

# デジタル図書館研究室

研究室紹介  
2017/4/25 メディアプロジェクト演習2

前田 亮

1

## 「デジタル図書館」とは？

- デジタル情報を収集・整理・保管し、提供するサービス全般
  - 対象となる情報の種類は限定しない
    - ・ 電子書籍・新聞・雑誌、古典資料、芸術作品、Webページ、画像・映像・音源、etc.
- Web上の情報は、従来の図書館のように体系付けられていない
- しかし、人類の知識を蓄積・提供する一種の巨大な**デジタル図書館**と捉えることができる



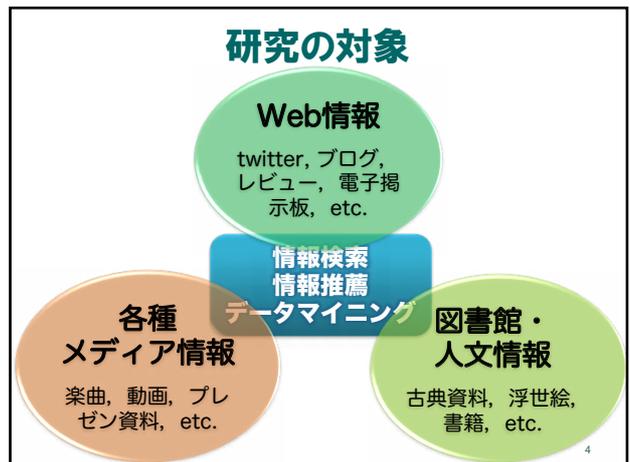
2

## デジタル情報の利点

- 「**検索**」ができる
- 情報の内容から自動的に「**分類**」できる
- 隠れた知識を「**発見**」(マイニング)できる
- 利用者に合った情報を「**推薦**」できる
- 様々な情報源を「**統合利用**」することができる
- これらの技術を「**情報アクセス技術**」と呼ぶ



3



## 主な研究テーマ

- **Web情報**
  - 文章中の**ネットスラング**の正しい表現への変換
  - スポーツ中継中の**Tweet**解析によるファン視点でのイベント検出
  - **多言語Wikipedia**記事への言語横断エンティティリンキング
- **各種メディア**
  - 図形群の意味や階層構造を用いた**プレゼンスライド**検索システム
- **古典史料**
  - **古典史料テキスト**からの情報抽出
  - 異言語の**浮世絵**データベースからの同一作品の同定<sup>5</sup>



5

## Web情報

6

## 文章中のネットスラングの正しい表現への変換

- 電子掲示板ではネットスラングが用いられる
  - 「草」「乙」「厨房」「密林」「リア充」「希ガス」...
- 形態素解析辞書と単語共起の統計情報を用いて、ネットスラングを正しい表現に変換

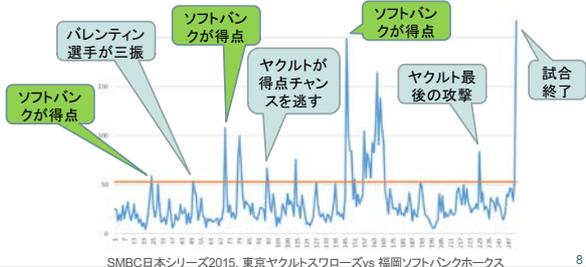
| 実際の書き込み                              | 変換結果  |
|--------------------------------------|---|
| ありがとうございますと赤を密林に頼んできましたwtkkしながら待ちますw | ありがとうございますと赤をAmazonに頼んでましたワクワクテカテカしながら待ちます(笑) |
| 本人と決めつける「自演乙！」                       | 本人と決めつける「自演お疲れ様！」                             |

赤字: ネットスラング      青字: 正しい日本語表現

7

## スポーツ中継中のTweet解析によるファン視点でのイベント検出

- Tweetの盛り上がり具合 (Tweet数, RT数, 選手名の出現数) を基に、ファン視点でのハイライトを検出



8

## 多言語Wikipedia記事への言語横断エンティティリンク

- 文書中のキーワードから、該当する知識ベース記事 (Wikipediaなど) に自動的にリンクを付与する手法

日本語文書

加齢黄斑変性?      日本語Wikipedia記事

- 用語によっては、母国語の記事がない、あるいは他言語の記事の方が充実している場合もある
  - 「Christendom (キリスト教界)」の日本語記事は存在しない

9

## 言語横断エンティティリンクの例

文章中の用語と記事名が異なっていてもリンク可能      英語版Wikipediaの該当記事

中国語版Wikipediaの該当記事      リーマン・ショック      日本語版Wikipediaの該当記事

言語間リンクがなくても該当記事にリンク

## 各種メディア

11

## 図形群の意味や階層構造を用いたプレゼンスライド検索システム

- プレゼンテーションスライド
  - 視覚的にわかりやすいように図が多用
- 図の作成は手間がかかる
  - 図の再利用のニーズ
- 従来のスライド検索手法では、図自体の検索が困難



図の作成が面倒だなあ...

以前作った図を再利用したい!

12

### 検索システムのユーザインタフェース

フォルダ選択  
C:\Users\seta\Documents\検索用フォルダ

クエリ入力

プレビュー

フォルダ選択

索引作成

図形群の意味  
流れ

図形群の要素の数  
3 ~

要素の形状  
四角形

さらに含まれる意味

データベース作成  
1. 文書群・アクションの抽出  
2. テスト項目の前処理  
3. メタデータの作成

特徴マッピング

インポートとターゲットマッピング

検索

検索結果

arrow\_group1流れ group1総リスト

preview

似ている図形群を探す

### 検索の例

入カクエリ  
図形群の意味: **重なり**  
図形群の要素の数: **4**  
図形群の要素の形状: **楕円**

トリプルのほうほうマーク

Outline of the system

第2章「キャラ」は何か

1位 score: 100

2位 score: 100

3位 score: 95

4位 score: 95

5位 score: 85

集合関係

集合関係

background

Algorithm of the inclusive judging

1位 score: 80

2位 score: 80

3位 score: 80

4位 score: 80

5位 score: 80

前の10件

次の10件

### 応用例:類似図形群の検索

入力図形群

検索された類似図形群

1位 score: 0

2位 score: 25

3位 score: 10

4位 score: 16.66666

5位 score: 17.91666

6位 score: 18.16666

7位 score: 18.16666

8位 score: 18.16666

9位 score: 18.16666

10位 score: 21.86666

前の10件

次の10件

### 古典史料

16

### 古典史料からの人物表現の自動抽出

平清盛

実名

別名

役職

鶏鳴清盛朝臣、義朝、義康等、軍兵都六百騎発向白河、六波羅亭、入道御座宰相中将被装束、次供奉行列被渡南庭、上皇女御殿御同車、自御前渡御馬場棧敷、前太政大臣以下歩行扈從、

人間関係の可視化などに役立つ

北条政子

後白河天皇

源頼朝

平清盛

源義経

武藏坊弁慶

那須与一

17

### 古典資料テキストからの人物情報の抽出と可視化①

古学

役者評判記には、歌舞伎役者の演技や容姿に対する評価(レビュー)が書かれている

役者評判記でレビューされている歌舞伎役者

位付(役者の評価)

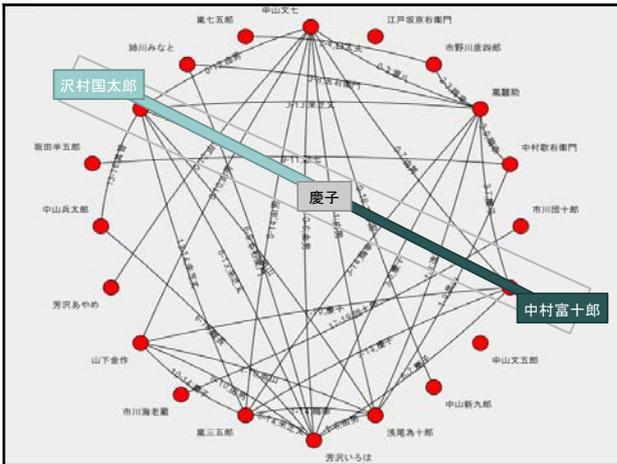
役者名

評価文

評価文中の人物呼称を抽出し、人物関係の可視化を行う

役者評判記『役者大極舞』  
立命館大学アート・リサーチセンター( <http://www.arc.nitsumei.ac.jp/index.html> ) 所蔵

18



### 古典資料テキストからの人物情報の抽出と可視化②

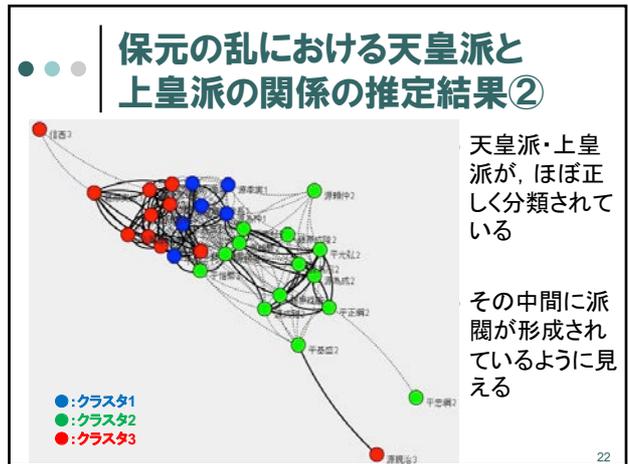
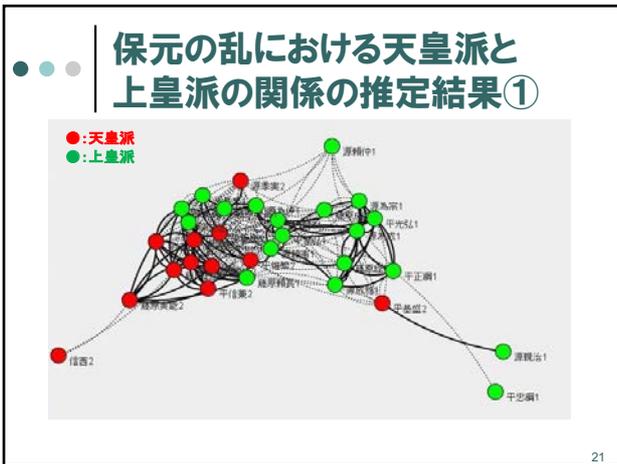
地名

人名

人名と地名が同段落に出現

段落

○ ある人名と地名の組が、同じ段落の中に現れた回数を集計することで、人名と場所の関連を推定



### 『東大寺要録』への注釈作業支援システム

- 12世紀に成立したとされる歴史資料
- 東大寺の歴史と当時の状況などが記されている
- 全十巻のうち一巻(約2万字)が電子テキスト化
- 東大寺要録研究会が書き下しと注釈作業中

電子テキスト化

注釈付き書き下し文

入力データ

注釈付けられていない文書

注釈の出現規則

注釈候補を抽出した文書

出力データ

● 注釈作業には、複数の分野の専門知識が必要

● 専門分野によって知識に差がある

### 機械学習による注釈候補の自動抽出

学習データ

既存の注釈

学習処理

注釈の出現規則

抽出処理

出力データ

## 異言語の浮世絵データベースからの同一作品の同定

- 浮世絵は、国内外の美術館・博物館で個々にデジタル化され公開
  - 同じ浮世絵作品が複数のデータベースに散在
  - データベースによってメタデータの言語・表記が異なる



25

## 浮世絵研究者からの要望

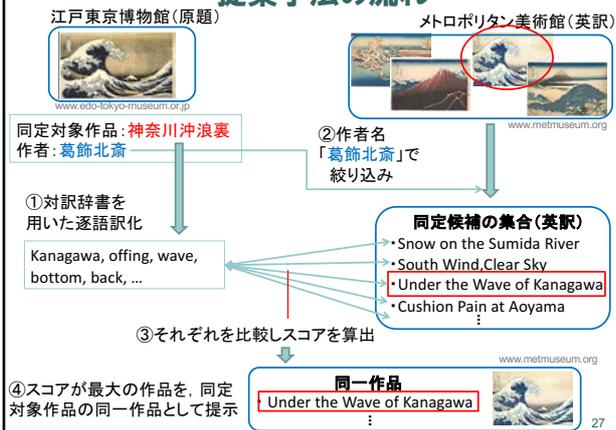
- データベース間で散在している**同一作品を見つけたい**
- 同一作品のメタデータを比較することにより、誤りの発見・修正、より詳細な情報の取得などが可能となる



- 本研究では**原題と英訳**を用いた同一作品の同定手法を提案

26

## 提案手法の流れ



27

## メディアプロジェクト演習2 作品制作のヒント (1)

- 「調査研究」
  - 新しい**Web技術**にはどのようなものがあり、それによって何が実現できるか
  - **デジタル図書館システム**としてどのようなものが開発され、技術的課題は何か
- 「未来創造」
  - 従来の**図書館とデジタル図書館**は、将来どのように共存していけばよいか
  - デジタル情報を**長期にわたって保存**していくには、どのようにすればよいか

28

## メディアプロジェクト演習2 作品制作のヒント (2)

- 「ソフトウェア」
  - 簡単な**情報検索システム**の作成
    - メディア情報学実験2「情報検索」の発展
    - 検索結果のランキング・分類・可視化など
  - 簡単な**デジタル図書館システム**の作成
    - 各種WebサービスAPIの利用
      - 国立国会図書館サーチAPI
      - 図書館API (カーリル)
      - 楽天ブックス書籍検索API
      - 版元ドットコム・書誌情報API

29

## デジタル図書館研究室について

- **現在のメンバー**
  - 前田亮 教授
  - ビルゲサイハン・バトジャルガル 専門研究員
  - 大学院 博士課程1名、修士課程7名
  - 学部生11名
- **共同研究**
  - 本学文学部、MOT研究科など
- **研究成果 (2016年度)**
  - ジャーナル1件、国際会議6件、国内会議5件、受賞1件

30



## おわりに



- 当研究室では、**デジタル図書館・情報アクセス技術**に関わる研究を幅広く行っている
- 研究室ホームページ
  - <http://www.dl.is.ritsumeai.ac.jp/>
- メールアドレス
  - [amaeda@is.ritsumeai.ac.jp](mailto:amaeda@is.ritsumeai.ac.jp)

31