 |
| Local | 1,338,053 | 14,032,449 | 401,913,172 | 587,034,329 |
| Total | 1,945,724 | 21,823,787 | 621,590,640 | 908,109,432 |

# ユースケース: LDAトピックモデルの解析結果例 (2011/01/01 〜 2012/03/31 :日本)

- **トピックモデル: 教師なし学習手法**
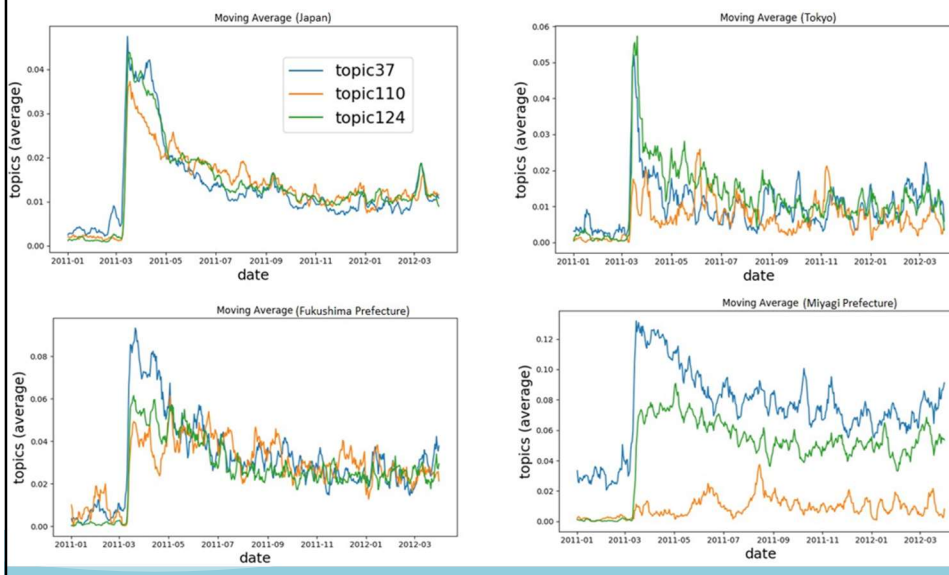- $p(\beta, \theta, z, w) = p(\beta)p(\theta)p(z|\theta)p(w|\beta, z)$ から

  $p(\beta, \theta, z|w) = \frac{p(\beta,\theta,z,w)}{p(w)}$ を求める問題

  計算コスト（時間）が高い(mdxに期待)

$\alpha \rightarrow \theta \rightarrow z \rightarrow w \quad \beta \leftarrow \eta$

N D K

Great East Japan Earthquake

← Date

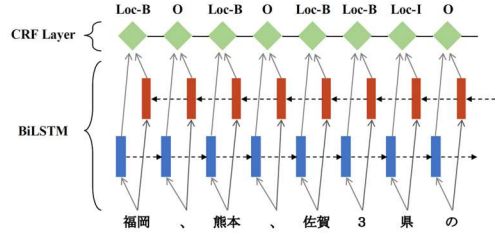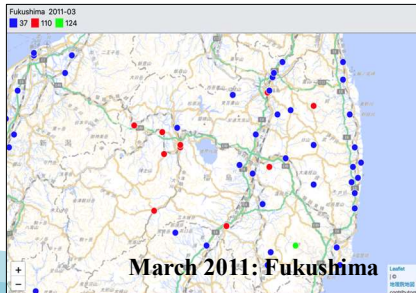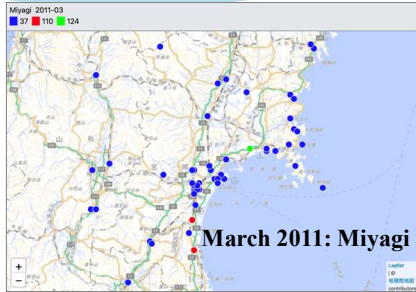Topics

Topic37

Topic110

Topic124

- Topic 37: 災害, 避難, 防災, 地震, 被害 (disaster, evacuation, disaster prevention, earthquake, damage)
- Topic 110: 原発, 電力, 発電, 原子, 放射 (nuclear power, electric power, power generation, atom, radiation)
- Topic 124: 震災, 被災, 地, 日本, 復興 (earthquake, damage, place, Japan, reconstruction)
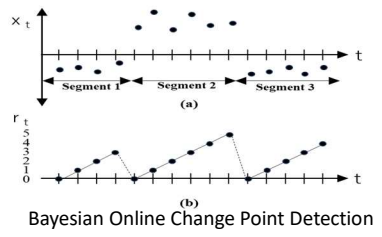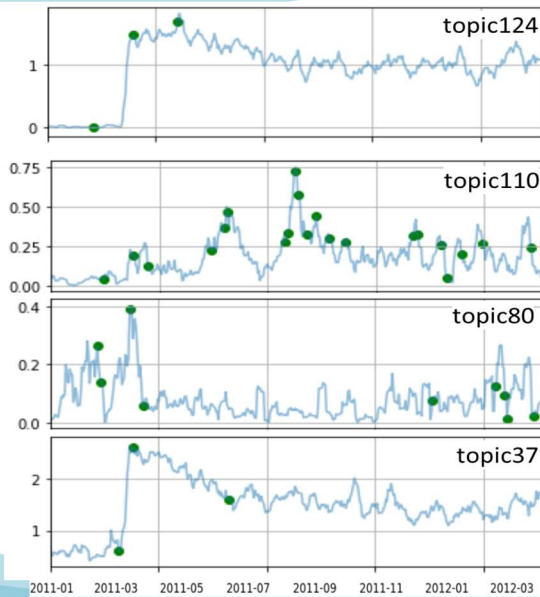
---

# ユースケース:トピックの時系列変化

Moving Average (Japan)

topic37
topic110
topic124

Moving Average (Tokyo)

Moving Average (Fukushima Prefecture)

Moving Average (Miyagi Prefecture)

# ユースケース:
# トピックの空間分布の変遷



March 2011: Miyagi

March 2011: Fukushima

CRF Layer

BiLSTM

福岡 、 熊本 、 佐賀 3 県 の

① BiLSTM-CRF:固有表現抽出。学習データから地名抽出法を学習する
② 地名辞書を検索し、地名に対応する緯度・経度の候補集合を取得する
③ 単純な戦略を用いて緯度・経度を一つに絞り込む
⇒ このアノテーションデータを共有する（何処にどのように蓄積する？）

■ Topic 37
■ Topic 110
■ Topic 124

---

# ユースケース:
# トピック変化点抽出の試み (宮城県)



topic124

topic110

topic80

topic37

Bayesian Online Change Point Detection

- Topic 124: 震災, 被災, 地, 日本, 復興 (earthquake, damage, place, Japan, reconstruction)
- Topic 110: 原発, 電力, 発電, 原子, 放射 (nuclear power, electric power, power generation, atom, radiation)
- Topic 80: 病院/医療/患者/看護/救急 (hospital, medical, patients, emergency )
- Topic 37: 災害, 避難, 防災, 地震, 被害 (disaster, evacuation, disaster prevention, earthquake, damage)

2011-01  2011-03  2011-05  2011-07  2011-09  2011-11  2012-01  2012-03

# ユースケース:トピック内容の変遷

Topic 80: 病院/医療/患者/看護/救急 (hospital, medical, patients, emergency )



# ツールボックスの開発（mdx）

# データ保存と共有:MyDatabase

- 多くは書誌的データベース（例：画像データと検索用メタデータ）だが、メタデータは研究者・資料毎に異なるので図書館システムには馴染まない
- 利用者のデータリテラシーはそれほど高くない
- 簡単に作れて変更も容易なデータベースシステムが必要



---

# データ保存と共有:MyDatabaseの例



- データベースの作成には時間と手間がかかるので数が増えない
- データ公開なら比較的容易であり、最近ではデータ公開が要求される
- 論文との紐付けが無いので図書館のデータレポジトリには馴染まない
- とりあえずMyDatabaseの機能拡張で対応する

## データ保存と共有:国際共同研究

**Research small Data Alliance in east and southeast asia (RsDA)**



**Member Institutes:** Academia Sinica (Taiwan), Chiang Mai University (Thailand), Chulalongkorn University (Thailand), International Research Center for Japanese Studies (Japa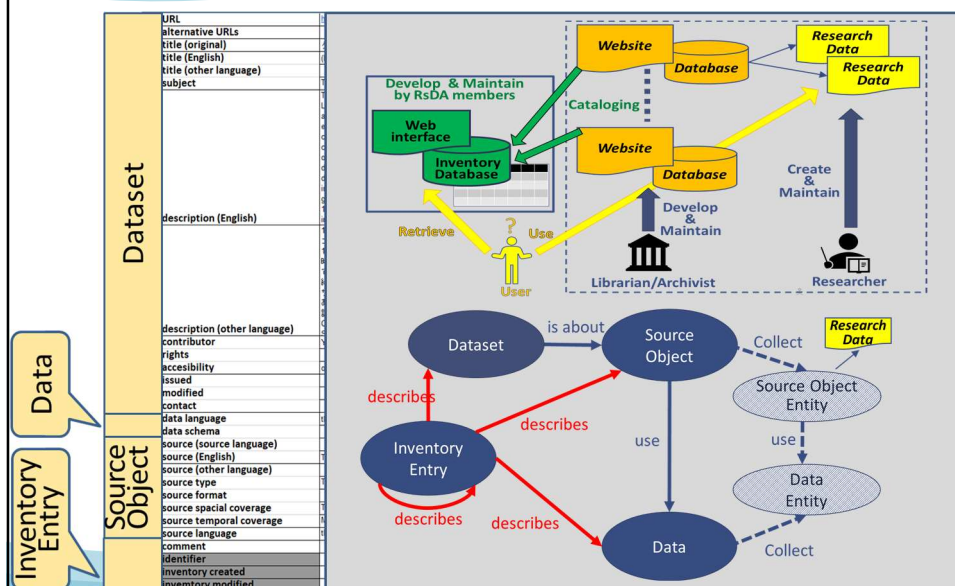n), Khon Kaen University (Thailand), National Museum of Japanese History (Japan), National Taiwan Normal University (Taiwan), National University of Singapore (Singapore), Royal University of Phnom Penh (Cambodia), Suranaree University of Technology (Thailand), Universitas Indonesia (Indonesia), Universitas Islam Riau (Indonesia), Universitas Lancang Kuning (Indonesia), Universiti Teknologi MARA (Malaysia), University of Colombo (Sri Lanka), University of the Philippines Diliman (Philippines), University of Tsukuba (Japan), Kyoto University (Japan)

1. ASEANの主要学術記憶機関における**学術デジタル資産**に関する現状の理解と共有
2. 研究データの共有と保存にどのように対処すべきか、利用可能な技術は何か、誰がこの問題に貢献すべきか等について、参加者が意見やアイデアを交換し共有する機会を提供。
3. コラボレーションのためのネットワークを構築するための実践的なアプローチについて議論

---

## データ保存と共有:メタデータ開発
### （MyDatabaseの機能拡張⇒ GakuNin RDMへ）

# 研究データ管理

1. どこにデータを蓄積したら良いのか（システムの問題）
   - 小規模の人文学系研究所（センター等）では恒久的な人員や機材の維持が困難
   - MyDatabaseは十数年に渡ってデータサービスを継続してきたが既に耐用年数を超えている
   - MyDatabaseでは収容できない大規模データの構築が進みつつある
   - 大学の研究データ管理システムあるいはGakuNin RDMへの期待
2. そもそもデータ管理はどのように行うのか（運用の問題）
   - 実はよく分かっていない（自己流：どのデータをどのように残す？）
   - GakuNin RDM利用するとしても、前提となる知識やスキルが必要ではないか（データとメタデータだけでデータの再利用は可能なのか？）
   - ⇒ 必要なスキルや評価の標準化が必要（例えば、京都大学学際融合教育研究推進センターのアカデミックデータ・イノベーションユニットが開発したBasic Rubric）
   - ユースケースを通じてRubricの利用と評価を試みる

---

# 研究データ管理:Rubricの役割と作成
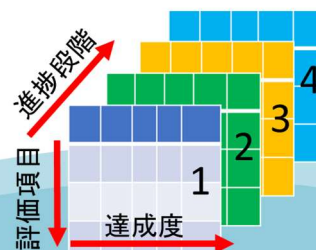
Rubricの役割
- 研究のための準備：Data Management Plans 提出等
- 研究進捗状況のチェック（データ取得・保管状況、処理・解析、文書化等）
- オープンサイエンス(データ)時代への対応

Rubricの役割
- 目標: 研究のライフサイクルの各段階でRDMに関する知識およびスキルの学習あるいは自己評価ができる仕組みの作成
- 分野: 様々な分野向けのルーブリックを作成。それらを基に、できる限り一般化する⇒Basic rubric
- 対象: 学生・若手研究者
- 特徴: 進捗段階に対応した4つのルーブリックに分割
    1. データに関する計画
    2. データの処理と整理
    3. データの解析
    4. データの共有または公開

ルーブリックの各項目毎に4段階の達成度を設定

## 研究データ管理:Basic Rubricの例
### (Planning for Data)

達成度

| PLANNING FOR DATA | Beginning | Developing | Enhancing | Completed |
|---|---|---|---|---|
| Data to Be Obtained | The data to be produced and the materials required to produce it have been determined. | The categories of data to be acquired and the relationship of the categories to specific research purposes have been identified. | The resolution, quality, and minimum amount of data, as well as necessary subjects for acquiring the data, have been identified. | The researcher is confident of the categories, resolution, quality, and minimum amount of data, and that needed subjects can be reliably acquired. |
| Methodology for Obtaining Data | The methodology required to obtain the data has been determined. | The instruments, data transfer plan, and storage device required to obtain and store the data have been identified. | Preparation or implementation of the instruments, data transfer arrangements, and storage device has begun. | Preparation or implementation of the instruments, data transfer arrangements, and storage device is complete. |
| Materials, Equipment, and Software for Data Acquisition | A plan for purchasing and/or creating the materials, equipment, and software has been determined. | The location(s) for collecting the data has been identified and arrangements for work in that location have been made. | Legal considerations associated with the materials, equipment, and/or software locally and internationally have been identified and observed. | Procedures for obtaining the necessary materials, equipment, and software for data acquisition have been documented and accomplished. |
| Team Members | Team members responsible for data acquisition and dissemination have been identified. | Appropriate roles have been assigned according to the skills and interests of individual team members. | A suitable research environment has been prepared to encourage information exchange, consultation, and cooperation among team members. | Team members fully understand the division of roles, as well as their own personal role, and are well prepared to conduct and disseminate the research. |
| Funding for Data Acquisition | Funding required for data acquisition has been estimated and available sources of funding identified. | Funding applications include team member expenses, purchase of equipment and materials, software for data acquisition, and legal expenses. | Requests for funding have been submitted to relevant funding sources. | Sufficient funding for all expenses related to data acquisition has been acquired. |
| Data Re-Use | Changes that may be required in ways that the research data is managed for re-use have been considered. | The availability of the data for re-use by other researchers has been determined. | Decisions have been made regarding the exclusive right of team members to have use of the data during an institutionally determined embargo period. | Plans for potential re-use of the data, including necessary changes to its management, are complete and accepted by all team members. |
| Open Versus Closed Data | Academic journals for potential publication have been identified and their conditions for data disclosure are understood. Sharing data in public repositories has been considered. | Decisions have been made as to which data (and materials) can be published and shared and which cannot. Relevant international and local legalities and proprietary issues have been considered. | The choice of licenses and repositories for publishing and sharing the data, as well as potential users of the data (and materials), have been considered. | The choices for publishing and sharing the data (and materials) have been documented and agreed to by all team members. |
| Planning for Publishing and Sharing Data | The requirements of funding agencies and research institutes being considered to manage the disclosure of the data (and materials) are understood. | Funding for disclosure of the data (and materials) has been estimated. It has been decided if the researcher's institution will be responsible for data disclosure or if the responsibility will be outsourced. | If needed, suitable repositories for outsourcing the publishing and sharing process have been identified and responsibilities for ongoing expenses have been determined. | A Research Data Management (RDM) plan for publishing and sharing the data (and materials) has been created, documented, and agreed upon by all team members and by the researcher's institution. |

評価項目

---

## 研究データ管理:ePortfolio

研究データ管理: ePortfolio



統合はこれから・・・