



学生・若手研究者のためのBERTワークショップ2 (固有表現タスク 1回目)

理化学研究所
革新知能統合研究センター
言語情報アクセス技術チーム

オープニング





本日の概要



- BERTおさらい（前回のふりかえり）（10分）
 - 前回のワークショップの資料やビデオはこちら
http://shinra-project.info/shinra2022/bert_workshop_shinra2022/
- BERT 体験(60分)
 - Google Colabを使って自分で動かしてみる
 - **Googleアカウントが必要**
 - コードの説明
- 森羅プロジェクトの説明、今後の説明（15分）

森羅Slack：shinra2022.slack.com

チャンネル：[#bert_workshop2022](https://shinra2022.slack.com)

招待リンク：https://join.slack.com/t/shinra2022/shared_invite/zt-14qkpf21i-lQNKIToalOU5We7xIZBqfQ



学生アルバイト募集

このWSのその先に



- 森羅プロジェクトでアルバイトを募集しています
 - 対象：B4~D1を対象
 - 形態：アルバイト（1年程度、週に15時間前後、リモート可）
インターン（2~3ヶ月程度集中的に。リモート相談）
 - 業務内容：森羅プロジェクトのシステム開発、データ開発
指導教官との相談の上、卒論、修論とすることも可能
理研のRIDEN（超巨大GPUマシン）を利用可能
- 社会人の兼業希望者も応相談

相談、希望は下記までメールを
satoshi.sekine@riken.jp



本日の講師：中山功太



学歴

2014年4月 - 2018年3月(学部), 2018年4月 - 2020年3月(修士課程)
豊橋技術科学大学 システムAIラボ(石田研)

2020年4月 - 現在(博士課程)
筑波大学 ヒューマンコンピューテーション研究室(馬場研)

職歴

2018年8月 - 2020年10月
理研 言語情報アクセス技術チーム(関根チーム) 研究パートタイマー

2020年11月 - 現在
同上 リサーチアソシエイト

論文(査読付き)

- [Kouta Nakayama](#), Shuhei Kurita, Akio Kobayashi, Yukino Baba, and Satoshi Sekine. 2021. Co-Teaching Student-Model through Submission Results of Shared Task. In Findings of the Association for Computational Linguistics: EMNLP 2021, pages 4525–4535, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Satoshi Sekine, [Kouta Nakayama](#), Maya Ando, Yu Usami, Masako Nomoto and Koji Matsuda, SHINRA2020-ML: Categorizing 30-language Wikipedia into fine-grained NE based on “Resource by Collaborative Contribution” scheme, In Proceedings of the 3rd conference on the Automated Knowledge Base Construction (AKBC 2021), 2021.
- Satoshi Sekine, Akio Kobayashi, [Kouta Nakayama](#), SHINRA: Structuring Wikipedia by Collaborative Contribution, In Proceedings of the 1st conference on the Automated Knowledge Base Construction (AKBC 2019), 2019.



クロージング





本日の概要



- BERTおさらい（前回のふりかえり）（10分）
 - 前回のワークショップの資料やビデオはこちら
http://shinra-project.info/shinra2022/bert_workshop_shinra2022/
- BERT 体験(60分)
 - Google Colabを使って自分で動かしてみる
 - **Googleアカウントが必要**
 - コードの説明
- 森羅プロジェクトの説明、今後の説明（15分）

森羅Slack : shinra2022.slack.com

チャンネル : [#bert_workshop2022](https://shinra2022.slack.com/channels/bert_workshop2022)

招待リンク : https://join.slack.com/t/shinra2022/shared_invite/zt-14qkpf21i-lQNKIToaIOU5We7xIZBqfQ



森羅プロジェクト

今回のWSのサンプルデータ



- Wikipediaの構造化プロジェクト

- 「説明できる人工知能」を作るために知識グラフを構築
- 協働による知識構築
 - リーダーボードを作り、参加者とともに実施

今回のタスク

- 3つのタスク

- Wikipediaページを約200種類のカテゴリーに分類
- ページの属性情報を抽出
- 抽出した情報をリンキング

- 進捗

- 2018年からタスクを実施。今年度、来年度は3つのタスク



3つのステップ

ステップ1 (分類)

各Wikipediaページを約220種類の拡張固有表現に分類
（「島崎藤村」は人名！）

ステップ2 (属性値抽出)

固有表現定義にある属性値をページから抽出
（「島崎藤村」の「作品」には「嵐」がある！）

ステップ3 (リンクの紐付け)

抽出した属性値を該当するWikipediaページに紐付け
（「嵐」はWikipediaページの「嵐（作品）」のこと！）



Wikipediaページ

(日本語は100万、英語は600万)



A collage of various Wikipedia article thumbnails in Japanese, including:

- アメリカ合衆国** (United States of America)
- 第二次世界大戦** (World War II)
- イリノイ州** (Illinois)
- シカゴ** (Chicago)
- シカゴ・オヘア国際空港** (Chicago O'Hare International Airport)
- ミッドウェー海戦** (Battle of Midway)
- 嵐 (グループ)** (Arashi)
- NHK紅白歌合戦** (NHK Red and White Song Festival)
- 大坂市** (Osaka City)
- 雲のじゅうたん** (Clouds of the Sky)
- 田中絹代** (Tanaka Kinuyo)
- 筑前県** (Chikugo Prefecture)
- 飯島県** (Iwajima Prefecture)
- 島崎藤村** (Shimazaki Tōson)
- 嵐 (歌)** (Arashi (song))
- Edward O'Hare**
- NHK**
- 日本放送協会** (NHK)

ステップ1 (分類)



A collage of various Wikipedia articles in Japanese, each with a colored circular label indicating its classification. The labels include:

- 国名** (Country Name): 日本 (Japan), アメリカ (USA)
- 戦争** (War): 第二次世界大戦 (World War II)
- 州** (State): イリノイ州 (Illinois)
- 都市** (City): シカゴ (Chicago)
- 空港** (Airport): シカゴ・ミッドウェー国際空港 (Chicago Midway International Airport)
- 人名** (Person Name): 坂井三郎 (Sakai Saburo), 田中絹代 (Tanaka Kinuko), 雲のじゅ (Kumogajyu)
- 組織** (Organization): NHK (NHK)
- 番組** (Program): NHK紅白歌合戦 (NHK Red and White Song Festival)
- 企業** (Company): 日本放送協会 (NHK)
- 船** (Ship): 大和 (Yamato)
- 名詞** (Noun): 艦 (Ship)
- 県** (Prefecture): 大阪府 (Osaka Prefecture)
- 都** (Metropolis): 東京都 (Tokyo Metropolis)
- 島** (Island): 島崎村 (Shimazaki Village)
- 教員** (Teacher): 坂東三津男 (Sakaito Mitsuotoko)
- 不** (Not): A yellow circle with the character '不' (Fu), likely representing a negation or a specific category.

The background features various images related to the articles, such as the American flag, a battleship, a city skyline, and a group of people.

ステップ2 (属性値抽出)

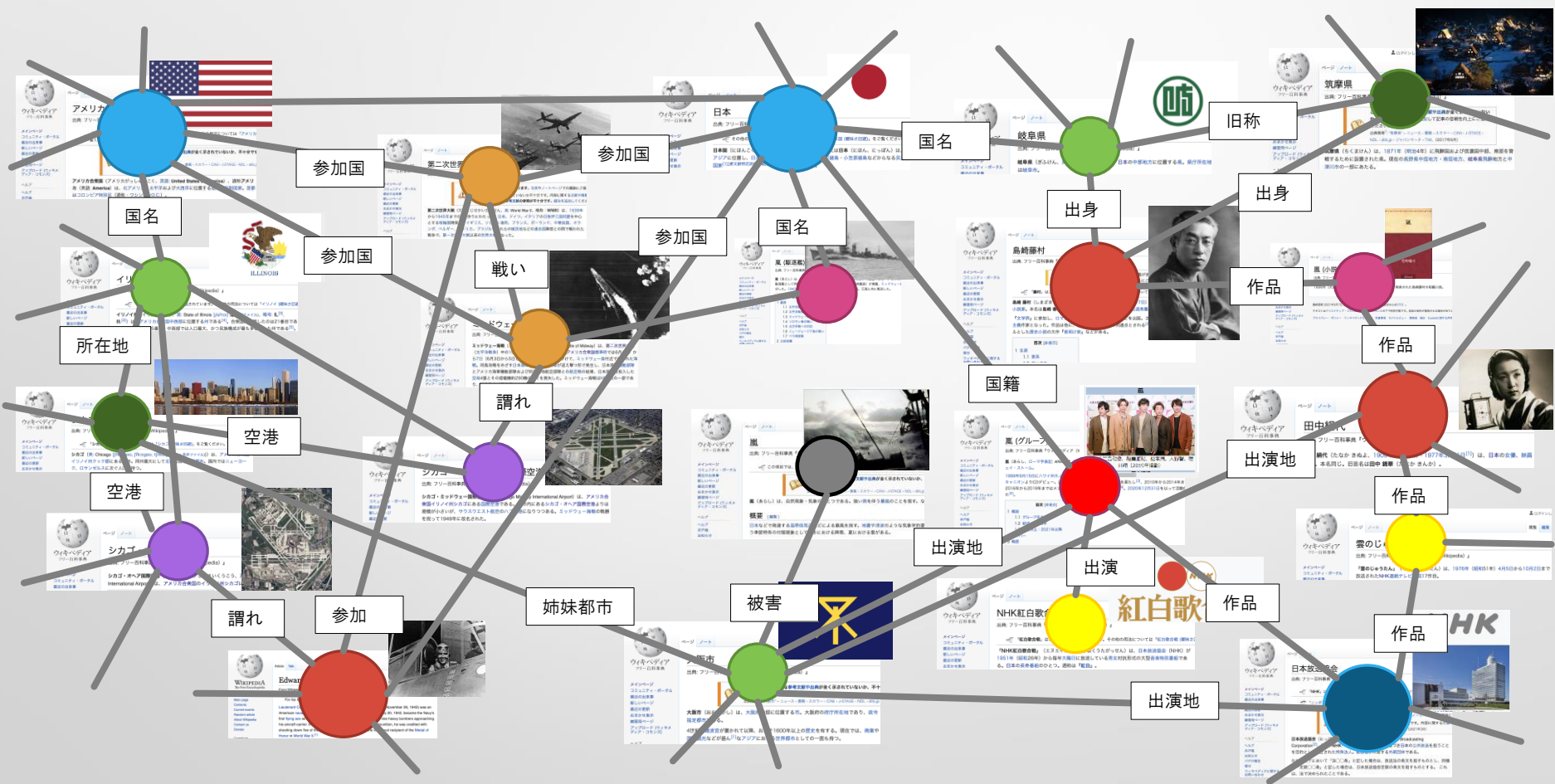


Background collage of various Wikipedia pages with highlighted categories: 国名 (Country Name), イベ (Event), 州 (State), 都市 (City), 空港 (Airport), 人名 (Person Name), 県 (Prefecture), 本名 (Real Name).

属性	属性値
国	アメリカ合衆国
所在地	イリノイ州
設立	1837年3月4日
人口	2695万人
人口データの年	2010年
空港	オヘア空港、ミッドウェー空港
...	...

属性	属性値
本名	島崎春樹
生年月日	1872年3月25日
国籍	日本
出身地	筑摩郡（現在の岐阜県）
地位職業名	小説家
作品	若菜集、嵐、夜明け前、。。
...	...

ステップ3 (リンキング)





物知り博士



第二次世界大戦に由来した名前を持つ2つの空港がある米国の都市はどこ？

第二次世界大戦の英雄に由来するオヘア空港と戦場名に由来するミッドウェー空港があるシカゴ！

雑談対話

おじいちゃんの好きな田中絹代が出演した1957年の「嵐」って映画の原作は、島崎藤村の小説なんだって



情報アクセス

教育支援

営業支援

認知症予防

介護福祉

多言語情報アクセス

他にも。。。

特定応用展開

観光

外国語教育

ビジネス

特許検索・分析

法律文書解析

オンライン医療

健康自己管理



WSの今後



- 2日目（10月末）
 - 質疑応答
 - Slackで挙げた質問なども含む
 - ベースラインを超えたシステムの紹介
 - 精度向上のためのヒントなど共有
- BERTによる他の属性抽出タスクの紹介
- 森羅プロジェクトの紹介



学生アルバイト募集

このWSのその先に



- 森羅プロジェクトでアルバイトを募集しています
 - 対象：B4~D1を対象
 - 形態：アルバイト（1年程度、週に15時間前後、リモート可）
インターン（2~3ヶ月程度集中的に。リモート相談）
 - 業務内容：森羅プロジェクトのシステム開発、データ開発
指導教官との相談の上、卒論、修論とすることも可能
理研のRIDEN（超巨大GPUマシン）を利用可能
- 社会人の兼業希望者も応相談

相談、希望は下記までメールを
satoshi.sekine@riken.jp