

Exposé

outlining a Master Thesis on:

Semantic approaches to scientific citation recommendation (tentative title)

Tarek Saier

Reviewer: Prof. Dr. Georg Lausen

Advisor: Dr.-Ing. Michael Färber

1. INTRODUCTION

This exposé will outline a prospective Master Thesis in the area of scientific citation recommendation and argue for its value. The approach will encompass the creation of a dataset and development of supervised learning methods with a focus on semantic analysis of citation contexts. Evaluation of the resulting implementation will follow the most prevalent methods in the field.

The remainder of this document is structured as follows. Section 2 will provide some theoretical background on relevant areas and give a quick overview of related work. A detailed description of the planned methodology and approach will be given in section 3; followed in section 4 by an outline of the planned evaluation. Section 5 and 6 conclude the exposé by listing the expected contributions of the Thesis and a proposed schedule.

2. BACKGROUND

2.1 Citation recommendation

The goal of citation recommendation is to provide adequate citations to a given input text. This can involve evaluating whether or not a given input text includes parts that are suitable to add citations to in the first place. For a given section of or position in an input text, the output recommendation can either be a single citation or a ranked list of multiple possible citations. A further distinction can be made concerning the granularity of text that a citation is recommended for. This can range from a complete document (global citation recommendation) to a specific point within a string of text (context aware/local citation recommendation). There are also approaches where citation markers—annotations in the text that mark the position of a citation—are left in the input text. In such a case the evaluation whether or not a citation should be recommended as well as the decision where exactly to put a citation are not necessary. In an ideal case, citation recommendation can even involve evaluating candidate documents in terms of their quality.

Given there are a lot of dimensions along which approaches can differ, section 2.3 will explain relevant terminology and section 2.4 will give an overview of these distinguishing dimensions. This will enable a more easily understandable overview of related work.

2.2 Semantic analysis

The idea of this thesis is to focus on semantic aspects of citation contexts. This means, rather than taking into account only syntactical aspects like n-grams, the analysis will go to a higher level of abstraction where the input's *meaning* is of importance. Because the focus of this analysis will most likely revolve around entities, claims and arguments, these terms will be defined in the following section.

2.3 Terminology

Citing/cited document. The former is the document making a reference while the latter the document being referenced. The contents of both can be taken into account when developing a citation recommendation approach, but in a considerable amount of approaches the *cited* documents' content is not.

Citation context. Within the citing document and concerning a single recommendation being made, this is the extend of text provided as input. Examples would be the citing documents abstract, a sentence containing a citation marker or a whole document.

Citation marker. A citation marker is an annotation in the input text (or a data set) that marks the location of a citation. In scientific publications this could, for example, be a [42]. When left in the input text for a recommendation process, the marker's association to its corresponding reference entry is, of course, removed (e.g. [42] could be changed to //, replaced by another type of annotation or the citing document's reference section could be made unavailable during the processing of the input).

Reference. For each citation marker there usually is a corresponding reference at the bottom of the page or near the end of the document. This reference identifies the cited document.

Citation function. The role of a citation or, put differently, the motivation that was behind putting a citation in a particular place. This can, for example, be just for referencing a data set that was used (by citing a data paper), backing up a claim or arguing for or against the overall proposition of a publication.

Metadata. In addition to a document's content, information *about* the document is also often taken into considera-

tion during the recommendation process. This is referred to as metadata.

Entity. A physical or abstract thing in the real world. Generally speaking entities like for example people, places, events and topics can be of interest.

Claim. In this setting a claim can be defined as an assertion which can be judged in terms of its factuality. While non-factual claims also exist (i.e. an opinion being stated), they do not need backing up by citations and are therefore not of interest for citation recommendation.

Argument. An argument can, in alignment with [1], be defined as being composed of a claim and one or more premises justifying the claim. To illustrate, this can take the form $\langle \text{premises} \rangle \langle \text{step(s) of deduction} \rangle \langle \text{claim} \rangle$ where the claim is the conclusion of the deduction.

2.4 Dimensions

To systematically categorize approaches to citation recommendation, distinctions can be made concerning the input and the output of a mechanism. With regards to the input, the dimensions *citation context* (length/position), *citation markers* (available or not) and *metadata* (available to what extent) can be used. In part, these can be further broken down as shown in the following table.

	learning	use
citing doc	$\langle \text{val} \rangle$	$\langle \text{val} \rangle$
cited doc	$\langle \text{val} \rangle$	$\langle \text{val} \rangle$

That is, citing and cited documents can be looked at separately, and a distinction can be made as to what is available during the learning phase and what needs to be provided as input during actual use of the resulting system. Note that for citation markers and context, only the *citing doc* row is applicable¹ and for metadata most likely only the *learning* column is². Because dimensions along tree axis are hard to visualize effectively, the distinction can be flattened to the following aspects:

- citation context (learning)
- citation context (use)
- citation markers (learning)
- citation markers (use)
- metadata (citing doc)
- metadata (cited doc)

¹The term "citation context" is used to refer to the context in the citing document. One could make a point, though, to furthermore introduce the notion of a context in the cited document. This could then be used, for example, to distinguish whether or not a mechanism outputs only a recommended document or also a specific section that is relevant; or to distinguish whether or not (parts of) cited documents are used during the learning phase.

²Although metadata aspects like the "date of the citing doc" could also be used in the online system. That is, given a newly written text without citations, an approach could interpret the input as a "recent citing doc" and recommend citations accordingly.

To give a concrete example, an approach could be trained on input with citation markers (citation markers learning), but be able to give useful output for input without markers (citation markers use) as well.

Above example also suggests, that there is a distinction to be made concerning an approach's output. A dimension *citation placement granularity* can be used to distinguish whether citation recommendations are given for a whole document, on a sentence level or if specific points within the text are identified.

2.5 Related work

In [2] Färber et al. give a comprehensive overview of the field of citation recommendation as well as a comparison of concrete approaches. Focus in the following will be works with distinct similarities or differences to the proposed approach (explained in section 3) which are therefore helpful in defining it.

Mishra et al. describe in [8] an approach to recommend news articles that can be used as references for Wikipedia articles describing historical events. Their goal is to offer readers an insight into the detailed view on and reporting of an event *at the time* as an addition to the more overarching representation on Wikipedia. This approach employs named entities as a key component to identify appropriate news articles to recommend. It is therefore similar in this regard to the first step in the Master Thesis' approach where the focus also will lie on recommendation based on entities. The domain, being Wikipedia and news articles, differs from scientific publications.

In [7], Levy et al. describe a method for claim detection using a cascade of classifiers. The detection of claims will also be necessary in the proposed Thesis' second step (citation recommendation based on claims). Levy et al. do, however, restrict their detection of claims to those related to a pre-defined topic and include claims that are statements of an opinion, which will most likely not be the case in the Master Thesis.

In a similar fashion Goudas et al. tackle argument extraction in [4], which will need to be done in the Thesis' third step (citation recommendation based on arguments). The document type being social media texts is, however, different.

3. METHODOLOGY AND APPROACH

```
foo bar
MAG[9][5][6][10]
entity[8]
claim[7]
argument[4]
```

data sets that were considered and why (benefits, drawbacks, ... (cite accordingly))

MAG start and arXiv start scenario (see wiki)

details of arXiv processing, challenges, etc. (MAG for evaluation where citation marker position is not relevant)

4. EVALUATION

```
foo bar
```

5. CONTRIBUTIONS

- apparently semantic stuff not very explored (cite survey if possible, look at tables) - creation of another nice (exact citation markers, large citation context, etc.) dataset like gold standard paper[3] - a nice dataset like gold standard paper[3] but not restricted to CS domain

6. SCHEDULE

7. REFERENCES

- [1] P. Besnard and A. Hunter. *Elements of Argumentation*. The MIT Press, 2008.
- [2] M. Färber and A. Jatowt. Citation Recommendation for Scientific Publications.
- [3] M. Färber, A. Thiemann, and A. Jatowt. A High-Quality Gold Standard for Citation-based Tasks. In *Proceedings of the 11th International Conference on Language Resources and Evaluation*, LREC 2018, 2018. r.
- [4] T. Goudas, C. Louizos, G. Petasis, and V. Karkaletsis. Argument extraction from news, blogs, and social media. In A. Likas, K. Blekas, and D. Kalles, editors, *Artificial Intelligence: Methods and Applications*, pages 287–299, Cham, 2014. Springer International Publishing. r.
- [5] D. Herrmannova and P. Knoth. An analysis of the microsoft academic graph. *D-Lib Magazine*, 22(9/10), 2016. r.
- [6] S. E. Hug, M. Ochsner, and M. P. Brändle. Citation analysis with microsoft academic. *Scientometrics*, 111(1):371–378, Apr 2017. r.
- [7] R. Levy, Y. Bilu, D. Hershcovich, E. Aharoni, and N. Slonim. Context dependent claim detection. In *Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers*, pages 1489–1500, Dublin, Ireland, August 2014. Dublin City University and Association for Computational Linguistics. r.
- [8] A. Mishra and K. Berberich. Leveraging semantic annotations to link wikipedia and news archives. In N. Ferro, F. Crestani, M.-F. Moens, J. Mothe, F. Silvestri, G. M. Di Nunzio, C. Hauff, and G. Silvello, editors, *Advances in Information Retrieval*, pages 30–42, Cham, 2016. Springer International Publishing. r (ch 1-3).
- [9] B. Paszcza. Comparison of microsoft academic graph with other scholarly citation databases, 11 2016. r (ch 1, ""3"").
- [10] A. Sinha, Z. Shen, Y. Song, H. Ma, D. Eide, B.-J. P. Hsu, and K. Wang. An overview of microsoft academic service (mas) and applications. In *Proceedings of the 24th International Conference on World Wide Web, WWW '15 Companion*, pages 243–246, New York, NY, USA, 2015. ACM. r.

Time frame	Task	Results
Oct 1 – Oct 21	Develop mechanism to generate dataset with citation markers from arXiv source dump	Dataset boilerplate (i.e. with citation markers but no semantic annotation)
Oct 22 – Oct 28	Write exposé	Thesis approval
Oct 29 – Nov 04	Add entity annotations to dataset	Dataset usable for supervised learning
Nov 05 – Nov 18	Develop entity based recommendation approach	-
Nov 19 – Nov 25	Add claim annotations to dataset	-
Nov 26 – Dec 09	Develop claim based recommendation approach	-
Dec 10 – Dec 22	Coordination with simultaneous tangential theses and integration into CiteRec system	-
Dec 23 – Jan 06	break/buffer	-
Jan 07 – Jan 13	Add argument annotations to dataset	-
Jan 14 – Jan 27	Develop argument based recommendation approach and start offline evaluation	-
Jan 28 – Feb 10	Offline evaluation	-
Feb 11 – Feb 24	Online evaluation	-
Feb 25 – Mar 17	Thesis writing	-
Mar 18 – Mar 31	buffer/paper writing	-