

# Scalable Ranking based on Web Pages and Semantic Search

S.Haseena<sup>1</sup>, R.Latha<sup>2</sup>, A.M.Barani<sup>3</sup>

<sup>1</sup>Department of Computer Application, St. Peters Institute of Higher Education & Research, Chennai, India

<sup>2</sup> Department of Computer Application, St. Peters Institute of Higher Education & Research, Chennai, India

<sup>3</sup> Department of Computer Application, St. Peters Institute of Higher Education & Research, Chennai, India

## ABSTRACT

Recommendation systems can take advantage of semantic reasoning-capabilities to overcome common limitations of current systems and improve the recommendations' quality. In this paper, present a personalized-recommendation system, a system that makes use of representations of items and user-profiles based on ontologies in order to provide semantic applications with personalized services. The recommender uses domain ontologies to enhance the personalization: on the one hand, user's interests are modeled in a more effective and accurate way by applying a domain-based inference method; on the other hand, the stemmer algorithm used by our content-based filtering approach, which provides a measure of the affinity between an item and a user, is enhanced by applying a semantic similarity method. Web Usage Mining plays an important role in recommender systems and web personalization. In this paper, we propose an effective recommender system based on ontology and Web Usage Mining. The first step of the approach is extracting features from web documents and constructing relevant concepts. Then build ontology for the web site use the concepts and significant terms extracted from documents. According to the semantic similarity of web documents to cluster them into different semantic themes, the different themes imply different preferences. The proposed approach integrates semantic knowledge into Web Usage Mining and personalization processes.

**Keywords:** *Browse Rank, Page Rank, Profile based results, Web search, thesauruses, search process.*

## I. INTRODUCTION AND RELATED WORK

Search Engine refers to an enormous info of web resources like sites, newsgroups, programs, images etc. It helps to find info on World Wide internet. User will seek for any info by passing question in style of keywords or phrase. It then searches for relevant info in its info and come to the user. Generally there are 3 basic elements of a look engine as listed below Web Crawler, Database and Search Interfaces.

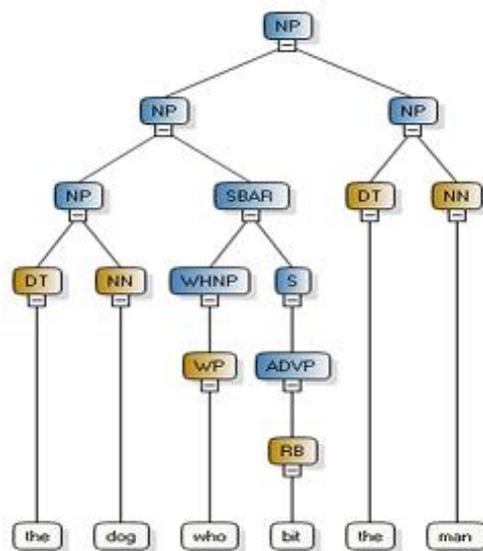
Web Crawler: It is conjointly referred to as spider or bots. It's a software system element that traverses the online to assemble info. Database: All the data on the online is hold on database. It consists of giant internet resources. Search Interfaces is associate interface between user and therefore the database. It helps the user to go looking through the database.

Web crawler, information and also the search interface are the main part of a groundwork engine that really makes program to figure. Search engines build use of Boolean expression AND, OR, to not prohibit and widen the results of a groundwork. Following are the steps that are performed by the search engine: Search appearance for the keyword within the index for predefined database rather than going on to the online to look for the keyword. It then uses software system to look for the data within the database. This software system part is thought as internet crawler.

Once internet crawler finds the pages, the Search Engine then shows the relevant websites as a result. These retrieved websites typically embody title of page, size of text portion, initial many sentences etc.

These search criteria might vary from one program to the opposite. The retrieved data is hierarchic in keeping with numerous factors like frequency of keywords, connection of data, links etc. User will click on any of the search results to open it.

The first downside considers the unit of study. Estimating the semantic similarity of ideas is a very important downside, well studied within the literature. Results of such studies area unit rumored in a very style of fields, including psychology, natural language processing [17], information retrieval [9], [12] language modeling [8], [14], and database systems [3]. Virtually all attempts to study the similarity of concepts model concepts as single words.



**Fig 1. Visualization of Treebank Parser**

We note, however, that the sequence of words in a phrase is important, and that some words contribute more to the meaning of a phrase than others. Consider two phrases: “the dog who bit the man” and “the man who bit the dog.” They contain the same words, so the method in [6] would consider them virtually equivalent, but they convey very different meanings. Thus, we need measures of both term similarities across two terms generally, and term importance, i.e., how critical the term is in the context of the phrase of which it is a part. If we can build a measure of similarity that is weighted by the importance of each term in the context of the phrase. Where it occurs, then we can leverage the similarity framework suggested in [6] with these importance measures.

We consider the words “roll” and “rolled” as conceptually equivalent for the purposes of concept matching. We do this by running each word through a standard stemming algorithm, e.g., the Porter stemmer [15], which reduces each word to its base form by removing common modifications for subject-verb agreement, or variation in parts of speech (e.g., “consider,” “considerably,” “considering,” “consideration”).

To generate the importance of each term, we can use a parser, e.g., OpenNLP [16], that can return the grammatical structure of a sentence. Given an input phrase, such a parser returns a parse tree, where the words in the input phrase that add the most to the meaning of the phrase appear higher in the parse tree than those words that add less to the

meaning of the input phrase. The details of this how this parser works are beyond the scope of this paper, but we note that the output representation is similar to a functional programming language with part-of-speech.

Useful knowledge discovery from Web usage data and satisfactory knowledge representation for effective Web-page recommendations are crucial and challenging. Existing system provide method to efficiently provide better Web-page recommendation through semantic enhancement by integrating the domain and Web usage knowledge of a website. Two new models are proposed to represent the domain knowledge.

The first model uses ontology to represent the domain knowledge. The second model uses one automatically generated semantic network to represent domain terms, Web-pages and the relations between them. Another new model, the conceptual prediction model, is proposed to automatically generate a semantic network of the semantic Web usage knowledge, which is the integration of domain knowledge and Web usage knowledge.

A number of queries have been developed to query about these knowledge bases. Based on these queries, a set of recommendation strategies have been proposed to generate Web-page candidates. The recommendation results have been compared with the results obtained from an advanced existing Web Usage Mining (WUM) method.

Existing recommendation systems are: cold-start, sparsely, overspecialization and domain-dependency. The performance of existing system depends on the sizes of training datasets. The bigger the training dataset size is, predicted pages are limited within the discovered Web access sequences. The domain ontology can be constructed manually by experts, or by automatically learning models is need to design and implement the learning models which can only be done by professionals at the beginning.

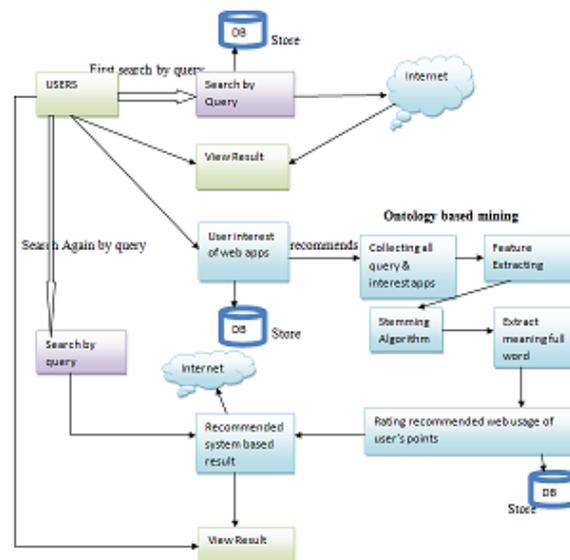
## **II. PROPOSED SOLUTION APPROACH**

In proposed system present a personalized-recommendation system, a system that makes use of representations of items and user-profiles based on ontologies in order to provide semantic applications with personalized services. The semantics method achieved by using two different methods. A domain-based method makes inferences about user's interests and a taxonomy-based similarity method is used to refine the item-user matching algorithm, improving overall results. The recommender proposed is domain-independent, is implemented as a Web service, and uses both explicit and implicit feedback-collection methods to obtain information on user's interests.

Proposed recommender system based on ontology and Web Usage Mining. The first step of the approach is extracting features from web documents and constructing relevant concepts. Then build ontology for the web site use the concepts and significant terms extracted from documents. According to the semantic similarity of web documents to cluster them into different semantic themes, the different themes imply different preferences.

## **III. SOLUTION ARCHITECTURE**

We now describe our implementation architecture, with particular attention to design for scalability. Integrating domain knowledge with Web usage knowledge enhances the performance of recommender systems using ontology-based Web mining techniques. The construction of this model is semi-automated so that the development efforts from developers can be reduced. The user-profile learning algorithm, responsible for expanding and maintaining up-to-date the long-term user's interests, employs a domain-based inference method in combination with other relevance feedback methods to populate more quickly the user profile and therefore reduce the typical cold-start problem. The filtering algorithm, which follows a stemming approach, makes use of a semantic similarity method based on the hierarchical structure of the ontology to refine the item-user matching score calculation.



**Fig. 2. Architecture for Real Time Efficient Web Page Searches Based on Ranking**

Ontology based mining Collecting the query form the user and remove the stop words then perform the stemming, Find the root word;

Step 1: First, if the word is in plural form, it is reduced to singular form. Then, any -ed or -ing endings are removed as appropriate, and finally, words ending in "y" with a vowel in the stem have the "y" changed to "i".

Step 2: Maps double suffixes to single ones when the second-to-last character matches the given letters. So "-ization" (which is "-ize" plus "-ation" becomes "-ize". Mapping to a single character occurrence speeds up the script by reducing the number of possible string searches. Note: for this step (and steps 3 and 4), the algorithm requires that if a suffix match is found (checks longest first), then the step ends, regardless, if a replacement occurred. Some (or many) implementations simply keep searching though a list of suffixes, even if one is found.

Step 3: Works in a similar strategy to step 2, though checking the last character

Step 4: Works similarly to steps 3 and 2, above, though it removes the endings in the context of VCVC (vowel-consonant-vowel-consonant combinations

Step 5: Removes a final "-e" and changes "-ll" to "-l" in the context of VCVC (vowel-consonant-vowel-consonant combinations).

### **Extract meaning full word:**

A user is unlikely enter specifically matching definition in the search, but should be able to enter something conceptually similar to the dictionary meaning. How the words are conceptually associated with each other. We tend to defines many forms of connectedness below.

Synonym set: A set of conceptually connected terms for t.  $W_{syn}(t) = \{t_1, t_2, t_3, \dots, t_n\}$ , is a synonym of t, as defined in the lexicon. For example,  $W_{syn}(\text{read})$  might consist of the set of words  $\{\text{learn, record, register}\}$ .

Antonym set: A set of conceptually negated terms for  $t$ .  $W_{syn}(t) = \{t_1, t_2, t_3, \dots, t_n\}$ , is a synonym of  $t$ , as defined in the lexicon. For example,  $W_{syn}(\text{simple})$  might consist of the set of words  $\{\text{complex, compound}\}$ .

Hypernym set: A set of conceptually a lot of general terms describing  $t$ .  $W_{hyp}(t) = \{t_1, t_2, t_3, t_j, \dots, t_n\}$ , wherever  $t_j$  may be a hypernym of  $t$ , as outlined within the lexicon. For example,  $W_{hyp}(\text{puppy}) = \{\text{dog, pup, domestic dog}\}$  Extract the meaning of the query from the extracted word. Perform the ranking and compare with the web search then produce the result.

#### IV. EXISTING METHOD OF WEB PAGE SEARCHES

##### A. *Conventional search*

The idea of archive has been characterized as any solid or emblematic sign, safeguarded or recorded, for reproducing or for demonstrating a wonder, regardless of whether physical or mental.[4] The developing thought of the archive among Jonathan Priest, Otlet, Briet, Schürmeyer, and alternate documentalists progressively stressed whatever worked as a report instead of conventional physical types of records. The move to computerized innovation would appear to make this refinement considerably more critical. Collect's insightful investigations have demonstrated that an accentuation on the innovation of advanced archives has obstructed our comprehension of computerized records as reports.

##### B. *Text Based Search*

In content recovery, full-content scan alludes to strategies for looking through a solitary PC put away archive or a gathering in a full-content database. The full-content inquiry is recognized from looks in light of metadata or on parts of the first messages stored in databases, (for example, titles, abstracts, chose segments, or bibliographical references).

##### C. *Multimedia Search*

Mixed media look empowers data to seek to utilize inquiries in numerous information sorts including content and other interactive media designs. Sight and sound hunt can be actualized through multimodal look interfaces, i.e., interfaces that permit submitting seek questions as printed asks for as well as through other media[7]. Pursuit is made utilizing the layers in metadata which contain data of the substance of a sight and sound record. Metadata look is less demanding, quicker and successful in light of the fact that as opposed to working with the unpredictable material, for example, a sound, a video or a picture, it seeks utilizing content.

##### D. *Conceptual Search*

An idea seek (or calculated pursuit) is a computerized data recovery strategy that is utilized to look electronically put away unstructured content (for instance, advanced documents, email, logical writing, and so on.) for data that is adroitly like the data gave in a hunt question[2]. As it were, the thoughts communicated in the data recovered in light of an idea seek inquiry are significant to the thoughts contained in the content of the question.

##### E. *Information System Search*

Data Systems Research is an associate inspected scholarly diary that spreads inquire about in the territories of data frameworks and data innovation, including subjective brain research, financial aspects, software engineering, operations look into, plan science, association hypothesis and conduct, human science, and key administration. It is distributed by the Institute for Operations Research and the Management Science and was as of late chosen as one of the main 20 proficient/scholastic diaries by Business Week.[4] Along with Management Information Systems Quarterly,

Information Systems Research is viewed as one of the two most renowned diaries in the data frameworks discipline.[10][11]

#### ***F. Personalised Search***

Customized seek suggests web look encounters that are custom-made particularly to a person's advantages by joining data about the person past particular question gave. There are two general ways to deal with customizing list items, one including adjusting the client's question and the other re-positioning hunt results [19].

#### ***G. Page Rank***

PageRank is a connection investigation calculation and it allows a numerical weighting to every component of a hyperlinked set of reports, for example, the World Wide Web, with the motivation behind measuring its relative significance inside the set. The calculation might be connected to any gathering of substances with complementary citations and references. The numerical weight that it allows to any given component E is alluded to as the PageRank of E and meant by PR (E). Different elements like Author Rank can add to the significance of an element.

#### ***H. Ranking (Information Retrieval)***

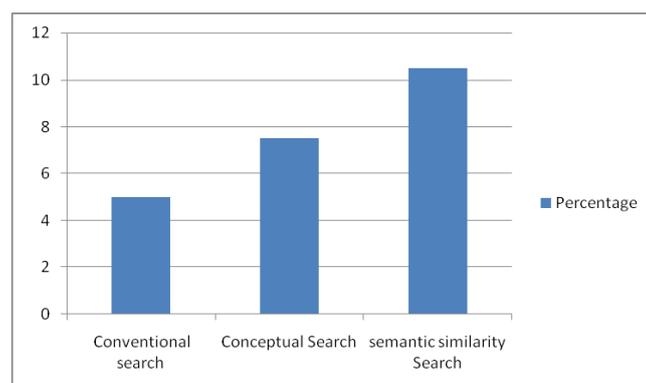
The positioning of inquiry comes about is one of the basic issues in data recovery (IR), the logical/building discipline behind web indexes. Given a question and answer gathering D of records that match the inquiry, the issue is to rank [19], that is, sort, the reports in D as indicated by some model so that the best results seem right on time in the outcome list showed to the client. Traditionally, positioning criteria is relating directly and significantly of records to be communicated in the inquiry.

#### ***I. Computing Search***

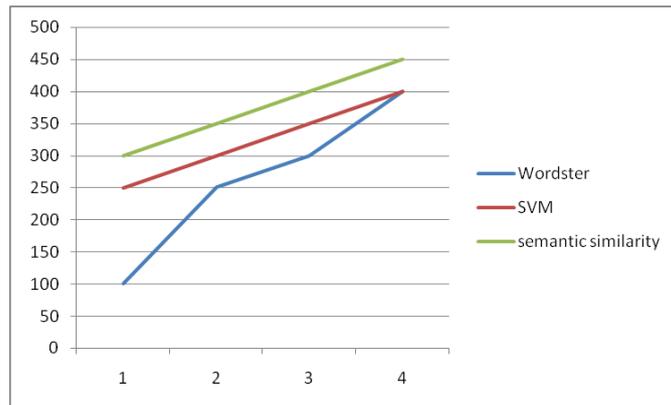
A web search tool is a data recovery framework intended to help discover data put away on a PC framework. The list items are typically exhibited in a rundown and are usually called hits. Web search tools help to limit the time required to discover data and the measure of data which must be counseled, similar to different strategies for overseeing data over-burden. The most open, noticeable type of a web search tool is a Web internet searcher which looks for data on the World Wide Web.

### **V. EXPERIMENTAL RESULTS**

This approach has access to all the word relationship data available and compares the input phrase and very huge impact changes in search rank.

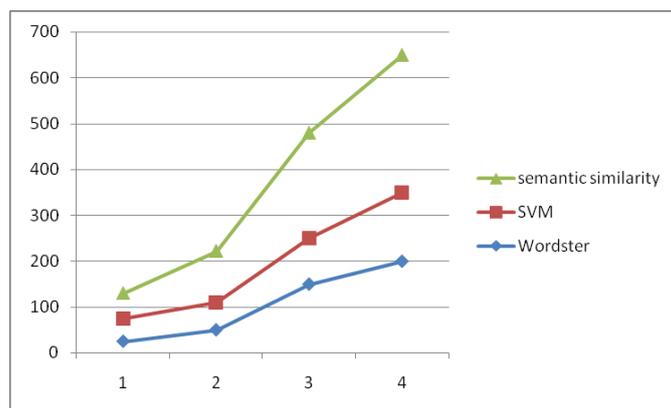


**Fig. 3 Impact changes of Semantic similarity Search and ranking**



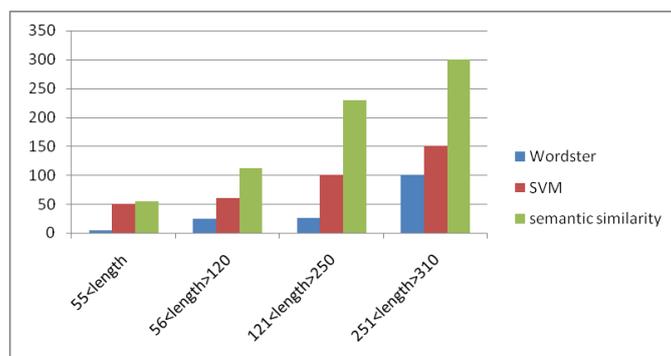
**Fig. 4 Scalability as CPU limits increases**

At runtime, we compared the user request with the model for each search. We ran each of the Wordster, SVM, Semantic similarity, thus the Semantic similarity scalability is increases.



**Fig. 5 Response time performance as request rate increases.**

The curves all show classic exponential growth with increasing load. However, the response time of the Wordster system is about an order of magnitude better than the response time for the SVM approach. Thus the Semantic similarity response time performance is very huge.



**Fig. 6 Accuracy as the value of alpha is varied**

The input search query length is increases. We progressively consider larger and larger numbers of possible output, thereby increasing the probability of including false positives, which causes precision to decrease as increases, while simultaneously increasing the probability finding expected results, which causes recall having an increasing trend as  $\alpha$  increase.

## VI. CONCLUSIONS

On account of an entire literary hunt, the initial phase in grouping of website pages is to discover a list of thing that may relate explicitly to the pursuit term. Earlier, web indexes started with a little rundown of URLs as a purported seed list, brought the substance, and parsed the connections on those pages for pertinent to data, which gave new connections. The procedure was very patterned and preceded until the point when enough pages were found for the searcher's utilization. Nowadays, a constant creep strategy is utilized rather than an accidental revelation in light of a seed list. The creep technique is an augmentation of previously mentioned revelation strategy. But there is no seed list on the grounds that the framework works constantly.

## REFERENCES

- [1] Adriaanse L and Rensleigh C (2013). Web of Science, Scopus and Google Scholar. The Electronic Library, 31(6), 727-744.
- [2] Aniko Hannak, Piotr Sapiezynski, Arash Molavi, Kakhki, Balachander Krishnamurthy, David Lazer, Alan Mislove and Christo Wilson (2013) "Measuring Personalization of Web Search (PDF)" Archived from the original (PDF) on April 25, 2013.
- [3] S. Berchtold, D.A. Keim, and H.-P. Kriegel, "Using Extended Feature Objects for Partial Similarity Retrieval," The VLDB J., vol. 6, no. 4, pp. 333-348, Nov. 1997.
- [4] Bornmann, L, Leydesdorff L and Mutz R (2013). The use of percentiles and percentile rank classes in the analysis of bibliometric data: Opportunities and limits. Journal of Informetrics, 7(1), 158-165.
- [4] D. Coltuc and J.-M. Chassery, "Very fast watermarking by reversible contrast mapping," IEEE Signal Process. Lett., vol. 14, no. 4, pp. 255-258, Apr. 2007. Christian Quast, Elmar, Priesse Pelin, Yilmaz Jan, Gerken Timmy, Schweer Pablo, Yarza Jörg, Peplies Frank and Oliver Glöckner (2013) "The SILVA ribosomal RNA gene database project: improved data processing and web-based tools" at Nucleic Acids Research, Volume 41, Issue D1, 1 January 2013, Pages D590-D596, <https://doi.org/10.1093/nar/gks1219>.
- [5] H. Cui, R. Sun, K. Li, M.-Y. Kan, and T.-S. Chua, "Question answering Passage Retrieval Using Dependency Relations," Proc. 28th Ann. Int'l ACM SIGIR Conf. Research and Development in Information Retrieval, pp. 400-407, 2005.
- [6] T. Dao and T. Simpson, "Measuring Similarity between Sentences," [http://opensvn.csie.org/WordNetDotNet/trunk/Projects/Thanh/Paper/WordNetDotNet\\_Semantic\\_Similarity.pdf](http://opensvn.csie.org/WordNetDotNet/trunk/Projects/Thanh/Paper/WordNetDotNet_Semantic_Similarity.pdf) (last accessed 16 Oct. 2009), 2009.
- [7] Jayanthi.J and Dr.K.S.Jayakumar (2011) "An integrated Page Ranking Algorithm for Personalized Web Search". In International Journal of Computer Applications (0975- 8887), Volume 12-No.11, January 2011.
- [8] J. Lafferty and C. Zhai, "Document Language Models, Query Models, and Risk Minimization for Information Retrieval," Proc. 24th Ann. Int'l ACM SIGIR Conf. Research and Development in Information Retrieval, pp. 111-119, 2001.

- [9] D. Lin, "An Information-Theoretic Definition of Similarity," Proc. Int'l Conf. Machine Learning, 1998.
- [10] Manikandan, R and Dr. R.Latha (2017). "A literature survey of existing map matching algorithm for navigation technology," International journal of Engineering Sciences & Research Technology", 6(9), 326-331. Retrieved September 15, 2017.
- [11] Mercy paul Selvan, A . Chandra Sekar and A. Priya Dharshini (2012) "Survey on Web page Ranking Algorithms" at International Journal of Computer Applications, Vol.41, No-19, ISSN 0975-8887.
- [12] R. Mihalcea, C. Corley, and C. Strapparava, "Corpus-Based and Knowledge-Based Measures of Text Semantic Similarity," Proc. Nat'l Conf. Artificial Intelligence, 2006.
- [13] Northcott, D. and Linacre S (2010). Producing spaces for academic discourse: The impact of research assessment exercises and journal quality rankings. Australian Accounting Review, 20(1), 38- 54.
- [14] Ponte and W. Croft, "A Language Modeling Approach to Information Retrieval," Proc. 21st Ann. Int'l ACM SIGIR Conf. Research and Development in Information Retrieval, pp. 275-281, 1998.
- [15] M. Porter, "The Porter Stemming Algorithm," <http://tartarus.org/martin/PorterStemmer/>, 2009.
- [16] O.S. Project "Opennlp," <http://opennlp.sourceforge.net/>, 2009.
- [17] P. Resnik, "Semantic Similarity in a Taxonomy: An Information- Based Measure and Its Application to Problems of Ambiguity in Natural Language," J. Artificial Intelligence Research, vol. 11, pp. 95- 130, 1999.
- [18] Sergey Brin and Lawrence Page (2012) "The anatomy of large scale hypertextual search engine" at Computer Networks Elsevier, vol.56, issue-18, pp3825-3833.
- [19] Mitsuo Yokokawa, Fumiyoshi Shoji, Atsuya Uno, Motoyoshi Kurokawa and Tadashi Watanabe (2011). The K computer: Japanese next-generation supercomputer development project. International Symposium on Low Power Electronics and Design (ISLPED). pp. 371–372. doi:10.1109/ISLPED.2011.5993668 (1–3 August 2011).