

Chapter 50

Tracing Related Scientific Papers by a Given Seed Paper Using Parscit

Resmana Lim, Indra Ruslan, Hansin Susatya, Adi Wibowo,
Andreas Handojo and Raymond Sutjiadi

Abstract The project developed a web site learning support for tracing scientific articles relating to a given input seed paper/article. The system will finds related articles that are listed on the references of the seed article. First, the reference list of an input seed article is extracted by utilizing Parscit citation extraction. Furthermore, the system searches the reference articles using Google Scholar and Mendeley API. Thus articles which are related to the seed article can be found. The system was built using PHP programming, it is utilizing Parscit modules, Google Scholar search and <https://www.mendeley.com/> API. Testing has been done by giving an input seed article. User will obtain the results of several articles related to the seed article.

Keywords Paper tracing · Citation extraction · Parscit · Google scholar · Mendeley.com

50.1 Introduction

To support literature review on a research, further reference search from an article/paper that we read is an important thing. When we read a scientific article, very often we are keen to discover more about the articles contained in the list of references. We can manually perform a search using Google scholar. But this will

R. Lim (✉) · H. Susatya
Electrical Engineering Department, Petra Christian University,
Surabaya, Indonesia
e-mail: resmana@petra.ac.id

I. Ruslan · A. Wibowo · A. Handojo
Informatics Engineering Department, Petra Christian University,
Surabaya, Indonesia

R. Sutjiadi
Computer Engineering Department, Institut Informatika Indonesia,
Surabaya, Indonesia

take substantial efforts when there are many references that we will find. Therefore we need a tool to find or download the references we find automatically given an input seed article.

An automated extraction of paper's bibliography metadata is a challenging task given the variety of different paper's layouts and formatting (citation style) of its reference strings. Several automated extraction approaches have been proposed which are using unsupervised methods and template matching. It finds that the most promising approach is utilizing supervised sequence classification such as Hidden Markov Models (HMMs) [1] or Conditional Random Fields (CRFs) [2]. ParsCit [3] is a popular reference extraction system that uses heuristics approach to detect and segment references within a scientific paper. ParsCit uses CRF to assign labels to the tokens within each reference string. ParsCit is an open-source CRF-based citation parser that has been successfully used by CiteSeerx [4] and scientific papers harvester system [5].

Our project contributes to the purpose of enriching a scientific article (we call the seed paper) by identifying its bibliography (list of references) and connecting them to the original resources (full text if available) by utilizing google scholar and Mendeley API. The project uses Parscit, to process the bibliography in scientific articles. Parscit generate reference's meta-data such as author name, article title, year, etc. The meta-data is then used to perform queries to Google Scholar and Mendeley. The system seeks to connect the bibliography with its source full text using Google Scholar or Mendeley.

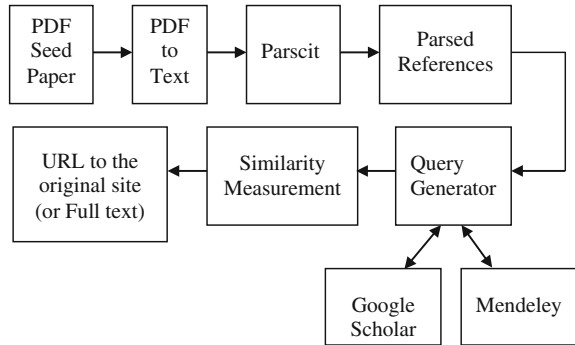
50.2 System Description

This web application is built by using PHP and HTML programming language and MySQL as database management system. The system uses Parscit as service provider to extract paper information data. All this components are run in a Linux based server. To accommodate function of parsing data to other data sources, PHP must be set to enable the cURL function.

In general, this system can be decomposed into parts of process as shown on Fig. 50.1. The first part is preprocessing. In this part, paper is prepared to be processed by Parscit by formatting PDF paper into TXT format. The second part is paper extraction process. This process is handled by Parscit to extract the paper meta-data such as: title, author(s), abstract, keywords. Parscit also do extract the list of references, which are used to make Google Scholar/Mendeley queries on the next phase. The third part is related paper searching process, where in this part, server makes query to other data sources (i.e. Google Scholar and <https://www.mendeley.com/>) to search related paper from references data. The last part is composing and presenting the information as the web output to the user.

As an input seed paper, user uploads a PDF paper into the system. Because Parscit can only process text based file, then before moving to the next part, unicode text must be extracted from PDF file. To do that, Apache PDFBox or similar

Fig. 50.1 System block diagram



application is used to extract unicode text from the input PDF file. The result from this process is a TXT file format.

This TXT file format is processed further by Parscit to grab paper meta-data (i.e. title, author, affiliation, abstract, etc.) and parsed its references/bibliography data. Before the meta-data extraction, Parscit must be trained first with some sample data using similar paper format and language to detect the sections of paper. Make sure that there are enough sample paper data to increase the accuracy of data extraction. The output from the Parscit extraction is an XML file format that already accommodate extracted information which is grouped by different XML tags as shown in Fig. 50.2.

In the next step, each parsed reference data is forwarded to query generator. In this part, query generator searches related paper data from 2 data sources, Google Scholar and Mendeley. The system initiates HTTP query using Document Object

```

<citation valid="true">
<authors>
<author>Adi Kusrianto</author>
</authors>
<title>Pengantar Desain Komunikasi Visual.</title>
<date>2007</date>
<publisher>Penerbit ANDI.</publisher>
<location>Yogyakarta:</location>
<marker>Kusrianto, 2007</marker>
<rawString>Kusrianto, Adi. (2007). Pengantar Desain Komunikasi Visual. Yogyakarta: Penerbit ANDI.</rawString>
</citation>
  
```

Fig. 50.2 Part of XML citation data

Model (DOM) parsing method to Google Scholar to search related paper using title and author data that already classified by Parscit XML result above. This parsing method is used because Google Scholar does not provide Application Programming Interface (API). If similar data are found in Google Scholar, than system saves the result into database to be processed further in the next step.

For Mendeley, system makes the same query to search related paper information, but in easier way because Mendeley provides PHP API to query the searching result. If similar data are found in Mendeley, than system saves the result into database to be processed further in the next step.

After getting the result, in the next part, system runs similarity checking to compare Google Scholar and Mendeley results with parsed references data from Parscit. Similarity checking is done by comparing the author name and title data between each reference data and then calculating the similarity value using formula 50.1 and 50.2 respectively as used in [6, 7].

$$AS = \sum_n \frac{LS}{LN} \quad (50.1)$$

where:

AS final similarity score of author name between paper and reference data.

LS number of matched characters between author name in paper and reference data.

LN number of total author name characters.

$$Sr = 1 - \frac{\|r1 - r2\|}{\|r1 + r2\|} \quad (50.2)$$

where:

Sr efficiency metric as similarity indicator of title between paper and reference data.

r1 vector of word sequence in paper data.

r2 vector of word sequence in reference data.

After similarity value has been known, system arranges the descending order of related paper data for each reference. This similarity checking process is to make sure that top listed is the most appropriate related paper.



Fig. 50.3 Result of related paper query

As the output as shown in Fig. 50.3, system shows the result in form of web page format. Each reference data of a paper is added with an URL link to download the full text paper from the result of data query above.

50.3 Result and Discussion

For testing, we have already trained the Parscit to recognize the English paper optimally using the same writing format with paper that used as a sample for this testing. As the input, we use seed English paper in PDF format and gives the output as shown in Fig. 50.4.

From extraction result above, system can parse the reference data optimally. This reference data than processed further to Mendeley and Google Scholar to search related papers. As the result, system can list the related paper successfully as shown in Fig. 50.5.

Application of a CC-VSI for Active Filtering and Photovoltaic Energy Conversion with a 1-to-1 MPPT controller

AUTHORS
Hanny H Tumbelaka, Masafumi Miyatake

ABSTRACT
This paper focuses on the implementation of a three-phase four wire current-controlled Voltage Source Inverter (CC-VSI) as both PV energy extraction and power quality improvement. For power quality improvement, the CC-VSI works as a grid current-controller shunt active power filter. Then, the PV array supported by a Look-up Table type of a MPPT controller is coupled to the DC bus of the CC-VSI. The output of MPPT controller is a DC voltage that determines the DC-bus voltage according to the PV maximum power. The computer simulation results show that the system works properly in steady state and dynamic condition.

EMAILS
tumbeh@petra.ac.id

AFFILIATION
1Department of Electrical Engineering Petra Christian University, Surabaya, Indonesia 2Department of Engineering and Applied Sciences Sophia University, Tokyo, Japan

KEYWORD
active power filter, MPPT, PV energy conversion

RELATED CONFERENCES
[Show Related Conferences](#)

RAW STRING

- # Borle, L., Zero Average Current Error Control Methods for Bidirectional AC-DC Converters, PhD Thesis, Electrical and Computer Engineering, Curtin University of Technology, Western Australia, 1999.
- # Cheri, Y., and Smedley, K.M., "A Cost-Effective Single-State Inverter with Maximum Power Point Tracking", IEEE Transactions on Power Electronics, 2004, 19(5): p. 1289-1294.
- # Castaner, L., and Silvestre, S., Modelling Photovoltaic System using Pspice, John Wiley & Sons, 2002.
- # Wanzeller, M.G. et.al., "Current Control Loop for Tracking of Maximum Power Point Supplied for Photovoltaic Array", IEEE Transactions on Instrumentation and Measurement, 2004, 53(4): p. 1304-1310.
- # El-Habrouk, M., M.K. Darwish, and P. Mehta, "Active power filters: a review", Electric Power Applications, IEEE Proceedings, 2000. 147(5): p. 403-413.
- # Tumbelaka, H.H., L.J. Borle, and C.V. Nayar. "Analysis of a Series Inductance Implementation on a Three-phase Shunt Active Power Filter for Various Types of Non-linear Loads", Australian Journal of Electrical and Electronics Engineering, Engineers Australia, 2005. 2(3): p. 223-232.
- # Tumbelaka, H.H., L.J. Borle, C.V. Nayar, and S.R.Lee, "A Grid Current-controlling Shunt Active Power Filter", Proceedings of ICPE'07, 2007. Daegu, Korea.
- # Grandi, G., Casadei, D., and Rossi, C., "Direct Coupling of Power Active Filters with Photovoltaic Generation System with Improved MPPT Capability", in IEEE Power Tech Conference, 2003. Bologna, Italy. Application of a CC-VSI for Active Filtering and Photovoltaic Energy Conversion

Fig. 50.4 Result of parsit extraction process



Fig. 50.5 List of reference related paper

50.4 Conclusion

We have developed a web site learning support for tracing scientific articles relating to a given input seed article. The system finds related articles that are listed on the references of the seed article so that user could track and download the paper easily.

Acknowledgments The authors would like to thank the Indonesian Directorate General of Higher Education, which provided the funds for the research project.

References

1. Hetzner, E.: A simple method for citation metadata extraction using hidden markov models. In: Proceedings of the 8th ACM/IEEE Joint Conference on Digital Libraries JCDL 08, vol. 18, no. 3, pp. 280–284 (2008)
2. Peng, F., McCallum, A.: Information extraction from research papers using conditional random fields. *Inf. Process. Manage.* **42**(4), 963–979 (2006)
3. Councill, I.G., Giles, C.L., Kan, M.: ParsCit: an open-source CRF reference string parsing package. In: Proceedings of LREC '08, pp. 661–667 (2008)
4. Teregowda, P.B., Uргаonkar, B., Giles, C.L.: Citeseerx: a cloud perspective. In: Proceedings of Second USENIX Work, Hot Topic in Cloud Computing (2010)
5. Sutjiadi, R., Lim, R., Wibowo, A., Handojo, A.: Mobile application for accessing paper citation with social network feature. *Adv. Sci. Lett.* **21**(7), 2179–2182 (2015)

6. Li, Y., Bandar, Z., McLean, D., O'Shea, J.: A method for measuring sentence similarity and its application to conversational agents. In: proceedings of FLAIRS Conference (2004)
7. Liliana, L., Lim, R., Kwan, E.: Voice conversion application (VOCAL). In: 2011 International Conference on Uncertainty Reasoning and Knowledge Engineering (URKE), vol. 1, pp. 259–262, IEEE August (2011)