

A Comparative Study of Reference Reliability in Multiple Language Editions of Wikipedia

Aitolkyn Baigutanova
KAIST, IBS
Daejeon, South Korea
aitolkyn.b@kaist.ac.kr

Diego Saez-Trumper
Wikimedia Foundation
Barcelona, Spain
diego@wikimedia.org

Miriam Redi
Wikimedia Foundation
London, United Kingdom
mredi@wikimedia.org

Meeyoung Cha
IBS, KAIST
Daejeon, South Korea
meeyoungcha@kaist.ac.kr

Pablo Aragón
Wikimedia Foundation
Barcelona, Spain
paragon@wikimedia.org

ABSTRACT

Information presented in Wikipedia articles must be attributable to reliable published sources in the form of references. This study examines over 5 million Wikipedia articles to assess the reliability of references in multiple language editions. We quantify the cross-lingual patterns of the *perennial sources list*, a collection of reliability labels for web domains identified and collaboratively agreed upon by Wikipedia editors. We discover that some sources (or web domains) deemed untrustworthy in one language (i.e., English) continue to appear in articles in other languages. This trend is especially evident with sources tailored for smaller communities. Furthermore, non-authoritative sources found in the English version of a page tend to persist in other language versions of that page. We finally present a case study on the Chinese, Russian, and Swedish Wikipedias to demonstrate a discrepancy in reference reliability across cultures. Our finding highlights future challenges in coordinating global knowledge on source reliability.

CCS CONCEPTS

• **Information systems** → *Wikis*.

KEYWORDS

Wikipedia, Verifiability, Information Credibility, Fake News, Misinformation

ACM Reference Format:

Aitolkyn Baigutanova, Diego Saez-Trumper, Miriam Redi, Meeyoung Cha, and Pablo Aragón. 2023. A Comparative Study of Reference Reliability in Multiple Language Editions of Wikipedia. In *Proceedings of Conference Acronym*. ACM, New York, NY, USA, 5 pages. <https://doi.org/10.1145/nnnnnnn>. nnnnnnn

1 INTRODUCTION

As a global knowledge encyclopedia, Wikipedia is one of the most popular websites worldwide. It is open access and maintained by the online community that collaboratively creates content following specific editing policies. A core content policy is verifiability [12], which requires that information included in Wikipedia articles is supported by reliable and relevant references from which claims are

derived. In order to encourage editors to avoid untrustworthy references, the Wikipedia community created a reliability index called the *perennial sources list* [13]. This list contains web domains and their labels (e.g., blacklisted, generally reliable), assigned based on a collective consensus through discussions among the editors [11]. Recent research has shown that these community efforts have led to improved reference quality in English Wikipedia [2].

A growing body of research explores the role of references in the encyclopedia. A study of sources cited in scientific articles has shown that Wikipedia relies heavily on prestigious journals from STEM fields, with sources from non-STEM disciplines being marginal but relevant in biographical content [14]. Recent research also found a moderate yet systematic liberal polarization in the selection of media and news sources [15]. While these studies are limited to English Wikipedia, there exist a few works that examined the popularity of sources across different languages of Wikipedia [5, 6], finding that some sources are common across editions but that each language often has its own distinct set of sources. However, the reliability of sources across a broader set of language editions of Wikipedia remains unexplored.

Understanding referencing behavior across languages is critical in multilingual ecosystems such as Wikipedia. The English Wikipedia edition is currently the largest in terms of community size and the number of articles, but it only covers a portion of the information found on other Wikipedias [3]. The impact of this dominant language edition on other editions remains unclear [9], and the cross-lingual effect on Wiki-related policies is still being discussed. One study has shown that even if large language editions share core editing rules, localized rule sets tend to become increasingly diverse over time [4].

Here, we examine the reliability of references in over 300 language editions of Wikipedia to study the cross-edition effects on reference quality: how do the English editor community's initiatives affect the referencing behavior of other language editions? We investigate the reliability of the most common sources and their presence in the same article across languages and find that non-authoritative references in the English edition tend to persist in other language versions of the same articles. We present our preliminary findings on the potential risks to content verifiability that may result from translating articles into less developed language editions.

Table 1: Descriptive statistics of *perennial source lists* in four language editions. Coverage represents the percentage of articles citing at least one source from the corresponding list.

Language edition	Number of domains	Overlap with the English list	Coverage	Coverage (English list)
English	1,156	100%	22.9%	22.9%
Chinese	211	17.5%	10.0%	11.5%
Swedish	51	54.9%	0.8%	2.4%
Russian	60	15.0%	5.7%	11.6%

While we hypothesized that a community-maintained list of untrustworthy sources in one language (i.e., English) can be a good starting point for improving reference quality in other smaller language editions, our data indicated that some domains deemed untrustworthy and hence banned in the English edition are still actively used in some language editions. One example is the Russian language online newspaper *lenta.ru*, which frequently appears in Russian Wikipedia. The daily tabloid *huanqui.com* which is under the auspices of the Chinese Communist Party, was another example that frequently appeared in Chinese Wikipedia. These web sources likely provide cultural- or regime-specific information and values that are not universally shared by all cultures. Wikipedia, however, is read by a global audience and serves as a foundation for many large language models and search engines. Therefore, our discovery of a discrepancy in reference reliability across cultures opens up a discussion about the need to account for the broad audience of this interconnected world when deciding “global knowledge”.

2 DATA

2.1 Reliability of Sources in Wikipedia

We use the internal labeling of source reliability decided by the Wikipedia community, called the *perennial sources list*. It comprises a collection of web domains in five categories: (1) blacklisted, (2) deprecated, (3) generally unreliable, (4) no consensus, and (5) generally reliable. Following the methodology of a recent study [2], we refer to the first two categories as non-authoritative sources. As of May 2023, a *perennial sources list* page exists in the following 12 language editions of Wikipedia: English, Russian, French, Persian, Swedish, Chinese, Portuguese, Greek, Lithuanian, Turkish, Vietnamese, and Nepali. Only the first six are actively maintained, with no updates made to the classification of the remaining lists this year. Furthermore, the Persian Wikipedia list includes only 14 web sources, while sources included in the French classification do not have an explicit reliability label.

For the remaining four lists in English, Chinese, Swedish, and Russian Wikipedia, we show descriptive statistics in Table 1, including the number of classified domains, the overlap of the sources in the list with English Wikipedia’s *perennial sources list*, and the coverage by the lists of the corresponding language edition and English edition. We define coverage as the percentage of pages in the corresponding language edition with at least one citation to the sources from a *perennial sources list*.

2.2 Dataset Description

We first retrieve the *perennial sources list* from English, Chinese, Swedish, and Russian editions using Python’s BeautifulSoup library to parse HTML documents. We use the online version of the lists from May 2023. Then, we manually filled in the missing data for the entries in the lists that did not include an explicit link to the source. This is the case for the Russian edition, where website domains are not listed for all the sources, but they can be inferred from the links to the discussion pages or the linked Wikipedia pages. As a result, we could obtain a table of sources along with their domain and category for each of the four editions. Only sources with explicit labels were included for the remainder of the analysis.

We used the following attributes from Wikidata dumps: item ID (a language-agnostic Wikidata page identifier), language edition ID, page ID in a given language edition, and the corresponding page title. We collect 2,182,103 unique Wikidata items across all language editions. We consider the pages in the article namespace of Wikipedia. We combine this data with the XML dumps for the most recent Wikipedia page versions (as of February 2023) to retrieve the raw text of each article. Using this text, we enrich the dataset with sources included in the English *perennial sources list* cited in a given article, along with the source’s category. This results in 5,189,606 articles across 314 editions that include at least one reference to a *perennial source*.

3 RESULTS

3.1 Reliability by Language Editions

To investigate the spread of English Wikipedia’s *perennial sources* across multiple language editions, we identify the proportion of articles in each edition that include at least one reference to these sources. Figure 1 shows the percentage of articles referencing reliable and non-authoritative sources in the 40 editions with the largest number of articles. Languages can be mapped from the Wikipedia code [10].

The plot shows outliers in the two directions of the confidence interval represented by the gray area. On the one hand, the English edition is located below the confidence interval, meaning the proportion of articles citing reliable domains is larger. This observation is consistent with recent research [2], as the community of English Wikipedia is more aware of the non-authoritative domains listed in the local *perennial sources list*. On the other hand, the outliers above the confidence interval appear to have a relatively larger proportion of articles citing deprecated or blacklisted domains. These are Russian (ru), Armenian (hy), Chinese (zh), French (fr), and Bulgarian (bg). We qualitatively assessed what non-authoritative sources are prevalent in these five editions. We found that for Russian and Armenian, the trend is primarily caused by the blacklisted domain *lenta.ru*, a Russian online newspaper. For the French edition, the US-based source *city-data.com* contributes the largest proportion. For the Chinese language edition, the trend is predominantly attributed to the Chinese tabloid *huanqiu.com*, and international media source *epochtimes.com*. Finally, the deprecated British tabloid *daily-mail.co.uk* contributes a significant proportion of all the editions.

We further examine specific domains that appear most frequently across multiple editions of Wikipedia. The blue dashed line in Figure 2 presents the number of editions for each category’s top ten

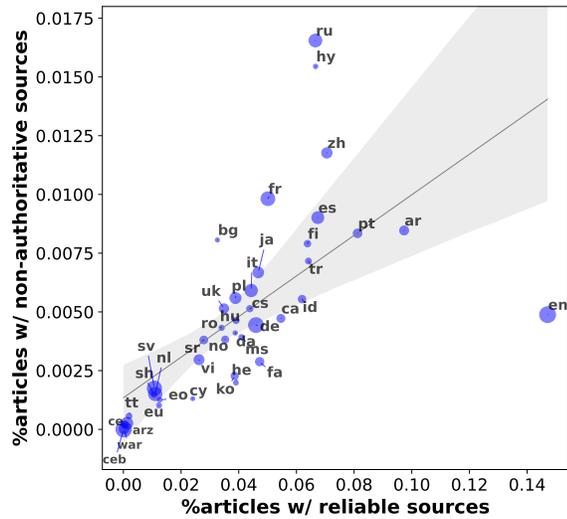


Figure 1: The proportion of articles referencing reliable and non-authoritative sources in 40 language editions. The circle size indicates the number of articles, while the shaded area indicates a 90% confidence interval.

most widely referenced domains. The most popular web domains across different languages are generally reliable *nytimes.com* and *bbc.co.uk*, generally unreliable *wikipedia.org*, and no consensus *britannica.com*, which are referenced in around 270 editions. We note that the most frequently used non-authoritative domain, namely *dailymail.co.uk*, is ranked 34th, cited in 181 language editions. This indicates that although deprecated and blacklisted sources are less popular, they are still commonly cited.

3.2 Reference Similarity with English Edition

Next, we explore the relationship between references in English Wikipedia and the other language editions. We measure how likely sources are to appear in a non-English version of an article when they are also present in the English one. That is, we analyze pages present in at least two editions of Wikipedia, where one is in English. First, we gather a collection of item IDs from each edition, referencing a specific source. Then, we compute the pairwise Jaccard similarity coefficient of the set of items citing a given source in English and other language editions. This coefficient is a commonly used statistic to determine the similarity between two sets and is computed as a ratio of two sets' intersection over their union.

Sources in the deprecated and blacklisted categories of the *perennial sources list* demonstrated a higher article set similarity with the English edition than generally reliable ($t=6.08, p<0.001$) and no consensus and generally unreliable ($t=6.16, p<0.001$) ones. In contrast, the generally reliable category also demonstrated a lower value on average than the middle two categories ($t=2.62, p<0.01$). Figure 2 displays the average article set similarity of non-English editions with the English Wikipedia per source. The scores are shown for each category's ten most commonly referenced domains in the *perennial sources list*. We observe that the highest similarity is exhibited by non-authoritative sources, such as *filmreference.com*,

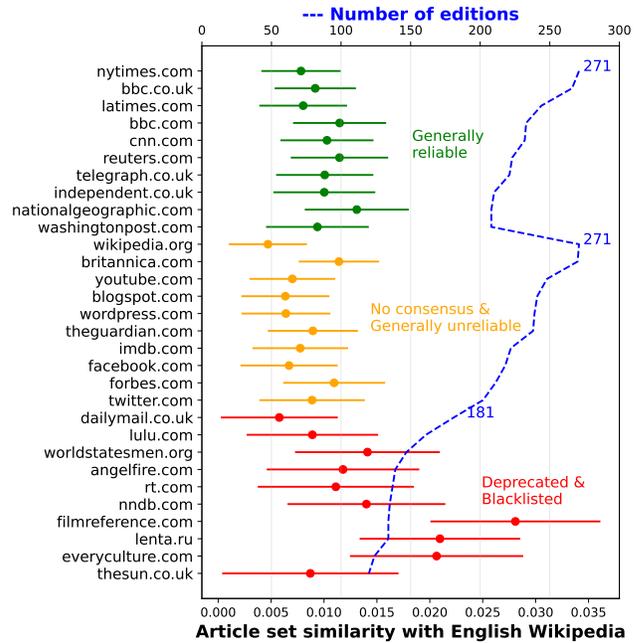


Figure 2: Average pairwise similarity of article sets citing the same source with English edition. The top ten most popular sources from each category are shown. The error bars indicate the standard error.

lenta.ru, and *everyculture.com*. These findings may indicate that non-authoritative domains found in the English edition are more likely to be present in other language editions.

There may be various reasons why references in one language may reappear in another. One such method is translation. The proportion of articles created through translation among all pages that currently reference at least one non-authoritative domain is found among the top 50 editions with the most articles. While less than 5% of these pages were created via translation in large language editions (e.g., Russian, Chinese, German, Italian, and Japanese), a significant proportion of pages with non-authoritative sources were created in this way in editions such as Uzbek (73%) or Hebrew (58%). These results suggest an important risk to content verifiability when translating articles into less developed Wikipedia language editions, as they have fewer resources to assess whether the original content included references to non-authoritative sources.

3.3 Local Source Reliability Labeling

Finally, we compare the impact of internal initiatives that maintain reference reliability in smaller editions to those implemented on English Wikipedia. In particular, we analyze the *perennial sources lists* from three language editions: Chinese, Swedish, and Russian. These lists are created by local editors through discussion and agreement, just like the English Wikipedia, making them comparable. Domains from the English list are frequently cited in articles written in all three languages. In the Russian edition, as shown in Table 1, domains in the English *perennial sources list* (11.6%) are cited in

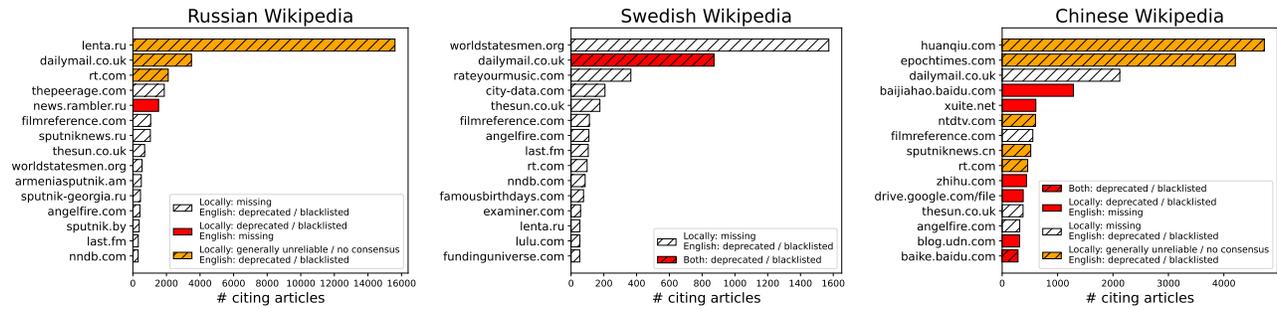


Figure 3: Top 15 non-authoritative sources (from the *perennial source* list of the local Wikipedia edition or the one of English Wikipedia) by the number of citations in Russian, Swedish, and Chinese Wikipedia editions.

more than twice as many pages as in the Russian *perennial sources list* (5.7%).

Next, we investigate the non-authoritative domain category of *perennial sources* for the three languages. Figure 3 shows the top 15 most frequently appearing deprecated and blacklisted sources from the *perennial source list* of the local language or English editions. A few domains are only included in the local lists (solid red bars). Some examples are *news.rambler.ru* in Russian Wikipedia and *xuite.com* in Chinese Wikipedia. Moreover, it is common that domains marked as deprecated and blacklisted in English Wikipedia are missing in the local list (white dashed bars) or have not reached consensus in the Russian and Chinese editions (orange dashed bars).

4 CONCLUSION

Implications We presented a comparative study of reference reliability on Wikipedia across multiple language editions to understand the status quo of knowledge integrity [1]. We started from a recent analysis of the community initiative labeling domain reliability (called *perennial sources list* in English Wikipedia) [2] and examined their prevalence in other language editions.

Our finding that trustworthy sources considered by English editors are also more frequently cited in other language editions suggests that overall reliability is well maintained. On the other hand, non-authoritative domains that somehow persisted in the English edition also appeared in the same articles in other languages, indicating that reliability risks in a major language edition will percolate to smaller language editions. Because smaller language editions lack the same level of human resources to maintain each page, these risks may have spread through translation.

These findings suggest the potential that language editions could co-share and co-develop the *perennial sources list*. This process will likely require global coordination among editors of various language editions, as well as agreement on what constitutes appropriate "global knowledge." Even in our case study, we found that cultural- or regime-specific domains were actively used in Wikipedia's local language edition, whereas the same domains were deemed unreliable in other language editions. One of the examples is a daily tabloid, *huanqiu.com* that supports the Chinese Communist Party. This domain frequently appears in the Chinese Wikipedia edition, yet it is labeled as untrustworthy in the English edition. Given the disparity in what constitutes reliable sources

across language editions, coordination among Wikipedia editors on deciding "global knowledge" may be challenging. Nonetheless, policy decisions are needed because Wikipedia content has become a key database for search searches and training large language models.

There could be a possibility of relying on external domain ratings (for example, see [15]) rather than maintaining a Wiki-specific list. According to one study, external reliability ratings have a high level of agreement among experts [7]. However, studies also acknowledge the challenges of handling culturally specific biases. For example, biases in Wikipedia citations to scholarly publications have been discovered, favoring authors affiliated with Anglosphere countries [16]. Because Wikipedia source reliability labels must result from a deliberative process among community members, discussions must include smaller language editions, particularly for local sources unique to that language.

Future Work We encountered several limitations that should encourage future work. First, we presented an analysis based on the most recent page versions without considering the articles' editing history. In Wikipedia, not only articles but also its rules evolve over time, including the decision to deprecate sources, as happened with the British tabloid *dailymail.co.uk* after an intense debate [8]. Further research could look into the temporal patterns of managing untrustworthy sources across languages, similar to previous work on reference quality in English Wikipedia [2]. Second, our analysis only considered the *perennial sources list* in English Wikipedia. This is because comparable lists in other language editions are much less developed. Some editions had a higher percentage of articles referencing deprecated or blacklisted domains. Therefore, we emphasize the need for more active discussions within non-English communities, and we hope that our work inspires local initiatives to improve the reference quality of the content.

Acknowledgements We thank the Wikipedia volunteers who contribute to content verifiability. This research was supported by the Institute for Basic Science (IBS-R029-C2) and the National Research Foundation of Korea (NRF) grant (RS-2022-00165347).

REFERENCES

- [1] Pablo Aragón and Diego Sáez-Trumper. 2021. A preliminary approach to knowledge integrity risk assessment in Wikipedia projects. *arXiv preprint arXiv:2106.15940* (2021).
- [2] Aitolkyn Baigutanova, Jaehyeon Myung, Diego Saez-Trumper, Ai-Jou Chou, Miriam Redi, Changwook Jung, and Meeyoung Cha. 2023. Longitudinal Assessment of Reference Quality on Wikipedia. In *Proceedings of the ACM Web Conference 2023* (Austin, TX, USA) (*WWW '23*). Association for Computing Machinery, New York, NY, USA, 2831–2839. <https://doi.org/10.1145/3543507.3583218>
- [3] Brent Jaron Hecht. 2013. *The mining and application of diverse cultural perspectives in user-generated content*. Ph. D. Dissertation. Northwestern University.
- [4] Sohyeon Hwang and Aaron Shaw. 2022. Rules and Rule-Making in the Five Largest Wikipedias. *Proceedings of the International AAAI Conference on Web and Social Media* 16, 1 (May 2022), 347–357. <https://doi.org/10.1609/icwsm.v16i1.19297>
- [5] Włodzimierz Lewoniewski, Krzysztof Węcel, and Witold Abramowicz. 2017. Analysis of References Across Wikipedia Languages. In *Information and Software Technologies*, Robertas Damaševičius and Vilma Mikašytė (Eds.). Springer International Publishing, Cham, 561–573.
- [6] Włodzimierz Lewoniewski, Krzysztof Węcel, and Witold Abramowicz. 2023. Companies in Multilingual Wikipedia: Articles Quality and Important Sources of Information. In *Information Technology for Management: Approaches to Improving Business and Society*, Ewa Ziemba, Witold Chmielarz, and Jarosław Wątróbski (Eds.). Springer Nature Switzerland, Cham, 48–67.
- [7] Hause Lin, Jana Lasser, Stephan Lewandowsky, Rocky Cole, Andrew Gully, David Rand, and Gordon Pennycook. 2022. High level of agreement across different news domain quality ratings. (2022).
- [8] Svrrir Steinsson. 2023. Rule Ambiguity, Institutional Clashes, and Population Loss: How Wikipedia Became the Last Good Place on the Internet. *American Political Science Review* (2023), 1–17.
- [9] Rodolfo Vieira Valentim, Giovanni Comarela, Souneil Park, and Diego Sáez-Trumper. 2021. Tracking knowledge propagation across wikipedia languages. In *Proceedings of the International AAAI Conference on Web and Social Media*, Vol. 15. 1046–1052.
- [10] English Wikipedia. 2023. Retrieved May 26, 2023 from ListofWikipedias
- [11] English Wikipedia. 2023. Reliable sources Noticeboard. Retrieved May 26, 2023 from https://en.wikipedia.org/wiki/Wikipedia:Reliable_sources/Noticeboard
- [12] English Wikipedia. 2023. Verifiability Policy. Retrieved May 26, 2023 from <https://en.wikipedia.org/wiki/Wikipedia:Verifiability>
- [13] English Wikipedia. 2023. Wikipedia: Reliable source, Perennial Source. Retrieved May 26, 2023 from https://en.wikipedia.org/wiki/Wikipedia:Reliable_sources/Perennial_sources
- [14] Puyu Yang and Giovanni Colavizza. 2022. A Map of Science in Wikipedia. In *Companion Proceedings of the Web Conference 2022*. 1289–1300.
- [15] Puyu Yang and Giovanni Colavizza. 2022. Polarization and reliability of news sources in Wikipedia. *arXiv preprint arXiv:2210.16065* (2022).
- [16] Xiang Zheng, Jiajing Chen, Erjia Yan, and Chaoqun Ni. 2023. Gender and country biases in Wikipedia citations to scholarly publications. *Journal of the Association for Information Science and Technology* 74, 2 (2023), 219–233.