

SCIENTIFIC REPORTS



OPEN

SciRide Finder: a citation-based paradigm in biomedical literature search

Adam Volanakis  & Konrad Krawczyk

There are more than 26 million peer-reviewed biomedical research items according to Medline/PubMed. This breadth of information is indicative of the progress in biomedical sciences on one hand, but an overload for scientists performing literature searches on the other. A major portion of scientific literature search is to find statements, numbers and protocols that can be cited to build an evidence-based narrative for a new manuscript. Because science builds on prior knowledge, such information has likely been written out and cited in an older manuscript. Thus, Cited Statements, pieces of text from scientific literature supported by citing other peer-reviewed publications, carry significant amount of condensed information on prior art. Based on this principle, we propose a literature search service, SciRide Finder (finder.sciride.org), which constrains the search corpus to such Cited Statements only. We demonstrate that Cited Statements can carry different information to this found in titles/abstracts and full text, giving access to alternative literature search results than traditional search engines. We further show how presenting search results as a list of Cited Statements allows researchers to easily find information to build an evidence-based narrative for their own manuscripts.

More than 60,000 articles are deposited in PubMed each month, making literature search an increasingly difficult task¹. A typical literature query consists of keyword-based search by services such as Google Scholar, PubMed, Scopus or Web of Science^{2–4}. The results typically consist of a list of titles and abstracts from documents that contain the query keywords. The scientist is then tasked with parsing through an extensive list of results, to extract information directly from titles/abstracts or to follow a link to the full document.

As such literature search can be burdensome, intelligent text mining of scientific publications has been seen as an alternative for extracting and organizing information from the ever-growing PubMed collection⁵. Sites such as iHOP or Chilobot mine field-specific knowledge by collating information regarding biomolecules from millions of PubMed publications^{6,7}. Less field-specific services such as COLIL, provide a service showing comments in more recent research on older manuscripts⁸. These tools demonstrate that strategic text mining and intelligent filtering can lead to new, more efficient tools for biomedical literature search.

Strategic text mining can be used to separate relevant information from tangential text. For instance, because of legal restrictions, typical literature search engines operate on the remit of copyright-available titles and abstracts alone, whereas full text contains more pertinent information⁹. For instance, tools such as Biotext or Yale Image Finder allow searches in Figure or Table captions alone in order to identify relevant information only^{10,11}. To understand what information is potentially irrelevant, it is necessary to identify portions of searchable documents that can be of more interest to the person performing literature search.

One major aim of literature search is to identify earlier papers to support the narrative presented in a new manuscript being written. Such narrative is constructed by citing findings, numbers, data and techniques from previous publications¹². Such pieces of text are easily identifiable in scientific manuscripts since they are annotated with references to prior peer-reviewed publications which support the statement being made^{12,13}. Therefore such statements in publications on previous literature, which we here call Cited Statements, offer succinct comments on prior art, whose information content is powerful enough to be used for article summarization¹⁴.

Here, we propose a simple strategy of improving text mining and literature search by creating a biomedical search corpus, which is constrained to such Cited Statements only. We show that Cited Statements can carry different information-retrieval data to these found in titles, abstracts and full text of documents they refer to, demonstrating that this methodology does not simply recapitulate information currently available in scientific

SciRide.org, Boston, MA, USA. Correspondence and requests for materials should be addressed to K.K. (email: konrad@sciride.org)

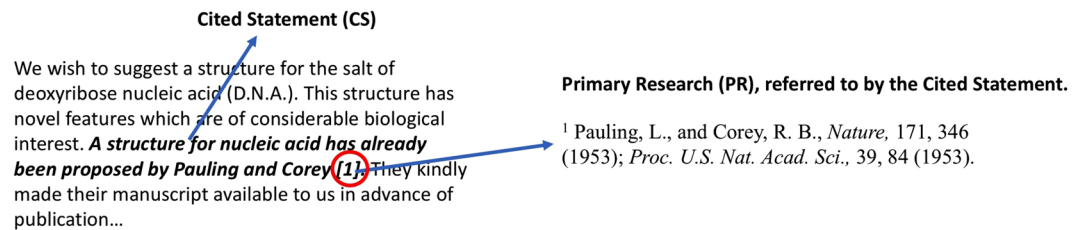


Figure 1. Example of a Cited Statement (CS) and Primary Research (PR). The CS is shown in bold in the excerpt on the left. PR which is referred to by the CS, is shown on the right. The text in the image was taken from the seminal paper by Watson and Crick in 1953, entitled 'A Structure of Deoxyribose Nucleic Acid'.

search engines. Furthermore, we show how presenting results in the form of Cited Statement text, can offer easy access to information in several literature-search scenarios. We hope that our service, available at finder.sciride.org will offer a streamlined way for biomedical scientists to build evidence-based narratives for their own manuscripts.

Results/Applications

SciRide Finder methodology. SciRide Finder offers an orthogonal literature search strategy to platforms such as PubMed or Google Scholar by focusing on Cited Statements only. In this manuscript, we refer to Cited Statement (CS) as any sentence from a peer reviewed publication containing citations to other manuscripts, which we refer to as Primary Research (PR) (Fig. 1).

We have extracted the CSs from all PubMed/Medline indexed documents where the copyright allowed for data mining and reproduction. To the best of our knowledge, the most suitable corpus for this task is the Open Access PubMed Central (OA PMC) dataset. It is a collection of open access journals from PubMed/Medline in standardized format. At the time of writing, there were approximately 1.7 m publications in the OA PMC dataset, which is 6% of a total of more than 26 m publications indexed in PubMed (or 15 m if only citations with abstracts are to be considered⁹).

The OA PMC dataset downloaded via the NCBI ftp service forms the core of our dataset. Nevertheless, the ~1.7 m OA PMC articles are only a subset of more than 4 m web-formatted documents available via PMC¹⁵. There are more than 2 m articles published after 1980 which are accessible via PMC 'eyes-only' subject to strict restrictions on machine access and heterogeneous publisher copyrights. We therefore extract such data manually if and only if the copyright situation is unambiguous.

We have set up a pipeline to collect data from the OA PMC and other publications in the public domain where copyright allows it (see Materials and Methods). At the time of writing, our data collection encompasses 1,786,322 peer-reviewed articles contributing 43,326,402 CSs. We make this corpus accessible via efficient Lucene-based search system as described in Materials & Methods. Here, we argue that our CS-based search system is a new literature search paradigm, distinct from traditional title/abstract and full text based methods.

SciRide Finder as an alternative biomedical literature search platform. Each search engine can be divided into (1) user input, (2) the search corpus, (3) retrieval methodology and (4) presentation of results. Here, we focus on demonstrating that creating a search engine whose corpus is constrained to CSs only can create a valuable service for the scientific community. As such we put full emphasis on (2) the search corpus and (4) the presentation of results. Input into SciRide Finder is comparable to any other search engine as it accepts keywords entered via an appropriate web interface. As such the retrieval methodology in our search engine is basic compared to major search engines (even though we cannot comment on specifics as they are typically not transparent about their algorithms²). However, use of only basic retrieval methodology reinforces the power of our approach by shifting the responsibility for the reliability of the service to where it differs significantly: constraining the search corpus to CSs only and presenting them as results.

The first major difference between traditional and CS-based search is the corpus employed to identify documents. In traditional systems, documents are retrieved if the query keywords are found within text of title, abstract or full text. In search engines such as Google Scholar or Semantic Scholar this can produce results where the keywords are widely separated (See Supplementary Examples 1–7). In contrast, the short nature of CSs enforces the proximity, focusing the top results (See Supplementary Examples 1–7). This is a conscious instrument that is supposed to increase the quality (precision) of the results, which typically reduces recall. Nevertheless, each CS comes from one document and carries at least one other document attached to it (the Primary Research). Therefore, searching by CSs identifies PR documents indirectly by text contained in other papers (Fig. 2). CS offers an alternative commentary on the PR, by scientists who were generally not involved in the original study. To prove this point, we demonstrate that CSs can hold alternative information to titles/abstracts and full text of PR documents, described in section 2.3.

The second major difference between traditional and CS-based search is the presentation of results. In traditional literature search systems, results are presented as titles/abstract, more seldom as full text excerpts. In contrast, CS-based search returns the concise pieces of text which cite other documents. In this capacity, it identifies the information which was used to build the evidence-based narrative for a manuscript: scientific statements, numbers, data and techniques, all supported by prior publications. In many scenarios, these are the pieces of text

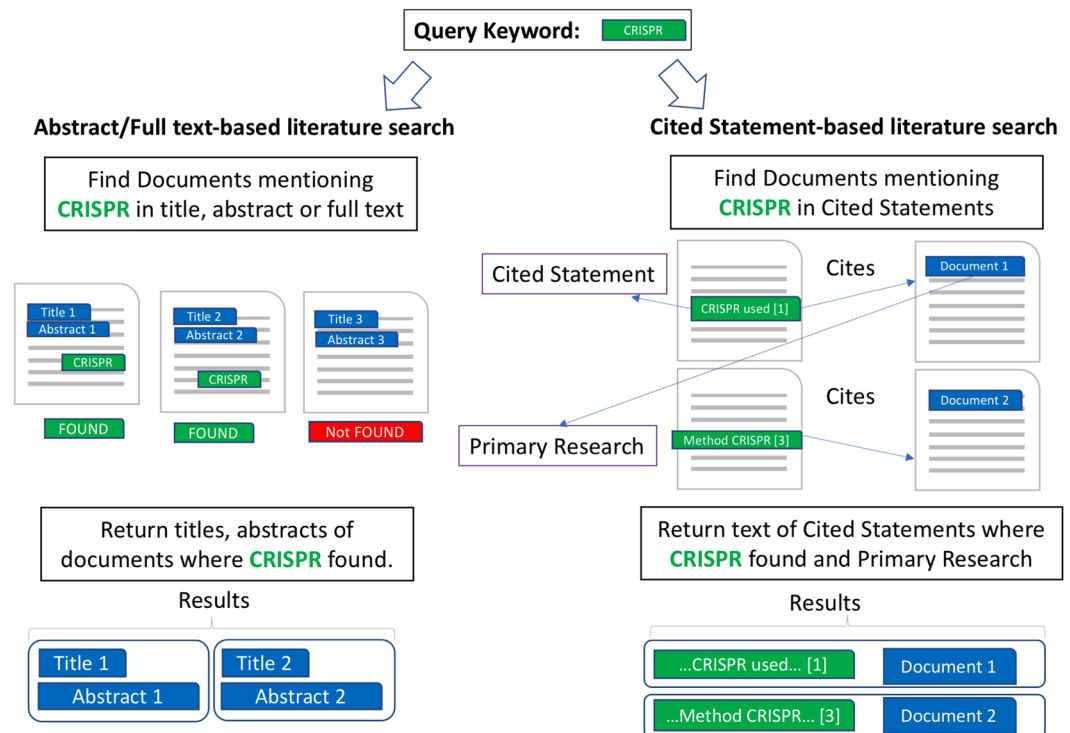


Figure 2. Contrasting the traditional literature search and Cited Statement-based literature search. Traditional literature search systems identify documents to be retrieved by keyword hits within them and present titles and abstracts as results (left). Cited Statement-based search identifies Primary Research documents by text from other publications and presents the citing text.

scientists look for in the first place to build an evidence-based structure for their own manuscripts. To exemplify this, we present possible applications of presenting results as CSs in section 2.4.

Cited Statements can hold different information on documents to titles/abstracts and full text.

We argue that CS-based search offers a novel way of retrieving documents, that can yield orthogonal results as compared to traditional search strategies. For this to be true, CSs must offer distinct information-retrieval data on the PR that would not normally be available by examining titles, abstract or even full text of PR document.

To quantify this, we identified 691,354 documents where we had CSs in our database referring to PR documents whose full text is available for text mining. For a given CS, we measured how many normalized words (stemmed, case-folded etc.) cannot be found in the title/abstract and full-text of the PR documents, which we refer to as ‘difference’. For each of 691,354 documents, we have identified the maximally different (as percentage) CS with respect to the title/abstract and full text of the PR document (Fig. 3). These results demonstrate that for 83% of our 691,354 PR documents, there exists a CS which is at least 50% different to the title/abstract of the PR document. For 61% of PR documents, there exists a CS which is at least 25% different to the PR document full text. Therefore, for a significant proportion of publications, there exists a CS which offers information that would not be available through a title/abstract or full-text search on the PR document.

We have also found that CSs tend to be different from titles/abstracts not only in the extreme as above, but also on average (see Supplementary Fig 1). These results demonstrate that CSs can contribute different information on the PR manuscripts than title/abstract and full text. Therefore, corpus constrained to CSs only, can provide orthogonal results to those offered by search engines which identify documents directly by their titles/abstracts and full text. This can offer a new paradigm in search, identifying documents indirectly by text contained in other documents, which can also provide a metric for the importance and the impact of this work in the scientific field.

Aside from the possible new paradigm in document retrieval, relationship between CS and PR text can offer valuable insight into our understanding of communication of knowledge. Each CS is supposed to be a statement being supported by evidence in PR. Natural language analysis of the relationship between CS and PR can therefore inform how evidence (PR) fuels making of scientific conclusions communicated via CS. For instance, the difference in content between CS and PR in certain cases can be indicative of whether PR was refuted or not¹⁶. Thus, ability to find whether a given statement was refuted or not is only one of the search scenarios where CS search can offer immediate insight and we describe the other orthogonal applications of this search paradigm in the next section.

Presenting the results as CS – search-scenario-driven comparison to other search engines.

Search results from SciRide Finder do not consist of titles, abstracts or full-text excerpts as is typically the case in other services, such as Semantic Scholar, PubMed or Google Scholar. Instead, in response to a query, we present

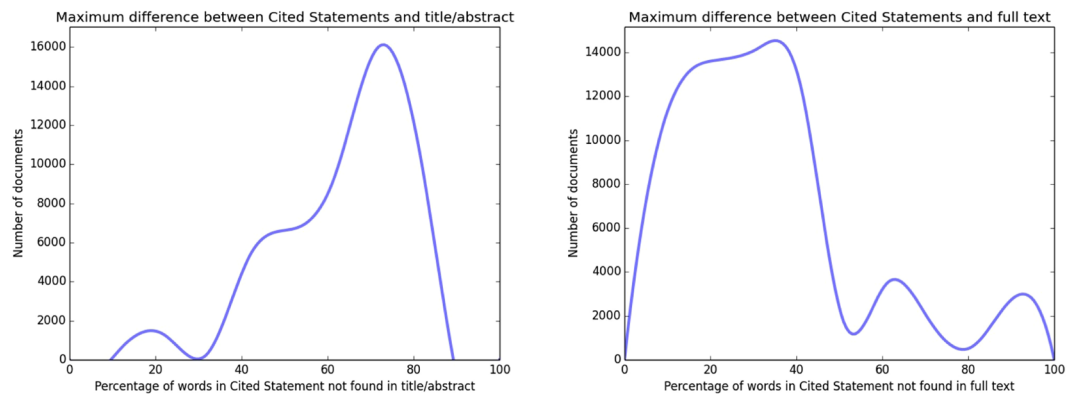


Figure 3. Maximum difference between CS and text of PR. For each of 691,354 documents, we identify the maximally different CS with respect to title/abstract (left) and full text (right) of PR. This stands to demonstrate that there are many PR documents, where there exists a significantly textually different CSs referring to them.

Source Paper: G3: Genes | Genomes | Genetics, Efficient CRISPR-Mediated Post-Transcriptional Gene Silencing in a Hyperthermophilic Archaeon Using Multiplexed crRNA Expression

"Cas9 was also elegantly engineered to repress or activate gene transcription in bacteria, and even eukaryotes, through the use of nuclease-deficient Cas9 variants that are guided to genomic sites to physically block regulatory enzymes or the transcription complex itself (;)."

Papers cited by the source:

- PMID: 25494202
- PMID: 23849981 (Cited by 2 result documents), 2nd most
- PMID: 23452860 (Cited by 5 result documents), 1th most

Figure 4. SciRide Finder search example for “Cas9 block” statements. Each result consists of the CS (B), the title of the paper that the statement appears in (A) and the PR sources that the statement is based on (C).

a list of CSs, PR documents and papers where the information was found (Fig. 4). This approach reduces the amount of text being presented on the screen, constraining information only to short statements. By nature of CS, such pieces of text refer to very specific information, such as numbers, datasets, protocols and concepts communicated by other papers. In this respect, CS corpus is a strategically constrained subset of biomedical literature as opposed to the big search engines such as Google Scholar, Semantic Scholar or PubMed that provide a universal literature search. However, one-size-fits-all methodologies might not provide optimal solutions¹⁷ and data rather than retrieval methodology is key¹⁸. Therefore we argue that our choice of search corpus combined with very basic retrieval methodology shows the power of the concept of CS-based search in a set of specific scenarios. We compile the detailed search scenarios in Supplementary Section 2 and we comment on them below:

Identifying citations supporting general knowledge. It is often problematic to identify citations supporting a well-known fact. For instance, the controversy surrounding the link between vaccines and autism is well known, but identifying studies discrediting it is not trivial. Searching SciRide Finder for “autism” “vaccine” and “discredited” would return results that debunk the notorious 1998 publication in Lancet by Wakefield and colleagues. Search scenarios of this type can be found in Supplementary data (examples 1 and 2).

Identifying numbers. Conducting research and writing scientific publications requires careful quantification that is only possible by finding the right numbers. Identifying the right numbers might be problematic if one does not know right away the publications they are reported in. Nevertheless, commonly used numbers are reported in the literature supported by citations – hence being found in CSs. To exemplify this, we offer examples of search scenarios identifying numbers for popular techniques and quantification of biological phenomena: error rates and read lengths in Illumina technology, or the current resolution of super-resolution microscopy (Supplementary Examples 3 and 4).

Identifying datasets. Datasets used in publications are rarely cited in titles and abstracts, rather being hidden in Methods sections. As an example, searching SciRide Finder for the terms “ChIP-seq” “HeLa” and “Pol II” would return the publications that have used datasets of RNA Polymerase II (Pol II) Chromatin Immunoprecipitation sequencing (ChIP-seq) experiments in HeLa cells, but also the original source of these datasets (Supplementary Figure 2), thus facilitating the retrieval of the datasets.

This scenario is applicable to any dataset that is widely commented on and thus cited in the literature. As an example, consider the usage of protein structures. These are typically downloaded from the Protein Data Bank, but to those not in the structural field it might not be obvious. We present example 5 in Supplementary material demonstrating how using SciRide Finder it is possible to identify the source of structural data.

Protocol/methodology technique identification. Similarly to datasets, specific techniques used are rarely available in abstracts and titles. Nevertheless, identifying publications which employ a given technique or software is indispensable for the retrieval of protocols. For example, the newly developed CRISPR/Cas9 genome editing method has an alternative usage as a block for gene transcription. A PubMed search for the terms “Cas9” and “Block” returns 47 publications (at the time of writing) and there is no way of knowing how and in what context this method was used in each paper without reading all manuscripts. The same search in SciRide Finder (Fig. 4) provides a list of publications where this technique was used in context. This allows us to identify publications describing the method, the theory behind it, protocols used, and the original research.

This methodology can be extrapolated to any common technique used and thus cited in the literature, providing immediate application context for the protocol. For instance, we present how SciRide Finder can be used to identify the context for the use of the standard drug-design methodology of Lipinski Rule of Five or how Deep Learning is used in biology (examples 6 and 7). In contrast, other search engines are more likely to return results on the respective techniques and methods rather their applications.

Mapping natural language connections between citing text and publications. SciRide Finder allows for searching for two or more terms appearing together, their context and the original research. For example, a CS-search for ‘mRNA export’ and ‘transcription’ would identify only the statements in which the two keywords appear together (Supplementary Figure 3) and provide original research context for the claim being made. Mapping such connections between CS and PR can be of particular interest for creating knowledge maps by the text mining community which among other things can be used to assess the veracity of claims being made in original publications¹⁶.

Conclusion

We have created a biomedical search service based on information content from CSs. These short pieces of text build the evidence-based narrative for a given manuscript and provide a reflection of knowledge contributed by previous publications. We have shown that constraining the search corpus to CS only, can be a viable alternative to conventional search methodologies as it provides different information from titles/abstracts and full text. Furthermore, presenting results as CSs, is beneficial in many areas of scientific literature search, whose major part is aimed at identifying evidence-based pieces of text to be used in future publications.

Established search methodologies, such as Google Scholar, aim to index all the information available on documents even if the publication itself is not in the public domain. On the contrary, our service indexes only a very well-defined subset of the full-text articles, namely the CSs. We currently extracted ~43 m CSs which contain comments on 34% (or 57%, if publications without abstracts are to be omitted⁹) of all of PubMed articles. This proportion should only increase as more publications become open access and repositories become legally and technically unified for systematic text mining^{15,19}. Despite the small scale of the corpus we used and the basic retrieval technology SciRide Finder offers useful search experience (previous section), offering argument in favor of CS-based search.

In summary, our system introduces an open-access, CS-only paradigm in literature search. Current manifestation of this paradigm, SciRide Finder, offers an orthogonal approach to reduce the burden currently associated with specific information retrieval in biomedical literature. We hope that our service will facilitate the efforts of researchers looking for Cited Statements, to build an evidence-based narrative for their own publications.

Materials and Methods

Data Collection for the base system – PMC Open Access Dataset. The OA PMC corpus was downloaded from the NCBI FTP website (<ftp.ncbi.nlm.nih.gov>) and divided into sentences using Natural Language Toolkit (nltk.org) and a custom set of heuristics, such as splitting text on terminal period ‘.’; removing the ‘.’ from short-hands such as ‘*et al.*’, ‘*ca.*’ and normalizing the scientific names (‘*H. Pylori*’). We identified the sentences containing citations as these having the <xref> tag with attributes pointing to references section (as opposed to non-bibliographic elements such as Tables and Figures). Rules were created for special cases where the citation pertaining to a sentence occurs after its terminal period. Each CS derived in this way contains the citing sentence, identifiers of cited articles (DOI or Pubmed ID) and the metadata on the manuscript it was derived from (journal title and article title). The system was set up to perform updates of this base dataset on a monthly basis.

Data Collection Beyond the PMC Open Access Dataset. We augment the information from the base-dataset manually from the ‘eyes-only’ documents where there was an unambiguous copyright situation on reproducing pieces of work in a normal citation scenario. Furthermore, it is sometimes possible to find the author-submitted PDF version of the document. These are documents available via platforms such as BioArxiv or author homepages. Whenever we could not identify a PMC version of an article, we attempted a PDF doi search. When a PDF document was found in such a way, we extracted information from it using PDFExtract tool from CiteSeer. PDFExtract is a utility which is capable of extracting portions of a PDF-formatted scientific document and present them in machine-readable plain text. Since information presented in such format is still very heterogeneous, we had to create different sets of rules to interpret the PDF-extracted plain text, which mostly involved detecting if the citations are number-based or author-name based.

Text Retrieval System. The Cited Statements are stored for rapid extraction in a Lucene-based system which was previously shown to be a robust search engine for biomedical applications²⁰. Since scientific documents are by and large written in English²¹ we have employed standard English analyser and stemmer as parameters of retrieval. We only perform searches on the text of the CS record, disregarding metadata of the full article it was retrieved from.

Documents are retrieved given a set of keywords to match the text of the CS. A post-processing step after document retrieval is introduced, where we count the number of shared citations between resulting documents. The documents are sorted in descending order firstly by the relevance score of the Lucene system (normalized to one decimal point) and secondly by the number of shared citations. This assures that statements on highly cited papers which are similar within the normalized value of the text-relevance score, are displayed first. The literature search service is available as a web service at <http://finder.sciride.org>. The text-mining of the CS corpus is available through an API which is described on the website.

Information content comparison. We measured how different the CSs are to PR document titles, abstracts and citations. Since a typical search engine operates on the remit of keywords, we have created a textual fingerprint for each CS, title/abstract and full text. Each fingerprint was a set of case-folded, stemmed and stop-word-free normalized words without duplicates.

For each PR document, we have collected three elements: its title/abstract, full text and a list of CSs referring to it. We have created a textual fingerprint for each title/abstract, full text and each CS, which was supposed to emulate a typical corpus employed by an information retrieval system.

To produce a fingerprint for a given piece of text, we split it into word tokens using the NLTK toolkit. We case-folded each word, and removed any punctuation (keeping special symbols such as Greek letters). We removed all stop-words (as defined by the NLTK corpus). Finally each word was stemmed so as to minimize mismatches in subtle inflection forms²². We did not keep word duplicates, thus for each text element (such as title/abstract), this resulted in a non-redundant list of normalized words.

A typical information retrieval algorithm can be expected to perform such text-normalization operations on a given document. Thus, it is reasonable to assume that if text-normalized fingerprints share many words, the information retrieval algorithm would treat them as contributing similar information and yield similar results. Therefore, the number of different normalized words between CS and PR title/abstract and full text was taken as a measure if CS contribute new information on PR.

Comparing two fingerprints (e.g. CS versus full-text) consisted of counting how many text-normalized words are found in one fingerprint but not the other.

Availability. The findings in this manuscript are available in three different guises: search engine, API and bulk data download. The search engine is designed to demonstrate a CS-based literature search in practice and is available at <http://finder.sciride.org>. Separately, we make the API available at endpoint located at <http://finder.sciride.org/api>. The API is designed for those who would like to use our current search-engine setup to perform text mining as it allows machine-readable retrieval of results in response to keywords. We also make the bulk of our data available formatted in json format for large-scale data mining experiments. The bulk data are available upon contacting authors at konrad@sciride.org.

References

- Neylon, C. & Wu, S. Article-level metrics and the evolution of scientific impact. *PLoS Biology* **7** (2009).
- Beel, J. & Gipp, B. Google Scholar's Ranking Algorithm: An Introductory Overview. *12th Int. Conf. Sci. Inf.* **1**, 230–241 (2009).
- Ostell, J. In *The NCBI Handbook* 1–6 (2002).
- Jacso, P. As we may search - Comparison of major features of the Web of Science, Scopus, and Google Scholar citation-based and citation-enhanced databases. *Current Science* **89**, 1537–1547 (2005).
- Beck, J. & Sequeira, E. In *NCBI Handbook* 1–17 (2013).
- Fernández, J. M., Hoffmann, R. & Valencia, A. IHOP web services. *Nucleic Acids Res.* **35** (2007).
- Chen, H. & Sharp, B. M. Content-rich biological network constructed by mining PubMed abstracts. *BMC Bioinformatics* **5**, 147 (2004).
- Fujiwara, T. & Yamamoto, Y. Colil: a database and search service for citation contexts in the life sciences domain. *J. Biomed. Semantics* **6**, 38 (2015).
- Westergaard, D., Stærfeldt, H.-H., Tønsberg, C., Jensen, L. J. & Brunak, S. Text mining of 15 million full-text scientific articles. *doi.org/162099*, <https://doi.org/10.1101/162099> (2017).
- Hearst, M. A. *et al.* BioText Search Engine: Beyond abstract search. *Bioinformatics* **23**, 2196–2197 (2007).
- Xu, S., McCusker, J. & Krauthammer, M. Yale Image Finder (YIF): A new search engine for retrieving biomedical images. *Bioinformatics* **24**, 1968–1970 (2008).
- Abu-Jbara, A. & Radev, D. Reference scope identification in citing sentences. *12 Proc. 2012 Conf. North Am. Chapter Assoc. Comput. Linguist. Hum. Lang. Technol.* 80–90 (2012).

13. Qazvinian, V. & Radev, D. R. Identifying non-explicit citing sentences for citation-based summarization. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics* 555–564, doi:Association for Computational Linguistics (2010).
14. Qazvinian, V. & Radev, D. R. Scientific Paper Summarization Using Citation Summary Networks. In *COLING '08 Proceedings of the 22nd International Conference on Computational Linguistics* 689–696, <https://doi.org/10.3115/1599081.1599168> (2008).
15. Piwowar, H. *et al.* The State of OA: A large-scale analysis of the prevalence and impact of Open Access articles. *PeerJ Prepr.* <https://doi.org/10.7287/peerj.preprints.3119v1> (2017).
16. Grabitz, P., Lazebnik, Y., Nicholson, J. & Rife, S. Science with no fiction: measuring the veracity of scientific reports by citation analysis. *bioRxiv* 172940, <https://doi.org/10.1101/172940> (2017).
17. Wolpert, D. *No free lunch theorems for search*. Technical Report SFI-TR-95-02-010, <https://doi.org/10.1145/1389095.1389254> (1995).
18. Banko, M. & Brill, E. Scaling to very very large corpora for natural language disambiguation. in. *Proceedings of the 39th Annual Meeting on Association for Computational Linguistics - ACL* **01**, 26–33, <https://doi.org/10.3115/1073012.1073017> (2001).
19. Piwowar, H. A., Day, R. S. & Fridsma, D. B. Sharing detailed research data is associated with increased citation rate. *PLoS One* **2**, (2007).
20. Yu, H. *et al.* Development, implementation, and a cognitive evaluation of a definitional question answering system for physicians. *J. Biomed. Inform.* **40**, 236–251 (2007).
21. Ferguson, G., Erez-Llantada, C. & Plo, R. O. English as an international language of scientific publication: a study of attitudes. *World Englishes* **30**, 41–59 (2011).
22. Porter, M. F. An algorithm for suffix stripping. *Program* **14**, 130–137 (1980).

Author Contributions

A.V. and K.K. designed the experiments wrote the manuscript and prepared the figures. K.K. wrote the text mining algorithms and created the finder.sciride.org website.

Additional Information

Supplementary information accompanies this paper at <https://doi.org/10.1038/s41598-018-24571-0>.

Competing Interests: The authors declare no competing interests.

Publisher's note: Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2018