# WIKIPEDIA AND OPEN ACCESS

**Puyu Yang**
Institute for Logic, Language and Computation (ILLC)
The University of Amsterdam
The Netherlands
p.yang2@uva.nl

**Ahad Shoaib**
School of Computer and Communication Sciences
École Polytechnique Fédérale de Lausanne (EPFL)
Switzerland
University of Waterloo
Canada

**Robert West**
School of Computer and Communication Sciences
École Polytechnique Fédérale de Lausanne (EPFL)
Switzerland

**Giovanni Colavizza**
Institute for Logic, Language and Computation (ILLC)
The University of Amsterdam
The Netherlands

## ABSTRACT

Wikipedia is a well-known platform for disseminating knowledge, and scientific sources, such as journal articles, play a critical role in supporting its mission. The open access movement aims to make scientific knowledge openly available, and we might intuitively expect open access to help further Wikipedia's mission. However, the extent of this relationship remains largely unknown. To fill this gap, we analyze a large dataset of citations from Wikipedia and model the role of open access in Wikipedia's citation patterns. We find that open-access articles are extensively and increasingly more cited in Wikipedia. What is more, they show a 15% higher likelihood of being cited in Wikipedia when compared to closed-access articles, after controlling for confounding factors. This open-access citation effect is particularly strong for articles with low citation counts, including recently published ones. Our results show that open access plays a key role in the dissemination of scientific knowledge, including by providing Wikipedia editors timely access to novel results. These findings have important implications for researchers, policymakers, and practitioners in the field of information science and technology.

***Keywords*** Wikipedia · Open Access · Open Science

## 1 Introduction

Open access (OA) publishing has emerged as a popular alternative to traditional subscription-based models, with the goal of making research more widely accessible to the public. This movement has gained momentum over the years, with many scholars recognizing the benefits of open access in promoting the dissemination of scientific knowledge and funding bodies adopting OA mandates [1, 2].

One of the most significant beneficiaries of the open-access movement is Wikipedia. As a dynamic platform for sharing and disseminating knowledge across the globe, Wikipedia is relied upon by millions of users every day to satisfy a wide range of information needs [3]. It has become a critical source of information for both the general public and academic researchers, and its impact is extending beyond the realm of general knowledge and into the academic sphere [4, 5, 6].

Citations to sources are critical to Wikipedia's mission of providing verifiable and reliable information to its readers[1]. Among several sources, academic and peer-reviewed publications are usually considered the most reliable[2]. Wikipedia's extensive use of citations makes it possible to analyze its reliance on academic publications, which is a central aspect

---

[1]https://en.wikipedia.org/wiki/Wikipedia:Core_content_policies.

[2]https://en.wikipedia.org/wiki/Wikipedia:Verifiability#Reliable_sources.

of our investigation. Seminal prior research has shown that open-access publications are more likely to be cited in Wikipedia than paywalled publications [7]. However, the extent to which open-access publications are relied upon by Wikipedia at the journal article level of granularity has not been examined in further detail.

In light of this, our study seeks to fill this gap by examining how open-access publications affect Wikipedia at the article-level granularity. Specifically, we aim to answer the following research questions:

1. RQ1: To what extent does Wikipedia rely on open-access publications? How has this been changing over time?
2. RQ2: To what extent does the open-access status of an article influence its likelihood of being cited in Wikipedia?

To address these questions, we will use descriptive statistics and regression analysis based on the *Wikipedia Citations* dataset [8]. To identify the information in article-level granularity such as the open access status of publications, we will use the OpenAlex and Scimago services. Our research contributes to our understanding of the role of open access in the dissemination of scientific knowledge and the impact of Wikipedia in this process, as well as informing policy and practice in the realm of open scholarly communication.

The remainder of the paper is structured as follows. Section 2 provides an overview of existing research in the field. Section 3 describes our dataset and methodology. Section 4 presents descriptive statistics of open-access publications in Wikipedia (RQ1), and then uses regression analysis to model the influence of open-access status on the likelihood of a paper being cited in Wikipedia (RQ2). Finally, Section 5 and 6 offer a discussion and conclusion of our findings.

## 2 Previous Work

### 2.1 Open Access in Science

The key idea behind open access (OA) is to provide unrestricted and free access to scientific outcomes, thus enhancing their visibility and reach regardless of financial or geographical constraints [9, 10]. The increasing popularity of OA in academic publications has generated extensive discussions among scholars in recent years. Empirical studies have shown that OA has had a significantly positive impact on the accessibility of scientific journal articles [11]. A comprehensive analysis of OA publications shows that at least 27.9% of the total 19 million scientific articles are OA [1]. In addition, studies report that around 55% of articles indexed by the Web of Science from 2009 to 2014 are OA, and more than 50% of scientific papers published since 2007 can be accessed freely [12, 13]. Among the various OA policies, Bronze OA is the most common type [1]. Although the distribution of OA varies across different fields, General Science, Technology, and Biomedical research have relatively higher OA rates, while Engineering and Arts&Humanities have lower rates [13, 12].

An "open access citation advantage" has also been a topic of ongoing debate. Some researchers have observed that a citation advantage linked to open access exists, although the effect magnitude varies based on the dataset and methods used. For example, OA articles have been found to receive 18% more citations than average based on Web of Science, while Scopus reports an even higher, positive 40% effect [1, 13]. Therefore, the effects of OA on citation patterns remain a topic of interest and active investigation.

### 2.2 Science and Wikipedia

With the rapid development of the internet, traditional peer review and journal publication can no longer meet the need for new ideas' development [14]. As one of the largest encyclopedias worldwide, Wikipedia aims to effectively and globally distribute information based on scientific findings[3], thereby making it a valuable altmetric source [15].

Previous research has indicated that the topical coverage of Wikipedia bears some similarity to science, with 13.44% of its citations being from open access journals [16]. STEM fields, especially biology and medicine, comprise the most prominently featured scientific topics in Wikipedia [17]. Specific fields such as Medicine and Psychology have a comparatively high number of citations to research papers on Wikipedia and are sometimes employed as a gateway to further academic research [18, 19]. Furthermore, journal articles cited from Wikipedia tend to be published in high-impact journals (e.g., by impact factor) and in open access more frequently than the average article [20, 7].

Science contributes a lot to Wikipedia, yet the influence goes both ways. Prior studies have established that Wikipedia can enhance the citation impact of an article it cites [21]. Additionally, Wikipedia has demonstrated its ability to rapidly and reliably incorporate novel scientific findings in response to ongoing public events or crises [22].

---

[3]`https://wikimediafoundation.org/about/mission/`

### 2.3 Citation Analyses of Wikipedia

The open release of citation datasets from Wikipedia has led to a surge of studies examining citation analysis on Wikipedia [23, 24, 8]. Among the articles on Wikipedia, 6.7% cite at least one journal article with an associated DOI [8], with the majority of cited journal articles being published in the past two decades [17]. [25] have conducted a study on the quality of citations in Wikipedia during COVID-19 and have found that Wikipedia mostly cites reliable sources and prefers open-access articles. Some researchers have focused on user behavior regarding reference usage on Wikipedia. [26] have found that engagement with citations on Wikipedia is generally low, but references are more frequently looked up when the information is not included.

Despite the increasing number of citation studies on Wikipedia, the relationship between open access and Wikipedia requires further exploration. Previous research has examined the effect of OA on Wikipedia, and found that articles with OA were 47% more likely to be cited than Closed Access articles when controlling for journal and research fields [7]. However, the study only used a smaller subset of journals and did not control for confounding factors such as citation counts. Furthermore, the various OA policies were not explored separately. To better understand the relationship between OA and Wikipedia, this study aims to improve on the previous results by using a more rigorous and comprehensive methodology with article-level granularity and controlling for more confounding factors.

## 3 Data

The data collection process followed a specific workflow as outlined below. Firstly, we obtain all citations from English Wikipedia to any source using the open dataset known as *Wikipedia Citations* [8]. Additionally, to identify journal articles, we relied on the classification and DOI provided by *Wikipedia Citations*. Secondly, to augment the primary dataset, we utilized the OpenAlex API to obtain information on the open-access status of journal articles. Finally, we employed data from Scimago to obtain relevant information for each journal. The subsequent sections provide a detailed description of the main datasets used in the study.

### 3.1 Wikipedia Citations

The primary dataset utilized in this research is *Wikipedia Citations*, a comprehensive dataset of over 29 million citations extracted from more than 6 million English Wikipedia articles as of May 2020 [8]. Out of these, around 2.5 million citations are classified as journal articles, with 1,705,085 of them containing a digital object identifier (DOI). To augment our analysis, we leverage the OpenAlex API[4] to obtain relevant article details such as OA status, publication date, publisher, and concepts, among others, for each DOI.

### 3.2 OpenAlex and Scimago

To examine the impact of open access (OA) articles, we utilized OpenAlex, a free and open platform that provides data on academic papers and researchers [27]. OpenAlex draws data from Microsoft Academic Service (MAG) and Crossref, among other sources, and contains more than 240 million academic works that can be utilized in the fields of bibliometrics, science and technology studies, and science of science policy [28, 29]. To obtain the necessary data for journal articles in *Wikipedia Citations*, we utilized the OpenAlex API[5]. After matching, we obtained article information from OpenAlex for 1,696,108 journal articles.

We considered five categories for the OA policy in our study, following the classification scheme proposed by [1]:

1. Gold: Refers to articles published in an open-access journal that is indexed by the Directory of open access Journals (DOAJ).

2. Green: Refers to articles that are available on the publisher's webpage but are also freely accessible in a pre-print repository.

3. Hybrid: Refers to articles that are available for free under an open license in a toll-access journal.

4. Bronze: Refers to articles that can be read for free on the publisher's webpage but do not have an easily identifiable license.

5. Closed: Refers to all other articles, including those shared only on an academic social network or in Sci-Hub.

---

[4]https://docs.openalex.org/api-entities/works.
[5]https://docs.openalex.org/how-to-use-the-api/get-single-entities.

Additionally, we collected journal information for conducting a regression analysis on the influence of OA. We obtained this information by downloading data from Scimago[6]. Scimago is an open-access resource that provides an internationally accepted journal rank indicator for analysis in the fields of scientometrics and infometrics [30, 31, 32]. We equipped each journal with an SJR score, H index, and other relevant information.

## 4 Results

We have augmented our analysis by incorporating additional metadata from OpenAlex and Scimago, which allowed us to obtain information for 99.4% (1,696,108) of the total 1,705,085 citations to a valid DOI. From these, we extracted 1,152,141 publications (unique DOIs) and an associated open access (OA) status. Our findings show that 42.3% (716,278) of the citations and 39.1% (450,277) publications were OA (i.e., not closed). These results align with the trend observed in the whole scientific community, where the percentage of OA articles has been increasing steadily, reaching 28% in 2018 [1]. Given that Wikipedia is an open and free online encyclopedia, it is not surprising that it cites OA articles more often. Additionally, our findings are consistent with the distribution of OA articles on Unpaywall, which reported a percentage of 47% [1].

### 4.1 Characterizing open access Articles within Wikipedia

We present our findings on the distribution of open access (OA) policies in Wikipedia citations. Our results demonstrate that the most commonly observed OA policy in Wikipedia citations is the bronze policy, which is again consistent with trends in scholarly literature [1]. The second most popular OA policy observed in Wikipedia citations is Green, which is significantly more prevalent than the Gold policy. This trend may be attributed to differences in the reference acquisition methods used by Wikipedia editors compared to researchers. Additionally, our comparison of OA policies in scientific articles and Wikipedia reveals a similar trend [1]. Interestingly, the percentage of Green policy in Wikipedia is much higher than in scientific articles, suggesting that pre-print repositories are a valuable source of information for Wikipedia editors.



(a) Over citations.

(b) Over unique journal articles.

Figure 1: Distribution of open access status by policy.

Figure 2 displays the change in OA status over the years. The blue line represents the proportion of OA status for each publication year, while the black line indicates the proportion of citations published each year. The figure indicates that there has been a steady increase in citing new OA articles over the past 40 years. This trend implies that OA articles are becoming more prevalent and could shape the future of science in Wikipedia. However, it should be noted that the percentage of OA may be subject to backfill bias, as articles may be labeled with their OA status after publication. Additionally, there has been a clear upward trend in the percentage of OA articles cited in Wikipedia from 2015 to 2019, with over 60% OA articles being cited during this period.

We proceeded to examine the breakdown of open access (OA) status and policies by journals in our dataset, which includes 40,191 journals. In order to visualize this information effectively, we determined the number of citations for each journal and selected the top 20 for analysis. Figure 7 displays the total number of citations for the top 20

---

[6]`https://www.scimagojr.com/aboutus.php`.

Figure 2: Fraction of OA citations by year (blue), and citations by publication year (black). The left y-axis serves for OA citations by year and the right y-axis for the fraction of a year's citations overall.

journals, with blue and orange representing OA and non-OA articles, respectively. As found in previous studies, some high-impact journals such as *Nature*, *PNAS*, and *Science* appear frequently on Wikipedia [20] and account for 5.7% of all citations. However, inferring article OA policy based on whether journals are classified as "open access" or "Closed Access" can be misleading [7], as there is a high variance in OA status among articles within journals. For example, while most articles in *Nature* and *Science* are OA, there are also non-OA articles. Therefore, it is inappropriate to categorize journals as solely "Open" or "Closed." It should be noted that our analysis reflects the OA policy distribution at a given point in time, and journals such as *PNAS* may have policies that automatically convert closed-access articles to OA after a certain period of time. As such, the effect of OA policies may only be observed in cases where the article is also classified as Green OA or is cited after the specified period.

To further investigate the distribution of OA policies among the top 20 journals, we plotted the data in Figure 8. Our analysis revealed a growing trend of bronze OA policies among journals. However, some journals that classify themselves as OA, such as "*Journal of Biological Chemistry*"[7] and "*PLOS one*"[8], have a significant proportion of articles classified as Hybrid or Gold OA. While there may be limitations in the classification of OA articles by Unpaywall, we accepted their classifications in our study.

We present an analysis of the distribution of open access (OA) status by OpenAlex concepts, as shown in Figure 3 and Figure 4. We used OpenAlex, a dataset containing 65k concepts and 19 root-level concepts. We employed fractional counting to determine the number of citations for each root-level concept. Given that 42.3% of citations on Wikipedia are OA, we used this percentage as a baseline for OA proportionality, represented by the black dotted line in Figure 3. Our analysis revealed substantial variance in field-specific OA proportions.

---

[7]https://www.elsevier.com/journals/journal-of-biological-chemistry/0021-9258/open-access-journal.

[8]https://journals.plos.org/plosone/s/journal-information#loc-open-access.

Figure 3: Distribution of OA status and count of citations by OpenAlex concept.



Figure 4: Distribution of OA policies by OpenAlex concept.

In Figure 3, we plotted the percentage of citations with OA status for each concept on the left, and on the right, we plotted the total number of citations of concepts, arranged from the largest to the smallest. Notably, only four concepts exhibited OA proportions higher than the average: Biology (53%), Physics (48%), Mathematics (46%), and Medicine (45%). These four concepts featured a roughly equal number of OA and closed access articles. By contrast, Psychology (30%), Art (30%), and History (22%) exhibited the lowest proportions of OA articles among cited Wikipedia articles. In general, our analysis suggests that fields with more OA articles cited on Wikipedia may indicate that the corresponding Wikipedia articles are accessible to a wider audience of editors and users. We suggest further research to investigate whether articles citing more OA articles on Wikipedia are indeed more popular among users, potentially motivating authors and publishers to release more studies as OA.

We also analyzed the distribution of OA policies in each concept. Our analysis indicates that bronze and green policies dominated most of the concepts in OA articles, with the exception of Art, in which the gold policy plays an important role.

## 4.2 OA Citation Advantage

In order to gain deeper insights into the influence of open access (OA) articles on scientific discourse in Wikipedia, we have developed a series of statistical models with varying formulations. Our aim is to examine the potential for an "OA citation advantage" effect in Wikipedia. To achieve this, we will compare citations in Wikipedia with a carefully

6

curated subset of scientific articles that are representative of those cited in Wikipedia. This subset has been selected using stratified sampling across citation count, journal, concept, and time of publication, based on the comprehensive set of scientific articles provided by OpenAlex. By conducting this analysis, we aim to provide a clearer understanding of the role that OA articles play in shaping scientific discourse in Wikipedia and to identify any potential advantages that may be associated with their use.

### 4.2.1 Model Specification

**Dependent variable**

To assess the potential advantage of OA articles in Wikipedia, we defined a binary dependent variable, denoted by *is_wiki*, that indicates whether an article has been cited in Wikipedia or not. Since our primary dataset consists solely of articles cited in Wikipedia, we use OpenAlex to obtain negative samples of articles not cited in Wikipedia, via stratified sampling.

**Independent variable**

To assess the impact of OA articles on their citation rates in Wikipedia, we analyze two types of variables: article-level and journal-level. At the article level, we consider the number of citations (*times_cited*), whether the article is OA or not (*is_oa*), the time of publication (*article_age*), and the field of research (*concept*). These features have been shown to have an influence on citation impact in previous studies [33, 34, 35, 36, 7, 20]. At the journal level, we primarily consider the Scimago journal rank (*SJR*). To account for changes in Scimago journal rank over time, we assign the rank corresponding to the year 1999 for articles published before 1999 and the rank assigned to journals for the same year an article was published.

Although these variables have been widely used to model citation impact in previous studies, little analysis has directly linked these indicators to whether an article is cited in Wikipedia or not, specifically with respect to different OA policies.

In this study, we use logistic regression as our model, which is usually used to analyze the relationship between a binary dependent variable and one or more independent variables. The logistic regression weights represent the size of the individual contributions of each predictor variable to the target variable. Figure 5 illustrates the assumed causal structure of Wikipedia's OA citation adoption effect, with a black line representing an assumed causal relationship between two variables. Specifically, we assume that the likelihood of a journal article being cited in Wikipedia is directly influenced by its article features, citation counts, and OA policy. At the same time, the OA policy can also influence the citation counts of the article, causing a further mediated effect on the adoption of this article in Wikipedia. With our models, we are interested in measuring both the (controlled) direct effect and the total effect of OA policy on being cited in Wikipedia. The former is shown as a black thick line in Figure 5, while the latter accounts for both the direct and the mediated (via citation counts) effects.



Figure 5: Causality of Wikipedia open access citation adoption effect.

### 4.2.2 Dataset construction

We aim to create a balanced dataset of journal articles that can be analyzed using regression analysis. The initial dataset was sourced from Wikipedia and contained 1,152,141 unique scientific articles. To initiate our regression analysis, we constructed two datasets. The first dataset was stratified by Journal and Year of publication, which will be utilized in our first two models. In this dataset, we finally get 708,156 articles. The second dataset was constructed by adding Journal, Year of publication and Concept as stratifying variables, and it will be used in the third model to account for the influence of ***concept***. For the second dataset, to restrict ourselves to root-level concepts and avoid any ambiguity, we filtered the citations to only include those with one associated concept, resulting in a set of 349,176 articles. We then aimed to construct a corresponding set of articles for these two datasets from OpenAlex, excluding those cited in Wikipedia, through a process of negative sampling.

To achieve this, we applied a stratified sampling strategy over the strata of Journal, Year of publication, and Concept, with the aim of obtaining a set of negative samples that were as similar as possible to the set of filtered articles. To reduce noise in the sampling strategy, we removed journals with no corresponding name available in Scimago and those with less than 20 citations. We also removed all articles published before 1900 to remove sparsely mentioned dates and accept a slight recency bias.

After pre-processing, we group the number of articles within each stratum and proceed as follows:

1. Filter the whole set of OpenAlex articles to those matching the fields in the strata.
2. If there are fewer articles in this filtered set, discard the strata and remove the respective articles from the dataset.
3. Otherwise, randomly sample the same number of articles from the filtered set and add it to the set of negative samples.

After repeating this process for all strata (156,354 for the first dataset and 73,353 for the second dataset), we obtain a final negative set size of 678,866 for the first dataset and 211,825 for the second one. Combining this with the remaining cited articles results in a total dataset size of 1,357,732 (first dataset) and 423,650 (second dataset).

To ensure the robustness of our sampling methodology, we repeated the process five times, resulting in five different sets that were used in the analyses. Although our method of matching strata to construct a set of negative samples is an approximation of the more rigorous method of propensity score matching (PSM), the discrete nature of our strata and the large population size contribute to the robustness of our analysis. A descriptive overview of this curated dataset can be found in Tables 7, 8, 9, 10 and 11 in the appendix.

### 4.2.3 Model results

We use logistic regression for interpretability and expressiveness, and we also use log transforms on the continuous variables. Our primary binary logit regression model is formulated as:

$$is\_wiki = article\_features + is\_oa \tag{1}$$

with

$$article\_features = ln(article\_age) + ln(SJR + 1)$$

a secondary formulation is:

$$is\_wiki = article\_features + is\_oa + ln(times\_cited + 1) \tag{2}$$

With each analysis being performed on each of our 5 samples, we consider a coefficient statistically insignificant if it is insignificant on the results of at least one sample. Below, we report the effects of the odd ratios based on the mean odd ratios across all 5 samples.

To assess the total effect of open access (OA) on citation adoption, we used the first formulation to examine the effect of *is_oa*. Our analysis of the data presented in Table 1 revealed that OA articles have a 15% higher odds of being cited in Wikipedia compared to closed access articles, and this finding was consistent across all five samples ($0.142 \pm 0.007$).

Table 1: Coefficients for overall OA adoption. Results for the first sample, model 1, $R^2 = 0.0007$.

| Feature | coef | odds_ratios | P>z |
|---|---|---|---|
| *is_oa* | 0.140 | 1.150 | 0.000 |
| $\ln(article\_age)$ | 0.025 | 1.025 | 0.000 |
| $\ln(1 + SJR)$ | -0.021 | 0.979 | 0.000 |

To investigate the direct effect of open access on the adoption of citations in Wikipedia while accounting for the relationship between citation count and open access policy, we utilized the second formula across five samples, as detailed in Table 2. After considering the citation count, the odds of an OA article being cited in Wikipedia is 12% higher than the closed-access article ($1.12 \pm 0.01$ across all samples).

Table 2: Coefficients for overall OA adoption. Results for the first sample, model 2, $R^2 = 0.012$.

| Feature | coef | odds_ratios | P>z |
|---|---|---|---|
| *is_oa* | 0.112 | 1.118 | 0.000 |
| $\ln(1 + times\_cited)$ | 0.162 | 1.176 | 0.000 |
| $\ln(article\_age)$ | 0.019 | 1.019 | 0.000 |
| $\ln(1 + SJR)$ | -0.214 | 0.807 | 0.000 |

In contrast, our analysis showed that other factors such as Scimago Journal Rank and publication age did not have a consistent impact on citation behavior in Wikipedia, which diverges from previous research on citation behavior in scientific fields. However, this inconsistency may be explained by the strong effect of citation count in Wikipedia, which had a significantly positive correlation with being cited, as indicated by its high coefficient and corresponding odds ratio.

To gain insight into the interaction between open access (OA) status and citation counts in Wikipedia, we use formula 2 to create a graph plotting the functions that contain these two variables and article features.



Figure 6: OA adoption effect at varying citation counts, based on model 2.

The graph, shown in Figure 6, displays the dependent variable, ***is_wiki***, on the y-axis and the citation counts on the x-axis. Articles are grouped according to their citation count (variable ***times_cited***). We plot the average model prediction calculated on each bin using the first data sample and provide 95% bootstrapped confidence intervals for each bin (faded color). The red line illustrates the trend of OA adoption by citation count under the condition that the OA status is closed, while the blue line shows the trend under the condition that the OA status is open. This graph reveals several insights. Firstly, when the citation counts are lower, there is a significant difference between OA and closed articles with the former getting a boost in their likelihood to be cited in Wikipedia. However, with increased citation

9

counts, this OA adoption effect becomes less distinguishable. In our previous work [17], we show that articles cited fewer than 100 times account for 70% of the total cited articles, and only about 3% of articles are cited 1,000 times or more. Therefore, most citations in Wikipedia benefit from this OA effect. We speculate that the OA adoption effect might be due to Wikipedia editors having an easier time finding and accessing open research results earlier, before they accumulate citations and therefore receive peer recognition.

Furthermore, we employed the first and second formulations to examine the impact of OA policy on citation adoption, with 'closed' as the baseline. The results, presented in Table 3, show that all OA policies significantly enhanced the overall adoption effect for OA articles. In Table 4, we used the second model to examine the indirect effect of OA policy and found a similar trend with the exception of the gold policy. To test the robustness of our findings, we conducted the regression in all 5 samples, and the results are reported in Tables 5 and 6.

Table 3: Coefficients for OA adoption by the policy. Results for the first sample, model 1, $R^2 = 0.001$.

| Feature | coef | odds_ratios | P>z |
|---|---|---|---|
| bronze | 0.101 | 1.106 | 0.000 |
| gold | 0.032 | 1.032 | 0.000 |
| green | 0.162 | 1.176 | 0.000 |
| hybrid | 0.190 | 1.210 | 0.000 |
| $\ln(article\_age)$ | 0.019 | 1.019 | 0.000 |
| $\ln(1 + SJR)$ | -0.036 | 0.965 | 0.000 |

Table 4: Coefficients for OA adoption by the policy. Results for the first sample, model 2, $R^2 = 0.01$.

| Feature | coef | odds_ratios | P>z |
|---|---|---|---|
| bronze | 0.110 | 1.116 | 0.000 |
| gold | 0.011 | 1.011 | 0.159 |
| green | 0.102 | 1.107 | 0.000 |
| hybrid | 0.159 | 1.172 | 0.000 |
| $\ln(article\_age)$ | 0.012 | 1.012 | 0.000 |
| $\ln(1 + SJR)$ | -0.214 | 0.807 | 0.000 |
| $\ln(1 + times\_cited)$ | 0.150 | 1.162 | 0.000 |

Table 5: Coefficients for OA adoption by the policy. Results across all 5 samples, model 1, $R^2 = 0.001$.

| Feature | coef | odds_ratios | P>z |
|---|---|---|---|
| bronze | 0.103 | 1.108 | 0.000 |
| gold | 0.033 | 1.033 | 0.000 |
| green | 0.164 | 1.178 | 0.000 |
| hybrid | 0.196 | 1.216 | 0.000 |
| $\ln(article\_age)$ | 0.019 | 1.019 | 0.000 |
| $\ln(1 + SJR)$ | -0.036 | 0.965 | 0.000 |

Table 6: Coefficients for OA adoption by the policy. Results across all 5 samples, model 2, $R^2 = 0.01$.

| Feature | coef | odds_ratios | P>z |
|---|---|---|---|
| bronze | 0.112 | 1.119 | 0.000 |
| gold | 0.012 | 1.012 | 0.128 |
| green | 0.104 | 1.109 | 0.000 |
| hybrid | 0.163 | 1.177 | 0.000 |
| $\ln(article\_age)$ | 0.012 | 1.012 | 0.000 |
| $\ln(1 + SJR)$ | -0.215 | 0.807 | 0.000 |
| $\ln(1 + times\_cited)$ | 0.150 | 1.162 | 0.000 |

## 5    Discussion

The popularity and growth of open access (OA) have significantly enhanced the dissemination of scientific knowledge. Our research demonstrates that Wikipedia editors regularly and increasingly cite OA articles, especially those published with bronze and green policies. While high-impact journals remain a preferred secondary source of Wikipedia, there are differences in the distribution of OA articles within journals. This finding underscores the importance of conducting research at the article level. Moreover, differences in OA Wikipedia citations among disciplines are observed, with biology, physics, and mathematics having higher OA citation rates, and disciplines in the social sciences and humanities having comparatively lower rates. Our study further reveals that the odds of an OA article being cited in Wikipedia over a closed-access article increases by 15% on average. On the one hand, Wikipedia appears to have an especially significant impact on the dissemination of OA articles with lower citation counts (e.g., those recently published), reflecting its ability to rapidly respond to new scientific advances. On the other hand, this effect might also reflect the fact that open-access articles are easier to find and use for Wikipedia editors who do not necessarily have access to journal subscriptions.

We acknowledge certain limitations in our study. For instance, our focus on articles with DOIs means that some conference papers and earlier literature were excluded, thus future research may consider additional sources. Additionally, while our regression model accounted for significant factors such as OA status, OA policy, and citation counts, other causal variables, such as article length, may also affect article citations on Wikipedia. Finally, our study did not consider time as an analytical dimension, and future research should examine the edit history of Wikipedia to extract specific data at the time when an article was cited. This would allow to further clarify the causal mechanisms at play when considering open access and Wikipedia.

## 6    Conclusion

This study examined the impact of open access (OA) on Wikipedia by analyzing article-level features. By utilizing a large dataset of citations from Wikipedia, coupled with OA-related metrics from OpenAlex and journal information from Scimago, we investigated the prevalence and significance of OA in Wikipedia. The results show that OA plays a crucial role in Wikipedia: OA articles are increasingly more cited over time, and have a higher chance of being cited in Wikipedia than similar closed-access articles. In particular, articles with low citation counts (e.g., those recently published) are substantially more likely to be cited in Wikipedia. These findings suggest that OA effectively facilitates the dissemination of scientific knowledge to the broader public, including via key platforms such as Wikipedia.

Our study provides a foundation for further research on Wikipedia and open science more broadly. Future studies should broaden source and variable coverage to better unpack the OA effect on Wikipedia. Furthermore, other forms of open science could be analyzed using a similar lens, for example, open research data and software. In conclusion, this study sheds light on the significance of OA in Wikipedia and potentially its broader impact. We believe that our findings will serve as a starting point for further research and contribute to the understanding of the impact and dissemination of OA.

## 7    Code and data availability statement

The code to replicate our work is made available online: `https://github.com/alsowbdxa/Open_access_and_wikipedia`. The Wikipedia Citations dataset is openly available [8], while access to OpenAlex can be requested through their portal. All other supporting datasets we used are openly available and referenced from the Data and Methods section.

## 8    Appendix

Our Appendix comprises three subsections: Figures, Tables, and Regression results supplements. Although these results offer a more extensive understanding of our research, they do not constitute the principal outcomes. Therefore, we have relocated them to the appendix for further reference.

### 8.1    Figures

Presented below are two figures depicting the distribution of OA status and policies among the top 20 journals. We have discussed it in the results part. This observation underscores the significance of conducting an article-level analysis for a more comprehensive understanding of the subject matter.

Figure 7: Distribution of OA status by top 20 journals.



Figure 8: Distribution of OA policies by top 20 journals.

## 8.2 Tables

The quality of our stratified samples is demonstrated through the descriptive statistics provided in Table 7, 8, 9 and 10. Additionally, Table 11 presents a count of articles by concepts within our dataset. The regression results for formulas 1 and 2 for the entire sample are displayed in Table 12 and 13.

Table 7: Descriptive statistics for the articles cited in Wikipedia (first dataset).

|       | times_cited | num_references | article_age | H index | is_oa   | $\ln(1 + SJR)$ |
|-------|-------------|----------------|-------------|---------|---------|----------------|
| count | 678,866     | 678,866        | 678,866     | 678,866 | 678,866 | 678,866        |
| mean  | 172.445     | 32.523         | 253.774     | 238.050 | 0.394   | 1.097          |
| std   | 792.548     | 45.811         | 196.048     | 273.388 | 0.489   | 0.737          |
| min   | 0           | 0              | 12          | 1       | 0       | 0.095          |
| 25%   | 14          | 7              | 130         | 80      | 0       | 0.556          |
| 50%   | 47          | 23             | 203         | 148     | 0       | 0.922          |
| 75%   | 134         | 42             | 306         | 271     | 1       | 1.466          |
| max   | 239,182     | 2,471          | 1,475       | 1,276   | 1       | 3.942          |

Table 8: Descriptive statistics for the articles not cited in Wikipedia. Average over all samples (first dataset).

|       | times_cited | num_references | article_age | H index | is_oa   | $\ln(1 + SJR)$ |
|-------|-------------|----------------|-------------|---------|---------|----------------|
| count | 678,866     | 678,866        | 678,866     | 678,866 | 678,866 | 678,866        |
| mean  | 110.742     | 28.405         | 253.803     | 238.050 | 0.364   | 1.097          |
| std   | 360.264     | 38.585         | 196.078     | 273.388 | 0.481   | 0.737          |
| min   | 0           | 0              | 12          | 1       | 0       | 0.095          |
| 25%   | 8           | 4              | 130         | 80      | 0       | 0.556          |
| 50%   | 38          | 21             | 203         | 148     | 0       | 0.922          |
| 75%   | 114         | 39             | 306         | 271     | 1       | 1.466          |
| max   | 81,780.80   | 3,100.40       | 1,475       | 1,276   | 1       | 3.942          |

Table 9: Descriptive statistics for the articles cited in Wikipedia (second dataset).

|  | times_cited | num_references | article_age | H index | is_oa | ln(1 + SJR) |
|---|---|---|---|---|---|---|
| n | 211,825 | 211,825 | 211,825 | 211,825 | 211,825 | 211,825 |
| mean | 165.813 | 34.901 | 265.953 | 275.844 | 0.456 | 1.249 |
| min | 0 | 0 | 13 | 1 | 0 | 0.095 |
| 25% | 16 | 9 | 141 | 92 | 0 | 0.647 |
| 50% | 51 | 26 | 216 | 177 | 0 | 1.061 |
| 75% | 138 | 45 | 313 | 331 | 1 | 1.730 |
| max | 148,954 | 1,976 | 1,475 | 1,276 | 1 | 3.942 |

Table 10: Descriptive statistics for the articles not cited in Wikipedia. Average over all samples (second dataset).

|  | times_cited | num_references | article_age | H index | is_oa | ln(1 + SJR) |
|---|---|---|---|---|---|---|
| n | 211,825 | 211,825 | 211,825 | 211,825 | 211,825 | 211,825 |
| mean | 153.128 | 31.303 | 265.975 | 275.844 | 0.437 | 1.249 |
| min | 0 | 0 | 12 | 1 | 0 | 0.095 |
| 25% | 12 | 7 | 141 | 92 | 0 | 0.647 |
| 50% | 49 | 25 | 216 | 177 | 0 | 1.061 |
| 75% | 152 | 43 | 313 | 331 | 1 | 1.730 |
| max | 47,776 | 2,211.600 | 1,475 | 1,276 | 1 | 3.942 |

Table 11: Count of articles by concepts in the final combined dataset (second dataset).

| Concept | n |
|---|---|
| Biology | 235,110 |
| Medicine | 73,120 |
| Chemistry | 32,400 |
| Physics | 16,014 |
| Psychology | 14,714 |
| Geology | 11,774 |
| Mathematics | 10,606 |
| Computer science | 7,958 |
| Philosophy | 3,654 |
| Political science | 3,390 |
| History | 2,798 |
| Art | 2,744 |
| Economics | 2,332 |
| Materials science | 2,116 |
| Geography | 2,000 |
| Business | 1,120 |
| Sociology | 828 |
| Environmental science | 596 |
| Engineering | 376 |

Table 12: Coefficients for overall OA adoption. Results across all 5 samples, model 1, $R^2 = 0.001$.

| Sample | Feature | coef | odds_ratios | P>z |
|---|---|---|---|---|
| 1 | is_oa | 0.140 | 1.150 | 0.000 |
| 2 | is_oa | 0.139 | 1.149 | 0.000 |
| 3 | is_oa | 0.143 | 1.154 | 0.000 |
| 4 | is_oa | 0.143 | 1.154 | 0.000 |
| 5 | is_oa | 0.144 | 1.155 | 0.000 |

## 8.3 Supplementary regression results

This section provides an in-depth analysis of OA citation adoption for each OpenAlex concept. To achieve this, we developed 19 distinct regression models, each dedicated to analyzing the adoption of OA citation for a single concept. We use the second formulation for each model, with data pertaining solely to the corresponding concept being considered in each case.

To gain insight into the effect of OA adoption on each concept, we present the coefficients for the *is_oa* variable in Table 14 and the coefficients for the $\ln(1 + times\_cited)$ variable in Table 15.

13

Table 13: Coefficients for overall OA adoption. Results across all 5 samples, model 2, $R^2 = 0.01$.

| Sample | Feature | coef | odds_ratios | P>z |
|---|---|---|---|---|
| 1 | *is_oa* | 0.112 | 1.118 | 0.000 |
| 2 | *is_oa* | 0.111 | 1.117 | 0.000 |
| 3 | *is_oa* | 0.115 | 1.122 | 0.000 |
| 4 | *is_oa* | 0.115 | 1.122 | 0.000 |
| 5 | *is_oa* | 0.116 | 1.122 | 0.000 |
| 1 | $\ln(1 + \textit{times\_cited})$ | 0.162 | 1.176 | 0.000 |
| 2 | $\ln(1 + \textit{times\_cited})$ | 0.161 | 1.175 | 0.000 |
| 3 | $\ln(1 + \textit{times\_cited})$ | 0.162 | 1.176 | 0.000 |
| 4 | $\ln(1 + \textit{times\_cited})$ | 0.162 | 1.175 | 0.000 |
| 5 | $\ln(1 + \textit{times\_cited})$ | 0.162 | 1.176 | 0.000 |

Table 14 reveals that OA articles in most concepts possess a statistically significant (p < 0.05) positive OA Wikipedia citation advantage and have greater odds of being cited in Wikipedia than closed-access articles. The top five concepts in terms of OA adoption advantage are Physics, Computer science, Chemistry, Environmental science and Biology, implying that STEM-related concepts tend to receive more attention on Wikipedia. Notably, only Business displays an average negative effect, albeit not statistically significant.

For $\ln(1 + \textit{times\_cited})$) in each concept, we find that citation counts moderate the positive effect of the relationship between OA status and adoption effect for OA articles. OA articles for several OpenAlex concepts, including Philosophy, Psychology, Mathematics, Medicine, Biology, Computer science, Geology, Chemistry, and Physics, show lower odds on average of being cited in Wikipedia than closed-access articles. OA articles in Business and Engineering, on the other hand, exhibit higher odds of being cited in Wikipedia than closed-access articles, but these results do not achieve statistical significance.

Table 14: Coefficients for OA adoption by concept for all samples (*is_oa*).

| concept | min OR | max OR | OR mean | Highest P-value | Mean R^2 |
|---|---|---|---|---|---|
| Political science | 1.494 | 1.700 | 1.598 | 0.001 | 0.019 |
| Philosophy | 1.699 | 1.868 | 1.788 | 0.000 | 0.022 |
| Economics | 1.211 | 1.404 | 1.290 | 0.399 | 0.123 |
| Business | 0.821 | 1.095 | 0.913 | 0.770 | 0.051 |
| Psychology | 1.524 | 1.748 | 1.632 | 0.000 | 0.067 |
| Mathematics | 1.894 | 2.015 | 1.972 | 0.000 | 0.074 |
| Medicine | 1.326 | 1.397 | 1.347 | 0.000 | 0.019 |
| Biology | 2.197 | 2.239 | 2.215 | 0.000 | 0.006 |
| Computer science | 2.478 | 2.744 | 2.665 | 0.000 | 0.064 |
| Geology | 1.984 | 2.264 | 2.121 | 0.000 | 0.020 |
| Chemistry | 2.335 | 2.547 | 2.435 | 0.000 | 0.010 |
| Art | 1.264 | 1.534 | 1.379 | 0.051 | 0.008 |
| Sociology | 1.453 | 1.874 | 1.609 | 0.260 | 0.032 |
| Engineering | 1.012 | 1.297 | 1.153 | 0.971 | 0.022 |
| Geography | 1.437 | 1.590 | 1.525 | 0.012 | 0.016 |
| History | 1.539 | 1.917 | 1.790 | 0.001 | 0.019 |
| Materials science | 1.660 | 2.264 | 1.943 | 0.033 | 0.037 |
| Physics | 3.008 | 3.357 | 3.208 | 0.000 | 0.012 |
| Environmental science | 2.076 | 2.585 | 2.257 | 0.009 | 0.025 |

Table 15: Coefficients for OA adoption by concept for all samples ($\ln(1 + \textit{times\_cited})$).

| concept | min OR | max OR | OR mean | Highest P-value | Mean R^2 |
|---|---|---|---|---|---|
| Political science | 0.835 | 0.917 | 0.878 | 0.100 | 0.019 |
| Philosophy | 0.823 | 0.847 | 0.836 | 0.003 | 0.022 |
| Economics | 0.961 | 0.994 | 0.979 | 0.899 | 0.123 |
| Business | 1.018 | 1.108 | 1.074 | 0.791 | 0.051 |
| Psychology | 0.913 | 0.944 | 0.924 | 0.010 | 0.067 |
| Mathematics | 0.856 | 0.885 | 0.866 | 0.000 | 0.074 |
| Medicine | 0.914 | 0.927 | 0.923 | 0.000 | 0.019 |
| Biology | 0.832 | 0.836 | 0.834 | 0.000 | 0.006 |
| Computer science | 0.854 | 0.891 | 0.865 | 0.000 | 0.064 |
| Geology | 0.817 | 0.855 | 0.837 | 0.000 | 0.020 |
| Chemistry | 0.816 | 0.835 | 0.826 | 0.000 | 0.010 |
| Art | 0.810 | 0.894 | 0.843 | 0.223 | 0.008 |
| Sociology | 0.820 | 0.934 | 0.867 | 0.556 | 0.032 |
| Engineering | 1.113 | 1.266 | 1.165 | 0.379 | 0.022 |
| Geography | 0.917 | 0.952 | 0.934 | 0.399 | 0.016 |
| History | 0.788 | 0.875 | 0.842 | 0.105 | 0.019 |
| Materials science | 0.892 | 0.962 | 0.918 | 0.483 | 0.037 |
| Physics | 0.783 | 0.804 | 0.791 | 0.000 | 0.012 |
| Environmental science | 0.765 | 0.844 | 0.805 | 0.071 | 0.025 |

# References

[1] Heather Piwowar, Jason Priem, Vincent Larivière, Juan Pablo Alperin, Lisa Matthias, Bree Norlander, Ashley Farley, Jevin West, and Stefanie Haustein. The state of OA: a large-scale analysis of the prevalence and impact of Open Access articles. *PeerJ*, 6:e4375, February 2018. Publisher: PeerJ Inc.

[2] Kim Holmberg, Juha Hedman, Timothy D. Bowman, Fereshteh Didegah, and Mikael Laakso. Do articles in open access journals have more frequent altmetric activity than articles in subscription-based journals? An investigation of the research output of Finnish universities. *Scientometrics*, 122(1):645–659, January 2020.

[3] Philipp Singer, Florian Lemmerich, Robert West, Leila Zia, Ellery Wulczyn, Markus Strohmaier, and Jure Leskovec. Why We Read Wikipedia. In *Proceedings of the 26th International Conference on World Wide Web*, WWW '17, pages 1591–1600, Republic and Canton of Geneva, CHE, April 2017. International World Wide Web Conferences Steering Committee.

[4] Taemin Kim Park. The visibility of Wikipedia in scholarly publications. 2011. Accepted: 2017-10-20T15:22:24Z Publisher: First Monday.

[5] Kayvan Kousha and Mike Thelwall. Are Wikipedia Citations Important Evidence of the Impact of Scholarly Articles and Books? *Journal of the Association for Information Science and Technology*, 68, November 2015.

[6] Fariba Tohidinasab and Hamid R. Jamali. Why and where Wikipedia is cited in journal articles? *Journal of Scientmetric Research*, 2:231–238, September 2013.

[7] Misha Teplitskiy, Grace Lu, and Eamon Duede. Amplifying the impact of open access: Wikipedia and the diffusion of science. *Journal of the Association for Information Science and Technology*, 68(9):2116–2127, 2017. _eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1002/asi.23687.

[8] Harshdeep Singh, Robert West, and Giovanni Colavizza. Wikipedia citations: A comprehensive data set of citations with identifiers extracted from english wikipedia. *Quantitative Science Studies*, 2(1):1–19, 2021.

[9] Jonathan P. Tennant, François Waldner, Damien C. Jacques, Paola Masuzzo, Lauren B. Collister, and Chris H. J. Hartgerink. The academic, economic and societal impacts of Open Access: an evidence-based review. Technical Report 5:632, F1000Research, September 2016. Type: article.

[10] LATINDEX Redalyc, REVENCIT Clase, and SERBILUZ IN-COM UAB. Berlin declaration on open access to knowledge in the sciences and humanities, 2003.

[11] Bo-Christer Björk, Patrik Welling, Mikael Laakso, Peter Majlender, Turid Hedlund, and Guðni Guðnason. Open Access to the Scientific Journal Literature: Situation 2009. *PLOS ONE*, 5(6):e11273, June 2010. Publisher: Public Library of Science.

[12] Alberto Martín-Martín, Rodrigo Costas, Thed van Leeuwen, and Emilio Delgado López-Cózar. Evidence of open access of scientific publications in Google Scholar: A large-scale analysis. *Journal of Informetrics*, 12(3):819–841, August 2018.

[13] Éric Archambault, Didier Amyot, Philippe Deschamps, Aurore Nicol, Françoise Provencher, Lise Rebout, and Guillaume Roberge. Proportion of Open Access Papers Published in Peer-Reviewed Journals at the European and World Levels—1996–2013. *Copyright, Fair Use, Scholarly Communication, etc.*, October 2014.

[14] Erik W. Black. Wikipedia and academic peer review: Wikipedia as a recognised medium for scholarly publication? *Online Information Review*, 32(1):73–88, February 2008.

[15] Cassidy R. Sugimoto, Sam Work, Vincent Larivière, and Stefanie Haustein. Scholarly use of social media and altmetrics: A review of the literature. *Journal of the Association for Information Science and Technology*, 68(9):2037–2062, 2017. _eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1002/asi.23833.

[16] Wenceslao Arroyo-Machado, Daniel Torres-Salinas, Enrique Herrera-Viedma, and Esteban Romero-Frías. Science through Wikipedia: A novel representation of open knowledge through co-citation networks. *PloS one*, 15(2):e0228713, 2020.

[17] Puyu Yang and Giovanni Colavizza. A map of science in wikipedia. In *Proceedings of The Web Conference*, 2022.

[18] Lauren A. Maggio, John M. Willinsky, Ryan M. Steinberg, Daniel Mietchen, Joseph L. Wass, and Ting Dong. Wikipedia as a gateway to biomedical research: The relative distribution and use of citations in the English Wikipedia. *PLOS ONE*, 12(12):e0190046, December 2017. Publisher: Public Library of Science.

[19] N. J. Schweitzer. Wikipedia and Psychology: Coverage of Concepts and Its Use by Undergraduate Students. *Teaching of Psychology*, 35(2):81–85, April 2008. Publisher: SAGE Publications Inc.

[20] Finn Aarup Nielsen. Scientific citations in Wikipedia, May 2007. arXiv:0705.2106 [cs].

[21] Neil Thompson and Douglas Hanley. Science Is Shaped by Wikipedia: Evidence From a Randomized Control Trial, February 2018.

[22] Giovanni Colavizza. COVID-19 research in Wikipedia. *Quantitative Science Studies*, 1(4):1349–1380, December 2020.

[23] Aaron Halfaker, Bahodir Mansurov, Miriam Redi, and Dario Taraborelli. Citations with identifiers in Wikipedia, 2018.

[24] Olga Zagorova, Roberto Ulloa, Katrin Weller, and Fabian Flöck. "I updated the <ref>": The evolution of references in the English Wikipedia and the implications for altmetrics. *Quantitative Science Studies*, pages 1–27, 12 2021.

[25] Omer Benjakob, Rona Aviram, and Jonathan Aryeh Sobel. Citation needed? Wikipedia bibliometrics during the first wave of the COVID-19 pandemic. *GigaScience*, 11:giab095, January 2022.

[26] Tiziano Piccardi, Miriam Redi, Giovanni Colavizza, and Robert West. Quantifying Engagement with Citations on Wikipedia. In *Proceedings of The Web Conference 2020*, WWW '20, pages 2365–2376, New York, NY, USA, April 2020. Association for Computing Machinery.

[27] Jason Priem, Heather Piwowar, and Richard Orr. OpenAlex: A fully-open index of scholarly works, authors, venues, institutions, and concepts, June 2022. arXiv:2205.01833 [cs].

[28] Laura Bredahl. Chapter 1. Introduction to Bibliometrics and Current Data Sources. *Library Technology Reports*, 58(8):5–11, November 2022. Number: 8.

[29] Hongtao Hao, Yumian Cui, Zhengxiang Wang, and Yea-Seul Kim. Thirty-Two Years of IEEE VIS: Authors, Fields of Study and Citations, August 2022. arXiv:2208.03772 [cs].

[30] Matthew E. Falagas, Vasilios D. Kouranos, Ricardo Arencibia-Jorge, and Drosos E. Karageorgopoulos. Comparison of SCImago journal rank indicator with journal impact factor. *The FASEB Journal*, 22(8):2623–2628, 2008. _eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1096/fj.08-107938.

[31] Jason Yuen. Comparison of Impact Factor, Eigenfactor Metrics, and SCImago Journal Rank Indicator and h-index for Neurosurgical and Spinal Surgical Journals. *World Neurosurgery*, 119:e328–e337, November 2018.

[32] Borja González-Pereira, Vicente P. Guerrero-Bote, and Félix Moya-Anegón. A new approach to the metric of journals' scientific prestige: The SJR indicator. *Journal of Informetrics*, 4(3):379–391, July 2010.

[33] Giovanni Colavizza, Iain Hrynaszkiewicz, Isla Staden, Kirstie Whitaker, and Barbara McGillivray. The citation advantage of linking publications to research data. *PLOS ONE*, 15(4):e0230416, April 2020. Publisher: Public Library of Science.

[34] Yassine Gargouri, Chawki Hajjem, Vincent Larivière, Yves Gingras, Les Carr, Tim Brody, and Stevan Harnad. Self-Selected or Mandated, Open Access Increases Citation Impact for Higher Quality Research. *PLOS ONE*, 5(10):e13636, October 2010. Publisher: Public Library of Science.

[35] Alfredo Yegros-Yegros, Ismael Rafols, and Pablo D'Este. Does Interdisciplinary Research Lead to Higher Citation Impact? The Different Effect of Proximal and Distal Interdisciplinarity. *PLOS ONE*, 10(8):e0135095, August 2015. Publisher: Public Library of Science.

[36] D. B. Struck, M. Durning, G. Roberge, and D. Campbell. Modelling the Effects of Open Access, Gender and Collaboration on Citation Outcomes: Replicating, Expanding and Drilling. *STI 2018 Conference Proceedings*, pages 436–447, September 2018. Publisher: Centre for Science and Technology Studies (CWTS).