

Journal of Theoretical and Applied Electronic
Commerce Research

E-ISSN: 0718-1876

ncerpa@utalca.cl

Universidad de Talca
Chile

Kim, Jeong-Su; Seo, Sang-Koo
Experiment and analysis for QoS of E-Commerce systems
Journal of Theoretical and Applied Electronic Commerce Research, vol. 1, núm. 3, december, 2006,
pp. 1-15
Universidad de Talca
Curicó, Chile

Available in: <http://www.redalyc.org/articulo.oa?id=96510302>

- How to cite
- Complete issue
- More information about this article
- Journal's homepage in redalyc.org

redalyc.org

Scientific Information System
Network of Scientific Journals from Latin America, the Caribbean, Spain and Portugal
Non-profit academic project, developed under the open access initiative

Experiment and Analysis for QoS of E-Commerce Systems

Jeong-Su Kim¹ and Sang-Koo Seo²

¹ Kwangwoon University, Department of Management Information Systems, jskim0723@paran.com

² Kwangwoon University, Department of Management Information Systems, skseo@kw.ac.kr

Received 20 September 2006; accepted 11 November 2006

Abstract

It is important for electronic commerce companies to understand the service quality of their systems, such as, the response time for users' interactions, the service delay zone and the number of appropriate users accessing the systems concurrently. The accurate and prompt management of the service quality can greatly help build and operate systems more efficiently. In this paper we present a methodology for the service quality measurement and the user capacity modeling for electronic commerce systems. While most previous researches on this issue have been based on the closed-LAN environment, we conduct experiments under real network environment using sample e-Commerce systems. Specifically, we measure the response times for e-Commerce transactions under Cable, DSL, and wireless networks, and analyze the delay zones in processing the users' service requests. The discrete event simulation and hybrid simulation are performed to estimate the maximum number of users using a response time limit as the service quality criterion. We also investigate the self-similar characteristics on the response time and the number of users, and the extensive results of the experiment and the simulations are described.

Key words: electronic commerce systems, QoS measurement and prediction, self-similarity traffic

1 Introduction

Since 1990s Internet e-business services have advanced and spread out tremendously. In these past years e-business solutions for enterprises have made rapid strides along with the growth of the Internet, and e-business computer systems have been upgraded or replaced with new systems in very short life cycles. Consequently, one of the main interests of e-business enterprises has been to make effective investments in building and operating electronic commerce services. How will the companies provide adaptable services accommodating more users in end-to-end electronic commerce systems? How much will the new services be better than existing ones? Where will be the major delay zones in communication and/or computation with the increased number of concurrent users of e-Commerce systems? If the answers to these questions can be predicted priorly before building actual systems or if they can be obtained promptly and precisely during system operations, electronic commerce companies will be able to build and manage their systems more effectively.

There has been a lot of previous work on this issue under the closed-Local Area Network (LAN) system environment [9], [10], [19], [23]. But the research on measurement and prediction of service quality considering business transactions in a real-life network environment is hard to find in the literature. In this paper, we try to measure the Quality of Service (QoS) of end-to-end high-speed Internet service on specific service areas. We analyze the maximum capacity of concurrent users from the source to the destination by performing various experiments on client response time using discrete event and hybrid simulations. This analysis is based on the actual measurement under end-to-end high-speed Internet service. We have also investigated the exponential distribution of self-similarity traffic characteristics, the pattern variations about traffic characteristics of Pareto distribution, the impact of variations on traffic rate, Hurst parameters, average response time, etc. The contribution of this research includes: firstly, we describe detailed methodology and procedure to measure QoS and to identify service delay zones for electronic commerce systems; secondly, we show how to analyze the impact of the number of users upon the service quality so as to provide services adaptively; finally, the paper contains various experiment results conducted under real network environment for sample e-Commerce systems.

The paper is organized as follows. Section 2 introduces the related work and background of end-to-end QoS measurement and self-similarity traffics. Section 3 describes sample e-Commerce systems and end-to-end QoS measurement. In Section 4 we define QoS parameters and prediction methodology, and analyze discrete event and hybrid simulations. Section 5 compares and analyzes various experiment results. Finally, in Section 6 we conclude the paper.

2 Related Work

2.1 End-to-end QoS Measurement

End-to-end QoS measurement and its prediction have been very interesting issues in the area such as network bandwidth of service provider, capacity design for application server, resource design for users, etc. One of the largest projects on this issue includes INTERMON [13]. The INTERMON defines an advanced architecture for inter-domain QoS analysis framework development. The advanced architecture encompasses the measurement of end-to-end resource control, traffic control between routers, and admission control from Communication Measurement (CM) toolsets. The CM toolsets support the QoS monitoring of application flows using the agent software, which is configured between end-to-end hosts. That is, the source and the destination hosts are interconnected into CM toolsets. The traffic control between routers is used to analyze protocol types and recognition of pattern behavior [1], [11], [12]. In addition, the integration of the INTERMON Toolkit User Interface and Policy Based Control Tool Interaction provides inter-domain QoS monitor, which provides an analysis function using databases access measurements. The architecture for the INTERMON Toolkit and the end-to-end QoS under inter-domain environment measurement configuration is shown in Figure 1.

The INTERMON Toolkit has simulation and prediction capabilities. The simulation function makes it possible to analyze Internet packets from routers after information is collected from traffic sources [4]. The prediction function estimates QoS parameters in the communication network of Autogressive Integrated Moving Average (ARIMA) model [20]. The traffic prediction forecasts future traffic QoS based on the pattern analysis of the past QoS behavior of abnormal traffics (e.g., route failures, fault operations, protocol variations, Denial of Service (DoS) attacks, configuration failures, etc) for the inter-domain network zone.

Another related study includes the prediction for data transmission performance of Wide Area Network (WAN), in which collected log data are analyzed and the performance evaluation is predicted from WAN [24]. In the study end-to-end performance information is collected from the past transmission, and the system provides prediction data

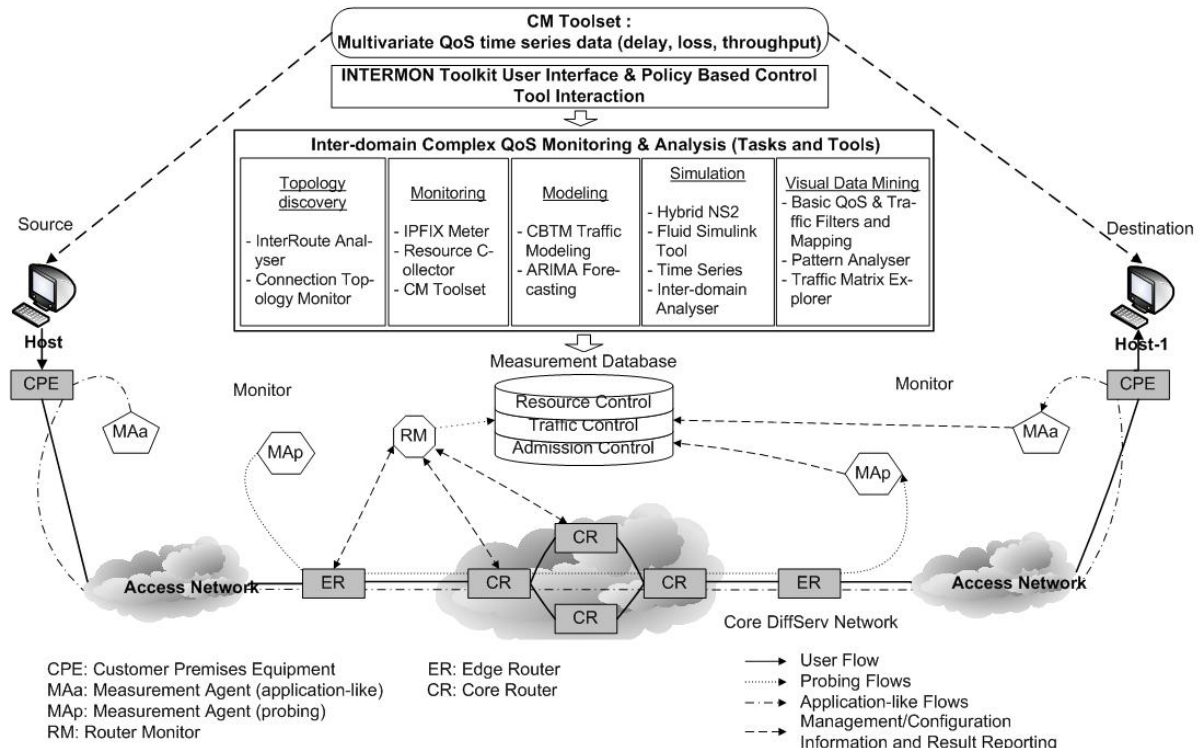


Figure 1: End-to-end QoS under Inter-domain Environment Measurement Configuration

results as the measurement of data transmission. Then the prediction framework is configured for users, and prediction data can be obtained from data transmission infrastructure.

Although both INTERMON and the data transmission performance study of WAN support prediction capabilities, the issue of the maximum capacity of concurrent users under the end-to-end real-life network environment has not been taken into account sufficiently. Our paper proposes a prediction methodology for estimating the allowable number of users, and describes in detail the simulation results for various QoS parameters. Also explained is the self-similarity traffic analysis for the average response time with respect to the number of the concurrent users and the comparison of the Hurst parameter with the Pareto distribution.

2.2 Self-Similarity Distribution

Telephone network mostly relies on the Poisson and exponential distribution because of the characteristics voice traffic. Until now, many experimental research have been based on the Poisson distribution, in which, however, it is hard to express the traffic characteristics for modern networks accurately. An alternative is to generate traffic with self-similarity characteristics. The self-similarity traffic is the property associated with one type of the fractal, an object whose appearance is unchanged regardless of the scale at which it is viewed. In the case of stochastic objects like time series, self-similarity is used in the distributional sense: when viewed at varying scales, the object's correlational structure remains unchanged. As a result, such time series exhibit bursts-extended periods above the mean at a wide range of time scales [5], [22], [25]. For this reason self-similarity has been suitable for analyzing Ethernet network traffics, the nature of World Wide Web (WWW) traffics, and the benchmark of network traffics. The self-similarity has intrinsic properties of a network system subject to self-similarity traffic conditions. For instance, self-similarity bursty network traffic comes about as a consequence of one of the most innocuous network activities - the transfer of files in a networked client/server system by a number of concurrent connections/sessions. However, the self-similarity traffic tends to get more variation when observed for a long time about overlapped traffic under small groups in a backbone network. Difference between self-similarity and Poisson are depicted in Figure 2.

When network designer plans network infrastructure, he/she must consider enough capacity by the burstiness. That is, the design must provide sufficient requirement levels for customer satisfaction. The various researches for self-similarity of data network have been performed. Crovella proposed measurement methods about the major parameters for heavy-tailed distribution in Web traffic [6]. The heavy-tailed distribution is appropriate to estimate tail weight in Web data (e.g., images, audio, video, text, archives, preformatted text and compressed files). A distribution is heavy-tailed if its tail asymptotically follows a power law. That is,

$$P_r\{X > x\} \sim x^{-\alpha}, \text{ as } x \rightarrow \infty, 0 < \alpha < 2 \quad (1)$$

The expression (1) implies that a random variable X follows a heavy-tailed distribution. One of the simplest heavy-tailed distributions is the Pareto distribution whose probability density function is given by

$$p(x) = \alpha k^\alpha x^{-\alpha-1} \quad (2)$$

where α is the shape parameter, $k > 0$ is the location parameter, and $x \geq k$. The distribution function has the form

$$F(x) = P_r\{X \leq x\} = 1 - (k/x)^\alpha \quad (3)$$

Heavy-tailed distributions have a number of properties that are qualitatively different from distributions more commonly encountered in networking research, in particular, the exponential distribution. If $\alpha \leq 2$, the distribution has infinite variance, and if $\alpha \leq 1$, the distribution has also infinite mean. Thus, as α decreases, a large portion of the probability mass resides in the tail of the distribution. In practical terms (also relating to our network model) a random variable that follows a heavy-tailed distribution can give rise to extremely large file size requests with non-negligible probability.

The measure of self-similarity can be expressed by Hurst parameter (H) with the range of $1/2 < H < 1$. The Hurst parameter value close to 1 means high self-similarity.

$$H = (3 - \alpha) / 2 \quad (4)$$

where the range of α is $1 < \alpha < 2$. Detailed description of the formulas can be found at [6].

Willinger presented in his study a physical explanation for the self-similar nature of today's packet network traffic [26]. Willinger provided mathematical results and validated his findings with detailed statistical analyses of two representative sets of high time-resolution traffic measurements from two different Ethernet LAN's. Traditional ON/OFF transmission models were assumed to execute exponential or geometric distributions, but each transmission source extended its range to get Noah Effect. The Noah Effect has characteristics about a wide range of time scales ("high-variability sources"), which is consistent with measured network traffic, thus, exhibiting the same self-similar or fractal properties as can be observed in the data.

It is argued that the self-similarity occurs because of the characteristics of traffic source generation [22]. The discrete event simulation and hybrid simulation conducted in our research may be in a sense confined to the limited data traffic generated by users' transactions in WAN environment. For this reason we investigate the self-similar characteristics on the response time and the number of maximum concurrent users of sample e-Commerce systems. We apply the Hurst parameter value for our simulation such that the Pareto distribution has the scale parameter value of 1 to get the heavy-tailed distribution.

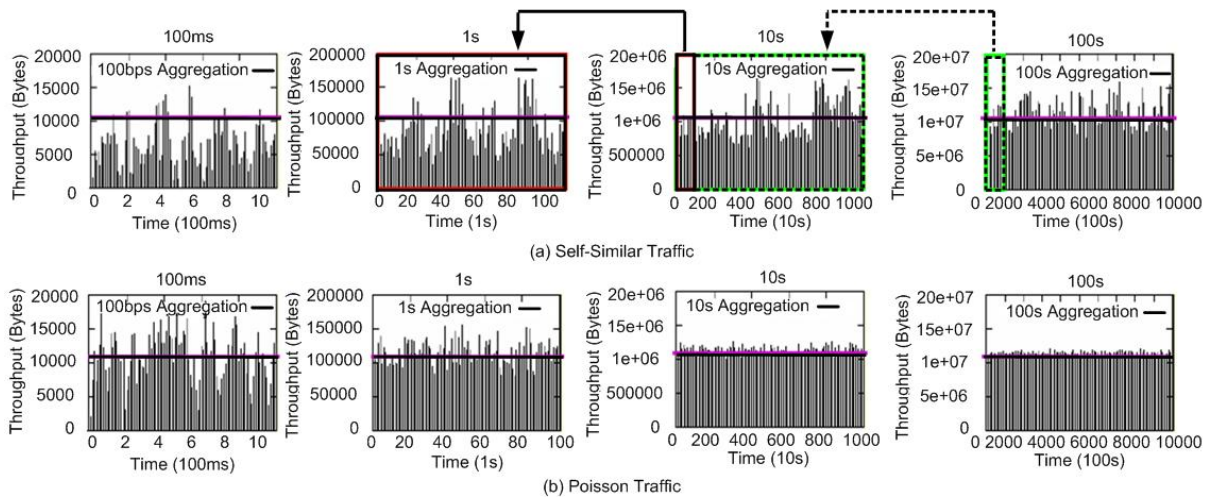


Figure 2: Self-similarity vs. Poisson

3 QoS Measurement Model for Electronic Commerce

3.1 Electronic Commerce Systems and Transaction Model

In order to investigate the end-to-end QoS measurement under real network environment we built a sample e-Commerce system. Since the transaction processing patterns play important roles in network and system performance, two kinds of Web systems are implemented: a shopping mall system and a ticketing system. We reflect the characteristics of e-Commerce systems and those two types are supposed to be most representative ones. The model for the e-Commerce transaction traffics in our research is based on the Layered Queuing Model (LQM), which is an extended concept of the Queuing Network Model (QNM) to estimate the response time between client and server computers in distributed systems [17]. The parameters and configuration of LQM are extended for user transactions under WAN environment instead of closed-LAN environment. A user transaction consists of a sequence of interactions between a user and the e-Commerce system from the web site connection through the end of an entire transaction.

Shopping mall users are assumed to follow the eight steps: connection to the main Web page, shopping items search, shopping items selection, putting the items into a shopping basket, deletion of items from the basket, addition of other shopping items into the basket, filling up an order sheet, and payment process. Each step involves transactional, non-transactional, or both types of works. A non-transactional work means the access of the http pages containing the static images and texts without requiring the database interactions. A transactional work means the access of the http pages that require the database interactions. For example, a shopping customer may put items into the shopping basket after filling up his/her customer ID and password into an http page. If the customer decides to purchase the items, all the shopping information, including that of the products and the customer, will be stored into the database. The ratio of non-transactional versus transactional works in the shopping mall system is configured as 20:80. The ratio might not exactly represent real shopping mall systems but it is meant to reflect database-centric applications. Likewise, the websites interactions for the ticketing system are assumed to be seven steps: connection to the main web page, searching movies, choosing a movie, deciding a theater, date and time selection, payment, and confirmation for the reservation. The ratio of non-transactional versus transactional works is configured as 60:40, which is to reflect more navigation-centric applications.

We make some additional assumptions for user interactions. A client is supposed to randomly select shopping items (or tickets) from the shopping mall (or the ticketing site) and the thinking time between the users' interactions is not considered. The thinking time may be incorporated safely without making serious impact on the methodology and result of our research. We exclude the process of the certification for the electronic payments (e.g., bank/card systems). Therefore, if a user's request for the shopping mall or the ticketing site is successfully done, these transactions are assumed to be stored safely in the database. Otherwise, the transactions will reset to the initial values without updating the database. The sample Web systems were built using the Active Server Page (ASP) under Windows 2000 servers. The databases for the shopping mall and the ticketing systems are maintained by the Microsoft Access 2003 and the Oracle 8.0.5, respectively.

3.2 Response Time Measurement

We conducted an experiment to measure the client response time for the shopping mall and the ticketing systems using the high-speed Internet services: Wireless LAN, ADSL, Cable, and VDSL. In the experiment Wireless LAN offers 11Mbps bandwidth (IEEE 802.11b) and the other services (ADSL, Cable, and VDSL) offer 100Mbps bandwidth each. The measurement configuration is shown in Figure 3.

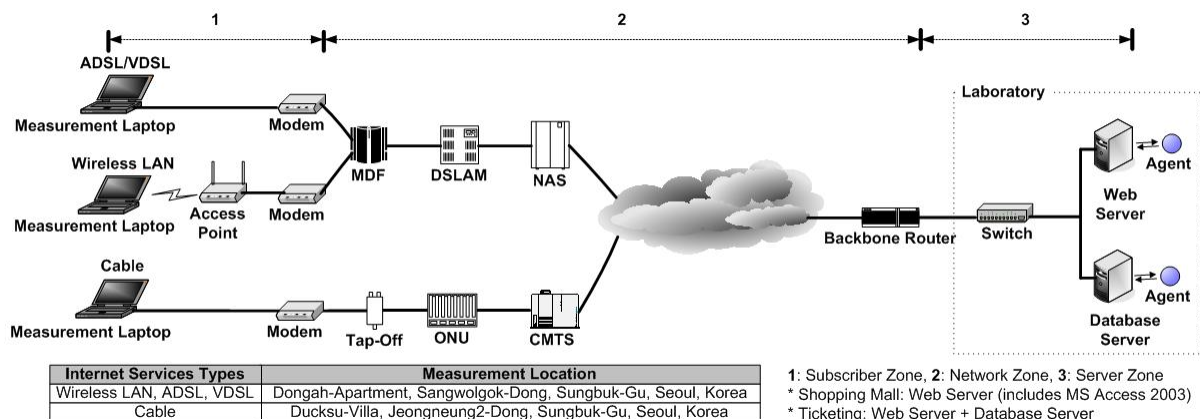


Figure 3: Measurement Configuration

We divide end-to-end (i.e., between client and server computers) delay zones into three regions: the subscriber zone, the network zone, and the server zone. The subscriber zone is from the personal computer of an e-Commerce user to modem, the network zone is from the back of the modem to the front of backbone router, and the server zone is from the back of backbone router to the server computer. The subscriber zone may show different delay time in each high-speed Internet service because of the differences in the equipments of the users' computers and the network environment. The three zones will be useful to identify the QoS responsibility among e-Commerce service providers and Internet Service Providers (ISPs).

As a measurement tool we use the IT Guru, version 10.5 of OPNET, with Application Characterization Environment (ACE) and ACE Decode Module [2], [14], [21]. The experiment steps for the end-to-end response time measurement are summarized in the following:

- **Step 1:** The agent software is installed at target servers to measure the response time.
- **Step 2:** High-speed Internet at a user's home is connected to the client computer.
- **Step 3:** Business transactions are submitted 10 times repeatedly.
- **Step 4:** Average response time is collected and analyzed.
- **Step 5:** Analyze and identify the delay zone of the end-to-end network.
- **Step 6:** Analyze the cause of the bottleneck.
- **Step 7:** Diagnose for the improvement.

We measure the response time for a whole (in another word, end-to-end) business transaction, which includes all the steps of interactions for the shopping mall and ticketing systems. The measurement results are shown in Figure 4. The packet information can be collected by the agents using the passive method. The results of packet analysis shows the major delay zones in the experiment are in the order of the subscriber zone > the network zone > the server zone for Wireless LAN, ADSL, and Cable connections. On the other hand, under VDSL, the order is the network zone > the server zone > the subscriber zone.

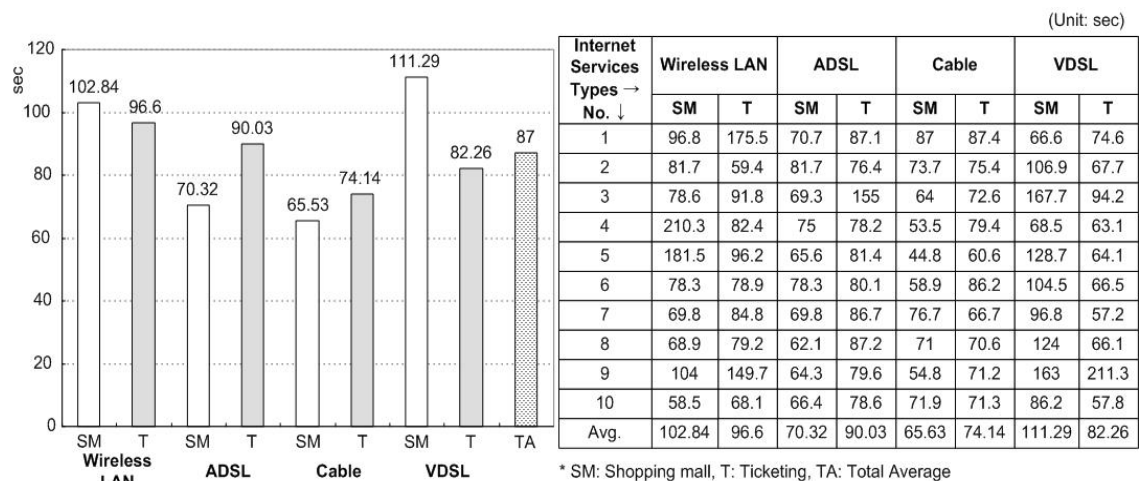


Figure 4: Measurement Results for End-to-end Business Transactions

The experiment results imply that the end-to-end QoS response time measurement can be used to find delay zones using the passive method between agents. While it is known that the delay zones for the case of LQM on closed-LAN were on the network and server zones, our experiment result indicates the subscriber zone as a major bottleneck. The reason is supposedly because the e-Commerce system is strongly based on a bi-directional packet transmission.

4 QoS Prediction and Analysis

4.1 QoS Parameters for Electronic Commerce

We define the parameters for QoS prediction as in Table 1, which are applied to our model for each high-speed Internet service. The simulation experiment for the maximum capacity of concurrent users consists of the discrete event simulation and the hybrid simulation. In the case of the discrete event simulation, input variables include the transaction type, the number of users, the system specification, the point-to-point bandwidth, etc. The hybrid simulation has additional variables for the LAN background utilization, the CPU background utilization, the link background utilization, the traffic flow (bi-directional), the device control, and so on. The dependent variables are the

CPU utilization (server and/or client), the throughput, the client response time, the LAN delay, the LAN inbound/outbound traffic, the traffic sent/received, the delay, the packet loss rate, etc.

Types	Parameters	Definition
Server	CPU Utilization	The utilization rate for the server's CPU
Client	Response Time	Total response time from the subscriber zone to the server zone for a whole business transaction.
	CPU	The utilization rate for client's CPU
	L A N Delay	Total transmission delay time for LAN on the network
	Inbound Traffic Outbound Traffic	Total amount of the traffic generated on LAN
Network	Traffic Sent	Total amount of sent bytes per second
	Traffic Received	Total amount of received bytes per second
	Queuing Delay	The queuing delay rate between each zone
	Packet Loss Rate	Bi-directional packet loss rate between source and destination
Commonness	Throughput	The network transmission rate between each zone
	Utilization	The network utilization rate between each zone

Table 1: QoS Parameters for Electronic Commerce

4.2 The Number of Users Prediction

The procedure for predicting the number of users consists of the following steps:

- **Step 1:** Each electronic commerce service defines the transaction type.
 - The shopping mall system and the ticketing system have eight and seven processes, respectively.
- **Step 2:** Business transactions are applied 10 times repeatedly from the end-to-end network environment.
- **Step 3:** Calculate the average response time.
 - Average response time is measured as 87 sec for the end-to-end Internet environment (See the Figure 4 for details).
- **Step 4:** Perform the discrete event simulation and the hybrid simulation.
 - Workload parameters are initialized.
 - Select a high-speed Internet service.
 - Measure the response time for business transactions.
 - Import the measurement values from the business transactions.
 - Measure the simulation time.
- **Step 5:** Simulate the assignment of the additional resources with respect to the average response time.
 - If the average response time doesn't meet the Service Level Agreements (SLA) criterion, additional resource (e.g., devices, nodes, and links) allocation is considered.
 - The input values for the parameters are re-set appropriately.
- **Step 6:** Prediction for the number of users
 - While the client response time is within the value of the criterion scope, the number of concurrent users at that moment is assumed suitable.
 - If the client response time is beyond the criterion scope, the input values are re-established.
- **Step 7:** Simulation finished when client response time is satisfied within the criterion scope.

The discrete event simulation and hybrid simulation allow estimating the delay zone and the maximum number of concurrent users. The previous research have considered various parameters to predict the number of users, including the response time, the throughput, the packet loss, the path-oriented differentiated service, etc [3], [7], [16], [18]. IBM initiated the Web Service Level Agreements (WSLA) project, which aimed at the creation and monitoring of SLAs in a Web services environment. The distributed monitoring framework can possibly manage to provide a different service with the service level agreements. Using the framework, service providers can manage their resources efficiently and flexibly to optimize the customer satisfaction [15]. These prediction capabilities include the extension of an abstract forecast type and new domain specific predicates [8]. We attempt to investigate the average response time with respect to the number of concurrent users with varying a criterion value, and the research on this aspect has not been published in the literature to our knowledge.

The QoS parameters are very important to predict the number of users. Predicting the number of users from the end-to-end environment requires considering the whole zones (the subscriber zone, the network zone, and the server zone). For instance, shopping customers want to process quickly and accurately. If it takes long delay time to

purchase goods, the customers will leave the e-Commerce site. Therefore, our research focuses on the maximum capacity of concurrent users with respect to the average response time. The experiment under closed-LAN environment appears that the 100Mbps network equipment may not reach its full capacity in the network throughput. That is, a service of 100Mbps bandwidth can support at most up to 70 ~ 80 Mbps bandwidth only. Also, the same is true for the 10/100Mbps network equipments. Furthermore, ISP companies are not willing to disclose the detailed information about their network bandwidths for specific regions and sub-nets. For this reason, while the network bandwidths of client and server zones can be measured, that of the other zone (i.e., network zone) is assumed to be unknown in our experiment.

4.3 Discrete Event and Hybrid Simulations

Discrete event simulation can predict various application performances, and generate the network model from the ACE module. For this purpose we exploit the model library of the IT Guru. On the other hand, the hybrid simulation combines with the explicit and background traffics. The processing in hybrid simulation is very similar to a real-life environment configuration because of the bi-directional traffic flow between the client and the application server. Specifically, the explicit traffic using the discrete event simulation is comprised of two parts: (1) the passive method based on device/node/link/traffic flow and (2) the use of the measured data from the ACE module. We choose the second method because it allows analyzing the maximum capacity of concurrent users under each high-speed Internet service. The background traffic is established in the passive method under the device link topology. These types include the router, the server, the workstation, and the LAN (designate group of clients). Moreover, the traffic flow is set up in the passive method with bi-directional traffics. The passive configuration and the workload parameters for simulation configuration are given in Table 2 and Table 3.

Client LAN (LAN Background Utilization)		Web Server & DB Server CPU Background Utilization		Cache Server CPU Background Utilization	
Time (sec)	Background Utilization (%)	Time (sec)	Background Utilization (%)	Time (sec)	Background Utilization (%)
10	10	10	10	10	90
20	20	20	20	20	80
30	30	30	30	30	80
40	40	40	40	40	70
50	50	50	50	50	80
60	60	60	60	60	90
Link Load, Background Load (Intensity(bps)), Commonness				Traffic Flow (Bi-directional)	
Direction: →		Direction: ←			
Sec	Bits/Sec	Sec	Bits/Sec	Sec	Bits/Sec
0.0	1,503,154.15	0.0	1,803,200.15	0.0	703,154.15
301	1,599,958.36	301	1,879,999.36	301	599,958.36
601	1,603,522.90	601	1,903,522.90	601	703,522.90
900	1,605,554.02	900	1,345,554.02	900	605,554.02
1,201	1,589,141.84	1,201	1,672,141.84	1,201	589,141.84
1,501	1,603,397.96	1,501	1,777,397.97	1,501	603,397.97
1,801	1,601,679.17	1,801	1,967,679.17	1,801	901,679.17
2,101	1,602,894.22	2,101	1,893,894.22	2,101	902,894.22
2,400	1,602,658.46	2,400	1,652,658.46	2,400	602,658.46
2,701	1,599,714.55	2,701	1,643,714.55	2,701	799,714.55
3,001	1,602,708.84	3,001	1,656,708.84	3,001	802,708.84
3,301	1,582,708.20	3,301	1,999,708.20	3,301	982,708.20

Table 2: Background Traffic Configuration

Category	Parameters
Packet Loss Rate and Latency	0
Number of Users	25 ~ 2900
Client Location	Remote
Simulation Time	1 hour
Random Seeds	128
Result of Statistics	100
Update Interval	500,000
Simulation Kernel	Optimize Mode

Table 3: Workload Parameters

4.4 Experiment Model

The experiment models are shown in Figure 5. In the Case 1, the wireless LAN is connected with 10Mbps bandwidth between the client and the remote switch, and ADSL, Cable, and VDSL are set up with 100Mbps bandwidth each. All high-speed Internet services are assumed to have 1.5Mbps bandwidth between the remote and the local routers. Delay zones are divided into three: the subscriber zone is from the user's PC to the remote switch; the network zone is from the back of the remote switch to the local router; and the server zone is supposed from the back of local router to the server computer.

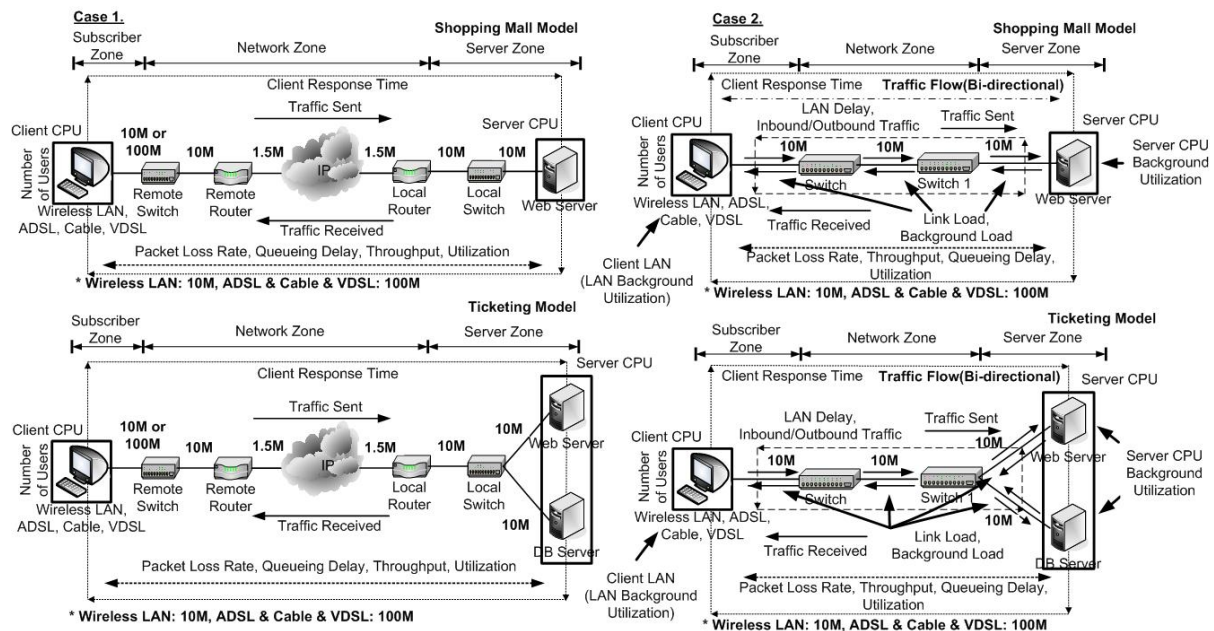


Figure 5: Experiment Model

The Case 2 in Figure 5 shows the hybrid simulation, which is done under 10Mbps bandwidth with two switches between the client and the server. The delay zones are also divided into three: the subscriber zone is from the user's PC to the previous switch; the network zone is from the switch to the switch 1; and the server zone is from the back of the switch1 to the server. Figure 6 shows the experiment results of the response time by increasing the number of concurrent users. In the case of shopping mall, the client response time has increased according to the increased number of concurrent users from the discrete event and the hybrid simulations. In the case of ticketing system, results are similar to the shopping mall in that the client response time is analyzed to be increasing with the increased users. The capacity of the concurrent users has shown much more overhead in the shopping mall system than in the ticketing system. The experiment results for the maximum capacity of the concurrent users for both systems are attached at Appendix 1. The parameters regarding the maximum capacity of concurrent users are summarized below.

1. Queuing delay
 - In Case 1 (the experiment results of the discrete event simulation), Each delay zone for the shopping mall and the ticketing systems has been analyzed to be the network zone.
 - In Case 2 (the experiment results of the hybrid simulation), all delay zones except ADSL of the shopping mall system have been on the network and the server zones.
2. Throughput
 - In Case 1, all of the Internet services of the ticketing system have shown more throughput in the server zone. But, the ADSL of the shopping mall system showed more throughput in the subscriber and the network zones.
 - In Case 2, the network and the server zones showed more throughput except for the ADSL of the shopping mall system and the Cable and the VDSL of the ticketing system.
3. Utilization
 - In Case 1, both shopping mall and ticketing systems showed more utilization in the network zone.
 - In Case 2, all the Internet services, except the ADSL of the shopping mall system and the VDSL of the ticketing system, have shown more utilization of the network and the server zones.

The other parameters are omitted due to the limited space. In Case 1, the maximum capacity of concurrent users in the shopping mall system is 70 in ADSL only. However, the other Internet services aren't satisfied for the given average response time. On the other hand, the maximum users of Wireless LAN, ADSL, Cable and VDSL for the

ticketing system are 645, 250, 440, and 340, respectively. The maximum capacity of concurrent users in the ticketing system in Case 2 are 950, 580, and 500 with the order of ADSL > VDSL > Cable, respectively. But the wireless LAN did not satisfy the given criterion scope. In the case of the ticketing system, the maximum capacity of concurrent users for Cable, VDSL, and ADSL are 2900, 2800, 2300, and 1800, respectively. In brief, the shopping mall system is supposed to have more business transactions according to the depth and the ratio of transactions. Therefore, the maximum capacity of concurrent users has decreased. The result implies that the application server and the network have shown higher utilization with the increased number of concurrent users. Delay zones are found to be the network and the server zones with the increased number of concurrent users.

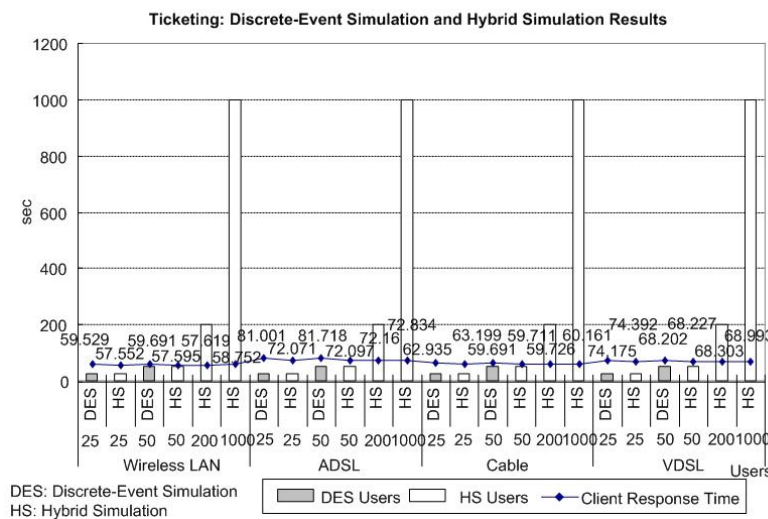
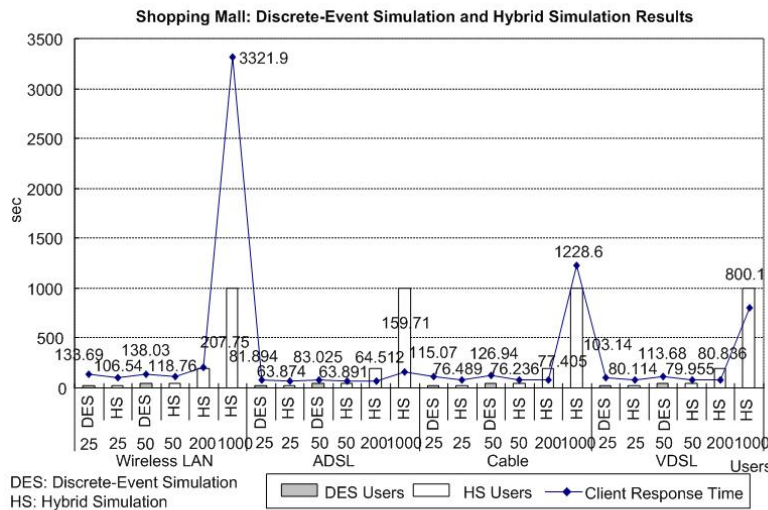


Figure 6: Client Response Time with Respect to the Number of Concurrent Users

Internet Services Types	Simulation Type	Number of Users	Client Response Time
Wireless LAN	Discrete-Event Simulation	25	133.69
		50	138.03
	Hybrid Simulation	25	106.54
		50	118.76
		1000	3321.9
ADSL	Discrete-Event Simulation	25	81.894
		50	83.025
	Hybrid Simulation	25	63.874
		50	63.891
		1000	159.71
Cable	Discrete-Event Simulation	25	115.07
		50	126.94
	Hybrid Simulation	25	76.489
		50	76.236
		1000	1228.6
VDSL	Discrete-Event Simulation	25	103.14
		50	113.68
	Hybrid Simulation	25	80.114
		50	79.955
		1000	800.1

Internet Services Types	Simulation Type	Number of Users	Client Response Time
Wireless LAN	Discrete-Event Simulation	25	59.529
		50	59.691
	Hybrid Simulation	25	57.552
		50	57.595
		1000	81.001
ADSL	Discrete-Event Simulation	25	81.001
		50	81.718
	Hybrid Simulation	25	72.071
		50	72.097
		1000	72.834
Cable	Discrete-Event Simulation	25	62.935
		50	59.691
	Hybrid Simulation	25	63.199
		50	59.711
		1000	60.161
VDSL	Discrete-Event Simulation	25	74.392
		50	68.202
	Hybrid Simulation	25	74.175
		50	68.227
		1000	68.993

5 Self-similarity Traffic

5.1 Traffic Modeling

As mentioned in the section 2.2, the discrete event simulation and the hybrid simulation allow estimating the maximum number of users using the transaction measurement done in a limited range of time. In order to analyze the experiment for larger time scale, we try to exploit the self-similarity characteristics of the network traffics for WAN environment. Figure 7 shows a configuration to support the Web services from users. The experiment compares Pareto and exponential distributions about the http method of the discrete event simulation. To assure the end-to-end network environment, we rely on empirically measured distributions for both client traces and the WWW server. Overall scenario is as follows.

- 1) The generation of the traffic rate from 10 to 100%
- 2) Comparison according to the variation of the Hurst parameter value in the major parameter of the Pareto distribution

3) Comparison of the average response time with respect to the number of concurrent users



Figure 7: Experiment Configuration

5.2 Analysis

The Pattern Variations about Traffic Characteristics of Pareto Distribution

The traffic characteristics show a certain difference between the Pareto distribution and the exponential distribution. At Figure 8(a), in the case of the Pareto characteristics, there occurs a large deviation from the exponential distribution. That is, the burstiness appears absolutely different. We fixed the Hurst parameter to be 0.82 in the analysis. Note that the Whittle estimator for the dataset yields an estimate of $H=0.82$ with a 95% confidence interval of (0.77, 0.87).

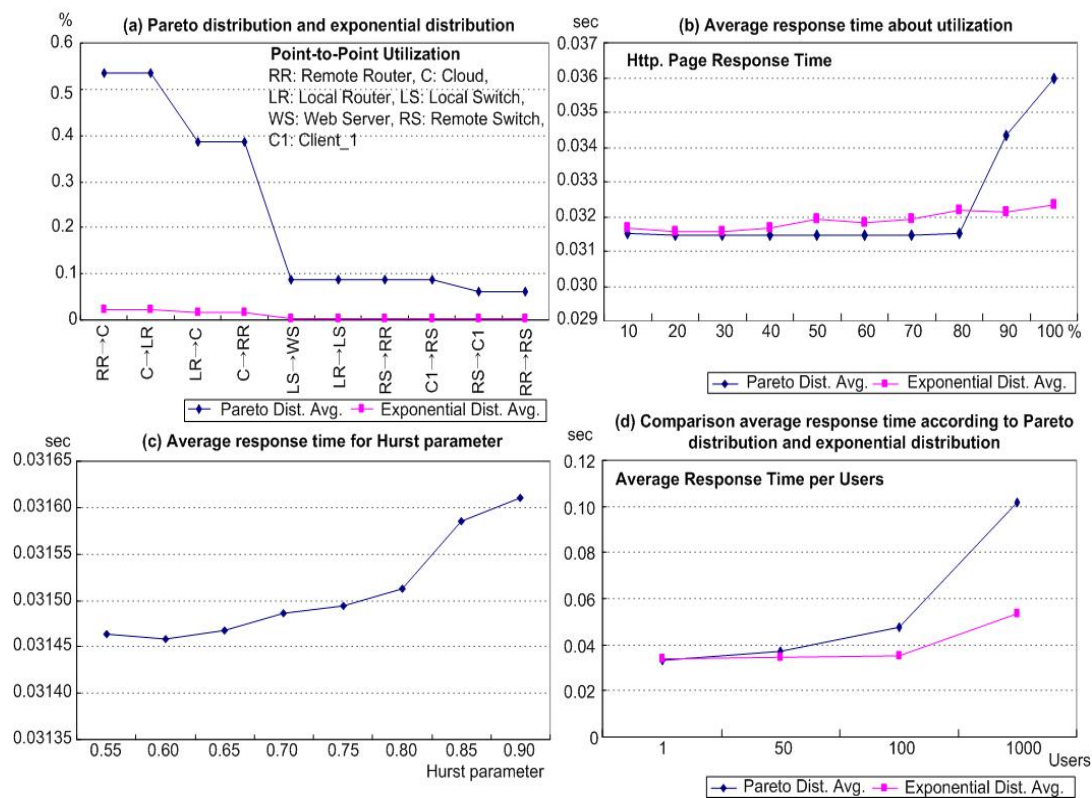


Figure 8: Simulation Result

The Average Response Time

The Pareto characteristics reached over 80% utilization rate, which shows the increase in the average response time. The result is shown in Figure 8(b).

The Hurst Parameter

As a major parameter of the Pareto function, the Hurst parameter has the range of $0.5 < H < 1$ and the self-similarity increases up to nearly 1 (one). Figure 8(c) shows the result of the average response time when the parameter value changes between 0.55 and 0.90.

Comparison of the Average Response Time According to Pareto and Exponential Distributions

Figure 8(d) shows the average response time about the number of concurrent users under the Pareto and the exponential distributions. In the case of a single user with the same value of both distributions, the average response time has increased for the Pareto distribution with increasing the number of concurrent users. When the number of concurrent users reaches about 1000, the response time is measured to be less than one second. We recognize through the experiment that it is hard to allow more capacity against the number of concurrent users.

6 Conclusion

The measurement and the prediction of the service quality are very important issues for the efficient implementation and the operation of e-Commerce systems. In this paper, we proposed a methodology for the measurement of the response time, identification of delay zones, and the user capacity modeling under real high-speed Internet environment instead of closed-LAN environment. We conducted extensive experiments by constructing sample e-Commerce systems and by using various simulation techniques. The experiment and simulation results indicate followings: firstly, the service delay zone with small number of users is mostly confined with the subscriber zone, and is shifting to the network and the server zones with the increasing number of concurrent users; secondly, we observe that the response time of each high-speed Internet services depends heavily on the sequence and the depth of the business transactions of e-Commerce systems and also on the ratio of their transactional versus non-transactional operations; finally, the result of the response time and the user capacity simulation is consistent with the self-similarity characteristics, leading the Hurst parameter value close to 1 (one).

Further work would include the investigation of various scenarios of e-Commerce systems with considering more advanced networking infrastructure and software packages, the analysis and the optimization for the link utilization between routers for the improvement of QoS, and the integration and the exploitation of advanced monitoring agents. We hope that the methodology and the experiment/simulation results in this paper be an aid for the future research on the QoS of more diverse types of e-Commerce systems, also serving as a useful data for the ISPs, the Internet Data Centers (IDCs), and e-Commerce companies.

References

- [1] AQUILA Consortium, (2002). Adaptive resource control for QoS using an IP-based layered architecture. [Online]. Available: <http://www-st.inf.tu-dresden.de/aquila/>.
- [2] ACE User Guide, IT Guru Product Documentation Release 10.5, OPNET Technologies, Inc., 2004.
- [3] M. F. Arlitt and C. L. Williamson, Internet web servers: workload characterization and performance implications, IEEE/ACM Transactions on Networking, vol. 5, no. 5, 1997.
- [4] F. Baumgartner, M. Scheidegger, and T. Braun, Enhancing discrete event network simulators with analytical network cloud models, International Workshop on Inter-domain Performance and Simulation (IPS), 2003, pp. 21-30.
- [5] M. E. Crovella and A. Bestavros, Self-similarity in world wide web traffic: evidence and possible causes, IEEE/ACM Transactions on Networking, vol. 5, no. 6, 1997.
- [6] M. E. Crovella, M. S. Taqqu, and A. Bestavros, Heavy-tailed probability distributions in the world wide web, in a Practical Guide to Heavy Tails, 1998, pp. 3-26.
- [7] Y. Cheng and W. Zhuang, Dynamic Inter-SLA resource sharing in path-oriented differentiated services networks, IEEE/ACM Transactions on Networking, vol. 14, no. 3, 2006.
- [8] A. Dan, H. Ludwig, and G. Pacifici, Web services differentiation with service level agreements, IBM Software Group, 2003.
- [9] J. C. Hu, S. Mungee, and D. C. Schmidt, Techniques for developing and measuring high-performance web servers over ATM networks, in INFOCOM '98 Conference, 1998, pp. 1222-1231.
- [10] J. C. Hu, I. Pyrali, and D. C. Schmidt, Measuring the impact of event dispatching and concurrency models on web server performance over high-speed networks, in Proceedings of the 2nd Global Internet Conference, 1997.
- [11] U. Hofmann, I. Miloucheva, T. Pfeifferberger, and F. Strohmeier, Evaluation of architecture for QoS analysis of applications in Internet environment, the 10th International Conference on Telecommunication Systems Modeling and Analysis Monterey, 2002.
- [12] U. Hofmann, I. Miloucheva, and T. Pfeifferberger, INTERMON complex QoS/SLA analysis in large Internet environment, ACM International Conference Proceeding, 2004, pp. 1-13.
- [13] INTERMON, (2001). Advanced architecture for Inter-domain quality of service monitoring, modeling and visualization. [Online]. Available: <http://www.ist-intermon.org/>.
- [14] IT Guru User Guide, IT Guru Product Documentation Release 10.5, OPNET Technologies, Inc., 2004.
- [15] IBM, (2003). Web Service Level Agreements (WSLA) Project. [Online]. Available: <http://www.research.ibm.com/wsla/>.
- [16] M. Jain and C. Dovrolis, End-to-end available bandwidth: measurement methodology, dynamics, and relation with TCP throughput, in Proceedings ACM SIGCOMM, 2002.
- [17] D. Krishnamurthy and J. Rolia, The Internet vs e-Commerce servers: when will server performance matter?, in Proceedings of CASCON '98, 1998, pp. 246-258.

- [18] Y. S. Kim, Response time simulation for 2-tier and 3-tier client/server system, Korea Society for Computer Science, vol. 9, no. 3, 2004.
- [19] J. C. Mogul, The case for persistent-connection HTTP, DEC Western Research Laboratory, Technical Report WRL 95/4, 1995.
- [20] I. Miloucheva, E. Muller, and A. Anzaloni, A practical approach to forecast quality of service parameter considering outliers, International Workshop on Inter-domain Performance and Simulation (IPS), 2003.
- [21] Methodologies and Case Studies, IT Guru Product Documentation Release 10.5, OPNET Technologies, Inc., 2004.
- [22] K. Park, G. Kim, and M. Crovella, On the effect of traffic self-similarity on network performance, in Proceedings SPIE International Conference on Performance and Control of Network Systems, 1997, pp. 296-310.
- [23] M. Sopitkamol, Ranking configuration parameters in multi-tiered e-Commerce sites, CAN SIGMETRICS Performance Evaluation Review, vol. 32, no. 3, pp. 24-33, 2004.
- [24] S. Vazhkudai, J. M. Schopf, and I. T. Foster, Predicting the performance of wide area data transfers, in 16th International Parallel and Distributed Processing Symposium, 2002.
- [25] W. Willinger, M. S. Taqqu, R. Sherman, and D. V. Wilson, Self-similarity through high-variability: statistical analysis of Ethernet LAN traffic at the source level, IEEE/ACM Transaction on Networking, vol. 5, no. 1, pp. 71-86, 1997.
- [26] W. Willinger, M. S. Taqqu, R. Sherman, and D. V. Wilson, Self-similarity through high-variability: statistical analysis of Ethernet LAN traffic at the source level, in Proceedings ACM SIGCOMM, 1995, pp. 100-113.

Appendix 1: Results for the Maximum Capacity of Concurrent Users

Case 1: Discrete Event Simulation for Shopping Mall System

IST	NoU	CPU Utilization				Client Response		Throughput											
		Client		Server		Time		Client RS	RS Client	RS RR	RR RS	RR Cloud	Cloud RR	Cloud LR	LR Cloud	LR LS	LS LR	LS WS	WS LS
		Avg. (%)	Max. (%)	Avg. (%)	Max. (%)	Avg. (sec)	Max. (sec)	Avg. (bits/sec)	Avg. (bits/sec)	Avg. (bits/sec)	Avg. (bits/sec)	Avg. (bits/sec)	Avg. (bits/sec)	Avg. (bits/sec)	Avg. (bits/sec)	Avg. (bits/sec)	Avg. (bits/sec)	Avg. (bits/sec)	Avg. (bits/sec)
ADSL	70	60,526	62,026	11,342	11,661	86,929	108,18	33,564	490,967	33,564	490,967	26,254	481,567	26,254	481,567	33,564	490,966	33,565	490,966
IST	NoU	Traffic Sent (bytes/sec)		Traffic Received (bytes/sec)		Utilization													
		Client	Server	Client	Server	Client RS	RS Client	RS RR	RR RS	RR Cloud	Cloud RR	Cloud LR	LR Cloud	LR LS	LS LR	LS WS	WS LS		
						Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	
ADSL	70	2,098	58,394	58,394	2,098	0,336	4,910	0,336	4,910	1,709	31,352	1,709	31,352	0,336	4,910	0,336	4,910		
IST	NoU	Queueing Delay														Packet Loss Rate			
		Client RS	RS Clients	RS RR	RR RS	RR Cloud	Cloud RR	Cloud LR	LR Cloud	LR LS	LS LR	LS WS	WS LS						
		Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)						
ADSL	70	0,000134	0,001091	0,000114	0,001091	0,000743	0,006964	0,000578	0,049475	0,000114	0,001091	0,000114	0,001091				0		

IST: Internet Services Types, NoU: Number of Users, WLAN: Wireless LAN, RS: Remote Switch, RR: Remote Router, LR: Local Router, LS: Local Switch, WS: Web Server, DB: Database Server

Case 1: Discrete Event Simulation for Ticketing System

IST	NoU	CPU Utilization				Client Response Time		Throughput													
		Client		Server				Client RS	RS Client	RS RR	RR RS	RS Cloud	Cloud RR	Cloud LR	LR Cloud	LR LS	LS LR	LS WS	WS LS	LS DB	DB LS
		Avg (%)	Max (%)	Avg (%)	Max (%)	Avg (sec)	Max (sec)	Avg (bits/sec)	Avg (bits/sec)	Avg (bits/sec)	Avg (bits/sec)	Avg (bits/sec)	Avg (bits/sec)	Avg (bits/sec)	Avg (bits/sec)	Avg (bits/sec)	Avg (bits/sec)	Avg (bits/sec)	Avg (bits/sec)	Avg (bits/sec)	Avg (bits/sec)
WLAN	645	96.811	99.758	49.228	50.699	86.341	141.23	152,682	1,192,510	152,682	1,192,509	127,324	1,166,395	127,324	1,166,398	152,683	1,192,594	161,735	1,205,363	12,770	9,052
ADSL	250	95.831	99.680	12.442	13.639	86.895	105.04	-	684,297	-	684,296	-	668,767	-	668,767	-	684,296	-	690,898	-	-
Cable	440	96.775	99.874	31.603	32.768	87.719	136.81	164,930	1,159,231	164,930	1,159,230	138,729	1,132,777	138,728	1,132,781	164,928	1,159,241	173,503	1,170,768	11,528	8,575
VDSL	340	96.712	99.904	12.802	13.934	85.789	117.72	134,382	944,409	134,382	944,409	113,219	923,186	113,219	923,189	134,383	944,465	141,947	953,827	9,364	7,564

IST	NoU	Traffic Sent (bytes/sec)		Traffic Received (bytes/sec)		Clients RS		RS Clients	RS RR	RR RS	RR Cloud	Cloud RR	Cloud LR	LR Cloud	LR LS	LS LR	LS WS	WS LS	LS DB	DB LS												
		Client	Server	Client	Server																Avg. (%)	Max. (%)										
WLAN	645	11,706	141,738	140,738	12,156	1,527	11,925	1,527	11,925	8,289	75,937	8,289	75,937	1,527	11,926	1,617	12,054	0,128	0,091													
ADSL	250	7,652	81,092	80,636	7,961	-	6,843	-	6,843	-	43,540	-	43,540	-	6,843	-	6,909	-	-													
Cable	440	12,970	137,345	136,548	13,426	1,649	11,592	1,649	11,592	9,032	73,749	9,032	73,749	1,649	11,592	1,735	11,708	0,115	0,086													
VDSL	340	10,626	112,014	111,341	11,075	1,344	9,444	1,344	9,444	7,371	60,103	7,371	60,104	1,344	9,445	1,419	9,538	0,094	0,076													

IST	NoU	Queueing Delay														Packet Loss Rate	
		Clients RS	RS Clients	RS RR	RR RS	RR Cloud	Cloud RR	Cloud LR	LR Cloud	LR LS	LS LR	LS WS	WS LS	DB LS	LS DB		
WLAN	645	0.00015	0.00096	0.00015	0.00096	0.00128	0.00612	0.00079	0.28794	0.00015	0.00096	0.00014	0.00089	0.00009	0.00012	0	
ADSL	250	-	0.000900	-	0.000900	0.001536	0.005721	-	0.085471	-	0.000908	-	0.000862	-	-	0	
Cable	440	0.00023	0.00092	0.00015	0.00092	0.00168	0.00587	0.00083	0.25970	0.00015	0.00092	0.00015	0.00086	0.00009	0.00012	0	
VDSL	340	0.00021	0.00093	0.00015	0.00093	0.00136	0.00595	0.00084	0.10694	0.00015	0.00093	0.00015	0.00087	0.00010	0.00012	0	

IST: Internet Services Types, NoU: Number of Users, WLAN: Wireless LAN, RS: Remote Switch, RR: Remote Router, LR: Local Router, LS: Local Switch, WS: Web Server, DB: Database Server

Case 2: Hybrid Simulation for Shopping Mall System

IST	NoU	CPU Utilization		Client Response Time		Throughput						Traffic Sent (bytes/sec)		Traffic Received (bytes/sec)		LAN		
		Clients	WS			Clients ↕ Switch	Switch ↕ Clients	Switch ↕ Switch1	Switch ↕ Switch	Switch1 ↕ WS	WS ↕ Switch1					Delay	Inbound Traffic	Outbound Traffic
		Avg. (%)	Avg. (%)	Avg. (sec)	Max. (sec)	Avg. (bits/sec)	Avg. (bits/sec)	Avg. (bits/sec)	Avg. (bits/sec)	Avg. (bits/sec)	Avg. (bits/sec)	Avg. (sec)	Avg. (bits/sec)	Avg. (bits/sec)				
ADSL	950	94,238	88,779	80.07	124.73	2,019,791	7,873,338	2,023,460	7,864,891	2,142,167	7,747,732	26,323	732,745	728,837	26,323	0.038785	20,273,134	14,560,736
Cable	500	93,742	86,392	86.00	122.00	1,873,282	7,259,701	1,874,798	7,267,262	1,874,798	7,256,906	13,415	661,339	656,338	13,415	0.030195	19,636,032	14,384,329
VDSL	580	96,589	84,820	86,724	99.912	1,905,377	7,364,357	1,908,081	7,363,548	1,905,886	7,365,611	16,753	666,868	666,851	16,753	0.022481	20,123,585	14,808,979
IST	NoU	Utilization						Queueing Delay						Packet Loss Rate				
		Clients ↕ Switch	Switch ↕ Clients	Switch ↕ Switch1	Switch ↕ Switch	Switch1 ↕ WS	WS ↕ Switch1	Clients ↕ Switch	Switch ↕ Clients	Switch ↕ Switch	Switch1 ↕ WS	WS ↕ Switch1						
		Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)					
ADSL	950	20.198	78.733	20.235	78.649	21.422	77.477	0.000216	0.012941	0.012175	0.000178	0.010417	0					
Cable	500	18.733	72.597	18.748	72.673	18.724	72.569	0.000167	0.008785	0.008654	0.000146	0.010964	0					
VDSL	580	19.054	73.644	19.081	73.635	19.059	73.656	0.0001752	0.0060294	0.0061826	0.0001534	0.0081762	0					

IST: Internet Services Types, NoU: Number of Users, WLAN: Wireless LAN, RS: Remote Switch, RR: Remote Router, LR: Local Router, LS: Local Switch, WS: Web Server, DB: Database Server

Appendix 1: Results for the Maximum Capacity of Concurrent Users (2)

Case 2: Hybrid Simulation for Ticketing System

IST	NoU	CPU Utilization			Client Response Time		Throughput								Traffic Sent (bytes/sec)			Traffic Received (bytes/sec)			LAN		
		Clients	Web Server	DB Server			Clients	Switch	Switch	Switch	Switch	Switch	WS	Switch							DB	Delay	Inbound Traffic
		Avg. (%)	Avg. (%)	Avg. (%)	Avg. (sec)	Max. (sec)	Avg. (bits/sec)	Avg. (bits/sec)	Avg. (bits/sec)	Avg. (bits/sec)	Avg. (bits/sec)	Avg. (bits/sec)	Avg. (bits/sec)	Avg. (bits/sec)	Avg. (bits/sec)	Avg. (sec)	Avg. (bits/sec)	Avg. (bits/sec)					
WLAN	1800	93.994	91.018	80.165	73.08	106.63	2,036,223	5,193,861	5,193,001	2,066,796	5,223,924	1,635,150	1,749,550	1,749,550	33,659	411,355	1,298	408,705	34,945	2,491	0.195620	4,871,327	1,847,714
ADSL	2300	94.252	83.946	83.318	80.519	91.711	2,378,540	7,348,616	7,354,449	2,423,116	7,409,127	1,654,725	1,775,040	1,775,040	62,551	666,006	2,572	662,221	65,069	3,782	0.044906	19,766,316	14,937,985
Cable	2900	97.084	95.014	87.420	78.08	120.36	2,666,723	9,302,150	2,671,675	2,886,536	9,210,134	1,842,386	1,656,171	1,656,171	84,436	893,537	2,975	888,397	87,409	5,039	0.159450	22,113,339	15,630,744
VDSL	2800	94.142	81.504	84.545	78.569	98.253	2,536,844	8,439,978	2,537,386	2,720,096	8,381,351	1,793,459	1,654,875	1,654,875	74,689	796,072	3,147	791,384	77,812	4,454	0.064161	20,841,853	15,061,779

IST	NoU	Utilization									Queueing Delay						Packet Loss Rate	
		Clients	Switch	Switch	Switch	Switch	WS	Switch	DB	Switch	Clients	Switch	Switch	Switch	Switch	WS		Switch
		Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)	Avg. (%)		Avg. (%)
WLAN	1800	20.362	51.939	20.471	51.930	20.668	52.239	16.352	17.495	0.0002198	0.0014083	0.0002129	0.0015210	-	0.0018806	0		
ADSL	2300	23.785	73.486	23.764	73.547	24.231	74.091	16.547	17.750	0.0003710	0.0051160	0.0002440	0.0061800	-	0.0217440	0		
Cable	2900	26.667	93.022	26.717	93.053	28.865	92.101	18.424	16.562	0.0004600	0.0234490	-	0.0649430	0.000257	0.0584950	0		
VDSL	2800	25.368	84.400	25.374	84.373	27.201	83.814	17.935	16.549	0.0004280	0.0120510	-	0.0190820	0.000247	0.0223300	0		

IST: Internet Services Types, NoU: Number of Users, WLAN: Wireless LAN, RS: Remote Switch, RR: Remote Router, LR: Local Router, LS: Local Switch, WS: Web Server, DB: Database Server