

Beyond Automatic Translation: Aligning Wikipedia Sections Across Multiple Languages

Bahodir Mansurov, Diego Sáez-Trumper,
Robert West, Leila Zia

Research Showcase
2018-03-21

**Create a data set with
section names
aligned across languages**

Why?

- Cross-lingual section recommendations.
- Improve the content translation tools.
- Towards generate an abstract ontology for section titles.

Награды
Russian

Биография
Russian

賞
Japanese

伝記
Japanese

Awards
English

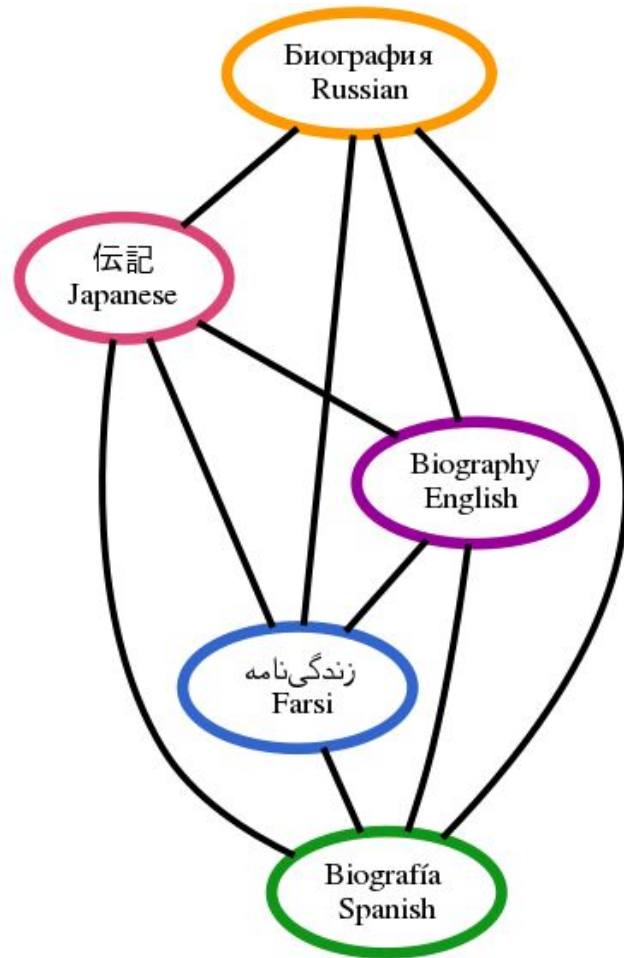
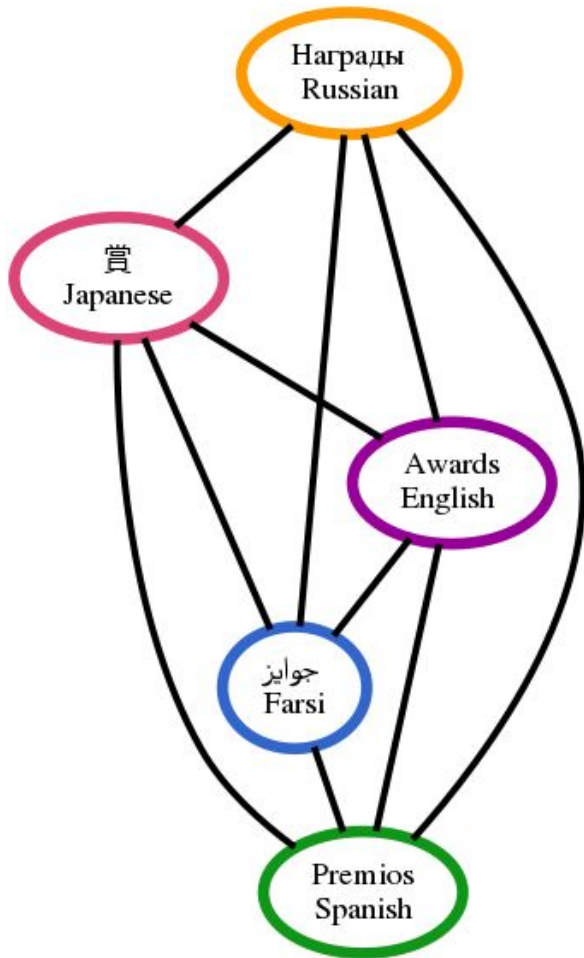
Biography
English

جوایز
Farsi

زندگی نامه
Farsi

Premios
Spanish

Biografía
Spanish



**Can we trust
Automatic Translation
services?**

Can we **trust** Automatic Translation services?

**Accuracy depends on language
(ex. we found very different results for English to
Spanish, and English to Farsi)**

Challenges & Constraints

- Keep the **style and conventions** of each wikipedia edition.
- **Equal importance** to all languages
- Ground truth is difficult to build (e.g. Latvian to Bengali)



Assets & Opportunities

- Large and active **community**.
- Self **reported language** skills (Babel Template).
- Entities links across languages
- Large set of pre-trained models using Wikipedia (e.g. **babylon project**).



Building a training dataset

- Select a set of diverse languages:
 - Different scripts
 - Different families
- اللغة العربية المعيارية الحديثة
- Русский язык
- Français
- Español
- English
- 日本語

Building a training dataset

- Select a set of diverse languages:
 - Different scripts
 - Different families
- **Modern Arabic**
- **Russian**
- **French**
- **Spanish**
- **English**
- **Japanese**

Wikipedia:Babel

en-5	This user is able to contribute with a professional level of English .
sv	Den här användaren talar svenska som modersmål .
no-4	Denne brukaren/brukeren meistar/behersker norsk på morsmålsnivå .
he-3	משתמש זה מסוגל לתרום ברמה מתקדמת של עברית .
lt-2	Šis vartotojas gali prisidėti prie projekto vidutinio lygio lietuvių kalba .
es-1	Este usuario puede contribuir con un nivel básico de español .
an-0	Iste usuario no repleca l' aragonés (u el repleca con prou dificultat).

```
SELECT user.user_name,  
       babel.babel_lang AS lang,  
       babel.babel_level AS level  
FROM babel  
LEFT JOIN user  
       ON user.user_id = babel.babel_user  
ORDER BY user.user_name ASC;
```

<https://quarry.wmflabs.org>

Wikipedia:Babel

en-5	This user is able to contribute with a professional level of English .
sv	Den här användaren talar svenska som modersmål .
no-4	Denne brukaren/brukeren meistarar/behersker norsk på morsmålsnivå .
he-3	משתמש זה מסוגל לתרום ברמה מתקדמת של עברית .
lt-2	Šis vartotojas gali prisidėti prie projekto vidutinio lygio lietuvių kalba .
es-1	Este usuario puede contribuir con un nivel básico de español .
an-0	Iste usuario no repleca l' aragonés (u el repleca con prou dificultat).

language pair	# of users
ar-en	139
ar-es	43
ar-fr	113
ar-ja	8
ar-ru	42
en-es	5991
en-fr	2796
en-ja	203
en-ru	5324
es-fr	2796
es-ja	48
es-ru	236
fr-ja	43
fr-ru	399
ja-ru	55
total	18236

```
SELECT us  
  ba  
  ba  
FROM babe  
LEFT JOIN  
  ON u  
ORDER BY
```

ng,
evel

babel_user

<https://>

[.org](https://)

Problem Definition

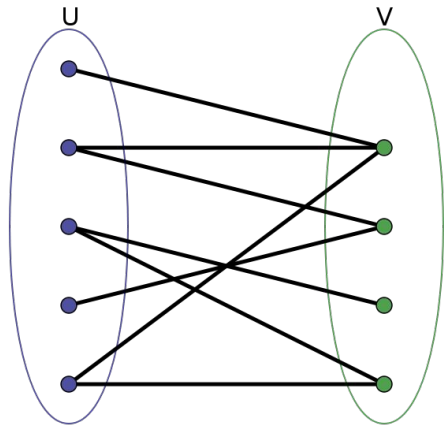
- Determine for each pair of sections $(S1, S2)$, where $S1$ and $S2$ are from different languages, whether $S1$ is the translation of $S2$, across multiple languages.

Problem Definition

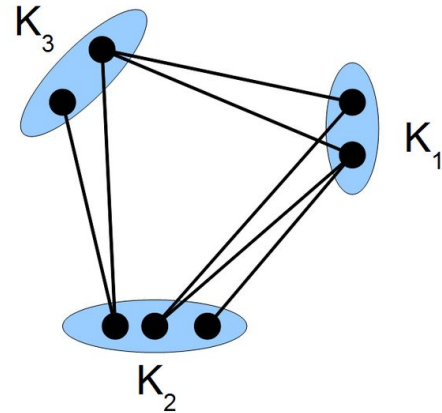
- Determine for each pair of sections (S_1 , S_2), where S_1 and S_2 are from different languages, whether S_1 is the translation of S_2 , across multiple languages.
- Consider a **link prediction** task, in **multipartite (k-partite) graph**, where each k group correspond to a language, nodes are section titles, and links represent translation.

Multipartite (k-partite) Graph

- k-partite graph is a graph whose vertices are or can be partitioned into k different independent sets.
- $k=2$ -> bipartite graph (ex. Actors and Movies)
- In our case k is equal to the number in languages in our dataset.



K = 2



K = 3

Features

- Automatic Translation
- Levenshtein distance
- Outcoming links
- **Word embeddings**
- **Co-occurrence counts**

Embeddings

"When some object X is said to be embedded in another object Y , the embedding is given by some injective and structure-preserving map $f : X \rightarrow Y$ "



<https://en.wikipedia.org/wiki/Embedding>

Distribution

Geography

Demographics

External links

Sources

References

Further reading

Bibliography

Filmography

Discography

Publications

Works

Gallery

Description

See also

History

Track listing

Notes

Cast

Biography

Early life and education

Personal life

Production

Family

Background

Early life

Life

Reception

Legacy

Death

Personnel

Overview

Charts

Plot

Education

Results

Awards

Honours

Career

Club career

Filmografía
Discografía
Publicaciones
Bibliografía

Bibliography
Filmography
Discography
Publications
Works

Aligned word embeddings

- Babylon Project:
 - Dictionary based alignment across languages
 - Linear transformation

https://github.com/Babylonpartners/fastText_multilingual



Bibliografía

Works **Publicaciones**
Obras **Publications**

Discografía **Bibliography**
Discography **Filmography** **Gallery**
Filmografía

Co-occurrence counts



- We use Wikidata to align articles, counting the co-occurrence of sections in articles about the same item across languages

Contents [\[hide\]](#)

- 1 Early life and education
 - 1.1 Family
 - 1.2 Primary and secondary school years
 - 1.3 Undergraduate years
 - 1.4 Graduate years
- 2 Career
 - 2.1 1966-1975
 - 2.2 1975-1990
 - 2.3 1990-2000
 - 2.4 2000-2018
- 3 Personal life
 - 3.1 Marriages
 - 3.2 Disability
 - 3.3 Disability outreach
 - 3.4 Plans for a trip to space
- 4 Death
- 5 Personal views
 - 5.1 Future of humanity
 - 5.2 Science vs. philosophy
 - 5.3 Religion and atheism
 - 5.4 Politics
- 6 Appearances in popular media
- 7 Awards and honours
 - 7.1 Stephen Hawking Medal for Science Communication
- 8 Publications
 - 8.1 Popular books
 - 8.1.1 Co-authored
 - 8.1.2 Forewords
 - 8.2 Children's fiction
 - 8.3 Films and series
 - 8.4 Selected academic works
- 9 Notes
- 10 References
 - 10.1 Sources
- 11 External links

Stephen Hawking (Q17714)



Índice [\[ocultar\]](#)

- 1 Biografía
 - 1.1 Primeros años y educación
 - 1.2 Carrera
 - 1.2.1 De 1962 a 1975
 - 1.2.2 De 1975 a 2018
 - 1.3 Fallecimiento
- 2 Obra
 - 2.1 Investigación del universo
 - 2.1.1 Investigación sobre el origen del universo
 - 2.1.2 Conjetura de protección de la cronología
 - 2.2 Pensamiento filosófico
 - 2.3 Creencias religiosas
- 3 Lucha personal contra la esclerosis lateral amiotrófica
- 4 Reconocimientos
 - 4.1 Principales premios y distinciones
- 5 Publicaciones
 - 5.1 Selección de obras de Stephen Hawking
 - 5.1.1 Científicas y divulgativas
 - 5.1.2 Ficción infantil
 - 5.1.3 Películas, documentales y series
 - 5.2 Vídeos musicales
 - 5.3 Literatura sobre Stephen Hawking
- 6 Véase también
- 7 Notas
- 8 Referencias
- 9 Bibliografía
- 10 Enlaces externos

**Please help us to create the
training dataset**

**[\[\[m:Research:Expanding_Wikipedia_articles_across_languages\]\]](#)
T183039**

Synonym Detection

Problem Definition

- Determine for each pair of sections (S_1 , S_2) from the same language whether they are synonyms of each other.

Features

- Word embeddings
- Levenshtein distance
- IsSubset
- Tf-Idf similarity (*)

Synonym Detection

- **Intuition:** Two sections are likely to be synonym if both tends to co-occur with similar sections, but (almost) never co-occurs among themselves.

Synonym Detection

- **Intuition:** Two sections are likely to be synonym if both tends to co-occur with similar sections, but (almost) never co-occurs among themselves.
- **Problem:** Most frequent sections tends to co-occur with almost all the sections.

Synonym Detection

- **Solution:**
 - Represent sections (S) as **TF-IDF vectors**, using the counting of co-occurrences as TF.
 - Measure cosine similarity of TF-IDF vectors, between pairs of sections S1 and S2
 - Divide by the number of co-occurrences (+1) between S1 and S2

Examples

S1	S2	tflidfSimilarity
Stud career	Stud record	0.99
Side effects	Adverse effects	0.99
Musical career	Music career	0.98
Home video	Home media	0.97
Line-up	Band members	0.96

Future Work & Tasks

- Gather labels for section alignment and synonyms.
- Test frameworks for the link detection.
- Improve word embedding alignments.

References

- K-partite graphs: https://en.wikipedia.org/wiki/Multipartite_graph
- Embeddings: <https://en.wikipedia.org/wiki/Embedding>
- Word Embeddings: https://en.wikipedia.org/wiki/Word_embedding
- Babylon Project:
https://github.com/Babylonpartners/fastText_multilingual
- Code: <https://github.com/digitalTranshumant/wmf-interlanguage>

Questions?