



# Types of Errors Hiding in Google Scholar Data

Romy Sauvayre

## ► To cite this version:

Romy Sauvayre. Types of Errors Hiding in Google Scholar Data. Journal of Medical Internet Research, In press. hal-03609870

**HAL Id: hal-03609870**

**<https://hal.science/hal-03609870>**

Submitted on 6 May 2022

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Sauvayre R., (in press) "What types of errors are hiding in Google Scholar data? A case study", *Journal of Medical Internet Research*.

## Viewpoint

# What types of errors are hiding in Google Scholar data? A case study

**Romy Sauvayre**

Université Clermont Auvergne, CNRS, LAPSCO, F-63000 Clermont-Ferrand, France.

Email: [romy.sauvayre@uca.fr](mailto:romy.sauvayre@uca.fr)

## Abstract

Google Scholar (GS) is a free tool that may be used by researchers to analyze citations, to find appropriate literature or to evaluate the quality of an author or a contender for tenure, promotion, a faculty position, funding or research grants. GS has become a major bibliographic and citation database. Following the literature, databases such as PubMed, PsycINFO, Scopus or Web of Science can be used in place of GS because they are more reliable. The aim of this study is to examine the accuracy of citation data collected from GS and provide a comprehensive description of the errors and miscounts identified. For this purpose, 281 documents that cited two specific works were retrieved via the Publish or Perish (PoP) software and examined. This work studied the false positive issue inherent in the analysis of neuroimaging data. The results reveal an unprecedented error rate: 279 of 281 the examined references (99.3%) contain at least one error. The nonacademic documents tend to contain more errors than the academic publications ( $U=5117.0$ ,  $P<.001$ ). This viewpoint article, based on a case study examining GS data accuracy, shows that GS data not only fail to be accurate but also potentially expose researchers who would use these data without verification to substantial biases in their analyses and results. Further work must be conducted to access the consequences of using GS data extracted by PoP.

**Keywords:** Reference accuracy; database reliability; false positive; academic publication; research evaluation; scientometrics; citation analysis.

## Introduction

Google Scholar (GS) has become a major bibliographic and citation database. Soon after its creation in 2004, GS received major criticism (see Orduna-Malea et al. [1] review), but subsequently, further studies described it more positively [2,3]. Indeed, the literature acknowledges the free access offered by GS [3-5] and the quality of its coverage [6-12]. The coverage of GS is considered better than that of both Web of Science (WoS) [12-15] and Scopus [9,10], which are GS's fee-based competitors. This is particularly true regarding its coverage of social sciences and humanities research [10,16,17], conference proceedings [10,14] and books [17]. The GS database has been substantially qualitatively [18] and quantitatively [10,19] improved in all scientific areas such that, according to De Winter et al. [18], it could supplant WoS.

However, "the automatic indexing of GS inevitably causes many errors" [20], such as duplicates [21] and false positive citations [18]. Most researchers generally claim that these errors are negligible [9,10,12,20,22-24], whereas others consider that data cleaning is necessary [16,19,25] but laborious [4,21]. Thus, some scholars have used GS without data cleaning [2,6,11,14,26,27], while others have identified and removed duplicates [4,9,12,17-19,24,28,29]. This removal was performed in 23 of 36 studies (41.8%) using GS data. **Erreur! Source du renvoi introuvable..** Furthermore, compared to the authors of related studies, these researchers less frequently identified false positives [17,18,30,31], missing values or omission errors [20,23], document type errors [18,32], author name errors [18,33], publication year errors [16,18,33], title errors [18,33], URL errors [16,32], citation miscounts [32], and inaccessible document errors [30]. None of these 36 studies mention any verification of journal names in their data cleaning process. Nevertheless, Haddaway et al. [4] attempted to explain the causes of duplicates, showing that they arise from typographical and capitalization errors occurring in journal names. Their findings were confirmed by a study conducted by Valderrama et al. [34] based on Scopus data.

However, the analysis of 36 articles published between 2008 and 2018 in journals with an impact factor from Journal Citation Reports (JCR) collected from WoS, GS and the relevant studies cited in the most cited research in this field shows that the data verification was not systematically followed by the calculation and the reporting of an error rate. Indeed, 14 studies of the 36 (38.9%) explicitly indicate a number or a rate of errors. A median error rate of 14.6% with a range from 0.04% to 53.5%, among corpora of citations ranging in size between 127 and 183,596 was calculated. Note that for those studies that were missing error rates but nevertheless had reported adequate results, the error rates were calculated and included. In addition, these studies report error data of a median of only 1 type of error (range: 0 to 6) and duplicates are the error types most frequently searched for in this sample of literature (23 of 36).

This median error rate therefore demonstrates that errors are recurrent in GS data. However, GS is a free tool that may be used by researchers to analyze citations, to find appropriate literature [35, 36] or to evaluate the quality or the influence [37] of an author, or a contender for tenure, promotion, a faculty position, funding or research grants [1,21]. Thus, the more an author is cited in a field, the more likely that person is to be considered a highly qualified researcher [38,39]. GS may also be used in research evaluations [23]. Thus, a comprehensive study of this failure of GS may be useful to the scientific community and researchers who wanted to use this database, whatever their field of study. However, as far as can be seen, no study reports and meticulously quantifies the different types of errors encountered in the GS data extracted by Publish or Perish (PoP) software, even though such a study would allow:

- better identification of the limitations of studies based on these data, as described by Hicks et al. [40] in the context of research evaluation;
- enrichment of the thoughtful methodological reflection on potential exposure to GS errors; and

Sauvayre R., (in press) "What types of errors are hiding in Google Scholar data? A case study", *Journal of Medical Internet Research*.

- development of appropriate methods to limit the negative effects of GS errors on the results produced.

The present case study aimed to examine the GS data extracted by PoP, provide a full count of the errors contained in these collected data, and present an epistemological reflection. By doing so, this study offers detailed categorizations of GS data that have not been provided by previous studies. The purpose is especially to address the following questions:

- 1) What types of GS errors could affect the data and results of researchers' studies?
- 2) What methodological problems may result from these errors?
- 3) How reliable can the citations of GS be without data cleaning?

## Methods

### Context

This GS study is part of broader research that aims to explore the diffusion process of a neuroimaging work that sought to alert the scientific community to the issue of false positives. Two references were examined: the first reference is a poster presented at the 15th Annual Meeting for the Organization for Human Brain Mapping (OHBM) [41] and the second is an article published in the short-lived *Journal of Serendipitous and Unexpected Results* (JSUR) [42]. The question was which researchers contributed to this diffusion or, in other words, who cited the OHBM poster or the JSUR article. The collection of those citation data then became necessary. Nevertheless, some full texts of the citing documents collected by GS did not cite either the OHBM poster or the JSUR article. From this, the present case study was born: the reliability of the GS data needed to be quantified to identify the limitations of the results produced with GS data before using these data in the diffusion study. This categorization of errors using these two references enables one to identify how GS works with literature not referenced by journal editors' websites. Google Scholar uses "automated software, known as 'parsers', to identify bibliographic data" [49] of documents available on the Internet. Then, the parser software "typically" collects the same data from full documents without metadata as the two used in the present case study.

### Data collection

To examine the reliability and accuracy of GS, the citations of both the OHBM poster and the JSUR article were analyzed. Note that GS was the only citation database available to collect the citation data for these two works because neither WoS nor Scopus have ever indexed them.

PoP version 5 was used to extract references that cited the poster. According to Harzing [43], this software provides a perfect collection of GS data: "Publish or Perish is as accurate or as inaccurate as Google Scholar itself." In addition, PoP is a common tool in scientometric studies using GS data [14,29,19,44].

The citation data were then collected from GS via PoP: the first author's name ("Bennett, Craig M") was entered without quotation marks, and the first part of the OHBM poster and JSUR article title ("Neural correlates of interspecies perspective taking in the post-mortem Atlantic Salmon") were entered with quotation marks into the "All of the words" software query box. As PoP's manual explains, the "All of the words" query "matches the search terms anywhere in the searched documents (author, title, source, abstract, references etc.)" [45], as GS does. Thus, this query was used to reproduce the same request with PoP and GS.

This title is so specific that only two results appeared:

- 127 references cited the JSUR article [42], and
- 154 references cited the OHBM poster [41]. In contrast, the reference that appeared in PoP and on GS was a paper supposedly published in a supplement of the famous *NeuroImage* journal and indexed by ScienceDirect. In reality, *NeuroImage* did not publish a journal article written by Bennett and colleagues in 2009 about the neuroimaging work. Bennett's *NeuroImage* paper does not exist. What this supplemental issue of *NeuroImage* does contain is the program of the OHBM conference. Therefore, when the citing documents cite the "*NeuroImage* ghost paper", they actually cite the OHBM poster.

Note that the JSUR article title is almost identical to the OHBM poster title—only the term "proper" in the second part of the title differs. The advantage of this strong similarity is the ability to evaluate the capacity of GS to manage citations of similar references.

A total of 281 references were extracted via PoP on October 6, 2017. Two CSV files were obtained (Multimedia Appendix 1)—one for each neuroimaging reference. In the present study, several columns that contained the following information were examined: authors, title of the citing document, publication year, publication or source, publisher, and web address of the citing document ("Article URL" as provided by GS). Each column was manually verified, and inaccuracies were counted and categorized in the following six steps:

- 1) The full text of accessible citing documents was downloaded and recorded;
- 2) The reference list of each citing document was consulted to verify and record the presence of the neuroimaging reference (OHBM poster, JSUR article, or both);
- 3) The document type was determined and recorded by reading it and searching for additional information on its source;
- 4) For each citing document, an accurate reference was elaborated for use as a standard and to determine whether GS data contain errors. An inductive and descriptive methodological approach was used to list and identify all the error types that occurred in the GS data. The reference accuracy literature served as a guide to avoid omitting the important errors in this field. A typology was elaborated and presented in the results section: 1) Data collection error (duplicates, reprints, translations, missing URLs, and inaccessible documents);

Sauvayre R., (in press) "What types of errors are hiding in Google Scholar data? A case study", *Journal of Medical Internet Research*.

- 2) Academic publication collection error (retrieval of types of documents other than journal articles, books, book chapters and conference proceedings);
  - 3) Citation error (false positive or citation counted by GS when the document does not cite the reference counted);
  - 4) Author error (missing authors, added authors, missing part of the author's name, initials errors);
  - 5) Title error (incorrect or incomplete title, spelling or typographical errors);
  - 6) Publication year error (erroneous or missing date of publication);
  - 7) Publication of source errors (journal name error identified in the "publication" column of GS);
  - 8) Publisher error (book editor name error identified in the "publisher" column of GS).
- 5) The GS errors found in each extracted column were listed; and
  - 6) The identified errors were aggregated by reference.

This collection, verification, and aggregation process required approximately 170 hours of work.

## Results

### Number of errors

A total of 755 errors were detected in 281 references retrieved from GS, for an average of 2.7 errors (range: 0-7) (Multimedia Appendix 2). Furthermore, 279 of 281 references (99.3%) contained at least one error.

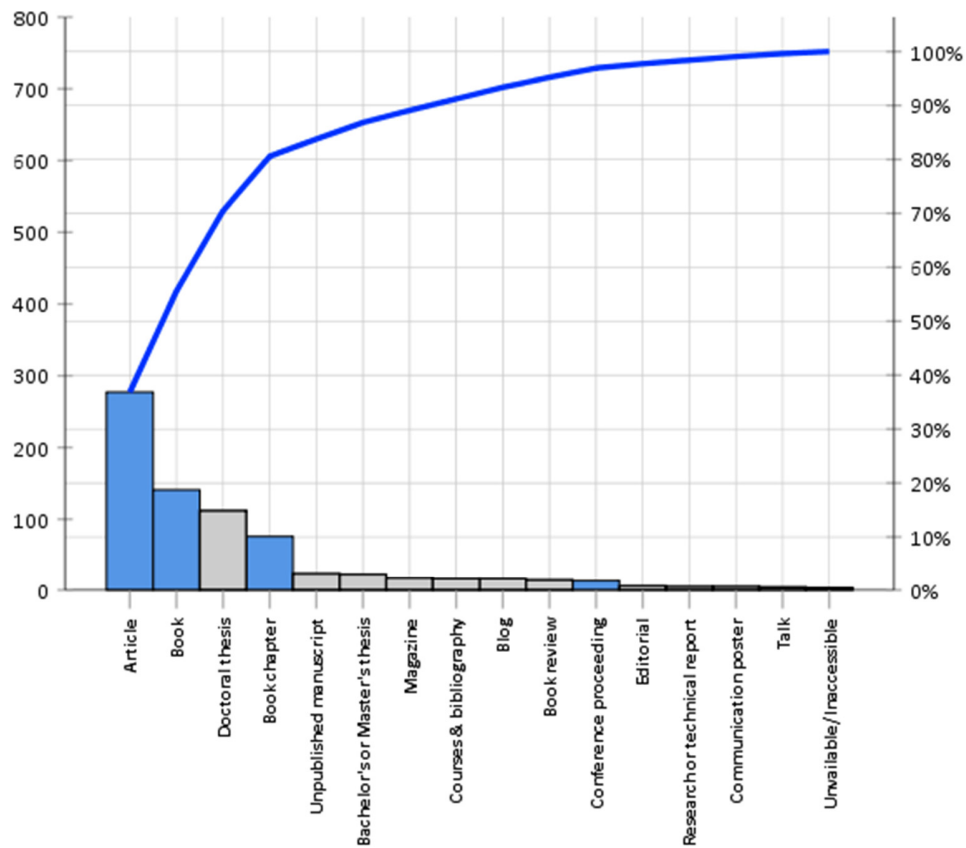


Figure 1. Pareto diagram: sum of errors detected (N = 755) as a function of document type. *Light blue*: academic publications; *Gray*: nonacademic documents. *Dark blue*: cumulative sum curve of errors detected.

### Typology of GS errors

After a manual examination of the references extracted from GS, eight types of errors were identified (Figure 2) and are thoroughly described above: 1) Data collection error; 2) Academic publication collection error; 3) Citation error (false positive); 4) Author error; 5) Title error; 6) Publication year error; 7) Publication error; 8) Publisher error.

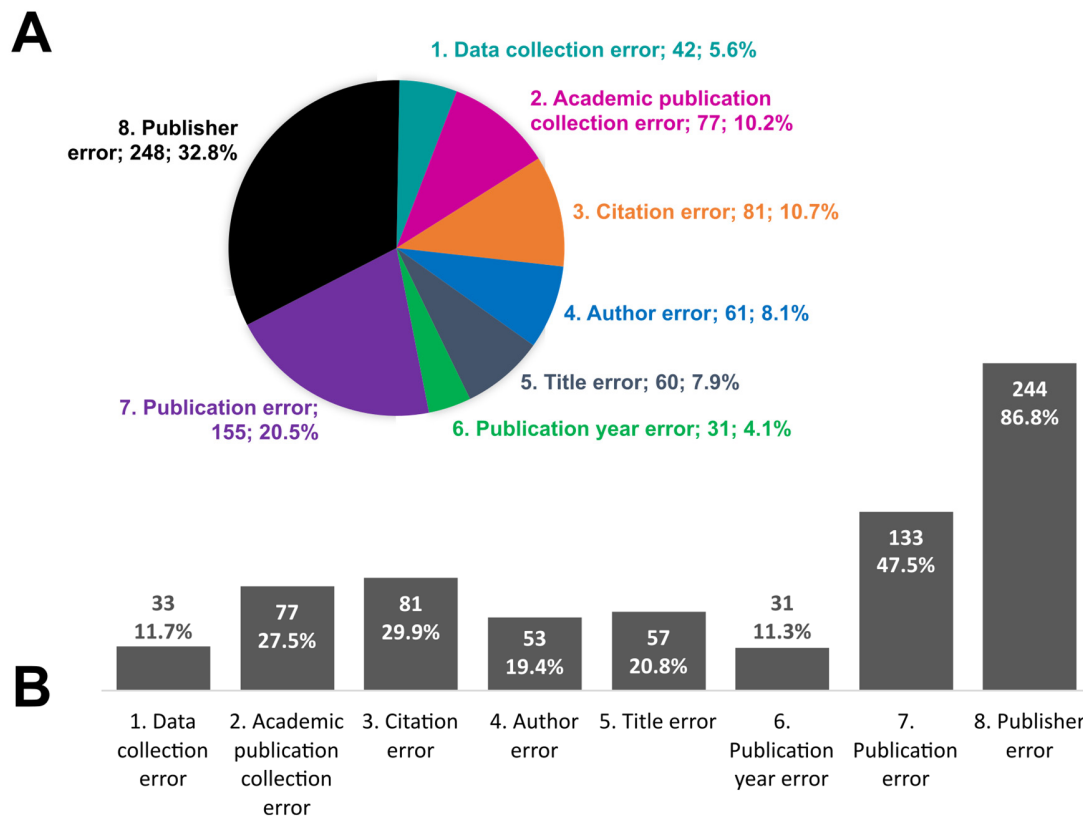


Figure 2. Typology of GS errors. Typology and proportion of errors identified (A) as a function of the number of valid references examined ( $271 < N < 281$ ) and (B) as a function of the total number of errors detected ( $N = 755$ ).

### Data collection errors

Data collection errors include duplicates, reprints, translations, missing URLs, and inaccessible documents (Multimedia Appendix 3). This type of error is identified in 33 of 281 references (11.7%), and among these errors, duplicates are detected in 16 of 281 references (5.7%). In addition, URL analysis indicates that none of GS data in any of the PoP extractions contain duplicate URLs. However, because 18 URLs were missing, a manual search for these references was then conducted to obtain and verify them. Among these missing-URL references, only 2 of the 18 citing documents were inaccessible, and 9 references were duplicates, translations, or reprints.

### Academic publication collection error

Some scientometric studies have used document type as a variable. Consequently, some researchers have focused exclusively on journal articles [3,6,29,30,46,47], whereas others have presented their collected citations per document type, including journal articles, books, book chapters and conference proceedings [25]. Furthermore, "grey literature" [4], such as theses and research reports [18], can also be included.

Sauvayre R., (in press) "What types of errors are hiding in Google Scholar data? A case study", *Journal of Medical Internet Research*.

Considering the diversity of this research method, it will be interesting to further explore the document types that GS is likely to retrieve and count.

GS describes itself as a database that "provides a simple way to broadly search for scholarly literature" [48]. But what does "scholarly literature" mean for GS? The definition provided by GS—and used in the document inclusion process—encompasses "journal papers, conference papers, technical reports, or their drafts, dissertations, pre-prints, post-prints, or abstracts" [49]. On another webpage, GS mentions that users "can search across many disciplines and sources: articles, theses, books, abstracts and court opinions" [48]. By contrast, GS excludes "news or magazine articles, book reviews, and editorials" [49] because they are "not appropriate" [49]. Nevertheless, there is no statement about the rejection of these undesirable documents from the GS index.

In the present study, the document type of each reference collected from GS was examined to determine whether the document in question were "academic publications". In this way, a document was considered an "academic publication" if and only if it was 1) an article that was published in a journal with an International Standard Serial Number (ISSN) or 2) a book, a book chapter or conference proceedings published with an International Standard Book Number (ISBN). All other document types (thesis, magazine, communication poster, bibliography, course, report, unpublished document)—so-called "nonacademic documents"—were classified as GS collection errors. Note that, according to this definition, a doctoral thesis is an academic work but not an academic publication.

As Multimedia Appendix 5 shows, GS retrieved 203 of 281 (72.5%) academic publications, but included 77 nonacademic documents in the corpus. The error rate reached 27.5% according to the definition given in the literature, whereas the GS definition led to a lower error rate (6.8%). In addition, because GS data are asymmetrically distributed, a nonparametric (Mann-Whitney) test was conducted with SPSS 25 and revealed that the nonacademic documents tend to contain more errors than the academic publications ( $U=5117.0$ ,  $P<.001$ ) (Multimedia Appendix 4).

#### **Citation error (false positive)**

The reference list of each citing document was examined to determine which of the two references (OHBM poster or JSUR article) had been cited. A total of 271 full documents were available and were read. This assessment reveals that 81 of the documents (29.9%) do not cite the reference retrieved from GS. In other words, 29.9% of the citations counted by GS are false positives. In 8 of the 271 cases (3.0%), neither of the two references are found. In 12 of the 271 cases (4.4%), the JSUR article reference is found instead of the OHBM poster one, that is, in the extraction of citations attributed by GS to the OHBM poster. Conversely, in 61 of the 271 cases (22.5%), the OHBM poster reference is found instead of the JSUR article one.

Additionally, these citation errors (false positives) affect eight times more OHBM poster references than JSUR article references (odds ratio = 7.77;  $4.4 < CI < 13.71$ ). Note that the OHBM poster reference is misreferenced in the citing documents more often than the JSUR article.

Sauvayre R., (in press) "What types of errors are hiding in Google Scholar data? A case study", *Journal of Medical Internet Research*.

#### **Author error**

As Multimedia Appendix 6 shows, 53 of 273 references (19.4%) contain at least one author error. For example, initials are removed or added. Authors are missing in 41 of 273 references (15.0%). Surprisingly, they are replaced by a journal name or by the title of either their own book or their own book chapter. Finally, 104 authors are missing, while 20 authors are improperly added. In sum, 124 authors of 565 (22.0%) are inaccurate.

#### **Title error**

A thorough examination of the "title" column extracted from GS shows that 57 of 274 references (20.8%) contain at least one error (Multimedia Appendix 7). The incompleteness of the title is the most common error identified. As a result of this error, some incomplete titles are similar to other publication titles. Furthermore, several errors are more questionable, such as replacement of a book title with a chapter title from this aforementioned book or with the title of a different chapter from another book by an author who contributed a chapter to this book. Other questionable title errors are the assemblage of two-chapter parts published in the same book and the replacement of the publication title by its editor's name or the domain name of the website that hosts it. Surprisingly, irrelevant parts are added to the publication title, such as an ISBN number, the price of the book, the name of the book collection, and an excerpt from the front page of a thesis ("a dissertation submitted for the degree of Doctor of Philosophy"). Lastly, the reference titles also contain typographical or spelling errors.

#### **Publication year error**

Publication year errors are detected in 31 of 274 references (11.3%). In most cases, the years are missing—they are replaced by "zero" in 22 references. In other cases, the actual publication year of the JSUR article or the OHBM poster is increased by one year or decreased by 1, 3, 7, or 100 years (Multimedia Appendix 8).

#### **Publication of source errors**

The "publication" or "source" column retrieved from GS via PoP shows inconsistencies that depend on the document type of references (Multimedia Appendix 9). Indeed, it contains journal names, books, edited book titles, conference proceedings titles, magazine names, publisher names, domain names of websites that host the citing documents, irrelevant parts of references, and even an author's address. Furthermore, a large number of missing values (i.e., "not provided" in Multimedia Appendix 9) are found in these publication data, affecting one in three references (32.0%). These missing publication data are observed most often for theses (bachelor's, master's and doctoral) and book references.

In total, 133 of 280 references (47.5%) contain errors (Table 1). These errors are mostly identified in conference proceedings, edited books, and journal articles. In

addition, only half of the citations counted by GS are usable as academic publication material (Table 1, "utility" column) because, for example, GS provides a domain name instead of the academic journal name. Among these usable data, 133 are inaccurate. Finally, in this corpus, only 60 of 280 references (21.4%) are proper usable data.

**Table 1.** Accurate and inaccurate content identified in the "Publication" column retrieved from GS via PoP (n=280).<sup>a</sup>

Type of errors, n (%)	Inaccurate publication	Accurate publication	Total references	Utility
Journal name	56 (51.9)	52 (48.1)	108 (100.0)	(+)
Magazine name	1 (50.0)	1 (50.0)	2 (100.0)	(-)
Book title	13 (100.0)	0 (0.0)	13 (100.0)	(-)
Edited book title	21 (72.4)	8 (27.6)	29 (100.0)	(+)
Conference proceedings title	5 (100.0)	0 (0.0)	5 (100.0)	(+)
Thesis title	2 (100.0)	0 (0.0)	2 (100.0)	(-)
Publisher name	2 (100.0)	0 (0.0)	2 (100.0)	(-)
Domain name	18 (90.0)	2 (10.0)	20 (100.0)	(-)
Preprint database name	1 (25.0)	3 (75.0)	4 (100.0)	(-)
Other	5 (100.0)	0 (0.0)	5 (100.0)	(-)
Missing value (not provided)	9 (100.0)	81 (90.0)	90 (100.0)	(-)
Total	133 (47.5)	147 (52.5)	280 (100.0)	

<sup>a</sup> The usable publication content for studies using academic publications is denoted by a "+". These errors were not easy to categorize because of the nonacademic documents. For instance, when the document type is a blog post or an unpublished draft, a journal name is not expected in the "publication" column and thus is counted as an inaccuracy. Nevertheless, this type of document had already been counted as a data collection error. Therefore, each document type was specifically analyzed to avoid falsely increasing the error count. However, the categorization was easier for other references, such as when the journal editor name was provided instead of the journal name. In addition, an examination of spelling and typographical errors, including capitalization errors, was conducted.

These source inconsistencies mainly occur in journal names as typographical errors, particularly capitalization errors (Multimedia Appendix 10). The second most frequent error is title and journal name incompleteness. Journal names are then heavily truncated, as shown in the following examples: "Journal of ..." instead of "Journal of Advertising Research" and "Rev ..." instead of "Revista de neurologia". The same type of inaccuracy was identified in the edited book titles: "... Routledge Handbook of ..." instead of "The Routledge Handbook of Neuroethics" and "... Imaging of the ..." instead of "Imaging of the Pelvis, Musculoskeletal System, and Special Applications to CAD". Furthermore, as several journal names begin with "Journal of" and several edited books begin with "Routledge Handbook of", the incompleteness of the GS data may cause difficulties.

### Publisher error

The "publisher" column retrieved from GS provides a variety of content (Multimedia Appendix 11): editor name (including journal editor), journal name, domain name of the website that hosts the citing document (e.g., 42 of the domain names are "books.google.com"), digital library (i.e., JSTOR), and missing values. The "publisher" column contains the highest error rate found in the GS data, which is 244 of 281 references (86.8%) (Table 2). Indeed, the 248 inaccuracies detected in this column constitute a third (32.9%) of the total errors identified. Journal editors and domain names are frequently inaccurate. The utility of this publisher data is then limited to studies using academic publication data. Only the editor names of books, book chapters, and conference proceedings are usable, but they actually represent 35 of the 281 references (12.5%). Furthermore, an error rate of 37.1% is found in these usable data. For example, an editor's name is replaced by an irrelevant one: The Penguin Press by Australia Books; and Palgrave Macmillan by Springer.

**Table 2.** Accurate and inaccurate content in the "Publisher" column retrieved from GS via PoP (N=281). The usable publication content for studies using academic publication data is denoted by a "+".

Type of errors, n (%)	Inaccurate publication	Accurate publication	Total references	Utility
Book & conference proceedings editor	13 (37.1)	22 (62.9)	35 (100.0)	(+)
Journal editor	51 (100.0)	0 (0.0)	51 (100.0)	(-)
Journal name	1 (100.0)	0 (0.0)	1 (100.0)	(-)
Digital library name	2 (100.0)	0 (0.0)	2 (100.0)	(-)
Domain name	167 (100.0)	0 (0.0)	167 (100.0)	(-)
Not provided	10 (40.0)	15 (60.0)	25 (100.0)	(-)
Total	244 (86.8)	37 (13.2)	281 (100.0)	

### Discussion

The aim of the present study was to examine the accuracy of citation data collected from GS via PoP and to provide a comprehensive description of the errors and miscounts identified. In fact, the extraction of raw data with inaccuracies from GS may generate incorrect results in several research areas, such as bibliometrics, scientometrics, and research evaluation. Despite the data cleaning performed by researchers (mainly duplicate removal), citation counts retrieved from GS are generally used without substantial caution. Furthermore, few comprehensive studies list the different types of GS errors, and no previous research seems to quantify the inherent problems of GS citations collected by PoP. The present study was therefore conducted to provide a meticulous analysis of GS data to anticipate the risk of errors that may affect data and the results of studies using them.

### Ranking of GS errors

The GS errors were analyzed using 281 documents that cited a neuroimaging work performed to raise awareness of false positive results in the scientific community. The present study reveals an unprecedented error rate: 279 of 281 examined references (99.3%) contain at least one error. The academic publications are not free from errors: they account for 503 of the 755 detected errors (67.0%). However, the nonacademic documents tend to contain more errors than the academic publications ( $U=5117.0$ ,  $P<.001$ ).

The cumulative error rate detected in the present study—99.3% of references contain at least one error—differs from the median rate (14.6%) reported in the literature over the past 10 years. This difference may be explained by several aspects of previous research:

- 1) an automatic approach was generally used to clean the data, while a manual examination was conducted in the present study;
- 2) a varied but low number of variables were examined in these studies: a median of one type of error was examined in previous studies, while eight types of errors were examined in the present study;
- 3) the usual purpose of these studies was to compare the coverage of GS, WoS, and Scopus; thus, the researchers mainly verified duplicates in an aggregated corpus drawn from these three databases; and
- 4) these studies do not cumulate the number of errors identified per reference.

These discrepancies make comparison difficult, but data provided by De Winter et al. [18]—"Online Supplementary material 5 Excel File"—make it possible. Through these data, an error rate cumulated by reference was calculated to compare what is comparable. However, as these researchers used four error types, the comparison was performed for academic publication collection errors, author errors, title errors and duplicates. All other things being equal, the present study reports an error rate three times higher than that reported by De Winter et al. [18] (64.8% and 20.5%, respectively). These findings suggest that citation counts and references extracted from GS are not fully reliable and may expose the researchers who use them to numerous errors. Note that the content of GS is the result of automatic indexing of websites by robots. The coverage depends on the indexed websites. Moreover, according to GS, "robots generally try to index every paper from every website they visit, including most major sources and also many lesser known ones" [50]. Thus, the reliability of GS is a type of "photography" of the reliability of authors and editors' websites. Since errors can happen, it is important to identify the possible impact of GS's lack of reliability with respect to research data.

### The impact of GS errors in research data

What is the probable impact of GS errors in the citation analysis or research evaluation areas when citation counts and references are used without data cleaning?

Sauvayre R., (in press) "What types of errors are hiding in Google Scholar data? A case study", *Journal of Medical Internet Research*.

#### **Publisher error**

The useful content that a researcher needs to find in the "publisher" column extracted from GS via PoP is the editor name for books, book chapters, and conference proceedings. However, this column mainly contains the domain name of the website hosting the citing document. Thus, only 7.8% of these collected data are free from errors and usable in an academic publication study. The "publication" column therefore requires meticulous examination before use. The first step is to determine the document type of each collected reference because GS still does not provide it.

#### **Publication error**

The PoP manual indicates that the "publication" column contains "journal name or similar"—"similar" is not explicitly defined—which is "not always available" and "sometimes wrong" [51]. However, the "publication" content is more disparate and incorrect than this. Indeed, it contains journal names, book titles, thesis titles, publisher names, and domain names, and one-third of this column is missing values. Only 21.4% of "publication" data are free from errors and usable in an academic publication study. In addition, GS errors can impact studies in the following ways. First, the GS error rate can negatively affect the evaluation of journal impact factors and the journal ranking. Second, missing values (32.4% of references) can alter relational database management [52]. Third, typographical errors (including capitalization errors) can lead to duplicates [4]. Lastly, note that there is a risk in using this "publication" data because of the large number of errors detected in the journal names, edited book titles, and conference proceedings titles.

#### **Citation error (false positive)**

The GS citation count is distorted by documents that do not cite the reference retrieved. This point is often reported in the literature—for instance, as a "phantom citation" or "false citation" [20]—but the no reported error rate [18,30,32,53,54] is as high as the rate found in the present study (29.9%). This difference can be explained by the highly similar titles of the two references examined (OHBM poster and JSUR article). This finding also demonstrates the difficulties of addressing this type of similarity in GS data. Consequently, researchers may use data samples that contain false positive citations and then may obtain biased findings.

#### **Academic publication collection error**

GS failed to retrieve only academic publications. Indeed, 27.5% of the citing documents are nonacademic publications, including doctoral theses, magazine articles, preprints, reports, courses, bibliographies, and blog posts. This error rate confirms previous findings [30]. However, if the GS definition of "scholarly literature" is applied, this error rate falls to 6.8%. This GS definition differs widely from the definition of "academic publication" used in the present study. Thus, GS seems to inaccurately report citation counts and references of academic publications, and

consequently, it does not accurately reflect the dissemination of published work. Therefore, the results of many scientometric studies using GS data to examine the publication activity of scientists, particularly in research evaluation, may be questionable when these data are not verified and cleaned (document type and false positive). The citation counts and h-index scores calculated by GS are also questionable. This raises questions about the reliability of studies that compare the coverage of GS, WoS and Scopus and conclude that GS collects significantly more citations [12,28,29] than its competitors. Further research should explore the citation counts of these databases to determine how comparable they are.

#### **Title error**

The main issue with the titles retrieved from GS is incompleteness, which causes problems such as false positive matches. The similarity between the OHBM poster title and the JSUR article title demonstrates this GS difficulty. Other errors—typographical and spelling—cause problems in database management. More unwelcome is the missing title error: instead of the title, 6.2% of references contain, for example, editor names, domain names, or ISBN numbers. These missing title errors raise several problems: 1) references cannot be retrieved by a search by title, and 2) duplicates can be more frequent and more difficult to detect.

#### **Author error**

The citing documents examined are cowritten by 565 authors. Nevertheless, 124 authors (22.0%) are either incorrect or missing. These errors can cause problems in studies of the structure of scientific collaborative networks, which are commonly used graphs. Indeed, a fifth of the collaborative networks built may be incorrect and thus may generate imperfect relationships. First, the missing authors may truncate an important share of all the authors involved. Second, the irrelevant added authors may create a bias that a graph's algorithms can reinforce. Consequently, researchers may overestimate a relationship or ignore another determinant one.

#### **Data collection error**

Duplicates, translations, and reprints are frequent in GS data. As collected data can be biased by duplicates, their detection is the first step implemented in studies using GS data. The duplicate, translation, and reprint rates found in the present study are similar to those in previous findings [30,32]. In addition, the URL address of citing documents is commonly used to detect duplicates and collect full text documents. Because a missing URL may cause difficulties, previous studies resolved this issue by automatically deleting a reference without a web address [16,32]. By contrast, in the present study, 6.4% of URL addresses are missing, but only 0.7% of them cannot be found with a manual search. Half of these found documents are usable, and half are

Sauvayre R., (in press) "What types of errors are hiding in Google Scholar data? A case study", *Journal of Medical Internet Research*.

duplicate, translation, and reprint references. Consequently, the duplicate search removes 7.8% of the references, whereas the irrelevant deletion of references without a URL address leads to the omission of 3.2% of the citing documents. Again, this may cause biased results.

#### **Publication year error**

Incorrect years have a lower frequency than missing years (3.3% and 8.0% of references, respectively). These missing values cause major problems in data collection. As GS limits the search results to the first 1000 citing references per query, certain researchers have collected data by publication year to obtain a larger corpus of GS citations [25] or to focus their analysis on a specific period of time [44]. Other researchers have removed references containing incorrect publication years [16]. Thus, these neglected references may lead to truncated data and biased results. Inherent to the failed indexation process of GS, this publication year error may cause sampling errors that affect the representativeness of findings.

#### **Data verification versus biased results**

The GS error rate seems to be negligible when types of errors are considered in isolation. These types of claims have been made about false positives [20], duplicates [10], and incorrect publication years [2]. By contrast, with regard to GS errors, Harzing [55] argues fatalistically that "bibliometrics is an inexact science and that any data source has its own flaws". However, Hicks et al. [40], in presenting the Leiden Manifesto, emphasize the importance of the quality of the data used in research evaluation. Conversely, when the GS error rates are observed as a whole, a worrying cumulative effect is revealed. Indeed, only 2 of 281 (0.71%) references collected from GS are free from errors. This raises a question about the reliability of the GS citation counts. In the present study, two neuroimaging works were cited 281 times according to GS. However, this citation count is incorrect. In fact, these works are cited 131 times in academic publications (i.e., excluding duplicates, reprints, translations, inaccessible documents, and false positives), which is 53.4% less than the GS claim. Thus, the full sample collected from GS (281 citations) can considerably differ from the proper sample (131 citations). There is thus a major risk of producing incorrect and biased results that do not accurately reflect the data examined.

Consequently, meticulous verification and cleaning of GS data are essential before using them. Considering this, several precautions should be taken to improve the reliability of GS data:

- 1) Detect and remove duplicate, translation, and reprint references and subsequently merge their citation counts;
- 2) Consult the full-text documents of the full sample to remove false positive matches; and

- 3) Verify the document type of each reference to exclude nonacademic publications.

Because results will be biased or wrong if these verification steps are not performed, is it possible to study a large-scale sample of GS citations (approximately several thousand)? It seems unlikely unless substantial resources are allocated for such verification. Indeed, Meho and Yang [21] were allowed 18 minutes per reference (3,000 hours of work for a 10,000 citation samples). In the present study, a work time of 32 minutes per reference was necessary to complete the verification (150 hours for 281 citations). What about automatization of the verification? Studies that cleaned large-scale data either in part or as a whole using an automatic cleaning process report a lower error rate and fewer error types than studies using a manual cleaning process. Therefore, it is reasonable to have doubts about the efficiency of this automatic cleaning.

Finally, studying a small sample of GS data seems more adequate than studying a large sample in terms of obtaining reliable data and accurate findings. Nevertheless, there is a need to conduct further research to develop statistical tools for weighting the correlation calculation in a large-scale sample of GS data, which are widely used in database coverage studies. However, these tools may not correct the collection issue inherent to the GS database.

Alternatively, according to the reference accuracy literature [9,10,12-15], databases such as WoS or Scopus can be used in place of GS because they are more reliable, though they have narrower coverage than GS. Indeed, WoS has an average error rate of 0.1% [4,18,31,32], and this rate is 1.0% for Scopus [10,31]. However, since GS is a free database [56], it may be the only possible way to conduct a study. However, knowing that all databases are likely to contain errors, verifying a sample of data is a useful precaution.

To conclude, the categorization of the errors encountered in the data extracted from GS provides researchers with methodological and epistemological reflections so that they become aware, with precision, of the probable errors that they are likely to encounter and can consequently adjust their methodological choice. For example, the number of citations obtained by GS may not be completely accurate, or the names of the authors mentioned may not be completely correct. With a sample of several thousand references, these errors can have a noticeable impact on the results.

## Conclusion

Almost all of the data retrieved from GS contain at least one error calling the reliability of GS data into question. Further, the reliability of studies using a large-scale sample without verification and data cleaning is also called into question. Moreover, studies using GS to evaluate research activity or compare the coverage of several databases (i.e., GS, WoS, Scopus) may be affected by substantial biases, including citation miscounts.

However, researchers able to spend a considerable amount of time on the meticulous verification of their small samples can obtain various references for journal articles, books, edited book chapters, and conference proceedings from GS. This ability can be

Sauvayre R., (in press) "What types of errors are hiding in Google Scholar data? A case study", *Journal of Medical Internet Research*.

especially useful in bibliometric studies based on material published in research areas in which journal articles are less predominant than other publication types.

### Limitations

Since the data used are limited and specific, the results obtained cannot be generalized. However, this case study provides a kind of "stress test" of GS to promote reflection on the limits of this free database.

### Conflicts of Interest

The author declares no competing interests.

### Abbreviations

**GS:** Google Scholar;

**JCR:** Journal Citation Reports;

**JSUR:** *Journal of Serendipitous and Unexpected Results*;

**JSUR article:** article published by Craig Bennett et al. [41] in the *Journal of Serendipitous and Unexpected Results*;

**OHBM:** Organization for Human Brain Mapping;

**OHBM poster:** poster presented by Craig Bennett et al. [39] at the 15th annual meeting for the Organization for Human Brain Mapping;

**WoS:** Web of Science;

**PoP:** Publish or Perish software;

**CI:** confidence interval

### Multimedia Appendix 1

Raw data.

### Multimedia Appendix 2

Number of errors detected per reference retrieved from GS

N° of errors, n, %	References	Errors
0	2 (0.7)	0 (0.0)
1	62 (22.1)	60 (7.9)
2	82 (29.2)	152 (20.1)
3	89 (31.7)	261 (34.6)
4	25 (8.9)	84 (11.1)
5	12 (4.3)	85 (11.3)
6	7 (2.5)	78 (10.3)
7	2 (0.7)	35 (4.6)
Total	281 (100.0)	755 (100.0)

Sauvayre R., (in press) "What types of errors are hiding in Google Scholar data? A case study", *Journal of Medical Internet Research*.

### Multimedia Appendix 3

Data collection error: types and rates.

Type of errors, n, %	N° errors	Error rate (%) / N° references	Error rate (%) / Total N° of errors
Duplicate	16 (38.1)	5.7	2.1
Translation/Reprint	6 (14.3)	2.1	0.8
Unavailable/Inaccessible	2 (4.8)	0.7	0.3
URL missing	18 (42.9)	6.4	2.4
Total	42 (100.0)	11.7	5.6

### Multimedia Appendix 4

Results of the Mann-Whitney test on the type of reference (academic publication or nonacademic document) and the number of GS errors and data visualization with a box plot.

#### Mann-Whitney test

Ranks				
Type of reference		N	Mean rank	Sum of ranks
Error type	Academic publication	203	127,21	25823,50
	Nonacademic document	77	175,54	13516,50
	Total	280		

#### Significance test<sup>a</sup>

	Nombre_erreurs
U of Mann-Whitney	5117,500
W of Wilcoxon	25823,500
Z	-4,599
Asymptotic significance (bilateral)	,000

a. Grouping variable: Type of reference

### Multimedia Appendix 5

Document types of references collected from GS as a function of three categories (academic journal, nonacademic journal, and GS non-"scholarly literature").

Sauvayre R., (in press) "What types of errors are hiding in Google Scholar data? A case study", *Journal of Medical Internet Research*.

Type of reference, n, %	Academic publications	Nonacademic publications	GS Non- "Scholarly literature"	Total references
Article	114 (40.7)	1 (0.4)		115 (41.1)
Bachelor's or master's thesis		8 (2.9)		8 (2.9)
Blog		3 (1.1)	3 (1.1)	3 (1.1)
Book	53 (18.9)			53 (18.9)
Book chapter	31 (11.1)			31 (11.1)
Book review		4 (1.4)	4 (1.4)	4 (1.4)
Conference proceeding	5 (1.8)			5 (1.8)
Communication poster		1 (0.4)		1 (0.4)
Courses & bibliography		4 (1.4)	4 (1.4)	4 (1.4)
Doctoral thesis		37 (13.2)		37 (13.2)
Editorial		3 (1.1)	3 (1.1)	3 (1.1)
Magazine		5 (1.8)	5 (1.8)	5 (1.8)
Research or technical report		2 (0.7)		2 (0.7)
Talk		1 (0.4)		1 (0.4)
Unpublished manuscript (preprint included)		8 (2.9)		8 (2.9)
Total	203 (72.5)	77 (27.5)	19 (6.8)	280 (100.0)

## Multimedia Appendix 6

Inaccurate content identified in the "Author" column retrieved from GS via PoP.

Type of errors, n, %	N° errors	Error rate (%) / N° references	Error rate (%) / Total N° of errors
Missing authors	41 (67.2)	15.0	5.4
Added authors	10 (16.4)	3.7	1.3
Missing part of the author's name	1 (1.6)	0.4	0.1
Initials errors	5 (8.2)	1.8	0.7
Replacement of authors by a book or a book chapter title	3 (4.9)	1.1	0.4
Replacement of authors by a journal name	1 (1.6)	0.4	0.1
Total	61 (100.0)	19.4	8.1

## Multimedia Appendix 7

Inaccurate content identified in the "Title" column retrieved from GS via PoP.

Type of errors, n, %	N° errors	Error rate (%) / N° references	Error rate (%) / Total N° of errors
----------------------	-----------	--------------------------------	-------------------------------------

Sauvayre R., (in press) "What types of errors are hiding in Google Scholar data? A case study", *Journal of Medical Internet Research*.

Spelling error	1 (1.7)	0.4	0.1
Typographical error	3 (5.0)	1.1	0.4
Incorrect title	5 (8.3)	1.8	0.7
Incomplete title	30 (50.0)	10.9	4.0
Replacement of the title by an editor	2 (3.3)	0.7	0.3
Replacement of the book title by one of the book chapter titles	11 (18.3)	4.0	1.5
Replacement of the thesis title by one of the thesis chapter titles	2 (3.3)	0.7	0.3
Replacement of the title by the domain name of the website hosting the document	2 (3.3)	0.7	0.3
Irrelevant part added	4 (6.7)	1.5	0.5
Total	60 (100.0)	20.8	7.9

### Multimedia Appendix 8

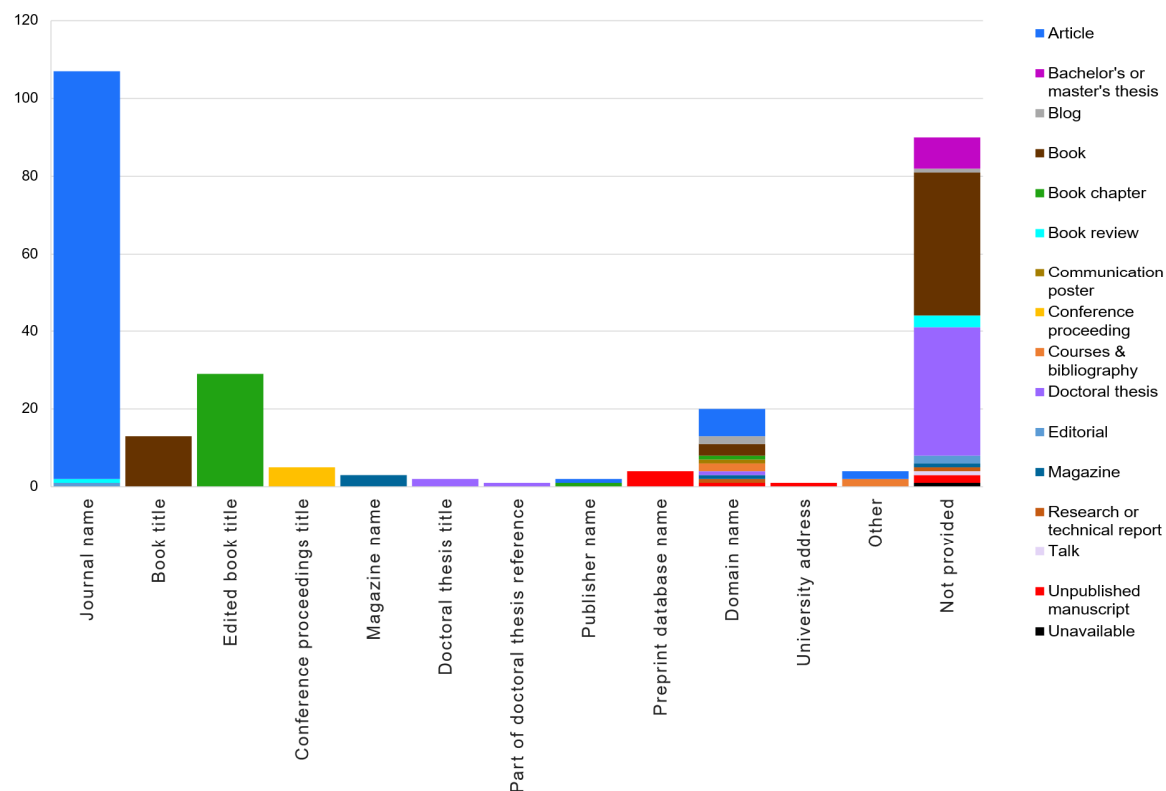
Inaccurate content identified in the "Year" column retrieved from GS via PoP.

Type of errors, n, %	N° errors	Error rate (%) / N° references	Error rate (%) / Total N° of errors
"0"	22 (71.0)	8.0	2.9
+ 1 an	2 (6.5)	0.7	0.3
- 1 an	4 (12.9)	1.5	0.5
- 3 ans	1 (3.2)	0.4	0.1
- 7 ans	1 (3.2)	0.4	0.1
- 100 ans	1 (3.2)	0.4	0.1
Total	31 (100.0)	11.3	4.1

### Multimedia Appendix 9

Content of the "Publication" column retrieved from GS via PoP as a function of reference document type.

Sauvayre R., (in press) "What types of errors are hiding in Google Scholar data? A case study", *Journal of Medical Internet Research*.



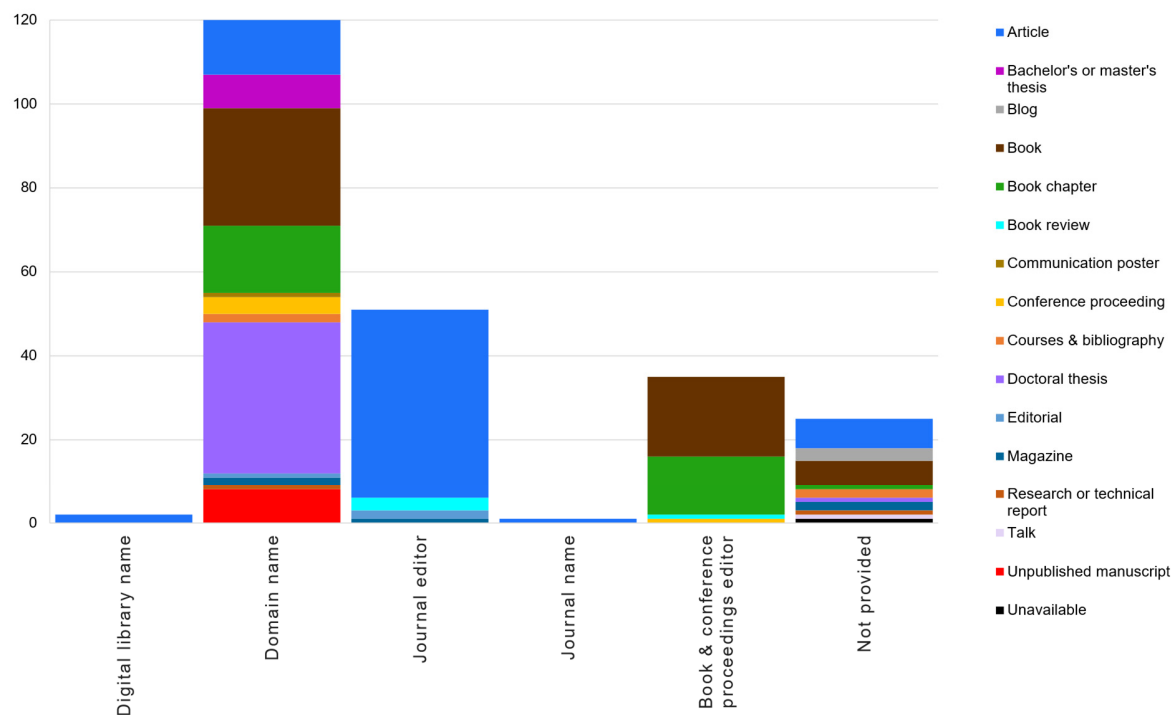
## Multimedia Appendix 10

Spelling and orthographical errors in the academic publications retrieved from GS via PoP: journal name, edited book title, and conference proceedings book title.

Type of errors, n, %	N° errors	Error rate (%) / N° references
<b>Journal name error</b>	<b>68 (100.0)</b>	<b>21.1</b>
Incorrect abbreviation	1 (1.5)	0.4
Capitalization error	34 (50.0)	12.1
Irrelevant part added	1 (1.5)	0.4
Incomplete name	31 (45.6)	11.1
Incorrect name	1 (1.5)	0.4
<b>Title error</b>	<b>33 (100.0)</b>	<b>11.8</b>
Incomplete edited book title	26 (78.8)	9.3
Incomplete conference proceeding book title or incomplete edited book title	5 (15.2)	1.8
Incorrect edited book title	2 (6.1)	0.7

## Multimedia Appendix 11

Content of the "Publisher" column retrieved from GS via PoP as a function of reference document type.



## References

1. Orduna-Malea E, Martín-Martín A, Delgado López-Cózar E. Google Scholar as a source for scholarly evaluation: a bibliographic review of database errors. *Rev Esp Doc Científica*. 2017;40(4):e185. doi:10.3989/redc.2017.4.1500
2. Martín-Martín A, Orduna-Malea E, Harzing A-W, Delgado López-Cózar E. Can we use Google Scholar to identify highly-cited documents? *J Informetr*. 2017;11(1):152-163. doi:10.1016/j.joi.2016.11.008
3. Falagas ME, Pitsouni EI, Malietzis GA, Pappas G. Comparison of PubMed, Scopus, Web of Science, and Google Scholar: strengths and weaknesses. *FASEB J*. 2008;22(2):338-342. PMID:17884971 doi:10.1096/fj.07-9492LSF
4. Haddaway NR, Collins AM, Coughlin D, Kirk S. The Role of Google Scholar in Evidence Reviews and Its Applicability to Grey Literature Searching. *PLoS ONE*. 2015;10(9):e0138237. PMID:26379270 doi:10.1371/journal.pone.0138237
5. Bornmann L, Marx W, Schier H, Rahm E, Thor A, Daniel H-D. Convergent validity of bibliometric Google Scholar data in the field of chemistry—Citation counts for papers that were accepted by *Angewandte Chemie International Edition* or rejected but published elsewhere, using Google Scholar, Science Citation Index, Scopus, and Chemical Abstracts. *J Informetr*. 2009;3(1):27-35. doi:10.1016/j.joi.2008.11.001

6. Abad-García M-F, González-Teruel A, González-Llinares J. Effectiveness of OpenAIRE, BASE, Recolecta, and Google Scholar at finding spanish articles in repositories. *J Assoc Inf Sci Technol*. 2018;69(4):619-622. doi:10.1002/asi.23975
7. Orduna-Malea E, Ayllón JM, Martín-Martín A, Delgado López-Cózar E. Methods for estimating the size of Google Scholar. *Scientometrics*. 2015;104(3):931-949. doi:10.1007/s11192-015-1614-6
8. Mingers J, Lipitakis EAECG. Counting the citations: a comparison of Web of Science and Google Scholar in the field of business and management. *Scientometrics*. 2010;85(2):613-625. doi:10.1007/s11192-010-0270-0
9. Harzing A-W, Alakangas S. Google Scholar, Scopus and the Web of Science: a longitudinal and cross-disciplinary comparison. *Scientometrics*. 2016;106(2):787-804. doi:10.1007/s11192-015-1798-9
10. Moed HF, Bar-Ilan J, Halevi G. A new methodology for comparing Google Scholar and Scopus. *J Informetr*. 2016;10(2):533-551. doi:10.1016/j.joi.2016.04.017
11. Gehanno J-F, Rollin L, Darmoni S. Is the coverage of google scholar enough to be used alone for systematic reviews. *BMC Med Inform Decis Mak*. 2013;13:7. PMID:23302542 doi:10.1186/1472-6947-13-7
12. Harzing A-W. A preliminary test of Google Scholar as a source for citation data: a longitudinal study of Nobel prize winners. *Scientometrics*. 2013;94(3):1057-1075. doi:10.1007/s11192-012-0777-7
13. Mingers J, Xu F. The drivers of citations in management science journals. *Eur J Oper Res*. 2010;205(2):422-430. doi:10.1016/j.ejor.2009.12.008
14. Franceschet M. A comparison of bibliometric indicators for computer science scholars and journals on Web of Science and Google Scholar. *Scientometrics*. 2010;83(1):243-258. PMID:33008629 doi:10.1007/s11192-009-0021-2
15. Harzing A-W. A longitudinal study of Google Scholar coverage between 2012 and 2013. *Scientometrics*. 2014;98(1):565-575. doi:10.1007/s11192-013-0975-y
16. Prins AAM, Costas R, Van Leeuwen TN, Wouters PF. Using Google Scholar in research evaluation of humanities and social science programs: a comparison with Web of Science data. *Res Eval*. 2016;25(3):264-270. doi:10.1093/reseval/rvv049
17. Kousha K, Thelwall M, Rezaie S. Assessing the citation impact of books: the role of Google Books, Google Scholar, and Scopus. *J Am Soc Inf Sci Technol*. 2011;62(11):2147-2164. doi:10.1002/asi.21608
18. De Winter JCF, Zadpoor AA, Dodou D. The expansion of Google Scholar versus Web of Science: a longitudinal study. *Scientometrics*. 2014;98(2):1547-1565. doi:10.1007/s11192-013-1089-2
19. Wildgaard L. A comparison of 17 author-level bibliometric indicators for researchers in Astronomy, Environmental Science, Philosophy and Public Health in Web of Science and Google Scholar. *Scientometrics*. 2015;104(3):873-906. doi:10.1007/s11192-015-1608-4
20. Franceschini F, Maisano D, Mastrogiacomo L. Empirical analysis and classification of database errors in Scopus and Web of Science. *J Informetr*. 2016;10(4):933-953. doi:10.1016/j.joi.2016.07.003

21. Meho LI, Yang K. Impact of data sources on citation counts and rankings of LIS faculty: Web of Science versus scopus and google scholar. *J Am Soc Inf Sci Technol*. 2007;58(13):2105-2125. doi:10.1002/asi.20677
22. Li J, Sanderson M, Willett P, Norris M, Oppenheim C. Ranking of library and information science researchers: comparison of data sources for correlating citation data, and expert judgments. *J Informetr*. 2010;4(4):554-563. doi:10.1016/j.joi.2010.06.005
23. Mingers J, O'Hanley JR, Okunola M. Using Google Scholar institutional level data to evaluate the quality of university research. *Scientometrics*. 2017;113(3):1627-1643. PMID:29200538 doi:10.1007/s11192-017-2532-6
24. Martín-Martín A, Orduna-Malea E, Delgado López-Cózar E. A novel method for depicting academic disciplines through Google Scholar citations: the case of Bibliometrics. *Scientometrics*. 2018;114(3):1251-1273. doi:10.1007/s11192-017-2587-4
25. Martín-Martín A, Orduna-Malea E, Ayllon JM, Delgado López-Cózar E. A two-sided academic landscape: snapshot of highly-cited documents in Google Scholar (1950-2013). *Rev Esp Doc Científica*. 2016;39(4):e149. doi:10.3989/redc.2016.4.1405
26. Pitol SP, De Groote SL. Google Scholar versions: do more versions of an article mean greater impact? *Libr Hi Tech*. 2014;32(4):594-611. doi:10.1108/LHT-05-2014-0039
27. Jamali HR, Nabavi M. Open access and sources of full-text articles in Google Scholar in different subject fields. *Scientometrics*. 2015;105(3):1635-1651. doi:10.1007/s11192-015-1642-2
28. Lasda Bergman EM. Finding citations to social work literature: the relative benefits of using Web of Science, Scopus, or Google Scholar. *J Acad Librariansh*. 2012;38(6):370-379. doi:10.1016/j.acalib.2012.08.002
29. De Groote SL, Raszewski R. Coverage of Google Scholar, Scopus, and Web of Science: a case study of the h-index in nursing. *Nurs Outlook*. 2012;60(6):391-400. PMID:22748758 doi:10.1016/j.outlook.2012.04.007
30. Granter SR, Laga AC, Larson AR. Calciphylaxis and the Persistence of Medical Misinformation in the Era of Google. *Am J Clin Pathol*. 2015;144(3):427-431. PMID:26276773 doi:10.1309/AJCPDMWVGKW9N1CU
31. Kulkarni AV, Aziz B, Shams I, Busse JW. Comparisons of Citations in Web of Science, Scopus, and Google Scholar for Articles Published in General Medical Journals. *J Am Med Assoc*. 2009;302(10):1092-1096. PMID:19738094 doi:10.1001/jama.2009.1307
32. García-Pérez MA. Accuracy and completeness of publication and citation records in the Web of Science, PsycINFO, and Google Scholar: a case study for the computation of h indices in psychology. *J Am Soc Inf Sci Technol*. 2010;61(10):2070-2085. doi:10.1002/asi.21372
33. Bar-Ilan J. Which h-index? — A comparison of WoS, Scopus and Google Scholar. *Scientometrics*. 2008;74(2):257-271. PMID:22748758 doi:10.1007/s11192-008-0216-y

34. Valderrama-Zurián J-C, Aguilar-Moya R, Melero-Fuentes D, Aleixandre-Benavent R. A systematic analysis of duplicate records in Scopus. *J Informetr.* 2015;9(3):570-576. doi:10.1016/j.joi.2015.05.002
35. Fiolet T, Guihur A, Rebeaud ME, Mulot M, Peiffer-Smadja N, Mahamat-Saleh Y. Effect of hydroxychloroquine with or without azithromycin on the mortality of coronavirus disease 2019 (COVID-19) patients: a systematic review and meta-analysis. *Clin. Microbiol. Infect.* 2021;27(1):19-27. PMID: 32860962 doi:10.1016/j.cmi.2020.08.022
36. Oyibo K, Sahu KS, Oetomo A, Morita PP. Factors Influencing the Adoption of Contact Tracing Applications: Protocol for a Systematic Review. *JMIR Res. Protoc.* 2021;10(6):e28961. PMID:33974551 doi:10.2196/28961
37. Silva LOJ e, Maldonado G, Brigham T, Mullan AF, Utengen A, Cabrera D. Evaluating Scholars' Impact and Influence: Cross-sectional Study of the Correlation Between a Novel Social Media-Based Score and an Author-Level Citation Metric. *J Med Internet Res.* 2021;23(5):e28859. PMID:34057413 doi:10.2196/28859
38. Cole S. Citations and the evaluation of individual scientists. *Trends Biochem Sci.* 1989;14(1):9-13. doi:10.1016/0968-0004(89)90078-9
39. Sauvayre R. De la diffusion scientifique aux normes de scientificité. Habilitation à diriger des recherches Dissertation. Sorbonne University; 2020.
40. Hicks D, Wouters P, Waltman L, De Rijcke S, Rafols I. Bibliometrics: The Leiden Manifesto for research metrics. *Nature.* 2015;520(7548):429-431. PMID:25903611 doi:10.1038/520429a
41. Bennett CM, Baird AA, Miller MB, Wolford GL. Neural correlates of interspecies perspective taking in the post-mortem Atlantic salmon: an argument for multiple comparisons correction. Poster session at the 15th Annual Meeting of the Organization for Human Brain Mapping. 2009 June 18-23, San Francisco, CA. <http://prefrontal.org/blog/2009/06/human-brain-mapping-2009-presentations/> [accessed Oct 2, 2017]
42. Bennett CM, Baird AA, Miller MB, Wolford GL. Neural correlates of interspecies perspective taking in the post-mortem Atlantic salmon: an argument for proper multiple comparisons correction. *JSUR.* 2010;1(1):1-5. [FREE Full text]
43. Harzing A-W. Accuracy: PoP vs GS for journals. Harzing.com. 2017 Oct 20. <https://harzing.com/resources/publish-or-perish/tutorial/accuracy/pop-vs-gs-for-journals> [accessed Feb 15, 2018]
44. Thelwall M, Kousha K. ResearchGate versus Google Scholar: which finds more early citations? *Scientometrics.* 2017;112(2):1125-1131. doi:10.1007/s11192-017-2400-4
45. Adams D. Google Scholar. Harzing.com. 2016 Dec 20. <https://harzing.com/resources/publish-or-perish/manual/using/data-sources/google-scholar> [accessed Mar 19, 2019]
46. Cole S, Cole JR. Scientific output and recognition: a study in the operation of the reward system in science. *Am Sociol Rev.* 1967;32(3):377-390. [PMID:6046811] doi:10.2307/2091085
47. Price DJ de S. Networks of scientific papers: the pattern of bibliographic references indicates the nature of the scientific research front. *Science.* 1965;149(3683):510-515. doi:10.1126/science.149.3683.510

Sauvayre R, (in press) "What types of errors are hiding in Google Scholar data? A case study", *Journal of Medical Internet Research*.

48. About. Google Scholar. Google Scholar. <https://scholar.google.com/intl/en/scholar/about.html> [accessed Jan 10, 2018]
49. Inclusion Guidelines for Webmasters. Google Scholar. <https://scholar.google.com/intl/en/scholar/inclusion.html#content> [accessed Jan 10, 2018]
50. Google Scholar Search Tips. Google Scholar. <https://scholar.google.com/intl/en/scholar/help.html#coverage> [accessed Jul 12, 2021]
51. Adams D. Results list. Harzing.com. 2016 Dec 20, 2016. <https://harzing.com/resources/publish-or-perish/manual/using/query-results/results-list> [accessed Feb 15, 2018]
52. Guagliardo P, Libkin L. Correctness of SQL Queries on Databases with Nulls. SIGMOD Rec. 2017;46(3):5-16. doi:10.1145/3156655.3156657
53. Baneyx A. "Publish or Perish" as citation metrics used to analyze scientific output in the humanities: international case studies in economics, geography, social sciences, philosophy, and history. Arch Immunol Ther Exp (Warsz). 2008;56(6):363-371. [PMID:19043670] doi:10.1007/s00005-008-0043-0
54. Bar-Ilan J. Citations to the "Introduction to informetrics" indexed by WOS, Scopus and Google Scholar. Scientometrics. 2010;82(3):495-506. doi:10.1007/s11192-010-0185-9
55. Harzing A-W. Accuracy. Harzing.com. 2016 Dec 20, 2016. <https://harzing.com/resources/publish-or-perish/manual/using/query-results/accuracy> [accessed Feb 15, 2018]
56. Basu A, Malhotra D, Seth T, Muhuri P. Global Distribution of Google Scholar Citations: A Size-independent Institution-based Analysis. J. Scientometr. Res.. 2019;8(2):72-78. PMID:23948488 doi:[10.5530/jscires.8.2.12](https://doi.org/10.5530/jscires.8.2.12)