# Automatic approaches to bridge knowledge gaps in Wikimedia projects

## WMF Research Team

**WIKIMEDIA**
FOUNDATION

Wikimedia Research Showcase October 2021

# The Research Team

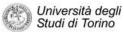Leila Zia
*Director, Head of Research*

Pablo Aragón
*Research Scientist*

Martin Gerlach
*Research Scientist*

Isaac Johnson
*Research Scientist*

Fabian Kaelin
*Senior Research Engineer*

Emily Lescak
*Senior Research Community Officer*

Miriam Redi
*Research Manager*

Diego Sáez-Trumper
*Senior Research Scientist*

UNIVERSITY OF CAMBRIDGE

NATIONAL CHENG KUNG UNIVERSITY 1931

Università degli Studi di Torino

UNIVERSITY OF MICHIGAN

Telefónica

UNIVERSITY OF MINNESOTA

EPFL
ÉCOLE POLYTECHNIQUE FÉDÉRALE DE LAUSANNE

UFES
UNIVERSIDADE FEDERAL DO ESPÍRITO SANTO

UNIVERSITY of WASHINGTON

香港浸會大學
HONG KONG BAPTIST UNIVERSITY

جامعة نيويورك ابوظبي
NYU ABU DHABI

ibs Institute for Basic Science
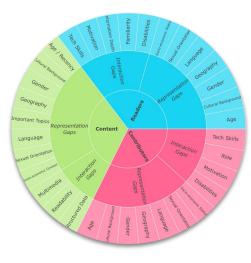
Learn more: https://research.wikimedia.org

# The strategic direction: Knowledge Equity
## [from Movement Strategy]

**Knowledge equity**: As a social movement, we will focus our efforts on the **knowledge and communities that have been left out by structures of power and privilege**. We will welcome people from every background to build strong and diverse communities. We will **break down the social, political, and technical barriers** preventing people from accessing and contributing to free knowledge.

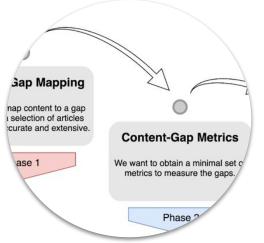# The research program: Addressing Knowledge Gaps

### Identify gaps

September 2020:

https://w.wiki/4DWU
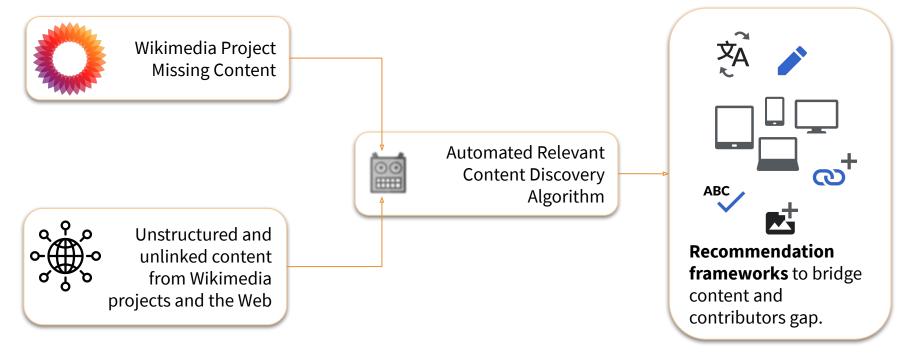


### Measure gaps

Ongoing:

https://w.wiki/4Exh



### Bridging gaps

Today:

https://w.wiki/4DWe

# Bridging knowledge gaps: Automated tools for content discovery and recommendation.



Wikimedia Project Missing Content

Unstructured and unlinked content from Wikimedia projects and the Web

Automated Relevant Content Discovery Algorithm

**Recommendation frameworks** to bridge content and contributors gap.

# Bridging knowledge gaps:
# Previous research work from our team

# Bridging gaps: Developing automatic tools

- Section alignment (Diego Sáez-Trumper) https://w.wiki/Bb6
  - In Collaboration with the Language team

- Link recommendation (Martin Gerlach) https://w.wiki/4BCP
  - In Collaboration with the Growth team

- Image recommendation (Miriam Redi) https://w.wiki/4DW8
  - In Collaboration with the Growth and Structured Data teams

- Equity in recommendations (Isaac Johnson) https://w.wiki/$4P

# Section alignment

Diego

# **Problem**: **Translation engines can not be directly used on Wikipedia**

- We want to compare sections and share content across languages.

- We can't compare sections across languages.

# Solution: Cross-lingual section alignment

- Create an ML system to align sections across languages.

# Step 1: Let's try to do better than Automatic Machine Translation

- Ask the community to help annotate data:

  - ar,es,fr,ja,ru,en (all pairs)

- Cross-lingual word-embeddings + Wikipedia specific features.



WikiAligments vs Google Translate

# Step 1: Let's try to do better than Automatic Machine Translation

✓

- Ask the community to help annotate data:

  - ar,es,fr,ja,ru,en (all pairs)

- Cross-lingual word-embeddings + Wikipedia specific features.

# Learn more and test the system!

[[m:Research:Expanding_Wikipedia_articles_across_languages/Inter_language_approach]]

https://secrec.wmflabs.org

# Step 2: Scale up to 100+ languages

- Add new languages and deploy.

- Test new embeddings and language models.

- Improve the recommender system.

# Link recommendation

**Martin Gerlach** (Research Team)
Growth Team (Marshall Miller, Rita Ho, Kosta Harlan, many more)
Djellel Difallah (NYU Abu Dhabi)

WIKIMEDIA
FOUNDATION

# Problem

# Editing is hard



**Technical**
*What is an infobox?*

**Conceptual**
*What is notability?*

**Cultural**
*Why are people so mean?*

# Structured task editing

- Break down editing into simpler tasks

  ○ Easier to: understand, do on mobile, get positive experience

- We picked the task of: "adding a link"

  ○ Well-defined, frequent, and attractive task type

- Machine-in-the-loop:

  ○ Generate recommendations using ML approaches.

  ○ Editors verify the output and validate the insertion.

# Link recommendation



**Hypatia** (born c. 350–370; died 415 AD) was a Hellenistic Neoplatonist philosopher, astronomer, and mathematician, who lived in Alexandria, Egypt, then part of the Eastern Roman Empire.

astronomer → Astronomy

astronomer → Astronomer

astronomer → --no link--

**?**

Entity-linking task

- Mention detection

- Link generation

- Link Disambiguation

# The Add-a-link Task in Wikipedia

- Language support

  - 300+ language version of Wikipedia; highest impact for smaller communities.

- Other considerations

  - Manual of style constraints

  - Prefer simpler over complex models (scalability, transparency, etc)

  - Utility: find a balance between precision over recall

# Step 1: Mention detection

- Build mention dictionary from all existing links
    - e.g. English Wikipedia: 7M anchors, 170M links
- String-matching sweeping window for all possible n-grams (N=1...10)
- Give preference to larger N



**Mention Dictionary**

**Step 1:** Mention detection

prominent thinker of the Neoplatonic school in Alexandria

# Step 2: Link generation

- From the anchor-dictionary extract all used links

- Drops links based on constraints (type-based)

- Drop links based on link probability heuristic:

  - Text to Link ratio < 6.5% (picked empirically, and supported by previous work)

# Step 3-a: Link disambiguation- features

- **N-gram size:** the number of tokens in the anchor (based on simple tokenization).

- **Frequency:** count of the anchor-link pair in the anchor-dictionary.

- **Ambiguity:** how many different candidate links exist for an anchor in the anchor-dictionary.

- **Kurtosis:** the kurtosis of the shape of the distribution of candidate-links for a given anchor in the anchor-dictionary

- **Levenshtein-distance:** a string similarity measure between the anchor and the link, e.g., the Levensthein-distance between "kitten" and "sitting" is 3.

- **Wiki2Vec Distance (entity embedding):** similarity between the article (source-page) and the link (target-page) based on the content of the pages.

# Step 3-b: Link disambiguation- classifier

- Extract fully linked sentences from the lead sections

  - Positive example: linked mention with correct link

  - Negative example: linked mention with incorrect link, unlinked mention

- Train a binary classifier (XGBoost)

# Evaluation

**Held-out testset** + **Manual evaluation**
(thanks: Bennoit Evellin, Habib Mhenni,
Martin Urbanec, Bluetpp, -revi)

**Tested Wikis**: Arabic, Bengali, Czech,
English, French, Vietnamese

**Precision**: 70% - 92%
How many suggestions are correct?

**Recall**: 30% - 66%
How many of the possible links captured?

# Link recommendation model

- Training pipeline for each language.

- Models/datasets published publicly

- Link recommendation API on kubernetes

## Get link recommendations

💬 Discussion   🕐 Updated 24 June 2021

`GET` `/service/linkrecommendation/v1/linkrecommendations/wikipedia/{language}/{title}`

Get a set of possible links that can be added to improve an article on Wikipedia.

https://api.wikimedia.org/wiki/API_reference/Service/Link_recommendation

# User interface

*Evaluate the suggestion*      *Feedback on algorithm*      *Edit summary*      *Next suggestion*

---

**✕  Suggestions**   **>**

Ke zdárnému růstu potřebuje sice tropické podnebí ✓, roste však i v místech se studenějším zimním počasí. S nástupem chladného počasí shodí listy a přečká ve stavu vegetačního klidu, dospělé rostliny přežijí i období s krátkodobými teplotami pod 0 °C. Preferuje vlhčí klima, špatně snáší suchý vzduch 🤖, obvykle prospívá až do nadmořské výšky 1200 m. Potřebuje hlubokou a výživnou půdu s dostatkem vlhkosti a dobrou drenáží, nesnáší delší sucho. Na nevyhovujícím stanovišti a v málo výživné půdě plodí slabě

🤖 1 / 7  ● ○ ○ ○ ○ ○ ○   ?

Text
**podnebí**

Should this link to the following article?

**Podnebí**
termín

‹  ✓ **Yes**   ✕ **No**   Next ›

---

**✕  Suggestions**   **>**

v indii jsou považovány za ovoce chudých.

Všechny části rostliny, nezralé plody i semena, obsahují alkaloidy, anonain, liriodenin, reticulin, které jsou toxické. Semena je možno bez

**Why is this not a good link?**
Your answers improve future suggestions.

● Almost everyone knows what it is

○ Linking to wrong article (e.g. linking "moon" to "planet")

○ Text should include more or fewer words (e.g. "palm tree" instead of "palm")

○ Other

**Done**

Text
**kůže**

Should this link to the following article?

**Kůže**
orgán pokrývající těla obratlovců

‹  ✓ Yes   ✕ **No**   Next ›

---

**‹  Save your changes**   **Publish**

## Summary

🤖 Suggestion                                      Linked?

podnebí:  🔗 Podnebí                              ✓

vzduch:  🔗 Vzduch                                 —

stálezelené:  🔗 Stálezelená rostlina             —

na dotek:  🔗 Na dotek                            ✓

období sucha:  🔗 Období sucha                   ✓

hnojiv:  🔗 Hnojivo                              ✓

kůže:  🔗 Kůže                                    ✕

By saving changes, you agree to the Terms of Use 🔗 and agree to release your contribution under the CC BY-SA 3.0 🔗 and GFDL 🔗 licenses.

**Review your changes**

---

☰  **WIKIPEDIA**   🔍   **52**

## *Color Out of Space* (film)

Page    Discussion

🗚    ★    🕑    ✏

🛈 This article's plot summary may be too long or excessively detailed. *(February 2020)*  Learn more

***Color Out of Space*** is a 2019 American science fiction cosmic horror [8] film directed

✓ You've published an edit. Thanks and keep going!

Richardson, Madeleine, Arthur, Olariska,

**Edit another suggestion:**

🖼   **Jonathan Gresham**
🤖 📊  Add links between articles

**View more suggested edits**

**Close and edit this article again**

# In practice

- Results from pilot-wikis (Arabic, Bengali, Czech, Vietnamese)
  - Newcomers prefer structured editing
  - Revert rate of structured edits is much lower (7.9% vs 25.5% for unstructured tasks)
  - Careless editing is rare
  - Reactions from community mostly positive
- Deployment
  - Currently: Arabic, Bengali, Czech, French, Hungarian, Persian, Polish, Romanian, Russian, Vietnamese
  - Planned: Catalan, Hebrew, Hindi, Korean, Norwegian, Portuguese, Simple English, Swedish, Ukrainian
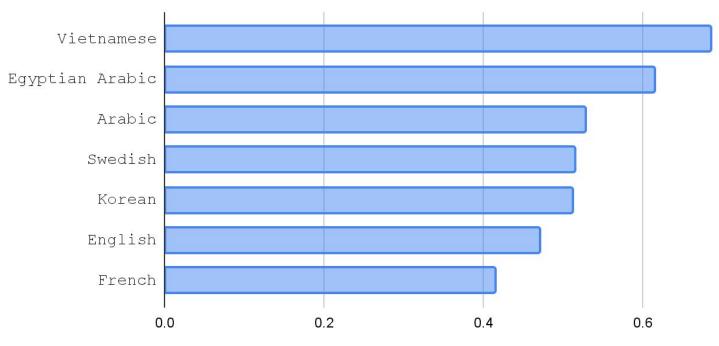
# Image recommendation

Miriam

**WIKIMEDIA**
F O U N D A T I O N

# Problem: **Wikipedia is missing images**

# Many Images!



Unillustrated Articles across Wikis

| Language | Percentage of unillustrated articles |
|---|---|
| Vietnamese | ~0.69 |
| Egyptian Arabic | ~0.61 |
| Arabic | ~0.53 |
| Swedish | ~0.51 |
| Korean | ~0.51 |
| English | ~0.47 |
| French | ~0.41 |

# "Add an Image" Structured Task

# "Add an Image" Algorithm Logic



Unillustrated Article

Wikimedia Commons

File:Hypatia.jpg

# The "Add an Image" Task in Wikipedia

- Manually finding good images for articles is hard:

  - The space of search for good free-licensed image matches can be huge, especially for newcomers

  - Sometimes finding matches is difficult: image search is not very effective in absence of metadata

- Automatically finding good images for articles is also hard:

  - We need to have "explainable" image-article matches

  - Due to the impact and visibility of images, accuracy/quality of matches is crucial

  - Computer vision technologies are not able to identify fine-grained objects, and simpler models are preferred

- Language support

  - Similar to links 300+ language version of Wikipedia

# A Simple Solution: **Leveraging Existing Matching Signals**

**Discovering unillustrated articles as articles without images or with *icons* only.**

Generate list of icons for each Wikipedia



Unillustrated Article

**Discovering related images from different sources**

Filtering

**Ranking images according to Relevance and Quality**

Task recommendation

## Borel set

In mathematics, a **Borel set** is any set in a topological space that can be formed from open sets (or, equivalently, from closed sets) through the operations of countable union, countable intersection, and relative complement. Borel sets are named after Émile Borel.

Unillustrated
article

**Dynamic unillustrated article detection based on Wiki size**

ويكيبيديا
الموسوعة الحرة

Has a
**Wikidata**
item?

NO

No recommendation

**Dynamic unillustrated article detection based on Wiki size**

WIKIPEDIA

Borel set

In mathematics, a **Borel set** is any set in a topological space that can be formed from open sets (or, equivalently, from closed sets) through the operations of countable union, countable intersection, and relative complement. Borel sets are named after Émile Borel.

Unillustrated article

ويكيبيديا
الموسوعة الحرة

Has a **Wikidata** item?

YES

**Discover image candidates:**
    **P18** (Image) of Wikidata item
    **P373** (Commons Category) of Wikidata  Item
    **Lead Images** of same article in other languages that are on Commons (i.e. can be used on all Wikis)

NO

**Filter out bad images:** placeholders, graphics, flags, maps **Remove articles about years, disambiguations, and lists**

**No recommendation**

NO

More than 0 candidates left?

**Dynamic unillustrated article detection based on Wiki size**

WIKIPEDIA

Borel set

In mathematics, a **Borel set** is any set in a topological space that can be formed from open sets (or, equivalently, from closed sets) through the operations of countable union, countable intersection, and relative complement. Borel sets are named after Émile Borel.

Unillustrated article

ويكيبيديا
الموسوعة الحرة

Has a **Wikidata** item?

YES

NO

**Discover image candidates:**
**P18** (Image) of Wikidata item
**P373** (Commons Category) of Wikidata Item
**Lead Images** of same article in other languages that are on Commons (i.e. can be used on all Wikis)

**Filter out bad images:** placeholders, graphics, flags, maps **Remove articles about years, disambiguations, and lists**

More than 0 candidates left?

NO

YES

**No recommendation**

**Rank images according to Reliability:**
- Prioritize sources according to their reliability
- Number of languages referring to the same image

**Assign image recommendations and confidence scores**

**Up to 3 recommendations**
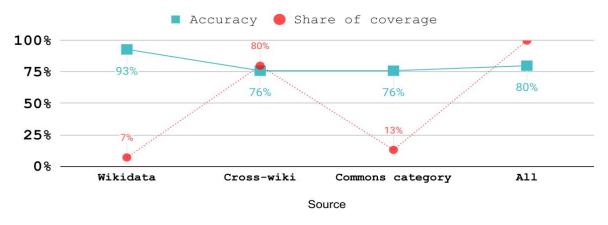
# Evaluation



## Android MVP "Train Image Algorithm"

Image Rec Task deployed on plwiki, eswiki, ruwiki, hewiki, ptwiki, fawiki, viwiki, frwiki, trwiki, itwiki, arwiki, enwiki, dewiki (many thanks to the Growth, Android and Platform Engineering Teams!!)
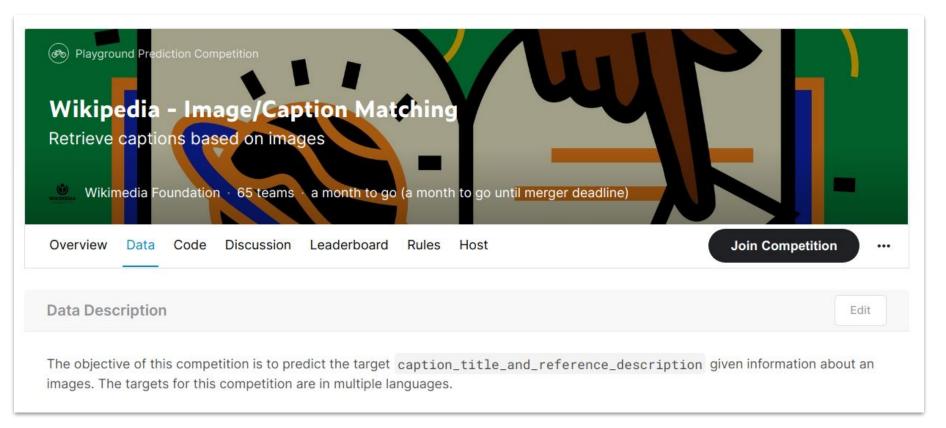
**Engagement is high**: average 9 edits/day

**Precision is high but depends on source**:



**Coverage**: 4% - 26% of unillustrated articles

# Gathering our Scientific Communities around the problem of image/text matching



Playground Prediction Competition

## Wikipedia - Image/Caption Matching
Retrieve captions based on images

Wikimedia Foundation · 65 teams · a month to go (a month to go until merger deadline)

Overview   Data   Code   Discussion   Leaderboard   Rules   Host

**Join Competition**   ...

### Data Description                                                                 Edit

The objective of this competition is to predict the target `caption_title_and_reference_description` given information about an images. The targets for this competition are in multiple languages.

# Equity in recommendations

Isaac

WIKIMEDIA
FOUNDATION

# Equity in Recommendations

- Recommender systems facilitate a growing number of Wikimedia edits:

  - Suggested Edits (300K edits Commons; 825K edits Wikidata)

  - Newcomer Tasks (190K edits across 70 Wikipedias)

  - Content Translation (1M articles created across all Wikipedias)

- How do they align with knowledge equity goals?

  - What languages are these tools available in?

    - Read more: https://w.wiki/4Hhg

  - What types of content receives the greatest benefits? should be prioritized?

    - Read more: https://w.wiki/$4P

# Newcomer Tasks -- Russian Wikipedia -- Topic Filters

Flow of geographic "bias" induced by algorithmic filters (Baseline -> Impressions) and editor behavior (Impressions -> Clicks -> Edits)

**Baseline**
54% of Russian Wikipedia articles are geographic

**Impressions**
Proportion of geographic impressions by country

**Clicks**
Proportion of geographic clicks by country
(5.7% of all impressions)

**Edits**
Proportion of geographic impressions by country
(56.5% of all clicks)



Geographic Articles

Russia

USA

Ukraine

France

Germany

United Kingdom

Japan

Other

# What challenges do you face in expanding section alignment to more Wikipedia languages?



Images of English and Vietnamese article for Modern Art showing how content can be aligned across Wikipedia articles through mapping wikilinks to their respective Wikidata IDs. More information: https://arxiv.org/abs/2103.00068

# How might we prioritize links to be added to Wikipedia articles?



Snapshot of English Wikipedia Modern Art article with links highlighted based on gender of the person they are about. Links to articles about men are often more prominent than links to articles about women or non-binary individuals. More information: Youtube:WikiConference North America 2021

# What special considerations do we make when considering the usage of recommendation or AI models for images?



Logo for VisibleWikiWomen, an annual campaign organized by Whose Knowledge? that focuses on adding images to articles about women.
CC-BY-SA 4.0. Whose Knowledge?

# Next steps ...

We need more and more research to Address Knowledge Gaps!