

Web workload generation challenges – an empirical investigation

Raoufhsadat Hashemian^{1,*}, Diwakar Krishnamurthy¹ and Martin Arlitt^{1,2}

¹The University of Calgary, Calgary, Alberta T2N 1N4, Canada

²HP Labs, Palo Alto, CA 94304, U.S.A.

SUMMARY

Workload generators are widely used for testing the performance of Web-based systems. Typically, these tools are also used to collect measurements such as throughput and end-user response times that are often used to characterize the QoS provided by a system to its users. However, our study finds that Web workload generation is more difficult than it seems. In examining the popular RUBiS client generator, we found that reported response times could be grossly inaccurate, and that the generated workloads were less realistic than expected, causing server scalability to be incorrectly estimated. Using experimentation, we demonstrate how the Java virtual machine and the Java network library are the root causes of these issues. Our work serves as an example of how to verify the behavior of a Web workload generator. Copyright © 2011 John Wiley & Sons, Ltd.

Received 6 October 2010; Revised 2 April 2011; Accepted 11 April 2011

KEY WORDS: workload generator; performance testing; benchmarking tools

1. INTRODUCTION

Web applications are used by many organizations to provide services to their customers and employees. An important consideration in developing such applications is the QoS that the users experience. This motivates the organizations to experimentally evaluate properties such as response time and throughput so that the service can be improved until a desired level of QoS is provided.

Practitioners and researchers rely on benchmarking systems such as RUBiS [1], TPC-W [2] and SPECweb [3] to evaluate the performance of their information technology (IT) infrastructure before putting it into production. These benchmarking systems typically contain a specially developed Web application. For example, RUBiS provides several different implementations of an online auction application using Web application platforms such as PHP [4] and Enterprise Java Beans (EJB) [5]. These benchmark systems also contain their own workload generator, designed to work specifically with that benchmark application. The *workload generator* is a tool that generates a synthetic workload to the benchmark application to emulate the behavior of the application's end users. The workload generator reports metrics such as response times for the emulated users and application throughput. These are typically used to assess the QoS provided by the system that executes the benchmark application. In addition to benchmark-specific workload generators, there are numerous general purpose Web request generation tools such as httpperf [6], S-Client [7] and JMeter [8]. However, a significant effort may be required to customize these tools to work with a specific application.

A desirable property of a Web workload generator is that it sends and receives requests to a specified Web server in a realistic manner. Although this is a simple property to describe, it can

*Correspondence to: Raoufhsadat Hashemian, Electrical and Computer Engineering, University of Calgary, Calgary, Alberta T2N 1N4, Canada.

†E-mail: rhashem@ucalgary.ca

be challenging to achieve in practice. Software bottlenecks, hardware bottlenecks, and software implementation features pertaining to the workload generation infrastructure can result in unrealistic workloads being generated or inaccurate measurements being recorded. To minimize the implications of measurement errors, a thorough examination of the workload generator is needed. Unfortunately, this aspect is often ignored in practice, thereby putting into risk the validity of the entire performance testing exercise.

In fact, the work described in this document was motivated by concerns about the validity of results from a benchmarking study of a multi-tier Web application system. Specifically, we found the results from an experimental testbed using the RUBiS benchmark were inconsistent across separate invocations of the benchmark. We decided to investigate the causes of such discrepancies through a controlled experimental study.

The details of our study are provided in the remainder of this document. The salient findings of our research are as follows:

- The workload generator can introduce significant errors in the measured end-user response times. In our experimental environment, we examined improper selection and tuning of Java virtual machine (JVM) as two examples of common mistakes, which can contribute to erroneous response time measurements.
- Multi-threading does not ensure a scalable workload generator. For example, we found the multi-threaded RUBiS workload generator, referred to as *RUBiS client* in this work, supported fewer concurrent users on a client host than the single-threaded *httperf* tool.
- An unexpected lack of realism in the generated workload can result in the scalability of the Web application being incorrectly estimated. For example, we found that the Transmission Control Protocol (TCP) connection management policy used (via a Java library) by RUBiS client is not appropriate for workload generation. The policy shared a single TCP connection across multiple emulated users. This results in an incorrect estimation of the overhead of TCP connection establishment, which a server experiences in practice.

Our work serves as an example of how to verify whether a Web workload generator is behaving as intended. Specifically, we offer a methodology to help recognize problems that could adversely impact the validity of performance testing exercises.

The remainder of this paper is organized as follows. Section 2 provides background information and examines related work. Section 3 describes our experimental environment, including the enhancements made to *httperf* and RUBiS client to support this work, and the network monitoring tools used to validate measurements from the workload generators. Section 4 explains our investigation of the challenges of Web workload generation. Section 5 discusses the implications of our findings. Section 6 summarizes our work and provides directions for future research.

2. BACKGROUND AND RELATED WORK

Workload generators for Web applications rely on the concept of an emulated user. An emulated user submits sequences of inter-dependent requests called *sessions* to a system under study. Dependencies arise because some requests in a session depend on the responses of earlier requests in the session. For example, an order cannot be submitted to an e-commerce system unless the previous requests have resulted in an item being added to the user's shopping cart. This phenomenon is known as an *inter-request dependency*. In each session, an emulated user issues a Web request, waits for the complete response from the system, and then waits for a period of time defined as the *think time* before issuing the next request. A workload generator typically emulates multiple users concurrently. The workload generator runs on one or more *client hosts*, depending on how many emulated users need to be generated.

A challenge for workload generation is ensuring that the synthetic workloads generated are representative of how real users use the system under study. In particular, one must carefully select realistic characteristics for workload attributes that can impact performance, such as the mix of different types of requests submitted to the system, the pattern of arrival of users to the system, and the

think times used within user sessions. Furthermore, the synthetic workload must preserve the correct inter-request dependencies to stress the system's application functions as intended. A discussion on the perils of selecting incorrect workload characterizations can be found in [9, 10]. Ferrari proposed several guiding principles to avoid common mistakes in designing synthetic workloads [11]. More recently, Barford and Crovella discussed the essential elements required for the generation of representative Web workloads [12].

Workload generators can differ in the type of user emulation they support. The most common type of user emulation employed is called the closed approach. With a *closed* approach, the number of concurrent users during any given test is kept constant. The next request in a user session is not generated until the previous request has been completed, and the think time has been observed. The load on the system can be controlled by manipulating the number of concurrent users in a test. With an *open* approach, users submit requests according to a specified rate, without waiting for the response of any of their previous requests that have not been completed in the expected time interval. This approach is useful for evaluating Web applications under overload conditions. However, it can violate inter-request dependencies. A hybrid approach combines aspects of the closed and open approaches. With a *hybrid* approach, user sessions are initiated at specified time instants. It is similar to the open approach in that a new session can be initiated before the previous sessions finish. However, similar to the closed approach, within each session, a request can only be issued after the response to the previous request in that session has been received. With the hybrid approach, the number of concurrent user sessions can change over the course of a test. Schroeder *et al.* argue that the hybrid approach is more representative of real systems than either the closed or open approaches [13].

Workload generators can also differ in the programming paradigms they employ. Most tools follow an approach that uses a combination of multi-threading and synchronous HTTP request-response handling. With this approach, each thread independently emulates the behavior of a user. An alternative approach uses an event-driven mechanism that relies on a single thread and asynchronous HTTP request-response handling. With this approach, a single thread switches between individual users as events such as HTTP requests and responses occur.

Numerous general purpose workload generation tools have been developed. Examples include S-Client [7], httpperf [6], SURGE [11], GEIST [14], WAGON [15], JMeter [8], and SWAT [16]. Typically, general purpose workload generators permit workload characteristics to be controlled in a fine-grained manner. This allows a system's performance to be studied under many different workloads of interest. Whereas GEIST and S-Client support only the open approach, httpperf and SWAT implement the open, closed, and hybrid approaches. WAGON employs the hybrid approach whereas JMeter and SURGE support only the closed approach. S-Client, httpperf, and SWAT (which is built on top of httpperf) follow the single-threaded, event-based design, whereas the other tools listed employ the multi-threaded paradigm.

As mentioned in section 1, practitioners and researchers use Web application benchmark systems such as SPECWeb [3], RUBiS [1], and TPC-W [2] extensively in performance-related studies. The workload generators that are bundled with these benchmark systems are not intended to be as flexible as general purpose tools with regards to fine-grained control over workload characteristics. However, they offer an advantage over general purpose tools in that they can be used "off-the-shelf" without the need for time-consuming customizations to handle inter-request dependencies of the benchmark applications.

In this study, we focus on the RUBiS client tool, as it was the workload generator we were using when we observed the inconsistent results. For comparison purposes, we enhanced httpperf so that it is capable of repeating the workload generated by RUBiS client. RUBiS client [17–20] and httpperf [21–25] have been used extensively to support experimental performance evaluation of virtualization techniques [18, 22, 24], multi-tier software systems [17, 19, 20, 23] and hardware platforms [21, 25]. The results of tests using these tools are often used in critical decision making. An unrecognized bottleneck caused by the tools or any other unintended behavior may affect the estimate of performance (i.e. how good is the experience of each user) as well as the estimate of the system's scale (i.e. how many concurrent users can be supported). Consequently, we believe that the outcome of our investigations is likely to be of interest to other researchers and practitioners.

Although benchmark systems are used extensively, there is very little work that focuses on systematically evaluating whether these tools function as intended. An exception is the work of Nahum [26], who compared the workload characteristics of the SPECWeb99 benchmark with the workloads observed at six production Web server systems. The study found that SPECWeb99 modelled certain characteristics such as file popularity well, but did a poor job of matching characteristics such as file size and HTTP response size. In another study, Nahum *et al.* show that not considering workload attributes such as network delays and packet losses during workload generation can result in overly optimistic estimates of server performance [27].

Our work differs from these studies in that it does not focus on validating the choice of workload characteristics or workload attributes used in benchmarks. Instead, we demonstrate and explain why some Java based workload generators can offer misleading performance results. We show that the results from such tools can lead to incorrect conclusions in performance studies.

3. EXPERIMENTAL ENVIRONMENT

In this section, we describe our testbed configuration, the network monitoring tools we used to validate the measurements reported by workload generators, and the subtle modifications we made to the workload generators to facilitate a direct comparison between them.

3.1. Testbed configuration

Our testbed contains three physical machines connected by a 1 Gb/s Ethernet switch. The specifications of these machines are described in Table I. A schematic of our testbed is shown in Figure 1. One of the three machines is used as the client host whereas the other two are used as Web and database servers, respectively. The application server is also installed on the Web server machine. The client machine runs a workload generator to exercise the RUBiS auction site installed on the two server machines. Several versions of RUBiS are available using three different technologies namely, PHP, Java Servlets and Enterprise Java Beans (EJB) [1]. Our testbed used the PHP version, which is installed as a Fast CGI module on either the Lighttpd (1.4.26) [28] or the Apache (2.0.63) [29] Web server. Lighttpd is a single process, event-driven, open-source Web server. The Apache Web server provides a group of Multi Processing Modules (MPMs) that allow it to run in a process-based, hybrid (process and thread) or event-hybrid mode. Only one of the Lighttpd or Apache Web servers runs at a time, as required by the specific experiment being run on the testbed. The MySQL database server is installed as the database tier. In our experiments, the Web server machine's CPUs were utilized more than the database server machine's CPUs.

Initially, the RUBiS client workload generator (version 1.4.3, released October 2004)[‡] [1] was used to exercise the system. However, as part of our investigation into why the benchmark results were not consistent across runs, we also used the *httperf* workload generator [6]. Switching to a different workload generator provided insights on causes of the inconsistencies, which we report on in section 4. To support tests with a large number of concurrent users, we followed the tuning procedure outlined by Brecht [31] that allowed the client machine to maintain a large number of open file descriptors.

3.2. Monitoring tools

To investigate the causes of the inconsistent benchmark results, we required three sets of data to be monitored. First, we needed the actual client-perceived response times for HTTP requests. In our study, we compare these response times to those reported by RUBiS client. Second, we wanted to record the TCP connection establishment and termination events in the system. Last, we needed the resource utilization on the client, Web server and database server.

To verify the HTTP response times as seen by the client, we required an independent monitoring tool, which could accurately measure the response times. We achieved this using *tcpdump* [32],

[‡]We note that the version of RUBiS we used corrected the bug in the previous version detected by Guitart *et al.* that can erroneously reduce the request generation rate [30].

Table I. Machine specifications.

Properties	Values
Number of processors	1
Number of cores	2
Processor model	Intel Core2 CPU 6400 @ 2.13GHz
Processor cache size L1	64 KB
Processor cache size L2	2048 KB
Memory total	2 GB
Memory swap	2 GB
Kernel	Linux version 2.6.9-55.0.9

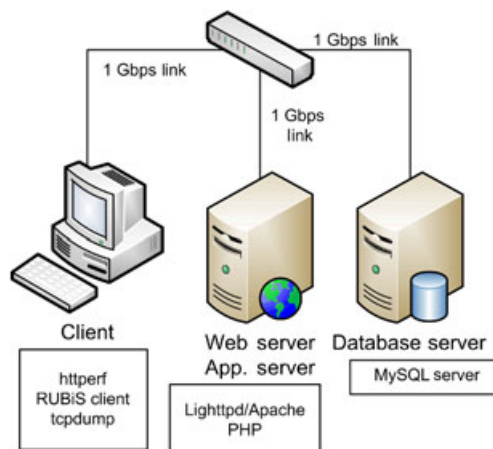


Figure 1. Testbed configuration.

Bro [33] and a Bro script (downloaded from: <http://bro-ids.org/bro-contrib/network-analysis/akm-imc05/>) that was developed for TCP delay analysis [34]. The script (*reduce.bro*) was modified (renamed as *response.bro*) to calculate the required timings in our experiments. In addition to the actual response time of the Web server, we used these tools to monitor the TCP connection establishment and termination events during each experiment. To obtain these, we modified the “*reduce.bro*” script (renamed as *connection.bro*) so that it tracks the TCP traffic between the client and the server.

To better understand how we use Bro, consider a typical HTTP request between a client and Web server, as shown in Figure 2. The first three packets construct the three-way handshake for establishing a TCP connection. The client sends the connection request to the server in the SYN packet (P1). The servers respond to the connection request by sending the SYN-ACK packet (P2). Then the client sends the ACK packet (P3) as a confirmation to the server, and the connection is established. Once the connection is created, the client sends the HTTP request (P4) to the Web server. The server then confirms the arrival of the TCP packet carrying the HTTP request by sending the acknowledgment (P5). After processing the HTTP request, the HTTP reply is sent in one or more TCP packets, depending on the size of the requested object. In Figure 2, the server sends the reply in three packets (P6, P7 and P8). Once the client receives all of the reply packets, it acknowledges receipt of the last packet, and the HTTP request is complete (P9). If the TCP connection is persistent, it can be used for additional HTTP requests. Any additional HTTP requests on this connection avoid the establishment overhead (i.e., the connection establishment time) as well as associated overheads (e.g., TCP slow-start). The HTTP 1.1 protocol allows persistent TCP connections to be exploited [35].

The time to complete the entire HTTP request can be broken into two parts, as shown in Figure 2. The first part is the difference between the time the first reply packet (P6) is received at the client and the time at which the client issued the HTTP request (P4). We refer to this time as the *reply time*. The second interval is the duration of time for the entire reply to be transferred to the client. We call

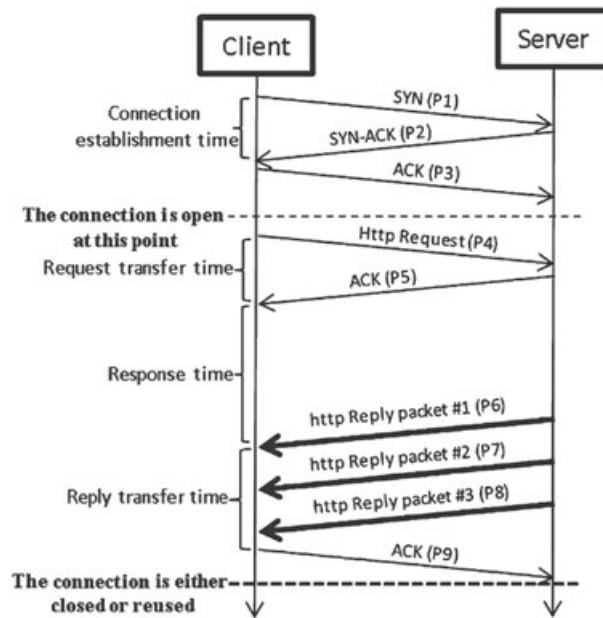


Figure 2. Typical HTTP communication.

this interval *transfer time*. We refer to the sum of the reply time and transfer time as the *response time*, which we use for evaluating the accuracy and representativeness of the workload generators.

Before each experiment, tcpdump is started and configured to capture all HTTP packets (TCP port 80). The tcpdump command line options we used to capture the network traffic are as follows:

- **tcpdump** ip host xxx.xxx.xxx.xxx and tcp port 80 -w *outputfile*
 - The “ip host” directive specifies that only IP packets “from” or “to” the specified xxx.xxx.xxx.xxx host address should be captured. For our tests “xxx.xxx.xxx.xxx” can refer to the IP address of either the Web server or the client machine.
 - The “-w” option writes the raw packets to *outputfile* rather than parsing and printing them to the console.

When each experiment ends, tcpdump is stopped, and its output file is given as the input to Bro. The primary purpose of Bro is as an intrusion detection system, but its powerful analysis capabilities make it attractive for our purposes. The modified Bro script (*response.bro*) calculates the reply times and transfer times of HTTP requests and records them in its own output file. Although Bro can run directly on live network traffic, we run it in an off-line manner to reduce the overhead on the client.

We describe the Bro command line options to run the *response.bro* script as follows:

- **bro** -r *outputfile* response.bro > *response.csv* 2 > *response.err*
 - The “-r” option specifies the name of the input tcpdump (pcap) format trace file to read and analyze.
 - “*response.bro*” is the name of Bro script that specifies the analysis to conduct.
 - “*response.csv*” is an output file that contains the HTTP response times calculated by Bro for each HTTP request submitted during the test.
 - “*response.err*” is an output file that captures any warnings generated by Bro.

The fields recorded by the output file (*response.csv*) of the modified Bro script are as follows:

- **response.csv** : “*End of transaction time stamp, Port, Connection ID, Request #, Pipelined flag, Reply time, Response time*”
 - *End of transaction time stamp*: The time stamp, which is recorded when the HTTP request is completed.
 - *Port*: The TCP port used by client to communicate with the server.

- *Connection ID*: An identifier for the TCP connection used to send the request.
- *Request #*: Number of HTTP requests, which have been sent through this TCP connection before current request.
- *Pipelined[§] flag*: Specifies whether the request pipelining used in this connection (1) or not (0).
- *Reply time[¶]*: The difference between the time at which the client issued the HTTP request and the time when first reply packet is received at the client.
- *Response time*: The sum of reply time and transfer time.

The *connection.bro* script can be executed similar to how the *response.bro* script is invoked. This script outputs a file called *connection.csv* that includes the following information:

- **connection.csv**: “*Event time stamp, Event flag, Number of open connections*”
 - *Event time stamp*: The time stamp, which is recorded at the completion of an event (TCP connection establishment and connection close).
 - *Event flag*: If set to “1” the event was “Connection Establishment”, and if set to “0” it was a “Connection Termination” event. The total number of connections established during the experiment can be calculated by counting number of “1” values in this column.
 - *Number of open connections*: The number of established connections, which are open at this instant.

There are several alternatives for executing tcpdump. For example, tcpdump can be run on a dedicated machine. That machine could receive a mirrored copy of the traffic from the Ethernet switch (e.g., from a switch that supports port mirroring). The main advantage of this approach is that there is no overhead on the client machine. A disadvantage is that it would have slightly less accurate measurements of when the client received each packet. For our study, we ran tcpdump on the client machine during our experiments. This enables us to get the best estimate of the response times as experienced by the client. However, it does place some load on the resources of the client machine. To ensure that the overhead of tcpdump does not affect the test results, a set of simple tests were conducted and repeated with and without tcpdump. We found that for our tests, tcpdump did not affect the accuracy of measured response times. However, tcpdump could affect measured response times in some situations (e.g., high network utilization), so care must be taken when using it.

Because the server and client machines used the Linux operating system, the *sysstat* [36] package was used to monitor the resource utilizations on each machine. The *sysstat* package contains the *sar*, *sadf*, *iostat*, *pidstat* and *mpstat* commands for Linux. The *sar* command collects and reports system activity information. The information collected by *sar* can be saved in a file in binary format for future inspection. We used the CPU, memory, network interface, and swap space metrics supported by the package to monitor how busy each machine was. The sampling interval was set to 1s for all the performance counters used by this study.

3.3. Workload generator modifications

As mentioned earlier, RUBiS client was the main target of this study. We also used httperf to independently validate results reported by RUBiS client. Although RUBiS client and httperf both emulate users’ HTTP transactions with a Web application, they have different capabilities and support different workload specifications. Therefore, we needed to make some minor changes in both applications so that their results would be directly comparable. In the following sections, the modifications made to each workload generator are discussed.

3.3.1. RUBiS client modifications. RUBiS client is a Java-based tool, which emulates user sessions using Java threads. The test specifications are defined in the “*rubis.properties*” file. We made the following four changes to RUBiS client:

[§]According to the HTTP protocol specification: “A client that supports persistent connections MAY “pipeline” its requests (i.e., send multiple requests without waiting for each response). A server MUST send its responses to those requests in the same order that the requests were received.” [32]

[¶]The response times in “response.bro” are measured with an accuracy of 10 microseconds.

1. The original RUBiS client application supports only closed session-based workloads. Because we wanted to have the option of submitting a more realistic workload to our Web server, we added the capability of creating sessions in a hybrid manner to RUBiS client. The inter-arrival times between successive new sessions are specified in a text file, and the file path is added to the *rubis.properties* file. Alternatively, our modifications also allow a tester to specify exponentially distributed session inter arrival times with a specified mean.
2. For our experiments, we required a record of each HTTP request and its response time (measured by RUBiS client), as well as a mapping of the request to the session that generated it. Because the original RUBiS client did not provide this feature, we customized the code to record this information in a log file. This file also records when each HTTP request gets submitted by RUBiS.
3. The original RUBiS client measures the response time in milliseconds. However, for some of the requests in our testbed, the response time was less than one millisecond. Consequently, the RUBiS output statistics were inaccurate in low load conditions, that is, many response times were reported as 0 ms. To calculate response time more accurately, we used *System.nanoTime()* [37] instead of the *System.currentTimeMillis()* [38] Java function. As the name implies, *System.nanoTime()* returns timestamps with nanosecond precision.
4. To simplify the experiments and knowing the fact that in most real systems a separate server is used as the “image server”, we disabled the image downloading part of the RUBiS code.

It is important to note that none of these modifications affect the general performance of the RUBiS client, whereas change #3 improves the accuracy of its response time measurements.

3.3.2. *httpperf* Modifications. *httpperf* is an event-based, single-threaded tool, which is capable of emulating different types of workloads. With *httpperf*, test parameters are specified as command line options. An example *httpperf* invocation is as follows:

- **httpperf** – server=www.example.com –port=80 –wssesslog 2000,0,*session_file* –period=*d*, 0.001

In this command, *httpperf* submits requests to “*www.example.com*” at TCP port “80”. The “*wssesslog*” option indicates that the workload to be generated is session-based. We have generated 2,000 sessions as per the session definitions found in the “*session_file*” text file. Each session definition in this file is a list of the URIs to be requested from the Web server with successive requests in the session definition separated by a think time. The “*period*” option specifies properties of the inter arrival time between successive new sessions generated by *httpperf*. The first parameter specifies the inter-arrival time distribution. The other parameter in the “*period*” option specifies the mean value for the session inter-arrival time distributions [39]. In the example shown previously, the inter-arrival times are set to a deterministic, that is, constant, value by selecting the “*d*” option, and this constant value is set to 1 ms. Closed workloads can be realized with *httpperf* by specifying a very low mean session inter-arrival time and a large number of requests within each session.

For each experiment in this study, we conduct a test with RUBiS client followed by a test with *httpperf*. To facilitate comparison, we require *httpperf* to submit the same workload to the system under test as was submitted by RUBiS client. In particular, the sequence of sessions submitted by RUBiS client and *httpperf* needs to be identical. This is achieved by extracting the session sequence and think time sequence within each session from the log file generated from the RUBiS client test (modification 2 of section 3.3.1) and saving that information to a “*session_file*” for use with *httpperf*. Furthermore, the sequence of inter-arrival times between successive new sessions generated by RUBiS client needs to be preserved during the *httpperf* test. This requirement caused us to modify *httpperf* to enhance the “*period*” option. The modification allows a user to specify a sequence of inter-arrival time values as input to *httpperf*. In the following example, *httpperf* uses a new “*s*” switch with the “*period*” option to submit sessions as per the sequence of inter-arrival times specified in the “*inter_arrival_file*”:

- –period=*s,inter_arrival_file*

In our experiments, we first extract the sequence of session inter-arrival times from a RUBiS test by parsing the log file generated by RUBiS client. We save this sequence in an “*inter_arrival_file*” and specify that file as input with the new “*period*” option. These modifications ensure that the sequence of sessions and the instants at which new sessions are generated (relative to the start of the test) are identical in the RUBiS client and httpperf tests in an experiment. Furthermore, a session generated by httpperf at a given time instant observes the same sequence of think times as the session generated by RUBiS client at the same time instant. To ensure that the Web server serves the requests from RUBiS client and httpperf in the same manner, we verified that the request header fields generated by RUBiS client and httpperf are identical.

We have also added a new option to httpperf for logging detailed information for individual HTTP requests submitted in a test. Specifically, the new “*-rfile_name=file_name*” option can be used to save the detailed information to the text file specified by “*file_name*”. The information recorded includes the time instant at which a request was submitted as well as the reply time and transfer time recorded by httpperf for that request. The values are recorded with 1 μ s precision.

4. INVESTIGATING WEB WORKLOAD GENERATION CHALLENGES

As mentioned in section 1, the goal of this study is to understand some inconsistencies observed in an earlier RUBiS benchmarking study. We started our exploration by attempting to validate the accuracy and correctness of the experiment results reported by RUBiS client. To achieve this, we needed to measure the client-observed response times independent of the RUBiS client. We did this by monitoring the network traffic and extracting the timing information of HTTP transactions, using the tools and methods described in section 3.2. In addition to verifying the accuracy of the response time measurements, we also wanted to verify whether RUBiS client emulates users in the expected manner. As mentioned in the previous section, we achieved this by conducting a side-by-side comparison between RUBiS client and a second workload generator (httpperf). The results of our validation are presented in section 4.1. We note that this general methodology can also be applied to validate other workload generators.

From our side-by-side comparison coupled with the independent response time monitoring, we observed some inaccuracies with RUBiS client. This motivated further investigation to determine the root causes of these problems. Because our initial investigation revealed that the problems occurred only with RUBiS client and not with the second workload generator, this second step is specific to RUBiS client. First, we examined external factors that could influence RUBiS client’s behavior. In particular, section 4.2 characterizes the impact of the JVM used by RUBiS client on the accuracy and scalability of the workload generator. Second, we considered internal factors (i.e., issues with RUBiS client’s implementation). Section 4.3 describes how the unrealistic TCP connection management policy used within RUBiS client can cause the workload generator to provide misleading insights into the scalability of the system under study. In general, performance debugging exercises tailored to the system under study (in this case, the workload generator) need to be carried out if problems are discovered in the validation phase.

4.1. Validating accuracy of measured end-user response times

To validate the accuracy of response times reported by RUBiS client, we conducted a set of experiments using the Lighttpd Web server. Because the original RUBiS client is only capable of generating closed workloads, we use a closed workload initially. The workload causes a mix of browse, buy, and sell transactions to be submitted to the RUBiS application. We set the think time distribution to a negative exponential [40] with a mean value of 7 s, consistent with the RUBiS specifications. The variable factor for these experiments is the number of concurrent user sessions in the system. We vary this parameter to achieve different utilization levels for the bottleneck resource in the server machines as well as the client machine. As described previously, in each experiment, we first conduct a test with RUBiS client. We then use httpperf to submit the same workload generated by RUBiS client using the approach described in section 3.3.2. We record the response times reported by each of the two workload generators and the actual Web server response times measured by Bro.

We define the difference between the mean response time reported by a workload generator and the actual mean response time reported by Bro as the *absolute error* of the workload generator.

We performed seven tests with each of the workload generators, varying the number of concurrent user sessions N from 500 to 4,000. Beyond 4,000 user sessions, the RUBiS client reported “out of memory” exceptions. This indicates that the JVM could not obtain enough memory from the operating system to spawn the requisite number of Java threads for user emulation.

Figures 3(a) and 3(b) show the absolute error of RUBiS client and httpperf, respectively, as a function of the Web server’s CPU utilization. It must be noted that Figure 3(a) and Figure 3(b) have very different scales for the Y -axis. Furthermore, both workload generators cause approximately the same utilization on the Web server’s CPUs, for a given value of N . This behavior is expected because both tools submit the same workload to the server.

Because Bro’s measurements are based on packet traces, we expect its response times to be slightly lower than what RUBiS client would see higher up the TCP/IP stack. However, the results in Figure 3(a) show that the increases are much larger than we would have anticipated. Even for light workloads ($N = 500$, mean Web server CPU utilization of 2.5%) the RUBiS client reports the mean response time 500 μ s higher than what Bro reports. With this absolute error, the RUBiS client mean response time estimate is 1.07 times the actual Bro measured mean response time. As N increases, the discrepancy increases significantly. For example at $N = 4,000$ the Web server CPU utilization is 22%, and the difference between the mean response times reported by RUBiS and Bro is around 150 ms. For this case, the RUBiS client mean response time is 2.6 times the Bro response time. In contrast, the absolute error values for httpperf do not change much with N and vary between 10 to 40 μ s. Most of this discrepancy is likely caused by Bro measuring in the lower layers of network protocol stack. The absolute error in the first experiment ($N = 500$) is more than 15 times higher in RUBiS client compared with httpperf. The gap between RUBiS client and httpperf in terms of error increases considerably as the workload intensity increases, and the Web server is more utilized. The httpperf results reveal that the RUBiS client errors are not caused by a common bottleneck like the network or the server. We note that the maximum mean Web server CPU utilization observed during the tests is only 22.2% at $N = 4,000$.

To better understand how often the discrepancies occur, we plot the cumulative distribution function (CDF) of the response times. Figures 4(a) and 4(b) show the CDF plots for both httpperf and RUBiS client tests for the $N = 3,000$ case. Figure 4(a) shows that there are significant discrepancies between the Bro measured response times (“Actual”) and RUBiS client reported response times (“Reported”) in both the high and low response time ranges. In contrast, Figure 4(b) shows that the two CDFs are almost indistinguishable for httpperf. We therefore conclude that RUBiS client significantly overestimates server response times whereas such a problem is not evident with httpperf.

We now investigate possible causes for the inaccurate RUBiS client response time measurements. Figure 5 shows the client machine’s mean CPU utilization when running the RUBiS client tests for different numbers of concurrent user sessions. The plot shows that the client machine’s CPU utilization increases rapidly with an increase in the number of concurrent user sessions. The figure

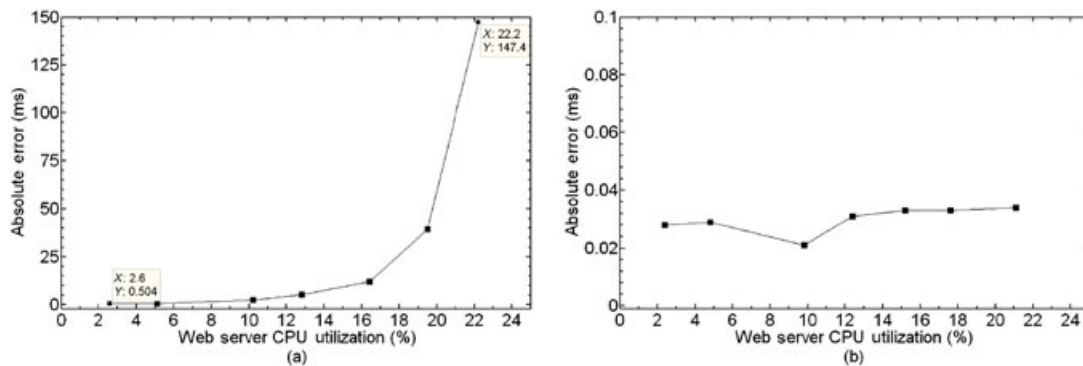


Figure 3. Absolute error (a) RUBiS client–GNU JVM (b) httpperf.

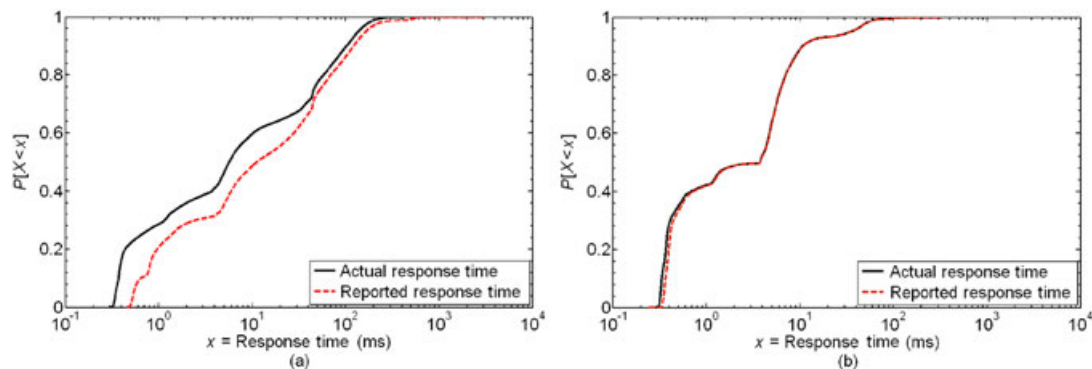


Figure 4. CDF of response times ($N = 3000$) - (a) RUBiS client-GNU JVM (b) httpperf.

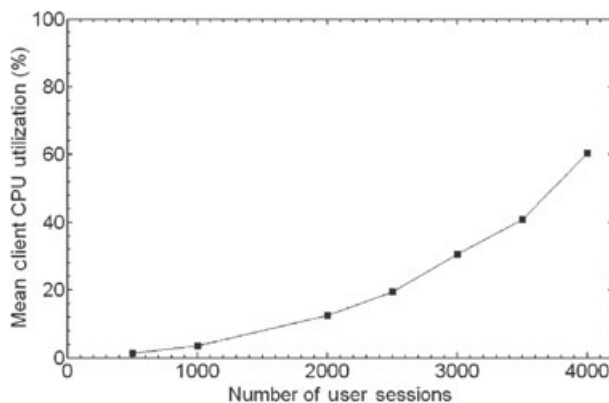


Figure 5. Mean client machine CPU utilization for RUBiS client-GNU JVM.

indicates that the client CPU utilization is strongly related to the number of threads created to emulate users. Figures 3 and 5 also jointly indicate that the CPU of the client machine is the bottleneck in the experimental setup. For example, the mean utilization of the Web server CPU with $N = 4,000$ is 22.2% (Figure 3(a)) whereas the mean utilization of the client machine's CPU at this setting is 60% (Figure 5). Therefore, it is likely that this bottleneck in the client side is the main source of the large absolute error values in the RUBiS client results.

4.2. Effect of Java virtual machine selection and tuning on the accuracy and scalability of RUBiS client

Because the performance results reported by RUBiS client were noticeably different from both httpperf and the independent measurement, our next step is to determine the cause(s). We begin by considering factors *external* to RUBiS client. The primary difference between the RUBiS client and httpperf is that the former is written in Java, whereas the latter is written in C. This motivates an investigation of the Java environment, to understand what effect (if any) it has on the response times reported by RUBiS client. Our approach to studying this is to use several different JVMs and configurations and to observe if any changes in behavior occur.

Although the effects of JVM selection and tuning are well-known in the Java community [41–44], these topics have not received much attention with respect to their use in Web workload generation. For example, we did not find any recommendations related to JVM selection and tuning in the RUBiS documentation. In the remainder of this section, we briefly explore their implications on Web workload generation.

The RUBiS client experiments described in section 4.1 used the default GNU JVM (GNU libgcj version 4.1.2 released in 2007) [45], which is bundled with some Linux distributions. Research in

the Java community suggested using a more recent JVM instead. Hence, we selected the latest version of Sun's JVM for Linux (version 1.6-18, released February 2010). A key difference between the two JVMs is that the Sun JVM includes the "HotSpot" [46] technology designed to improve performance.

Figure 6 compares the actual response times measured by Bro for the two JVMs. Figure 6 shows that the Sun HotSpot JVM allowed us to emulate more than the 4,000 users possible with the GNU JVM. Furthermore, we noticed that the accuracy of RUBiS client improves when using Sun HotSpot JVM. For RUBiS client using the Sun HotSpot JVM, the absolute error was less than 3 ms up to $N = 4,800$. We verified that the ability to support more concurrent users and the improved accuracy are caused by the better performance of the Sun HotSpot JVM. In particular, the Sun HotSpot JVM utilizes the client machine's CPUs considerably less than the GNU JVM. For instance, whereas the CPU in the client machine is on average 30% utilized for $N = 3,000$ in the GNU JVM, this value is decreased to 5% with Sun HotSpot JVM.

With the default configuration of the Sun HotSpot JVM, we were able to emulate up to 5,100 users. Beyond this, an out-of-memory exception was encountered by RUBiS client. Following known best practices for JVM tuning, we used the Java command line option "`-Xmx=heapsize`" to increase the maximum heap size for the JVM from 64 MB^{||} to 512MB. As shown in Figure 6, the new heap size increased the capacity of RUBiS client from 5,100 to 6,300 users. Beyond 6,300 threads, we encountered the out-of-memory exception again.

We conclude this section by comparing the scalability of RUBiS client and `httperf` based on the results presented so far. Figure 6 plots the actual mean response times for various values of N for RUBiS client and `httperf`. The maximum number of users that RUBiS client could emulate on the client node was 6,300 with the Sun HotSpot JVM and 512 MB heap size setting. At this setting, the mean utilization of the Web server CPUs was found to be 36%. A memory bottleneck at the client machine prevents RUBiS client from stressing the server further (i.e., with a larger JVM heap size). In contrast, we were able to conduct tests with up to $N = 8,000$ concurrent users with `httperf` as shown in Figure 6. At this setting, `httperf` was able to drive the utilization of the Web server CPU to up to 90%. Furthermore, the mean response time reported by `httperf` was still very accurate for this case. For more than 8,000 concurrent users, the server was severely loaded leading to unstable behavior such as very long request response times and a large number of connection resets and timeouts.

The above results indicate that multi-threading does not automatically ensure a scalable workload generator. We note that `httperf` was configured to exploit only one of the two available processor cores of the client machine in our experiments. In contrast, RUBiS client used both cores. This suggests that the single-threaded, event-driven `httperf` tool can support significantly more number of emulated users on a single host than RUBiS client.

4.3. Effect of Transmission Control Protocol connection management policy on the validity of performance results

In addition to considering external factors that may influence RUBiS client's behavior, we also examined internal factors; that is, those related to RUBiS client's implementation. We examined the source code of RUBiS client, to identify how it generated the workload. This led us to investigate the effect of the TCP connection management policy on the performance results reported by each workload generator. For this comparison, we consider only the Web server response times measured by Bro, to facilitate a direct comparison between the workload generators. Figure 6 shows that the RUBiS client with GNU JVM, RUBiS client with Sun HotSpot JVM, and `httperf` yield very different mean server response times, even though they have been configured to submit the same workload. For example, with $N = 3,000$ users the mean server response times are 31.6, 8.3, and 4.2 ms, respectively, with RUBiS client (GNU JVM), RUBiS client (Sun HotSpot JVM), and `httperf`. The remainder of this section explains the reason for these differences and their implications to server performance evaluation.

^{||}This setting pertains to the JVM optimized for "client" code.

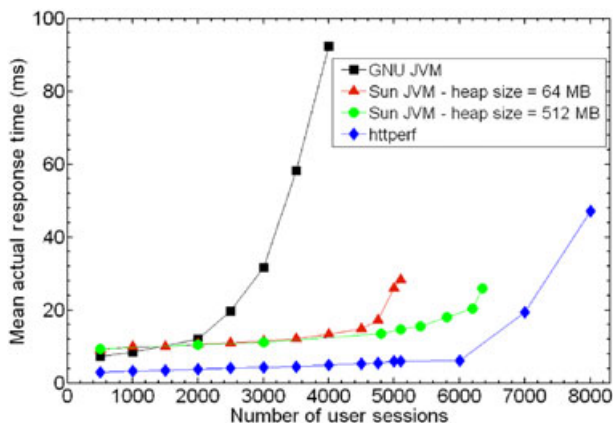


Figure 6. Mean actual response time of Web server for all workload generators.

After a detailed examination of the behaviors of RUBiS client and httperf, we realized that the primary difference between them was caused by differences in the TCP connection management policy used with each application. On paper, both claimed to open a single (dedicated) persistent TCP connection per user session (see RUBiS documentation [47] and httperf man page). However, our test results revealed significant differences in the way httperf and RUBiS client handle persistent connections. Although we found that httperf exhibits its documented behavior, we observed that RUBiS client handles the TCP connections quite differently. These differences are caused by its use of certain functions in the Java network library. Specifically, a thread first creates a *URL* object, which then calls the *URL.openStream()* [48] Java function. This function returns an *InputStream* object whose methods are called by the thread to read the HTTP response pertaining to the *URL*. After reading the response, the thread closes the *InputStream* object by calling its *close()* method. Java documentation states that this method releases any system resources associated with the stream [49]. The ensuing results show that the use of these functions causes RUBiS client to deviate from its intended behavior of using one dedicated TCP connection per user session.

To illustrate the differences between the TCP connection management policies, we present in Figure 7 a time series of the number of concurrent connections observed during a test with both httperf and RUBiS client for one of the experiments presented in the earlier sections ($N = 2,000$). This data were obtained by using the *connection.bro* script described in section 3.2. In the RUBiS client tests (Figures 7(a) and 7(b)), the total number of concurrent connections fluctuates from 0 to 75 for the GNU JVM and 0 to 250 for the Sun HotSpot JVM. These maximum values are much less than the number of concurrent user sessions ($N = 2,000$). In contrast, the httperf test (Figure 7(c)) sees the number of concurrent TCP connections jump from 0 to 2,000 at the beginning of the experiment**. After around 300 s, the number of concurrent connections starts to decrease as sessions begin to complete.

Figure 8 provides further evidence of the differences in the way persistent connections are handled by RUBiS client and httperf. Figures 8(a) and 8(b) show the number of requests sent through each unique TCP connection used during RUBiS client and httperf tests, respectively. The *x* axis of these figures show connections sorted based on the time they were closed. Figures 8(a) and 8(b) show that both workload generators used persistent TCP connections because each connection was used to issue multiple requests. With httperf, 2,000 connections were opened during the test, and each connection submitted 50 requests as shown in Figure 8(b). The number of requests submitted per connection corresponds to the number of requests per session, which was configured to be 50 for both workloads. In contrast, with RUBiS client the number of connections used during the test was 2,912, which is greater than the number of concurrent sessions. This indicates that there are more connection establishment and connection shutdown activities in the RUBiS client workload. From

**Because this example uses a closed workload, all the sessions start at the beginning of the experiment.

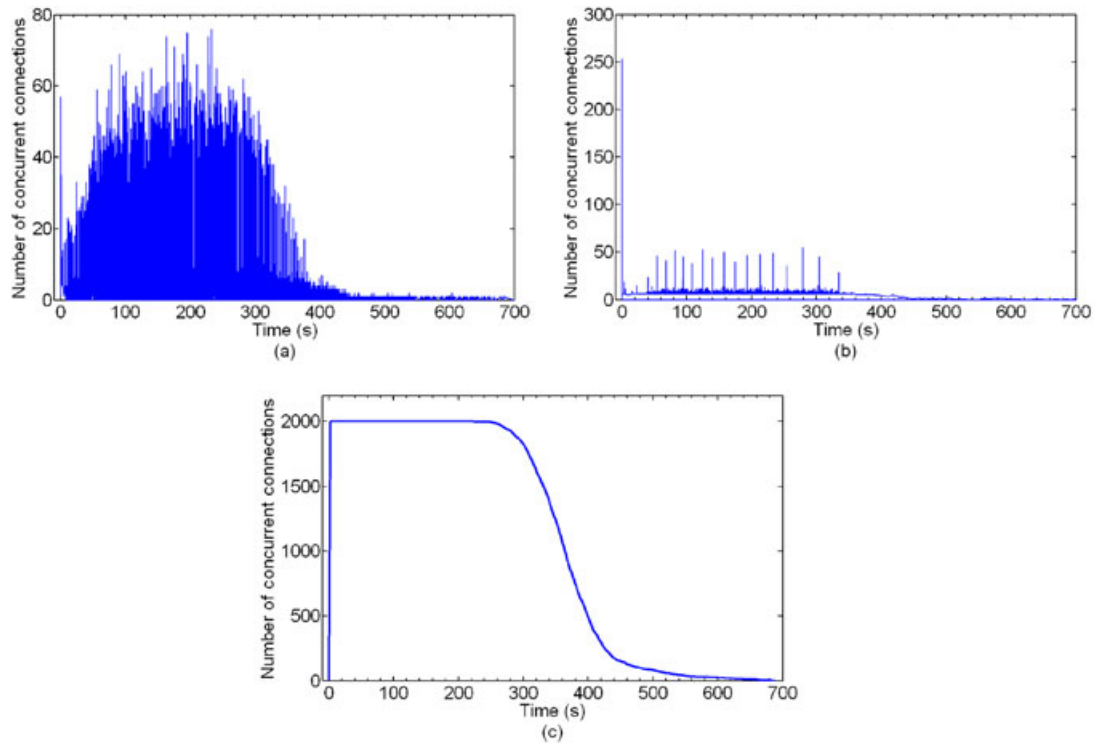


Figure 7. Number of concurrent connections - (a) RUBiS client-GNU JVM (b) RUBiS client-Sun JVM (c) httpperf.

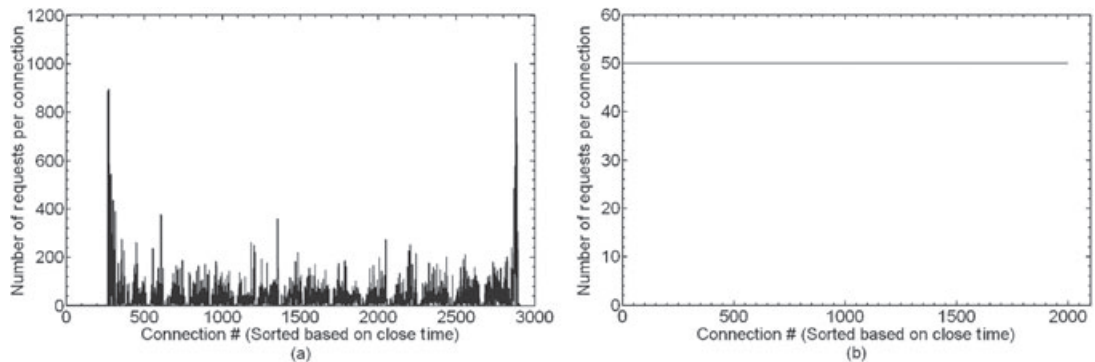


Figure 8. Number of requests per connection ($N = 2,000$) - (a) RUBiS client-Sun JVM (b) httpperf.

Figure 8(a), at the beginning of the RUBiS experiment a large number of connections (around 250) are opened to send the first requests for the 2,000 users. These connections are closed after one or two HTTP requests were submitted through them. However, the number of requests submitted per connection for a vast majority of connections is significantly greater than 50, the number of requests per session. The number of requests submitted in a connection was as high as 1,000. Figure 7 and Figure 8 together establish that RUBiS client causes multiple user sessions to use a single connection. In effect, a small number of concurrent TCP connections are shared across a large number of emulated users.

The higher server response times observed with RUBiS client in Figure 6 are likely caused by the higher connection establishment and connection shutdown overheads in RUBiS client. Specifically, the highest response time in Figure 6 was caused by GNU JVM, which established the most connections. The Sun HotSpot JVM established fewer connections, which caused requests to incur

lower response times at the server. httpperf established the fewest TCP connections, causing the least stress on the server.

The sharing of TCP connections across users in RUBiS client is not representative of the behavior observed in real systems, because clients do not initiate requests from the same TCP connection. We conducted additional experiments to better understand how the connection management policies used by the tools impact server behavior. Specifically, we constructed hybrid workloads to emulate a “flash crowd” scenario, that is, a sudden increase in the rate of arrival of sessions, using httpperf and RUBiS client. We then studied the behavior of the Lighttpd and Apache servers under these workloads. We configured Apache to use the “prefork” module [50] for request processing. Recalling from section 3, Lighttpd is an event-based server that uses an asynchronous mechanism to handle HTTP requests. It is lightweight in that it is designed to use just a single process per processor. For all our experiments, we configured Lighttpd to use two processes [51]. The prefork module of Apache maintains a pool of worker processes with each process handling an incoming connection in a synchronous manner. When the number of incoming connections exceeds the number of worker processes, Apache spawns additional processes to handle the increased load. For our experiments we used the default *Apache* setting where the initial size of the worker process pool is 15.

Figure 9(a) shows the non-bursty and the bursty, that is, flash crowd workloads used for Lighttpd. From Figure 9(a), in the non-bursty workload, the sessions arrive with a mean rate of 10 sessions per second. In the bursty workload, the session arrival rate increases suddenly from 10 to 333 causing a burst of sessions over a 4-s time interval. Table II shows the mean response times measured by Bro for both workloads while using httpperf and RUBiS client. Under both the workload generators, Lighttpd is able to handle the bursty workload with only a marginal increase in mean response time. However, the mean response times with httpperf are lower than the corresponding response times with RUBiS client. This is consistent with the behavior observed in the previous experiments, which also used Lighttpd. The extra connection establishment and connection disconnection overheads imposed by RUBiS client seem to dominate this scenario.

Figure 9(b) shows the non-bursty and bursty workloads used for Apache. Because the Apache and Lighttpd servers have dissimilar overheads for handling requests from new TCP connections,

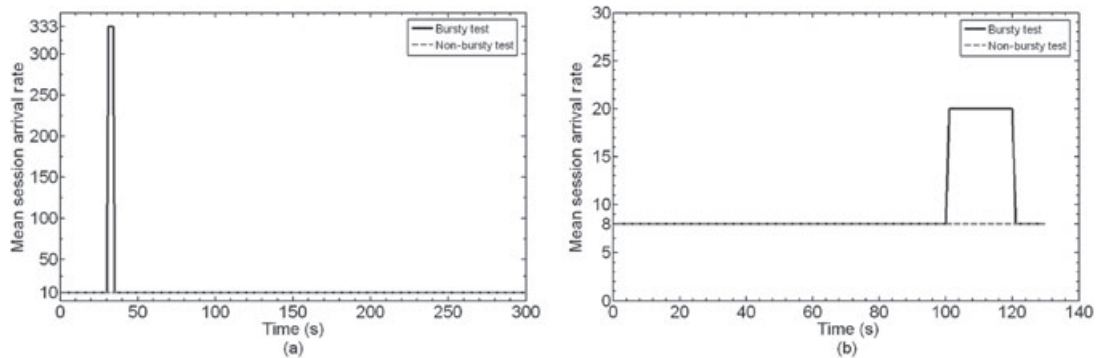


Figure 9. Session arrival rate (Sessions/Second) - (a) Lighttpd (b) Apache.

Table II. Bursty test results.

		Apache response time (ms)			Lighttpd response time (ms)		
		Median	Mean	95 percentile	Median	Mean	95 percentile
RUBiS client	<i>Non-bursty</i>	2.118	4.333	23.334	2.243	11.003	44.995
	<i>Bursty</i>	2.135	4.571	23.755	2.411	12.369	48.384
httpperf	<i>Non-bursty</i>	2.164	11.537	27.878	2.150	4.061	20.588
	<i>Bursty</i>	2.233	58.647	33.457	2.165	4.643	20.488

we had to apply a workload with a different burst specification. Figure 9(b) shows that in the non-bursty workload, sessions arrive at the rate of eight sessions per second. In the bursty workload, the session arrival rate suddenly increases from eight to 20 causing a burst of sessions to arrive over a 20-s time interval.

From Table II, data obtained from RUBiS client tests on Apache show that there is very little increase in the mean, median, and 95 percentile of response times from the non-bursty to bursty case. This may lead a tester to reach the conclusion that the Apache Web server is scalable with respect to handling bursts in session arrivals. However, the use of *httperf* provides a diametrically opposing viewpoint regarding the server's scalability. While using *httperf*, the mean response time for the bursty workload is more than five times the mean response time of the non-bursty workload. The 95th percentile of response time is also higher for the bursty workload whereas the median response time is not affected much. These observations imply that a small fraction of requests in the bursty workload encountered very high response times. The results show that burstiness has a significant detrimental impact on Apache's performance. Interestingly, whereas RUBiS client places lesser stress than *httperf* on Apache, the situation is reversed for *Lighttpd*. This suggests that a factor other than connection establishment and connection shutdown overheads dominates server performance in this scenario.

The reason for the long response times encountered while using *httperf* to generate a bursty workload to Apache can be explained as follows. To generate the burst of sessions shown in Figure 9(b), *httperf* initiates a large number of concurrent connections, one per each new session in the burst, to the server. This causes the total number of concurrent connections issued by *httperf* to Apache to increase beyond the worker process pool size. As a result, a number of sessions in the burst encounter significant delays related to the time Apache takes to spawn new worker processes to handle the increase in load. From Table II, this type of performance degradation was also observed with the non-bursty workload, although to a smaller extent. RUBiS client does not cause the bottleneck in worker process pool size to be exposed caused by its sharing of connections across multiple users. The maximum number of concurrent connections issued by RUBiS client during the test was always less than the Apache worker process pool size. Table II reveals that *Lighttpd* is less sensitive to burstiness under both *httperf* and RUBiS client. Because of its event-driven, asynchronous request processing architecture, it avoids overheads related to spawning new processes to handle a burst of incoming connections.

In summary, the choice of unrealistic TCP connection management policies can provide misleading insights into Web server performance and scalability. For example, the policy of the Java library used by RUBiS client provided pessimistic estimates of performance for *Lighttpd* but overly optimistic estimates for Apache with respect to the more realistic one dedicated connection per session policy used by *httperf*. We note that whereas *httperf* uses a more realistic policy than that of the Java library used by RUBiS client, many modern browsers use more than one connection per session. For example, Souders reports that Web browsers like Internet Explorer 8 and Firefox 3 use up to six parallel TCP connections to transfer HTTP transactions [52]. Workload generators must consider such complexities to ensure the validity of performance testing exercises.

5. DISCUSSION

Benchmarking computer systems is a challenging task, as there are many possible mistakes that can be made. One common mistake listed by Jain is "not validating measurements" [40]. The solution to this is to cross-check the measurements, which is an approach we used in this work. By implementing such an approach based on network traffic monitoring, we revealed that the RUBiS client was incorrectly reporting the performance and scalability of the Web server under test. We then explored internal and external factors. These explorations revealed that the undesirable behaviors were caused by limitations of the Java networking library used by the tool and by the JVM originally used in the test bed.

Our specific implementation of the network monitoring based validation methodology is straightforward to apply in other Web server benchmarking studies. Whereas our approach can alert users

to possible problems with the workload generator, it cannot automatically identify the root causes of those problems. Resolving the problems contributing to response time inaccuracies requires manually searching for root causes, such as looking into workload generator specific configuration issues, for example, JVM tuning in our case, and going through the source code as we did to identify the cause of the connection management issue. It is important to note that the overhead of the validation approach should be quantified (as we did) in each case, to avoid another common benchmarking mistake [40].

Code reuse is a common practice in software development, as it can dramatically reduce the time (and therefore cost) to develop an application. However, a disadvantage of reusing source code is that any problems that exist with the initial code can propagate to other applications. This issue is relevant to our work, as the TPC-W workload generator shares a similar implementation to the RUBiS client generator. In particular, both are implemented in Java, support a closed workload approach, follow the multi-threaded paradigm, and employ the same TCP connection management policy (via a common Java library) described in section 4.3. As a result, we expect studies that have used either RUBiS or TPC-W to benchmark a Web server may have incorrectly estimated the performance or scalability of the server, for the reasons discussed in section 4. In contrast, another popular Java-based tool, Apache JMeter [8], does not use this library and hence is unlikely to suffer from the connection management related problems we observed in this study.

We believe that the need to validate workload generators becomes even more crucial with the advent and widespread use of Web 2.0. Cormode and Krishnamurthy [53] note that there are significant differences between Web 1.0 and Web 2.0 with respect to the technologies involved and the types of user interaction supported. The increased complexity of Web 2.0 implies greater challenges in ensuring the correctness and realism of workload generators targeted for such applications.

6. CONCLUSION

This paper described our experience in validating the performance and scalability results reported by the RUBiS client Web workload generator. After observing inconsistent benchmarking results with RUBiS client, we implemented a method to cross-check the results. This uncovered two root causes for the inconsistent results: the JVM and the Java network library used by the generator. We also showed that a multi-threaded workload generator is not necessarily more scalable than an efficiently implemented, event-based, single-threaded generator.

Because of the importance of Web workload generation, we believe that similar validation work should be conducted for other common workload generators. In particular, we plan a similar study for the SPECWeb workload generator, as it is commonly used in industry to benchmark new generation Web servers. Because the results of such studies are used for purposes such as purchase decisions, the ramifications for inaccurate benchmark results are potentially more significant.

Source code for the modified Bro monitoring script can be accessed at: <http://bro-ids.org/bro-contrib/network-analysis/hka-spe11/>. Our workload generator enhancements and further details on our experimentation methodology can be found at: <http://people.ualgary.ca/~dkrishna/SPE>

ACKNOWLEDGEMENTS

This work was financially supported by Natural Sciences and Engineering Research Council (NSERC) Canada and Hewlett-Packard (HP). The authors would like to thank Jerry Rolia of HP Labs and the anonymous reviewers for their helpful suggestions and constructive feedback.

REFERENCES

1. RUBiS–Homepage. <http://rubis.ow2.org/> [01 October 2010].
2. TPC–W– Homepage. <http://www.tpc.org/tpcw/> [01 October 2010].
3. SPEC– Benchmarks. <http://www.spec.org/benchmarks.html#web> [01 October 2010].
4. PHP: Hypertext Processor. <http://www.php.net/> [01 October 2010].

5. Enterprise Java Bean Technology. <http://www.oracle.com/technetwork/java/index-jsp-140203.html> [01 October 2010].
6. Mosberger D, Jin T. httpperf: A tool for measuring web server performance. *ACM SIGMETRICS Performance Evaluation Review* 1998; **26**(3):31–37.
7. Banga G, Druschel P. Measuring the capacity of a web server under realistic loads. *World Wide Web* 1999; **2**(1):69–83.
8. Apache JMeter. <http://jakarta.apache.org/jmeter/> [01 October 2010].
9. Feitelson D. The forgotten factor: facts on performance evaluation and its dependence on workloads. *Proceeding of Int. Euro-Par Conference 2002*; **2400**:49–60.
10. Paxson V, Floyd S. Wide area traffic: the failure of Poisson modeling. *IEEE/ACM Transactions on Networking* 1995; **3**(3):226–244.
11. Ferrari D. On the foundations of artificial workload design. *ACM SIGMETRICS Performance Evaluation Review* 1984; **12**(3):8–14.
12. Barford P, Crovella M. The surge traffic generator: generating representative web workloads for network and server performance evaluation. *Proceedings of the ACM SIGMETRICS* 1998:151–160.
13. Schroeder B, Wierman A, Harchol-Balter M. Open versus closed: a cautionary tale. *Proceedings of the 3rd Conference on Networked Systems Design & Implementation* 2006; **3**:18–18.
14. Kant K, Tewari V, Iyer R. GEIST: Generator of ecommerce and internet server traffic. *Proceedings of Int. Symposium on Performance Analysis of Systems and Software* 2001:49–56.
15. Liu Z, Niclausse N, Jalpa–Villanueva C. Traffic model and performance evaluation of web servers. *Performance Evaluation* 2001; **46**(2-3):77–100.
16. Krishnamurthy D, Rolia J, Majumdar S. A Synthetic Workload Generation Technique for Stress Testing Session-Based Systems. *Proceedings of the IEEE Transactions on Software Engineering* 2006; **32**(11):868–882.
17. Guitart J, Carrera D, Torres J, Ayguadé E, Labarta J. Tuning dynamic web applications using fine-grain analysis. *Proceedings of the 13th Euromicro Conference on Parallel, Distributed and Network-Based Processing* 2005:84–91.
18. Padala P, Zhu X, Wang Z, Singhal S, Shin K. Performance evaluation of virtualization technologies for server consolidation. *Technical Report HPL-2007-59* 2007.
19. Malkowski S, Hedwig M, Pu C. Experimental evaluation of N-tier systems: observation and analysis of multi-bottlenecks. *Proceedings of the IEEE International Symposium on Workload Characterization* 2009:118–127.
20. Sicard S, Boyer F, De Palma N. Using components for architecture-based management: the self-repair case. *Proceedings of the 30th International Conference on Software Engineering* 2008:101–110.
21. Guitart J, Carrera D, Beltran V, Torres J, Ayguadé E. Dynamic CPU provisioning for self-managed secure web applications in SMP hosting platforms. *Computer Networks* 2008; **52**(7):1390–1409.
22. Wood T, Cherkasova L, Ozonat K, Shenoy P. Profiling and modeling resource usage of virtualized applications. *Proceedings of the 9th ACM/IFIP/USENIX International Conference on Middleware* 2008:366–387. DOI: 10.1007/978-3-540-89856-6_19.
23. Rolia J, Krishnamurthy D, Casale G, Dawson S. BAP: a benchmark-driven algebraic method for the performance engineering of customized services. *Proceedings of the First Joint WOSP/SIPEW International Conference on Performance Engineering* 2010:3–14.
24. Kusic D, Kephart JO, Kandasamy N, Jiang G. Power and performance management of virtualized computing environments via lookahead control. *Cluster Computing* 2009; **12**(1):1–15.
25. Ramamurthy P, Sekar V, Akella A, Krishnamurthy B, Shaikh A. Remote profiling of resource constraints of web servers using mini-flash crowds. *USENIX Annual Technical Conference on Annual Technical Conference* 2008:185–198. DOI: 10.1.1.145.4897.
26. Nahum E. Deconstructing SPECweb99. *7th International Workshop on Web Content Caching and Distribution (WCW)* 2002.
27. Nahum E, Rosu M, Seshan S, Almeida J. The effects of wide-area conditions on WWW server performance. *Proceedings of the ACM SIGMETRICS International Conference on Measurement and Modeling of Computer Systems* 2001:257–267.
28. Lighttpd fly light. <http://www.lighttpd.net/> [01 October 2010].
29. The Apache HTTP Server Project. <http://httpd.apache.org/> [01 October 2010].
30. Guitart J, Carrera D, Torres J, Ayguadé E, Labarta J. Successful experiences tuning dynamic web applications using fine-grain analysis. Research report number: UPC-DAC-2004-3 / UPC-CEPBA-2004-2, 2004.
31. Brecht T. Linux: increasing the number of open file descriptors. Online tutorial available at: <http://www.cs.uwaterloo.ca/~brecht/servers/openfiles.html> [01 October 2010].
32. TCPDUMP/LIBPCAP public repository. <http://www.tcpdump.org/> [01 October 2010].
33. Paxson V. Bro: a system for detecting network intruders in real-time. *Computer Networks* 1999; **31**(23-24): 2435–2463. DOI: 10.1016/S1389-1286(99)00112-7.
34. Arlitt M, Krishnamurthy B, Mogul JC. Predicting short-transfer latency from TCP arcana: a trace-based validation. *Proceedings of the 5th ACM SIGCOMM Conference on Internet Measurement* 2005:19–19.
35. Hypertext Transfer Protocol – HTTP/1.1. <http://www.w3.org/Protocols/rfc2616/rfc2616.html> [01 October 2010].
36. SYSSTAT. <http://pagesperso-orange.fr/sebastien.godard/> [01 October 2010].
37. Java Documentation Link. <http://download.oracle.com/javase/1.5.0/docs/api/Java/lang/System.html#nanoTime%28%29> [01 October 2010].

38. Java Documentation Link. <http://download.oracle.com/javase/1.4.2/docs/api/java/lang/System.html#currentTimeMillis%28%29> [01 October 2010].
39. httpperf Manual. <http://www.hpl.hp.com/research/linux/httpperf/httpperf-man-0.9.pdf> [01 October 2010].
40. Jain R. *The Art of Computer Systems Performance Analysis*. John Wiley & Sons: New York, 1991.
41. Tong L, Lau FCM. Exploiting memory usage patterns to improve garbage collections in Java. *Proceedings of the 8th International Conference on the Principles and Practice of Programming in Java (PPPJ '10)*. ACM, New York, NY, USA, 2010; 39–48, DOI: <http://doi.acm.org/10.1145/1852761.1852768>.
42. Georges A, Buytaert D, Eeckhout L. Statistically rigorous Java performance evaluation. *Proceedings of the 22nd Annual ACM SIGPLAN Conference on Object-Oriented Programming Systems and Applications*. ACM, New York, NY, USA, 2007; 57–76, DOI: <http://doi.acm.org/10.1145/1297027.1297033>.
43. Brecht T, Arjomandi E, Li C, Pham H. Controlling garbage collection and heap growth to reduce the execution time of Java applications. *ACM Transactions on Programming Languages and Systems* 2006; **28**(5):908–941.
44. Blackburn SM, Cheng P, McKinley KS. Myths and realities: the performance impact of garbage collection. *Proceedings of the Joint International Conference on Measurement and Modeling of Computer Systems (SIGMETRICS '04/Performance '04)*. ACM, New York, NY, USA, 2004; 25–36.
45. GCJ: The GNU Compiler for Java-GNU Project - Free Software Foundation (FSF). <http://gcc.gnu.org/java/> [01 October 2010].
46. Oracle (Sun) Java HotSpot Technology. <http://www.oracle.com/technetwork/java/javase/tech/index-jsp-136373.html> [01 October 2010].
47. Amza C, Chanda A, Cox A, Elnikety S, Gil R, Rajamani K, Cecchet E, Marguerite J. Specification and implementation of dynamic Web site benchmarks. *Proceedings of WWC-5: IEEE 5th Annual Workshop on Workload Characterization 2002*:3–13.
48. Java Documentation Link. <http://Java.sun.com/j2se/1.4.2/docs/api/Java/net/URL.html#openStream%28%29> [01 October 2010].
49. Java Documentation Link. <http://Java.sun.com/j2se/1.4.2/docs/api/Java/io/InputStream.html#close%28%29> [01 October 2010].
50. Prefork - Apache HTTP Server. <http://httpd.apache.org/docs/2.0/mod/prefork.html> [01 October 2010].
51. Lighttpd - Server.max-workerDetails - lighty labs. <http://redmine.lighttpd.net/projects/lighttpd/wiki/Server.max-workerDetails> [01 October 2010].
52. Souders S. Roundup on Parallel Connections. Available at: <http://www.stevesouders.com/blog/2008/03/20/roundup-on-parallel-connections/> [01 October 2010].
53. Cormode G, Krishnamurthy B. Key differences between Web1.0 and Web2.0. *First Monday* 2008; **13**(6). <http://firstmonday.org/htbin/cgiwrap/bin/ojs/index.php/fm/article/view/2125/1972>.