

## Chapter 13

# Exploring Semantic Characteristics of Socially Constructed Knowledge Repository to Optimize Web Search

**Dengya Zhu**

*Curtin University, Australia*

**Heinz Dreher**

*Curtin University, Australia*

### **ABSTRACT**

*Short-term queries preferred by most users often result in a list of Web search results with low precision from a user perspective. The purpose of this research is to improve the relevance of Web search results via search-term disambiguation and ontological filtering of search results based on socially constructed search concepts. A Special Search Browser (SSB) is developed where semantic characteristics of the socially constructed knowledge repository are extracted to form a category-document set. kNN is employed with the extracted category-documents as training data to classify Web results. Users' selected categories are employed to present the search results. Experimental results based on five experts' judgments over 250 hits from Yahoo! API demonstrate that utilizing the socially constructed search concepts to categorize and filter search results can improve precision by 23.5%, from Yahoo's 41.7% to 65.2% of SSB based on the results of five selected ambiguous search-terms.*

## INTRODUCTION

The introduction and subsequent explosion of the Web has dramatically changed our approach to access and use of information. Internet is becoming a part of life for most people in the world. However, as indicated by Baeza-Yates and Ribeiro-Neto (1999), most users have difficulties in expressing their information needs in search-term format: they prefer short queries instead of the Boolean expressions (Jansen & Spink, 2006). To address this issue, most search engines encourage users to enter very short search terms as queries, and then return a list of search results which are ranked by technologies such as traditional information retrieval models and PageRank (Page, et al., 1998), according to the relevance degree of the results with respect to a given query. However, as the volume of information on the Web is becoming unbelievably huge, short search terms based Web search usually leads to search engines return a list of thousands, even millions of search results. Searchers are thus frustrated when facing such a long list of results especially when half of the search results are irrelevant to their information needs (Gauch, Chaffee, & Pretschner, 2003). It is now commonly recognized that information search services are far from perfect. The challenges of search engines are summarized in Table 1.

The first challenge for search engines is the *Information overload* (C1). What does “1-10 of 55,400,000 for jaguar” mean? Can we really access 55,400,000 information items about jaguar? Are all of the items are relevant to my information need? Oh my God, how can manage to read all of them. It seems that search engines prefer to present a number of search results as huge as this for short queries. However, as research indicated (Jansen & Spink, 2006), the tendency is fewer results pages are browsed. Therefore, millions of search results are a kind of information overload for users.

The second challenge of search engines is *Mismatching hits* (C2). The performance of an

information retrieval system is usually measured by precision and recall. Precision is an evaluation of how retrieved results of the information retrieval system are relevant; whereas the recall is an evaluation of how the relevant results are retrieved (Baeza-Yates & Ribeiro-Neto, 1999). While millions of search results being retrieved, meaning that the recall is very high, precision remains low, because most the retrieved results are irrelevant. For example, the top ten search results of a Chinese name “Wei Liu” are about ten different persons (Zhu, 2007), and you are lucky if the person you are searching for is among the top ten results; or on the other hand, the information you are looking for is located on the tenth page.

The third issue of search engines is that search results are presented in the form of a *flat list* (C3). The flat list of results is suitable for a small amount of results because it provides an easy and quick way to locate a relevant information item. However, when there are hundreds or even thousands of retrieved results returned, users have to do re-search amongst the returned items; they need to go through page by page to pick up useful information items. Finding relevant information

Table 1. Challenges of web search engines (Zhu, 2007)

	Challenges	Phenomena
C1	Information overload	Millions of Web hits
C2	Mismatching hits	High recall, low precision, many irrelevant results
C3	Flat list of results	Results are presented in a flat list, users have to pick up useful items among the list, like finding a needle in haystack
C4	Mismatching mental model	Automatically formed hierarchy used to re-organize Web hits usually mismatches human mental model
C5	Homogeneity	Search engines present “the same for all” hits, not personalized
C6	Low recall of Web navigation	Web navigation is more accurate, but the recall is very low

20 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage:

[www.igi-global.com/chapter/exploring-semantic-characteristics-socially-constructed/71859](http://www.igi-global.com/chapter/exploring-semantic-characteristics-socially-constructed/71859)

## Related Content

---

### An Efficient Lightweight Network Based on Magnetic Resonance Images for Predicting Alzheimer's Disease

Boan Ji, Huabin Wang, Mengxin Zhang, Borun Mao and Xuejun Li (2022). *International Journal on Semantic Web and Information Systems* (pp. 1-18).

[www.irma-international.org/article/an-efficient-lightweight-network-based-on-magnetic-resonance-images-for-predicting-alzheimers-disease/313715](http://www.irma-international.org/article/an-efficient-lightweight-network-based-on-magnetic-resonance-images-for-predicting-alzheimers-disease/313715)

### Research Synthesis and Thematic Analysis of Twitter Through Bibliometric Analysis

Saleha Noor, Yi Guo, Syed Hamad Hassan Shah, M. Saqib Nawaz and Atif Saleem Butt (2020). *International Journal on Semantic Web and Information Systems* (pp. 88-109).

[www.irma-international.org/article/research-synthesis-and-thematic-analysis-of-twitter-through-bibliometric-analysis/256548](http://www.irma-international.org/article/research-synthesis-and-thematic-analysis-of-twitter-through-bibliometric-analysis/256548)

### Semi-Automatic Knowledge Extraction to Enrich Open Linked Data

Elena Baralis, Giulia Bruno, Tania Cerquitelli, Silvia Chiusano, Alessandro Fiori and Alberto Grand (2013). *Cases on Open-Linked Data and Semantic Web Applications* (pp. 156-180).

[www.irma-international.org/chapter/semi-automatic-knowledge-extraction-enrich/77204](http://www.irma-international.org/chapter/semi-automatic-knowledge-extraction-enrich/77204)

### Modeling for Research and Communication

Gilbert Paquette (2010). *Visual Knowledge Modeling for Semantic Web Technologies: Models and Ontologies* (pp. 439-465).

[www.irma-international.org/chapter/modeling-research-communication/44943](http://www.irma-international.org/chapter/modeling-research-communication/44943)

### An Idea Ontology for Innovation Management\*

Christoph Riedl, Norman May, Jan Finzen, Stephan Stathel, Viktor Kaufman and Helmut Krcmar (2009). *International Journal on Semantic Web and Information Systems* (pp. 1-18).

[www.irma-international.org/article/idea-ontology-innovation-management/41734](http://www.irma-international.org/article/idea-ontology-innovation-management/41734)