# ACM-CR: A Manually Annotated Test Collection for Citation Recommendation

Florian Boudin

florian.boudin@univ-nantes.fr

LS2N, Université de Nantes

Nantes, France

## ABSTRACT

Citation recommendation is intended to assist researchers in the process of searching for relevant papers to cite by recommending appropriate citations for a given input text. Existing test collections for this task are noisy and unreliable since they are built automatically from parsed PDF papers. In this paper, we present our ongoing effort at creating a publicly available, manually annotated test collection for citation recommendation. We also conduct a series of experiments to evaluate the effectiveness of content-based baseline models on the test collection, providing results for future work to improve upon. Our test collection and code to replicate experiments are available at https://github.com/boudinfl/acm-cr

## KEYWORDS

citation recommendation, test collection, digital libraries

## 1 INTRODUCTION

Citing has always been an integral part of academic research, whether it is for backing up claims or for referring to previous work of relevance. Yet, citing properly is becoming increasingly difficult and time-consuming as the volume of published research continues to grow exponentially [5]. Researchers are under constant pressure to publish their findings, but have less and less time to browse the literature to retrieve relevant papers to cite. This has motivated an active line of research on citation recommendation systems (see [4] for a recent survey), whose goal is to relieve the researchers from performing this task by recommending appropriate citations for a given text.

More precisely, citation recommendation is the task of finding relevant citations for a given *citation context*, that is, a text passage (e.g. a sentence or a paragraph) within a document (See Figure 1).[1]

---

[1]This is not to be confused with paper recommendation, which is the task of recommending documents to the user that are worth reading [1].

Citation recommendation, also referred to as *context-aware citation recommendation* or *local citation recommendation* in previous work, is often viewed as a retrieval task where, given a citation context (query), the task is to retrieve the most relevant citations (documents) from a collection of scientific texts. The benefits of doing so are two-fold: 1) well-known retrieval models can be readily applied to the task at hand, and 2) the effectiveness of citation recommendation systems can be evaluated offline through *test collections*.
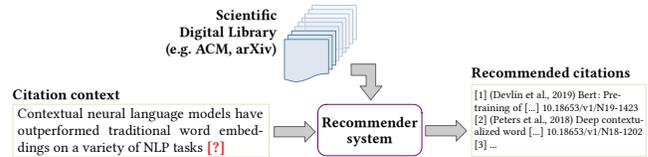


**Figure 1: Illustration of the citation recommendation task.**

Existing test collections for citation recommendation are built automatically by extracting citation contexts (queries) and cited reference(s) (relevant judgments) from a collection of scientific papers. As a result, they are known to be quite noisy and unreliable due to errors in citation parsing and/or PDF to text conversion [4]. Other reported problems include muddled citation contexts (roughly defined as fixed-size windows of characters around citations), and incomplete metadata information of papers. In this work, we focus on addressing these problems and present our ongoing effort to create ACM-CR, a new publicly available, manually annotated test collection for citation recommendation.

## 2 THE ACM-CR TEST COLLECTION

As a first step, we assembled a sizeable collection of documents by collecting bibliographic records of scientific papers (BIBTEX entries) from the ACM Digital Library. Our document collection currently holds 114,882 records of scientific papers on topics related to Information Retrieval (IR) and adjacent research fields (e.g. machine learning, databases or data mining).[2] Each record carries a Digital Object Identifier (DOI), and most of them (91%) provide paper abstracts which are the primary units of indexing in scientific literature search engines [6].

To create queries and relevance judgments, we collected 50 open-access scientific papers published at top-tier IR conferences in 2020[3], from which we manually extracted citation contexts and cited references. For citation contexts, we extracted full paragraphs that include citations from the introduction and related work sections of each paper (other sections being less suitable for the task [7]).

---

[2]We use the SIGs IR, KDD, CHI, WEB and MOD sponsored conferences and related journals as a means to filter articles.

[3]Venues are SIGIR, CHIIR, ICTIR and WSDM.

We then took a step further and removed end-of-line hyphens, split paragraphs into sentences and anonymized citations that reveal their author's identity (i.e. "Smith et al. [1] proposed" is reformulated as "[1] proposed"). For cited references, we mapped each reference to its DOI by querying Google Scholar[4]. We paid a particular attention to preprint references in case they have been superseded by refereed publications. For those without DOIs but available online, we indicated the URLs to get their PDF files for later use. The papers we collected have 31.8 cited references on average, half of which occur in our document collection. Annotating a single paper took about an hour on average, all annotations being encoded in XML format.

Table 1 presents some statistics about the test collection. Overall, we extracted 341 citation contexts (paragraphs) from which 269 can be used as queries for evaluation because they have cited references that appear in our document collection (underlined in Table 1). It should be noted that our citation contexts are on average twice as long as the 400 characters used in previous work, and thus we might expect better retrieval results. On a lower level of granularity, we annotated 837 out of the 1,800 sentences that make up the 341 citation contexts, each of which has 2 cited references on average. Again, a smaller portion of these annotated sentences (552) have cited references that appear in our document collection and can be used as queries. Here, the idea is to provide two different types of queries (paragraphs and sentences) to investigate the effect of citation context length on retrieval performance.

| Context | #nb | #queries | #tokens | #char. | #cit. |
|---------|-----|----------|---------|--------|-------|
| Paragraphs | 341 | 269 | 145.4 | 809 | 5.1 |
| ∟ Sentences | 837 | 552 | 30.9 | 164 | 2.1 |

**Table 1: Statistics of the test collection. The average number of tokens, characters and citations per context are reported.**

## 3 EXPERIMENTS

We report a series of experiments to evaluate the effectiveness of content-based citation recommendation models on the introduced test collection. It serves two main purposes: 1) to perform an initial sanity check on the test collection, and 2) to provide baseline results for future work to improve upon. The first model we consider is BM25, a standard *ad-hoc* retrieval model that is commonly used as baseline in citation recommendation. Specifically, we use the implementation of BM25 from the Anserini open-source IR toolkit [9] with the default parameters. For the second model, we adopt a two-stage neural document ranking approach inspired by [8], and use SciBERT [2] to re-rank the top-20 documents retrieved by BM25. We apply the uncased model[5] without fine tuning and re-rank documents against citation contexts by computing the cosine similarity between their hidden representations. We evaluate effectiveness of the models in terms of recall and nDCG at the top 10 recommendations as recommended in [4]. We use the Student's paired t-test to assess statistical significance of our retrieval results at $p < 0.05$.

Results are presented in Table 2. We observe noticeably higher scores for the BM25 model in contrast to prior work (e.g. $\approx$ +10%

---

[4]https://scholar.google.com/
[5]https://github.com/allenai/scibert

| Model | Paragraphs | | Sentences | |
|-------|------------|--|-----------|--|
| | R@10 | nDCG@10 | R@10 | nDCG@10 |
| BM25 | 33.58 | 28.27 | 39.03 | 31.56 |
| BM25+SciBERT | 34.40 | 29.37[†] | 40.30[†] | 32.31[†] |

**Table 2: Retrieval effectiveness of content-based citation recommendation models. † indicates significance over BM25.**

in comparison to [3]), which we attribute to the high quality of the manually extracted citation contexts. Re-ranking using SciBERT increases the retrieval effectiveness, with statistically significant improvements in three out of four cases. The impact of re-ranking is more important on sentences, which seems reasonable since short queries are prone to vocabulary mismatch issues.

## 4 CONCLUSION AND FUTURE WORK

In this paper, we described our progress in creating ACM-CR, the first manually annotated test collection for citation recommendation. We hope that this resource will provide a useful benchmark for future studies, and help us gain better insights on the effectiveness of citation recommendation models. For future work, we plan to collect and annotate more papers, while also exploring two directions for improving our test collection. The first one will be to obtain the full-text of all the open-access articles in our document collection, and to construct the underlying citation graph on which many citation recommendation models rely on. The second direction concerns the incomplete nature of the extracted relevance judgments, i.e. that do not include uncited, yet relevant papers. Here, we will investigate how co-citation instances can be used to overcome the sparseness of the relevance judgments.

## ACKNOWLEDGMENTS

## REFERENCES

[1] Joeran Beel, Bela Gipp, Stefan Langer, and Corinna Breitinger. 2016. Paper recommender systems: a literature survey. *International Journal on Digital Libraries* 17, 4 (2016), 305–338. https://doi.org/10.1007/s00799-015-0156-0

[2] Iz Beltagy, Kyle Lo, and Arman Cohan. 2019. SciBERT: A Pretrained Language Model for Scientific Text. In *Proceedings of EMNLP*. 3615–3620. https://doi.org/10.18653/v1/D19-1371

[3] Michael Färber and Ashwath Sampath. 2020. HybridCite: A Hybrid Model for Context-Aware Citation Recommendation. In *Proceedings of JCDL*. 117–126. https://doi.org/10.1145/3383583.3398534

[4] Michael Färber and Adam Jatowt. 2020. Citation recommendation: approaches and datasets. *International Journal on Digital Libraries* 21 (2020), 375–405. https://doi.org/10.1007/s00799-020-00288-2

[5] Nature Publishing Group (Ed.). 2009. Credit where credit is due. *Nature Cell Biology* 11, 1 (2009), 1–1. https://doi.org/10.1038/ncb0109-1

[6] Chien-yu Huang, Arlene Casey, Dorota Głowacka, and Alan Medlar. 2019. Holes in the Outline: Subject-Dependent Abstract Quality and Its Implications for Scientific Literature Search. In *Proceedings of CHIIR*. 289–293.

[7] David Jurgens, Srijan Kumar, Raine Hoover, Dan McFarland, and Dan Jurafsky. 2018. Measuring the Evolution of a Scientific Field through Citation Frames. *Transactions of the Association for Computational Linguistics* 6 (2018), 391–406. https://doi.org/10.1162/tacl_a_00028

[8] Rodrigo Nogueira, Zhiying Jiang, Kyunghyun Cho, and Jimmy Lin. 2020. Evaluating pretrained transformer models for citation recommendation. In *BIR 2020 Workshop on Bibliometric-enhanced Information Retrieval*. 89–100.

[9] Peilin Yang, Hui Fang, and Jimmy Lin. 2017. Anserini: Enabling the Use of Lucene for Information Retrieval Research. In *Proceedings of SIGIR*. 1253–1256. https://doi.org/10.1145/3077136.3080721