

# Information retrieval for children based on the Aggregated search paradigm

Sergio Duarte  
University of Twente

## 1 Introduction

The Internet is increasingly being used by children for information and entertainment purposes. Moreover, it has recently shown that children are often trusted to search the Internet on their own [12]. Unfortunately, most of the current Information Retrieval (IR) systems are designed for adults and previous studies have shown that the information needs and search approaches of children and adults differ substantially. Particularly, children have been found to be less focused during the search process. They often follow a non-linear navigational style, in which previous searches and clicked hyper-links are reactivated frequently [4]. This behavior along with the lack of logical progression in the exploration of results expose the disorientation children experience searching the Web and their difficulty to decide which information is relevant [3]. Children also have difficulty constructing meaning from the results, specially in the case of complex information needs that required pieces of information for different sources [3]. All these search characteristics hamper the search success and search experience of children. It is important to mention that although these observations were drawn from case-studies, we recently were able to observe this search behavior on a large-scale by identifying queries retrieving information for children based on the DMOZ kids & teens category and the domains clicked on a commercial search engine query log [6].

The aim of this research is to develop information services for children by expanding and adapting current Information retrieval (IR) technologies according to the search characteristics and needs of children. Concretely, we will employ the aggregated search paradigm as theoretical framework [11]. Aggregated search refers to the selection of results from diverse sources and the presentation of these results on an single result page by organizing them in a coherent way beyond the classic result list provided by current search engines [10, 7].

We consider that this paradigm is promising to address the search characteristics of children because the presentation of different information sources can highly reduce the cognitive load of children to find information. The coherent presentation of results using different information types (e.g., wiki, images, videos) can also help children to have a better understanding of the information searched. Additionally, parents or moderators can preselect several trusted sources using this paradigm.

Sushmita et al.[13] provided some evidence to support that aggregated search interfaces help users to find relevant information faster and to increase the quan-

tity and diversity of items retrieved per information request. We intend to provide evidence that this paradigm also helps children to search from the Web. For this purpose we will address the following central research question:

- What are the most appropriate information types and sources that an IR system should select given a child’s information request and how to present this information to ease and improve his/her search experience?

We identify 4 challenges that need to be addressed to answer this research question: (1) Evaluate the performance of the verticals of existing search engines and propose methods to improve these results for children’s information needs. (2) Explore methods to dynamically select the most relevant verticals given the user’s request. (3) Study methods to present the information from these verticals to improve the search for children. (4) Explore mechanism to discover verticals to evaluate our methods on complex scenarios. The research questions and research methods associated to each one of these challenges are described as follow.

## 2 Evaluation of current search engines from the children IR perspective

We plan to explore the following research questions:

- How well do commercial search engines satisfy children’s information needs and how appropriate are the results for users aged 8 to 12 years old?
- How effective are simple strategies to enrich queries (e.g., as adding *for kids* to the queries) to improve the performance of search engines to retrieve relevant and appropriate information for children?

The motivation of this research is to evaluate the adequacy of the verticals of current search engines to satisfy children’s information needs in terms of the relevance and appropriateness of the results. More importantly, we want to explore methods to query these verticals to improve the quality of the results considering that the users are children. We speculate that simple query expansion methods using cue words (e.g., adding *for kids* to the query) have high potential to improve the performance of current search engines to retrieve children-friendly content. Concretely, we will evaluate the web, video, images, news, blogs and wiki verticals of two major search engines. We will also evaluate the results from search engines oriented for children as *Yahoo! Kids*.

*Evaluation:* For this purpose a set of topics reflecting realistic children’s information needs will be created based on the queries and sessions we identified in [6] given that these queries are a reasonable approximation of children’s information needs on the web. A test collection will be created based on these topics to evaluate the system in a similar fashion as it is done in TREC.

## 3 Vertical selection in IR for children

The second challenge is to design and evaluate methods to dynamically select the most relevant verticals given the user’s request. The research question that we will address is:

- Does the presentation of results from suitable verticals improve the search experience of children given his/her information request?

Arguello et al. [1] presents a classifier based on query string, query log and vertical-corpora features to select the best vertical for an input query. We will explore the use of social bookmarking systems as an external source of evidence for vertical selection. Concretely, graph-based expert finding methods will be employed to identify users annotating resources for children. The annotations from these users (e.g., bookmarks tags) will be used to learn aspects and the types of information that are found in the Web for the most popular children topics, which can be achieved using clustering methods [2] on the annotations provided by these users. Bridging these annotations to the verticals will allow IR systems to decide which verticals the system should utilize given the topics involved in the user's request. The extraction of query features is a feasible way to solve this task given that it has been shown that the tags used to describe Web resources commonly overlap with popular query terms [8].

*Evaluation:* We hypothesize that we can improve the search experience of children by displaying information from the verticals in the same result page. To verify this hypothesis, we will perform two test studies. In the first test-study the accuracy of the vertical selection method and the performance gain obtained from mining social bookmarking systems will be evaluated by gathering manual assessments, in a similar fashion as in [1]. Once we show that these resources are useful for the vertical selection task, we will perform a second case-study with children to determine if the IR systems displaying information from these verticals on the same result page actually improve the search experience of children.

## 4 Content aggregation in IR for children

In the *organization phase* of the aggregated search paradigm the results of the verticals are reorganized using a meaningful criteria [10]. In this research the challenge is to organize the results to improve the search experience of children. The research question will be addressed is:

- What is the most effective way to organize the information from the verticals to reduce the cognitive load of children searching the web?

We speculate that the alignment of results from different verticals using probabilist content models [9] can successfully provide *multi-modal views* of the information retrieved. This can be achieved by placing together (e.g., contiguous blocks) those results from the verticals that are better associated by these models. This is a promising approach to guide the search and improve the understanding of the information requested, specially when the information need require different content types to be satisfied.

*Evaluation:* The evaluation will be performed using a case-study with children to verify which organization scheme is more effective and if the methods to align results is beneficial to reduce the search cognitive load.

## 5 Vertical discovery for children users

As a follow up we will address the following research questions:

- Is the mining of social bookmarking systems beneficial to find safe web search services (verticals) for children?
- How well perform the methods we proposed under complex scenarios with a greater number of verticals?

We plan to discover verticals that contain information relevant for children. We will carry out this task by extending expert finding techniques [5, 14] to find users that share information suitable for children in social bookmarking systems. The motivation is to exploit the social knowledge of thousands of users to find web resources that can be accessed by IR systems using search standards as OpenSearch, which is a set of formats for sharing search results<sup>1</sup>. The motivation of this research is twofold: (1) To provide semi-automatic mechanisms to scale-up the *trusted* resources of IR systems for children.(2) Evaluate the scalability of the methods proposed for vertical selection and information presentation under richer scenarios.

*Evaluation:* The evaluation of the two research questions can be carried out by employing the evaluation methodology developed for the first and second challenge described in this proposal, respectively.

## 6 Evaluation challenges

A key point in our research is the evaluation of the search success and experience of children using IR systems. Currently, there is research gap in evaluation methodologies to evaluate the retrieval performance of systems when the final users are children. Thus, there is an salient need in the research community to discuss and to agree in the criteria and methodology to perform this evaluation, specially on a large-scale. This observation has also been drawn by Bilal [4], who pointed out that the understanding children have from the results is an important factor to include in the evaluation besides measuring the visualization of relevant results. We also pointed out that standard measures of success in query logs need to be revised to better capture the search success of children [6].

## 7 Conclusions

The contributions of this PhD proposal can be summarized as follow: (1) We will measure how well verticals from state-of-the-art search engines satisfy children's information needs in terms of appropriateness and relevance. We will create novel methods to enrich queries to improve the performance of these verticals given the information needs of children. The evaluation of the search engines and our method will also represent a contribution since we will create a test collection to evaluate the performance of IR systems for these users. To the best of our knowledge the creation of test collections to evaluate IR for children has

---

<sup>1</sup><http://www.opensearch.org/>

not been addressed before. (2) We will provide new sources of evidence for the task of vertical selection based on social bookmarking systems. For this purpose we will expand current expert finding techniques to mine the social knowledge of users sharing safe information for children.(3) We will explore innovative methods to present the information from the verticals to reduce the cognitive load of children searching the web and to increase their understanding on the topics being searched. Evaluation methodologies to measure the efficiency of these methods will also contribute to the understanding of children search behavior. (4) We will design automatic vertical discovery methods to provide safe mechanisms to scale-up *trusted* resources in IR systems for children.

## References

- [1] J. Arguello, F. Diaz, J. Callan, and J.-F. Crespo. Sources of evidence for vertical selection. In *SIGIR '09*, pages 315–322, New York, NY, USA, 2009. ACM.
- [2] J. Aslam, K. Pelekhev, and D. Rus. Static and dynamic information organization with star clusters. In *In Proceedings of the 1998 Conference on Information Knowledge Management*, pages 208–217, 1998.
- [3] D. Bilal. Children’s use of the yahoooligans! web search engine: Ii. cognitive and physical behaviors on research tasks. *J. Am. Soc. Inf. Sci. Technol.*, 52(2):118–136, 2001.
- [4] D. Bilal. Children’s use of the yahoooligans! web search engine. iii. cognitive and physical behaviors on fully self-generated search tasks. *J. Am. Soc. Inf. Sci. Technol.*, 53(13):1170–1183, 2002.
- [5] C. S. Campbell, P. P. Maglio, A. Cozzi, and B. Dom. Expertise identification using email communications. In *In CIKM '03: Proceedings of the twelfth international conference on Information and knowledge management*, pages 528–531. ACM Press, 2003.
- [6] S. Duarte, D. Hiemstra, and P. Serdyukov. Query log analysis in the context of information retrieval for children. *To appear in SIGIR 2010*.
- [7] K. Gyllstrom and M. Moens. A picture is worth a thousand search results: Finding child-oriented multimedia results with collage. *To appear in SIGIR 2010*.
- [8] P. Heymann, G. Koutrika, and H. Garcia-Molina. Can social bookmarking improve web search? In *WSDM '08: Proceedings of the international conference on Web search and web data mining*, pages 195–206, New York, NY, USA, 2008. ACM.
- [9] M.-F. M. Koen Deschacht1. Finding the best picture: Cross-media retrieval of content. In *Advances in Information Retrieval*, pages 539–546. Springer Berlin / Heidelberg, 2008.
- [10] A. Kopliku. Aggregated search: From information nuggets to aggregated documents. In *CORIA*, pages 495–502, 2009.

- [11] V. Murdock and M. Lalmas. Workshop on aggregated search. *SIGIR Forum*, 42(2):80–83, 2008.
- [12] Ofcom. Uk children’s media literacy: Research document, March 2010.
- [13] S. Sushmita, H. Joho, and M. Lalmas. A task-based evaluation of an aggregated search interface. In *SPIRE*, pages 322–333, 2009.
- [14] J. Zhang, J. Tang, and J. Li. Expert finding in a social networks. In *Proc. of Database Systems for Advanced Applications (DASFAA ’2007)*.