

Singapore Management University

Institutional Knowledge at Singapore Management University

Research Collection School Of Computing and Information Systems

School of Computing and Information Systems

4-2014

CeleBrowser: An example of browsing big data on small device

Song TAN

Chong-wah NGO

Singapore Management University, cwngo@smu.edu.sg

Jun XU

Yong RUI

Follow this and additional works at: https://ink.library.smu.edu.sg/sis_research



Part of the [Databases and Information Systems Commons](#), and the [Graphics and Human Computer Interfaces Commons](#)

Citation

TAN, Song; NGO, Chong-wah; XU, Jun; and RUI, Yong. CeleBrowser: An example of browsing big data on small device. (2014). *Proceedings of the 4th ACM International Conference on Multimedia Retrieval, ICMR 2014, Glasgow, United Kingdom, April 1-4*. 514-517.

Available at: https://ink.library.smu.edu.sg/sis_research/6496

This Conference Proceeding Article is brought to you for free and open access by the School of Computing and Information Systems at Institutional Knowledge at Singapore Management University. It has been accepted for inclusion in Research Collection School Of Computing and Information Systems by an authorized administrator of Institutional Knowledge at Singapore Management University. For more information, please email cherylids@smu.edu.sg.

CeleBrowser: An example of browsing big data on small device

Song Tan
Department of Computer
Science
City University of Hong Kong
Kowloon, Hong Kong
sophysongtan@gmail.com

Chong-Wah Ngo
Department of Computer
Science
City University of Hong Kong
Kowloon, Hong Kong
cscwngo@cityu.edu.hk

Jun Xu
Department of Automation
University of Science and
Technology of China
Hefei, China
junx1992@gmail.com

Yong Rui
Microsoft Research
Beijing, 100080, China
yongrui@microsoft.com

ABSTRACT

In this demonstration, we demonstrate a mobile-based celebrity video browsing system called CeleBrowser. Using this system, users can interactively switch among four views: people-centric, timeline-centric, month-centric and topic-centric, for browsing celebrity-related hot videos. A peculiarity of the demonstration is to highlight the advantage of multi-perspective information organization and presentation in engaging users for exploratory browsing of large number of Web videos on a device with small screen. Technology-wise the demonstration shows how query logs collected for six months from two vertical search engines are leveraged for mining hot events and videos of celebrities.

Categories and Subject Descriptors

H.3.3 [Information Search and Retrieval]: Retrieval models

Keywords

Video browsing, log analysis, event extraction, visual summarization

1. INTRODUCTION

Searching the right videos to watch is always difficult. Today's search engines mostly return a long list of videos, where the relationship among videos are not always clear. It thus becomes the users' responsibility to pick the right videos from the bunch of somewhat related and noisy videos. Searching in such a scenario on smart phones in particular is tedious due to the limited displayed space. A vivid way of presenting is by summarizing the available information, for

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the Owner/Author.

Copyright is held by the owner/author(s).
ICMR '14, Apr 01-04 2014, Glasgow, United Kingdom
ACM 978-1-4503-2782-4/14/04.
<http://dx.doi.org/10.1145/2578726.2582614>

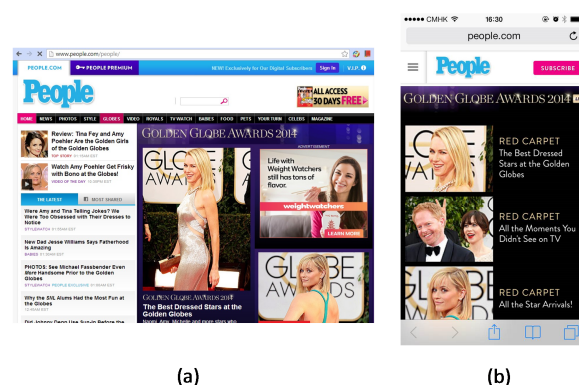


Figure 1: The desktop (a) and mobile phone (b) versions of the website People.com.

example through metadata and browsing history, to provide multi-perspective access of videos.

This paper presents a celebrity video browser named CeleBrowser, which organizes celebrity-related Web videos from multiple perspectives for fast access and browsing of videos. According to Bing search volume statistics, celebrity search consumes around 10% of Internet traffic. Due to huge demand, there are professional-edited websites such as TMZ [3] and PopSugar [2], as well as magazines such as US weekly [4] and People [1], featuring the breaking news, photos, videos, movies and gossips about celebrities. Producing such websites, nevertheless, requires significant manual work and is thus expensive. The large varieties of information across time are also difficult to be summarized, and as a consequence most websites do not provide the retrospective view of celebrities unless certain milestone events such as "death of Michael Jackson" occurs. Furthermore, to squeeze the rich set of information for small screen display, the information available for desktop version is often simplified, as shown in Figure 1. In such a case, exploratory browsing of celebrities could sometimes become even more difficult.

CeleBrowser is different from the celebrity-based websites, where the summaries are generated automatically based on the analysis of user query logs from two vertical search en-

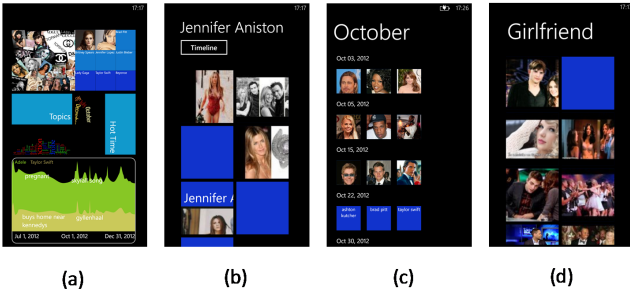


Figure 2: A celebrity video search system: (a) home page, (b) list view of milestone videos related to Jennifer Aniston, (c) list view of hot videos in October, and (d) videos of topic “girlfriend”, developed on Windows Phone 8.

gines. The summaries are multi-perspective, providing quick access to the hot videos of a celebrity, a month and a topic, as shown in Figure 2. Retrospective views of important videos and topics are also supported for visualizing the level of importance for videos and topics across time. CeleBrowser is developed on Windows Phone 8. Capitalized on the animated and resizable tiles, multi-perspective of information is naturally displayed with dynamic content in tiles of different sizes to reflect the event importance.

2. TECHNOLOGIES

The technologies behind CeleBrowser are the mining of milestone periods, topics, and videos from user query logs. The processing pipeline is content-free, and thus the underlying algorithms are extremely efficient in mining. Query logs have been considered by commercial search engines as one of the most effective ways for improving search results [9], and have been leveraged for tasks such as query reformulation [6] and video recommendation [5].

2.1 Hot time detection

The detection of milestone periods is by thresholding the search volume of queries about a celebrity. Basically the periods that above the thresholds are regarded as “hot times” of a celebrity. With the intuition that most queries about an event are issued on the day of occurrence, thresholding of search volume is generally a practical strategy, which are also adopted in works like [7]. Nevertheless, picking a right threshold is not always easy and will render different sets of hot times. In our studies, nevertheless, a more critical issue than picking a right threshold is the selection of a search volume for analysis.

In the celebrity domain, both the vertical search engines of news and videos receive a large volume of celebrity-related queries. The distributions of queries across time between the two search engines, nevertheless, are quite different. For video search, the queries tend to spread over few days after an event occurs. While for news search, the queries are mostly peaked at the days of event occurrence. We speculate that the very first queries most users input to are the news search engine, for the reasons that news articles are used to be regarded as the “first hand” information, which are professionally edited and originated from the authority sources. On the other hand, there is usually a lag of time between video upload and the actual date of event occurrence. The

queries to video search engine are also more diverse than that of news search engine. For example, when the news of Adele pregnant leak, the top queries in news log are all relevant to this event. However, there are a large portion of queries issued to search for the MTVs and performances of Adele in the video log.

In our system, the hot time detection relies only on news search log. The setting of threshold in detection is not particularly sensitive because the queries of hot events are naturally peaked on the day of occurrence in news domain. There are few exception, however, such as the event like Tom Cruz’s divorce, the peak happens a few days later after the burst of the news. This is because the peripheral events following the divorce indeed attract more attention than the original event.

2.2 Mining hot queries and videos

The mining is grounded on the intuition that queries referring to the milestone events often result in the clicks of representative documents. Similarly, documents describing milestone events are accessed through queries relevant to the events. Thus, a straightforward solution for this problem is to discover the mutual relationship between queries and documents through the user click through data. We employ reinforcement algorithm for mining the milestone topics and videos. Specifically, the former refers to hot queries of users, and the latter is the set of representative videos for hot queries.

Nevertheless, direct applying the reinforcement algorithm on query logs from video search engine does not perform well in our empirical studies. This is mainly because some queries are loosely related to the event of occurrence. The search for porn videos, particularly, is very frequent and evenly distributes throughout any period of time. Thus, we integrate query logs from both news and video search engines as a tripartite graph for reinforcement learning, as shown in Figure 3. Basically, the queries from the news and videos engines form an interfacing layer for updating the importance of videos as well as news articles.

Let \vec{q} , \vec{v} , \vec{s} denotes the query, video and news layers respectively. The query layer \vec{q} is formed by a set of queries, where each query is represented as a node. Similarly, \vec{v} and \vec{s} are formed by a set of videos and news articles being clicked by a query q_i in \vec{q} . The \mathbf{M} and \mathbf{N} denote the query-by-video and query-by-news matrix respectively. Each element in the matrix denotes the number of clicks from a query to a video or news. Each column of the matrix is normalized so that the summation of column elements is equal to 1.

Each node in the bipartite graph carries a weight. Using \vec{q} as the interfacing layer, the scores in the video layer \vec{v} can be propagated to the scores in the news layer \vec{s} through the following equation:

$$\vec{s} = \mathbf{N}^T \mathbf{M} \vec{v} \quad (1)$$

Conversely, \vec{s} can be propagated to \vec{v} through the following equation:

$$\vec{v} = \mathbf{M}^T \mathbf{N} \vec{s} \quad (2)$$

The initial scores of video and news, denoted v_0 and s_0 respectively, are set to 1. Through the aforementioned propagation process, the scores for videos and news articles will

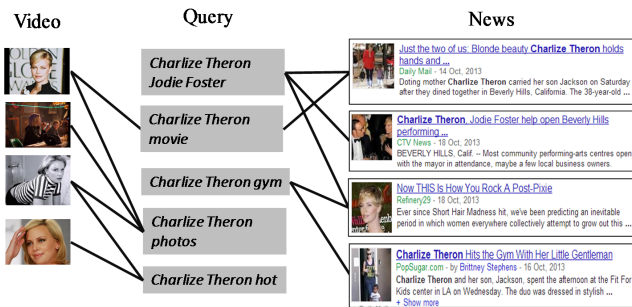


Figure 3: Video-Query-News reinforcement: an example showing the videos and news articles clicked by the queries collected from two vertical search engines.

be as follows after n -iteration:

$$\vec{v} = (\mathbf{M}^T \mathbf{N} \mathbf{N}^T \mathbf{M})^{n-1} \mathbf{M}^T \mathbf{N} \vec{v}_0 \quad (3)$$

$$\vec{s} = (\mathbf{N}^T \mathbf{M} \mathbf{M}^T \mathbf{N})^n \vec{s}_0 \quad (4)$$

Because the matrices $\mathbf{M}^T \mathbf{N} \mathbf{N}^T \mathbf{M}$ and $\mathbf{N}^T \mathbf{M} \mathbf{M}^T \mathbf{N}$ are symmetric, the scores of videos and news will be converged to the eigenvectors corresponding to the largest eigenvalues of the matrices. The scores of video and news reflect their relevant importance. In the implementation, the user queries containing the names of celebrities as the keywords are used to build the bipartite graph. The reinforcement algorithm is conducted in a temporal window with day as the unit. Only days which are regarded as “hot” (as described in section 2.1) will be considered for reinforcement learning.

3. MULTI-PERSPECTIVE ACCESS

With the milestone topics (or queries) and videos mined from the click-through data, these information are organized into celebrity-centric, timeline-centric, month-centric and topic-centric views for multi-perspective information access. Each view is visualized by an animated tile as shown in Figure 4.

The celebrity-centric view shows a list of celebrities as icons (Figure 4a). Touching an icon brings to next page that offers two options of showing the milestone videos of the celebrity in the list-wise or timeline manner.

The topic-centric view lists milestone topics of celebrities sorted in alphabetical order. The topics are pooled from the candidate topics mined from hot days. Basically, we consider only the top-10 hot queries of each hot day. Some examples of popular topics are “divorce”, “baby”, “bikini”, “wardrobe malfunction”. To visualize the popularities of topics, a tag cloud is generated for analytics as shown in Figure 4b.

The month-centric view displays the list of milestone videos sorted by months. To provide vivid visualization, each month is represented by an animated tile, with the size of a tile proportional to the number of milestone videos in a month (Figure 4c). When a tile is touched, a list of milestone videos sorted by the date of a month will be shown.

There are two kinds of timeline-centric views. The first type is shown in Figure 4a, where a trend-line extracted from the statistics of news search volume is shown for each celebrity. The trend-line indicates the hot degree of each

day. The thumbnails of milestone videos are attached to the trend-line for quick understanding of topic evolution over time for a celebrity. Figure 4a shows an example of timeline view for browsing the milestone videos of Jennifer Anniston. A video is represented by a thumbnail plus video title. The video will be displayed if the corresponding icon is touched. The second type of timeline view is shown in Figure 4, which offer the concurrent view of two celebrities across time [8]. The interface is designed as like two rivers of different width running across time, where the relative width between the rivers indicates the degree of hotness at a particular time. Along the rivers, salient topics are plotted such that users can have a quick overview as well as comparison of major topics between two celebrities.

4. DEMONSTRATION

The demonstration will show the mining results for 38 celebrities who are listed in the Forbes Celebrities 100 list. The celebrities in the list is widely recognized as the world’s “most powerful celebrities”. The user queries and click-through data used for mining are obtained from the Microsoft Bing news and video search engines. These logs were collected for a total duration of six months from July to December of year 2012. In total, there are 495,622 queries, 323,792 news and 540,374 videos being clicked.

5. USER STUDY

We compare three very different systems for celebrity browsing: People.com, YouTube and CeleBrowser. People.com is a professionally edited website, while the information presented by CeleBrowser is automatically generated based on query logs analysis. YouTube provides search function, while People.com and CeleBrowser only present summaries for browsing. A total of 12 human subjects are invited to evaluate the three systems. All the evaluations are done on smart phones. Except Celebrowser which should be run on Windows Phone platform, we do not restrict the types of smart phones to be used for evaluation. Basically, the subjects are asked to use their own mobile phones to complete the evaluation.

The subjects are split into two groups. The first group compares People.com and CeleBrowser, while the second group evaluates YouTube and CeleBrowsers. Each subject is requested to answer two questions using a system. To prevent potential bias due to the question difficulty, the questions assigned to evaluators are designed in such a way that each question will be answered by different subjects using different systems. The questions being asked are: 1) Who did Taylor Swift date?, 2) Find the videos about Katy Perry’s date at pool?, 3) Among the hollywood stars who have divorce issue? and 4) Who shoot Jennifer Aniston and Justin Theroux’s first photo?. The evaluators are informed of the answers happened within the period of July to December of 2012, and are requested to complete all the questions within 20 mins. There are more than one answer for question 1 and 3. The answers for all the questions cannot be easily located by querying search engines such as Google and Bing.

Each subject fills in a questionnaire after complete the evaluation. The questionnaire includes two following criterions in the scale of 1 (not acceptable) to 7 (highly preferable):

- Presentation: To what degree the organization and presentation of celebrity-related information help in

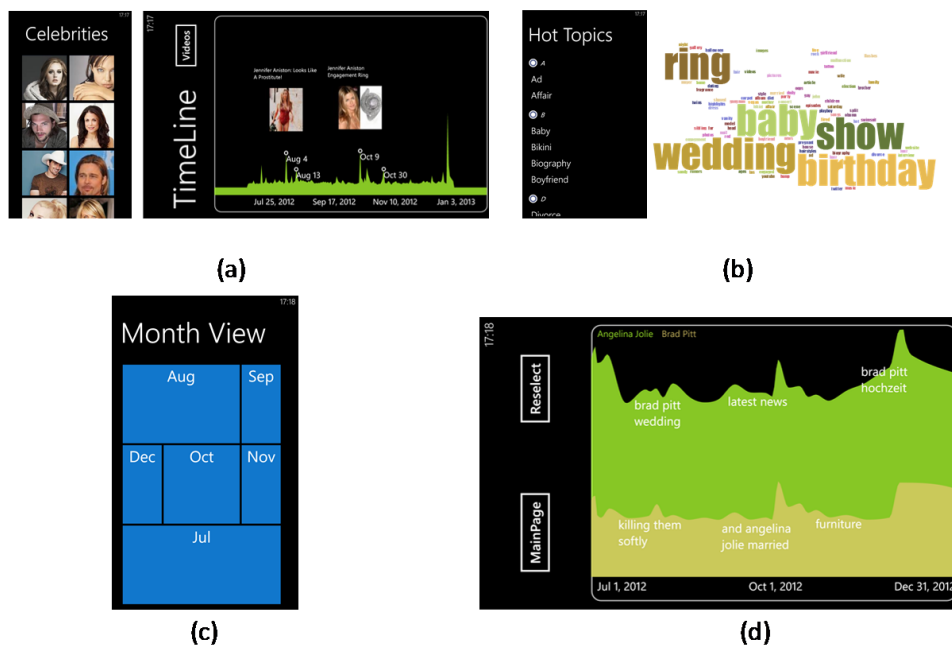


Figure 4: Multi-perspective organization of information: (a) celebrity-centric view and timeline visualization of milestone videos with Jennifer Anniston as an example, (b) topic-centric and its tag cloud generated from hot topics, (c) matrix view of hot month, and (d) concurrent timeline view of two celebrities Brad Pitt and Angelina Jolie.

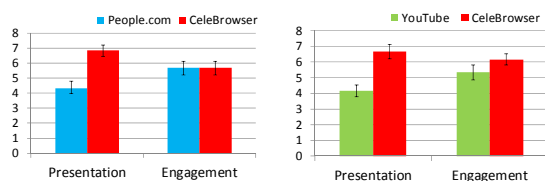


Figure 5: Subjective evaluation. Pairwise comparison of three systems using two criteria.

browsing?

- Engagement: To what degree do you enjoy using the system?

Figure 5 shows the result of user evaluation. Basically using CeleBrowser answers the most number of questions, followed by YouTube and People.com. From the questionnaires, the rating of CeleBrowser for the criterion Presentation is close to the perfect score (7) compared to two other systems which achieves a score of slightly more than 4. The result basically indicates that most evaluators agree that multi-perspective organization of information as presented by CeleBrowser is relatively helpful. For engagement criterion, CeleBrowser achieves higher score than YouTube and same score as People.com. The fact that People.com offers attractive and interesting content by professional editing results in high score of Engagement criterion. CeleBrowser, on the other hand, receives high rating because of the design of multi-perspective information access for supporting exploratory browsing.

6. ACKNOWLEDGMENTS

The work described in this paper was fully supported by a grant from the Research Grants Council of the Hong Kong Special Administrative Region, China (CityU 119610), and a grant from Microsoft Research Asia Windows Phone Academic Program (FY12-RES-OPP-107).

7. REFERENCES

- [1] <http://www.people.com/>.
- [2] <http://www.popsugar.com/>.
- [3] <http://www.t TMZ.com/>.
- [4] <http://www.usmagazine.com/>.
- [5] S. Baluja, R. Seth, D. Sivakumar, Y. Jing, J. Yagnik, S. Kumar, D. Ravichandran, and M. Aly. Video suggestion and discovery for youtube: taking random walks through the view graph. In *Proceedings of the 17th international conference on World Wide Web, WWW '08*, pages 895–904, New York, NY, USA, 2008. ACM.
- [6] Y. Hu, Y. Qian, H. Li, D. Jiang, J. Pei, and Q. Zheng. Mining query subtopics from search log data. In *Proceedings of the 35th international ACM SIGIR conference on Research and development in information retrieval, SIGIR '12*, pages 305–314, New York, NY, USA, 2012. ACM.
- [7] S. Tan, C.-W. Ngo, H.-K. Tan, and L. Pang. Cross media hyperlinking for search topic browsing. In *ACM Multimedia*, pages 243–252. ACM, 2011.
- [8] F. Wei, S. Liu, Y. Song, S. Pan, M. X. Zhou, W. Qian, L. Shi, L. Tan, and Q. Zhang. Tiara: A visual exploratory text analytic system. In *Proceedings of the 16th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '10*, pages 153–162, New York, NY, USA, 2010. ACM.
- [9] G.-R. Xue, H.-J. Zeng, Z. Chen, Y. Yu, W.-Y. Ma, W. Xi, and W. Fan. Optimizing web search using web click-through data. In *Proceedings of the thirteenth ACM international conference on Information and knowledge management, CIKM '04*, pages 118–126, New York, NY, USA, 2004. ACM.