

# Wikipediaから文化の広がりを理解する

村山太一  
(大阪大学産業科学研究所)



## 背景

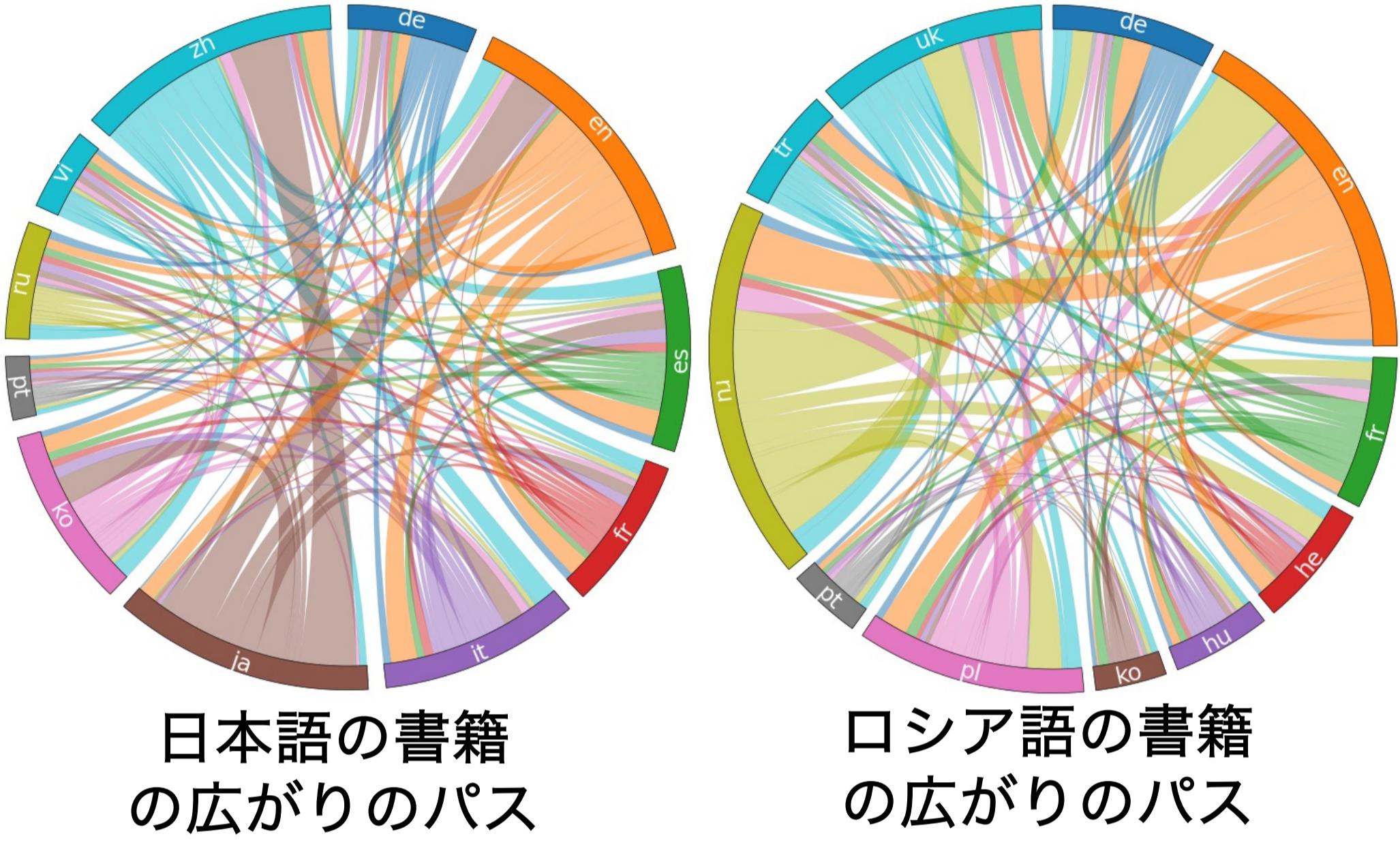
- WikipediaはOpen Knowledge Graphの一種
- Wikipediaには300を超える言語版があり、5,570万記事、2,040万のEntityが登録されており、毎月25万記事が作成
- Wikipediaの各記事はボランティアなどによって編集され、各言語の文化や人気が反映される

## Question・問題提起

- 各言語のWikipediaの編集履歴やEntityの登録履歴を分析することで、
- どのような文化が広がっているのか？
  - 言語を超えて広がる文化はどのような特徴があるのか？
  - 文化が広がるためにキーとなる言語が存在するのか？
  - コンテンツに含まれる要素の広がりはどのようなものか？

上記のことがわかると、**文化と言語の関係が明らかになるとともに、商品やアイテムを世界に展開したい場合に必要となってくる要素が明らかになる**

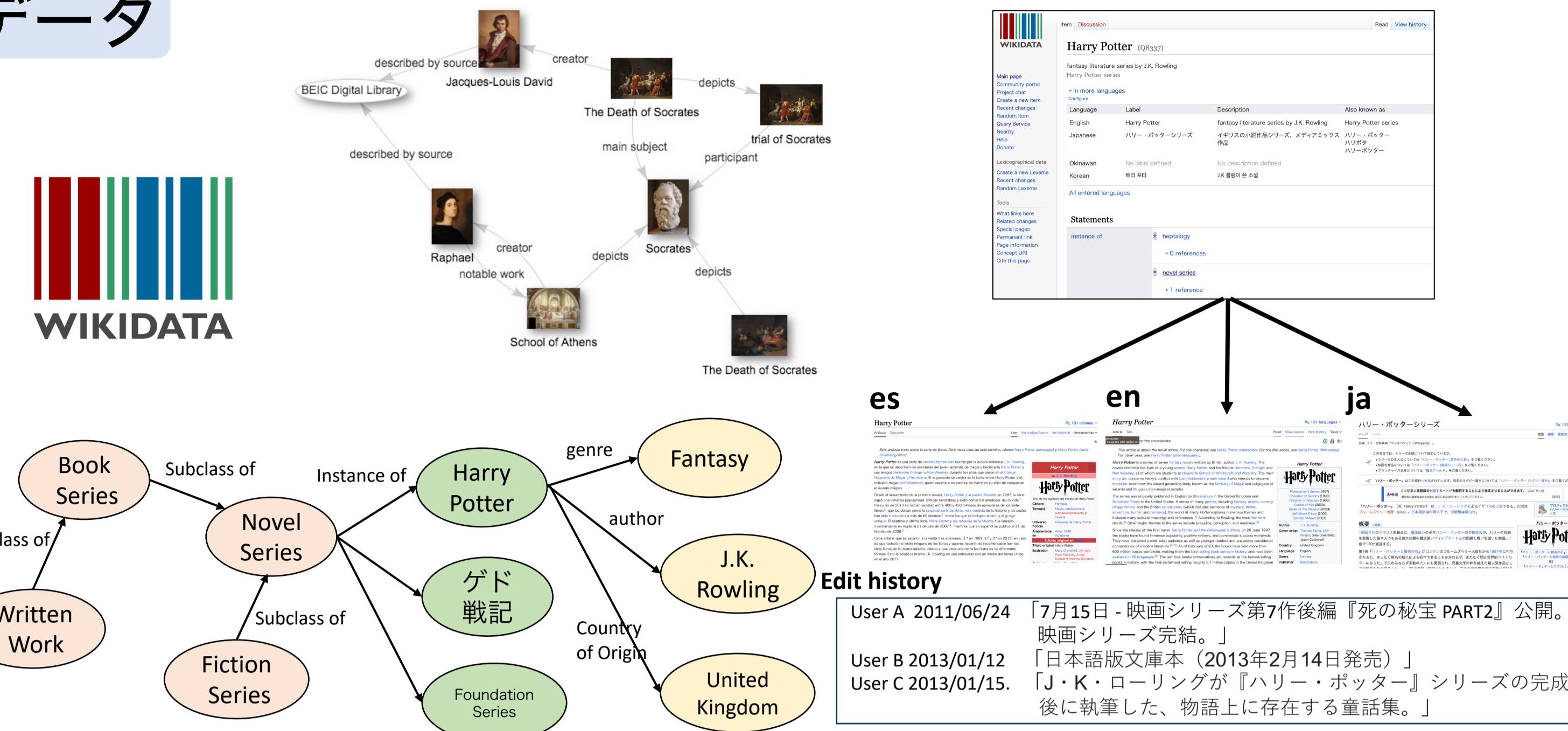
## マクロの観点での取り組み



### 5年後に複数言語にページが作成されているか各特徴のオッズ比

Business:	6.23
Crime Novel:	5.29
encyclopedia:	4.20
Isekai (異世界物):	4.08
アイルランド語:	4.41
カタルーニャ語:	3.32
ラテン語:	2.92
ハンガリー語:	2.63
オフィシャルサイト:	1.70
FreeBase登録:	1.60

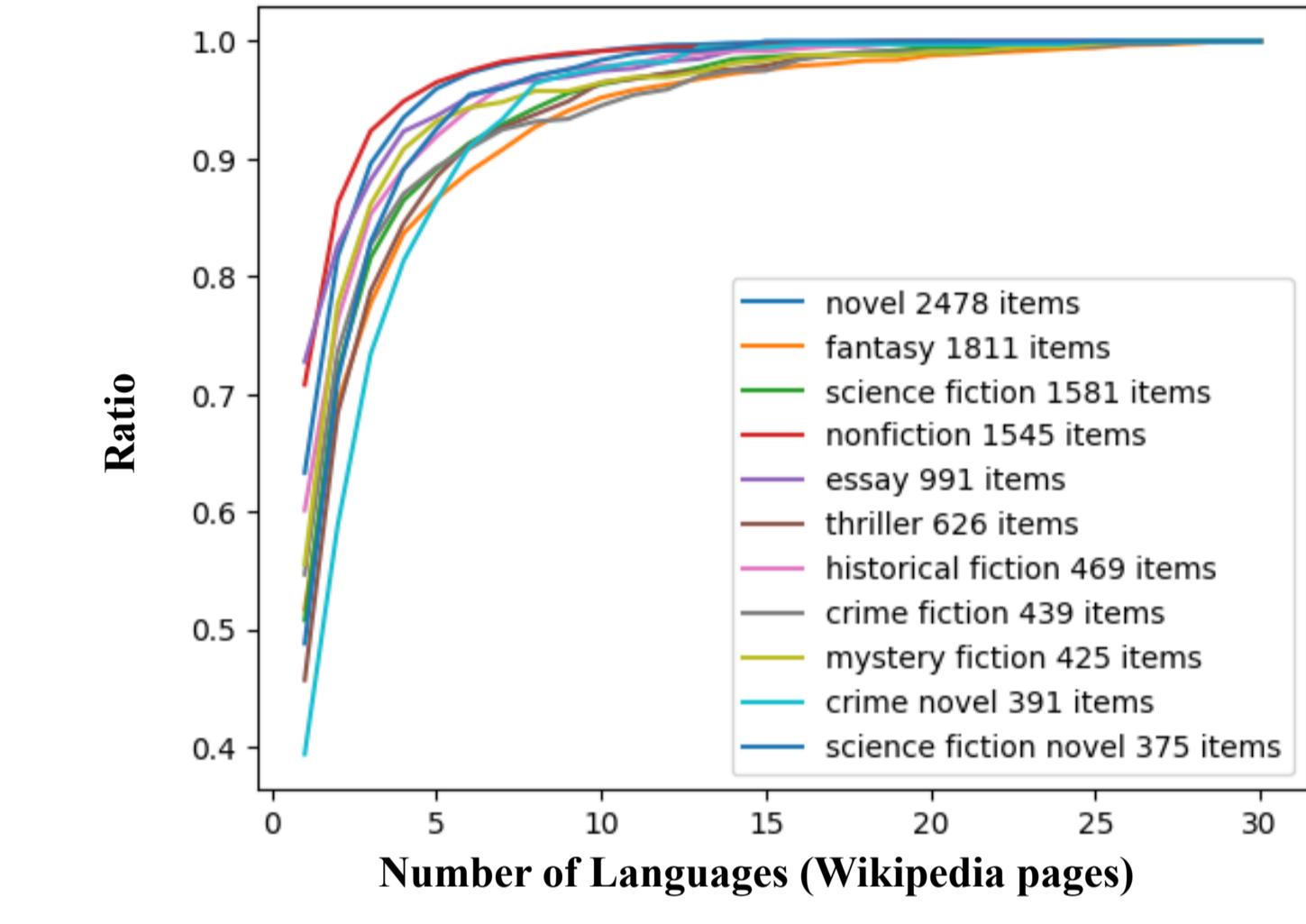
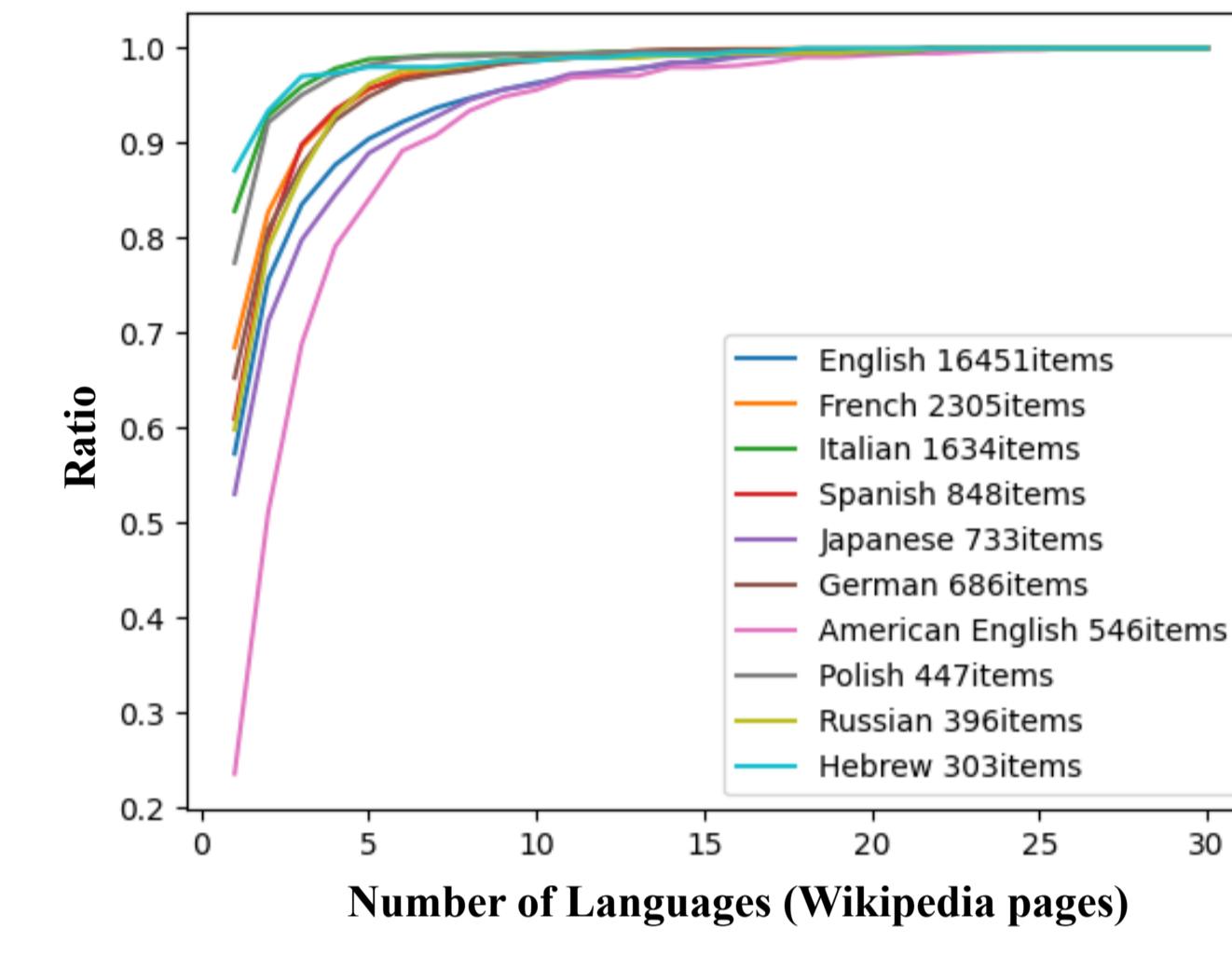
## データ



- Wikidataは階層的なオントロジー構造で、各データが相互にリンクされている多言語対応のKnowledge Database
- 特定のクラスに属する、Wikipediaの編集履歴を獲得

### Written Workより下のSubclassに属する2000年以降に出版されたWikidata (書籍)が対象 (全43,058 Entity)

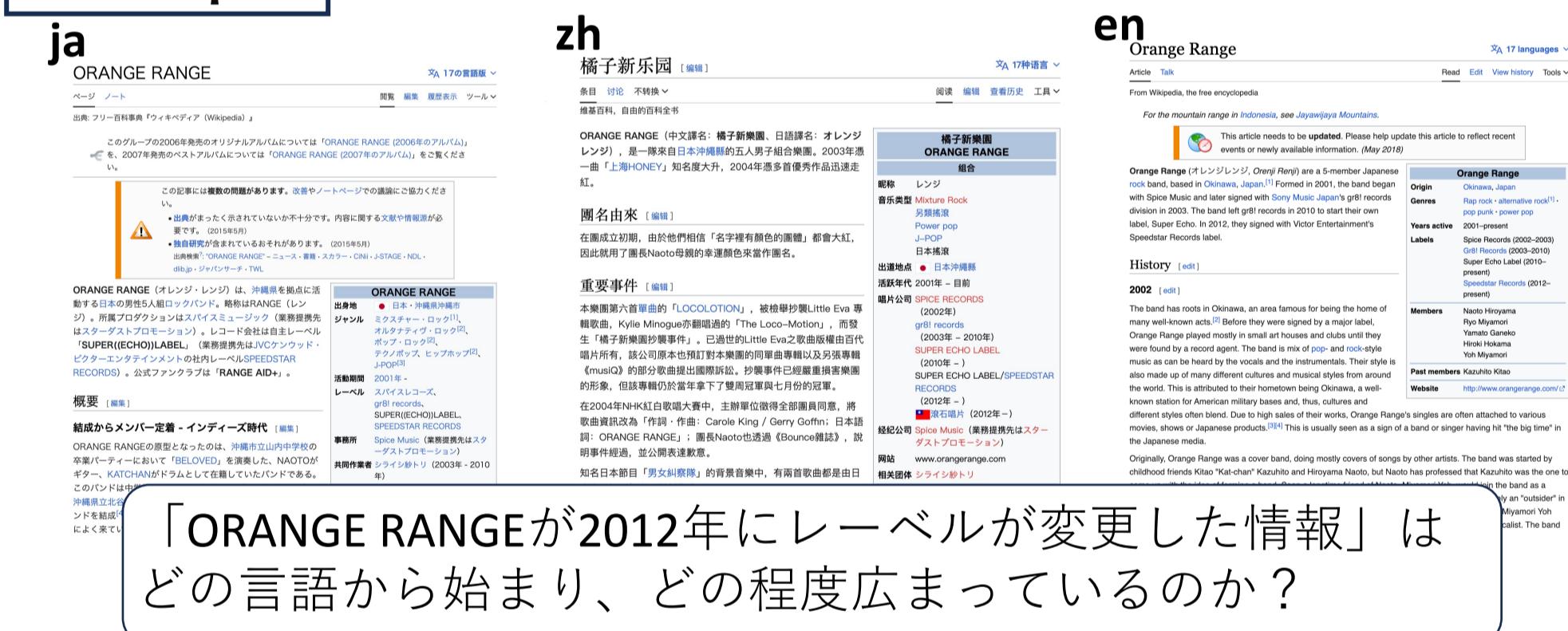
各Entityがどれだけの言語のWikipediaにページが作成されているかを調査



## ミクロの観点での取り組み

各言語のWikipediaの編集履歴を見ることで、コンテンツの各要素がどこから生まれ、どの要素が世界に広まっているかが理解できるのではないか？

### Example



### 手法の概要

- 各言語のWikipediaで同じ意味の文を抽出し、どの言語版でその意味の文書が生まれたのか？どれくらい広がっているかについて示す
- 文の埋め込み表現の抽出にはMultilingual Sentence-BERTを採用 (huggingfaceで提供のdistiluse-base-multilingual-cased-v2)
- 文のクラスタリングにはコサイン類似度を採用
- 各言語のWikipediaに含まれる文を3つのクラスに分類
- その言語ではじめて生まれ、世界的に広がっている文
- その言語で生まれたが、世界に広がっていない文
- 他の言語で生まれ、輸入してきた文

## アルゴリズム

### ① Wikipediaの編集履歴の前処理

i : 編集履歴をSentenceごとに分割

2011年6月24日追加  
「8番目の物語（=事実上の第8巻）」と銘打たれている文はト書きと台詞で構成されており、第1巻 - 第7巻のような小説のスタイルはとられていない。またジョン・ティファニー、ジャック・ソーンとローリングの共著名義である。

分割

2011/06/24 「8番目の物語（=事実上の第8巻）」と銘打たれている文はト書きと台詞で構成されており、第1巻 - 第7巻のような小説のスタイルはとられていない。またジョン・ティファニー、ジャック・ソーンとローリングの共著名義である。

ii : 言語に関わらず時系列順に並び替え

Date	Lang	Sentence
2011/06/24	ja	またジョン・ティファニー、ジャック・ソーンとローリングの共著名義である。
2011/06/24	es	Además, EA produjo en 2003 un juego de simulación de Quidditch: Harry Potter: Quidditch World Cup.
2011/06/24	en	Some of the translators hired to work on the books were well-known authors before their work on Harry Potter, such as Viktor Golyshov, who oversaw the Russian translation of the series' fifth book.
.	.	.
2019/02/23	fr	Le trio retourne donc au château, très vite attaqué par Voldemort.
2019/02/24	pt	Arrecadando em bilheterias aproximadamente 7,7 bilhões de dólares;
2019/02/24	zh	學者與記者的解讀更多元，有的還包括政治。大致上，對於小說系列主體的解讀有

### 時系列順に各文章を適用

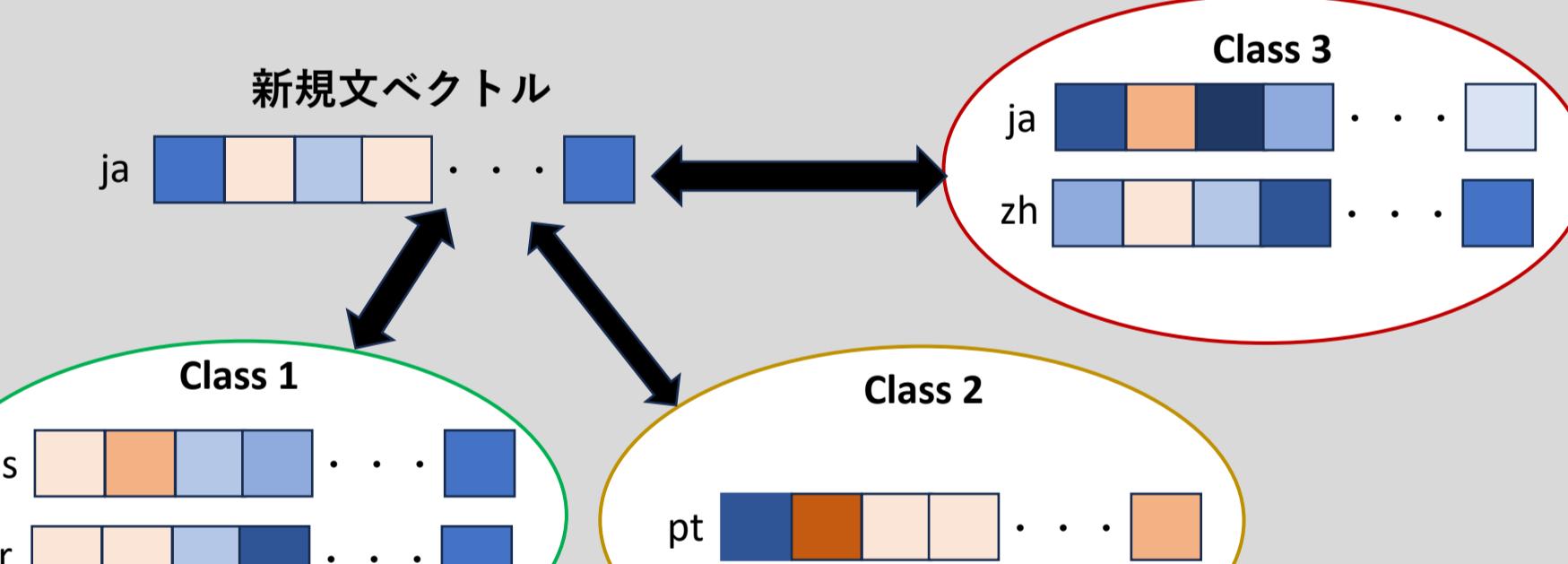
### ② 文章の埋め込み表現獲得

またジョン・ティファニー、ジャック・ソーンとローリングの共著名義である。

Multilingual Sentence Transformers

### ③ 文章のクラスタリング

i 以前出現した文とのコサイン類似度



閾値以上の類似度

ii 既存クラスへの追加

ja

zh

es

fr

pt

Class 2

Class 3

閾値以下の類似度

ii 新クラスの追加

ja

zh

es

fr

pt

Class 3

Class 4

課題/質問

- 情報がない場合は、そもそも別ページに作成されている可能性がある

在籍中のメンバー

吉岡 聖恵 (よしおか きよえ, 1984年2月29日 - ) (39歳) : ポーカル

水野 良樹 (みずの よしき, 1982年12月17日 - ) (40歳) : ギター・ピアノ

旧メンバー

・山下 鶴穂 (やました ほたか, 1982年8月27日 - ) (40歳) : ギター・ハモニカ

Current members [edit source]

Kyoko Honda (吉岡 聖恵, Kyoko Honda, born February 29, 1984) sings vocals.<sup>111</sup> She was born in the city of Shirooka, and later moved to Atsugi, Kanagawa when she was five years old. She went to City Minamisuna High school in Kanagawa Prefecture, and graduated from Showa University College of Music. Her brother was a high school classmate of Mizono and Yamashita. She became one of the first members of the band in 2004, and has been performing with the band ever since. She has contributed to some of their songs like "Kimigayo" and "GOLDEN GIRL", which are A-side songs, "Mira! Wakuusei" (from "Hajimete no Uta" album), Tokyo (from "Nirai" album), and "Shiro Diary" (from "NEWTRAL" album).<sup>112</sup>

Former members [edit source]

Kyo Honda (吉岡 聖恵, Kyoko Honda, born February 29, 1984)

He was born in Hamamatsu, Shizuoka Prefecture, but he later moved to the city Ebina, Kanagawa Prefecture, at an early age.

It was announced that he registered his marriage on August 17, 2013, to a woman who has no connection to the media industry. He has released no further information about her identity.

2005年3月26日に初のホールでのワンマンライブを厚木市文化会館小ホールで開催した。

● 他に良い手法が無いか？

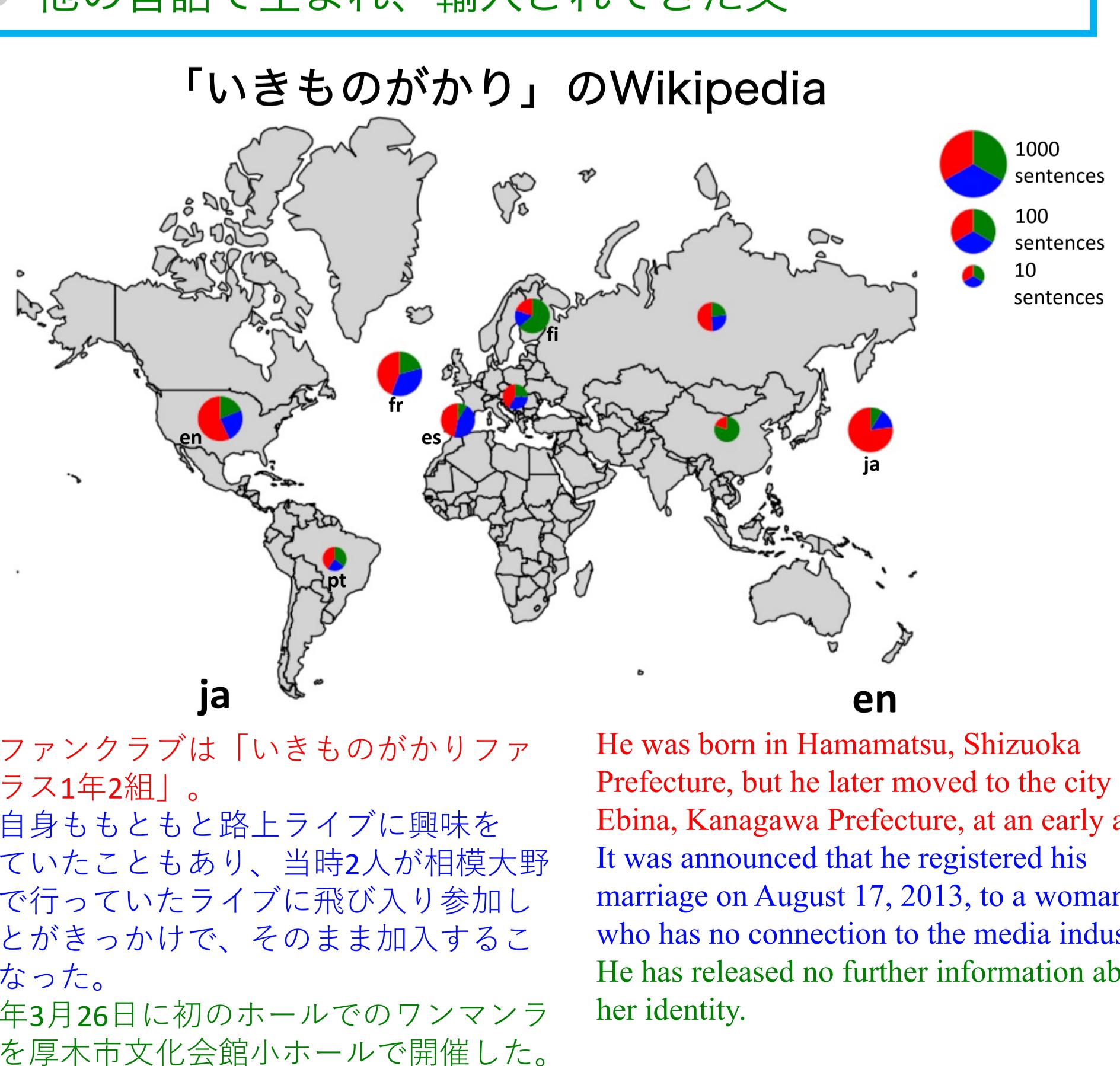
- Entity Extractionは検討したが、表記揺れの問題で多言語のEntityを一致できない

いきものがかり = Ikimonogakari = 生き物係？

= Ikimonogakari = 生き物係？

● 評価をどうするか？どのようにして論文として落とし込むか？

- Information Gapの解決という切り口はあるが…



公式ファンクラブは「いきものがかり」

クラス1年2組」。

吉岡自身もともと路上ライブに興味を

持っていたこともあり、当時2人が相模大野

駅前で行っていたライブに飛び入り参加し

たことがきっかけで、そのまま加入するこ

ととなった。

2005年3月26日に初のホールでのワンマンラ

イブを厚木市文化会館小ホールで開催した。