# E-learning Document Search Method with Supplemental Keywords Derived from Keywords in Meta-Tag and Descriptions which are Included in the Header of the First Search Result

Kohei Arai

Graduate School of Science and Engineering
Saga University
Saga City, Japan

Herman Tolle

Software Engineering Department
Brawijaya University
Malang, Indonesia

*Abstract*— **Optimization method for e-learning document search with keywords which are derived from the keywords and descriptions in the meta-tag of web search results together with thesaurus engine is proposed. 15 to 20% of improvement on hit rate of search performance is confirmed with the proposed search engine.**

*Keywords- Search engine; e-learning content; thesaurus engine.*

## I. INTRODUCTION

When a word or words are typed in search engines, a list of web sites that contain those words is displayed. The words you enter are known as a query [1]. Baeza-Yates and Ribeiro-Neto linked Information Retrieval to the user information needs which can be expressed as a query submitted to a search engine [2]. Search engines were also known as some of the brightest stars in the Internet investing frenzy that occurred in the late 1990s [3]. Although search engines are programmed to rank websites based on their popularity and relevancy, empirical studies indicate various political, economic, and social biases in the information they provide [4],[5].

Based on our previous experiment [6], our system can detect client mobile browser and provide proper format for document and mobile markup language. We propose a new method and approach for developing a new search engine for helping people search E-Learning document on the Internet. We develop a new system based on the open search engine API (Application Protocol Interface) like Google and Yahoo. We create an e-learning specific search engine with improvement in the efficiency and effectiveness in searching document file format comparing with just using original Google or Yahoo. The new system is also accessible through mobile browser on mobile devices for support recent and future technology in the mobile area [7] Method for e-learning contents search engine of ELDOXEA is proposed already. ELDOXEA allows search e-learning contents with a single keyword. Hit rate of search performance of the ELDOXEA is not good enough. In order to improve hit rate, supplemental keywords are required in addition to the firstly input keyword which is referred to primary keyword hereafter.

In order to choose appropriate supplemental keywords for improvement of hit rate, several attempts are performed. First one is to use keywords and descriptions in meta-tag of the header of the first web search result. Keywords are used to be in the meta-tag. Also, there are some keywords in the descriptions in the header of the first web search result. Therefore, these supplemental keywords in the meta-tag and descriptions in the header are applicable to add to the primary keyword.

Second one is to use keywords which are derived from thesaurus engines. Thesaurus engine provides similar words to the primary keyword with their priority. Therefore, these are used for supplemental keywords.

Third one is to use twitter for gathering suggestions of supplemental keywords from twitters. Reliability of the twitters can be evaluated with the previously proposed method.

Fourth one is to use Bulletin Board System: BBS for gathering suggestions of supplemental keywords from the community members. Firstly, the user has to send a message which is requesting supplemental keywords to BBS system with the primary keyword. Then the user has reply message with supplemental keywords. These are to be candidates of supplemental keywords.

Experiments with some queries, "Linear Algebra" and "Hazardous Materials Handler examination" are conducted for searching e-learning contents. By adding supplemental keywords to the primary keyword based on the aforementioned four methods, we confirm their efficiency, hit rate improvements.

The second section describes the proposed methods for choosing supplemental keywords while the third section describes some experimental results followed by some concluding remarks and some discussions.

## II. PROPOSED METHODS

### A. Search Engine

There are three types of search engines, (1) Directory type,

(2) Information Collection Robot type, and (3) Hybrid type.

(1)  Directory type

Content in the Web directories is searched and examined by operators so that it is reliable. Information content is not so much. There is some response delay after updating web sites.

(2)  Information Collection Robot: ICR type

ICR is collecting from the web sites so that information contents are rich. On the other hands, classification items are not so many as appropriate for a wide variety of search purposes.

Ex) Google[1], Inforseek[2], etc.

(3)  Hybrid type

It has the aforementioned both benefits.

Ex) Yahoo[3], MSN[4], ELDOXEA[5], etc.

The proposed search engine for e-learning document search is based on the ELDOXEA of Hybrid type.

*B.  Efficiency Improvement on Searching Process of E-Learning Document on the Internet*

To improve the efficiency of the searching process of E-Learning document, we design a new process for searching and display the document. One of the most problems on efficiency while we search a document file is how to get the appropriate files in a fast way. People usually have to check on each document files of the results set, start from the first results.

Most of document files format is not able to display directly on browser without additional plug-in or application. So, in the conventional way, we should download the document file first, and then open it in our PC, for example, we open PPT files using *Microsoft PowerPoint*. After check the content then we can decide to keep this file or not. This process takes time if the file size is large and we should wait for download process. Another problem with this process is storage problem for download too many unrelated files.

We design a new process that help user to preview document, before they decide to save (download) it. The preview processes displaying the document file in the same results area, so users still stay in the same page while checking each result. This will  lead the user  to control which document is related or not related and decide to save or skip it. The proposed algorithm as follows:

1.  Get *SearchKey* from user input
2.  Get *RelatedKey* and *NotRelatedKey* from a Database based on user's *SearchKey*
3.  Create the *CompleteKeyword* (1)
4.  Search using the *CompleteKeyword*
5.  Get results and display it
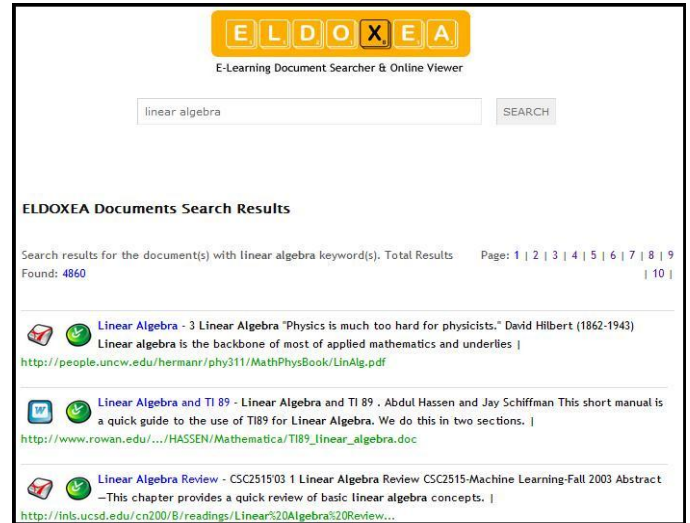6.  If user click one of the results, preview the document files

in same page
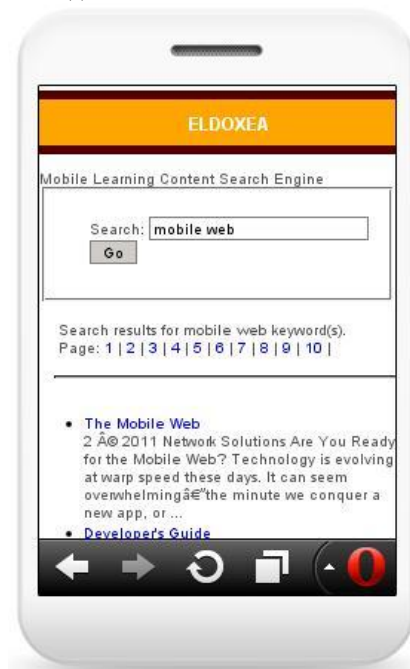7.  Preview display, user can choose Download/Save or close button
8.  If user click Download, then save the files

Example can be seen through http://www.eldoxea.com as shown in Figure 1 ((a) for Internet terminals and (b) for mobile phone).



(a)ELDOXEA for internet terminals



(b) ELDOXEA for mobile phone

Figure 1.   Top page of ELDOXEA

*C.  Supplemental Keyword Selection with the Keywords in the Meta-Tag and the Descriptions in the First Search Result*

In order to choose appropriate supplemental keywords for improvement of hit rate, several attempts are performed. First one is to use *keywords* and *descriptions* in meta-tag of the first web search result. Keywords are used to be in the meta-tag.

---

[1] http://en.wikipedia.org/wiki/Google_Search
[2] http://en.wikipedia.org/wiki/Infoseek
[3] http://en.wikipedia.org/wiki/Yahoo!_Search
[4] http://en.wikipedia.org/wiki/Bing
[5] http://b.hatena.ne.jp/entry/www.eldoxea.com/

Also, there are some keywords in the descriptions in the first web search result. Therefore, these supplemental keywords in the meta-tag and descriptions are applicable to add to the primary keyword. This approach using the assumption that related website should contain the similar keyword, while not related website containing another not related keyword. So, we try to find the intersection of keyword in meta-tag of page header between webs in the search results set.

Figure 2 shows an example of header. In the header, there are meta-tag and descriptions. In these meta-tags and the descriptions, there are some keywords. We could use these keywords as supplemental keywords for search.



Figure 2.    shows an example of header

Figure 3 shows the first three search results based on Google search with the keyword "Linear Algebra". When I visit the first URL of Wikipedia in Japanese, then Figure 4 appears. Then source code can be displayed as shown in Figure 5. Although we cannot get any keyword in the meta-tag sometime or description in the header, it used to be appeared in the header. If we repeat the same keyword twice as keyword for search, then we get the other search results as shown in Figure 6



Figure 3.    First three search results based on Google search with the keyword "Linear Algebra".
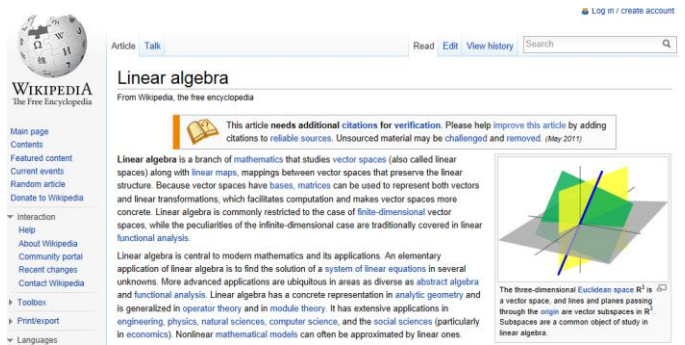


Figure 4.    Top page of the first URL of Wikipedia of Linear Algebra in Japanese



Figure 5.    The source code of the first URL of Linear Algebra.



Figure 6.    Other search results we used to get with double keyword for search (in this case "Linear Algebra is refrained twice)

## D. Supplemental Keyword Selection from Thesaurus Engine

Second one is to use *keywords* which are derived from thesaurus engines. Thesaurus engine provides similar words to the primary keyword with their priority. Therefore, these are used for supplemental keywords. An example of the search results of thesaurus engine with "Linear Algebra" is shown in Figure 7. As shown in Figure 7, similar words to "Linear Algebra" are listed in the order of priority. Also, you can refer to the URLs as results of thesaurus engine as shown in Figure 8. Figure 9 shows an example of top page of the first priority of URLs of the search results of thesaurus engine, *Weblio*[6]. Then we can get the keywords in the meta-tag and the descriptions in the header when you check the source word of the web pages as shown in Figure 10.



Figure 7. An example of search results of thesaurus engine of Weblio with the keyword "Linear Algebra" in Japanese.
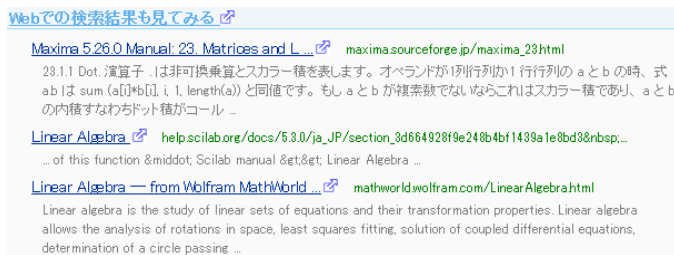


Figure 8. URLs can be referred as results of thesaurus engine (in this case with the keyword of "Linear Algebra").
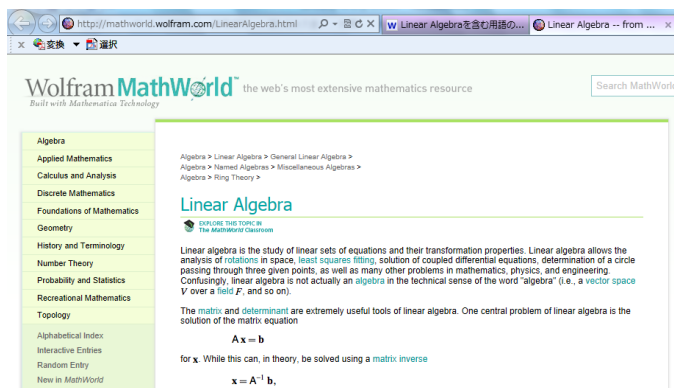


Figure 9. Example of top page of the first priority of URLs of the search results of thesaurus engine, *Weblio*.
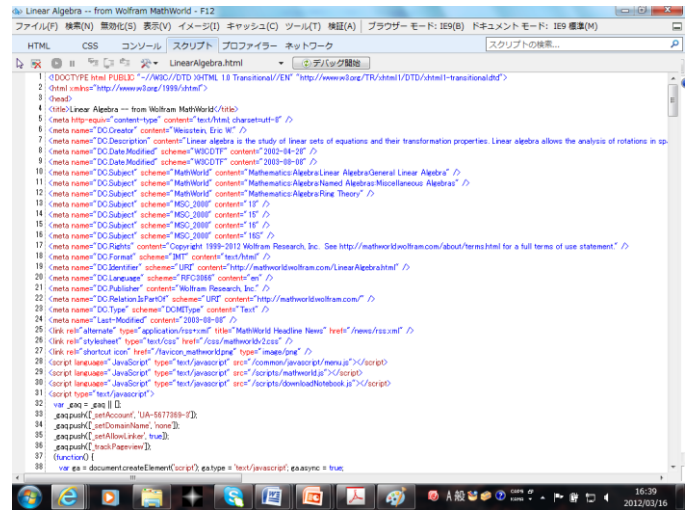


Figure 10. The keywords in the meta-tag and the descriptions in the header when you check the source word of the web pages

After that, we could determine the supplemental keywords in accordance with their priority.

## E. Supplemental Keyword Selection from Twitters

Third one is to use twitter for gathering suggestions of supplemental keywords from twitters. Reliability of the twitters can be evaluated with the previously proposed method[7]. The *Weblio* of thesaurus engine provides the gate for twitter. From this gate, we can bet valuable information of URLs related to the primary keyword. Reliability of the twitter has to be checked though.

## F. Supplemental Keyword Selection from BBS

Fourth one is to use Bulletin Board System: BBS and chatting for gathering suggestions of supplemental keywords from the community members. Firstly, the user has to send a message which is requesting supplemental keywords to BBS system as well as chat with the primary keyword. Then the user has reply message with supplemental keywords. These are to be candidates of supplemental keywords. Most of Learning Management System like *Moodle*[8] provides BBS and chatting capabilities. Using these functions, we can get valuable information relating to the primary keyword.

### III. EXPERIMENTS

## A. Method for Experiments Conducted

We conduct the experiment for our propose methods in two ways. Subjectively search through Yahoo search engine and objectively creating a new search engine using Yahoo API. We check the results of a search engine that only using primary keyword, and then comparing the results while using combination of primary key and supplementary keyword. First, correct answer of Yahoo search with queries of "Linear Algebra" and "Hazardous Material Handler test" is determined. Yahoo search is a kind of hybrid type of search engine that is

---

[6] http://thesaurus.weblio.jp/content/エンジン

[7] http://www.readwriteweb.com/archives/twazzup_a_better_twitter_search_engine.php
http://www.govloop.com/profiles/blogs/twitters-reliability-an-issue
[8] http://moodle.org/

same as ELDOXEA. The first 50 candidates of URLs are selected. All of the 50 sites is visited and evaluated subjectively. Then these sites are divided into "*appropriate*" or "*not appropriate*". The term "*appropriate*" means that the web content is containing e-learning materials that related to the input keyword. The term "*not appropriate*" means that the web content may not related to the input keyword in the area of e-learning.

The first four keywords for the primary keywords, "Linear Algebra" and "Hazardous Material Handler Examination" are as follows,

(1)  Linear Algebra: Linear Algebra, Mathematics, Matrix

  Algebra/Geometry

(2)  Hazardous Material Handler Examination: Hazardous Material Handler Examination, Gasoline Handler, Hazardous, National Certificate

Figure 12 shows hit rate of the proposed method and of the search results with same primary keywords repeatedly. Square denotes the hit rate with the proposed method while the upside down triangle denotes search result with using the same primary keyword "Linear Algebra" repeatedly. In the case, of usage of the same primary keyword of "Linear Algebra" repeatedly, search result shows that hit rate is saturated at the number of supplemental keyword is 1, which is corresponding to the search with two same primary keyword results in maximum hit rate.

On the other hands, the first two supplemental keywords show the maximum hit rate for the "Hazardous Material Handler Examination" case while repeated usage of the same primary keyword does work for the "Hazardous Material Handler Examination" case.

Second, keywords in the meta-tag and the descriptions are extracted from the header of these sites. Then the keywords are sorted with priority depending on their frequency. These processes are automatically done by our search engine system. Figure 11 shows the screenshot of our meta-tag keyword extractor for automatically extract *keywords* from the meta-tag of a search results.


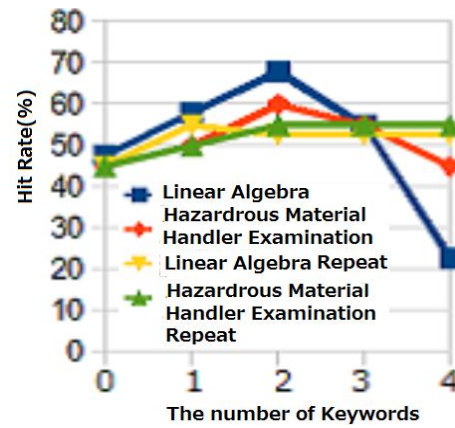
Figure 11. Our automatic Meta Tag Keywords extractor system



Figure 12.  Hit rate of the proposed method and of the search results with same primary keywords repeatedly.

### B.  Supplemental Keyword Selection with Thesaurus Engine, Weblio

As shown in Figure 7, similar words to "Linear Algebra" are listed in the order of priority. There, however, is less related keyword from the Weblio of thesaurus engine. Therefore, it would not be worked for finding supplemental keyword at all.

### C.  Applicability of the Proposed Search Engine

The proposed search engine is applicable to the other primary keywords. Other than "Linear Algebra", the proposed search engine is applied to the chemistry, mechanics, etc. Improvements of hit rate for these primary keywords are evaluated. Figure 13 shows the improvements of hit rate for searching the primary keyword with the primary and supplemental keywords.
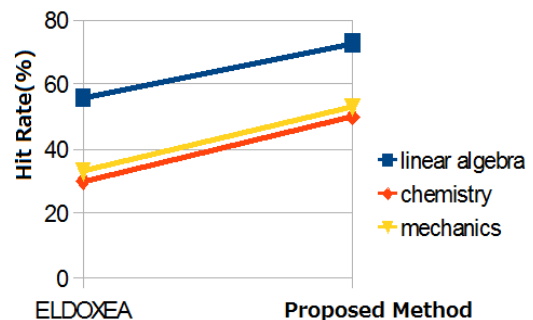


Figure 13.  Improvement of hit rate of the proposed method in comparison to the previously proposed ELDOXEA of e-learning content search engine.

17 to 20% of improvements are confirmed for the primary keywords, Linear Algebra", Chemistry" and "Mechanics".

Other two methods by utilizing twitter and BBS system will be discussed in the other paper in the near future.

### IV.  CONCLUSIONS

The proposed search method uses not only one single primary keyword, but also supplemental keywords which are derived from the keywords in the meta-tag and descriptions in the page header of a website appeared in the first search result.

Hit rate is defined as matching accuracy between subjectively determined success search and the current search results of URLs. The hit rate of the proposed method is compared to that of the search method with the same primary keyword with repeatedly used supplemental keyword. Depending on the number of supplemental keyword, hit rate is increasing. Improvement of the hit rate of the proposed search method is 15 to 20 % while that of the search method with repeatedly used supplemental keyword which is same as primary keyword is around 10 %. It is also found that the proposed search method is applicable to the other primary keywords.

### ACKNOWLEDGMENT

### REFERENCES

[1] S. Brin and L. Page. The anatomy of a large-scale hypertextual Web search engine. In WWW, 1998.

[2] Ricardo Baeza-Yates and Berthier Ribeiro-Neto. Modern Information Retrieval. Addison Wesley, May 1999.

[3] Gandal, Neil (2001). "The dynamics of competition in the internet search engine market". *International Journal of Industrial Organization* **19** (7): 1103–1117

[4] Segev, Elad (2010). Google and the Digital Divide: The Biases of Online Knowledge, Oxford: Chandos Publishing

[5] Vaughan, L. & Thelwall, M. (2004). Search engine coverage bias: evidence and possible causes, Information Processing & Management, 40(4), 693-707

[6] Kohei Arai, Herman Tolle "Module Based Content Adaptation of Composite E-Learning Content for Delivering to Mobile Learners", International Journal of Computer Theory and Engineering (IJCTE), Vol 3, No. 3, pp. 381-386, June 2011

[7] Kohei Arai, Herman Tolle, Efficiency improvements of e-learning document search engine for mobile browser, International Journal of Research and Reviews on Computer Science, 2, 6, 1287-1291, 2011.

### AUTHORS PROFILE

**Kohei Arai,** He received BS, MS and PhD degrees in 1972, 1974 and 1982, respectively. He was with The Institute for Industrial Science, and Technology of the University of Tokyo from April 1974 to December 1978 also was with National Space Development Agency of Japan from January, 1979 to March, 1990. During from 1985 to 1987, he was with Canada Centre for Remote Sensing as a Post Doctoral Fellow of National Science and Engineering Research Council of Canada. He moved to Saga University as a Professor in Department of Information Science on April 1990. He was a councilor for the Aeronautics and Space related to the Technology Committee of the Ministry of Science and Technology during from 1998 to 2000. He was a councilor of Saga University for 2002 and 2003. He also was an executive councilor for the Remote Sensing Society of Japan for 2003 to 2005. He is an Adjunct Professor of University of Arizona, USA since 1998. He also is Vice Chairman of the Commission "A" of ICSU/COSPAR since 2008. He wrote 30 books and published 322 journal papers.

**Herman Tolle,** He graduated Bachelor degree in Electrical Engineering from Brawijaya University, Malang in 1998, also graduated Master degree in Telecommunication Information System from Bandung Institute of Technology (ITB), Bandung in 2002. He is with Engineering Faculty of Brawijaya University from 2002 to present. He is now a Doctoral student in Department of Information Science, Faculty of Science and Engineering, Saga University Japan. He has a major concern of research in image analysis, multimedia, content adaptation and web engineering.