

International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirset.com

Vol. 6, Issue 6, June 2017

Effective Page Refresh Policies for Green Web Crawling

Pallavi Ghodke, Prof. P.S.Desai

ME Student, Dept. of Computer Network, SKNCOE Pune., Savitribai Phule Pune University Pune India

Professor, Dept. of Computer Engineering, SKNCOE Pune. Savitribai Phule Pune University Pune India

ABSTRACT: On web large no. of requests goes to server. Due to large no. of HTTP requests to web servers increase the energy consumption and carbon footprint of the web servers and for that computational resources are used while serving the requests. In existing system the problem of Web crawlers that maintain local copies of remote Web pages for Web search engines. In this context, remote data sources (Websites) do not notify the copies (Web crawlers) of new changes, so we need to periodically poll the sources to maintain the copies up-to-date. The proposed work is mainly motivated by the need to manage updated web data. In this Web base system, it stores a significant portion of the web in local repository for web searching. This web search engines also maintain copies and/or indexes of Web page, and they need to periodically visit the pages to maintain them up-to-date. Here page to be refreshed is selected based on a metric that considers the page's staleness, its size, and the greenness of the energy consumed at the web server premises. If page is updated then increase its greenness and not updated then increase its staleness. It avoids no. of frequently sending http request to server. Personalized profile is maintained for each user.

KEYWORDS: Crawling, carbon footprint, greenness, staleness, Web Dynamics, Web repository.

I. INTRODUCTION

Web crawling is responsible for traversing the hyperlink structure among the web pages to discover and download the content in the Web, as well as for refreshing already downloaded pages in the web repository [2]. Frequently occurrence of HTTP request increases carbon emission. To avoid large no. of HTTP request one needs to maintain local copies of remote data sources for better performance or availability. In this paper Web search engine copies a significant subset of the Web and maintain copies and/or indexes of the pages to help users' access relevant information. Part of the local copy may get out-of-date because changes at the sources are not immediately propagated to the local copy [9]. Therefore, it becomes important to design a good refresh policy that maximizes the "freshness" of the local copy [3]. On web updated pages are not displaying top and request goes to only one server frequently so its energy consumption get increased. Here we consider staleness and greenness of web pages. Updated contain will reflect immediately by checking pages last updated time and changing it in local copy and displaying them to user. Here define policy depending on server greenness and pages staleness.

1. Background

Web crawling is responsible for traversing the hyperlink structure among the web pages to discover and download the content in the Web, as well as for refreshing already downloaded pages in the web repository [2]. Frequently occurrence of HTTP request increases carbon emission. To avoid large no. of HTTP request one needs to maintain local copies of remote data sources for better performance or availability. In this paper Web search engine copies a significant subset of the Web and maintain copies and/or indexes of the pages to help users' access relevant information. Part of the local copy may get out-of-date because changes at the sources are not immediately propagated to the local copy [9]. Therefore, it becomes important to design a good refresh policy that maximizes the "freshness" of the local copy [3]. On web updated pages are not displaying top and request goes to only one server frequently so its energy consumption get increased. Here we consider staleness and greenness of web pages. Updated contain will

International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirset.com

Vol. 6, Issue 6, June 2017

reflect immediately by checking pages last updated time and changing it in local copy and displaying them to user. Here define policy depending on server greenness and pages staleness.

2. Motivation

We devise a web repository refreshing technique that takes into accounts both the greenness and staleness concepts when scheduling the download of web pages. This technique aims to reduce the total staleness of pages in the web repository while constraining web servers' total carbon footprint resulting from the activities of the crawler. We evaluate the performance of our technique using a large, real-life web server data set. Beyond creating an environment-friendly crawler, our work has implications for large-scale web search engines, which should comply with regulations about carbon footprint reduction.

Purpose:

- 1) To minimize http request to web server.
- 2) Getting user expected result.
- 3) Minimize carbon foot print.
- 4) Works use the update history of related pages to capture the actual update likelihood of a page better. Works in another line devise page refresh strategies that aim to minimize the negative impact on users due to stale search results

II. LITERATURE SURVEY

4) This paper proposed a general, optimization-based framework to minimize datacenter costs in the presence of different carbon footprint reduction goals, renewable energy characteristics, policies, utility tariff, and energy storage devices (ESDs). Advantage is on-site renewable can help lower costs due to their ability to reduce the peak datacenter power draw from the utility and disadvantage is On-site and off-site these hybrid combination is the most cost-effective across the spectrum [9].

5) Proposed two joint dynamic speed scaling and traffic shifting schemes, one sub gradient-based and the other queue-based. Our schemes assume little statistical information of the system, which is usually difficult to obtain in practice. Its advantage is a joint dynamic speed scaling and traffic shifting scheme. But the developed QTF leads to a higher cost [2].

6) This Proposed system decides how to schedule Web pages for selective (re)downloading into a search engine repository. The scheduling objective is to maximize the quality of the user experience for those who query the search engine. Advantage is it maximizes user experience. But it is effectively impossible to compute the exact quality value of a large repository of Webpages [10].

7) Proposed system identified challenges through an architectural classification, starting from single node search system and moving towards hypothetical multi-site web search architecture. Its main advantage is it present the current scalability challenges in large-scale web search engines. Drawback is lack of large hardware infrastructures and datasets make conducting scalability research quite difficult [1].

2) This paper solve the problem related to conflicting requirements of maximizing quality of services delivered by the cloud services while minimizing energy consumption of the data center resources. Proposed the concept of inception of data center energy-efficiency controller that can consolidate data center re-sources with minimal effect on QoS requirements. Its main Advantage is present inception of data center energy-efficiency controller that can consolidate data center resources with minimal effect on QOS requirements. Its drawback is energy-efficient hardware devices such as SSDs and optical Interconnects do not fit into current cloud paradigm due to their high cost. [5]

9) Proposed various refresh policies and studies their effectiveness. We first formalize the notion of "freshness" of copied data by defining two freshness metrics, and we propose a Poisson process as the change model of data sources. Based on this framework, we examine the effectiveness of the proposed refresh policies analytically and experimentally. Advantage is proposed refresh policies improve the "freshness" of data very significantly. After drawback is it increases local data storage [4].

1) This paper proposed a system that makes efficient use of communication resources in approximate replication environments. Advantage is when data changes rapidly, good performance. Disadvantage is that It doesn't ensure exact consistency [7].

International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirset.com

Vol. 6, Issue 6, June 2017

8) This proposed system focuses the role of communication fabric in data center energy consumption and presents a scheduling approach that combines energy efficiency and network awareness, named DENS. The DENS methodology balances the energy consumption of a data center, individual job performance, and traffic demands. Advantage: Balances the energy consumption of a data centre, individual job Performance and traffic demands. Disadvantage: produce heavy data streams directed to the end-users [3].

3) It proposed dynamic placement of virtual machines (VMs) with deterministic and stochastic demands. In order to ensure a quick response to VM requests and improve the energy efficiency, a two-phase optimization strategy has been proposed, in which VMs are deployed in runtime and consolidated into servers periodically. Its advantage is that it gives quick response to VM requests, but Deadlock may occur in the actual migration process [6].

10) Propose methods for selecting related pages, and have shown that related pages that correlate over time in their content change contain valuable information for page content change prediction. Using page content as well as related pages significantly improves prediction accuracy and compares it to common approaches. It presents a drawback large scale prediction tasks is that storing snapshots from multiple iterations is prohibitively expensive [8].

III. EXISTING SYSTEM APPROACH

On web user query for his requirement. It is hard to find for crawler to notice changes newly made to web Pages. It takes more time. The New updated pages were not found on web. A web crawler can do to reduce the energy consumption it incurs to web servers without sacrificing the coverage or freshness of its web repository. This is because the amount of energy consumed on web servers depends only on factors related to the hardware and software resources that are not managed by the search engine company. Previous system schedules a Web crawler to improve freshness. The model for the Web crawler and freshness are very different. It is important to note “importance” or the “weight” of a page is proportional to the change frequency of the page[3]. While this assumption makes analysis simple, but it lead hard to discover the fundamental trends. On web digital data becomes available rapidly, it will be increasingly important to collect it effectively. A crawler cannot refresh all its data constantly, so it must be very careful in deciding what data to poll and check for freshness.

IV. PROPOSED SYSTEM APPROACH

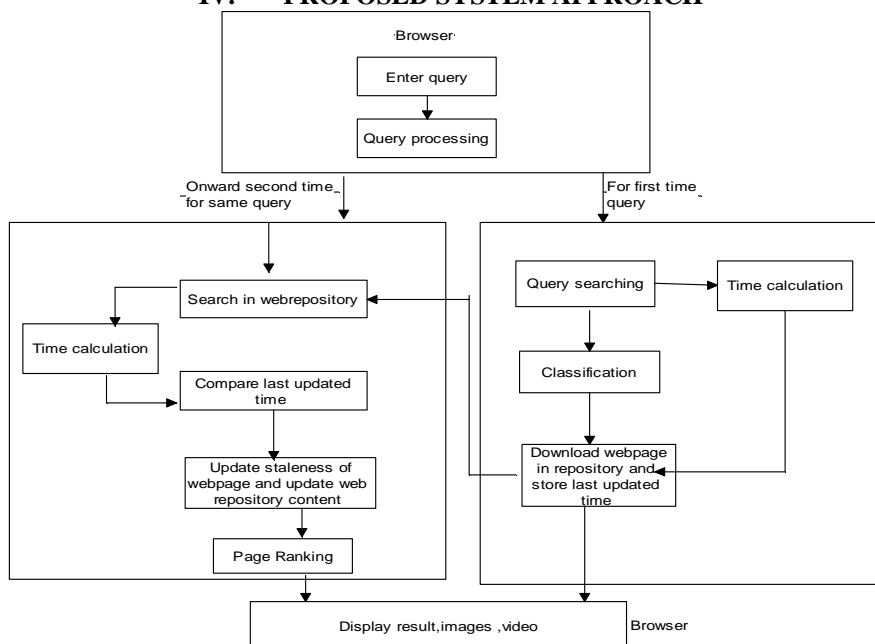


Fig. No. 01) System Architecture of Green crawler

International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirset.com

Vol. 6, Issue 6, June 2017

EXPLANATION

Web searching:

1) Seed collection:

2) Document classification: Web pages are classified for that for which website these query is entered and document are classify.

3)Greenness and staleness updating: When downloading a document system checks last updated time of that page if updated time is latest then get that data and store in own repository. pageto get updated content GreedyIL algorithm ids used and store in web repository's page and display it to user. If page is latest updated then increase its greenness and if not then its staleness will increased. In that way policy is performed for each request.

4)Page Ranking: To reduce large no. of request for particular one server links are displayed to user according to user hit count. If hit count is exceed some value then again it reduces the hit count and display it to use at top. That helps to give chance to equal ranking.

ADVANTAGES

1) It reduces no. of HTTP request to servers.

2)This technique aims to reduce the total staleness of pages in the web repository while constraining web servers' total carbon footprint resulting from the activities of the crawler.

3) It minimizes total staleness of a web repository is related to the problem of minimizing the total waiting time.

4) Personalization is maintained for each user.

V. MATHEMATICAL APPROACH

NOTATION-

1) $j(t)$: Staleness of page.

2) t : Time to take decision.

3) j : Page to download

3) W :set of web pages.

4) $x_j(t)$: Time to download web page.

EQUATION-

$$1) s_j(t+1) = \begin{cases} 0 & \text{if } x_j(t) = 1 \\ s_j(t) + 1 & \text{if } x_j(t) = 0. \end{cases} \quad (1)$$

$$2) \min_{\sum_{t=0}^T \sum_{j \in W} s_j(t)} \quad (2)$$

$$3) s.t \sum_{j \in W} x_j(t) \forall t = 0, \dots, T-1 \quad (3)$$

$$4) \sum_{t=0}^{T-1} \sum_{j \in W} x_j(t) e_j \leq G \quad (4)$$

$$5) \sum_{j \in W} s_j(t+1) = \sum_{j \in W} s_j(t) + (|W| - 1) s_{j^*(t)}(t) \quad (5)$$

VI. ALGORITHM-GREEDYIL

Input: seed Documents

Output: updated document

Step 1: Let priority queue Q of seed Document, I be the iteration index to 0 Create Q for seedDocument.

Step 2: if no. of seed Document < no. of offline document

Step 3: update Rank

Step 4: else

Step 5: parseSimilarDocument

International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirset.com

Vol. 6, Issue 6, June 2017

Step 6: while risk(q,documents)>threshold

Step 7: do set s ← parameter(t,G)

Step 8: if t has no sibling then

Step 9: Insert(S,IL(s)) to Q

Step 10: else if t has sibling then

Step 11: merge t into shadow sibling

Step 12: if no operation on sibling in Q then

Step 13: Increase staleness

Step 14: else

Step 15: increase greenness

Step 16: Update rank

Step 17 : return updatedDocument

VII. MODULE

a. Web searching:

In this system user enters the query to search result. According to query stop words are removed and remaining word are passed to search web pages to crawler. If query is entered same time which is previously entered for that system search web pages in his own repository.

b. Document classification: Web pages are classified for that for which website these query is entered and document are classify.

Document download: When web pages are available according to query they are downloaded in web repository.

c. Greenness and staleness updating: When downloading a document system checks last updated time of that page if updated time is latest than already stored page then parse the latest updated page and take only updated data of that page and store that updated data in already downloaded page. To check page content similarity it perform greedily algorithm and take only updated content and store in web repository's page and display it to user. If page is latest updated then increase its greenness and if not then its staleness will increased. In that way policy is performed for each request. Result will display to user which pages have better greenness.

VIII. EXPERIMENTAL SET UP

Result table I:

time	EDD Like policy	Optimal policy
10	0.7	0.7
20	0.7	0.8
30	0.7	0.9
40	0.7	1
50	0.7	1.1

Table 1: The performance of the optimal policy in comparison with that of the EDD-like policy in terms of total staleness as a function of λ .

Result table 3:

	Overall		com		Gov.	
	Fressness	Age	Fressness	Age	Fressness	Age
1	0.12	400days	0.7	386days	0.69	19days
2	0.57	5days	0.38	8days	0.82	2days
3	0.62	4days	0.44	7days	0.85	1.3days

International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirset.com

Vol. 6, Issue 6, June 2017

Table: Freshness

- 1) Performance measure used:
Graph:

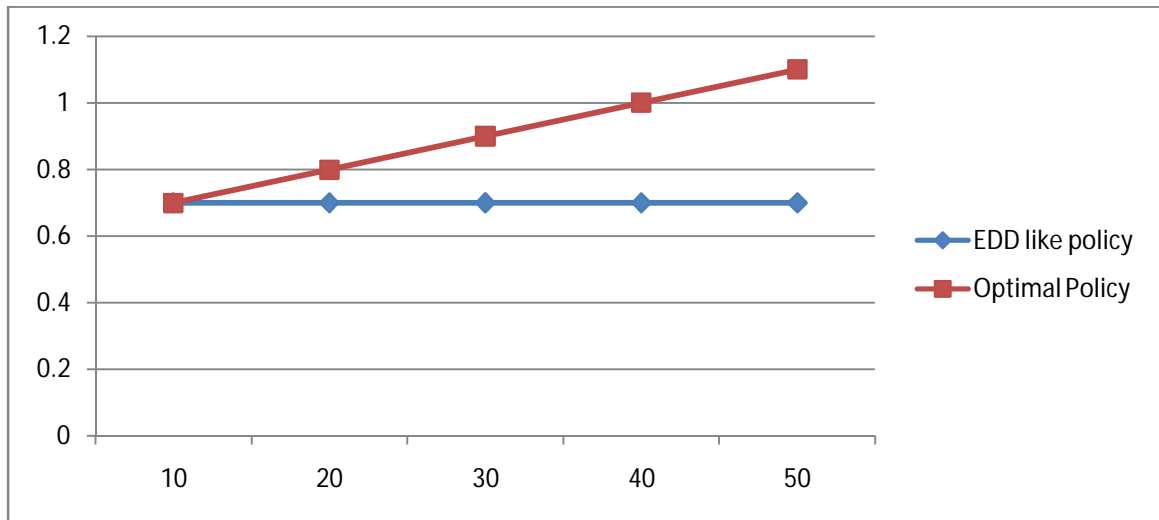


Fig 2: x-axis λ (min/g) Y-Axis Total staleness(min)

The performance of the optimal policy in comparison with that of the EDD-like policy in terms of total staleness as a function of λ . As λ increases, the servers' desire to keep their carbon footprint low increases as well and thus, the crawler is obstructed from minimizing the Staleness of pages.

Proposed system design an optimal policy, which can be implemented in an online manner, based only on the instantaneous values of page staleness and greenness indicators of the servers. Here we stored page content in our repository so that helps to assign staleness to pages. Here we are minimizing sending http request to server.

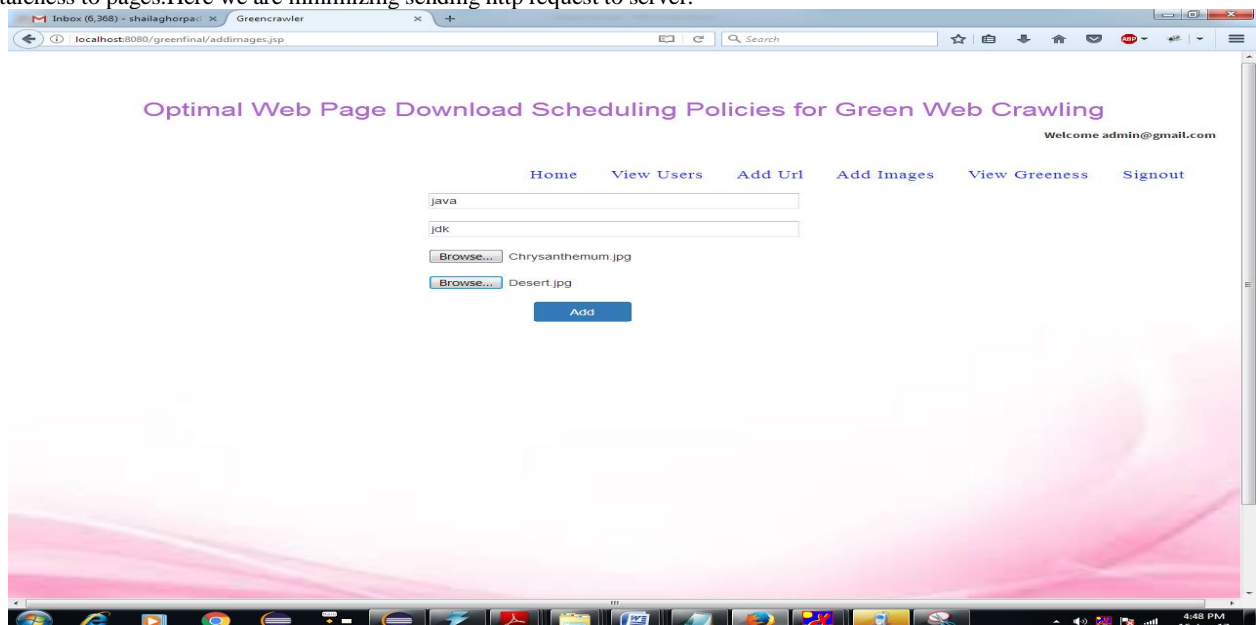


Fig 3 Greenness.jsp:

International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirset.com

Vol. 6, Issue 6, June 2017

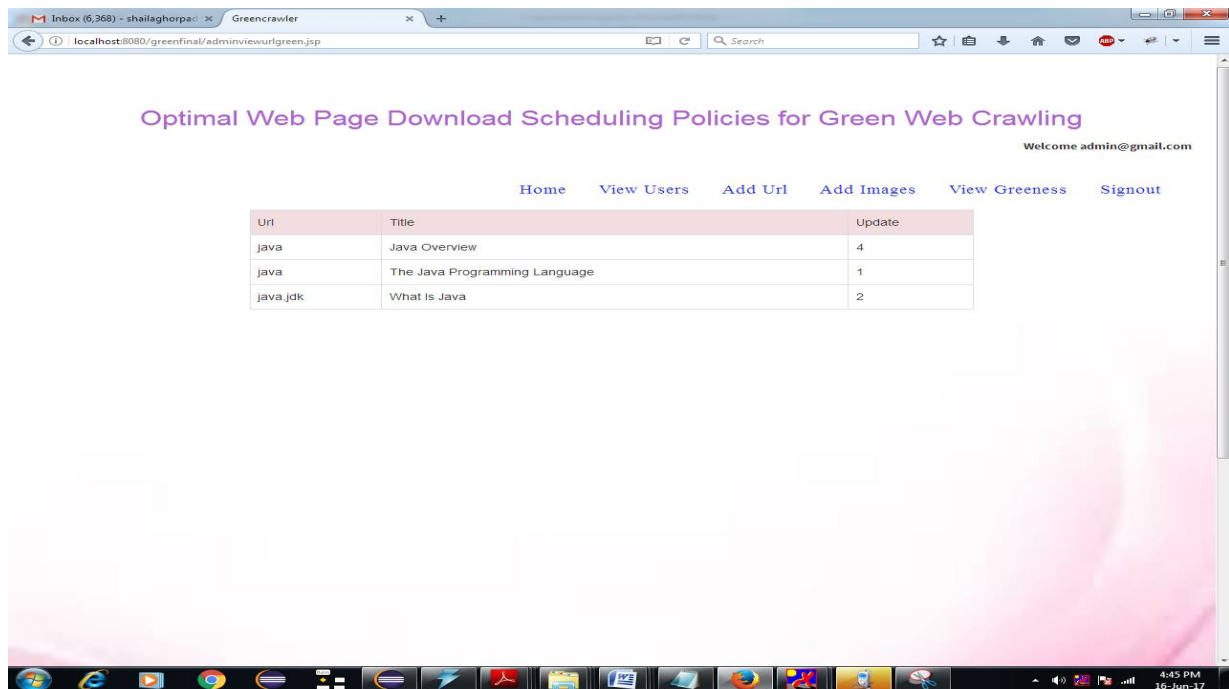


Fig 4 Greenness of query word:

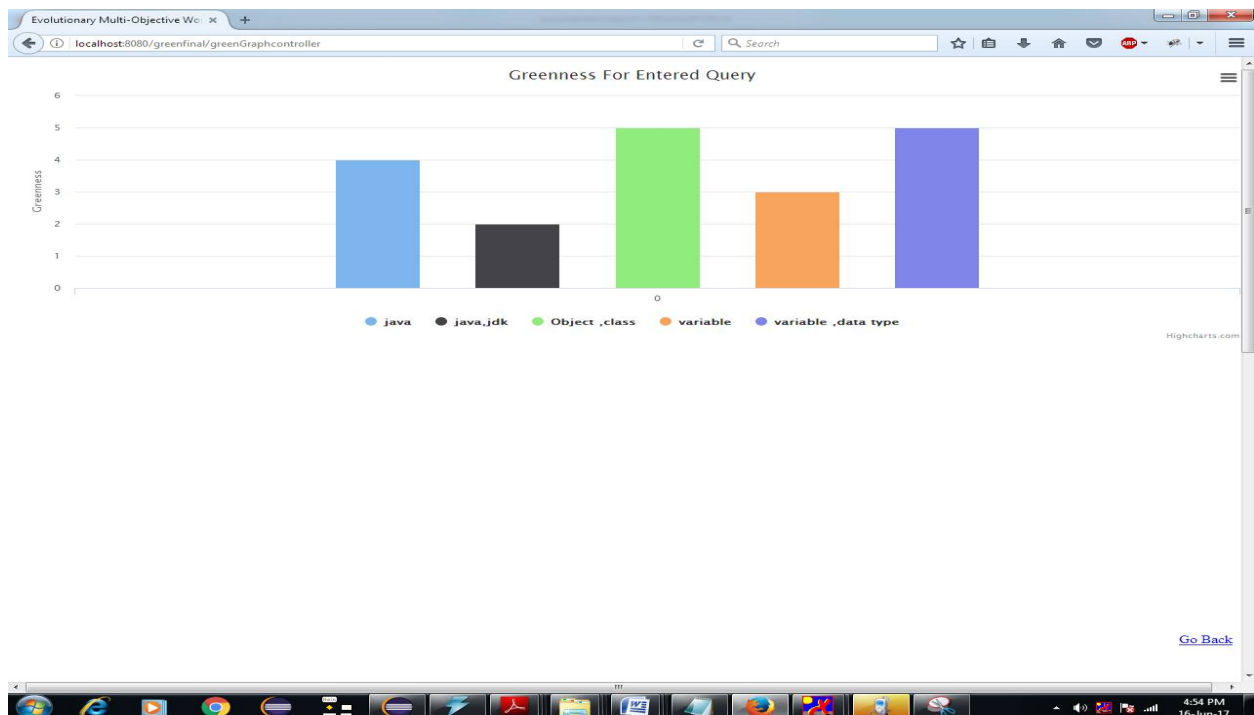


Fig 5 Greenness for Entered Query.

International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirset.com

Vol. 6, Issue 6, June 2017

IX. CONCLUSION

Proposed system design an optimal policy, which can be implemented in an online manner, based only on the instantaneous values of page staleness and greenness indicators of the servers. It works on Stored page content in repository so that helps to assign staleness topages. Here it minimizing sending http request to server. . It avoids no. of frequently sending http request to server. Personalized profile is maintained for each user.

REFERENCES

- [1] Scalability challenges in web search engines, in Synthesis Lectures on Information Concepts, Retrieval, and Services. San Mateo, CA, USA: Morgan, 2015, B. B. Cambazoglu and R. A. Baeza-Yates,
- [2] Geographic trough filling for internet datacenters, in Proc. INFOCOM, 2012, pp. 2881–2885, D. Xu and X. Liu,
- [3] DENS: Data center energy efficient network-aware scheduling, Cluster Comput., vol. 16, no. 1, pp. 65–75, Mar. 2013, D. Kliazovich, P. Bouvry, and S. U. Khan.
- [4] Effective page refresh policies for web crawlers, J. Cho and H. Garcia-Molina ACM Trans. Database Syst., vol. 28, no. 4, pp. 390–426, Dec. 2003.
- [5] Survey of techniques and architectures for designing energy-efficient data centers, IEEE Syst. J., pp. 1–13, Jul. 2014, J. Shuja et al.
- [6] Dynamic placement of virtual machines with both deterministic and stochastic demands for green cloud computing,” Math. Prob. Eng., vol. 2014, 11 pp., Jul. 2014, W. Yue and Q. Chen
- [7] Efficient Monitoring and Querying of Distributed, Dynamic Data via Approximate Replication, 2005 IEEE Christopher Olston, Jennifer Widom.
- [8] Predicting Content Change on the Web, 2013 ACM 978-1-4503-1869-3/13/02, Kira Radinsky, Paul N. Bennett
- [9] Carbon-Aware Energy Capacity Planning for Datacenters, 2012, Chuangang Ren, Di Wang, Bhuvan Uргаonkar, and Anand Sivasubramaniam.
- [10] User-centric web crawling, in Proc. 14th World Wide Web, 2005, pp. 401–411, S. Pandey and C. Olston,