

unarXive 2022: All arXiv Publications Pre-Processed for NLP, Including Structured Full-Text and Citation Network

Tarek Saier
tarek.saier@kit.edu
Karlsruhe Institute of Technology
Karlsruhe, Germany

Johan Krause
johan.krause@student.kit.edu
Karlsruhe Institute of Technology
Karlsruhe, Germany

Michael Färber
michael.farber@kit.edu
Karlsruhe Institute of Technology
Karlsruhe, Germany

ABSTRACT

Large-scale data sets on scholarly publications are the basis for a variety of bibliometric analyses and natural language processing (NLP) applications. Especially data sets derived from publication's *full-text* have recently gained attention. While several such data sets already exist, we see key shortcomings in terms of their domain and time coverage, citation network completeness, and representation of full-text content. To address these points, we propose a new version of the data set unarXive. We base our data processing pipeline and output format on two existing data sets, and improve on each of them. Our resulting data set comprises 1.9 M publications spanning multiple disciplines and 32 years. It furthermore has a more complete citation network than its predecessors and retains a richer representation of document structure as well as non-textual publication content such as mathematical notation. In addition to the data set, we provide ready-to-use training/test data for citation recommendation and IMRaD classification. All data and source code is publicly available at <https://github.com/lllDence/unarXive>.

CCS CONCEPTS

• **Information systems** → **Information retrieval**; • **Computing methodologies** → **Information extraction**; *Language resources*; **Knowledge representation and reasoning**.

KEYWORDS

scholarly data, information extraction, citation network, \LaTeX

ACM Reference Format:

Tarek Saier, Johan Krause, and Michael Färber. 2023. unarXive 2022: All arXiv Publications Pre-Processed for NLP, Including Structured Full-Text and Citation Network. In *Proceedings of . ACM*, New York, NY, USA, 5 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 INTRODUCTION

Large data sets derived from the full-texts of academic publications are of ever-increasing importance. Beyond large-scale metadata, which is the basis for bibliometric analyses, research output quantification [9], and various applications such as trend detection [3], data sets reflecting the *full-text* content of papers have recently

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2023 Association for Computing Machinery.
ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00
<https://doi.org/XXXXXXX.XXXXXXX>

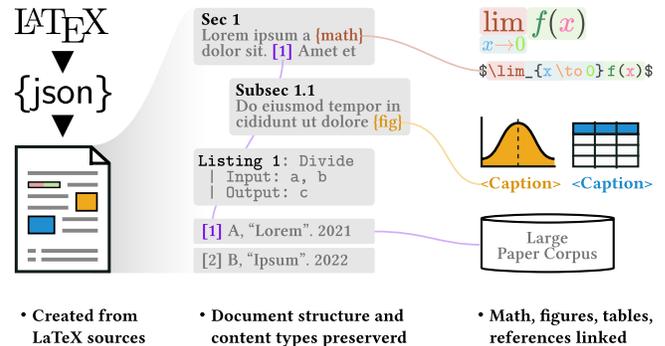


Figure 1: Schematic of our data set.

Created from arXiv.org \LaTeX sources, our data set preserves document *structure* (sections, subsections, ...) and *content types* (paragraphs, listings, ...). In-text positions of mathematical notation, figures, tables and citation markers are linked to \LaTeX math content, figure/table captions, and bibliographic references respectively. Bibliographical references are linked to the large paper corpus OpenAlex.

enabled more sophisticated analyses and applications, such as scientific document summarization [14], claim verification [26], and knowledge graph generation [13].

Key aspects of such data sets are (1) basic measures such as quality, size, and temporal as well as disciplinary coverage, (2) their citation network, and (3) handling of non-textual content. (1) Quality is affected by the source material (e.g. PDF or \LaTeX) and parsing method. (2) The citation network is important to allow for bibliometric analyses. (3) Non-textual content such as tables, figures, and mathematical notation often contain important information.

Across these key aspects, we see significant shortcomings in currently available data sets, as shown in Table 1. For example, (1) limited size (SciXGen), (2) omission of a citation network (arXMLiv), and (3) no or limited handling of mathematical notation (S2ORC, unarXive 2020).

To address these issues, we propose a new version of the data set unarXive, which comprises 1.9 M publication across several disciplines, includes a more complete citation network than its predecessors, and retains structured mathematical notation as well as table and figure captions (see Figure 1). Apart from the data set itself, we furthermore provide ready-to-use training and test data for two NLP tasks. Overall, we make the following contributions.

- We provide a 1.9 M document scholarly data set, containing structured full-text, annotated in-text citations, linked table and figure captions, structured mathematical notation, and a high quality citation network.

Table 1: Comparison of large data sets derived from paper full-texts

[†]Cit. network completeness is reported in two ways. “general”: the whole data set; not directly comparable. “compare”: for arXiv.org data from 1991–2020; directly comparable.

[‡]References in the PMC-OAS are partially linked to a mixed set of IDs (PubMed, MEDLINE, DOI) [7]. Therefore there is no single, comprehensive number for its completeness.

Data Set	Source		Citation Network [†]		Structured			Disciplines	Purpose
	Data	Format	general	compare	Doc.	Math.	# Docs		
CORE [17]	multiple	PDF	0%	-	×	×	>100 M	various	general NLP
S2ORC (PDF) [11]	multiple	PDF	69.4%	-	✓	×	12 M	various	general NLP
unarXive 2020 [20]	arXiv.org	L ^A T _E X	42.6%	42.6%	×	×	1.2 M	phys., maths, CS	general NLP
S2ORC (L ^A T _E X) [11]	arXiv.org	L ^A T _E X	31.1%	31.1%	✓	✓	1.5 M	phys., maths, CS	general NLP
arXMLiv [6]	arXiv.org	L ^A T _E X	0%	0%	✓	✓	1.6 M	phys., maths, CS	maths linguistics
SciXGen [4]	arXiv.org	L ^A T _E X	41.6%	-	✓	✓	205 k	CS	text generation
PMC-OAS [15]	PubMed	XML	mixed [‡]	-	✓	✓	3.3 M	biomedical	not NLP specific
unarXive 2022 (ours)	arXiv.org	L ^A T _E X	44.4%	44.4%	✓	✓	1.9 M	phys., maths, CS	general NLP

- We provide ready-to-use training/test data for the development and evaluation of approaches to two NLP tasks, namely citation recommendation and IMRaD classification.
- We distribute our data in accordance to the FAIR principles [27] and share our source code freely available under a permissive license.

2 RELATED WORK

In Table 1 we give an overview of related work. Excluded are data sets that are either just sets of PDFs, or only contain metadata.

CORE [17], while being very large, does not contain a citation network, nor is document structure preserved. S2ORC (PDF) [11] is second in size and, while not directly comparable due to different publications covered, has the most complete citation network. However, mathematical notation is only partially preserved as plain-text. unarXive 2020 [20] has the second highest citation network completeness in direct comparison, but lacks structured content.

The bottom part of the table are data sets with both document structure preserved and structured mathematical notation. S2ORC (L^AT_EX) [11] is a discontinued¹ subset of S2ORC and has a limited citation network, arXMLiv [6] offers the highest level of structure but no citation network, and SciXGen [4] is limited in size. The PMC-OAS [15] is comparable to unarXive 2022 in size and structure, but has a partial and mixed citation network.

Overall, unarXive 2022 has the most complete citation network as far as direct comparison is possible, preserves document structure as well as structured mathematical notation, and is the largest data set covering physics, mathematics and computer science.

3 APPROACH

We base our data set creation approach in part on S2ORC (L^AT_EX) and in part on unarXive 2020. This is motivated as follows.

As shown in Table 1, the majority of related data sets is based on paper’s L^AT_EX sources—which is less noise-prone than parsing PDFs [1]. Among these, S2ORC (L^AT_EX) provides well structured full-text content usable for a wide variety of applications (see Section 4.2), while arXMLiv and SciXGen are optimized for special

purposes. We therefore base our structured document representation on S2ORC (L^AT_EX). Regarding the citation network, however, unarXive 2020 achieves the most high quality results in direct comparison among existing data sets. We therefore base our citation network creation on unarXive 2020.

Regarding both S2ORC (L^AT_EX) and unarXive 2020, we don’t just copy, but also improve upon the existing work. To furthermore provide an up-to-date data set, we use as source data all papers on arXiv.org up until the end of 2022.

Conceptually, our overall data set creation process can be broken down into two major steps, namely document parsing and reference linking. In the following these are described in more detail.

3.1 Document Parsing

To convert the L^AT_EX source of a paper into a format that is well suited for NLP applications and analyses, we follow S2ORC (L^AT_EX) and unarXive 2020 and perform the following three steps. First, we flatten the paper’s L^AT_EX source into a single .tex document using latexpand.² Next, we use the tool Tralics³ to convert the L^AT_EX source into XML. In the last step, we create an easy to handle JSON structure from the XML.

We adapt and extend the JSON structure of S2ORC as shown in Table 2. Adding paper metadata facilitates easier analyses (e.g. for specific or across disciplines). Including information on section numbers and types reflects the document structure more closely (e.g. the nesting structure is not lost). Retaining URLs from embedded links helps with reference linking (see Section 3.2).

We mark the position of citation markers, tables, figures, and mathematical notation within the running text, and link citations markers to their references, tables and figures to their captions (i.e., textual surrogates of their content), and mathematical notation to its original L^AT_EX content.

3.2 Reference Linking

To add a citation network to the data set, bibliographical references—which at this point are just raw strings of text—need to be associated with the cited documents they’re referencing. We follow the

¹Last release including the L^AT_EX subset is 2019-09-28, see <https://github.com/allenai/s2orc> (accessed 2023/02/12).

²See <https://ctan.org/pkg/latexpand>.

³See <https://www.sop.inria.fr/marelle/tralics/>.

Table 2: Extension of S2ORC format

Entity	S2ORC data	Added data
Paper	<ul style="list-style-type: none"> • ID • abstract • full-text (list of paragraphs) • bibliographic references 	<ul style="list-style-type: none"> • Metadata (title, list of authors, discipline, license, version history)
Paragraph	<ul style="list-style-type: none"> • Section title • text 	<ul style="list-style-type: none"> • Section number • Section type (e.g. <i>section</i>, <i>subsection</i>) • Content type (e.g. <i>paragraph</i>, <i>listing</i>, <i>proof</i>)
Bibliographic reference	<ul style="list-style-type: none"> • Parsed reference • ID of cited document 	<ul style="list-style-type: none"> • Raw reference string • List of contained arXiv IDs • List of embedded links (i.e. URLs of clickable links not rendered as text when viewing the document)

methodology of unarXive 2020 and link references to a large corpus of publication metadata. To do this, references are first parsed to determine the contained information (title, authors, year, venue, etc.), which is then matched against the paper records in the large metadata corpus. For these two steps, we make the following changes and improvements of the unarXive 2020 approach.

Parsing. unarXive 2020 utilizes the tool Neural Parscit [18] for reference parsing and furthermore uses a heuristic procedure to determine identifiers such as DOIs or arXiv IDs found within reference string. We use GROBID [12], a more commonly used and actively developed tool. Additionally, we extend the identifier determination heuristics to be more robust and versatile by refining matching patterns and extending them to more citation styles.

Matching. unarXive 2020 matches references to paper records in the Microsoft Academic Graph (MAG) [23], which is no longer publicly available. Instead of the MAG, we use OpenAlex [19], the MAG’s open successor provided by the nonprofit organization OurResearch.⁴ Choosing OpenAlex allows us to also match references to recent papers, which would not be contained in legacy versions of the MAG. Additionally, the fact that OpenAlex paper records contain a variety of identifiers (e.g. DOI and PubMed ID) facilitates combined and comparative analyses of our data with others. Furthermore, OpenAlex has been deemed better suited for bibliographic analyses than the MAG [22].

4 RESULTS

In the following, we first present key statistics of our proposed data set. Following that, we explain how the data set can be used for analyses as well as the development of NLP applications, and introduce training/test data for two NLP tasks. Lastly, we describe how the data set is distributed to facilitate easy adoption by the community of researchers and practitioners.

⁴See <https://ourresearch.org/>.

4.1 Data Set

Our data set comprises *1,881,346 papers*, which contain a combined *182,586,547 paragraphs*, *63,367,836 references* and *133,744,613 in-text citation markers*. The distribution across disciplines is 57% physics, 20% mathematics, 17% computer science, and a combined 5% for others. We are able to link 28,135,565 references (44.4%) and 64,547,944 (48.3%) in-text citation markers to OpenAlex. As shown in Table 1, this makes our citation network more complete than that of existing data sets.

In Listing 1 we show an excerpt of our document representation for one paper, showcasing the extracted plain text and structured content.

```

/* ----- example paper (arXiv:2105.05862) ----- */
{ "paper_id": "2105.05862",
  "metadata": { ... },
  "abstract": { ... },
  "body_text": { ... },
  "ref_entries": { ... },
  "bib_entries": { ... }
/* ----- one of the sections in body_text ----- */
{ "section": "Memory wave form",
  "sec_number": "2.1",
  "sec_type": "subsection",
  "content_type": "paragraph",
  "text": "The gauge choice leading us to this solution does not fix
completely all the gauge freedom and an additional constraint
should be imposed to leave only the physical degrees of freedom.
This is done by projecting the source tensor  $\{\{formula:7fd88bcd-9013-433d-9756-b874472530d9\}\}$  into its transverse-traceless (TT)
components (see for example  $\{\{cite:80dbb6c8b9e12f561a8e585faceac5f4e104d60d\}\}$ ). Doing this and without loss of generality, we will
use the following very well known ansatz for the source term
proposed in  $\{\{cite:bc9a8ca19785627a087ae0c01abe155c22388e16\}\}$ \n"
/* ----- ref_entries entry for  $\{\{formula:7fd88\dots\}\}$  ----- */
{ "latex": "$_S{\nu\mu\nu}$",
  "type": "formula"
/* ----- bib_entries entry for  $\{\{cite:80dbb\dots\}\}$  ----- */
"bib_entry_raw": "R. Epstein, The Generation of Gravitational Radiation by Escaping Supernova Neutrinos, Astrophys. J. 223 (1978) 1037.",
  "contained_links": [
    { "url": "https://doi.org/10.1086/156337",
      "text": "Astrophys. J. 223 (1978) 1037.",
      "start": 87,
      "end": 117 }
  ],
  "ids" { ... }

```

Listing 1: Data example.

In Figure 2 we show the number of papers across all disciplines over all years covered. We can see that yearly arXiv.org submissions in computer science are likely to surpass those in physics in 2023. As a simple showcase of the use of structured full-text content, we show in Figure 3 how the average number of bibliographic references per paragraph developed over time for the three major disciplines represented in the data set. Dividing by paragraphs is done to account for variation in paper length. We can see that the density of references is increasing more rapidly in physics and computer science, than it is in mathematics.

4.2 Applications

As is evident by the past use of our data set’s predecessors unarXive 2020 and S2ORC, large-scale scholarly data sets created with NLP research in mind have broad applicability. Example uses are analyses of citation behavior across languages [21] or disciplines [24] and the development of models for claim verification [26], document retrieval [16], summarization [14], or information extraction [25].

Due to its similarities in structure and contained information, unarXive 2022 is equally suitable for the applications named above. Beyond these, we provide data for two NLP tasks on unarXive 2022, namely content based citation recommendation and IMRaD classification, which are described in the following.

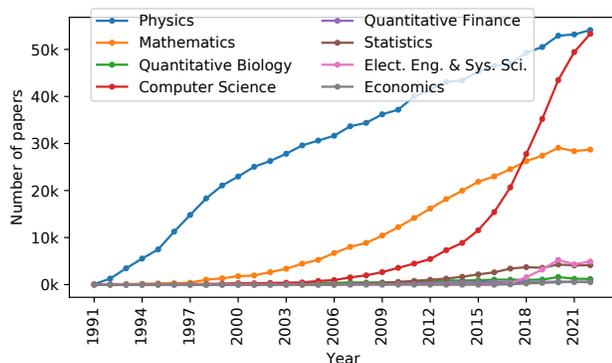


Figure 2: Number of papers per year.

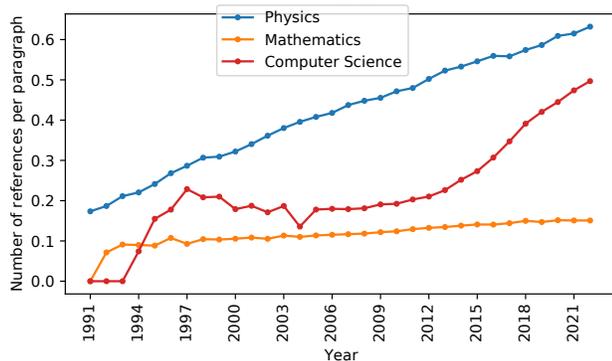


Figure 3: Reference density per year.

Content Based Citation Recommendation. Given a piece of text and a citation-marker position, the task of content based citation recommendation entails identifying publications which are suitable to cite in the given text at the given position [2, 5]. Large full-text corpora of publications with a citation network provide a rich source for supervision of machine learning (ML) models for this task. That is, human made citations are used as training examples, or for evaluating models in a citation re-prediction setting. From the premissively licensed papers in our data set we use all in-text citation markers with a linked reference cited at least three times, to allow splitting into train, dev, and test data. The result is 2.5 M items consisting of (1) a paragraph and citation marker position (model input), and (2) the ID of the cited document (desired model output).

IMRaD Classification. Scientific publications are usually structured into sections commonly summarized as “Introduction, Methods, Results, and Discussion” (IMRaD). Classifying sections of scientific text into these four classes is done, for example, in fine-grained citation classification. Because conventions differ between disciplines, we prepare data for this task for computer science papers only. To aforementioned four classes we add the common “Related Work” section as a fifth class. From the premissively licensed computer science papers in our data, we use those that are unambiguously assignable to one of the five classes. The result is 530 k items consisting of (1) the paragraph text (model input), and (2) the class (desired model output). An exemplary application scenario

for a model trained on this data is a paper writing assistant that can detect parts in a manuscript, which might be better placed in a different section (e.g. discussion rather than results).

4.3 Distribution

Under consideration of the FAIR principles, we chose the following well established distribution channels and licenses for our data set, aforementioned NLP task data, as well as our source code.

- The **data set** is distributed on Zenodo.
 - <https://doi.org/10.5281/zenodo.7752615> (open subset)
 - <https://doi.org/10.5281/zenodo.7752754> (full)
 In accordance with the licensing terms of our source data, we share our data set in two versions.
 - (1) The subset generated from permissively licensed source data (165 k publications, 9%) is openly accessible.
 - (2) The full data set, generated partially from source data under arXiv.org’s “non-exclusive license to distribute,”⁵ is accessible through Zenodo’s “restricted access” policy, making it possible to grant access to the data on request given the intended use is in accordance with the license terms.
- The **NLP task data** is provided on the Hugging Face Hub.
 - https://huggingface.co/datasets/saier/unarXive_citrec
 - https://huggingface.co/datasets/saier/unarXive_imrad_clf
 This facilitates easy access and use by the NLP community.
- The **source code** for creating the data set is shared on GitHub under the MIT License.
 - <https://github.com/IMDepence/unarXive>
 Sharing the code openly and permissively licensed allows anyone to freely modify and extend the code to their needs. This makes, for example, integration into other NLP projects such as benchmarks and frameworks possible.

5 CONCLUSION

We propose unarXive 2022, a data set generated from 1.9 M \LaTeX paper sources and suitable for a wide variety of analyses and NLP applications. We base our approach to data set creation and format on existing works, while also addressing their shortcomings. Improving upon these tried and tested predecessors, unarXive 2022 offers the most complete citation network and most structured content compared to existing data sets, and is surpassed in size only by the PMC-OAS, which covers a different set of disciplines.

With our data set we provide data for two NLP tasks, content based citation recommendation and IMRaD classification, to facilitate its usage. We furthermore distribute our work under consideration of the FAIR principles, sharing it through well established channels and permissively licensed, thereby ensuring proper accessibility, easy use, and possibilities for adaption and extension.

We plan to incrementally update our data set with new arXiv.org submissions. For future developments, we note the importance of mathematical notation in academic publications, as reflected by recent SemEval tasks in 2021 and 2022 [8, 10]. Similar to existing projects,⁶ we plan to investigate novel analyses and applications based on the combination of our data set’s citation network and structured mathematical notation.

⁵See <http://arxiv.org/licenses/nonexclusive-distrib/1.0/>.

⁶See <https://github.com/PierreSenellart/theoremkb>.

AUTHOR CONTRIBUTIONS

Tarek Saier: Conceptualization, Data curation (lead), Formal analysis, Methodology, Software (lead), Visualization, Writing – original draft, Writing – review & editing. Johan Krause: Data curation (support), Software (support). Michael Färber: Writing – review & editing.

ACKNOWLEDGMENTS

This work was partially supported by the German Federal Ministry of Education and Research (BMBF) via [KOM,BI], a Software Campus project (01IS17042). The authors acknowledge support by the state of Baden-Württemberg through bwHPC. We thank Johannes Reber for supporting early stages of the software development.

REFERENCES

- [1] H. Bast and C. Korzen. 2017. A Benchmark and Evaluation for Text Extraction from PDF. In *2017 ACM/IEEE Joint Conference on Digital Libraries (JCDL)*. 1–10. <https://doi.org/10.1109/JCDL.2017.7991564>
- [2] Chandra Bhagavatula, Sergey Feldman, Russell Power, and Waleed Ammar. 2018. Content-Based Citation Recommendation. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1*. Association for Computational Linguistics, 238–251. <https://doi.org/10.18653/v1/N18-1022>
- [3] Chaomei Chen. 2006. CiteSpace II: Detecting and visualizing emerging trends and transient patterns in scientific literature. *Journal of the American Society for Information Science and Technology* 57, 3 (2006), 359–377. <https://doi.org/10.1002/asi.20317>
- [4] Hong Chen, Hiroya Takamura, and Hideki Nakayama. 2021. SciXGen: A Scientific Paper Dataset for Context-Aware Text Generation. In *Findings of the Association for Computational Linguistics: EMNLP 2021*. 1483–1492. <https://doi.org/10.18653/v1/2021.findings-emnlp.128>
- [5] Michael Färber and Adam Jatowt. 2020. Citation recommendation: approaches and datasets. *International Journal on Digital Libraries* 21, 4 (Dec. 2020), 375–405. <https://doi.org/10.1007/s00799-020-00288-2>
- [6] Deyan Ginev. 2020. arXivLiv:2020 dataset, an HTML5 conversion of arXiv.org, hosted at <https://sigmathling.kwarc.info/resources/arxvmliv-dataset-2020/>. SIG-MathLing – Special Interest Group on Math Linguistics.
- [7] Bela Gipp, Norman Meuschke, and Mario Lipinski. 2015. CITREC : An Evaluation Framework for Citation-Based Similarity Measures based on TREC Genomics and PubMed Central. In *iConference 2015 Proceedings*. <http://hdl.handle.net/2142/73680>
- [8] Corey Harper, Jessica Cox, Curt Kohler, Antony Scerri, Ron Daniel Jr., and Paul Groth. 2021. SemEval-2021 Task 8: MeasEval – Extracting Counts and Measurements and their Related Contexts. In *Proceedings of the 15th International Workshop on Semantic Evaluation (SemEval-2021)*. 306–316. <https://doi.org/10.18653/v1/2021.semeval-1.38>
- [9] Jorge E Hirsch. 2005. An index to quantify an individual’s scientific research output. *Proceedings of the National academy of Sciences* 102, 46 (2005), 16569–16572.
- [10] Viet Lai, Amir Pouran Ben Veyseh, Franck Deroncourt, and Thien Nguyen. 2022. SemEval 2022 Task 12: Symlink - Linking Mathematical Symbols to their Descriptions. In *Proceedings of the 16th International Workshop on Semantic Evaluation (SemEval-2022)*. 1671–1678. <https://doi.org/10.18653/v1/2022.semeval-1.230>
- [11] Kyle Lo, Lucy Lu Wang, Mark Neumann, Rodney Kinney, and Daniel Weld. 2020. S2ORC: The Semantic Scholar Open Research Corpus. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, 4969–4983.
- [12] Patrice Lopez. 2009. GROBID: Combining Automatic Bibliographic Data Recognition and Term Extraction for Scholarship Publications. In *Research and Advanced Technology for Digital Libraries*. 473–474.
- [13] Yi Luan, Luheng He, Mari Ostendorf, and Hannaneh Hajishirzi. 2018. Multi-Task Identification of Entities, Relations, and Coreference for Scientific Knowledge Graph Construction. In *Proc. Conf. Empirical Methods Natural Language Process. (EMNLP)*.
- [14] Yuning Mao, Ming Zhong, and Jiawei Han. 2022. CiteSum: Citation Text-guided Scientific Extreme Summarization and Domain Adaptation with Limited Supervision. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*. 10922–10935. <https://aclanthology.org/2022.emnlp-main.750>
- [15] Bethesda (MD): National Library of Medicine. [n. d.]. PMC Open Access Subset. <https://www.ncbi.nlm.nih.gov/pmc/tools/openfst/> [Internet]. B. 2003 - [cited 2023 Feb 7].
- [16] Mathias Parisot and Jakub Zavrel. 2022. Multi-objective Representation Learning for Scientific Document Retrieval. In *Proceedings of the Third Workshop on Scholarly Document Processing*. Association for Computational Linguistics, Gyeongju, Republic of Korea, 80–88. <https://aclanthology.org/2022.sdp-1.9>
- [17] Nancy Pontika, Petr Knoth, Matteo Cancellieri, and Samuel Pearce. 2016. Developing Infrastructure to Support Closer Collaboration of Aggregators with Open Repositories. *LIBER Quarterly* 25, 4 (April 2016), 172–188. <http://oro.open.ac.uk/45935/>
- [18] Animesh Prasad, Manpreet Kaur, and Min-Yen Kan. 2018. Neural ParsCit: A Deep Learning Based Reference String Parser. *International Journal on Digital Libraries* 19 (2018), 323–337. <https://doi.org/10.1007/s00799-018-0242-1>
- [19] Jason Priem, Heather Piwowar, and Richard Orr. 2022. OpenAlex: A fully-open index of scholarly works, authors, venues, institutions, and concepts. <https://doi.org/10.48550/ARXIV.2205.01833>
- [20] Tarek Saier and Michael Färber. 2020. unarXiv: a large scholarly data set with publications’ full-text, annotated in-text citations, and links to metadata. *Scientometrics* (March 2020). <https://doi.org/10.1007/s11192-020-03382-z>
- [21] Tarek Saier, Michael Färber, and Tornike Tsereteli. 2021. Cross-Lingual Citations in English Papers: A Large-Scale Analysis of Prevalence, Usage, and Impact. *International Journal on Digital Libraries* (Dec. 2021). <https://doi.org/10.1007/s00799-021-00312-z>
- [22] Thomas Scheidsteger and Robin Haunschild. 2022. Comparison of metadata with relevance for bibliometrics between Microsoft Academic Graph and OpenAlex until 2020. <https://doi.org/10.48550/arXiv.2206.14168>
- [23] Arnab Sinha, Zhihong Shen, Yang Song, Hao Ma, Darrin Eide, Bo-June Paul Hsu, and Kuansan Wang. 2015. An Overview of Microsoft Academic Service (MAS) and Applications. In *Proceedings of the 24th International Conference on World Wide Web (Florence, Italy) (WWW’15)*. 243–246.
- [24] Michele Delli Veneri, Rafael S. de Souza, Alberto Krone-Martins, E. E. O. Ishida, M. L. L. Dantas, Noble Kennamer, and for the COIN collaboration. 2022. How Have Astronomers Cited Other Fields in the Last Decade? *Research Notes of the AAS* 6, 6 (jun 2022), 113. <https://doi.org/10.3847/2515-5172/ac74c7>
- [25] Vijay Viswanathan, Graham Neubig, and Pengfei Liu. 2021. CitationIE: Leveraging the Citation Graph for Scientific Information Extraction. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. 719–731. <https://doi.org/10.18653/v1/2021.acl-long.59>
- [26] David Wadden, Shanchuan Lin, Kyle Lo, Lucy Lu Wang, Madeleine van Zuylen, Arman Cohan, and Hannaneh Hajishirzi. 2020. Fact or Fiction: Verifying Scientific Claims. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 7534–7550. <https://doi.org/10.18653/v1/2020.emnlp-main.609>
- [27] Mark D. Wilkinson, Michel Dumontier, IJsbrand Jan Aalbersberg, Gabrielle Appleton, Myles Axton, Arie Baak, Niklas Blomberg, Jan-Willem Boiten, Luiz Bonino da Silva Santos, Philip E. Bourne, Jildau Bouwman, Anthony J. Brookes, Tim Clark, Mercè Crosas, Ingrid Dillo, Olivier Dumon, Scott Edmunds, Chris T. Evelo, Richard Finkers, Alejandra Gonzalez-Beltran, Alasdair J. G. Gray, Paul Groth, Carole Goble, Jeffrey S. Grethe, Jaap Heringa, Peter A. C. ’t Hoen, Rob Hooft, Tobias Kuhn, Ruben Kok, Joost Kok, Scott J. Lusher, Maryann E. Martone, Albert Mons, Abel L. Packer, Bengt Persson, Philippe Rocca-Serra, Marco Roos, Rene van Schaik, Susanna-Assunta Sansone, Erik Schultes, Thierry Sengstag, Ted Slater, George Strawn, Morris A. Swertz, Mark Thompson, Johan van der Lei, Erik van Mulligen, Jan Velterop, Andra Waagmeester, Peter Wittenburg, Katherine Wolstencroft, Jun Zhao, and Barend Mons. 2016. The FAIR Guiding Principles for scientific data management and stewardship. *Scientific Data* 3, 1 (2016), 160018. <https://doi.org/10.1038/sdata.2016.18>