# Uncovering Alternative Metrics:
## Data Mining Wikipedia for Evidence of Public Engagement and Impact

Lars Alvik, Peter Neish, Sally Tape

# Acknowledgment
# of country

# What we will cover

- **Background on Wikipedia ecosystem**

- **References in Wikipedia and disinformation**

- **Metric and Alternative Metrics**

- **Data mining Wikipedia**

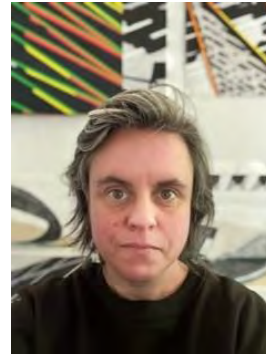- **Where to from here**

# Introduction, who are we?

**Peter Neish**

**Program Manager
Stewardship and Open Research**
University of Melbourne
Wikimedia Australia

**Sally Tape**

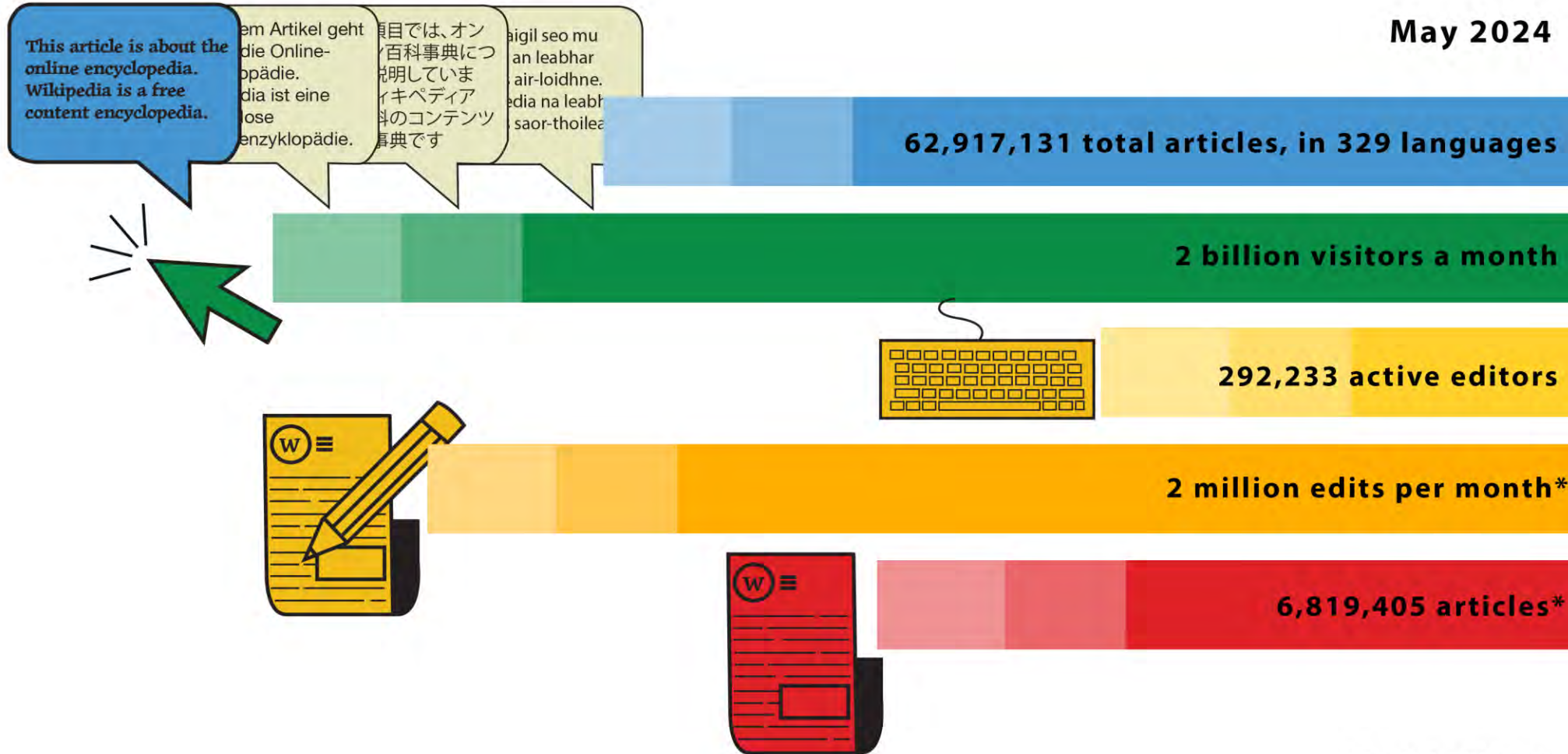**Open Research Support Specialist**
University of Melbourne

**Lars Alvik**

**Digital Curation Technical Specialist**
University of Melbourne

# Background on Wikipedia



This article is about the online encyclopedia. Wikipedia is a free content encyclopedia.

May 2024

62,917,131 total articles, in 329 languages

2 billion visitors a month

292,233 active editors

2 million edits per month*

6,819,405 articles*

*English Wikipedia

https://meta.wikimedia.org/wiki/List_of_Wikipedias

https://en.wikipedia.org/wiki/Wikipedia:Size_of_Wikipedia

# Background on Wikipedia



**Highly discoverable**



**FAIR**



**Easy to read
& understand**



**Impact &
integration**

# Background on Wikipedia

- **Structured knowledge banks**
- **Verifiable information**
- **Integration**



Wikimedia Foundation, CC BY-SA 3.0, via Wikimedia Commons

# The role of references in Wikipedia and disinformation

- **Ideally, everything in Wikipedia is verifiable by a reliable published\* secondary source eg:**
  - **Scholarly works: monographs, articles, textbooks**
  - **Well-established News organisations**

**\*made available to the public in some form (could be behind a paywall)**



The fin de siècle newspaper proprietor 1894: Frederick Burr Opper, Public domain, via Wikimedia Commons

# The role of references in Wikipedia and disinformation

- **Wikipedia is not a for original research**

- **Wikipedia is versioned and any vandalism or misinformation can be readily rolled back**

- **Bots and people can monitor and review content**

- **Studies continually find that Wikipedia is equal or more reliable than scholarly textbooks and published Encyclopaedias**
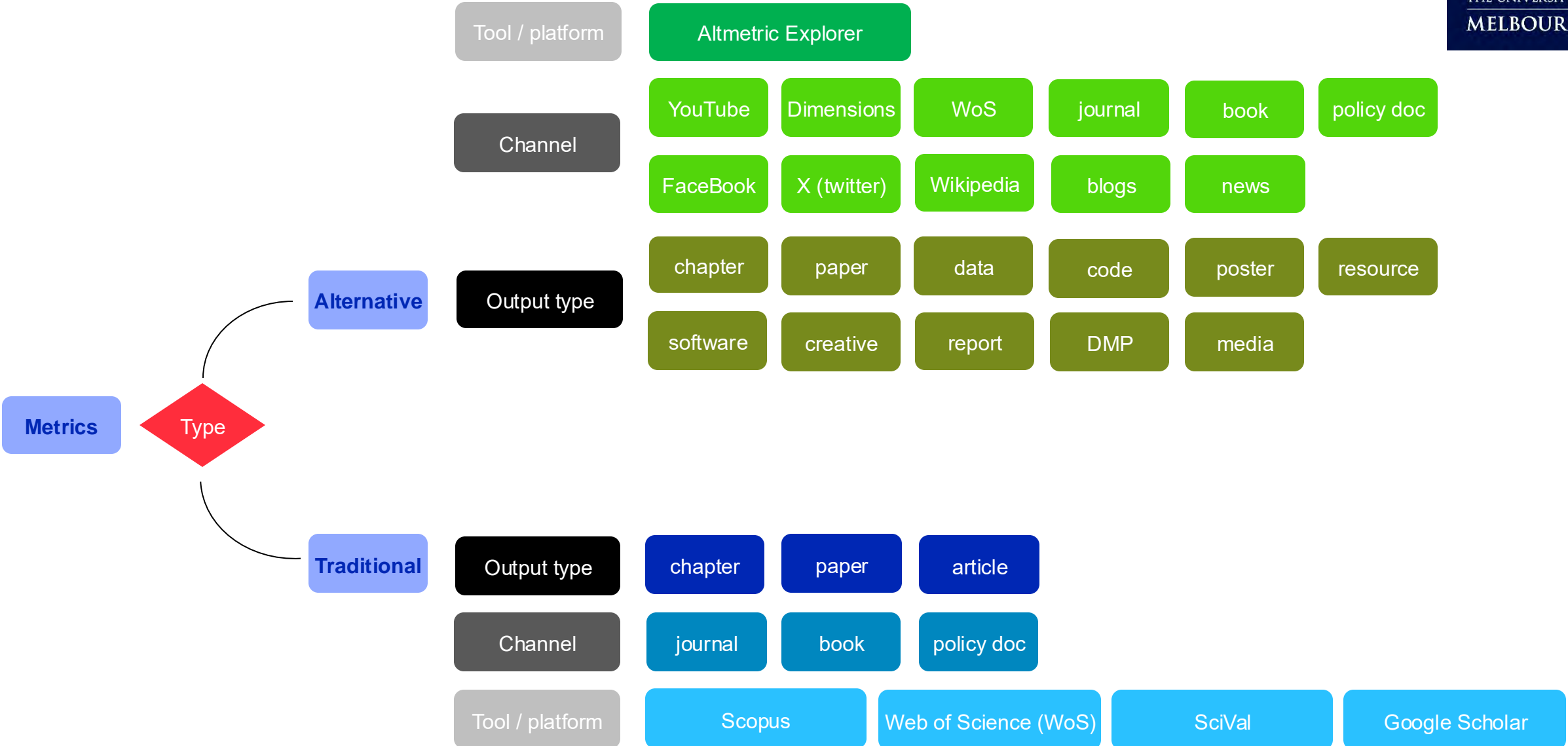


Joi Ito., CC BY 2.0 via Wikimedia Commons

# Why is Wikipedia interesting from a University perspective?



- **Impact through wide readership outside of scholarly community**

- **Contribute to Open Knowledge**

- **Content related to institution alumni or areas of expertise**

# Background on Metrics

# Background on Metrics

- Metrics from Wikipedia are tracked through 'mentions' on-line

- Wikipedia articles are published in 300+ different languages

- Metrics are currently recorded on the Altmetric Explorer
for 34 Wikipedia languages





https://www.altmetric.com/explorer/outputs?identifier=10.5194%2Fhess-11-1633-2007&scope=all&show_details=3112735

# What we did

- **Initially we wanted to look at the whole Wikipedia and use their tools to get all links to doi.org from articles within Wikipedia. (this didn't work)**

- **The alternative was to download the parsed html of each article and scan through for links to doi.org. This worked, but we had to cut down on the number of articles.**

- **We narrowed it down to articles within the first 3 category levels of a handful of Australian science categories.**

# Scripts

1. **Runs through categories and saves page names to a list.**
2. **Gets the parsed text from the Wikipedia API and saves DOIs to a list.**
3. **Checks doi.org if the DOI is a Crossref or DataCite DOI.**
4. **Gets metadata about the DOI from the relevant API.**
5. **Print stats.**



Category name: eg. Biotechnology

Script recursively going through the category

Articles

Parsing wikitext and search for DOIs

DOIs in each article

Get DOI metadata from Crossref API

# Findings (1 / 2)

## History of Australia

- **1255 articles / 243 with DOIs (19%)**
- **1526 authors listed in the DOIs**
  - **446 with affiliations (29%)**
  - **108 with ORCIDs (7%)**

## Politics of Australia

- **2768 articles / 279 with DOIs (10%)**
- **1228 authors listed in the DOIs**
  - **290 with affiliations (23%)**
  - **82 with ORCIDs (7%)**

# Findings (2 / 2)

## Geology of Australia

- 888 articles / 283 with DOIs (32%)
- 3649 authors listed in the DOIs
  - 980 with affiliations (27%)
  - 368 with ORCIDs (10%)

## Biota of Australia

- 5961 articles / 2641 with DOIs (44%)
- 15267 authors listed in the DOIs
  - 2691 with affiliations (17%)
  - 954 with ORCIDs (6%)

# Problems

- **The lack of affiliations and affiliations being a free text field makes it hard for us to gauge the impact of our researchers.**

- **only a subset of Wikipedia, so we can't get a complete picture**

- **only DOIs, doesn't take into account references that don't use a persistent identifier**

- **Long process, many steps, and time consuming.**

# Future!

- **WikiData - https://www.wikidata.org/wiki/Q21090057**
- **Scholia - https://scholia.toolforge.org/**
- **Wikicite - http://wikicite.org/**