



学生・若手研究者のためのBERTワークショップ（1回目）
～BERTを使った自然言語処理の開発を
2回のワークショップで体験！～

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



- 今、流行りのBERTって何？ (30分)

- レクチャーによる説明

事前知識がない方は
難しいと思っても大丈夫！
ビデオで再学習できます

- BERTをとりあえず動かしてみよう(60分)

- Google Colabを使って自分で動かしてみる
 - **Googleアカウント、Google Drive上に600MB必要**
- コードの説明

- リーダーボード、コミュニティーの紹介 (15分)

Slackで連絡を
行います。
只今、zoomの
チャットで招待
リンクを共有

- 森羅プロジェクトの説明、今後の説明 (5分)

森羅Slack : shinra2022.slack.com

チャンネル : #bert_workshop2022

招待リンク : https://join.slack.com/t/shinra2022/shared_invite/zt-1dlrhtpty-eCBovUEWpJ1k4Uhugh4Fw



チームメンバー募集

このWSのその先に



- (BERTを使って) 森羅の研究に取り組みたい学生を募集
 - 最大3名 (だいたいB4~D1を対象)
 - チームにアルバイト、インターンとして参加
 - 期限はまず今年度末 (相談の上、来年度の延長も可能)
 - 2回目 (8月10日) に応募方法をアナウンス

本日の講師の中山功太さんも、
2017年のインターン学生

- 社会人の希望者があれば、応相談

リーダーボード作成等の門脇
さんも、社会人参加者



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.

学生・若手研究者のためのBERTワークショップ（1回目）
～BERTを使った自然言語処理の開発を
2回のワークショップで体験！～

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

クロージング





本日の概要



- 今、流行りのBERTって何？ (30分)
 - レクチャーによる説明
- BERTをとりあえず動かしてみよう(60分)
 - Google Colabを使って自分で動かしてみる
 - Googleアカウント、Google Drive上に600MB必要
 - コードの説明
- リーダーボード、コミュニティーの紹介 (15分)
- 森羅プロジェクトの説明、今後の説明 (5分)



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

- 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

今回のタスク

- 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)
- ミッドウェー海戦** (Battle of Midway)
- シカゴ** (Chicago)
- シカゴ・ミッドウェー国際空港** (Chicago Midway International Airport)
- シカゴ・オヘア国際空港** (Chicago O'Hare International Airport)
- Edward O'Hare**
- 日本** (Japan)
- 岐阜県** (Gifu Prefecture)
- 島崎藤村** (Ishikawa Futenjū)
- 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 Heaven)
- 日本放送協会** (NHK)

ステップ1 (分類)



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

- 国名** (Country Name): Labels for USA (アメリカ), Japan (日本), and Osaka Prefecture (大阪府).
- 戦争** (War): Labels for World War II (第二次世界大戦), the Battle of Midway (ミッドウェー海戦), and the atomic bombing of Nagasaki (原子爆弾投下).
- 州** (State): Label for Illinois (イリノイ州).
- 都市** (City): Labels for Chicago (シカゴ) and Chicago-Midway International Airport (シカゴ・ミッドウェー国際空港).
- 空港** (Airport): Label for Chicago-Midway International Airport.
- 人名** (Person Name): Labels for Uchiyama Gudō (内田 良平), Uchiyama Gudō (内田 良平), and Uchiyama Gudō (内田 良平).
- 船** (Ship): Label for the USS Arizona (阿佐ヶ浦丸).
- 名詞** (Noun): Label for the term '艦' (Ship).
- 組織** (Organization): Label for the group '嵐 (グループ)' (Arashi).
- 番組** (Program): Labels for NHK Red and White Song Battle (NHK紅白歌合戦) and NHK Red and White Song Battle (NHK紅白歌合戦).
- 企業** (Company): Label for NHK (NHK).
- その他** (Other): Labels for '都' (Capital), '本' (Original), and '不' (Not).

ステップ2 (属性値抽出)

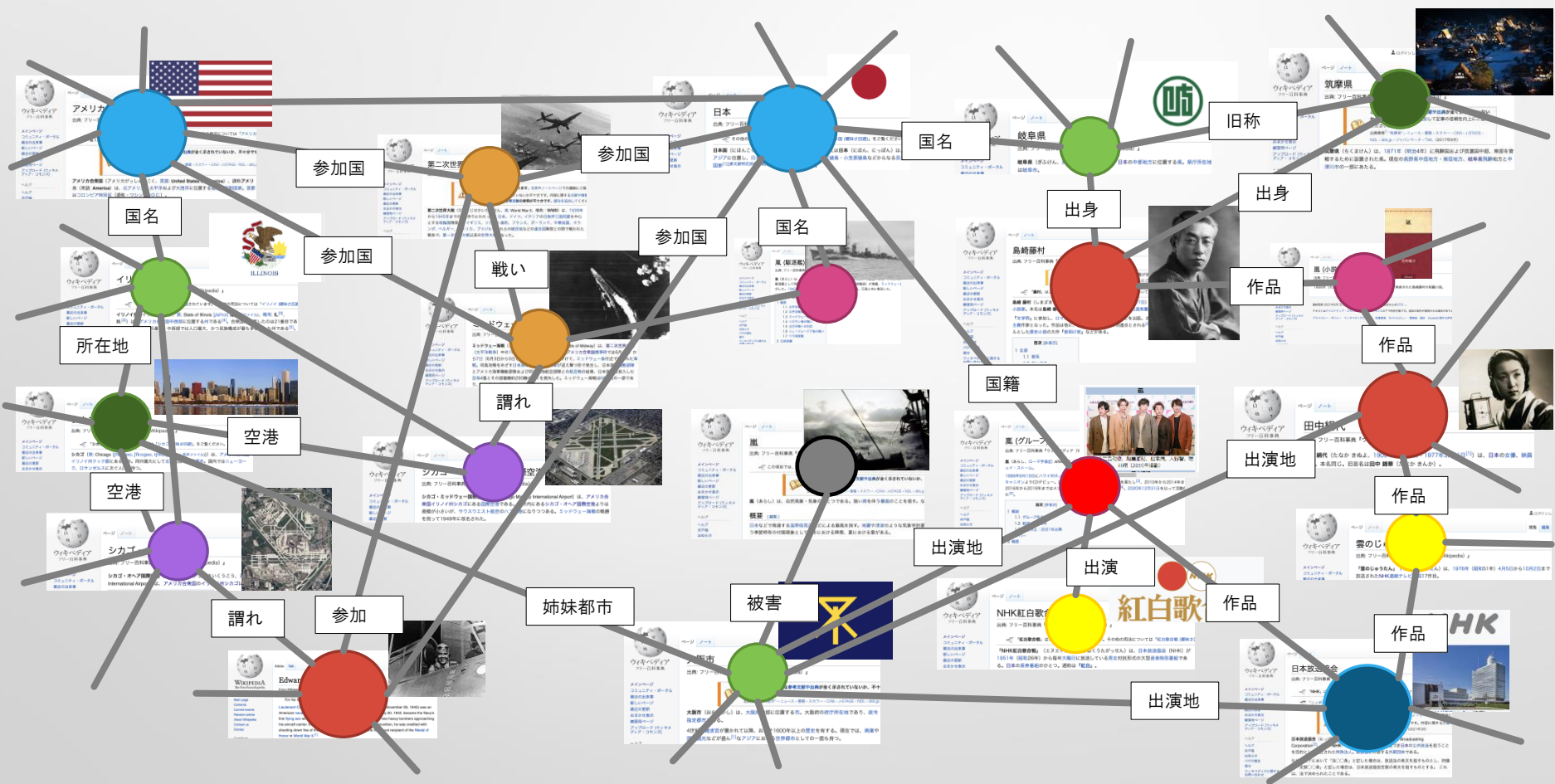


Background collage of various Wikipedia pages with highlighted categories: 国名 (Country), 州 (State), 都市 (City), 空港 (Airport), イベ (Event), 岐阜県 (Gifu Prefecture), 島崎藤村 (Ishikawa Fumoto), 本名 (Real Name).

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

他にも。。。

多言語情報アクセス

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



- 2日目（8月10日）13時

Doorkeeper: <https://c5dc59ed978213830355fc8978.doorkeeper.jp/events/141143>

- 質疑応答
 - Slackで挙がった質問なども含む
 - 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による他の自然言語処理タスクの紹介
- その後：BERTを使った色々なタスク（予告）
 - 今回は「分類」をBERTで実現
 - 次回は「属性値抽出」をBERTで実現
 - 次次回は。。。



森羅チームメンバー募集

このWSのその先に



- 森羅のテーマに研究として取り組みたい学生を募集
 - 最大3名（だいたいB4~D1を対象）
 - チームにアルバイト、インターンとして参加
 - 期限はまず今年度末（相談の上、来年度の延長も可能）
 - 2回目（8月10日）に応募方法をアナウンス

本日の講師の中山功太さんも、
2017年のインターン学生

- 社会人の希望者があれば、応相談

リーダーボード作成等の門脇
さんも、社会人参加者