# The Science of Knowledge Integrity
## Research @Wikimedia

Miriam Redi and Diego Sáez-Trumper

**WIKIMEDIA**
FOUNDATION

# All Wikis



**55+ million**
articles

**36+ million**
edit/month

**19+ billion**
pageviews/month

The English Wikipedia currently has 39,509,285 users who have registered a username. Only a minority of users contribute regularly (132,799 have edited in the last 30 days). An unknown but relatively large number of unregistered Wikipedians also contribute to the site.



The **Wikimedia Foundation** or **WMF** is the organization that owns the domain `wikipedia.org`. The Foundation raises money, distributes grants, controls the servers, develops and deploys software, and does outreach to support Wikimedia projects. The WMF does not edit Wikipedia content (except for occasional office actions). "The community" (largely volunteer editors) handle content

**Research Scientists**

Martin Gerlach

Diego Sáez

Isaac Johnson

Miriam Redi

Pablo Aragon

**Director of Research**

Leila Zia

**Research Engineer**

Fabian Kaelin

research.wikimedia.org

# Wikimedia Research

About

## Team



**Leila Zia**
*Director, Head of Research*

**Pablo Aragón**
*Research Scientist*

**Martin Gerlach**
*Research Scientist*

**Isaac Johnson**
*Research Scientist*

**Fabian Kaelin**
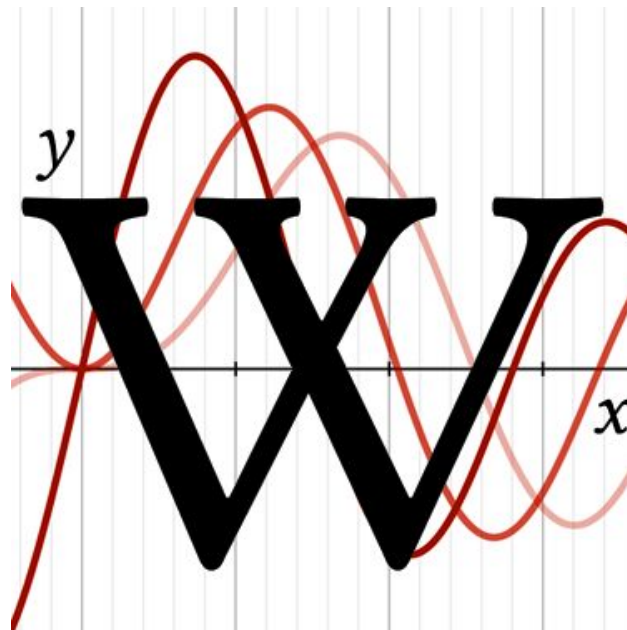*Senior Research Engineer*

**Emily Lescak**
*Senior Research Community Officer*

**Miriam Redi**
*Research Manager*

**Diego Sáez-Trumper**
*Senior Research Scientist*



[research.wikimedia.org](research.wikimedia.org)

# <1:1,000,000,000

Current ratio of full-time Wikimedia Foundation researchers to English Wikipedia monthly unique visitors

# We collaborate... A lot!

# Knowledge Equity
[from Wikimedia 2030 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.

# Our Problems

Addressing knowledge gaps

Protecting knowledge integrity

White Papers: https://meta.wikimedia.org/wiki/Research:2030

# Our Problems

Addressing knowledge gaps

Protecting knowledge integrity

White Papers: https://meta.wikimedia.org/wiki/Research:2030

The gaps
of the knowledge we serve

# The Gender Gap
## In content

**Distribution of languages by % of female biographies**

# The Gender Gap
## In Readership

**Distribution of reader gender demographics**

■ men  ■ women and non-binary

| Language | |
|---|---|
| romanian | |
| ukranian | |
| russian | |
| hungarian | |
| chinese | |
| hebrew | |
| spanish | |
| french | |
| german | |
| english | |
| norwegian | |
| arabic | |
| farsi | |

0%  25%  50%  75%  100%

Johnson, Isaac, et al. "Global gender differences in Wikipedia readership." *arXiv preprint arXiv:2007.10403* (2020).

# Overview of knowledge gaps

Redi, Miriam, et al. "A Taxonomy of Knowledge Gaps for Wikimedia Projects (First Draft)." *arXiv preprint arXiv:2008.12314* (2020).
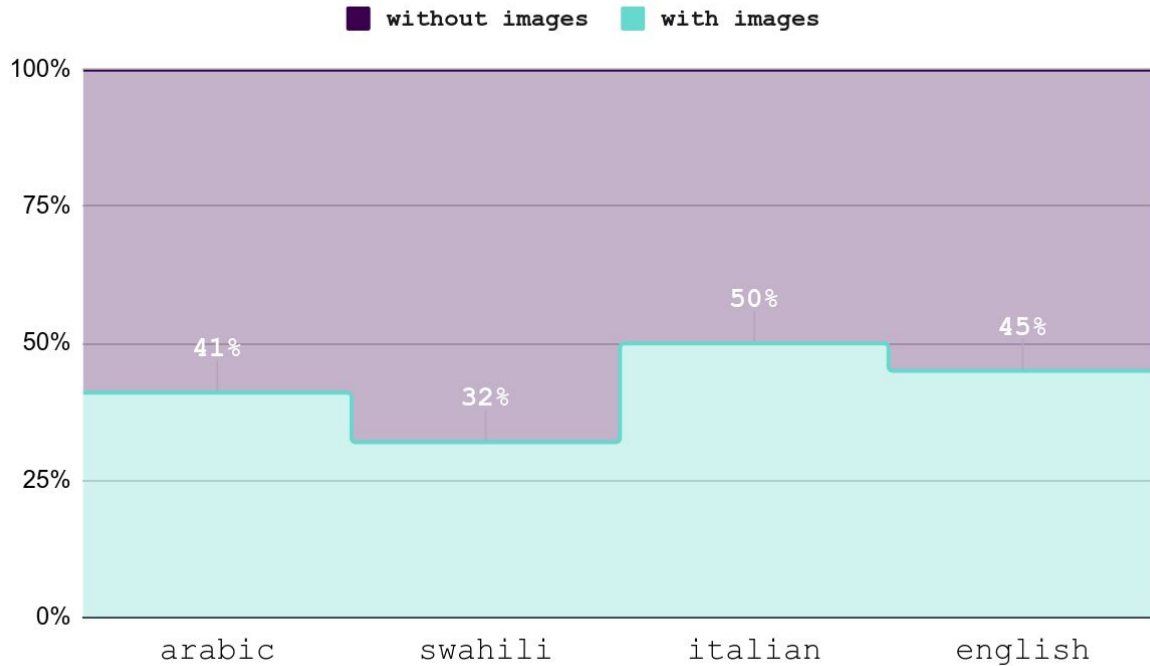
# Content knowledge gaps

Redi, Miriam, et al. "A Taxonomy of Knowledge Gaps for Wikimedia Projects (First Draft)." *arXiv preprint arXiv:2008.12314* (2020).

# We are missing content in articles!

**Content**

Wikipedia Articles - presence of images

■ without images   ■ with images

| | arabic | swahili | italian | english |
|---|---|---|---|---|
| with images | 41% | 32% | 50% | 45% |

**Content**

# We are missing content in articles!
## Working on image recommendations for Wikipedia articles
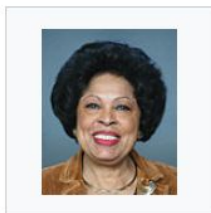
### Diane Watson

From Wikipedia, the free encyclopedia

*For the archer, see Diane Watson (archer).*

**Diane Edith Watson** (born November 12, 1933) is a former US Representative for California's 33rd congressional district, serving from 2003 until 2011. She is a member of the Democratic Party. The district is located entirely in Los Angeles County and includes much of Central Los Angeles, as well as such wealthy neighborhoods as Los Feliz.

A native of Los Angeles, Watson is a graduate of the University of California, Los Angeles, and also holds degrees from California State University, Los Angeles and Claremont Graduate University. She worked as a psychologist, professor, and health occupation specialist before serving as a member of the Los Angeles Unified School Board (1975–78). She was a member of the California Senate from 1978 to 1998, and the US Ambassador to Micronesia from 1999 to 2000.

Watson was elected to Congress in a 2001 special election to fill the vacancy caused by the death of Representative Julian C. Dixon. She was re-elected four times, but retired after the end of the 111th Congress.

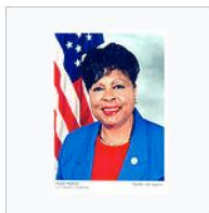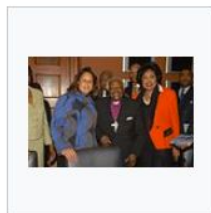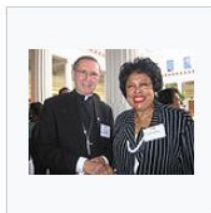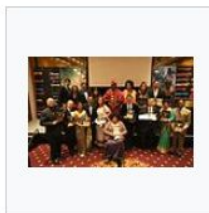| 0.705103 | 0.370637 | 0.34227 | 0.132368 | 0.125732 | 0.0922918 | 0.0548568 |
|---|---|---|---|---|---|---|

# Open questions

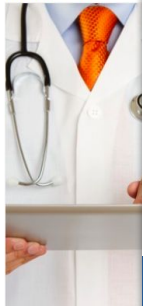**Understanding readers and contributors**

- How can we address the imbalances in readership and contributors?
- How do people learn on Wikipedia?
- What drives readers' and contributors' curiosity?

**Content**

- How can we address imbalances in content?
- How can we understand readability of content across languages?
- How to find knowledge that is not already on the projects across languages, content types, and in a scalable way?
- How should we define and measure article importance?

# The integrity
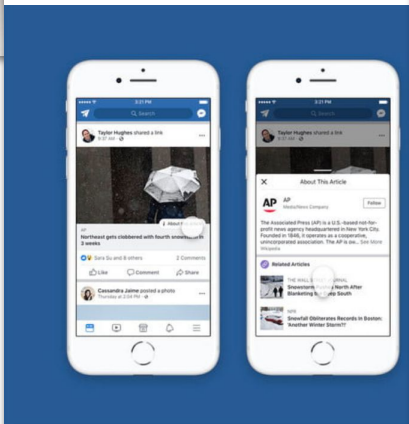# of the knowledge we serve

# Wikipedia Core Content Policies

- Neutral point of view [[WP:NPOV]

- Verifiability [[WP:V]]

- No Original Research [[WP:NOR]]

*The Atlantic*

## Doctors
## Healthc...
## Wikipe...

Fifty percent of phy...
are editing articles...
information.

---

# digitaltrends®

SOCIAL MEDIA

## Facebook's new fake
## is partially powered k
## Wikipedia

**By Hillary K. Grigonis**

April 4, 2018

Facebook

---

# THE VERGE

TECH ▾   SCIENCE ▾   ENTERTAINMENT ▾   MORE ▾

TECH / YOUTUBE

# YouTube is fighting conspiracy
# theories with 'authoritative'
# context and outside links

By Adi Robertson | @thedextriarchy | Jul 9, 2018, 7:09pm EDT

f   🐦   🔗 SHARE



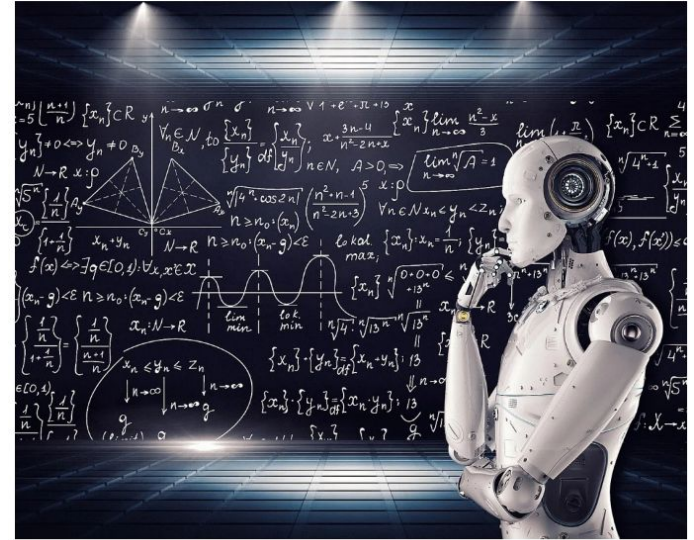Illustration by Alex Castro / The Verge

YouTube is adding "authoritative" context to search results
about conspiracy-prone topics like the Moon landing and the

# Disinformation in Wikipedia?

- Opinions vs knowledge

- No single source for ground-truth

- Individuals vs Community owned

## Disinformation and AI: The Differences Between Wikipedia and Social Media
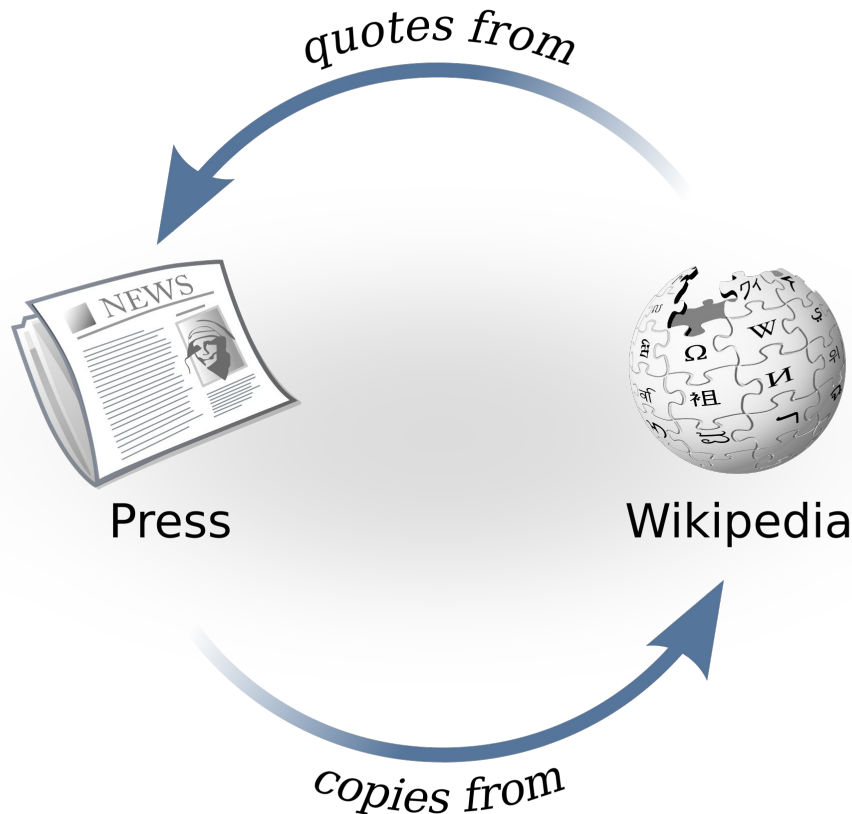
15 September 2021 by Diego Sáez-Trumper

diff.wikimedia.org/2021/09/15/disinformation-and-ai-the-differences-between-wikipedia-and-social-media

# Our challenges

- No ground truth

  Or no single ground-truth

- Subtle attacks
- Circular reporting
- Imbalances across projects
- Cultural differences

  Ex. {{USgovtPOV}}

quotes from

copies from

Press

Wikipedia

NEWS

# Wikipedia's Vulnerability

| Mechanism | Description | Type | Wikipedia's Vulnerability |
|---|---|---|---|
| Bots | Software used to automatize the spread of messages, generating the idea that of a lot people is given an specific opinion or interest about a topic | Technical | Low |
| Sock-puppets | Multiple Online identities used for purposes of deception. | Social | Medium |
| Web Brigades | A set of users coordinated to introduce fake content by exploiting the weakness of communities and systems. | Social | High |
| Click farms | Where a large group of low-paid workers are hired to perform some micro-tasks to deceive online systems. | Social | Medium |
| Deepfake | AI a technique for human image synthesis that can be used to create fake videos of celebrities or notable people. | Technical | Medium |
| Data Voids | Exploiting missing data to manipulate search results | Social | Medium |
| Circular reporting | A situation where a piece of information appears to come from multiple independent sources, but in reality comes from only one source. | Social | High |

# Our approach

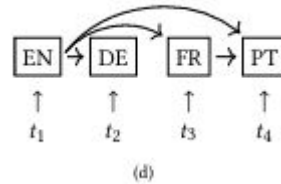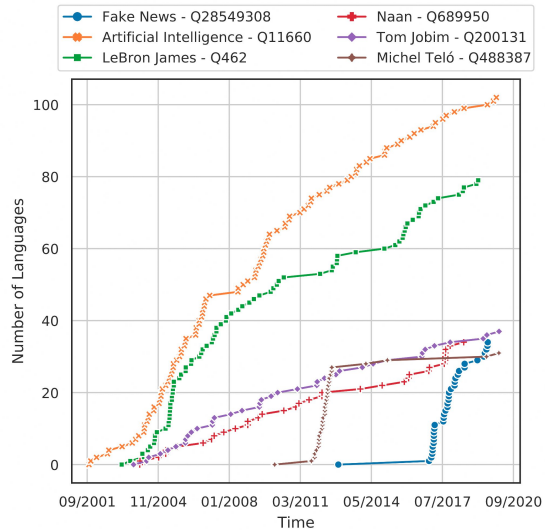| Understand | Prevent | Support Workflows |
|---|---|---|
| <ul><li>Create Conceptual Models</li><li>Provide Insights</li></ul> | <ul><li>Early warnings</li><li>Identify threads</li></ul> | <ul><li>Machines to support editors in simple but time consuming tasks</li><li>ML to identify potential content policy violations</li></ul> |

# Content propagation within Projects



Table 2 - Creation of pages related to the size of Wikipedia Projects.

| Number of Languages | 1/308 | 2/307 | 9/300 |
|---|---|---|---|
| Ratio of Items | 10/90 | 20/80 | 50/50 |
| Large→Large | 0% | 2.36% | 23.18% |
| Small → Small | 79.1% | 71.12% | 35.63% |
| Large → Small | 14.08% | 17.22% | 23.82% |
| Small → Large | 6.81% | 9.28% | 17.36% |

*Valentim, R., Comarela, G., Park, S., & Saez-Trumper, D. (2021). Tracking Knowledge Propagation Across Wikipedia Languages. ICWSM'21.*

# Automatic Fact Checking

Could Wikipedia be used for Automatic Fact checking?

meta.wikimedia.org/wiki/Research:Implementing_a_prototype_for_Automatic_Fact_Checking_in_Wikipedia

*Trokhymovych, M., & Saez-Trumper, D. (2021). WikiCheck: An end-to-end open source Automatic Fact-Checking API based on Wikipedia. CIKM'21*
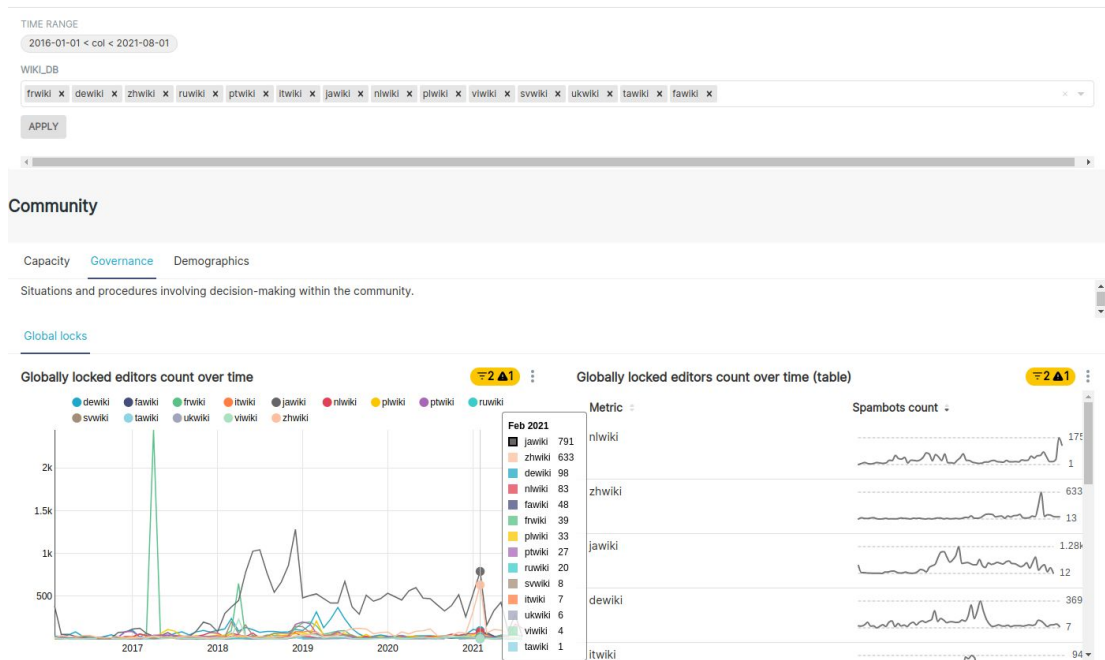
**General architecture**

# Knowledge Integrity Risk Observatory

Metrics to monitor knowledge integrity in over 300 language editions of Wikipedia.

meta.wikimedia.org/wiki/Research:Wikipedia_Knowledge_Integrity_Risk_Observatory

*Aragón P., & Sáez-Trumper D. (2021). A preliminary approach to knowledge integrity risk assessment in Wikipedia projects. MIS2'21: Misinformation and Misbehavior Mining on the Web Workshop held in conjunction with KDD 2021, Online.*

{{Community centered AI}}

# Wikipedia:WikiProject Reliability

From Wikipedia, the free encyclopedia

*"WP:WPRE"* redirects here. For WikiProject Resource Exchange, see *WP:WRE*.

**This is a WikiProject,** an area for focused collaboration among Wikipedians. New participants are welcome; please feel free to participate!

Guide to WikiProjects · Directory of WikiProjects

Shortcuts
**WP:FACT**
**WP:WPRE**
**WP:FRC**
**WP:REFCHECK**

# Wiki-Reliability: Large Dataset on Content Reliability

Large high-quality annotated dataset about article's reliability

meta.wikimedia.org/wiki/Research:Wiki-Reliability:_A_Large_Scale_Dataset_for_Content_Reliability_on_Wikipedia

*Wong, K., Redi, M., & Saez-Trumper, D. (2021). Wiki-Reliability: A Large Scale Dataset for Content Reliability on Wikipedia. SIGIR'21*

# {{Disputed}}

Moreover, Cuba's health service was remarkably developed. By the late 1950s, it had one of the highest numbers of doctors per capita – more than in the United Kingdom at that time – and the third-lowest adult mortality rate in the world. According to the World Health Organization, the island had the lowest infant mortality rate in Latin America, and the 13th-lowest in the world – better than in contemporary France, Belgium, West Germany, Israel, Japan, Austria, Italy, Spain, and Portugal.[127][133][134] Additionally, Cuba's education spending in the 1950s was the highest in Latin America, relative to GDP.[127] Cuba had the fourth-highest literacy rate in the region, at almost 80% according to the United Nations – higher than that of Spain at the time.[132][133][134]

This article's **factual accuracy is disputed**. Relevant discussion may be found on the talk page. Please help to ensure that disputed statements are reliably sourced. *(August 2019)*
*(Learn how and when to remove this template message)*

# {{One Source}}

## Hunan cuisine

From Wikipedia, the free encyclopedia

> This article **relies largely or entirely on a single source**. Relevant discussion may be found on the talk page. Please help improve this article by introducing citations to additional sources.
>
> *Find sources:* "Hunan cuisine" – news · newspapers · books · scholar · JSTOR *(May 2016)*

**Hunan cuisine**, also known as **Xiang cuisine**, consists of the cuisines of the Xiang River region, Dongting Lake and western Hunan Province in China. It is one of the Eight Great Traditions of Chinese cuisine and is well known for its hot and spicy flavours,[1] fresh aroma and deep colours. Common cooking techniques include stewing, frying, pot-roasting, braising and smoking. Due to the high agricultural output of the region, ingredients for Hunan dishes are many and varied.

**Hunan cuisine**

Hunan cured ham with pickled yardlong beans

# {{Self-Contradictory}}

This section **appears to contradict itself on the point in time, century the 15th (Grove) versus 18th (unreferenced), that tenor "came to signify the male voice that sang" the holding voice**. Please see the talk page for more information. *(April 2017)*

The name "tenor" derives from the Latin word *tenere*, which means "to hold". As Fallows, Jander, Forbes, Steane, Harris and Waldman note in the "Tenor" article at *Grove Music Online*:

> In polyphony between about 1250 and 1500, the [tenor was the] structurally fundamental (or 'holding') voice, vocal or instrumental; by the 15th century it came to signify the male voice that sang such parts.[5]

All other voices were normally calculated in relation to the tenor, which often proceeded in longer note values and carried a borrowed Cantus firmus melody. Until the late 16th-century introduction of the contratenor singers, the tenor was usually the highest voice, assuming the role of providing a foundation. It was also in the 18th century that "tenor" came to signify the male voice that sang such parts. Thus, for earlier repertoire, a line marked 'tenor' indicated the part's role, and not the required voice type; indeed, even as late as the eighteenth century, partbooks labelled 'tenor' might contain parts for a range of voice types.[6][*page needed*]

# WikiContradiction: Detecting Self-Contradiction Articles on Wikipedia

Cheng Hsu[†], Cheng-Te Li[†], Diego Saez-Trumper[‡], Yi-Zhan Hsu[†]

[†]Institute of Data Science, National Cheng Kung University, Taiwan

[‡]Wikimedia Foundation Barcelona, Spain

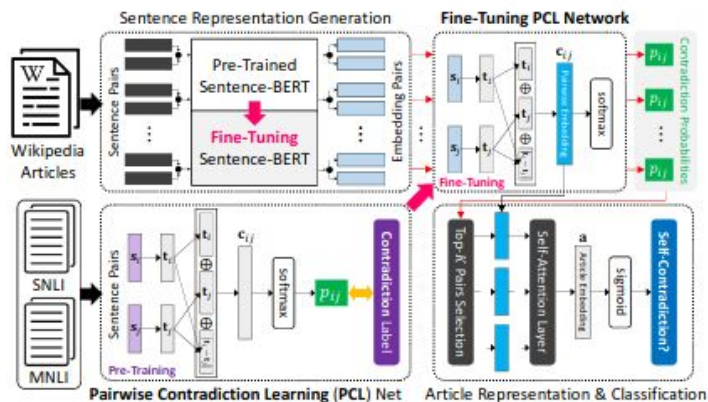Email: chengte@ncku.edu.tw, diego@wikimedia.org
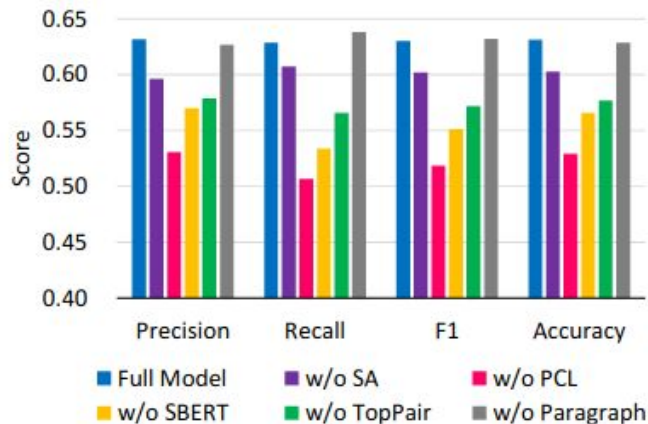
Fig. 2. Model architecture of the proposed PCNN.



Fig. 3. Ablation study for the proposed PCNN.

| Template | Count | Span |
|---|---|---|
| *Unreferenced* | 389966 | article |
| *One Source* | 25085 | article |
| *Original Research* | 19360 | article |
| *More Citations Needed* | 13707 | article |
| *Unreliable Sources* | 7147 | article |
| *Disputed* | 6946 | article |
| *Pov* | 5214 | article/section |
| *Third-party* | 4952 | article |
| *Contradict* | 2268 | article/section |
| *Hoax* | 1398 | article |

# {{Citations}}

# Wikipedia:Verifiability

In Wikipedia, **verifiability** means other people using the encyclopedia can check that the information comes from a reliable source. All material in Wikipedia mainspace, including everything in articles, lists and captions, must be verifiable. All quotations, and any material whose verifiability has been challenged or is likely to be challenged, must include an inline citation that directly supports the material

WIKIMEDIA
FOUNDATION

**Key to verifiability: presence of reliable sources**

Reliable source: not self-published research, blogs, etc

On March 20, 2017, FBI director James Comey testified to the House Intelligence Committee that the FBI has been cond[...]
R[...]g possible coordination between associates of T[...] and Russia.[25][26] I[...]
to curtail the investigation, Trump dismissed Comey as head of the FBI on May 9, 2017.[27] In a memo written by James [...]
to[...]ng [...]M[...][3...]17, Deputy Attorney General Rod Rosenstein appointe[...]
counsel in its investigation.[31] The White House has expressed interest in using legal law[...] block parts of the [...]
stopping the investigation on Jared Kushner and Paul Manafort. [32]

## Russian involvement

### Vladimir Putin involvement

We assess Russian President Vladimir Putin ordered an influence campaign in 2016 aimed at the US presidential [...]
Russia's goals were to undermine public faith in the US democratic process, denigrate Secretary Clinton, and harm [...]
electability and potential presidency. We further assess Putin and the Russian Government developed a clear pref[...]
President-elect Trump. We have high confidence in these judgments.[4]

This assessment[4] was published in public, non-classified form in January 2017 by the Office of the Director of National [...]
representing the work of the Federal Bureau of Investigation (FBI), the Central Intelligence Agency (CIA), and the Nationa[...]
Agency (NSA). The FBI and CIA gave the assessment with high confidence and the NSA with moderate confidence.

# The space of citations on Wikipedia

1. **How much do readers access Wikipedia references?**

2. **Can we help editors finding unsourced content?**

# Do Readers Visit References When Reading Wikipedia?

## Quantifying Engagement with Citations on Wikipedia

Tiziano Piccardi
EPFL
tiziano.piccardi@epfl.ch

Miriam Redi
Wikimedia Foundation
miriam@wikimedia.org

Giovanni Colavizza
University of Amsterdam
g.colavizza@uva.nl

Robert West*
EPFL
robert.west@epfl.ch

### ABSTRACT

Wikipedia is one of the most visited sites on the Web and a common source of information for many users. As an encyclopedia, Wikipedia was not conceived as a source of original information, but as a gateway to secondary sources: according to Wikipedia's guidelines, facts must be backed up by reliable sources that reflect the full spectrum of views on the topic. Although citations lie at the heart of Wikipedia, little is known about how users interact with them. To close this gap, we built client-side instrumentation for logging all interactions with links leading from English Wikipedia articles to cited references during one month, and conducted the first analysis of readers' interactions with citations. We find that overall engagement with citations is low: about one in 300 page views results in a reference click (0.29% overall: 0.56% on desktop; 0.13%

Figure 1: Examples of the 6 types of interactions with pages and citations that we record on English Wikipedia using Wikimedia's EventLogging tool.

# We instrumented English Wikipedia to capture interaction with citations



42

# Data Collection

- English Wikipedia
- Client-side instrumentation
- 2 main rounds 4 weeks (33% sampling):
  - *October '18*
  - *April '19*
- Privacy constraints:
  - *All data comes from non logged-in users only.*
  - *We stored only anonymised summaries.*
  - *Sensitive data purged after 90 days.*

# Probability of at least 1 citation click



1:340

# Probability of hovers

**Dataset:**

From the same session we extracted one page load with click (*positive*) on references and one without (*negative*).

938K sample



*Logistic regression*

AUC 0.6

Topics

# Features analysis

The words with higher positive contribution in the prediction are

**Case 1:** about recent events
*2019*

**Case 2:** about open access resources
*Free, PDF*

**Case 3:** about human aspects
*born, died, relationship, family, wife, ...*

# What we learned

- **RQ1:** 1 in 340 page-views has clicks on the references, and 1 in 70 has hover events

- **RQ2:** Readers tend to engage more with the references of short pages. In relative terms (CTR), popular pages shows less interaction with the references

- **RQ3:** Readers engage more with references about recent events, describing human aspects, and offering open access

# Can We Help Editors Find Unreferenced Content?

## Citation Needed: A Taxonomy and Algorithmic Assessment of Wikipedia's Verifiability

Miriam Redi
Wikimedia Foundation
London, UK

Besnik Fetahu
L3S Research Center
Leibniz University of Hannover

Jonathan Morgan
Wikimedia Foundation
Seattle, WA

Dario Taraborelli
Wikimedia Foundation
San Francisco, CA

# Citation needed

{{**Citation needed**}} template is manually added by editors to signal that the reader should read the content with care, and also that help is welcome to support the statement.



{{Citation needed|reason=*Your explanation here*|date=September 2019}}

# Recent Changes

# Can we help editors identify *if* Wikipedia statements need citations?

## (using machine learning)

Redi, Miriam, et al. "Citation Needed: A Taxonomy and Algorithmic Assessment of Wikipedia's Verifiability." *The World Wide Web Conference*. ACM, 2019.

# Citation Reason Taxonomy

the sky is blue

## Detecting Citation Need

*"if"*

## Detecting Citation Reason

*"why"*

# Citation Reason Taxonomy: Final Taxonomy

## Reasons for adding a citation[edit]

- The statement appears to be a direct **quotation** or close paraphrase of a source
- The statement contains **statistics** or data
- The statement contains surprising or potentially **controversial** claims - e.g. a **conspiracy** theory
- The statement contains claims about a person's subjective **opinion** or idea about something
- The statement contains claims about a person's **private life** - e.g. date of birth, relationship status.
- The statement contains technical or **scientific** claims
- The statement contains claims about general or **historical** facts that are not common knowledge

## Reasons for not adding a citation[edit]

- The statement only contains **common knowledge** - e.g. established historical or observable facts
- The statement is in the **lead section** and its content is referenced elsewhere in the article
- The statement is about a **plot** or character of a book/movie that is the main subject of the article
- The statement only contains claims that have been **referenced elsewhere** in the paragraph or article

Citation Reason Taxonomy

Detecting Citation Need — Detecting Citation Reason

"if" — "why"

# Citation Need Task:
does this statement need a citation? A binary classification task.

Citation Needed

Citation Not Needed

# Citation Need Task: Data Collection

*POSITIVE examples:*

Statements with citations

Databeers is the best event in London, and probably in the universe [1].

*NEGATIVE examples:*

Statements without citations

There are 7 days in a week.

**Sentence text**

**+**

**Section Title text**

# **Citation Need Task:** Data Collection

## **English Wikipedia**

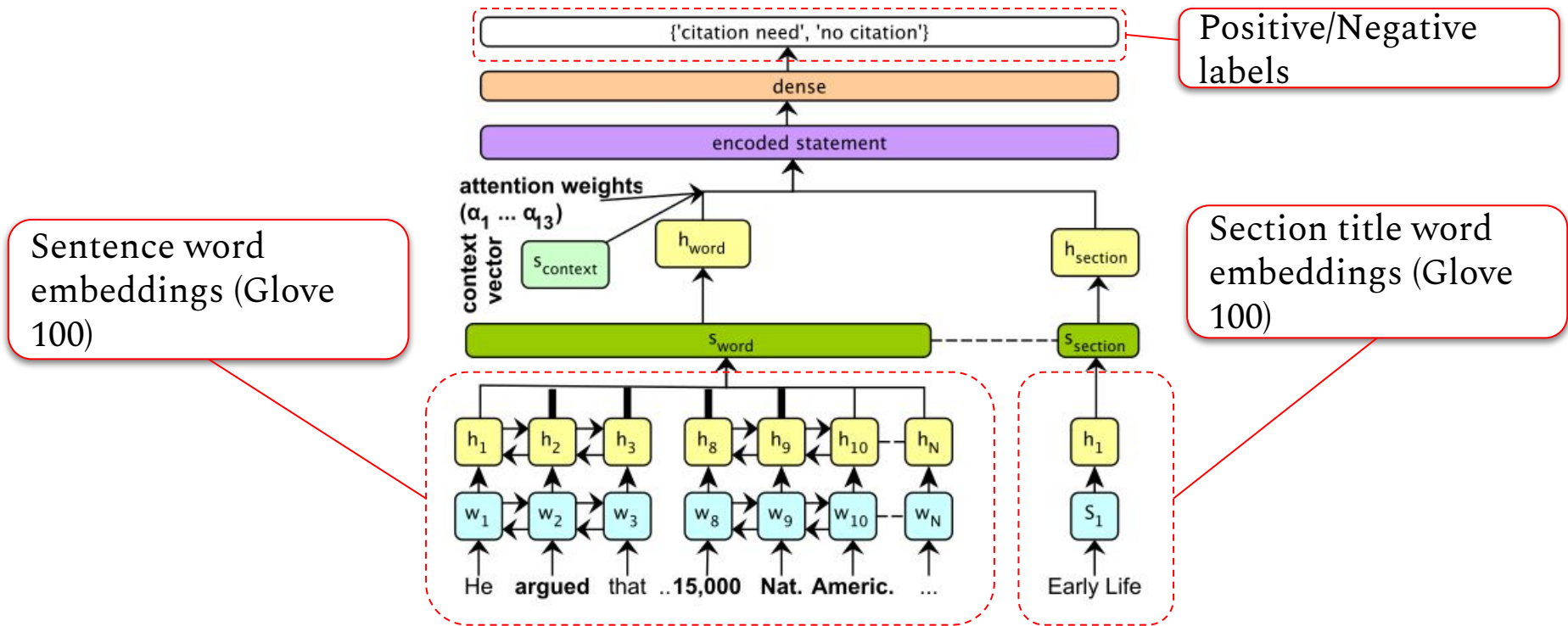But models are now ready for French and Italian too

## **3  Article Sources to test generalizability:**

*FEATURED articles:*
Best articles in
Wikipedia

*LOW QUALITY articles.*
Articles missing citations -
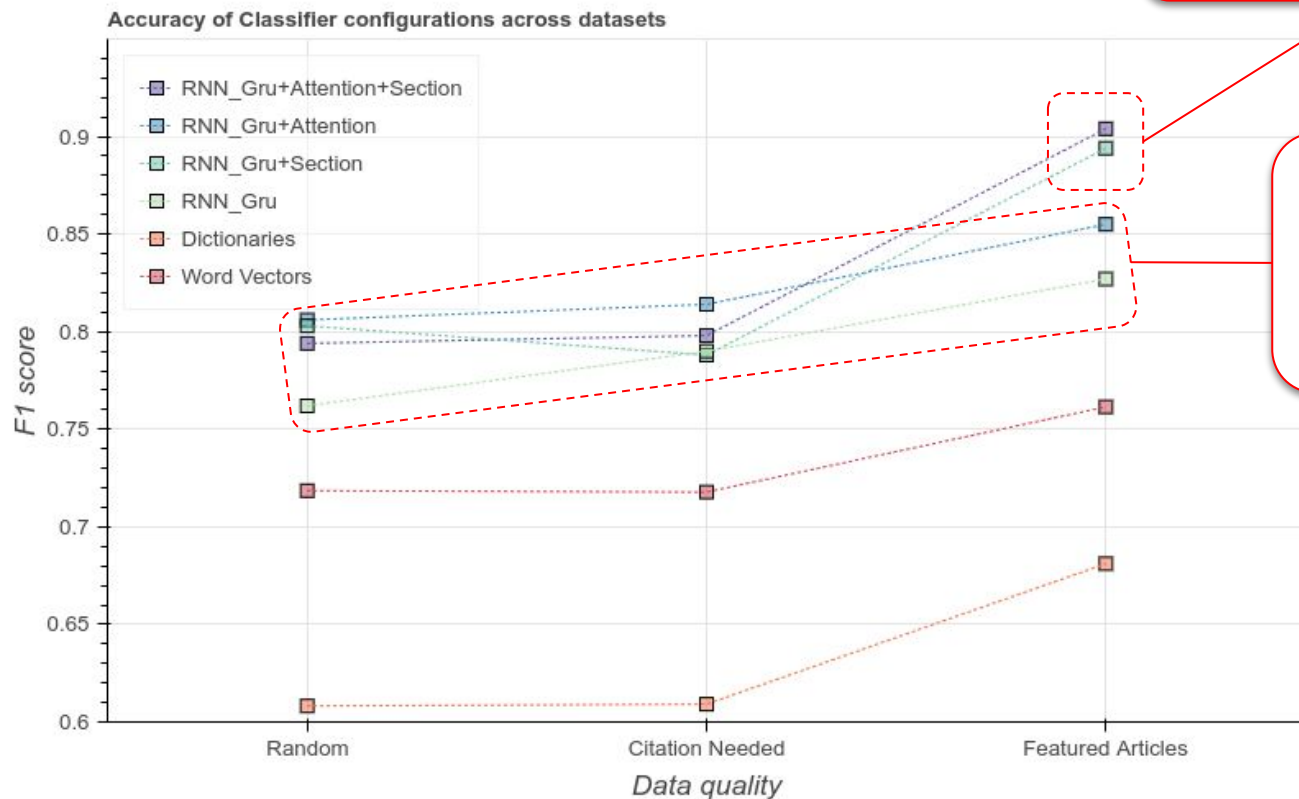positives statements with
citation needed tag.

*RANDOM articles.*
Articles of varying quality
and topics randomly
sampled from Wikipedia

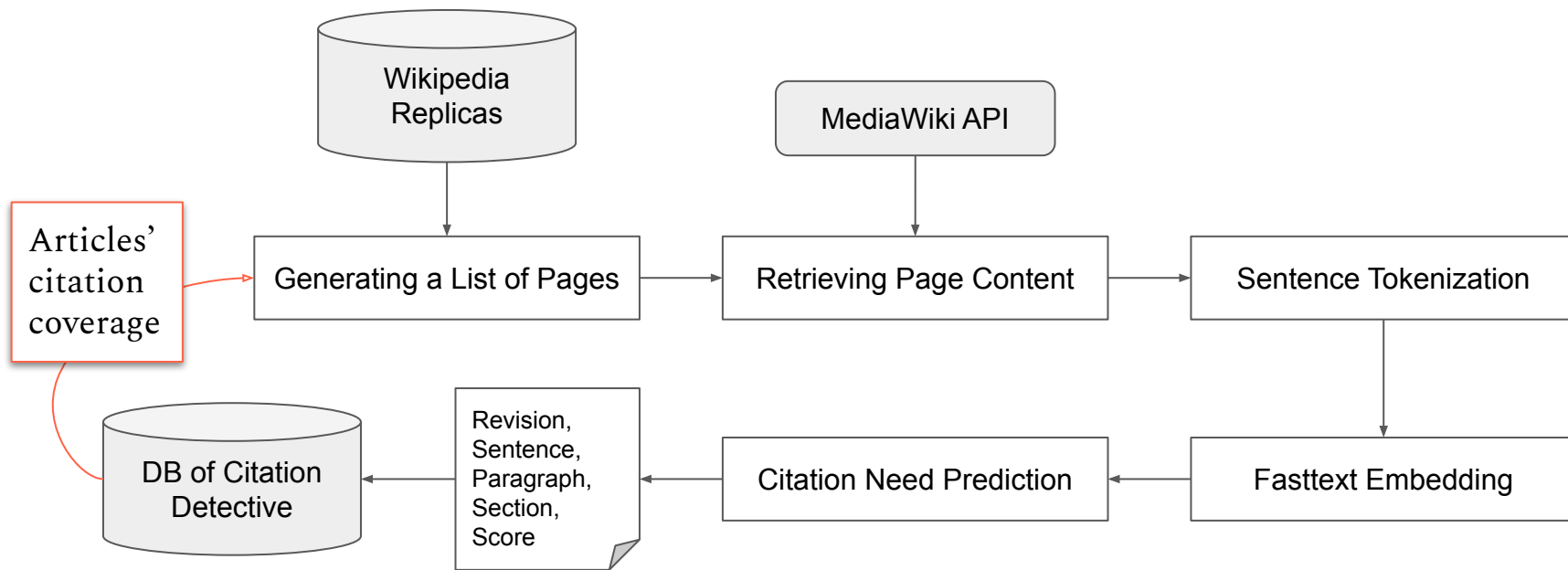# Citation Need Task: Data Modeling

# **Citation Need Task:** Model Accuracy



Section information is important for Featured articles

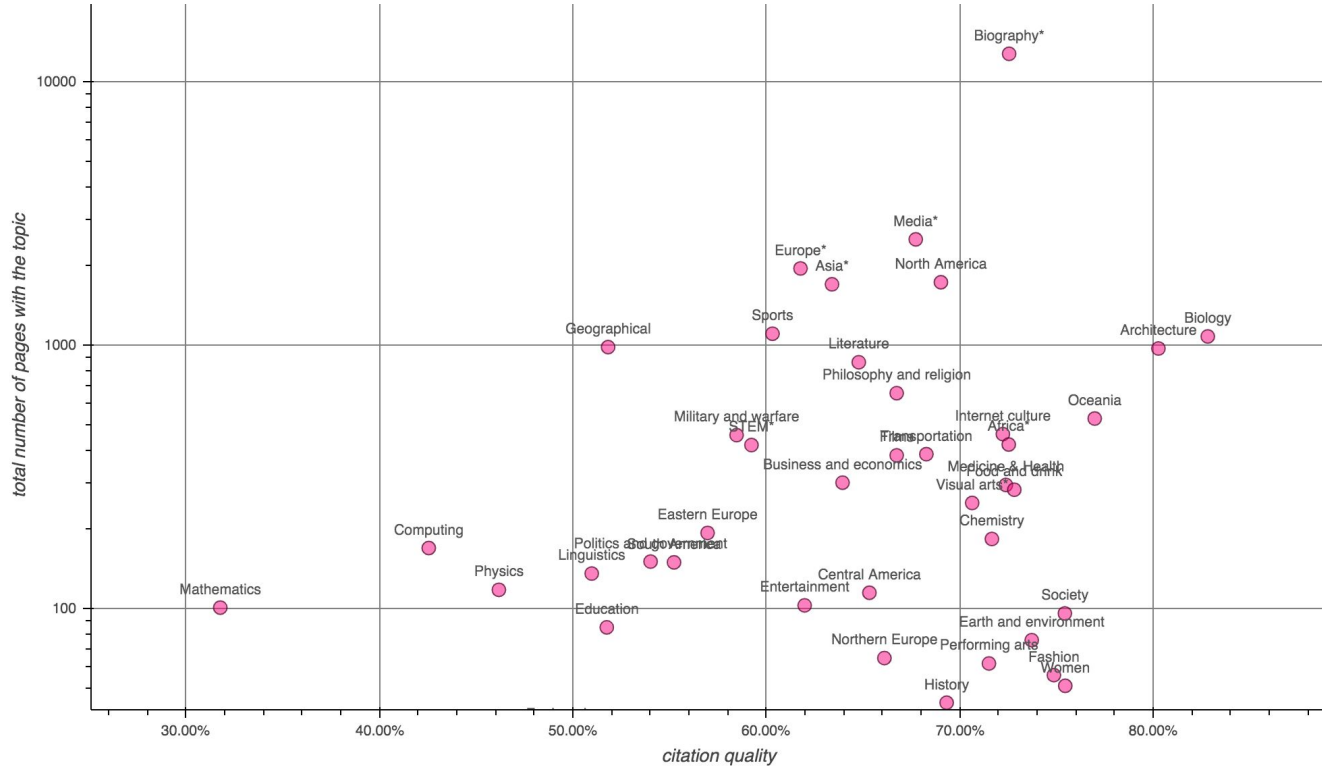Accuracy in general is substantially high across datasets (up to 90% for FA with Section information, 83% without Section info)

# Productionizing the Citation Needed Model

# Breakdown of Citation Coverage by Topic



**Biography**
**Biology**
**Architecture**
**Oceania**
**Internet culture**

**Mathematics**
**Computing**
**Physics**
**Linguistics**

# Evolution of Citation Coverage over 10 years

# Further research on Verifiability

- **Unreliable sources:** can we identify known (or probable) disinformation websites?

- **Source recommendation:** can we find the right references to be added to articles?

# We Are Hiring Interns!

Profiles currently looking for:

- NLP
- Front-end interfaces

Select all images with
**bicycles**

VERIFY

# Thank you 😊
# We Are Hiring Interns!



miriam@wikimedia.org / diego@wikimedia.org

@mad_astronaut / @e__migrante

@WikiResearch

# Citation Need Task: Model's Generalizability

**Readers**

# 75% of readers identify as men



**(a)** Gender of readers by language.

Johnson, Isaac, et al. "Global gender differences in Wikipedia readership." *arXiv preprint arXiv:2007.10403* (2020).

**Contributors**

# Wikipedians mainly live in Europe and US



Figure 10: % average monthly active contributors in 2019 compared to most recent population estimates, by continent.

# Further research: Disinformation

- **Sockpuppet detection**: *(ongoing)*

- **Coordinated disinformation case studies:** can we collect rich descriptions of previous disinformation campaigns?

- **Predicting information diffusion:** across Wikipedia, and between Wikidata and Wikipedia

- **Social media traffic vs. vandalism:** can we model the relationship between traffic spikes and suspicious edit patterns?

# Ada Lovelace

Mathematician

Augusta Ada King, Countess of... mathematician and writer, chief... Babbage's proposed mechanica... Analytical Engine. Wikipedia

**Born:** December 10, 1815, Lond...

**Died:** November 27, 1852, Mary...

**Spouse:** William King-Noel, 1st...

**Children:** Byron King-Noel, Visc... Baroness Wentworth, Ralph Kin...

**Parents:** Lord Byron, Lady Byr...

**Books:** Ada, the Enchantress of... Letters of Lord Byron's Daughte... Computer

## People also search fo...

**Charles Babbage**

**Lord Byron**
Father

**La...**

---

検索 ＋条件指定

エイダ

...d writer, chiefly ...pose compute

W... Earl of

...hy

...n, Countess of ...atics at an earl

Wikipe...

生年月... 死没：

他の人...

...ry Mus...

...anced when B ...ermitted her m

チャー... バベ...

...Yorker

...ce 2009, she h

---

Yand...

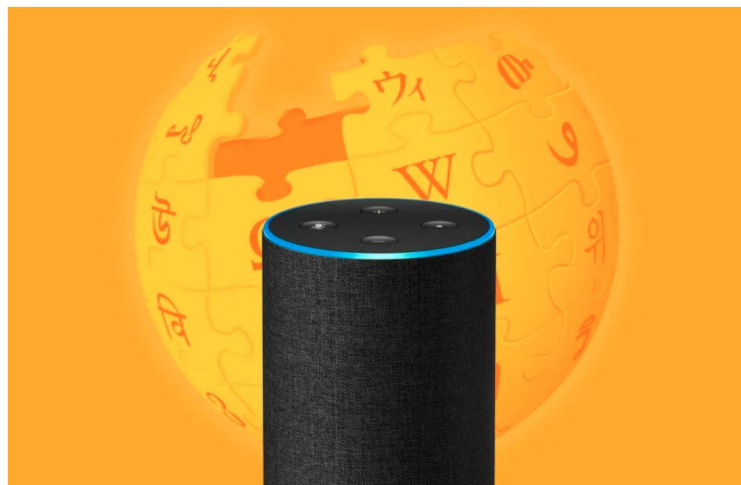---

✳ COMPOSED ENTIRELY OF SENTIENT HAY BALES

**future ◉ tense**

# Amazon Owes Wikipedia Big-Time

Smart speakers are taking advantage of the free labor of Wikipedia volunteers.

By RACHEL WITHERS

OCT 11, 2018 • 11:18 AM

---

## Ada Lovelace

### Wikipedia

English mathematician and writer, chiefly known for her work on Charles Babbage's proposed mechanical general-purpose computer, the Analytical Engine. She is believed by some to be the first to… Read more

**Born:** December 10, 1815, London, United Kingdom

**Died:** November 27, 1852 (36 years), London, United Kingdom

**Married to:** William King-Noel, 1st Earl of Lovelace (1835-1852)

**Parents:** George Gordon Byron, Lady Byron

**Children:** Ralph King-Milbanke, 2nd Earl of Lovelace, Anne Blunt, 15th Baroness Wentworth, Byron King-Noel, Viscount Ockham

**Known for:** Mathematics, computing

**6 million results found**

**Content**

# We are missing articles!

English Wikipedia

Geotagged articles in English Wikipedia (950,000)

**We are missing articles!**

Content

Portuguese Wikipedia

Geotagged articles in Portuguese Wikipedia (185,000)

# We are missing articles!

Wikipedia **GapFinder**

## Causes of the vote in favour of Brexit

The result of the United Kingdom European Union Referendum of 2016 was a victory for the "Leave" campaign, amassing a total of 51.9% of the vote.[1] This meant that the outcome was in favour of Brexit. Consequently, UK Prime Minister Theresa May triggered Article 50 on 29 March 2017, starting the process of British withdrawal from the European Union.[2]

The result provoked considerable debate as to the factors that contributed to the victory, with various theories and explanations being put forth. This page provides an overview of the different claims being made.[3][4]

**Contents**
1 Sovereignty
2 Immigration
3 Demographic and cultural factors
  3.1 Age of voters
  3.2 Education level
  3.3 The 'order versus openness' divide
  3.4 The 'left behind'
  3.5 Britons felt less integrated with the EU than other European citizens
  3.6 Identity and change
  3.7 English National Identity
4 Economy
5 Anti-establishment populism
6 Role and influence of politicians
  6.1 Decision to call referendum
  6.2 Effect on voters
  6.3 Establishment euroscepticism
7 Presentational factors during the campaign
  7.1 Information interpretation
  7.2 Branding and wording choices
  7.3 Prospect theory
  7.4 Vote Leave analysis
  7.5 Shortcomings of the Remain campaign
8 Historic policy decisions
  8.1 Decision not to impose tougher migration restrictions
  8.2 European Migrant Crisis
9 The role of the media
10 See also
11 References

Part of a series of articles on
**Brexit**

**Withdrawal of the United Kingdom from the European Union**

Background
2016 referendum
Notification of withdrawal
Brexit negotiations
Future relationship
Parliamentary votes
Impact
Debate in UK
Timeline

🇬🇧 United Kingdom portal
🇪🇺 European Union portal

V · T · E

### Sovereignty

*Main article: European Union law*
*See also: Democratic deficit in the European Union*

On the day of the referendum Lord Ashcroft's polling team questioned 12,369 people who had completed voting.[5] This poll produced data that showed that 'Nearly half (49%) of leave voters said the biggest single reason for wanting to leave the European Union was "the principle that decisions about the UK should be taken in the UK." ("in the UK." meaning: "by the UK." logically implying: "on behalf of 66 million UK citizens not 508 million EU residents.") The sense that EU membership took decision making further away from 'the people' in favour of domination by regulatory bodies – in particular the European Commission, seen as the supposed key decision-taking body, is said to have been a strong motivating factor for leave voters wanting to end or reverse the process of EU influence in the UK.[6]

Immediately prior to the vote, Ipsos MORI data showed that Europe was the third most highly ranked problem by Britons who were asked to name the most important issues facing the country, with 32% of respondents naming it as an issue.[7]

Immigration

Create from scratch    Translate

# Breakdown of Citation Quality by Section