# Towards Visual Web Search: Interactive Query Formulation and Search Result Visualization

Marian Dörk
mdoerk@ucalgary.ca

Carey Williamson
carey@cpsc.ucalgary.ca

Sheelagh Carpendale
sheelagh@cpsc.ucalgary.ca

Department of Computer Science, University of Calgary
2500 University Drive NW, Calgary, AB, Canada T2N 1N4

## ABSTRACT

Search on today's Web is influenced by the early Web that was primarily text-based: search parameters are typically entered as text queries and the resulting resources are mostly displayed as textual lists. To gain an overview of the information space and the retrieved resources, the information seeker has to issue several search queries and skim many search results. In this paper we show how visualization widgets (VisGets) are a viable way to query and visualize multiple types of Web data. We have applied VisGets to three prominent resource types on the Web—hypertext, syndicated content, and Semantic Web data—and discuss the limitations of the VisGets prototype and the challenges of this approach.

## Keywords

Information retrieval, information visualization, World Wide Web, exploratory search.

## 1. INTRODUCTION

As we approach the 20th anniversary of the World Wide Web, it is appropriate to reconsider the interfaces used to query and navigate Web-based information. The early Web was mostly unstructured and text-based, so text-based search mechanisms were appropriate and effective. However, today's Web is much richer, featuring a vast amount of information, with a wide range of formats (e.g., news feeds, photos, music, and videos). Keyword searches, with long lists of ranked results, may not provide the best way to navigate large and diverse Web-based information collections. While text-based searches are effective for precise information needs, they can be frustrating for less specific information needs. In particular, it is difficult to acquire a general overview of an information collection or data repository.

In this paper, we show how recent visualization advances can improve query and navigation of increasingly structured Web data. The underlying motivation for this work is two-fold. First, the rapid evolution of the Web has produced a vast, rich, and diverse information repository requiring better interfaces for information exploration on the Web. Second, the Web browser itself is now highly suitable for interactive visualizations that can be used for search query formulation and result set exploration. While visualization widgets (VisGets) [8] allow for interactive query formulation for navigating RSS feeds, we illustrate in this work that such interactive visualizations are viable for a larger range of Web data by extending VisGets to three types of Web resources (publications, photos, and famous people) revealing strengths, limitations, and areas for further research.

## 2. RELATED WORK

Early information navigation in online systems has been described as *berrypicking* [4] referring to a notion of information seeking comparable to dynamic roaming between different sets of documents. As searchers gather multiple information bits they learn more about the content and may change their information needs and search queries. Such open-ended information seeking has been recently revisited as *exploratory search* that evolves from an initially vague information need into a learning process [11]. However, these approaches involve low-level navigation. The onus is on the information seeker to gain an overview of the information space and orientation within it by formulating multiple queries and considering many search results. This however, becomes problematic as the number of resources increases, along with the volume of possibly relevant information to navigate.

A large body of prior work has shown that interactive visualizations can considerably improve the exploration of databases. Instead of entering text-based queries, information seekers can explore databases by interacting with graphical user interface elements [1]. Furthermore, results can be augmented with visualizations representing, for example, document-query similarities [9]. *Coordinated views* can provide richer insights into interdependencies between different dimensions within an information space [3]. For example, interactive visualizations that feature maps, topics, and timelines have been developed for the exploration of video libraries [7]. While coordinated visualizations have been successful when applied to static databases, the use of multidimensional information visualization for search within dynamic information spaces such as the Web invites further research.

There has been a significant increase in visual and interactive interfaces on the Web over the last few years. Partly due to improved support for Web standards, simple interactive visualizations are now possible within the Web browser. Furthermore, emerging data on the Web is increasingly structured and semantically organized [5]. A large portion of research on visualizing hypertext and semantic data focused on navigation hierarchies [12] and node-link diagrams [13]. However, hierarchies and graphs representing Web data emphasize low-level data representation rather than viewers' information needs. In contrast, geographic maps, time sliders, and other interactive widgets have been introduced to ease authoring and exploration of small-scale semantic data sets within the Web browser [10]. Large-scale information repositories such as Wikipedia, can be transformed into semantic information repositories [2], which enable sophisticated queries, but still lack accessible search and browse interfaces. While RSS (Really Simple Syndication) feeds containing structured metadata can be explored with spatial and temporal filters [6], these controls do not provide visual overviews along multiple dimensions.

# 3. VISGETS

Based on earlier work on RSS feed exploration using visualization widgets (VisGets) [8], we created a system that allows the visual search for three different types of Web data (see next section). In general, a VisGet is an information visualization widget that combines interactive query formulation with visual summarization of search results along a particular dimension. In doing so, it can be thought of as using this dimension to create a facet or view into the data. Since Web resources are often structured along time, location, and tags, the first three VisGets use these relatively common types of Web-based information. The following briefly summarizes VisGet functionality (for further details see [8]).

**Time VisGet.** Interactive bar charts indicate the temporal distribution of information items and allow formulation of temporal queries (Figure 1: Left). This can be done either by selecting an individual bar or moving sliders to select a temporal sub-region.

**Location VisGet.** The spatial distribution of information items is indicated in the location VisGet (Figure 1: Middle), with circles placed on regions of a map. Zooming into the map increases the spatial resolution, whereas zooming out aggregates close regions to clusters. Selecting circles, changing zoom levels, or panning the map also modifies the spatial query constraints.

**Tag VisGet.** The topical composition of the information space is summarized with an interactive tag cloud, in which font sizes indicate relative numbers of information items (Figure 1: Right). Selecting tags defines filters in the topical domain.
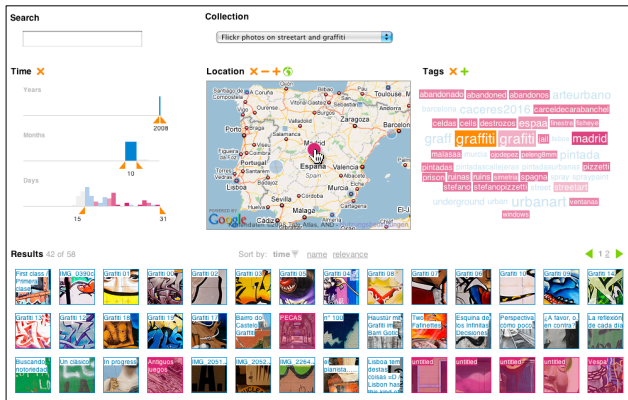


**Figure 1: VisGets allow exploration of resources along multiple facets. Hovering over elements highlights relations.**

VisGets allow for two types of coordinated interactions: multidimensional query refinement and weighted brushing. The combination of query constraints as defined in VisGets allows multidimensional query formulation. On the other hand, weighted brushing constitutes a more transient interaction technique that only requires hovering the mouse pointer over a visual element causing related elements in the interface to be highlighted. In addition, a text search (Figure 1: Left, top) complements the query parameters set with VisGets. Thus, the searcher can integrate text and visual search or switch between these different ways of information seeking.

The information items complying with the multidimensional query are displayed in the area below the VisGets (Figure 1: bottom). The results are displayed as squares with the title of the information source and, if present, an image embedded. Hovering over a result item shows a preview. As the searcher changes query parameters with the VisGets, information items are removed and added through animated transitions. The VisGets interface utilizes common browser features, for example, to provide an interaction history and bookmarking of search results (see Figure 2).
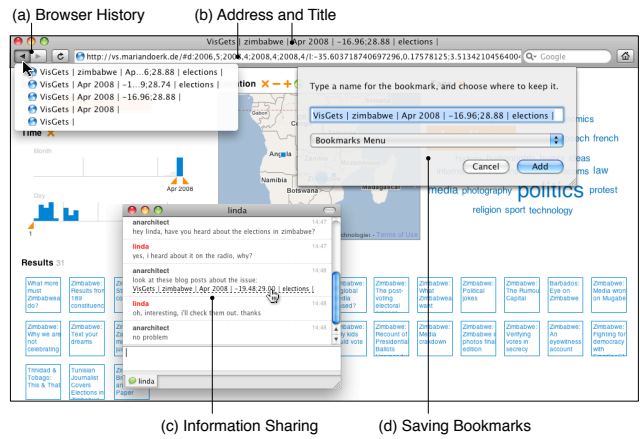


**Figure 2: VisGets use conventional Web browser functions.**

# 4. CASE STUDIES

To better understand the potential of coordinated visualizations for Web-based query formulation and result navigation, we have applied VisGets to three resource types: academic publications from the ACM digital library, geo-referenced photos from Flickr, and famous people on Wikipedia. These examples represent data formats that are widely available on the Web, namely hypertext, syndicated content, and Semantic Web data.

## 4.1 Hypertext: Academic Publications

Our hypertext example considers academic publications in the ACM digital library, specifically those from the past eight years of the WWW conference. These data provide a good example of rich hypertext content, with many implicit dimensions that require either sophisticated analysis or manual extraction. We used a semi-automatic approach, in which we identified relevant structural entities on each publication page for batch data processing. We have chosen the date of publication for the temporal dimension, the institutional affiliation of the first author for the location, and the paper's indicated `General Terms` and `Keywords` as tags. Furthermore, we extracted additional metadata per publication (e.g., title, authors, abstract) for detail-on-demand display in the results list.

The data extraction was carried out using regular expressions defined based on the relevant information context on the page. The extracted information snippets were then integrated into our database schema. For location information, we used the affiliation string, by stripping the first part (institution) and using the rest (city, state, country) to query the GeoNames Web service [15] for geographic (latitude and longitude) information. We automated this process, relying on the consistent presentation of structural information in ACM's digital library.

Exploring publications in digital libraries, such as ACM's, traditionally involves entering search terms, refining text-based queries, and following hyperlinks. VisGets facilitate new ways to explore academic publications, by providing visual overview and summarization information about the document collection. For example, a research interest (e.g., "search") can be entered as a text-based query (see Figure 3). The search interface provides not only a set of results, but also the temporal, spatial, and topical distribution of publications in the VisGets and mechanisms to refine the query along those dimensions. The time VisGet provides a visual presentation of publication trends over the years, while additional brushing shows how the spatial distribution of research institutions and prominent research topics change over time. The location VisGet highlights influential research institutions; brushing these shows
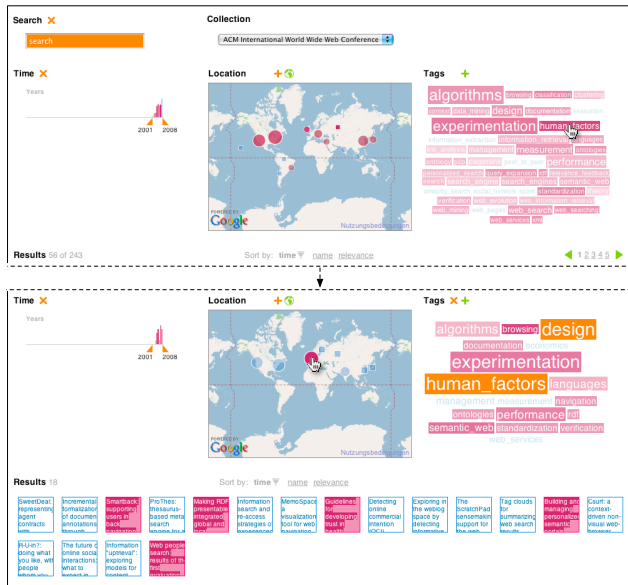
**Figure 3: Academic publications can be explored with visual widgets by setting query filters and brushing visual elements.**

their research topic specializations and a temporal characterization of their activity. Providing a broad overview of the paper themes, the tag VisGet allows the searcher to refine a text-based query with multiple tags as filters (e.g., "design" and "human factors"), explore those research themes that occur within the selected publications and the region being currently activated (see Figure 3).

## 4.2 Syndicated Content: Photos

There are many types of resources on the Web that are made available via RSS feeds. As an example for this category, we have selected photos published on Flickr. With the growing use of location-aware devices, photos and other Web resources increasingly have geographic metadata. Flickr provides GeoRSS feeds that include geographic information associated with the shared photographs. The dimensions chosen for this information space are the time of creation, the location where the photo was taken, and the attached keywords. These all fit well with our existing set of VisGets. The extraction of the corresponding structural properties (timestamp, latitude and longitude, and tags) is straightforward, as described in our earlier work [8].

Shared photos on the Web can be explored in many different ways: entering search terms, going through photo listings, using interactive geographic maps, or accessing tag clouds. While these methods facilitate browsing along specific dimensions, the separated display of either thumbnails, map markers, or tags alone does not provide an adequate overview of the rich, multi-faceted information space. With VisGets, multiple dimensions of a photo collection can be visualized and made accessible interactively at the same time. As an example, consider a traveller with interests in stencil graffiti and street art who is planning a trip to the Iberian Peninsula. Selecting the target location on the geographic map provides an easy filter on the information space, with a resulting gallery of photos displayed. Selecting "graffiti" via the tag cloud or with a text-based entry narrows the information collection as well (see Figure 1). The time slider can be used to restrict the photos to recent ones, or to a specified time of year. Different locations on the map could then be explored through weighted brushing indicating different types of graffiti in different cities or regions.

## 4.3 Semantic Web: Famous People

Semantic Web data constitutes a particularly compelling type of Web content, since it provides the most sophisticated data modelling. While the Semantic Web is still arguably in its infancy, future information on the Web will increasingly be semantically structured. As an example of semantic data, we selected a subset of Wikipedia's content: people born after 1900. We have chosen birthdate, place of birth, and occupation as the primary conceptual dimensions, to align with the temporal, spatial, and topical VisGets. These properties (and many others) are made available by the DBpedia project as Linked Data [2]. This means that structural information is made explicit as RDF and available through the semantic query language SPARQL.

To aggregate information about famous people from Wikipedia, the VisGets system requests data subsets from the DBpedia endpoint. In the current prototype, data is cached in a MySQL database, since SPARQL queries are often too slow for responsive interaction. Furthermore, the SPARQL standard does not yet provide query methods that can be used for visual aggregations. Because DBpedia does not provide a cache invalidation mechanism, cached entries are regularly updated. Extracting birth dates and locations can be done by using DBpedia's ontology. For occupations, our VisGets system uses the YAGO ontology [14]. In this ontology, general concepts such as "politician" or "football player", are associated with people via the `rdf:type` relation. This information can then be used in the tag VisGet.



**Figure 4: Using weighted brushing to explore relationships between occupation and age of famous people on Wikipedia.**

Typically, one navigates information on the Wikipedia by entering a search query and navigating hyperlinks from one entry to another. Our navigation experiments with the Wikipedia data set considered famous personalities such as athletes, politicians, singers, and writers. The time VisGet provides interesting visual overviews of these data collections, highlighting the temporal distributions for different types of famous people. For example, most famous football (soccer) players are young, while most famous politicians are not—both observations can be made by brushing the corresponding tags (see Figure 4).

## 5. ONGOING WORK

Having applied our VisGets-based prototype to three different information spaces with different data formats has sparked further research addressing a wide range of challenges.

**Interface Customization.** Our current system is based on a fixed set of dimensions (time, location, tags) used for the initial set of VisGets. Currently, when any of these dimensions is absent, the resource is excluded from the interface. However, the notion of filtering based on certain properties could also be utilized to query certain types of resources based on different sets of VisGets

that could be activated by the searcher. That is, the information seeker could invoke VisGets spontaneously to explore, for example, spatio-temporal resources using only the time and location VisGets. This way the visual search interface could be customized based on the information needs a searcher wishes to pursue.

**User Studies.** While we believe that VisGets have the potential to support a new type of information seeking, it is necessary to better understand how search behaviour is affected by visualization widgets. Do VisGets enrich or impede Web search? In particular, it would be interesting to investigate which kind of search tasks may be facilitated by visualization widgets and when conventional text-based queries are more efficient for satisfying an information need. To find answers to these and further questions regarding visual information exploration on the Web we intend to conduct studies deploying visual search interfaces in real-world contexts.

**Additional VisGets.** Besides time, location, and tags, many other dimensions remain to be considered, especially in domain-specific contexts. For example, price is a relevant dimension for some resources on the Web, since it could be used as a discriminating factor by an information seeker wishing to explore products. A price VisGet could provide a visual overview by indicating the distribution of prices and at the same time allow query formulation using price intervals. In combination with temporal and spatial VisGets the search interface could help making purchase decisions. The challenge is to extract and visualize a greater variety of dimensions as VisGets for other common types of concepts found on the Web, such as text documents, multimedia content, and products. For this, several existing visualization techniques can be appropriated to support both visual summarization and query formulation along many possible dimensions.

**Structure and Analysis.** VisGets rely on meaningful and extractable dimensions. There are two main strategies for achieving semantically structured resources on the Web: either provide semantic data models or employ data mining techniques to derive meaning from unstructured representations. While our prototype system currently employs simplistic data extraction of Web resources with differing degrees of structure, the success of this approach ultimately relies on both semantic data modelling and data mining. With advances in both research domains, 'explorable' dimensions can be derived from increasing amounts of Web data. The semantic structures emerging from more advanced data modelling and analysis will require and enable more sophisticated query formulation and search result presentation.

**Search Engine Integration.** So far, our prototype browses relatively small sets of Web resources compared to the scale of the Web as a whole. There are significant challenges that are to be addressed to integrate VisGets with general Web search engines. Providing visual summaries and responsive interactivity for millions of web resources requires advanced large-scale indexing architectures that go beyond the conventional text-based query-response paradigm. Can visual overviews be pre-computed when there are multiple dimensions with virtually endless combinations of query parameters? How can visualization generation be integrated with traditional keyword indexing of Web search engines?

## 6. CONCLUSION

In this paper, we have discussed the role of visualization for Web-based query formulation and search result presentation. For this, we have extended and applied VisGets to three different Web-based information collections with diverse data types, namely academic publications (hypertext), geo-referenced photos (syndicated content), and person data (Semantic Web). We have described the use of VisGets for these collections, and discussed the efforts undertaken for data extraction. These case studies illustrate both potential and limitations of VisGets while also identifying significant challenges related to the immense scale and diversity of the Web. These challenges are the drivers for our ongoing research on visual exploration of Web resources.

## 7. ACKNOWLEDGMENTS

## 8. REFERENCES

[1] C. Ahlberg and B. Shneiderman. Visual information seeking: Tight coupling of dynamic query filters with starfield displays. In *CHI '94: Conference on Human factors in computing systems*, pages 313–317. ACM Press, 1994.

[2] S. Auer, C. Bizer, G. Kobilarov, J. Lehmann, R. Cyganiak, and Z. Ives. DBpedia: A nucleus for a web of open data. In *ISWC '08: International Semantic Web Conference*. Semantic Web Science Association, 2007.

[3] M. Q. W. Baldonado, A. Woodruff, and A. Kuchinsky. Guidelines for using multiple views in information visualization. In *AVI '00: Advanced Visual Interfaces*, pages 110–119. ACM Press, 2000.

[4] M. J. Bates. The design of browsing and berrypicking techniques for the online search interface. *Online Information Review*, 13(5):407–424, 1989.

[5] T. Berners-Lee, J. Hendler, and O. Lassila. The semantic web. *Scientific American*, 284(5):28–37, 2001.

[6] Y.-F. Chen, G. D. Fabbrizio, D. Gibbon, R. Jana, S. Jora, B. Renger, and B. Wei. GeoTracker: geospatial and temporal rss navigation. In *WWW '07: World Wide Web Conference*, pages 41–50. ACM Press, 2007.

[7] M. Christel. Accessing news video libraries through dynamic information extraction, summarization, and visualization. *Visual Interfaces to Digital Libraries*, 2539:98–115, 2002.

[8] M. Dörk, S. Carpendale, C. Collins, and C. Williamson. VisGets: Coordinated visualizations for web-based information exploration and discovery. *IEEE Trans. on Visualization and Computer Graphics*, 14(6):1205–1212, 2008.

[9] M. A. Hearst. Tilebars: Visualization of term distribution information in full text information access. In *CHI '95: Conference on Human factors in computing systems*, pages 59–66. ACM Press, 1995.

[10] D. F. Huynh, D. R. Karger, and R. C. Miller. Exhibit: Lightweight structured data publishing. In *WWW '07: World Wide Web Conference*, pages 737–746. ACM Press, 2007.

[11] G. Marchionini. Exploratory search: from finding to understanding. *Comm. of the ACM*, 49(4):41–46, 2006.

[12] D. Nation, C. Plaisant, G. Marchionini, and A. Komlodi. Visualizing websites using a hierarchical table of contents browser: WebTOC. In *Proc. 3rd Conference on Human Factors and the Web*, 1997.

[13] L. Reeve, H. Han, and C. Chen. *Visualizing the Semantic Web: XML-Based Internet and Information Visualization*, chapter Information Visualization and the Semantic Web, pages 19–44. Springer, 2006.

[14] F. Suchanek, G. Kasneci, and G. Weikum. YAGO: a core of semantic knowledge. In *WWW '07: World Wide Web Conference*, pages 697–706. ACM Press, 2007.

[15] M. Wick. Geonames. http://www.geonames.org/ (Retrieved 2008-10-31).