

# Improving Performance of Internet Services Through Reward-Driven Request Prioritization

Alexander Totok and Vijay Karamcheti  
Department of Computer Science  
Courant Institute of Mathematical Sciences  
New York University, New York, NY, USA  
Email: {totok,vijay}@cs.nyu.edu

**Abstract**—Meeting client QoS expectations proves to be a difficult task for the providers of modern Internet services, especially when web servers experience overload conditions (steady or transient), which cause increased response times and request rejections, leading to user frustration, lowered usage of the service and reduced revenues. In this paper, we propose a server-side request scheduling mechanism that addresses these problems. Our *Reward-Driven Request Prioritization (RDRP)* algorithm gives higher execution priority to client web sessions that are likely to bring more service *profit* (or any other application-specific *reward*). The method works by predicting future session structure by comparing its requests seen so far with aggregated information about recent client behavior, and using these predictions to preferentially allocate web server resources. Our experiments using the TPC-W benchmark application and CBMG-based web workloads, with an implementation of the RDRP techniques in the JBoss application server, show that RDRP can significantly boost reward attained by the service, while providing better QoS to clients that bring more reward.

## I. INTRODUCTION

In recent years, the role of the Internet has undergone a transition from simply being a data repository to one providing access to a variety of sophisticated Internet services such as e-mail, banking, on-line shopping, and entertainment. Typical interaction of users with such services is organized into *sessions*, a sequence of related requests, which together achieve a higher level user goal. An example of such interaction is an on-line shopping scenario for an e-Commerce web site, which involves multiple requests that search for particular products, retrieve information about a specific item, add it to the shopping cart, initiate the check-out process, and finally commit the order. The *success of the whole session* now becomes the ultimate QoS goal, which contrasts with the *per-request success* performance metrics of the early Internet.

This shift in QoS metrics is beginning to manifest itself in how an Internet service schedules the execution of client requests. Traditionally, a web server hosting an Internet service would process user requests on a first-come-first-served (FIFO) basis. Although this approach provides fair access to the service for all clients, it does not work as well if the service operates under overload conditions, whether such conditions are steady or transient (as a result of bursty client behavior). In such situations, clients see increased response times and their requests (and the containing sessions) may get rejected, which

leads to user frustration, and as a consequence, to lowered usage of the service and reduced service revenues.

Consequently, modern-day services may contain one or more server-side mechanisms to deal with such overload situations. Session-based admission control (SBAC) [1] admits only as many sessions as can be served by the service. More complex service differentiation mechanisms have also been used to provide stable QoS guarantees (e.g., request throughput, response times) to different client groups, based on prenegotiated Service-Level Agreements (SLAs). Common to such schemes is the consideration that QoS received by a client is determined upfront by his association with a client group or by his service membership status.

Although they offer better performance than FIFO scheduling, the above schemes fall short of delivering the best performance in many important scenarios. In particular, a service provider often encounters situations where it makes sense to differentiate among clients based on the (dynamic) activities these clients perform in a session, rather than on their (static) identity, in order to boost service revenues, or for other application-specific goals. Let's consider the following three examples.

- In the online shopping scenario introduced earlier, the service provider might be interested in giving a higher execution priority to the sessions that have placed something in the shopping cart (potential *buyer* sessions), as compared to the sessions that just *browse* product catalogs, making sure that clients that buy something (and so – bring profit to the service) receive better QoS.
- For Internet services, some of whose web pages contain third party-sponsored advertisements, the service provider's profits may (directly) depend on the number of visits to these pages. Consequently, the service provider may wish to provide better QoS to the sessions that visit these pages more often.
- For many services whose client interactions involve different length sessions, service profit may be defined in terms of the number of sessions that visit a distinguished "success" page. In such cases, the service provider may prefer shorter sessions visiting the success page over longer ones, because more of them could be served.

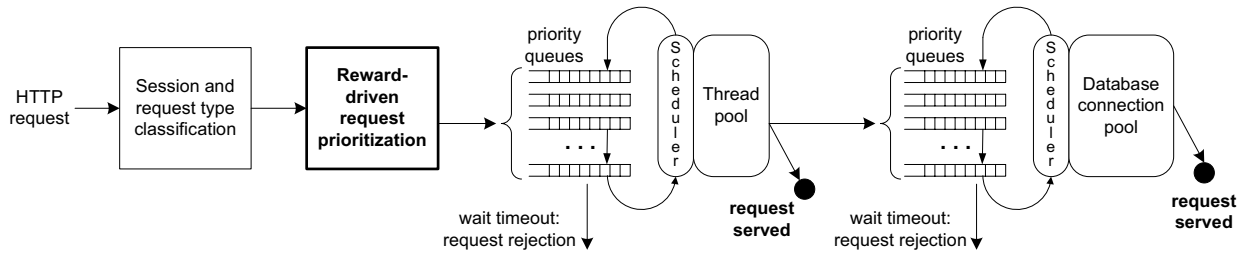


Fig. 1. The model of web application server architecture.

These examples are unified by the idea that the service may benefit from providing better QoS to sessions, that bring more *profit* (give more *reward*), where the notion of profit or reward is defined in an application-specific fashion. What is important is that the information about the client’s possible usage of a service (and its associated contribution to service reward) is not encoded in any static profile, so application-logic-independent SLA-based service differentiation approaches are not as beneficial here.

Instead, to be able to provide better QoS to the sessions that bring more reward, the service provider now needs to predict the behavior of a client. If the client is a returning customer and his identity can be determined (e.g., using cookies), then decisions on QoS provided to this client can be based on the history of his service usage (e.g., history of previous purchases). However, the success of this *per-client history-based approach*, is, not unexpectedly, highly dependent on the correlation between the past and the future behavior of a client, and may not work well if such a correlation is absent or weak.

Instead of focusing on individual client behavior, we advocate the approach of predicting a session’s activities by associating it with aggregated client behavior or broader service usage patterns, obtained for example through online request profiling. Specifically, we propose *Reward-Driven Request Prioritization (RDRP)* mechanisms that try to maximize reward attained by the service, by dynamically assigning higher execution priority values to the requests whose sessions are likely to bring more reward. Our methods compare the sequence of a session’s requests seen so far with aggregated information about client behaviors, and use a Bayesian inference analysis to statistically predict the future structure of a session, and so – the reward the session will bring and the execution cost it will incur. The predicted reward and execution cost values are used to compute each request’s priority, which is used in scheduling “bottleneck” server resources, such as server threads and database connections, to incoming client requests.

We have implemented our proposed methods as a set of middleware mechanisms, which are seamlessly and modularly integrated in the open-source Java web application server JBoss [2]. A *Request Profiling* module performs automatic real-time monitoring of client requests to extract parameters of service usage and to maintain the histories of session requests. It also performs fine-grained request profiling to

identify execution times for different service request types. The *RDRP module* uses the information gathered by the Request Profiling module to compute and assign request priorities, that in turn influences queueing behavior for various application server resources.

We evaluate our approach on the TPC-W benchmark application [3], emulating an online store selling books, and compare it with both the session-based admission control and per-client history-based approaches. Our experiments show that RDRP techniques yield benefits in both underload and overload situations, for both smooth and bursty client behavior. In underload situations, the proposed mechanisms give better response times for the clients that bring more reward, which is crucial for ensuring return customers.<sup>1</sup> In overload situations, when some of the requests get rejected, the mechanisms ensure that sessions that bring more reward are more likely to complete successfully and that the aggregate profit attained by the service increases compared to other solutions. Additionally, we show that the history-based approach matches the performance of our RDRP mechanisms on the amount of reward attained and response times only if the correlation between the clients’ past and future behavior is 75% or greater, and 50% or greater respectively.

The rest of this paper is organized as follows. Section II presents models and assumptions used throughout the paper. Section III describes the reward-driven request prioritization techniques. Section IV presents our testing methodology and experimental results. In Section V we discuss related and future work, and we conclude in Section VI.

## II. MODELS AND ASSUMPTIONS

### A. Web application server architecture

We work with Internet services implemented on top of modern middleware platforms, such as the J2EE component framework [6]. Such services are usually built as complex (and often distributed) software systems, consisting of several logical and physical tiers (e.g., web tier, application tier, and database tier) and accessing multiple backend data sources. We present our request prioritization algorithms in a centralized setting however, to focus on the benefits of the proposed request prioritization techniques. We expect that the methods

<sup>1</sup>Several independent studies showed that the bulk of a service customers are returning clients, and that providing good QoS to long-time customers is a key factor in service success [4], [5].

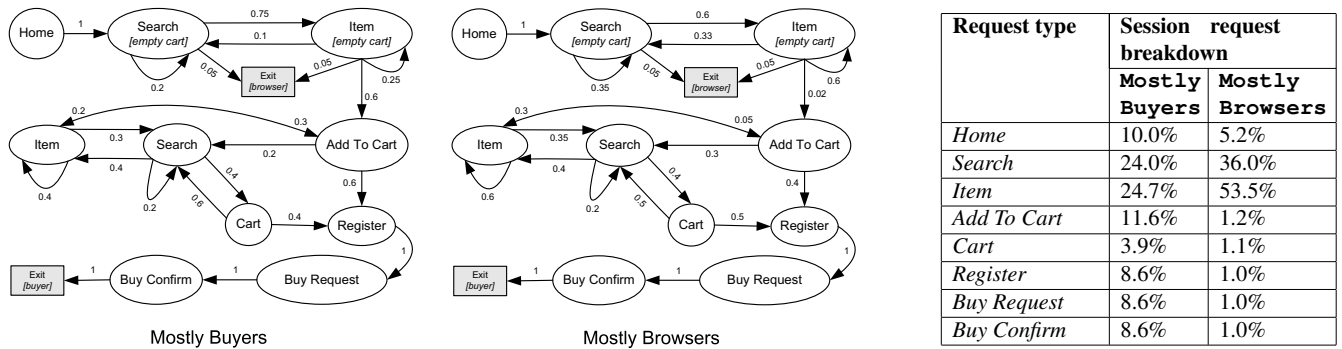


Fig. 2. Two CBMGs used for TPC-W web workload: “Mostly Buyers” (left) and “Mostly Browsers” (middle). Average breakdown of sessions by request types (right).

will show their utility in a distributed setting as well, where they can be independently applied at every system resource contention point that sees concurrent requests competing for server resources.

Our work uses middleware-level mechanisms for server performance optimization, specifically control over request scheduling policies. We adopt this approach, because the middleware itself often does not have control over the low-level OS resources (e.g., CPU and memory), and uses higher-level mechanisms, such as request scheduling, component pool management, transaction demarcation, etc., to improve server performance. It is often the case that middleware server performance is limited by several “bottleneck” resources, that are *held exclusively* by a request for the whole duration or some significant portion of it (as opposed to low-level *shared* OS resources), such as server threads, or database (DB) connections. The default allocation policy of these resources to requests is FIFO. In the absence of application errors, failing to obtain such a resource is the major source of request rejection.

We advocate and use a request execution model, where a request is rejected (with an explicit message) if it fails to obtain a critical server resource within a specified time interval. This approach is shared by a vast majority of robust server architectures that bound request processing time in various ways (e.g., by setting a deadline for request completion), as opposed to a less robust approach, where a request is kept in the system indefinitely, until it is served (or is rejected by lower-level mechanisms such as TCP timeout). The former approach has the advantage of more efficiently freeing up server resources held by requests whose processing cannot be completed because of server capacity limitations.

Fig. 1 illustrates this application server architecture and the flow of a request through the system. Requests compete for two critical exclusively-held server resources: server threads and database (DB) connections; these resources are pooled by the web server and the application server respectively. Scheduling of requests to available threads and DB connections is done according to the request priority set by the RDRP module. The request with the highest priority is served first, with FIFO used as a tiebreaking policy. Timeout values for obtaining a thread and a DB connection are set to be 10s. If

this timeout expires, the request is rejected with an explicit message. Note that some requests do not require database access, so they can be successfully served just by acquiring a server thread.

### B. TPC-W Application

To test the benefits of the proposed request prioritization mechanisms, we use the TPC-W transactional web e-Commerce benchmark [3], that emulates an on-line store that sells books. The TPC-W specification describes in detail the application data structure and the 14 web invocations (WI) that constitute the web site functionality, and defines how they change the application data stored in the database. A typical TPC-W session consists of the following requests: a user starts web site navigation by accessing the Home page, searches for particular products (Search), retrieves information about specific items (Item Details), adds some of them to the shopping cart (Add To Cart), initiates the check-out process, registering and logging in as necessary (Register, Buy Request), and finally commits the order (Buy Confirm). We use our own implementation of the TPC-W benchmark, realized as a J2EE component-based application [7].

### C. Web workload model

In this study, as in numerous other web server performance studies, we use synthetic web workloads, which are injected into a working application server environment using a load generator machine. Utilizing synthetic workloads is a common and widely adopted way to evaluate web server performance. Although not as realistic as using real web traces, this approach is more convenient for controlled exploration of the range of client behaviors.

Several session-based web workload models have been proposed, based on a detailed analyses of real web traces [8]–[12]. A dominant fraction of these models [8], [9], [13]–[16], as well as workload generators of web server performance benchmarks, such as TPC-W [3], use first or higher-order Markov chains to model session structure.

Our study follows this trend and adopts the *Customer Behavior Model Graph* (CBMG) [8] approach for session structure modeling. CBMG is a state transition graph (i.e., a first-order Markov chain), where states denote results of

service requests (web pages), and transitions denote possible service invocations. Transitions in CBMG are governed by probabilities  $p_{i,j}$  of moving from state  $i$  to state  $j$  ( $\sum_j p_{i,j} = 1$ ). It was shown that web workloads consisting of several different CBMG session structures can approximate any given sequence of user requests (web request log) as “close” as desired, by appropriately choosing the model parameters (i.e., the number of CBMGs and their transition probabilities), which in turn can be obtained from the web request log by the proposed clustering algorithm [8]. The greater the number of CBMGs in the workload model, the closer such an approximation can be made. In web workloads consisting of several CBMGs, each one of them represents a typical navigational pattern exhibited by the service users.

**Session structure.** Our workload for the TPC-W application is a 50%/50% mix of the two CBMGs shown in Fig. 2. We use simplified user session structures, which use only a subset of the TPC-W request types, but are rich enough to include essential application activities and represent requests with a wide range of functional and execution complexity. Both of the CBMGs in Fig. 2 use the same state transition structure, but with different transition probabilities. The “Mostly Buyers” CBMG produces user sessions that tend to buy products, while the “Mostly Browsers” CBMG produces more browsing-biased sessions. This results in different frequencies of requests being invoked by the two kinds of sessions (see Fig. 2, right). Note that not all “Mostly Buyers” sessions result in a purchase, and analogously, not all “Mostly Browsers” sessions just browse the product catalog. The 50%/50% mix of the given “Mostly Buyers” and “Mostly Browsers” sessions results in approximately 52% of sessions finishing with a purchase. This value may be higher than most retail e-Commerce web sites see in real life, however such client behavior may be more characteristic of web sites providing online brokerage services, where a greater portion of user sessions results in completion of reward-bringing transactions of selling and buying stocks. We introduce this bias towards purchasing sessions to highlight the benefits of our request prioritization approach. However, we expect our methods to exhibit the same *relative* improvements even in workloads with fewer purchasing sessions.

**Timing parameters.** We model session inter-request user think times as exponentially distributed with mean 5s for the “Mostly Buyers” sessions and 10s for the “Mostly Browsers” sessions. If not stated otherwise, the flow of incoming new sessions is modeled as a Poisson process with arrival rate  $\lambda$ , which determines the overall load produced on the system (average request rate received by the service is  $\lambda \cdot N$ , where  $N$  is the average session length, in requests). The Poisson process produces relatively smooth sequence of events, and fails to model inherently bursty and self-similar traffic typically observed at web sites [17]. To better model the latter, we additionally use the B-model [17], which has been shown to produce synthetic traces with burstiness matching that of real traffic. We use this model to produce load with different *degrees of burstiness* (determined by the  $b$ -parameter of the

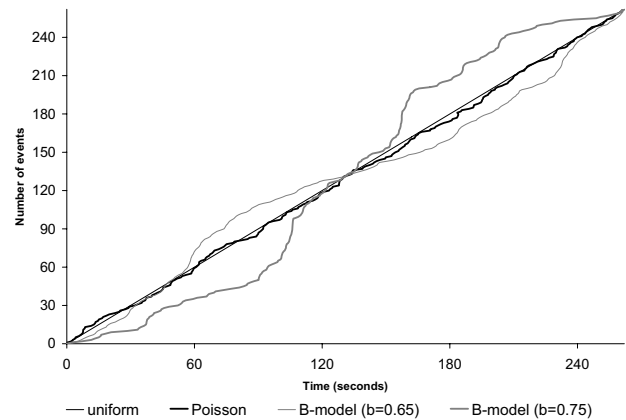


Fig. 3. Event arrival patterns for the three processes: Poisson ( $\lambda = 1$ ) and B-model ( $b=0.65$  and  $b=0.75$ ).

B-model), and do it in a way to only imitate local (short-lived) burstiness to avoid substantial shifting of massive event clusters to short time intervals. Specifically, we model two types of bursty load, one with  $b = 0.65$  and another with  $b = 0.75$ ,<sup>2</sup> and refer to these as “low-bursty” and “high-bursty” load respectively. In contrast with these two methods, we refer to the Poisson arrival model as “smooth.” Fig. 3 shows the event arrival patterns for a Poisson process ( $\lambda = 1$ ), and for the two B-model processes ( $b = 0.65$  and  $b = 0.75$ ) with the same average event arrival rate (1 event/s). This graph helps to visually assess the degree of event arrival burstiness produced by the different models.

#### D. Session reward and request cost specification

It is the service provider responsibility to define the *reward* function associated with the session. The model we adopt in this study is simple yet general enough to encompass several possible applications: a reward value is defined for every request type of the service. The reward of the session is the sum of rewards of the requests in the session. The reward “counts” only if the session completes successfully.

To illustrate the reward formulation, let us revisit the example scenarios presented in Section I. In the on-line shopping scenario the profit of the service is reflected by the volume of items sold. So one way to define a reward function for the on-line store service is by assigning a reward value of 1 for the **Add to Cart** request – the shopping cart will contain as many items in it as the number of times the **Add to Cart** request was executed. In the example of third party-sponsored advertisements, the idea of reward specification is straightforward – assign each such web page a reward value based on the agreement between the service and the third party, e.g., based on how much the latter pays the former for a client’s click on this page. In the example where the service provider wants to maximize the number of successful sessions the reward specification is done by assigning a reward value

<sup>2</sup>In the original B-model study [17], the authors analyzed real web traces and inferred that the  $b$ -parameter for that traces ranged from 0.6 to 0.8, so we felt that values 0.65 and 0.75 would be reasonably representative.

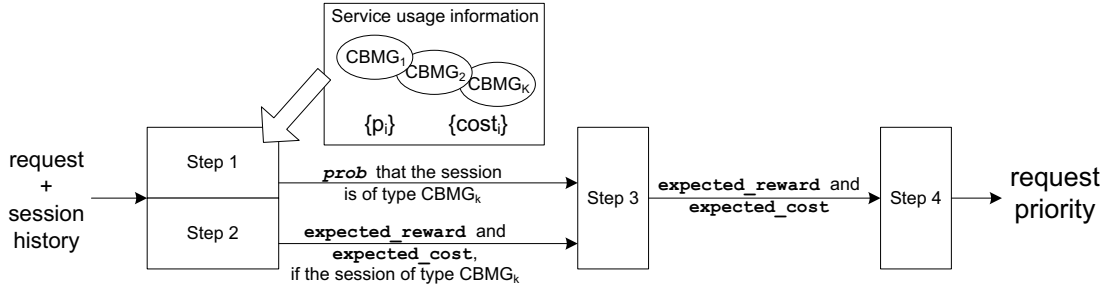


Fig. 4. Logical steps of the RDRP method.

of 1 to the request corresponding to a visit to the “success” page.

The cost of a session is similarly specified in terms of the *relative request execution cost* for each request type. This choice has the following rationale. Processing times for individual requests in typical services, including our TPC-W application, can vary widely by as much as two-to-three orders of magnitude. However, there tends to be much more variation across request types than for requests within the same type but with different request parameters [15], [18]. Information about the relative execution costs permits the RDRP Bayesian inference algorithm to be able to make adequate predictions of future server resource consumption by a session. Note that request execution costs can either be specified by the service provider (in abstract cost units), or determined by online request profiling as the average processing time of requests of a particular type.

### III. REWARD-DRIVEN REQUEST PRIORITIZATION

Our reward-driven request prioritization (RDRP) algorithms work with the assumption that information about the aggregate structure of the user load is known. Specifically, we assume that the workload consists of  $K$  CBMGs:  $CBMG_1, CBMG_2, \dots, CBMG_K$ . The probability of a session having the structure of  $CBMG_k$  is  $p_k$ ,  $\sum_{k=1}^K p_k = 1$ . For the session of structure  $CBMG_k$ , the probability of transition from state  $i$  to state  $j$  is  $p_{i,j}^k$ . We also assume that as stated earlier, that for each request type  $i$ , its *relative execution cost* –  $cost_i$  – is known. We note that CBMG structures can be extracted from web server request logs through offline or online cluster analysis [8], and the various probabilities and per-request type execution costs can be updated at runtime through request profiling.

Given the above, the RDRP mechanism works in the following way. For every incoming request, it looks at the sequence of requests already seen in the session and compares this sequence with the known CBMG structures of the session types comprising the user load. A Bayesian inference analysis estimates the probability that the given session is of type  $CBMG_k$ , for each  $k = 1, \dots, K$  (step 1). For each session type  $CBMG_k$ , the algorithm computes the values of *expected reward* and *execution cost*, resulting from the *future* requests of the session, assuming it had the structure  $CBMG_k$  (step 2). This information is used to get the non-conditional values

of expected reward and execution cost of the future session’s requests (step 3). These values are used then to define the *priority* of the request (step 4), which governs the scheduling of available server threads and DB connections to incoming requests (see Fig. 1). The logical sequence of the RDRP algorithm steps is depicted in Fig. 4 and is explained in detail below.

**Step 1.**  $\Pr\{CBMG_k \mid \text{req hist}\}$ , the Bayesian estimate that the session is of certain type  $CBMG_k$  (for a given history of session requests) is given by the following formula:

$$\Pr\{CBMG_k \mid \text{req hist}\} = \frac{\Pr\{\text{req hist} \mid CBMG_k\} \cdot p_k}{\sum_{i=1}^K \Pr\{\text{req hist} \mid CBMG_i\} \cdot p_i} \quad (1)$$

where  $\Pr\{\text{req hist} \mid CBMG_k\}$  is the probability of having a certain sequence of  $L$  requests  $\{i_1, i_2, \dots, i_L\}$  in a session of type  $CBMG_k$  and is determined as

$$\Pr\{\text{req hist} \mid CBMG_k\} = \prod_{j=1}^{L-1} p_{i_j, i_{j+1}}^k \quad (2)$$

**Timing parameters.** In the basic Bayesian analysis of distinguishing among possible session types (equations (1) and (2)), we took into account only the CBMG state transition information. However, if session inter-request user think times differ for sessions of various CBMG types, than this additional information can be used in an attempt to make the Bayesian inference analysis more accurate. Imagine that we know the distribution of user think times for each CBMG session type comprising the load, in particular – their PDF functions,  $PDF_k(x)$ ,  $k = 1, \dots, K$  ( $PDF_k(x) = \Pr\{\text{time} < x\}$ , for  $CBMG_k$ ), and that the observed session inter-request times are  $t_1, \dots, t_{L-1}$  ( $L$  is the number of requests seen in the session). Then equation (2) can be substituted by the following one:

$$\begin{aligned} \Pr\{\text{req hist}, t_1, \dots, t_{L-1} \mid CBMG_k\} = \\ = \prod_{j=1}^{L-1} p_{i_j, i_{j+1}}^k \cdot \prod_{j=1}^{L-1} PDF_k(t_j) \cdot (\Delta t)^{L-1} \end{aligned} \quad (3)$$

where the infinitesimal time interval  $\Delta t$  appears in the equation, because the session inter-request times have supposedly continuous distributions. When equation (3) is substituted in equation (1), the infinitesimal value  $(\Delta t)^{L-1}$  appears in both the numerator and the denominator, and cancels each other out. We call the RDRP scheme that involves inter-request timing



to maintain histories of session requests. The low-level request processing information is used to periodically update the values of relative request execution cost  $cost_i$  in the RDRP algorithm, defined in our experiments as the average request processing time, without the time spent waiting for a thread or a DB connection.

**TPC-W application.** As stated earlier, this study uses our own implementation of the TPC-W benchmark, realized as a J2EE component-based application [7]. The TPC-W application parameters (e.g., for database population) are chosen so as to achieve diverse execution complexity for different request types involved in the simulated sessions. The average response times for the TPC-W request types, when executed in isolation (only one request is processed by the server at a time), range from 5 ms for the **Register** request (which does not require database access) to 450 ms for the **Search** request (which performs execution of complex database queries). When executed concurrently, the requests see larger response times, because of queuing delays for critical server resources (threads and DB connections) and possible database contention.

**Client load.** A separate workstation is used to produce client load and to gather statistics. The client load simulates requests according to the CBMG structures discussed in Section II-C. The maximum sustainable request rate of the server configuration under the resulted request mix is approximately 20 req/s, with the bottleneck being the MySQL database server.<sup>3</sup> The overall load produced on the system is determined by  $\lambda$  – the arrival rate of new sessions. We use different values of this parameter to generate server overload as well as underload conditions and report the load measured as a percentage of the system processing capacity.

**Reported metrics.** For each experiment, we measure the reward attained by the service (i.e., number of items bought by successfully completed sessions) and average request response times for sessions bringing different reward. We do not report the reward metric for the server underload situation, because in this situation all sessions complete successfully, and each request scheduling algorithm produces the same reward value. Where absolute values of reward are reported, they are counted *per incoming user session*. This is done to show how close the employed algorithms are to the ideal situation, when all the *buying* sessions complete successfully, which brings the average per-session reward of 0.7 (this value is determined by the mix and the structure of the involved CBMGs shown in Fig. 2). In some of the experiments we show attained reward measured as a percentage of the reward value produced by the default FIFO request scheduling algorithm. This is done to emphasize the relative benefits that employed methods bring, compared with the default web application server policies.

<sup>3</sup>This seemingly low server throughput is attributed, first, to the underprovisioned one generation old machines we were using for the experiments, and second, to the fact that we did not perform scrutinized database and TPC-W application tuning. However, we expect the relative performance improvements achieved by the RDRP methods to be similar in more powerful server environments.

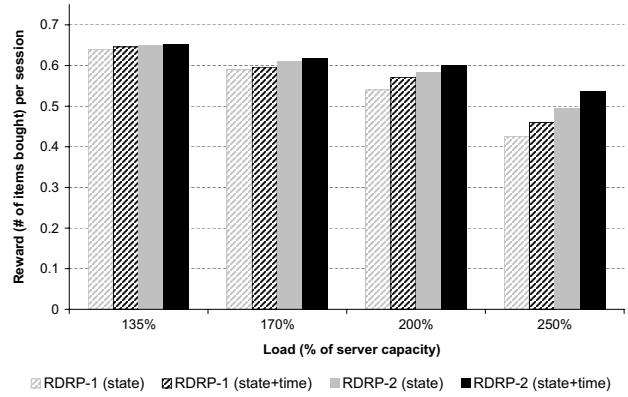


Fig. 6. Comparison of benefits brought by the two flavors of the RDRP method. RDRP-1 takes into account the execution cost of all session’s requests (seen and expected). RDRP-2 only takes into account the cost of current and the future requests of the session.

### B. Comparison of two priority schemes

We ran a set of experiments comparing the performance of the RDRP-1 and the RDRP-2 methods, corresponding to the two priority formulations in equations (4) and (5), under various load conditions. Figure 6 compares the performance of the two methods under different amounts of server overload, for “smooth” session arrivals. The RDRP-2 method outperforms RDRP-1 in all scenarios, but especially under high client load. To informally understand this outcome, consider a session that is just one or two steps away from its completion (e.g., it is in the **Register** state in the CBMG of Fig. 2). The RDRP-2 method, according to equation (5), gives this request a higher priority than RDRP-1 (equation (4)), because it does not count the cost already incurred by the session. Consequently, under RDRP-1, the request might get rejected due to a low priority value, which will waste all the effort it took to bring the session to its nearly complete state. A careful examination of the logs produced during the experiments supports this explanation: the primary reason for the poor performance of RDRP-1 is the fact that some sessions are rejected with one or two requests left to complete the session, a phenomenon that never happens with RDRP-2. By ignoring the cost already incurred by the session, the RDRP-2 method appears to increase the likelihood of session completion as compared to its RDRP-1 counterpart. This also agrees with economics theory, which argues that *sunk costs* (i.e., costs that have already been incurred and which cannot be recovered, like  $cost_{incurred}$  in equation (4)), should not be taken into account when making rational decisions [19]. In the rest of the experiments we used only the RDRP-2 algorithm, and refer to it from now on as simply the RDRP method.

### C. Imitating the “history-based” approach

As stated in Section I, an alternative method to prioritize client requests to improve service reward is by using a per-client history-based approach. Broadly considered, such an approach models the behavior of any application-specific technique in which all requests of a session are assigned a constant

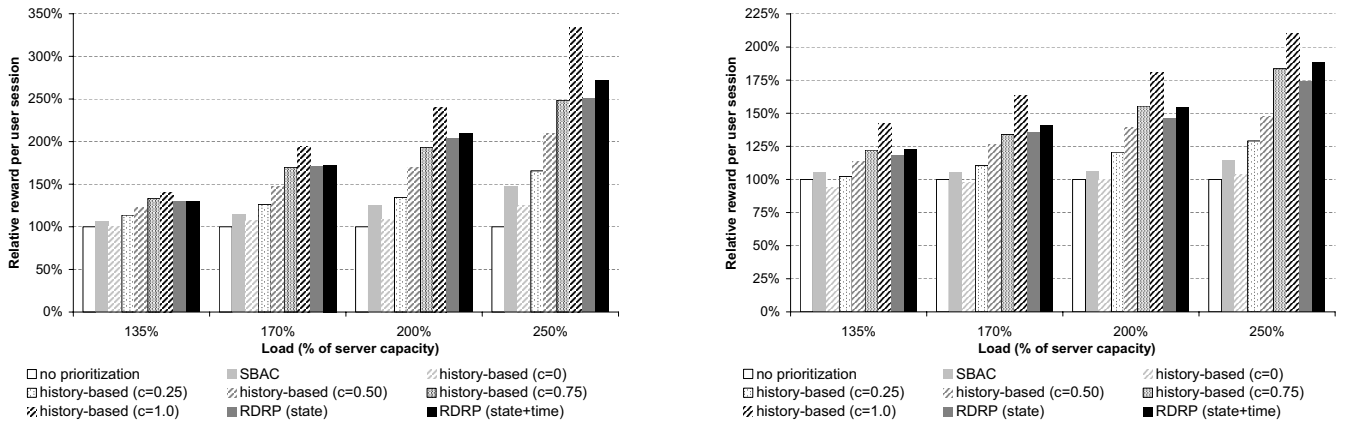


Fig. 7. Reward (number of items bought by successfully completed sessions, per user session), relative to the default no-prioritization (FIFO) scheme, for the “smooth” Poisson (left) and “high-bursty” (right) client loads.

priority value and are scheduled according to this priority. The priority assignment can have arbitrary logic, for example, it can be done in an attempt to predict the client’s future behavior based on the history of the client’s previous purchases, or it can be determined solely by the client’s membership status. The success of such approaches is determined, of course, by how good they are in predicting the client’s behavior or, more precisely, the statistical correlation between assigned session priority and the actual reward brought by this session. To the best of our knowledge, prior work on workload characterization has not addressed such correlation in behavioral patterns (especially with the information that we need). We therefore employed the following scheme for producing a predefined correlation between the assigned session priorities and the actual rewards brought by the sessions. Each session announces in advance the reward it would bring, enabling the session prioritization mechanism to set the session’s priority so that the statistical correlation (parameter  $c$ ) between the assigned priorities and the sessions’ rewards meets the predefined value. A value of  $c = 1.0$  brings the best achievable performance because the prioritization algorithm always assigns to requests from the session, a priority value in direct correspondence with the reward the session will bring.

#### D. Performance of RDRP

We compare the relative costs and benefits of RDRP mechanisms against the following alternative server-side request scheduling and overload protection methods:

- Default FIFO request scheduling with no request prioritization.
- Session-Based Admission Control (SBAC), which admits approximately as many sessions as can be processed by the server capacity; all of the admitted sessions are allowed to complete successfully. This method is used only in the server overload situation.
- The per-client “history-based” approach described in Section IV-C. We run five sets of experiments with  $c = 0, 0.25, 0.5, 0.75, \text{ and } 1.0$ .
- Our RDRP(state) and RDRP(state+time) methods, described in Section III.

1) *Server overload*: First, we evaluate the behavior of the methods in server overload situations. We run four sets of experiments, modeling loads of 135%, 170%, 200%, and 250% of server capacity, for both the “smooth” (Poisson) and “high-bursty” (B-model with  $b = 0.75$ ) client loads (see Section II-C for details). Fig. 7 shows the reward attained by the service, relative to the performance of the default FIFO request scheduling mechanism. Figs. 8 and 9 show average request response times for sessions bringing different reward, for the “smooth” and “high-bursty” client loads. Several conclusions can be drawn from the results of these experiments.

**Reward attained.** As expected, the default FIFO request scheduling policy shows the worst performance, because a request may get rejected anywhere in the session, which results in low successful session throughput. The SBAC method works better, because it at least allows the sessions that have started to complete successfully, however it does not try to necessarily admit those sessions that bring the greatest reward. The history-based approach shows an increase in reward attained with an increase of the correlation between assigned session priorities and sessions’ rewards. Note that even with values of  $c = 0.25$ , this method already outperforms the SBAC algorithm. Finally, both RDRP methods significantly boost reward attained by the service. The RDRP(state+time) method works slightly better than RDRP(state), because it takes into account the inter-request time differences between more-profitable “Mostly Buyers” sessions and less-profitable “Mostly Browsers” sessions and better distinguishes between them. The theoretically best history-based ( $c = 1.0$ ) method, of course, shows the best performance, however the history-based approach matches the performance of the RDRP algorithms, only for values of  $c \geq 0.75$ .<sup>4</sup> The performance of all algorithms goes down, when the client load experiences bursty behavior, because under bursty conditions the queues for critical server resources are more susceptible to rapid build-ups, which results in higher rates of request rejections. However, the relative advantages of RDRP over the other

<sup>4</sup>Whether such good prediction is possible in real life, remains an open question, due to the lack of publicly available information with such statistics.

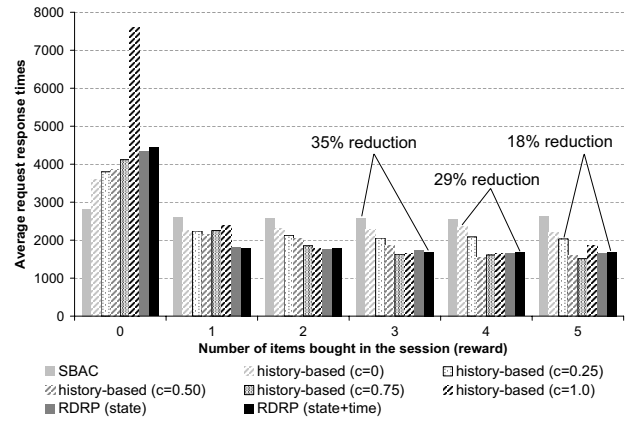
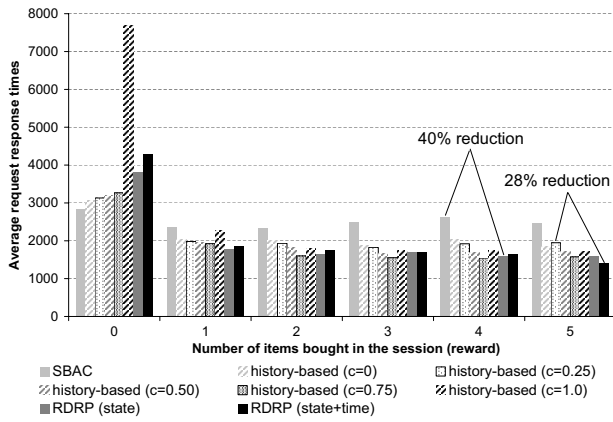


Fig. 8. Average request response times for sessions that bring different reward, for “smooth” traffic, for the 135% server capacity (left) and the 170% server capacity (right) overload situations.

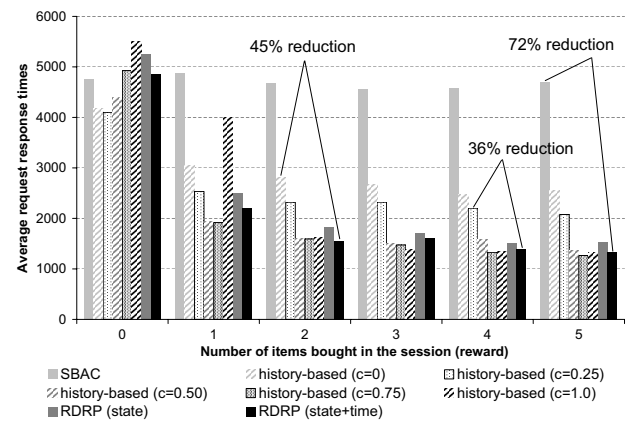
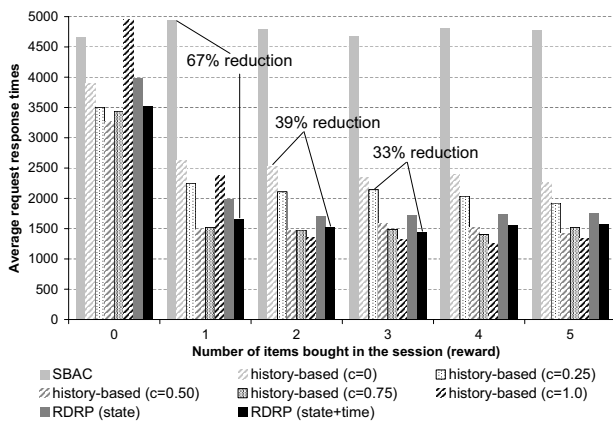


Fig. 9. Average request response times for sessions that bring different reward, for “high-bursty” traffic, for the 135% server capacity (left) and the 170% server capacity (right) overload situations.

methods stay the same.

**Request response times.** All algorithms that perform request/session prioritization, and manage to correctly guess (at least to a certain degree) the session’s reward, decrease request response times for sessions that bring non-zero reward, as compared to the SBAC method. Both RDRP methods perform on par with the history-based approach for values of  $c \geq 0.5$ . For the “smooth” client load, the RDRP algorithms reduce response times by up to 40% compared to SBAC, and show up to 28% lower response times than the history-based approach with  $c = 0$  and  $c = 0.25$ . For bursty client load, the difference is more pronounced: response times from the RDRP methods are lower than that from SBAC and the history-based approach with  $c = 0$  and  $c = 0.25$  by up to 72%, 45%, and 36%, respectively.

Note, that for “smooth” client load, the sessions with zero reward (i.e., *browsing* sessions) see significantly increased response times, when the history-based approach with  $c = 1.0$  is applied (Fig. 8). This happens, because with the history-based approach, all browsing sessions (48% of all sessions, see Section II-C for an explanation) get the same (zero) priority, because the priority is defined as the session’s reward, while the remaining 52% of sessions get a higher execution priority.

Being all stuck in a single lowest-priority queue (with a FIFO tiebreaking policy), browsing sessions see higher rates of request rejections. This in turn produces higher response times for the session because a rejected request spends at least 10s in the system (before experiencing a timeout). Interestingly, this effect is reduced with bursty session arrivals (Fig. 9).

2) *Server underload:* For server underload situations, we ran experiments with both “smooth” (Poisson) and bursty client loads. The Poisson-modeled web workload generated such a smooth flow of request arrivals, that all request scheduling algorithms showed more or less the same performance. This happened because the server queues for the critical resources (threads and DB connections) almost never built up.

Experiments with the bursty client load showed very different behavior. Figures 10 and 11 show average request response times for sessions bringing different reward, for the two bursty client loads.<sup>5</sup> Several conclusions can be drawn from the results of these experiments.

<sup>5</sup>The experiments labeled as “100% of server capacity” were actually ran at a rate slightly lower than the system capacity, which experienced slight variations because of the non-deterministic behavior of the web application server. This ensured that our experiment did not slide into the overload mode of server operation.

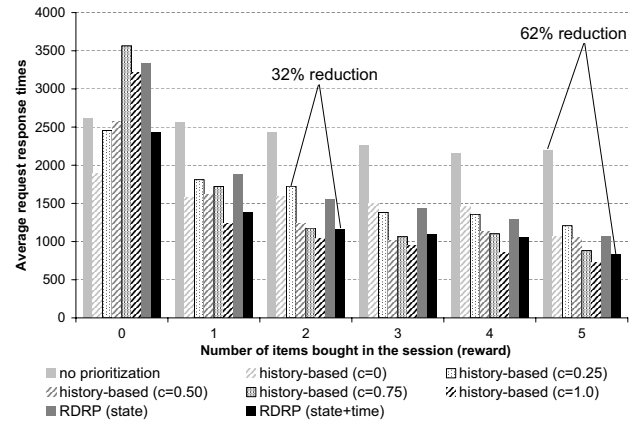
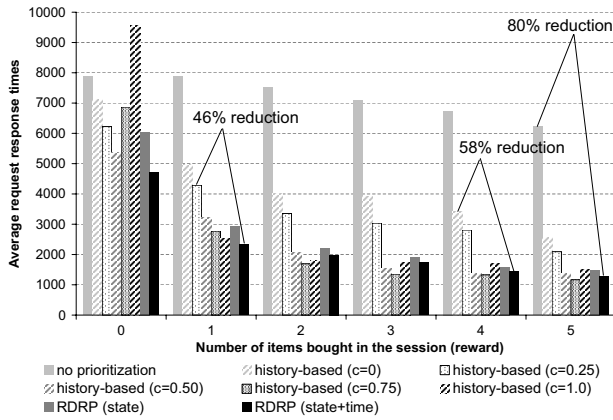


Fig. 10. Average request response times for sessions that bring different reward, for ‘high-bursty’ traffic, for 100% server capacity (left) and the 80% server capacity (right) underload situations.

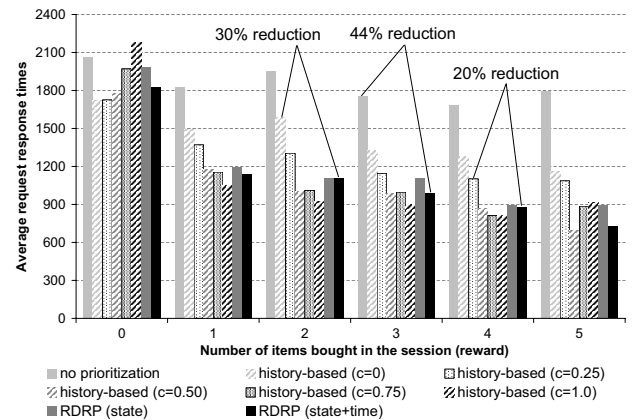
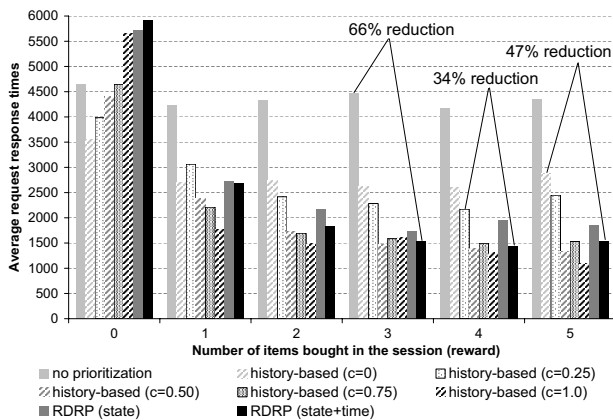


Fig. 11. Average request response times for sessions that bring different reward, for ‘low-bursty’ traffic, for the 100% server capacity (left) and the 85% server capacity (right) underload situations.

The RDRP methods (as well as the history-based approaches) decrease request response times for the sessions that bring non-zero reward. This happens because with bursty arrivals (unlike the smooth arrival case described above), the queues for the critical server resources (server threads and DB connections) occasionally build up, and the request prioritization mechanisms minimize the queueing delays seen by the sessions that bring more reward by assigning their requests higher priorities. For ‘high-bursty’ traffic, the effects of request prioritization are visible for loads above approximately 70% of server capacity (in Fig. 10 we show experiments with the load of 100% and 80% of server capacity), while for the ‘low-bursty’ traffic, the effects are visible for the load in the range of 85%–100% of server capacity.

As in the server overload situation, the performance of the RDRP methods is matched by the history-based approach only for values of  $c \geq 0.5$ . Under ‘high-bursty’ traffic, RDRP outperforms the history-based method by up to 58% (for  $c = 0$ ) and 46% (for  $c = 0.25$ ). This advantage of RDRP over the history-based approach diminishes a bit under ‘low-bursty’ traffic conditions (Fig. 11). The default FIFO method performs worst of all. It is interesting to note that even the history-based approach with  $c = 0$ , which is not supposed to

ever correctly guess the session’s reward, gives lower response times (for all reward values, including 0) than the default FIFO request scheduling scheme.

To our understanding, this behavior happens for the following reason. The request scheduling algorithm we adopt to imitate a history-based approach with  $c = 0$  works by uniformly assigning priorities to sessions as integer values in the range of 0 to 100 (this process does not correlate with the session reward, therefore corresponds to  $c = 0$ ). Some sessions get higher priorities than the other, and all sessions are uniformly sorted into a discrete number of priority buckets. Unlike the FIFO scheduling case, where all requests have to wait in one long queue produced by a traffic burst, the uniform session prioritization scheme permits some sessions to sneak ahead of other sessions. This perturbs the waiting times seen by requests sufficiently so as to achieve an average response time lower than that seen by the FIFO case.

## V. RELATED WORK AND DISCUSSION

The notion of a web session, as a structural organization of client communication with Internet services, was first investigated in [1] and [20]. Since then several other studies have explored session characterization of web workloads [8]–[12],

[21]. The work in [22] and [23] acknowledged that *service usage patterns* affect the performance of Internet services, and that understanding the nature of user workloads is crucial for properly designing and provisioning web servers. Our study follows this trend, by using aggregated information about client usage of an Internet service to boost the service profit or other application-specific reward that clients bring to the service provider. Our work combines information about the workload structures seen at a service with Bayesian inference analysis to guide the scheduling of bottleneck server resources to incoming client requests.

The idea of profiling client requests and gathering service usage statistics, which later can be used for various purposes, is itself not new [24], [25]. An array of data mining techniques has been proposed to extract information from web access logs for myriad applications, including one most relevant for this work: to predict user access sequences [26]–[28]. Our use of this information to influence scheduling of server resources is somewhat novel, differing from past applications of such information that have ranged from reducing user-perceived latencies for personalized web sites [29] to improving web server caching and prefetching behavior [30].

Of more direct relevance to the work described in the paper is the large body of prior work in the areas of web server QoS, overload and admission control, and request scheduling. Due to space limitations, we focus on a few representative approaches most relevant to our work, and skip discussion of other, somewhat orthogonal approaches for improving service performance such as service distribution, load balancing and content adaptation.

Various forms of *admission control* have been used to prevent services from being overwhelmed in the presence of persistent or transient overload. Among these, Session-Based Admission Control (SBAC) [1] is suitable for session-oriented client loads. An overloaded service can experience a severe loss of throughput measured in completed (successful) sessions while still maintaining its throughput measured in requests per second. This happens because a request can be rejected anywhere in the session, even if the session has already had a lot of its requests served and is very close to completion. The SBAC method works by admitting as many sessions as can be processed by the service, trying to make sure that if a client starts a session with the service, it will be successfully completed. However, this approach is oblivious to any application-specific information, e.g., to the reward brought by different sessions.

In recent years, several researchers have investigated the effects of request *scheduling* and *prioritization* on web server performance, for both web sites providing static content [31]–[33] and database-driven dynamic web sites [15], [34]. It has been shown that request response times and server throughput can be improved by employing such scheduling algorithms as Shortest Job First (SJF) [15], [31], [32] and Shortest Remaining Processing Time first [33]. Some of these studies used request scheduling algorithms combined with admission control policies [15], [18], [34], [35]. Our work shares the

same goals as these efforts, but its use of more sophisticated scheduling policies is somewhat constrained by the hooks exposed by the underlying middleware. In web studies focusing on static content, the *cost* of servicing a job is usually approximated by the size of the downloaded file. For web sites serving dynamic content, it was noticed that request processing times depend primarily on the *request type* rather than on the parameters of the request [15], [18]. Our work shares the same observation, using fine-grained request profiling to determine first absolute and then relative request processing times for different request types. An analogous technique is used in [15].

A notable difference from our work is that for the most part, the request scheduling studies above do not pursue the goal of increasing likelihood of session completion (even if they take into account session-oriented nature of client workloads). An exception is the work of Chen et al. [16] on the Dynamic Weighted Fairing Sharing request scheduling algorithm (DWFS), which, among other goals, tries to avoid processing of requests that belong to sessions that are likely to be aborted in the near future. In trying to increase session completion rate, this study shares a commonality with our work, which is oriented towards completion of sessions that bring more service reward.

Several previous studies have also proposed admission control and request scheduling techniques that take into account application-specific information about the reward or profit brought to the Internet service by individual request types, and try to maximize this reward. In [36], a Profit Aware QoS policy (PAQoS) is developed, aimed at maximizing the web site's profit under SLA constraints. The authors of [13] propose using queuing of requests based on their types, where a reward function corresponding to the service provider's objective is maximized using techniques for nonlinear optimization. However, there are conceptual differences between these studies and our approach. In both [36] and [13], the authors assume that reward is brought by individual requests, rather than by completion of the whole session. They also assume a Generalized Processor Sharing (GPS) model for request execution, rather than a model of prioritized scheduling of requests to *exclusively-held* resources, such as server threads and DB connections. In our opinion, the latter model is a closer match to existing web application server architectures.

Finally, although this paper has demonstrated our work on a sample TPC-W application emulating an online store selling books, the underlying ideas and the approach are broadly applicable to all services accessed in a session-oriented fashion. This is in contrast to a recent study [14], which focuses specifically on improving web server performance for the web shopping scenario. In that work, the authors propose a combined LIFO-Priority scheme for overload control of a retail e-Commerce web site, where all service requests are divided among browser and (in their terminology) *revenue-generating transaction* requests (i.e., reward-bringing requests in our terminology). LIFO scheduling is applied to the browser requests, while the revenue-generating requests are given the highest priority. The experimental results show that such scheduling

discipline improves the overall web server throughput while also increasing the completion rate of the revenue-generating requests. Note however, that while the study shows successful execution of revenue-generating requests, unlike our work, it does not explicitly address the issue of ensuring successful completion of *buying* sessions that contain such requests.

## VI. CONCLUSION

In this paper we have proposed *Reward-Driven Request Prioritization* (RDRP) mechanisms, which maximize the reward attained by an Internet service by dynamically assigning higher execution priorities to the requests whose sessions are likely to bring more profit (or any other application-specific reward) to the service. We have implemented the proposed methods as pluggable middleware mechanisms in the J2EE application server JBoss [2], and tested them on the TPC-W benchmark application [3] using CBMG-based web workloads. Our experiments showed that RDRP techniques yield benefits in both underload and overload situations, for both smooth and bursty client behavior, against state-of-the-art alternatives such as session-based admission control and history-based session prioritization approaches. In the situation of service underload the proposed mechanisms gave better response times for the clients that brought more reward. In the situation of service overload, the mechanisms ensured that sessions that brought more reward were more likely to complete successfully and that the aggregate profit attained by the service increased compared to other solutions. Additionally, we showed that the history-based approach matched performance of our RDRP mechanisms only if the correlation between the clients' past and future behavior reached the mark of 75% for the reward attained, and 50% – for the request response times.

## ACKNOWLEDGMENT

This research was sponsored by DARPA agreement N66001-01-1-8929; by NSF grants CAREER:CCR-9876128, CCR-9988176, and CCR-0312956; and Microsoft. The U.S. Government is authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright annotation thereon. The views and conclusions contained herein are those of the authors and should not be interpreted as representing the official policies or endorsements, either expressed or implied, of DARPA, Rome Labs, SPAWAR SYSCEN, or the U.S. Government.

## REFERENCES

- [1] L. Cherkasova and P. Phaal, "Session-based admission control: A mechanism for peak load management of commercial web sites," *IEEE Trans. Comput.*, vol. 51, no. 6, pp. 669–685, 2002.
- [2] JBoss J2EE Application Server. <http://www.jboss.org>.
- [3] Transaction Processing Performance Council. TPC-W Benchmark. <http://www.tpc.org/tpcw/>.
- [4] D. K. Pecaut, M. J. Silverstein, and P. Stanger, "Winning the online consumer: Insights into online consumer behavior," Boston Consulting Group, March 2000. [Online]. Available: <http://www.bcg.com>
- [5] S. VanBoskirk, C. Li, and J. Parr, "Keeping customers loyal," Forrester Research, May 2001. [Online]. Available: <http://www.forrester.com>
- [6] Sun Microsystems. Java 2 Enterprise Edition. <http://java.sun.com/j2ee/>.
- [7] TPC-W-NYU: a J2EE implementation of the TPC-W benchmark. <http://www.cs.nyu.edu/~totok/professional/software/tpcw/tpcw.html>.

- [8] D. Menascé, *et al.*, "A methodology for workload characterization of e-commerce sites," in *Proc. ACM Conf. on Electronic Commerce*, November 1999.
- [9] —, "In search of invariants for e-business workloads," in *Proc. ACM Conf. on Electronic Commerce*, October 2000.
- [10] M. Arlitt, "Characterizing web user sessions," in *Proc. Performance and Architecture of Web Servers Workshop*, June 2000.
- [11] M. Arlitt, D. Krishnamurthy, and J. Rolia, "Characterizing the scalability of a large web-based shopping system," *ACM Trans. Internet Technol.*, vol. 1, no. 1, pp. 44–69, 2001.
- [12] D. Menascé, *et al.*, "Characterizing e-Business workloads using fractal methods," *J. Web Engineering*, vol. 1, no. 1, pp. 74–90, 2002.
- [13] J. Carlstrom and R. Rom, "Application-aware admission control and scheduling in web servers," in *Proc. IEEE INFOCOM'02*, June 2002.
- [14] N. Singhmar, *et al.*, "A combined LIFO-priority scheme for overload control of e-Commerce web servers," in *Proc. IEEE RTSS International Infrastructure Survivability Workshop*, December 2004.
- [15] S. Elnikety, *et al.*, "A method for transparent admission control and request scheduling in dynamic e-Commerce web sites," in *Proc. WWW'04*, May 2004.
- [16] H. Chen and P. Mohapatra, "Session-based overload control in QoS-aware web servers," in *Proc. IEEE INFOCOM'02*, June 2002.
- [17] M. Wang, *et al.*, "Data mining meets performance evaluation: Fast algorithms for modeling bursty traffic," in *Proc. ICDE'02*, February 2002.
- [18] X. Chen, H. Chen, and P. Mohapatra, "An admission control scheme for predictable server response time for web accesses," in *Proc. WWW'01*, May 2001.
- [19] H. R. Varian, *Intermediate Microeconomics: A Modern Approach*, 7th ed. W. W. Norton & Company, 2005.
- [20] D. Krishnamurthy and J. Rolia, "Predicting the performance of an e-Commerce server: Those mean percentiles," in *Proc. ACM SIGMETRICS Workshop on Internet Server Performance*, June 1998.
- [21] W. Shi, *et al.*, "Workload characterization of a personalized web site – and its implications for dynamic content caching," in *Proc. International Workshop on Web Caching and Distribution*, August 2002.
- [22] P. Barford, *et al.*, "Changes in web client access patterns: Characteristics and caching implications," *WWW Journal*, vol. 2, no. 1, pp. 15–28, 1999.
- [23] D. Llambiri, A. Totok, and V. Karamcheti, "Efficiently distributing component-based applications across wide-area environments," in *Proc. ICDCS'03*, May 2003.
- [24] A. Joshi, K. Joshi, and R. Krishnapuram, "On mining web access logs," in *Proc. ACM SIGMOD Workshop on Research Issues in Data Mining and Knowledge Discovery*, May 2000.
- [25] J. Srivastava, *et al.*, "Web usage mining: Discovery and applications of usage patterns from web data," *SIGKDD Explorations*, vol. 1, no. 2, pp. 12–23, 2000.
- [26] J. Pitkow and P. Piroli, "Mining longest repeating subsequences to predict World Wide Web surfing," in *Proc. USITS'99*, 1999.
- [27] A. Nanopoulos, D. Katsaros, and Y. Manolopoulos, "Effective prediction of web-user accesses: A data mining approach," in *Proc. WEBKDD'01 Workshop*, August 2001.
- [28] E. Frias-Martinez and V. Karamcheti, "A prediction model for user access sequences," in *Proc. WEBKDD'02 Workshop*, July 2002.
- [29] —, "Reduction of user perceived latency for a dynamic and personalized web site using web-mining techniques," in *Proc. WEBKDD'03 Workshop*, August 2003.
- [30] Q. Yang, *et al.*, "Mining web logs to improve web caching and prefetching," in *Proc. Asia-Pacific Conf. on Web Intelligence*, October 2001.
- [31] L. Cherkasova, "Scheduling strategy to improve response time for web applications," in *Proc. International Conf. on High Performance Comput. and Networking (HPCN-Europe)*, April 1998.
- [32] M. Crovella, R. Frangioso, and M. Harchol-Balter, "Connection scheduling in web servers," in *Proc. USITS'99*, October 1999.
- [33] B. Schroeder and M. Harchol-Balter, "Web servers under overload: How scheduling can help," Computer Science Dept. Carnegie Mellon Univ., Tech. Rep. CMU-CS-02-143, May 2002.
- [34] N. Bhatti and R. Friedrich, "Web server support for tiered services," *IEEE Network*, vol. 13, no. 5, pp. 64–71, 1999.
- [35] P. Bhoj, S. Ramanathan, and S. Singhal, "Web2K: Bringing QoS to web servers," HP Laboratories, Tech. Rep. HPL-2000-61, May 2000.
- [36] Q. Zhang, E. Smirni, and G. Ciardo, "Profit-driven service differentiation in transient environments," in *Proc. IEEE MASCOTS'03*, October 2003.