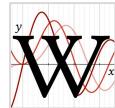


# Open, collaborative, reproducible research to support free knowledge

*I Congresso Científico Brasileiro da Wikipédia* • October 14, 2016

Dario Taraborelli  
@readermeter



# Outline

Research at the Wikimedia Foundation

Wikimedia Research: current priorities

A platform for open, collaborative, reproducible research

# Outline

Research at the Wikimedia Foundation

Wikimedia Research: current priorities

A platform for open, collaborative, reproducible research

# Outline

Research at the Wikimedia Foundation

Wikimedia Research: current priorities

A platform for open, collaborative, reproducible research

# Research at the Wikimedia Foundation



**WIKIMEDIA**  
FOUNDATION

# A brief history of research at WMF

## Early years

- 2010 Research committee (RCom), the Research Index
- 2011 First research roles
- 2011 Wikimedia Summer of Research
- 2012 Research microteams at the Foundation
- 2013 Research and Analytics
- 2015 Research department

# A brief history of research at WMF

## Early years

- 2010      Research committee (RCom), the Research Index
- 2011      First research roles
- 2011      Wikimedia Summer of Research
- 2012      Research microteams at the Foundation
- 2013      Research and Analytics
- 2015      Research department

# A brief history of research at WMF

## Early years

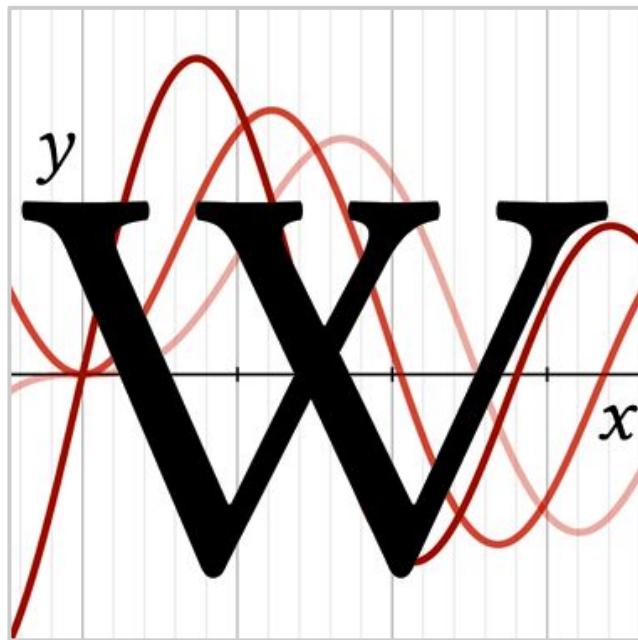
- 2010 Research committee (RCom), the Research Index
- 2011 First research roles
- 2011 Wikimedia Summer of Research
- 2012 Research microteams at the Foundation
- 2013 Research and Analytics
- 2015 Research department

# A brief history of research at WMF

## Early years

- 2010 Research committee (RCom), the Research Index
- 2011 First research roles
- 2011 Wikimedia Summer of Research
- 2012 Research microteams at the Foundation
- 2013 Research and Analytics
- 2015 Research department

# Wikimedia Research



## Research



**Dario Taraborelli**  
Director, Head of Research



**Nathaniel Schaaf**  
Software Engineer  
(International)

## Research and Data



**Aaron Halfaker**  
Principal Research Scientist



**Morten Warncke-Wang**  
Research Fellow



**Robert West**  
Research Fellow



**Ellery Wulczyn**  
Data Scientist



**Erik Zachte**  
Data Analyst (International contractor)



**Leila Zia**  
Senior Research Scientist

## Design Research



**Abigail Ripstra**  
Lead Design Research Manager



**Samantha Becker**  
Participant Recruiter  
(Contractor)



**Jonathan Morgan**  
Senior Design Researcher

# altmetrics

## OPINION

### Measuring the Impact of Altmetrics

When it comes to ranking academic influence, PageRank is just the start

By PAUL MCFEDRIES / AUGUST 2012



**TRIAL BY TWITTER**

Blogs and forums are bypassing peer-reviewed journals as the go-to spot for academics to share their research.

**S**ince the early days of the Internet, academics have been using blogs and forums to share their research. But until recently, they were doing so largely outside the academic system. Now, however, that's changing, as more and more scholars are turning to social media to share their work. This shift is having a significant impact on how research is perceived and used. In this article, we'll explore the reasons behind this change and what it means for the future of academic communication.

**Illustration:** Jesse Lefkowitz

## ReaderMeter

About | FAQ | News

DUNCAN J WATTS

H<sub>R</sub>-Index: 16  
G<sub>R</sub>-Index: 26

Most read publication: 145

Total number of publications: 57

Total bookmarks: 767

[what does this mean?]

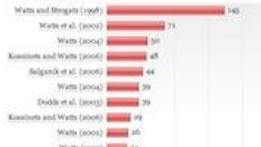
Permalinks

HTML: [http://readermeter.org/Watts.Duncan\\_J](http://readermeter.org/Watts.Duncan_J)  
JSON: [http://readermeter.org/Watts.Duncan\\_J.json](http://readermeter.org/Watts.Duncan_J.json)

Powered by  MENDLEY

ReaderMeter is a free service based on data from Mendeley (Terms of Use). Our data is not governed. Use of ReaderMeter's contents is free under a CC license.

### Top 10 publications by readership (?)



### Duncan J Watts's coauthors

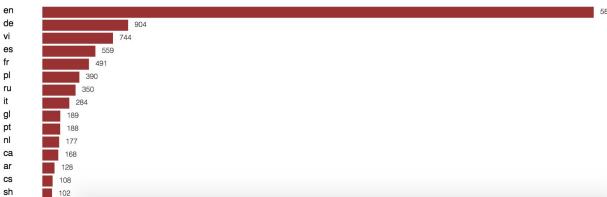
## WikiPedia Cite-o-Meter

### Statistics for Public Library of Science (10.1371)

[stats] [data] [json] [about]

Total number of citations across Wikipedia: 12110  
Project with the largest number of citations: en (6910 citations)  
Total number of citations in Wikimedia Commons: 16055  
Data last updated: 2015-12-20 03:35:44

### Citations for Public Library of Science (10.1371) in the top 100 Wikipedias



The altmetrics manifesto  
<http://altmetrics.org/manifesto/>

# Wikimedia Research: Current priorities

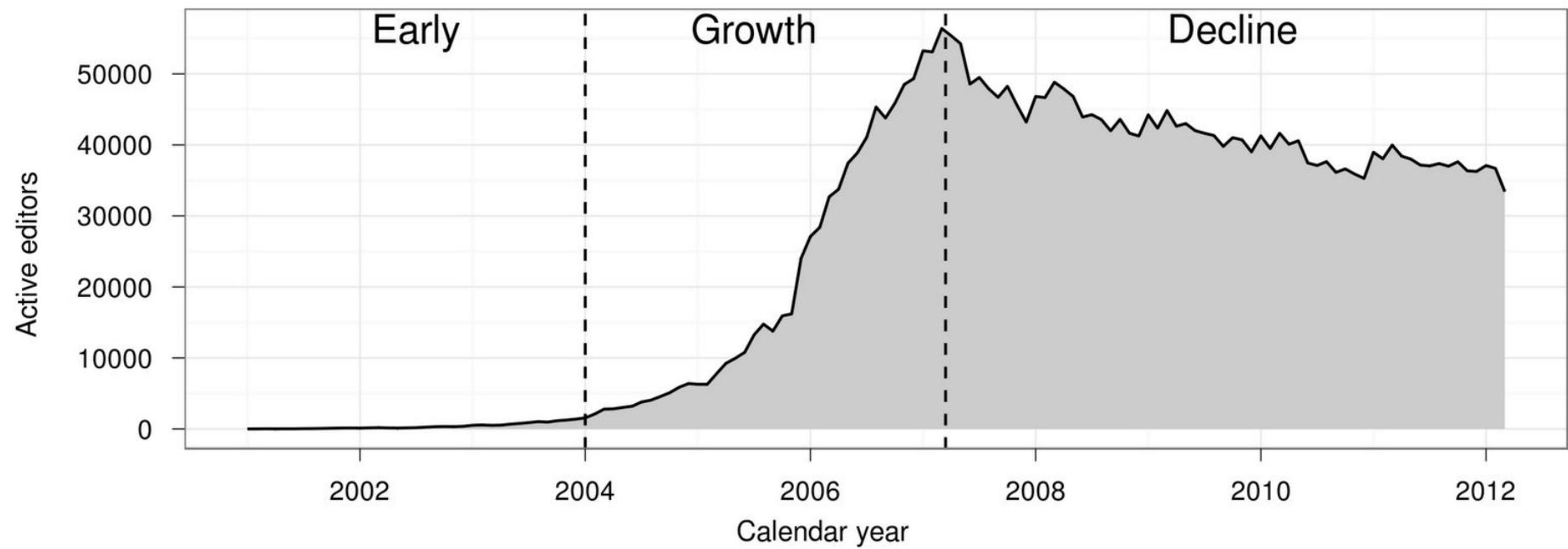


FOUNDATION

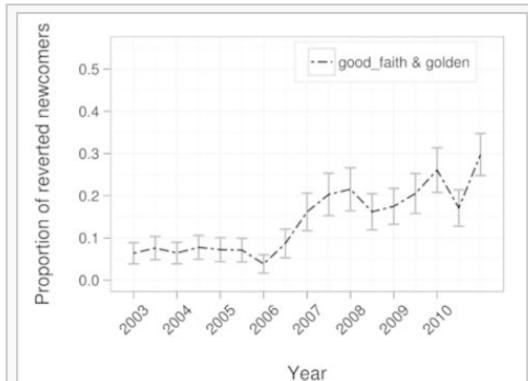
# Artificial intelligence as a service



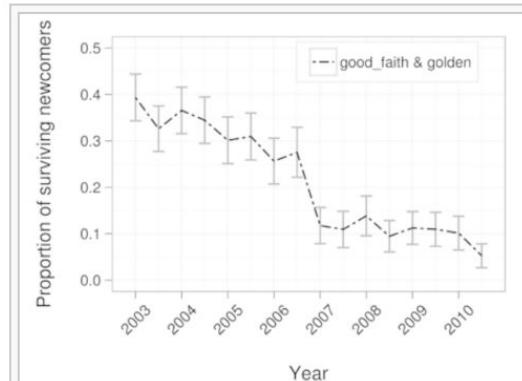
# 1. *Socially aware* quality control



first-time editors are *as good as they were* in 2004-05, but their *retention* sharply declined



**Rejection of desirable newcomers.** The proportion of desirable newcomers who are reverted in their first session of editing is plotted over time. □



**Survival of desirable newcomers.** □  
The proportion of desirable newcomers who edit for at least two months is plotted over time.

Wikipedia is a firehose

Bad edits must be reverted

Cost of quality control work must be minimized

*What about new contributors?*

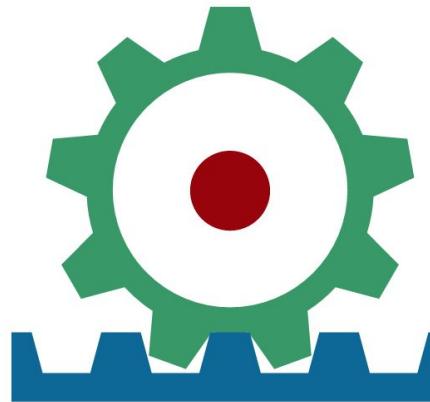
*How to build socially-aware quality control systems?*

# Correcting the side effects of quality control

The Objective Revision Evaluation Service (ORES) is a web service that provides machine learning as a service for Wikimedia Projects. The system is designed to help automate critical wiki-work -- for example, vandalism detection and removal. This service is developed as part of the [Revision scoring as a service](#) research project.

## Scores API

ORES is intended to be used as a source of information by tool developers. To access ORES, scores, a simple RESTful API is provided. There are two versions of the scoring API that differ slightly in their behavior. Version 2 provides access to model info in a scoring request. Version 1 is preserved for backwards compatibility.



<http://ores.wikimedia.org>

<https://blog.wikimedia.org/2015/11/30/artificial-intelligence-x-ray-specs/>

# The ORES API

Revision as of 21:52, 12 January 2015 (edit) (undo)

Blank123456789 (talk | contribs)

Next edit →

Line 11:

==Evaluation of intelligence==

"Intelligence" is hard to define, whether in dogs, other animals, or humans. The ability to learn quickly might be taken as a sign of [[intelligence (trait)|intelligence]], but such evidence must be interpreted with care, because learning speed may be affected by such things as the effectiveness of the rewards used in training or the motivation or activity level of the dog. For example, some breeds, such as [[Siberian Husky|Siberian Huskies]], are said to be not particularly rewarded by pleasing their owners, but quickly learn to escape from yards or catch small animals, often using ingenious ways of doing both. **LLAMAS GROW ON TREES**

```
{  
  "642215410": {  
    "prediction": true,  
    "probability": {  
      "false": 0.07561753062984156,  
      "true": 0.9243824693701584  
    }  
  }  
}
```

<https://ores.wikimedia.org/scores/enwiki/damaging/642215410>

<https://en.wikipedia.org/w/index.php?diff=prev&oldid=642215410>

# The ORES API



<http://ores.wikimedia.org>

<https://meta.wikimedia.org/wiki/ORES>

Image credits: Mun May Tee-Galloway

context	edit quality			article quality
	damaging	goodfaith	reverted	wp10
arwiki Arabic Wikipedia			✓	
cswiki Czech Wikipedia			✓	
dewiki German Wikipedia			✓	
enwiki English Wikipedia	✓	✓	✓	✓
enwiktionary English Wiktionary			✓	
eswiki Spanish Wikipedia			✓	
etwiki Estonian Wikipedia			✓	
fawiki Persian Wikipedia	✓	✓	✓	
frwiki French Wikipedia			✓	✓
hewiki Hebrew Wikipedia			✓	
huwiki Hungarian Wikipedia			✓	
idwiki Indonesian Wikipedia			✓	
itwiki Italian Wikipedia			✓	
nlwiki Dutch Wikipedia	✓	✓	✓	
nowiki Norwegian Wikipedia			✓	
plwiki Polish Wikipedia	✓	✓	✓	
ptwiki Portuguese Wikipedia	✓	✓	✓	
ruwiki Russian Wikipedia	✓	✓	✓	✓
svwiki Swedish Wikipedia			✓	
trwiki Turkish Wikipedia	✓	✓	✓	
ukwiki Ukrainian Wikipedia			✓	
viwiki Vietnamese Wikipedia			✓	
wikidatawiki Wikidata	✓	✓	✓	

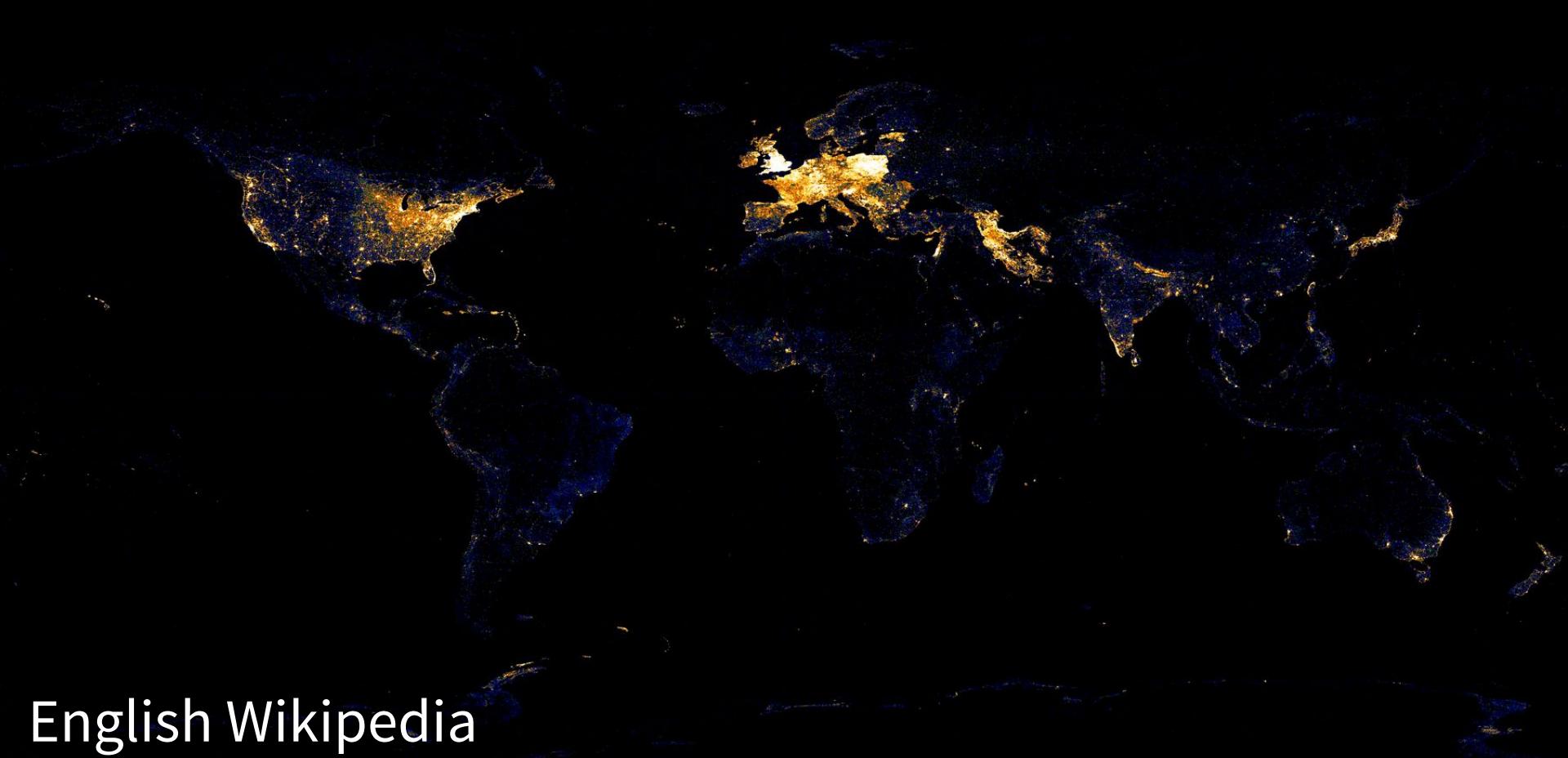
## 2. Content growth via recommendations



# How complete is Wikipedia?

Geotagged articles across all Wikipedia languages (2 million)

<https://iccl.inf.tu-dresden.de/web/Wikidata/Maps-06-2015/en>



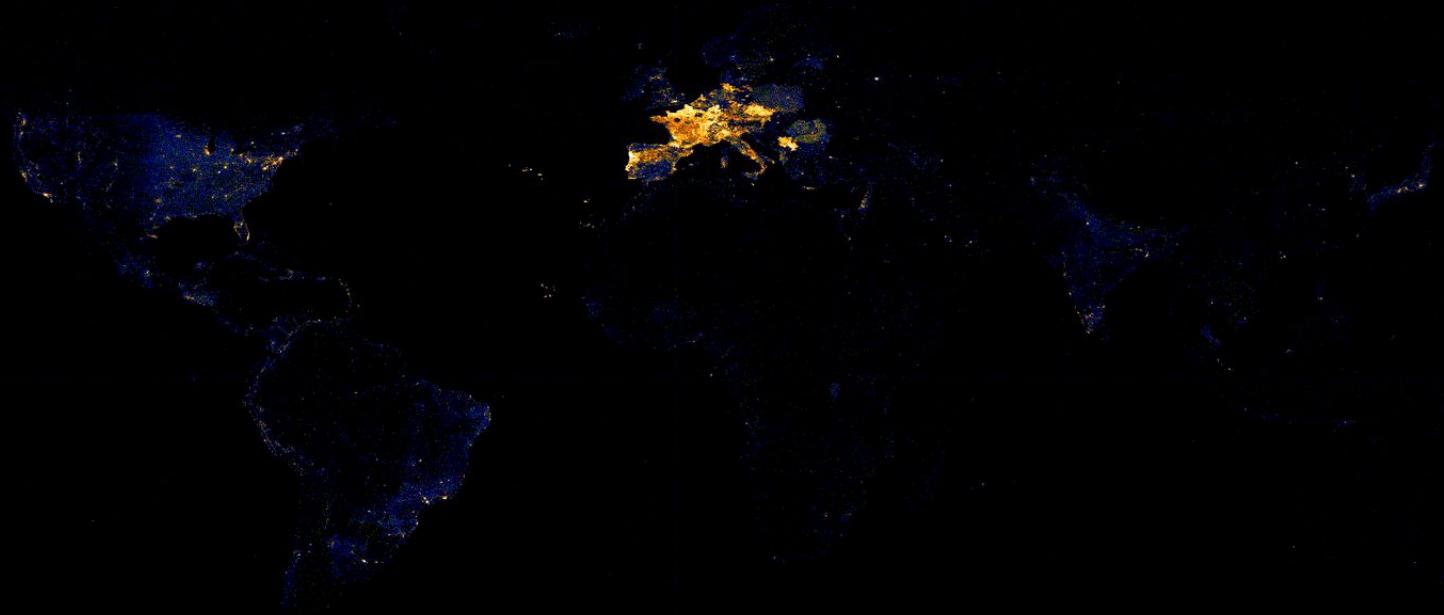
# English Wikipedia

Geotagged articles in English Wikipedia (950,000)



# Spanish Wikipedia

Geotagged articles in Spanish Wikipedia (261,000)



# Portuguese Wikipedia

Geotagged articles in Portuguese Wikipedia (185,000)



# Arabic Wikipedia

Geotagged articles in Arabic Wikipedia (87,000)

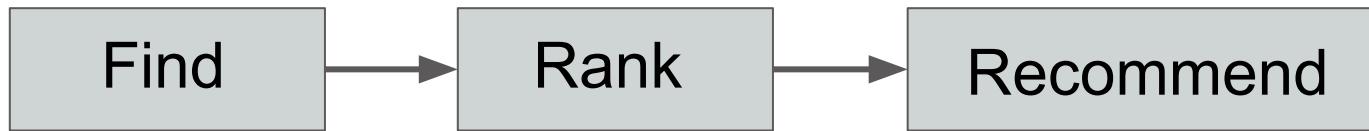
Depending on the language you speak, you can only learn about a small set of topics in Wikipedia.

*How to fill these content gaps?*

<http://blog.wikimedia.org/2016/04/27/article-recommendation-system/>

# Article recommendations

**Goal:** Recommend important missing articles to editors



In collaboration with:

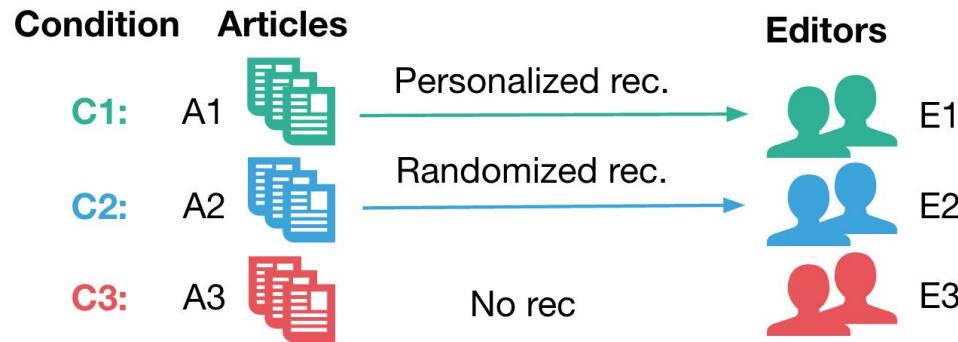


Stanford  
University

<http://blog.wikimedia.org/2016/04/27/article-recommendation-system/>

# Article recommendations

**Experimental results:** article creation recommendations can triple article creation rate, while controlling for quality



português ▾

العربية ▾

 Rio de Janeiro

Polícia Militar do Estado...

2k visualizações recentes



Volta Redonda

5k visualizações recentes

Resende (Rio de Janei...  
município localizado no sul do  
estado do Rio de Janeiro, no  
Brasil

4k visualizações recentes



Geografia do Rio de Ja...

2k visualizações recentes



Petrópolis

6k visualizações recentes



Campeonato Carioca d...

6k visualizações recentes



Carioca

2k visualizações recentes



Guanabara

3k visualizações recentes

# Understanding editing culture



# 1. Toxic comments and personal attacks

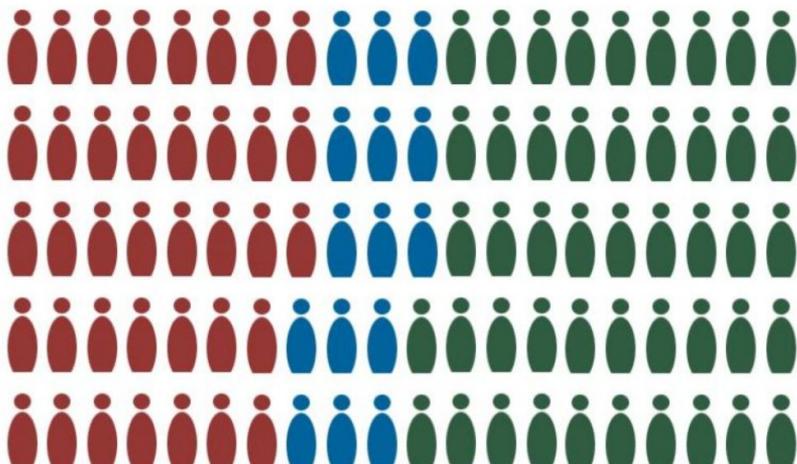
If quality control has caused new good-faith contributors to leave at increasingly higher rates

*what about toxic comments and personal attacks?*

<https://meta.wikimedia.org/wiki/Research:Detox>

# Harassment is common in Wikipedia

Respondents were asked if they had **personally experienced harassment**. Out of 2,495 that responded to this question :

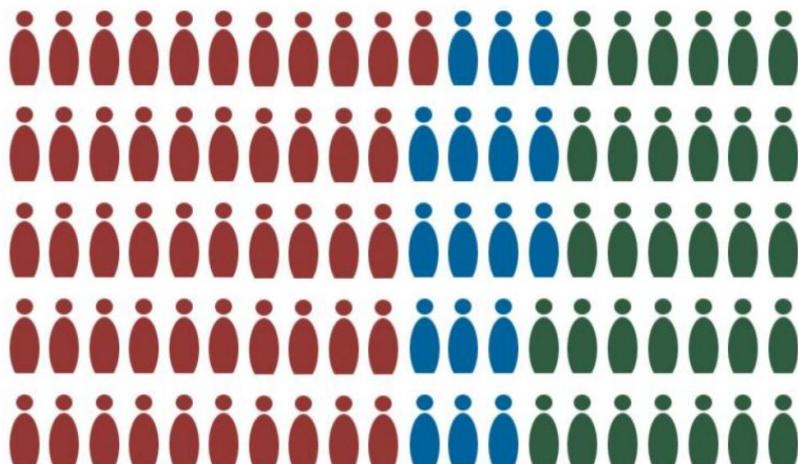


38% said yes

16% were unsure

47% said no

Respondents were asked if they had **witnessed the harassment of others**. Out of 2,078 that responded to this question:

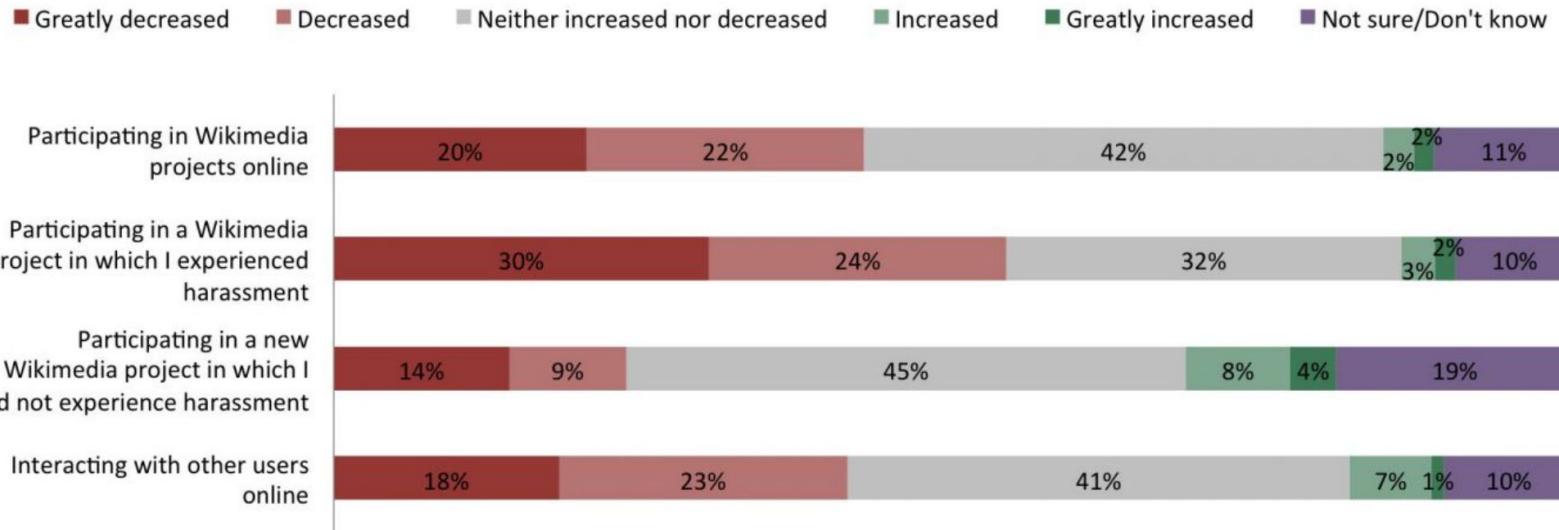


51% said yes

17% were unsure

32% said no

# Victims of harassment are less likely to contribute



# Research:Detox

## Goals:

1. Design algorithms to detect personal attacks on Wikipedia
2. Use these algorithms to analyze and identify harassing behavior
3. Analyze the impact of harassment on editor retention

In collaboration with:



Jigsaw

# Wikipedia DeTox

This page contains experimental machine learning classifiers for detecting toxicity in Wikipedia discussions. See [Wikimedia Detox Research](#) for more information.

Select a model:

- Attack
- Aggression

Select Input Type:

- Text
- Revision ID

Your point of view is pretty stupid and irrelevant.

Score

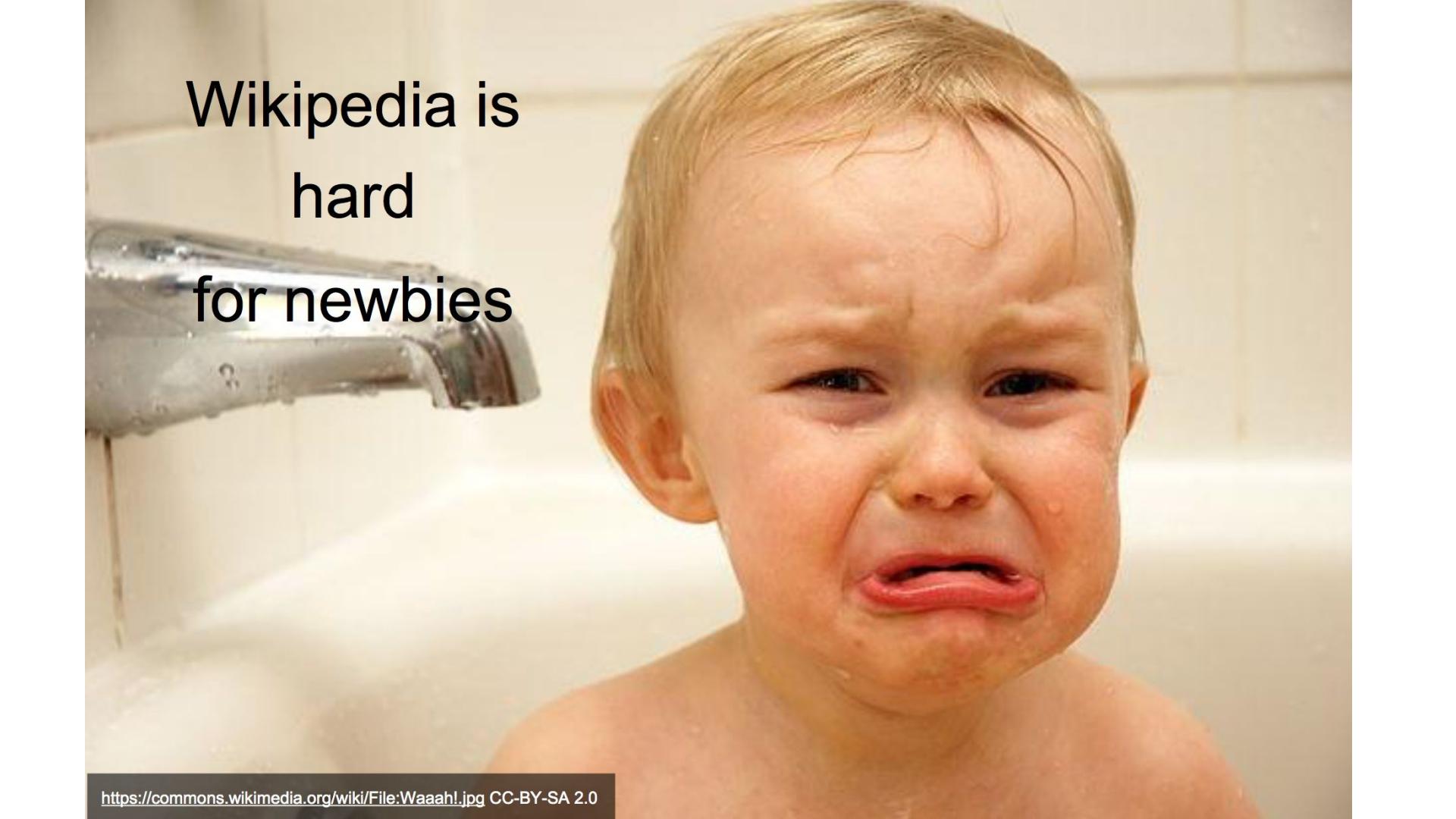
Results:

aggressive: 0.95

neutral: 0.05

friendly: 0.00

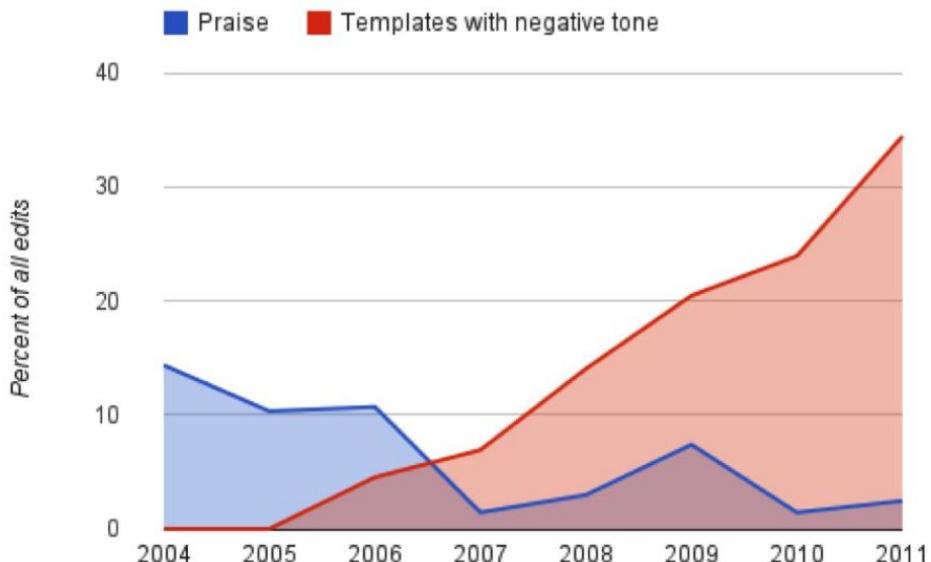
## 2. Socializing newcomers

A close-up photograph of a young child with light blonde hair, crying in a bathtub. Water is dripping from a chrome faucet on the left. The child's face is wet and expressive. The background shows the white tiled wall of the bathtub.

Wikipedia is  
hard  
for newbies

# Wikipedia is hard for new editors.

*Does a more welcoming experience increase new editor retention?*



# Teahouse: newcomer socialization experiments

Short term 1+ edit survival: +13%

Long term 5+ edit survival: +16%



Hi **RyanHarmy**! Thanks for contributing to Wikipedia. Be our guest at [the Teahouse](#)! The Teahouse is a friendly space where new editors can ask questions about contributing to Wikipedia and get help from peers and experienced editors. I hope to see you there! Doctree (I'm a Teahouse host)

[Visit the Teahouse](#)

This message was delivered automatically by your robot friend, [HostBot \(talk\)](#) 17:22, 20 October 2015 (UTC)

# Understanding Wikimedia readers



# 1. A taxonomy of Wikipedia readers

# Characterizing readers

*Daily pageviews:* 250 million (EN), 13 million (PT)

*Monthly unique devices:* 602 million (EN), 45 million (PT)

*Active editors:* 30,000 (EN), 1,500 (PT)

([analytics.wikimedia.org](https://analytics.wikimedia.org) [stats.wikimedia.org](https://stats.wikimedia.org), September 2016)

*Readers are a very large part of Wikipedia's ecosystem that are often overlooked*

# Characterizing readers

## Webrequest logs

- IP address, user agent, page requested, etc.

- A wealth of data

- We can't learn "why" readers read Wikipedia from webrequest logs alone

## Surveys

- Surveys help understand participant motivation

- Can't be meaningfully compared with quantitative data

*Combine survey and webrequest data*

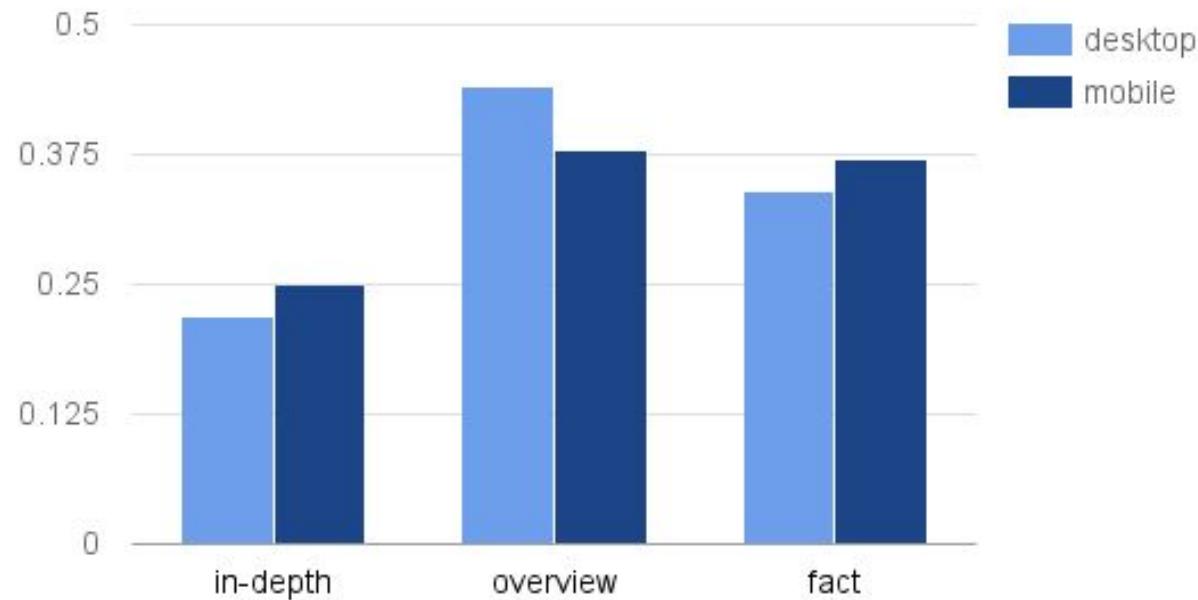
# A large scale reader survey

**Goal:** Collect data on three dimensions of reader experience, link responses to webrequest log data via opt-in tokens

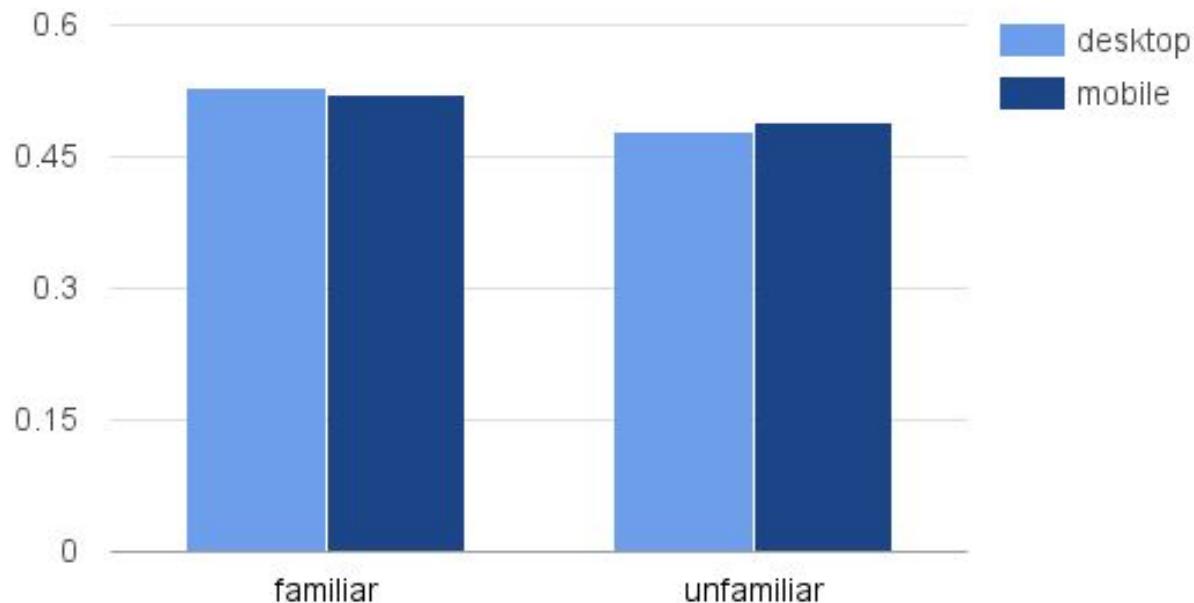
**Questions:** Depth of knowledge / Familiarity / Motivation

**Data:** 36K responses collected from English Wikipedia readers  
(Desktop+Mobile)

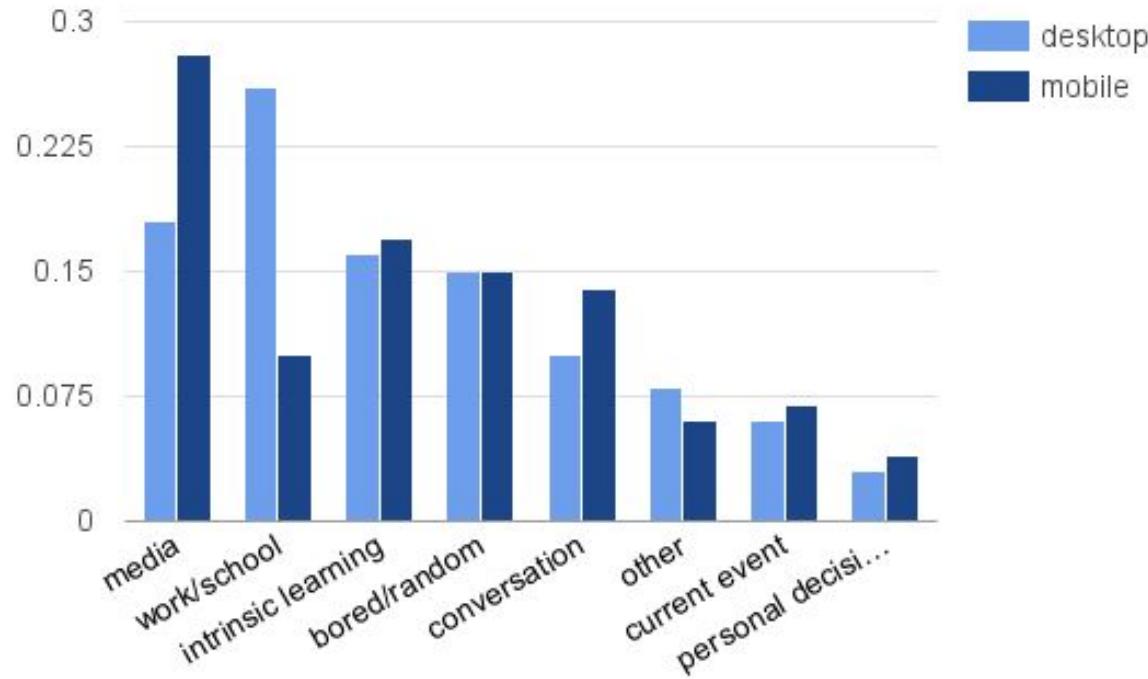
# Information Depth



# Prior knowledge



# Motivation



## 2. New readers

# New reader research

## Countries



**Mexico**



**India**



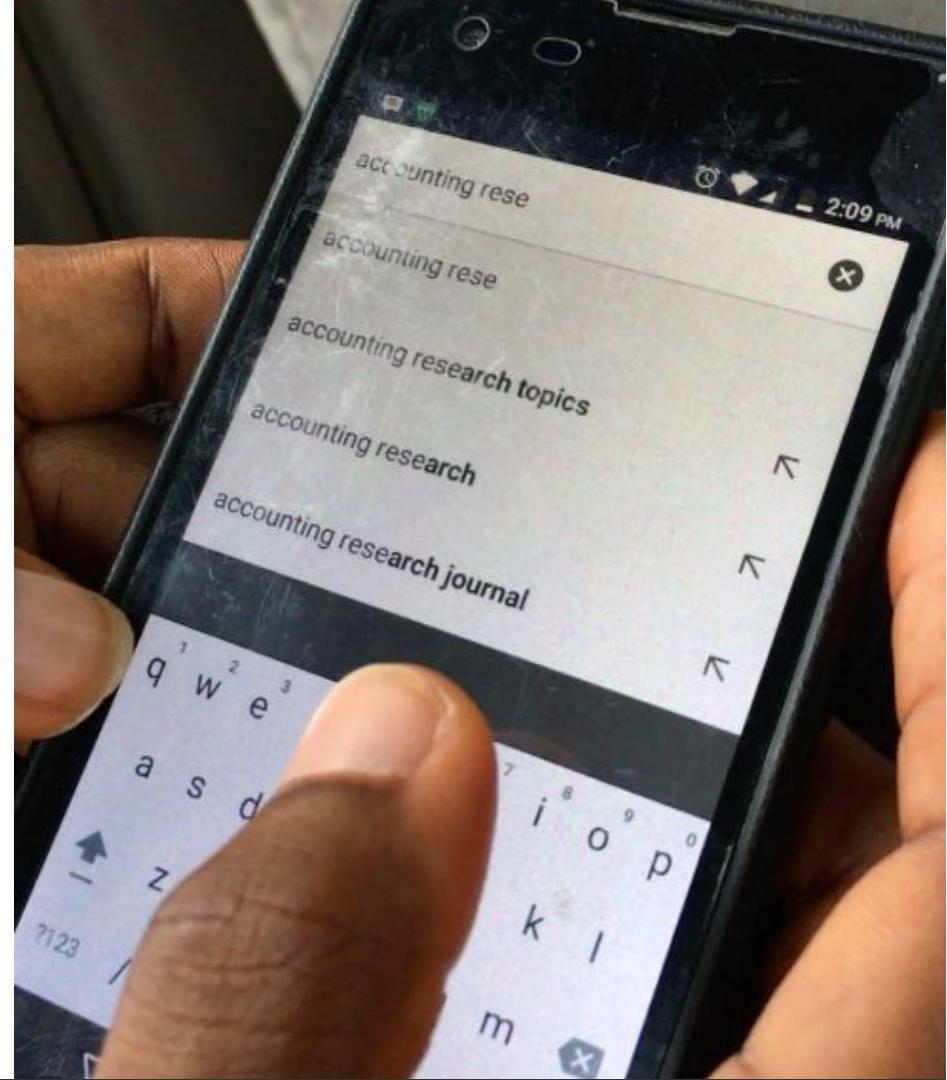
**Nigeria**

## Themes

1. Information seeking
2. Accessing the internet
3. Understanding the internet
4. Using the internet
5. Getting information online
6. Wikipedia Awareness
7. Wikipedia usage

**22. People confuse Wikipedia with a search engine or social media platform. This can create unrealistic expectations of its functionality.**

Learn more: [Research deck, slide 73](#)



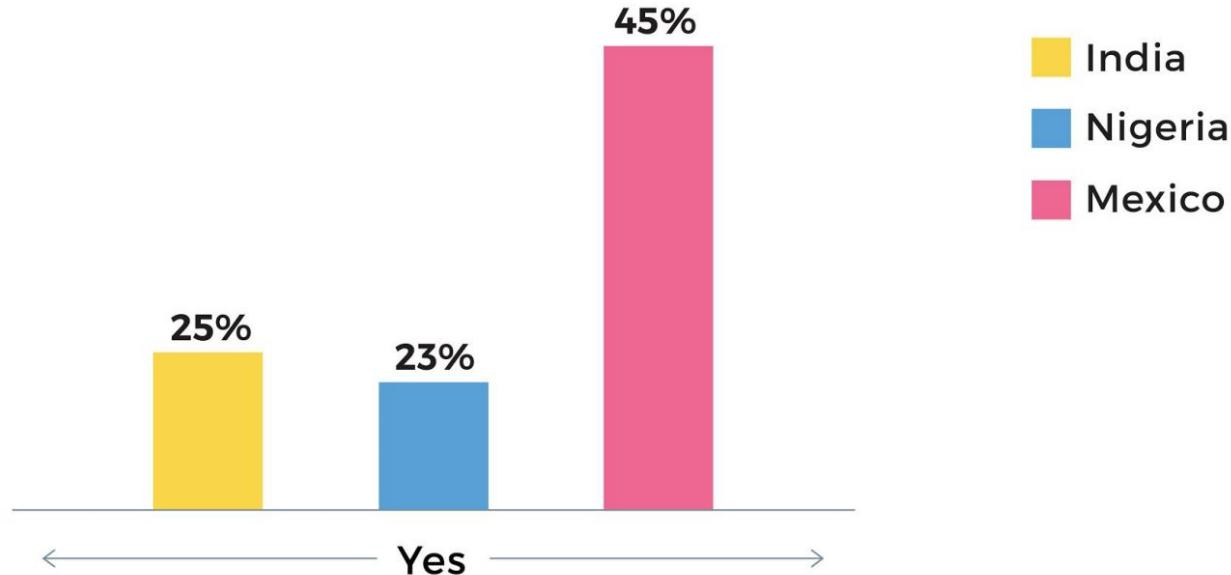
**“Wikipedia is a 'poor cousin' of Google. It is the lesser model.”**

**“Google and Wikipedia are similar. Google is more distributed; Wikipedia is more analytical and comprehensive.”**

**“Wikipedia is a social network. You'd use it if a friend in the US was on it and you wanted to connect with them.”**

## Phone survey findings

### Have you ever heard of Wikipedia?



# Wikimedia Research: Current priorities

## Artificial Intelligence as a service

*Socially aware* quality control

Content growth via recommendations

## Understanding editing culture

Toxic comments and personal attacks

Socializing newcomers

## Understanding readers

A taxonomy of Wikipedia readers

New readers

# Wikimedia Research: How to scale?

**1 : 100,000,000**

Current ratio of full-time Wikimedia researchers to Wikimedia monthly unique visitors  
(legacy comScore data)

A platform for open,  
collaborative, and reproducible  
research



FOUNDATION

# Formal collaborations



*Carnegie Mellon University*

*Dresden University of Technology*

*GESIS - Leibniz Institute for the Social Sciences*

*Institute for Scientific Interchange*

*Jigsaw*

*Stanford University*

*University of Minnesota*

*University of Washington*

*West Virginia University*

*Wikimedia Germany*

[https://www.mediawiki.org/wiki/Wikimedia\\_Research/Collaborators](https://www.mediawiki.org/wiki/Wikimedia_Research/Collaborators)

# Open, reusable research



# Open access policy

---

The Wikimedia Foundation's **mission** is to disseminate open knowledge effectively and globally. In keeping with this mission, the Wikimedia Foundation supports research in areas that benefit the Wikimedia community. We aim to make any work produced with our support **openly available** to the public and reusable on Wikimedia projects.

## 1. Expectations

---

Researchers will need to provide unrestricted access to and reuse of all their research output if their research receives support from the Wikimedia Foundation in the form of:

- funds;
- letter of endorsement;
- equipment, hosting, or office space;
- access to non-public data or special API privileges; or
- other support under an agreement between researchers and the Wikimedia Foundation.

[https://wikimediafoundation.org/wiki/Open\\_access\\_policy](https://wikimediafoundation.org/wiki/Open_access_policy)

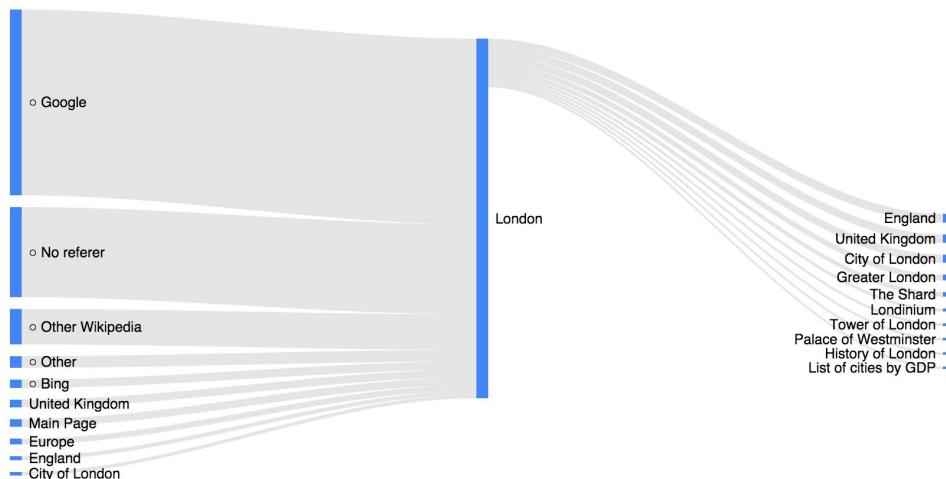
# Open data



# *Growing free knowledge through open data*

By [Dario Taraborelli](#), Wikimedia Foundation

March 13th, 2015



<https://blog.wikimedia.org/2015/03/13/open-data-sets/>

# An open notebook infrastructure



# PAWS

PAWS: A Web Shell (PAWS) hosts python notebooks and a terminal accessible with a browser.



<https://paws.wmflabs.org/>

<https://wikitech.wikimedia.org/wiki/PAWS>

Sample notebook by Brian Keegan

## Identify the most well-connected nodes

Compute the directed degree centralities. The network is constructed such that users contribute to articles. It should be impossible for articles to contribute to articles, or users to contribute to users, or articles to contribute to users.

Following the enforced direction, in-degree values for articles should be non-zero and reflect the number of users who contributed to them while it should be zero for users. The out-degree values for users should be non-zero and reflect the number of articles they contributed to while it should be zero for pages.

```
in_degree_d = {node:int(centrality)[len(collab_noboot_g)-1] for node,centrality in nx.in_degree_centrality(collab_noboot_g).items()}
out_degree_d = {node:int(centrality)[len(collab_noboot_g)-1] for node,centrality in nx.out_degree_centrality(collab_noboot_g).items()}
```

Convert the in- and out-degree distributions from above into a DataFrame called "cg\_degree\_df" for collaboration degree centrality DataFrame. Also create new columns in "degree\_df" that label whether the row is a page or user with filtering later on.

```
cg_degree_df['Page'] = pd.Series(cg_degree_df.index, index=cg_degree_df.index).str.contains('pr')
cg_degree_df['User'] = pd.Series(cg_degree_df.index, index=cg_degree_df.index).str.contains('us')
```

Look at the nodes with the highest in-degree: these are articles and the number of unique users who contributed to them.

```
cg_degree_df['In'].sort_values(ascending=False).head(10)
```

Look at the nodes with the highest out-degree: these are users and the number of unique articles (in the set) they contributed to.

```
cg_degree_df['Out'].sort_values(ascending=False).head(10)
```

Are there any articles with only a single editor?

```
cg_degree_df.head()
```

	In	Out
us:Editor:23	118	0
us:Wikimania	134	0
us:WikiTuring	99	0
us:WRC2012	99	0
us:Community	94	0
us:Neutrality	90	0
us:WikiProject	85	0
us:Jelasooy	79	0
us:CommonsEditor	78	0
us: Wiki Dot, Jupyter	184	0

Look at the articles that only have a single contributor to the set of hyperlinked articles. Some of these look like major topics, but they aren't the main page themselves, but rather sections that are placeholders pointing to another page elsewhere. For instance, an article like "Timeline 11/2001 attacks" exists but

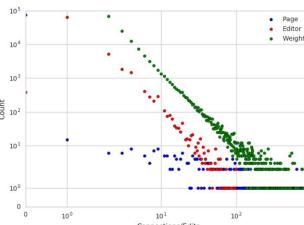
actually points to "September 11 attacks". The former page is an example of a redirect that has its own, but typically extremely short, revision history.

```
cg_degree_df.query("In == 1 & Page != 'In'")
```

```
Article_Short
prCampaign_finance_reform
prClimate_change_investigation
prClimate_change_investigation
prFort_McMurray
prGMO_greens
prHBO
prHacker_Farms
prHealthcare_health_care_plan
prHumanitarian_aid
prIslam_(person)
prReader_A_iphone_iPhone
prWar_in_Afghanistan_(2001-present)
prWar_in_Afghanistan_(2001-present)
Name: In, dtype: int64
```

```
in_degree_dist = cg_degree_df['In'].value_counts().reset_index()
out_degree_dist = cg_degree_df['Out'].value_counts().reset_index()
revise_weighted_edits = pd.merge(in_degree_dist, out_degree_dist, how='left').edit/(in_degree_dist['In']*out_degree_dist['Out'])
```

```
fig = plt.subplots(1,1)
in_degree_dist.plot.scatter(x='index',y='In',ax=ax, c='blue',label='Page')
out_degree_dist.plot.scatter(x='index',y='Out',ax=ax, c='red',label='Editor')
revise_weighted_edits.plot(x='index',y='edit',ax=ax,c='green',label='Weight')
ax.set_xlabel('Connections/Edits')
ax.set_xscale('log')
ax.set_ylabel('Count')
ax.set_yscale('log')
ax.set_xlim(1,1000)
ax.set_ylim(0,1e05)
```



*learn more about*  
**Open data, APIs and  
research tools**



workshop at 2pm, auditorium PPGI

# Conclusion



**WIKIMEDIA**  
FOUNDATION



# [m:Research:Newsletter](#)

**Wikimedia Research Newsletter**

Vol: 6 • Issue: 05 • May 2016

[CONTRIBUTE] [ARCHIVES]

*English as Wikipedia's lingua franca; deletion rationales; schizophrenia controversies*

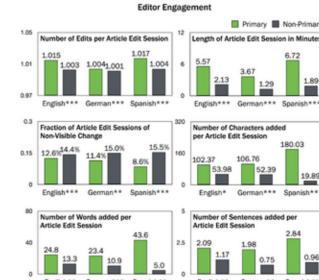
With contributions by: Morten Warncke-Wang, Piotr Konieczny, Federico Leva, Steve Jankowski and Tilman Bayer

## English still the lingua franca of Wikipedia

Reviewed by [Morten Warncke-Wang](#)

Many of the more active Wikipedia contributors are [multilingual](#). In the April 2011 Wikipedia Editors Survey, [\[supp 1\]](#) 72% of respondents said they read Wikipedia content in more than one language, and 51% said they contributed to multiple Wikipedias. Research has estimated that approximately 15% of active Wikipedians are multilingual. [\[supp 2\]](#) These contributors are important as they can enable knowledge transfer between different language editions of Wikipedia, yet little is known about who they are and what they do.

A recent paper published in [PLOS ONE](#) by researchers at [KAIST](#) and [OII](#), titled "Understanding Editing Behaviors in Multilingual Wikipedia"[\[1\]](#), adds to our knowledge of multilingual contributors by investigating their engagement level, topic interests, and language proficiency. The paper uses a dataset spanning a month of Wikipedia contributions in July–August 2013 and defines a *multilingual editor* as one who make contributions to multiple languages. Overall the dataset contains 12,577



Primary and non-primary editors show different engagement levels (figure 4 from the paper)

**@WikiResearch**

TWEETS: 2,559 FOLLOWING: 479 FOLLOWERS: 4,018 LIKES: 414

**Tweets** Tweets & replies Media

Pinned Tweet  
WikiResearch @WikiResearch · Sep 19  
W The next @WikiResearch showcase will be live-streamed this Wednesday 9/21 at 11:30am Pacific Time youtube.com/watch?v=tTdkVe... (1/3)

262 Photos and videos

Banner Configuration?

WikiMedia Research Showcase - September 2016 The September 2016 Wikimedia Research Showcase: [https://www.mediawiki.org/wiki/Wikimedia\\_ResearchShowcase](https://www.mediawiki.org/wiki/Wikimedia_ResearchShowcase)

2 4 \*\*\*

WikiResearch @WikiResearch · 18m  
W I Congresso Científico Brasileiro da Wikipedia kicks off #ccbwiki

#CCBWiki

New to Twitter?  
Sign up now to get your own personalized timeline!

[Sign up](#)

You may also like - Refresh



Worldwide Trends

Bob Dylan	1,000,000 tweets
#Nobel	8,816 tweets
#GloryDays	46,5K tweets
ノーベル文学賞	58,4K tweets
#FelizJuventud	20,3K tweets
ムハマド	11,9K tweets
#祷りの歌と歌の祷り	56,0K tweets
Dario Fo	49,2K tweets
さだやま	17K tweets
Francisco Correa	12,7K tweets

© 2010 Twitter About Help Terms Privacy Cookies Ads Info

# Thank you

## Acknowledgments

*The Wikimedia Research team, our fellows and collaborators:*

Samantha Becker, Grage Gellerman, Aaron Halfaker, Jonathan Morgan, Yuvi Panda, Abbey Ripstra, Amir Sarabadani, Nathaniel Schaaf, Morten Warncke-Wang, Robert West, Erik Zachte, Leila Zia

*Our formal collaborations:*

[mw:Wikimedia\\_Research/Collaborators](#)