

WebBeholder: A Source of Community Interests and Trends based on Cooperative Change Monitoring Service on the Web

Santi Saeyor

Mitsuru Ishizuka

Department of Information and Communication Engineering

School of Engineering, The University of Tokyo

7-3-1 Hongo, Bunkyo-ku, Tokyo 113-8656, JAPAN

{santi,ishizuka}@miv.t.u-tokyo.ac.jp

Abstract

Besides the target information itself, the changes upon the previously released information are significant and worth being notified to those who perceived the out of date information as soon as possible. The changes made in Web repositories occur at unpredictable rates. Unfortunately, stock type information source has no means to inform its prospective users about the changes. While the stock type information source occupies a large percentage of sources on the Web, it is necessary to have a system that monitors changes on the Web, and provides comprehensive presentation to the prospective users. In this paper, we propose a mechanism that provides change monitoring and presentation service for a large group of users by coalition among service agents. The service agents keep improving the overall utilization factor by several schemes based on the decision made by game analysis. We apply a cost model to the service mechanism in order to study the cooperative behavior of the service agents. The reduction of cost is designed to comply with the level of cooperation among service agents. The WebBeholder community is a promising source of current trends and interests.

1 INTRODUCTION

The explosive growth of the World Wide Web (WWW) brings about overwhelming information, in addition, it is supposed to be changed dynamically without any prior notification. A large number of information sources are stock type. The users access this type of information in pull mode, mostly by Web Browsers. These information sources have no mechanism to bring information of the changes to prospective users. The users have to deal with the matter by themselves. Browsing through the sites for new updates is not only time

consuming task but also vain in case that there is no change made on the sites once visited. This puts a significant load to the users besides exploring brand new information. We need some representatives to do such burdensome and tedious jobs for us. Furthermore, we would like to know when the changes occurred and how they look. That means not only tracking tools but notification and presentation issues are also taken into account.

This paper considers the evolution of mechanism that detects and evaluates changes on the Web, provides it in comprehensive form, and pushes the information to prospective users. At overall level, these operations induce the flow of information, change and difference instances, from the information sources to the users. With this system, The ubiquitous stock type information sources on the Web have no need to provide any effort to convey their updates to the users. As a result, the information is seemingly transferred in push mode.

We incorporate shared resource management in our system in order to enable the framework a larger scale of service. The shared resource management plays an important role to make the push mode transfer of changes and differences practical. The system would not be practical if the available resources are used to provide service to a large group of users without an effective resource management processing.

The service attempts to increase the utilization factor of the system by several schemes. Each scheme tries to increase identical services in the monitoring process. These identical or virtually identical requests give a significant impact on the utilization of our service. As long as we can implement the service with reasonable resource allocation, more users can have access to our service. At this point, we let the service interact with user. However, the interaction must be made minimal in order to maintain the level of

automatism of the service. We consider the schemes that deal with user as games. The games are played by users and the service agent. The service agent decides its moves based on experience in the past and game analysis. A useful tool we use here is the game theory. We apply the tool from the view point of our service. Overall, the service agent improves the utilization factor of the service and gets the most out of limited interaction with the users.

2 ARCHITECTURE OF THE SERVICE

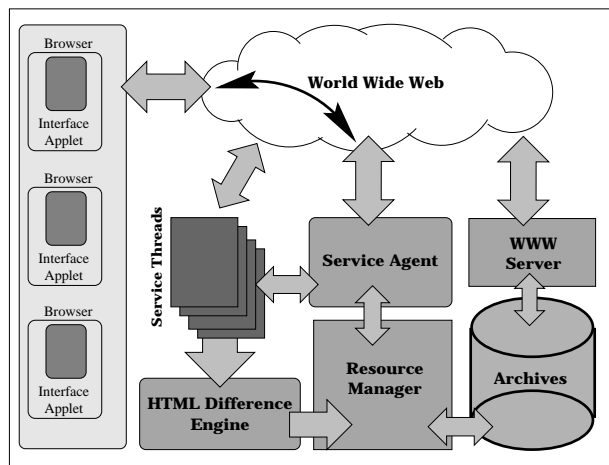


Figure 1: The architecture of change monitoring service.

The architecture of the change monitoring service is shown in [Fig.1]. The service is provided openly on the Web. The user access the service via the WWW by using any browser with Java Virtual Machine. Requests can be made directly to the *Service Agent* which is the front end of the service.

The functions of main modules can be listed as follows:

- **Resource Manager:** All resources for the service are handled by this module. The results from the HTML Difference Engine will be kept in the archives by Resource Manager. It keeps improving the utilization factor of the system by matching all identical or virtually identical requests.
- **Service Threads:** Each Web page monitoring request will be handled by a service thread. The thread keep monitoring and comparing revisions of that Web page.

- **Service Agent:** This is the heart of the service that interacts with other modules in order to retrieve and compare Web pages. It takes requests from the users and consults the Resource Manager to allocate resource for incoming requests. It activates the service threads to start services. When the users query their service profiles, it works as a service broker that retrieve the profiles for the users. The Service Agent handles all responses made by users when the service need to know opinions of the users. The responses from the users are proceeded to the Resource Manager in order to decide how to improve the service.

- **HTML Difference Engine:** The service threads implement the Difference Engine in order to compare the content of updated pages and see whether there are significant changes in them. The old and new versions of HTML documents are compared by running the Difference Engine. The results from Difference Engine are very important for the agent to classify the changes. At the same time, it will summarize the updated information into another HTML document by an innovative algorithm described in the *Difference and Display* subsection.

- **WWW server:** The page archives contain the old and new version of Web pages together with summary pages constructed by the HTML Difference Engine. When the users are notified by the Service Agent, they can view the summary pages with their browsers via the WWW server.

3 MONITORING THE CHANGES

The monitoring service keep monitoring the changes in Web repositories and making comprehensive presentation of the changes available to the users. Once the Service Agent found significant changes, it notifies the users. The considerations for each task is described in detail in following topics.

HTML Document Stream Filtering The result from the differentiator is fed into the HTML constructor and the user's interest based filter. The filtering process is perform right in this filter. The filter implements the three categories of user interest as described above. The existence of contents can be checked by finding whether the old contents are still in the document. At the same time, the filter scans the document whether any new contents are added to the document. In the same manner, the filter checks whether

anything changed with the tags that control the appearance of the document. Once the filter found any changes that fall into these two categories, it evaluates a score of the changes.

$$Score_{content} = \sum_{i=1}^{N_c} w_i(\phi) \quad (1)$$

$$Score_{topic} = \sum_{i=1}^{N_k} \sum_{j=1}^{N_i} 2^{j+5} \quad (2)$$

$$Score_{total} = Score_{content} + Score_{topic} \quad (3)$$

where:

N_c = number of content and appearance changes

N_k = number of key words

N_i = number of occurrences of the i^{th} key word.

ϕ = category of the change

$w_i(\phi)$ = weight of the ϕ category

The topics in user interests can be found by checking the key words that user specified. The filter checks the existence of those key words in the context of difference in the stream then evaluates another score for this kind of user interest by the equation above. The more the occurrences of a key word, the closer to the topic in user interest we can expect. Finally the total score is summed up and then determined by the decision maker. If the score exceeds the specified threshold, the user will be notified of the changes.

4 SHARED RESOURCE MANAGEMENT AND COOPERATION BASED ON COST MODEL

When serving a large number of users, we expect to have some identical or virtually identical requests. These requests can share the resource. The more identical or virtually identical requests, the better utilization factor we can get from the service. The Resource Manager deals dynamically with the request matching. The conditions of virtually identical requests change dynamically upon the the changes in Web repositories and the parameters of the request made by the users.

Unlike the service for individual usage, we attempt to provide personal service while maintaining the efficiency of resource usage. In the case of pushing changes and difference information, we push the information to prospective users. This implies that each user has different degree of interest and attitude against the detected changes. [Fig. 2] shows that among different service pages, there are some identical requests.

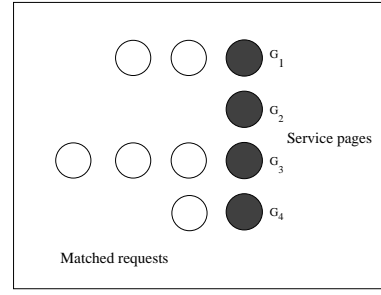


Figure 2: Resource sharing of matched requests.

We can express the utilization of resource as following equations.

$$N = \sum_{i=1}^M G_i \quad (4)$$

$$M = P_{diff} N \quad (5)$$

$$\psi = \frac{N - M}{M} = \frac{1 - P_{diff}}{P_{diff}} \quad (6)$$

where:

N = number of all requests

M = number of different kinds of requests

G_i = number of matched requests for the i^{th} group

P_{diff} = Probability of having different kinds of requests

ψ = Utilization factor

We can see obviously from the equations that if we share the resource among users, we are likely to get more profit than serving each user separately. The utilization factor, finally, depends on the P_{diff} which ranges from $\frac{1}{N}$ to 1. The range tells us that our utilization factor ranges from 0 to $N - 1$.

The amount of identical or virtually identical requests can vary dynamically. This is the case when some requests among currently identical requests are satisfied by the changed conditions but some are not. For examples, we decide to push changes information to the user if we found that the change score is higher than specific threshold points. Suppose we have 2 users who specified the score threshold for an identical page at 1500 and 1000 points respectively. Both requests are considered identical if the change score is 2000 points. However, if the change score falls between 1000 and 1500 points, the requests are no longer identical. The Resource Manager analyzes the characteristic of changes on the Web pages. A parameter to check is the notification threshold. In some cases,

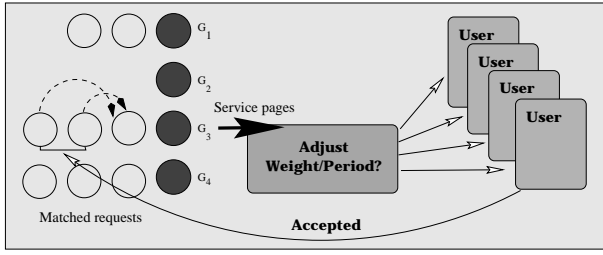


Figure 3: *Increasing matched requests by asking users to modify appropriate parameters.*

some users specified high thresholds with high frequency of monitoring. If the Resource Manager found that the change rates of those pages are relatively slow, it may ask the users to adjust threshold weight of notification or monitoring frequency. Adjusting these parameters has probability to increase more matched requests as shown in [Fig. 3]. Meanwhile, the Re-

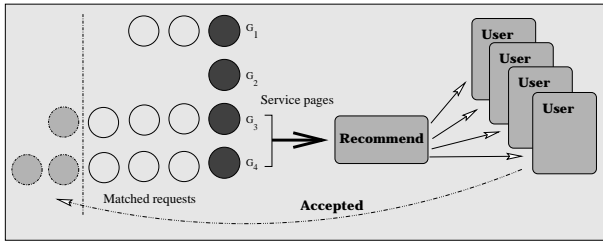


Figure 4: *Increasing matched requests by recommending hot requests to the users.*

source Manager detects hot requests shared by a large number of users. The hot requests trend to be interesting pages. The Resource Manager may recommend these requests to other users as shown in [Fig. 4]. If some users accept the recommendation, the utility factor of the service trends to be increased according to the increasing matched requests.

However, the Resource Manager has to make decision based on facts and experience in the past whether it should ask the users to adjust some parameters or recommend some hot requests to the users. The Resource Manager makes a decision by choosing the most profitable choice from the pay-off matrix of a game. We will consider how to make decision based on our service games as follows:

We define the utility:

Table 1: Expected pay-off matrix for the games of the service. Pay-off $\langle x_1, x_2 \rangle$ indicates that the Service Agent has an expected payoff of x_1 (where an improvement is worth 1, and no improvement is worth 0 for the Service Agent) and the user has an expected payoff of x_2 . In the case of our service $x_1 + x_2$ must be 1.

		User		
		No Adj.	Adj.	
		No Adj.	$\langle p_0, q_0 \rangle$	$\langle p_1, q_1 \rangle$
Service Agent	Adj.	$\langle p_2, q_2 \rangle$	$\langle p_3, q_3 \rangle$	

$$utility(ServiceAgent, 1) \leftarrow improvement$$

$$utility(ServiceAgent, 0) \leftarrow \neg improvement$$

$$utility(User, 1) \leftarrow improvement$$

$$utility(User, 0) \leftarrow \neg improvement$$

Improvement occurs when:

$$improvement \leftarrow ServiceAgent(D) \wedge User(D) \wedge improve_if_follow(D)$$

$$improvement \leftarrow ServiceAgent(Adj.) \wedge User(NoAdj.) \wedge improve_if_sa_un$$

$$improvement \leftarrow ServiceAgent(NoAdj.) \wedge User(Adj.) \wedge improve_if_sn_ua$$

Right here, the *improve_if_sa_un* is the p_2 and the *improve_if_sn_ua* is the p_1 . Suppose that the Service Agent is to choose a strategy with $p_a = P_{ServiceAgent}(Adjust)$ and the user is to choose a strategy with $p_u = P_{User}(Adjust)$. In this setup, the probability of having improvement $P(improve)$ is defined by

$$P(improve) = p_a p_u p_3 + (1 - p_a)(1 - p_u)p_0 + (1 - p_a)p_u p_1 + p_a(1 - p_u)p_2 \quad (7)$$

In a randomized equilibrium from the Service Agent view, the payoff for *Adjust* and *NoAdjustment* must be equal. The payoff for asking for adjustment is the above formula with $p_a = 1$, the payoff for asking no adjustment is the formula with $p_a = 0$. These are equal, so we have

$$(1 - p_u)p_0 + p_u p_1 = p_u p_3 + (1 - p_u)p_2 \quad (8)$$

Similarly for the user

$$(1 - p_a)p_0 + p_a p_2 = p_a p_3 + (1 - p_a)p_1 \quad (9)$$

Solving for p_u and p_a we derive

$$p_u = \frac{p_2 - p_0}{p_1 + p_2 - p_0 - p_3} \quad (10)$$

and

$$p_a = \frac{p_1 - p_0}{p_1 + p_2 - p_0 - p_3} \quad (11)$$

Substitute in equation (7), we derive

$$P(improve) = \frac{p_1 p_2 - p_0 p_3}{p_1 + p_2 - p_0 - p_3} \quad (12)$$

The service agent checks whether the $P(improve)$ is above 0.5 which means the system has probability to improve the service more than 0.5, if it asks the user for adjustment. In our service, the probability p_3 can be derived by a function that evaluates how significant a request for adjustment is. The probability p_2 comes from the experience in the past which is, in other words, how much the user refuse the suggestion. The probability p_1 comes from the improvement made when the user adjust the service profile without suggestion from the Service Agent. Finally, the probability p_0 comes from the self-improvement rate occurred as the conditions of changes in Web repositories vary in time domain.

We can see obviously that the variables used in the game analysis above can be evaluated at ease from the statistic of the service. The Service Agent can make decision to deal or not to deal with the user based on this analysis. Moreover, each decision trends to be more exact as the experience of the Service Agent increases. A cost model is applied to the service mechanism to investigate cooperative behavior of service agents in the community. The invited users under the invitation of preference adjustment requests and the new users gather together to get assigned to appropriate service groups. This implies that they desire to play in an assignment game in order to improve profits. Right in this process, the requests will be clustered into groups of prospectively identical requests. These groups are applied to a cost model for finding acceptable matches. The cost for a service is defined on the concept that the more effective identical requests, the lower the cost of each service in the identical group.

We define C_i , the cost per service for each user in group i^{th} base on F , the full cost per user, as follows:

$$C_i = F(1 - k \frac{e^{\frac{g_i}{N}} - 1}{e - 1}) \quad where \ 0 \leq k \leq 1 \quad (13)$$

The cost reduction increases exponentially according to the number of identical requests in the group. This create a persuasion for the users to join a large group even the request for that group is not exactly match what they want since this can reduce their costs. Another effect of the cost model is that the service costs in the community become lower remarkably as the number of identical requests increases. As a result, the competitiveness of the community becomes higher in case that we consider the service in multiple communities. The above is the evaluation of value of the services on the service provider side. The evaluation of services on the user side can be derived from the evaluation of distance between existing services and the service in need. We define the appreciation value of the user j^{th} against the service of the group i^{th} by u_{ij} . Then, the difference value for the service of group i^{th} viewed by the user j^{th} is $a_{ij} = u_{ij} - C_i$. If we have a set of service groups $G = \{g_1, g_2, \dots, g_n\}$, and a set of users $H = \{h_1, h_2, \dots, h_m\}$, then we can find the maximum profit combination by finding the maximum value of $\sum a_{ij}$ among the combination space of G and H .

5 IMPLEMENTATION

The prototype of the system has been implemented locally in our laboratory. The users can access from anywhere on the Web, provided that they can access our interfacing JAVA applet. Currently, we are serving up to 20 users and monitoring over 200 pages. The followings are some of our results. Once the users are

```

Subject: Page Update Information
Date: Wed Jan 27 03:36:48 GMT+09:00 1999
From: WebBeholder@shizuoka.ac.jp
To: santi@miv.t.u-tokyo.ac.jp

Report for http://www.miv.t.u-tokyo.ac.jp/~santi:

3 changes for page link
  http://www.cs.colorado.edu/home/mcbryan/WWWWWW.html was added
  http://ycos.cs.cmu.edu/ was added
  http://home.netscape.com/home/internet-search.html was added.
87 characters changed
The summary can be found at
http://onvx.miv.t.u-tokyo.ac.jp:8081/webbeholder/~santi@miv.t.u-tokyo.ac.jp/http://onvx.miv.t.u-tokyo.ac.jp:8081/ht
Total changed weight: 811

```

Figure 5: Changes information pushed from a Service Thread via Email.

notified via Email, they can easily access the summary pages on the Service Agent side by clicking on the link

11:00 – 11:45	Hand Writing Detection and Optical Character Recognition: Japanese as the Leading Developer Dr. Kittu Kosavisutte
11:45 – 12:00 12:00 – 13:00	Lunch
Setting up the Computers for a Multilingual Capacity: Desktop Processing and Telecommunications	
13:15 – 14:15	Demonstration on a UNIX machine W/Out-See-Tune Yamada Jun

Figure 6: A look of presentation of changes.

attached to the notification mails. The HTML Difference Engine is able to indicate the changes precisely down to the cell level of table structure in the document. The deleted text is displayed in a stroke font and the new text is displayed in an underline font. The deleted links or images will be marked by small foot icons and the inserted links or images will be marked by peg icons. The users can jump from one change to another quickly since the changes are displayed noticeably.

6 Summary

In this paper, we presented the mechanism of Web repositories change monitoring service that notifies the users about changes in Web repositories and creates comprehensive summary pages of the changes. The improvement of overall utilization factor is derived from the resource management within each service unit and the cooperation among service units. The Matchmaking Agent has potential to balance the services among service units. Decision making process of both Matchmaking Agent and Service Agent in each service unit relies on the game analysis. The expected payoff values based on experience in the past have direct impacts to the decision. The users that are engaged to the service assignment game are grouped to appropriate service group based on cost model. At the same time, the mechanism serves well as a source of extracting trends and users' interest.

References

[1] Santi Saeyor and Mitsuru Ishizuka: WebBeholder: A Revolution in Tracking and Viewing Changes on The Web by Agent Community, in *proceedings of WebNet98, 3rd World Conference on WWW and Internet*, Orlando, Florida, USA, Nov. 1998.

- [2] Fred Douglass, Thomas Ball, Yih-Farn Chen and Eleftherios Koutsofios: The AT&T Internet Difference Engine: Tracking and viewing changes on the web *World Wide Web* Volume 1 Issue 1, 1998. pp. 27-44
- [3] Fred Douglass: Experiences with the AT&T Internet Difference Engine 22nd International Conference for the Resource Management & Performance Evaluation of Enterprise Computing System (CMG96), December, 1996.
- [4] F. Douglass, T. Ball, Y. Chen, E. Koutsofios. Webguide: Querying and Navigating Changes in Web Repositories. In *Proceedings of the Fifth International World Wide Web Conference*, Paris, France, May 1996. pp. 1335-1344.
- [5] Brian Starr, Mark S. Ackerman, Michael Pazzani: Do-I-Care: A Collaborative Web Agent *Proceeding of ACM CHI'96*, April, 1996.
- [6] Imma Curiel: Cooperative Game Theory and Applications, Kluwer Academic Publishers, 1997
- [7] Theo Driessen: Cooperative Games, Solution and Applications, Kluwer Academic Publishers, 1988