# QNS – AN ONLINE SYSTEM FOR THE STUDY OF QUEUING MODELS

Hemamalini Sankar and Jozo J. Dujmović
Department of Computer Science
San Francisco State University

*QNS is an online system for the self-study of queuing models. QNS is publicly available on the Internet. It offers a diverse educational support including theoretical presentation of material, queuing theory models, graphical simulator with queuing network animation, laboratory experiments based on numerical solver, quiz subsystem with automatic grading, control system with GUI, and remote access support. The paper presents the structure and the use of QNS, as well as the basics of its implementation.*

## 1. Introduction

Queuing models are fundamental for understanding performance of computer systems. Our goal is to provide cooperating online systems for the self-study of performance models based on queuing theory. Our online educational support includes the following three systems:

- QNAS – Queuing network animation and simulation system [PAN02, PD02, QNAS]
- RAND – Random number generation and testing system [LIU03, RAND]
- QNS – Queuing networks online study system [SAN05, QNS]

QNAS is based on a discrete event simulator that is used both for deriving numerical performance indicators and for animation of queuing phenomena. The animation subsystem is used for visualization of dynamic behavior of queuing networks. It shows customers visiting service centers and receiving service from servers. The animation subsystem can operate at various speeds of animation. In the fastest mode, animation is suspended and the discrete event simulator is used for creating numerical results.

All events that occur in queuing systems are random. Good understanding of the nature of randomness is a prerequisite for the study of queuing models. RAND is a self-study system for random number generation and testing.

In this paper we present QNS, the third component of our educational support. QNS is based on the online course model presented in Fig. 1. At the beginning and at the end of all activities the user consults the table of contents which serves as a roadmap for selecting topics and assessment of progress in the study of presented material.
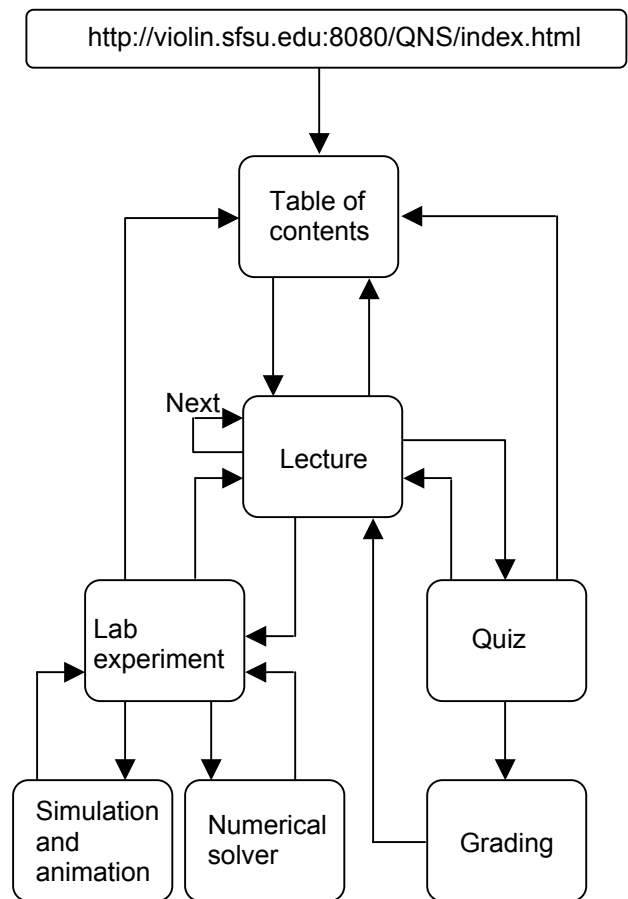


Figure 1. The online course model

Theoretical study of queuing systems is very difficult unless mathematical models come on top of a clear intuitive understanding of the dynamics of queuing phenomena. This is a reason why the lab is a very important component of our model. It is designed

to support two goals:

- Building intuitive understanding of queuing phenomena through animation of the most important queuing models.
- Numerical experiments based on analytic models. They are necessary to complement and verify the progress of theoretical study.

The animation subsystem is important because it provides visual basis for intuitive understanding of queuing phenomena, such as *random variations of queue length, service times, response times, and wait times, as well as service disciplines, cooperation of parallel servers, bottlenecks and saturation phenomena, load dependent phenomena*, etc. It also helps to visualize and understand the mean value performance indicators, such as *the mean utilization of a server, throughput, response time, queue length*, etc.

The animation is included in QNS to make sure that users clearly differentiate mean values of random variables from transient phenomena. Queuing models can be used without having substantial mathematical background, but cannot be used without a thorough intuitive understanding of random phenomena and dynamics of random events. This was the reason for developing QNAS, and this was the reason for including animation elements in QNS.

Quiz and grading subsystems are provided as tools that enable users to control their progress. Since QNS is designed primarily for the self-study, the role of quiz and grading subsystem is not to provide global grades and issue certificates.

## 2. The structure of QNS tutorial

The structure of QNS tutorial is shown in Fig. 2. It consists of six sections that cover the most frequently used queuing models. The first section is an introduction that provides description of queuing systems, and an exhaustive list of definitions of all performance indicators. It also provides notation and symbols used in queuing analysis.

Our presentation includes a spectrum of queuing models and examples of their use, but *not* their mathematical derivation. The majority of derivations require sophisticated mathematical background, and QNS is designed for practitioners, who regularly have no interest in mathematical derivations, and sometimes have no background in probability theory and calculus, that is indispensable for understanding the derivations of various formulas. Instead of focusing on origins of analytic models we suggest that QNS users experiment with QNAS to develop strong intuitive background. Simultaneous use of QNAS and QNS is useful to understand differences between

analytic and simulation models, and to benefit from experimenting with queuing models.

---

1. **Introduction**
   1.1 Basic queuing concepts
   1.2 A simple queuing system
   1.3 Classifications and notations of queuing systems
   1.4 Graphical representation of queuing systems
2. **Poisson arrival process**
3. **Service station models**
   3.1 Survey of service station models
   3.2 M/M/1 queuing model
   3.3 M/G/1 queuing model
   3.4 M/D/1 queuing model
   3.5 GI/G/1 queuing model
   3.6 M/M/k queuing model
4. **Open queuing networks**
   4.1 Stochastic modeling of open queuing networks
   4.2 Elementary queuing network models
   4.3 Operational modeling of open queuing networks
5. **Closed queuing networks**
   5.1 Stochastic modeling of closed queuing networks
   5.2 Operational modeling of closed queuing networks
6. **Mean value analysis (MVA)**
   6.1 Load-independent MVA
   6.2 Load-dependent MVA

---

Figure 2. The structure of QNS queuing model tutorial

The process of independent and random arrivals of customers is provided in the second section. The Poisson flow of service requests is a fundamental prerequisite for understanding material presented in other sections.

The simplest queuing models are the single service station models described in the third section. This presentation shows the relationship between the statistic properties of the arrival and service processes and the complexity of analytic models. It starts with the simplest M/M/1 model that reflects the arrival process with exponential interarrival times, and exponentially distributed service times. Then we introduce more complex cases including the cases of general independent arrival process and general
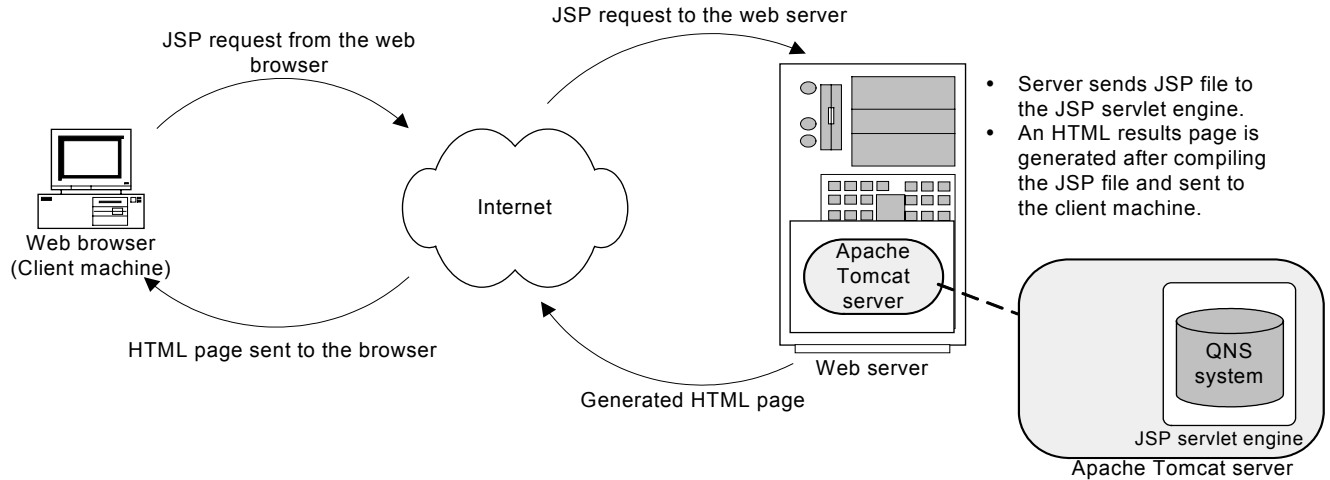
Figure 3. Generation of dynamic HTML pages using JSP technology

service time distribution (approximate GI/G/1 model). At the end of the third section we study service centers with multiple servers (the M/M/k model).

Models of open queuing networks are introduced in the fourth section. We present simple serial, parallel, and cyclic network models, and then use them to develop general models of queuing networks. Included are both stochastic and operational models. Stochastic models study systems over infinite observation intervals, while operational models focus on finite (and usually short) observation intervals during which it is possible to count frequencies of service events.

Closed queuing networks are presented in the sixth section. This presentation introduces the role of feedback in queuing systems and shows general stochastic and operational models.

The mean value analysis is a topic presented in the sixth chapter. Multiple parallel servers, caching, and operating system optimization are main sources of load dependent service behavior of real servers in computer systems. Consequently, we present both the load-independent and the load-dependent mean value analysis. These algorithms are important for analytic modeling of many real-life computer systems.

### 3. Internal structure of QNS system

QNS is implemented using Java, Microsoft FrontPage (HTML), JSP, and Tomcat. The animation subsystem and the numerical solver are developed using Java and are activated by a web browser as Java applets.

Being a web-based tool, the presented system can be used by anyone with a computer and an Internet connection. The system is useful for users of various

skill levels. Beginners can learn by reading the course materials and working with the experimental tools, whereas advanced users, who are already familiar with the concepts of queuing networks, can select a topic of their choice to learn or go directly to the experimental tools and learn by working with it. Moreover, this tool helps a user to learn at his/her own pace. The user can go over the materials repeatedly or s/he can randomly select a topic to read, without having to go through the material in a sequential manner.

Once the QNS system is activated, the graphical animation and simulation tools are automatically downloaded to the local system and run completely within the client machine, reducing network traffic between the client machine and the server. The Java *.class* files are compressed in to a *.jar* file hence decreasing the download time. In addition, the *.jar* file is downloaded in a single HTTP transaction thereby minimizing the need for opening new connection for each file.

The Java applet consists of a main frame, which shows the simulation/animation of the respective queuing model and a side frame, which is an input panel for changing simulation/animation parameters. The input panel includes a results button that activates the QNS numerical solver, and displays a window containing numerical results. The simulation/animation system is programmed using Java where the animation is created using Swing and Java AWT toolkit and the simulation is done using respective distributions (exponential, deterministic, etc.) for inter-arrival and service times. The numerical results are generated using Java programs that compute respective output parameters for that particular simulation/animation system. More detailed

simulation and animation studies can be performed using the QNAS system [PD02, QNAS].

The tutorial materials are developed as web pages using Microsoft FrontPage. They are static HTML pages. The user can easily navigate between the chapters, simulation/animation tools, quiz and grading subsystem, and numerical solver. The navigation system points out the exact location in the tutorial and provides hyperlinks to other chapters, respective quizzes and experimental tools.

The quiz and grading subsystem are implemented using the JSP (Java Server Pages) technology [HALL01, HB04]. The user input for quizzes is processed by the quiz and grading subsystem and JSP is used to generate a dynamic HTML page containing results. The results page gives a side-by-side comparison of user input and correct answers, and displays a percentage score. The quiz subsystem is organized as multiple choice questions using form layout in HTML. Once the user clicks the SUBMIT button in the quiz page, the user input is collected from the quiz page and is forwarded to the JSP servlet engine. Then the user's answers are compared with the correct answers and a dynamic HTML page is generated by the JSP servlet engine and is sent to the client. Finally, the results are displayed at the client side.

Apache Tomcat is used as a web server for the HTML pages and Java Server Pages [TOMCAT, HB04]. Fig. 3 illustrates the interaction between the client and the web server for generating dynamic HTML pages using JSP technology.

The entire application, which includes the HTML files (with file extension *.html*), Java Server Pages (with file extension *.jsp*), compressed Java class files (with file extension *.jar*) and image files (with file extension *.jpeg* and *.gif*), is packaged as a web using Microsoft FrontPage and is stored under the Apache Tomcat's *ROOT* directory. The default HTTP connector port for Tomcat is set to 8080. Hence, to access a webpage, say with file name, index.html, stored in the server with domain name violin.sfsu.edu, the URL will be http://violin.sfsu.edu:8080/index.html.

## 4.   Lab and animation subsystems

There are many phenomena that can be observed and intuitively understood using an animation system. They include:
- The job arrival pattern
- Random interarrival times
- Random think and serve times
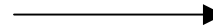- Random queue length
- Throughput

- Server utilization
- Server idling
- Server saturation
- Bottleneck resources
- Average number of jobs in a system
- Probabilistic branching and routing
- Effects of feedback in queuing networks
- Load balancing

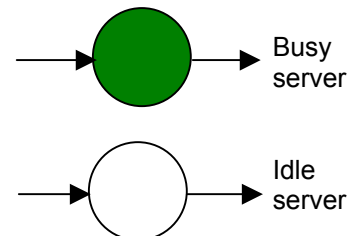The animation subsystem uses the following components:
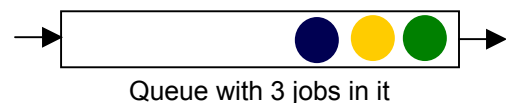- Each job is represented as a circle and is color-coded. E.g.:



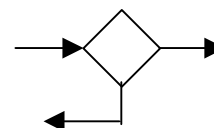- Direction of the job flow is denoted using arrows.



- Each server is represented as a bigger circle: a busy server is shown with the color of the job that is currently served and an idle server is shown empty:



Busy server

Idle server

- A queue is denoted by a rectangle. The jobs currently in the queue are shown as occupying specific positions in the rectangle:



Queue with 3 jobs in it

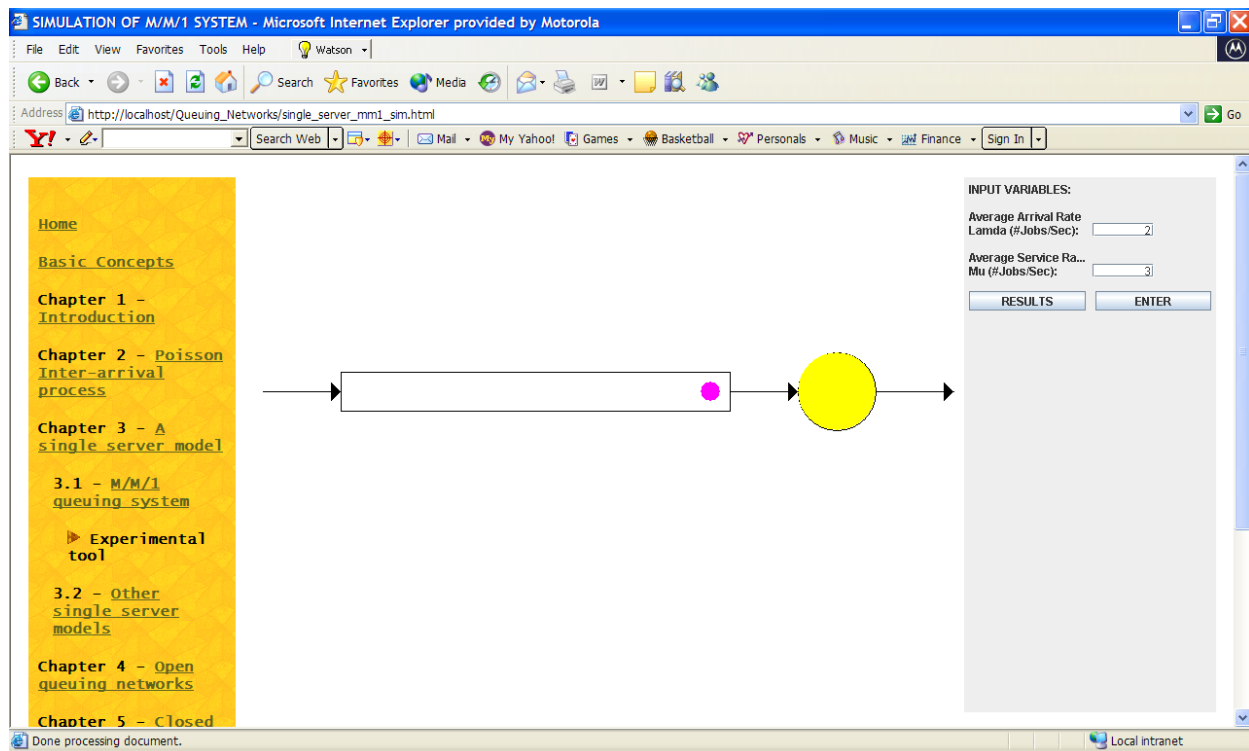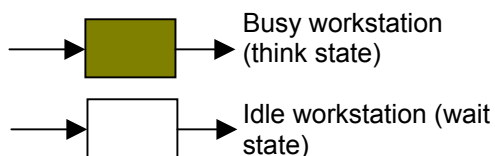- A probabilistic branching element is represented using a diamond:

Figure 4. An example of the M/M/1 queuing system

- Interactive workstations are denoted by a small rectangle. Similar to servers, if a workstation is busy then it shows the color of the job that currently resides in the workstation; otherwise, the workstation is empty:



Busy workstation (think state)

Idle workstation (wait state)



Figure 5. Numerical results for the M/M/1 queuing system example (Fig. 4)

The simulation/animation tools are developed for the following systems:
- M/M/1 queuing system.
- M/D/1 queuing system.
- Open queuing network.
- Closed queuing network with I/O devices
- Interactive, closed queuing network.
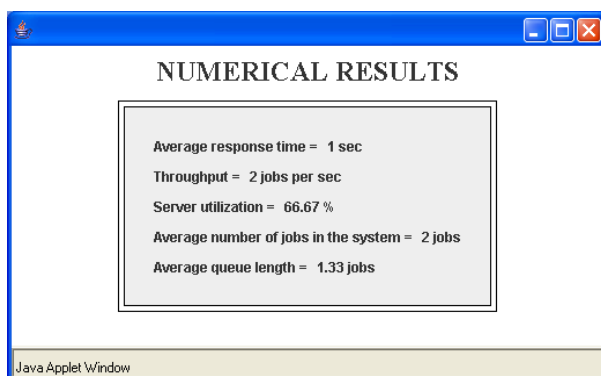- Load-dependent and load-independent closed queuing network with I/O devices.

M/M/1 queuing system. A sample M/M/1 system is shown in Fig. 4 and its numerical results are shown in Fig. 5. The M/M/1 queuing system has a single server and a FCFS (first-come-first-served) queue. The arrival of jobs follows Poisson arrival process. The arrival rate, $\lambda$, and the service rate, $\mu$, are exponentially distributed. The animation of M/M/1 queuing system can be viewed for any $\lambda$ and $\mu$ value. The numerical solver can be activated by selecting the 'Results' button. The numerical solver displays the output values such as the average response time, system throughput, server utilization, the average number of jobs in the system, and the average queue length.
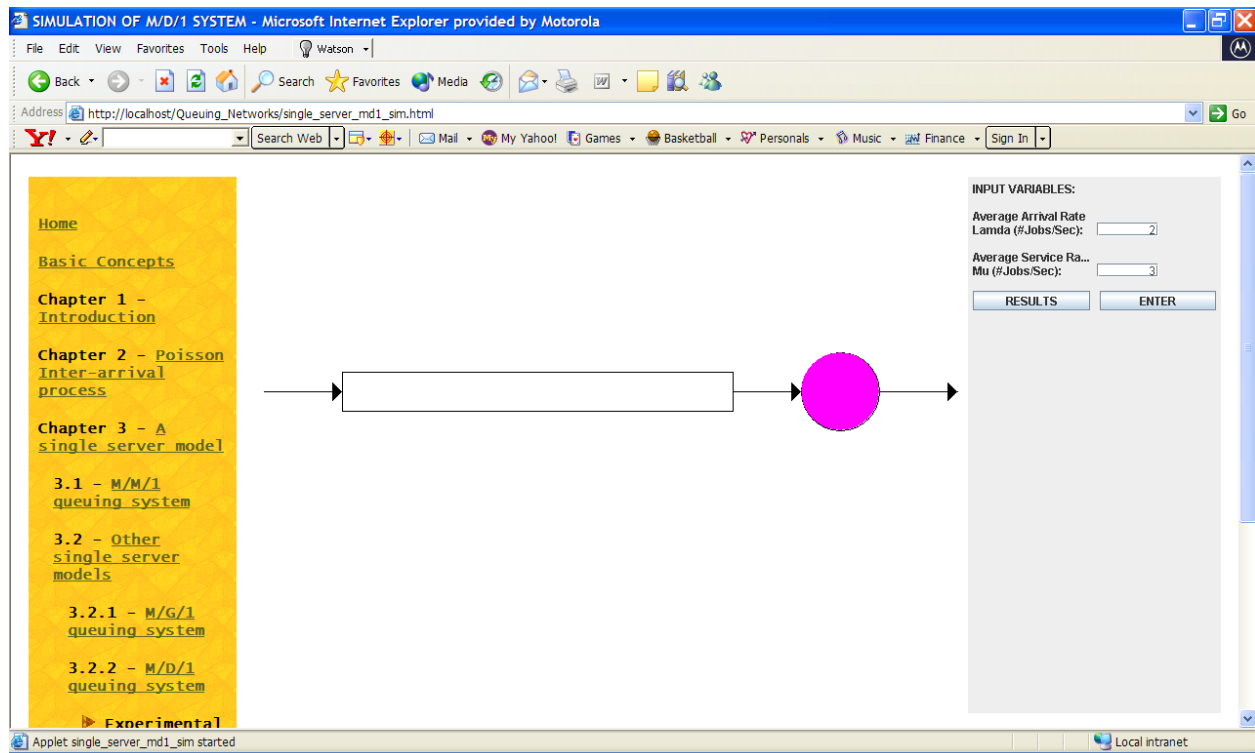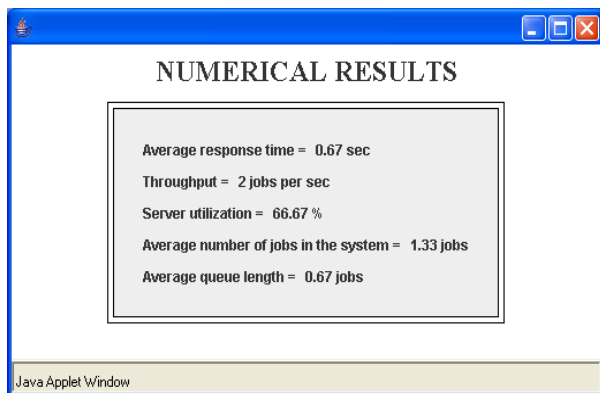
Figure 6. An example of the M/D/1 queuing system



Figure 7. Numerical results for the M/D/1 queuing system example (Fig. 6)

M/D/1 queuing system. Similar to the M/M/1 queuing system, the M/D/1 queuing system is also a single server system, i.e. a server and a queue. The arriving of jobs follows the Poisson arrival process. The arrival rate, $\lambda$, is exponentially distributed, whereas the service rate, $\mu$, is a constant. Similar to M/M/1 system, the queue-scheduling algorithm used in this system is FCFS. The animation of M/D/1 queuing system can be viewed for any $\lambda$ and $\mu$ value. The numerical solver displays the values of basic performance indicators: the average response time, throughput, server utilization, average number of jobs in the system and average queue length. A sample M/D/1 system is shown in Fig. 6 and its numerical results are shown in Fig. 7.

Open queuing networks. An open queuing network is a network of interconnected serial, parallel, and feedback service centers. The network has one or more inputs and one output. QNS uses a sample feedback open queuing network shown in Fig. 8. The number of service centers can range from 1 to 5. After the user enters the number of servers and activates the 'Enter' button, the input panel displays the default values of all the input variables. The user can alter the input variable values to view the simulation under different input scenarios. In cases of 2 or more servers, a job leaving a server can visit any other server with equal probability. The output parameters calculated in this case include: the average response time, throughput, server utilization, average number of jobs in the system and average queue length. An open queuing model with five servers and
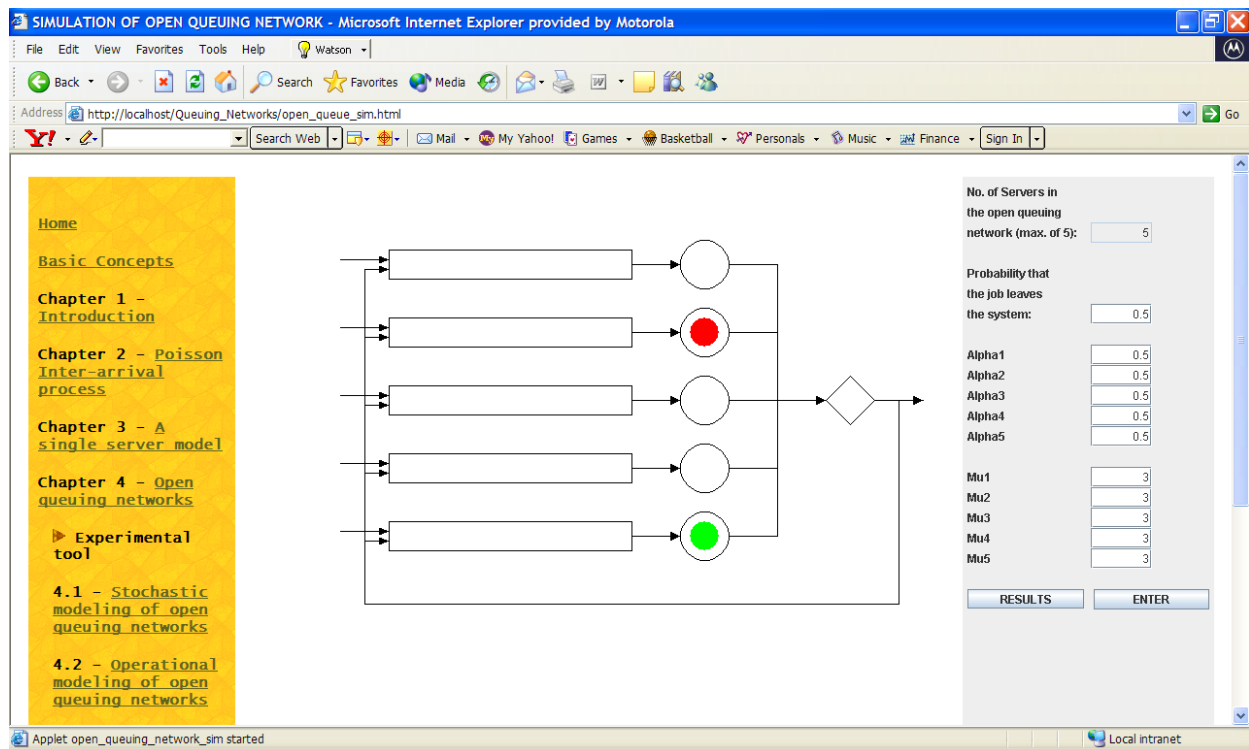
Figure 8. An example of open queuing network

its numerical results are shown in Fig. 8 and Fig. 9, respectively.
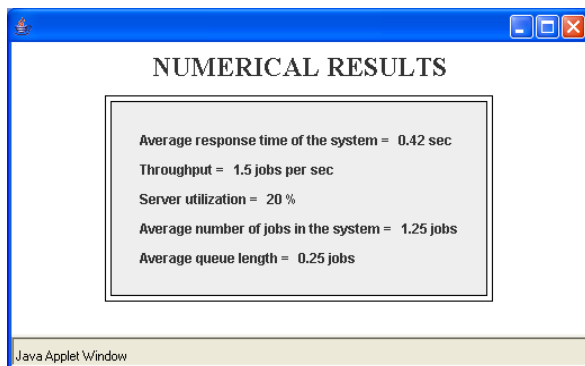


Figure 9. Numerical results for the open queuing network example (Fig. 8)

Closed queuing networks with I/O devices. A closed queuing model with five servers and five disks, and its numerical results are shown in Fig. 10 and Fig. 11, respectively. A closed queuing network is a queuing network, which has no input and no output. The number of jobs in the network is assumed to be a constant. The experimental tool developed in QNS is a closed queuing network with file I/O devices. In this network the number of servers can range from 1 to 5 and the number of file I/O devices (disks) can range

from 1 to 5. Similar to open queuing network, when the user activates the 'Enter' button, after selecting number of servers, number of disks and number of jobs, the input panel displays the default values of all input variables. The user can alter the input variable values to view the simulation under different input scenarios. The output performance indicators for this model are: overall residence time of the server(s), residence time of each disk, average response time of the system, throughput, server utilization, utilization of each disk, and average queue length in each disk queue.

Closed queuing networks - interactive systems. An interactive system is a closed queuing system that has one or more processors, a processor queue and one or more workstations. Workstations are nothing but delay centers, which has a think time, Z, and it is assumed that exactly one job is generated from each workstation. Therefore, if there are $n$ workstations, then there will be exactly $n$ jobs in the system. In this system, the number of servers can range from 1 to 5 and the number of workstations can range from 1 to 10. The output parameters for this model are: the average response time of the system, throughput, server utilization, and critical number of users. The interactive closed queuing model with five
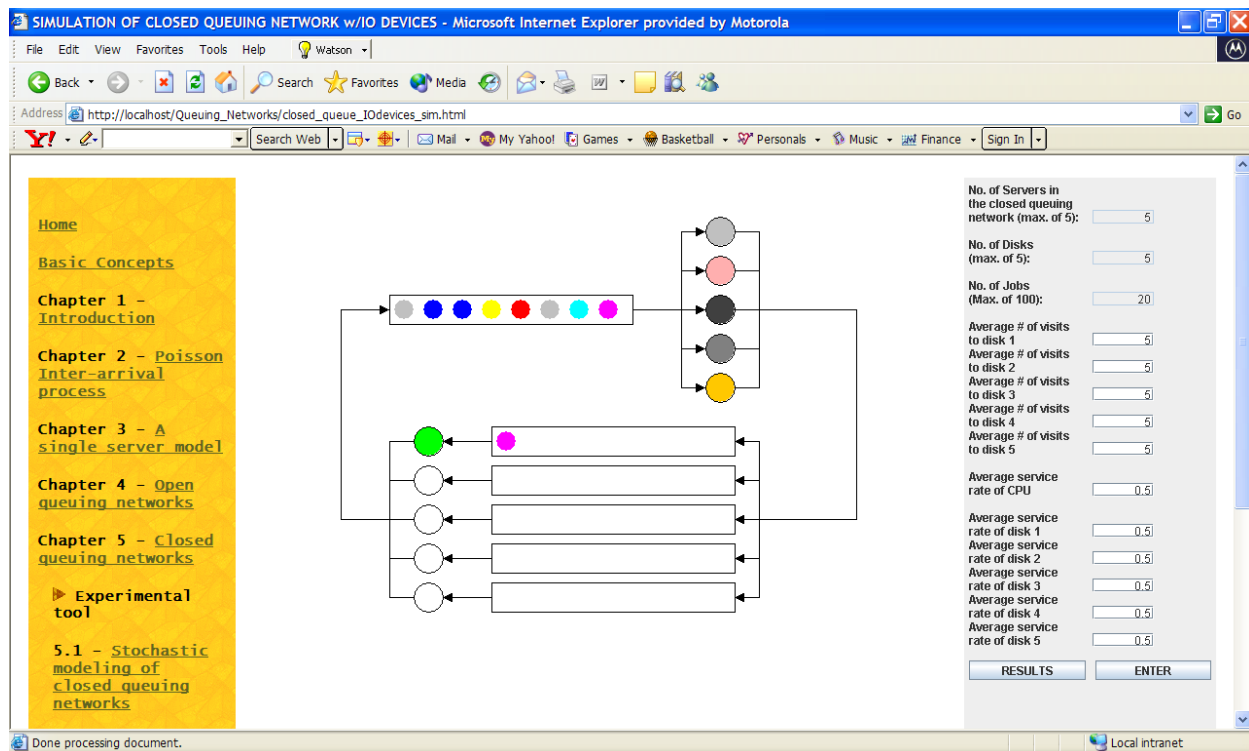
7

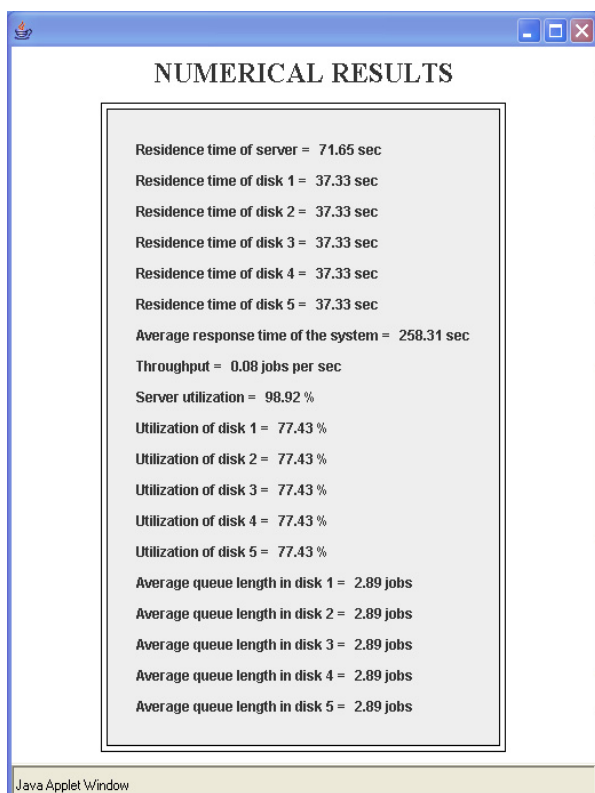Figure 10. An example of closed queuing network with I/O devices



Figure 11. Numerical results for the closed queuing network with I/O devices (Fig. 10)

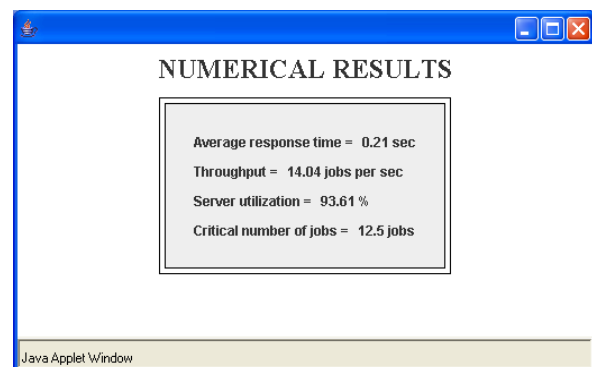servers and ten workstations, and its numerical results are shown in Fig. 13 and Fig. 12, respectively.



Figure 12. Numerical results for the closed queuing network model of an interactive system (Fig. 13)

Closed queuing networks – MVA. This system is a closed queuing system that has one or more processors, a processor queue, one or more disks and one or more workstations. This system is a combination of closed queuing system with I/O devices and interactive system, since it has both the workstations and disks. The jobs are generated from the workstation(s). Therefore, if there are $n$ workstations then there will be exactly $n$ jobs in the
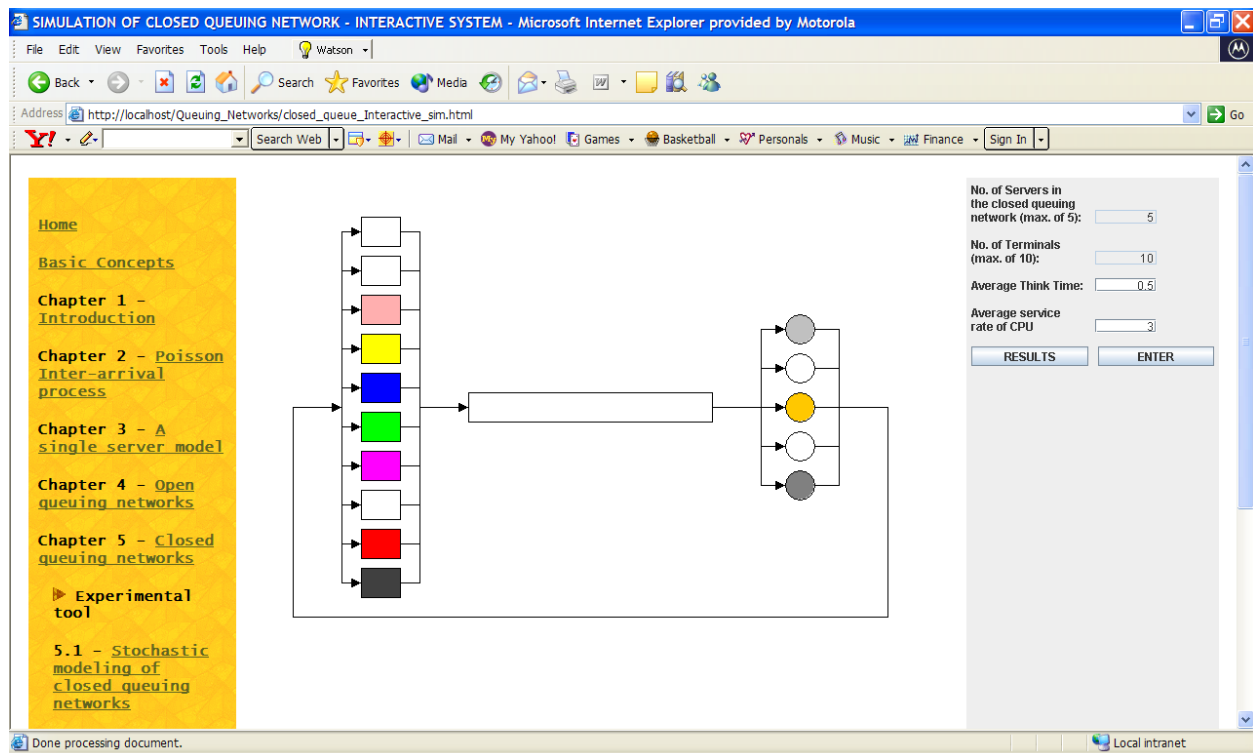
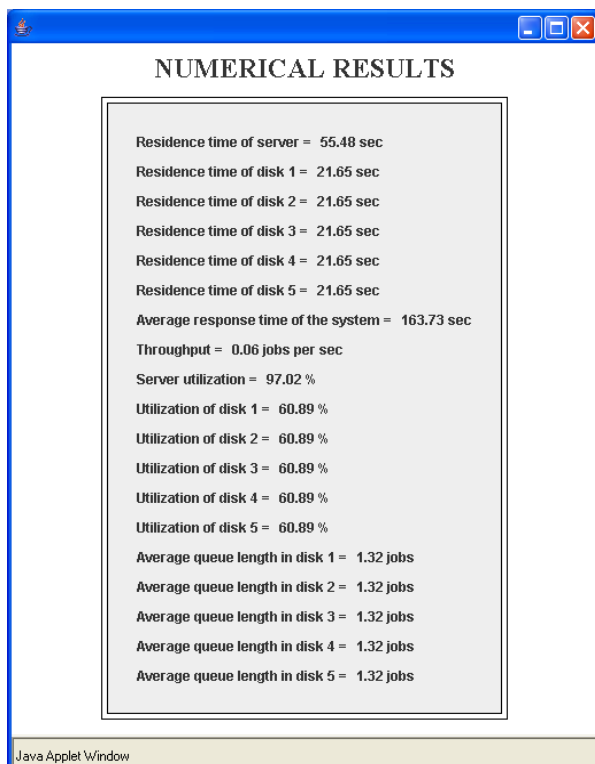Figure 13. A closed queuing network model of an interactive system



Figure 14. Numerical results for the closed queuing network model based on mean value analysis (Fig. 15)

system. In our sample system the number of servers can range from 1 to 5, the number of disks can range from 1 to 5 and the number of workstations can range from 1 to 10. The output parameters for this model are: the overall residence time of the server(s), the residence time of each disk, the average response time of the system, throughput, server utilization, utilization of each disk, and average queue length of each disk queue. The MVA model with five servers, five disks and ten workstations, and its numerical results are shown in Fig. 15 and Fig. 14, respectively.

## 5.  Quiz and grading subsystem

The quiz subsystem is used for testing the knowledge of the user in a selected area. The chapter quiz is organized into multiple-choice questions in which the user can select one of the several options listed for a question. After selecting the answers for all the quiz questions, the user should select the 'Finish' button. The user can reset the quiz and start over by clicking the 'Reset' button. After selecting the 'Finish' button, the corresponding JSP page is invoked from the server. The user's selected answers are passed as parameters to the JSP page. The user's input is compared to the correct answers and the quiz results page is generated and sent to the client machine.
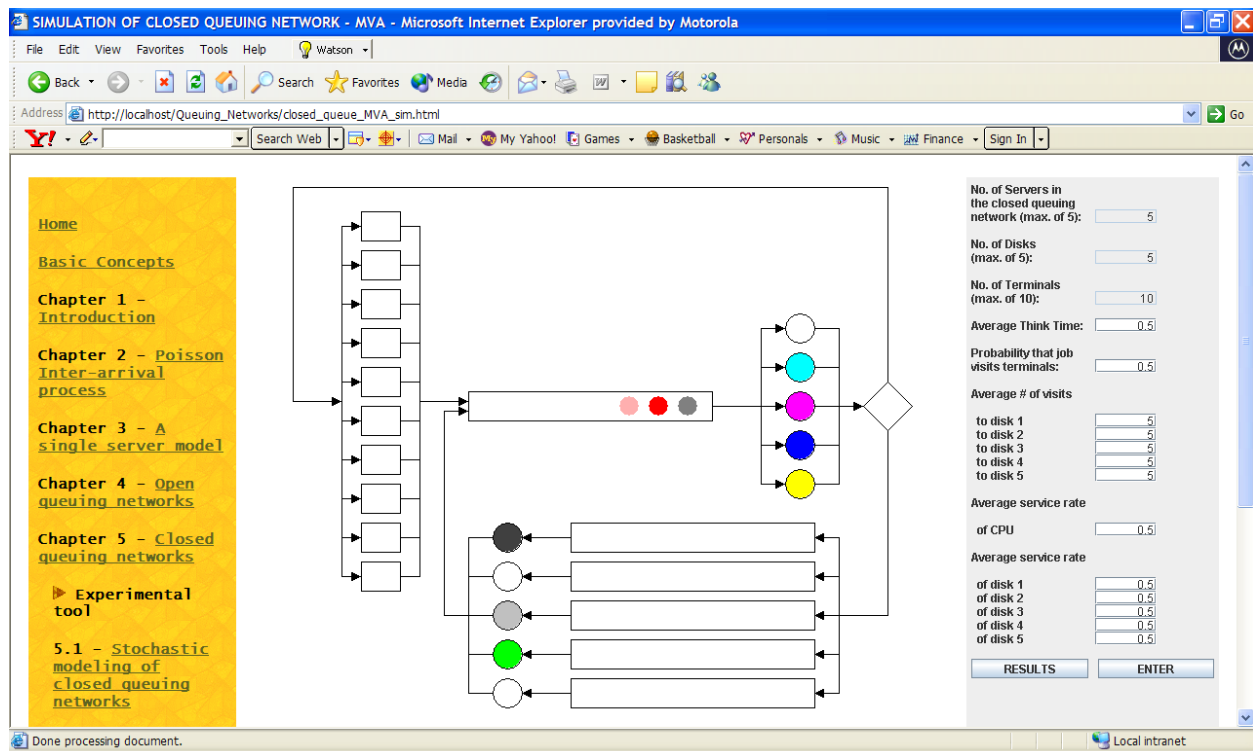
Figure 15. An example of closed queuing network model based on mean value analysis

The user has an option of taking the quiz anytime while reading the chapter. This enables navigation in a random way instead of going through all the material and quiz in a sequential way. Further, the quiz subsystem consists of both theoretical questions and numerical problems.
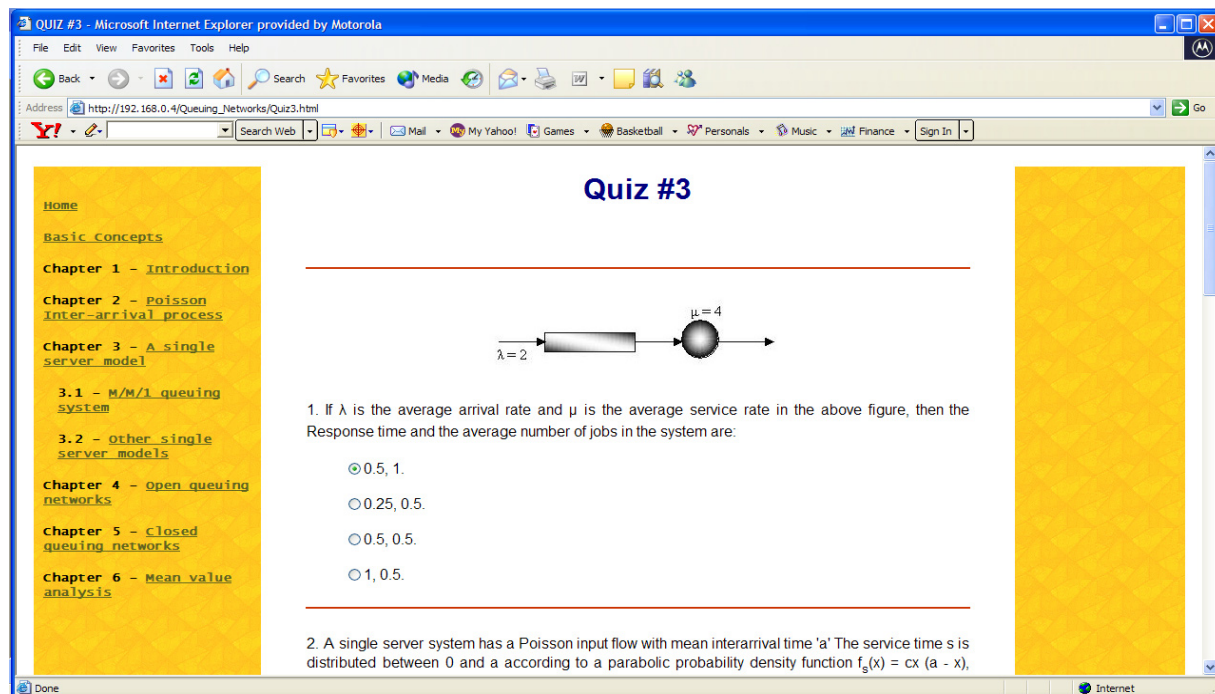


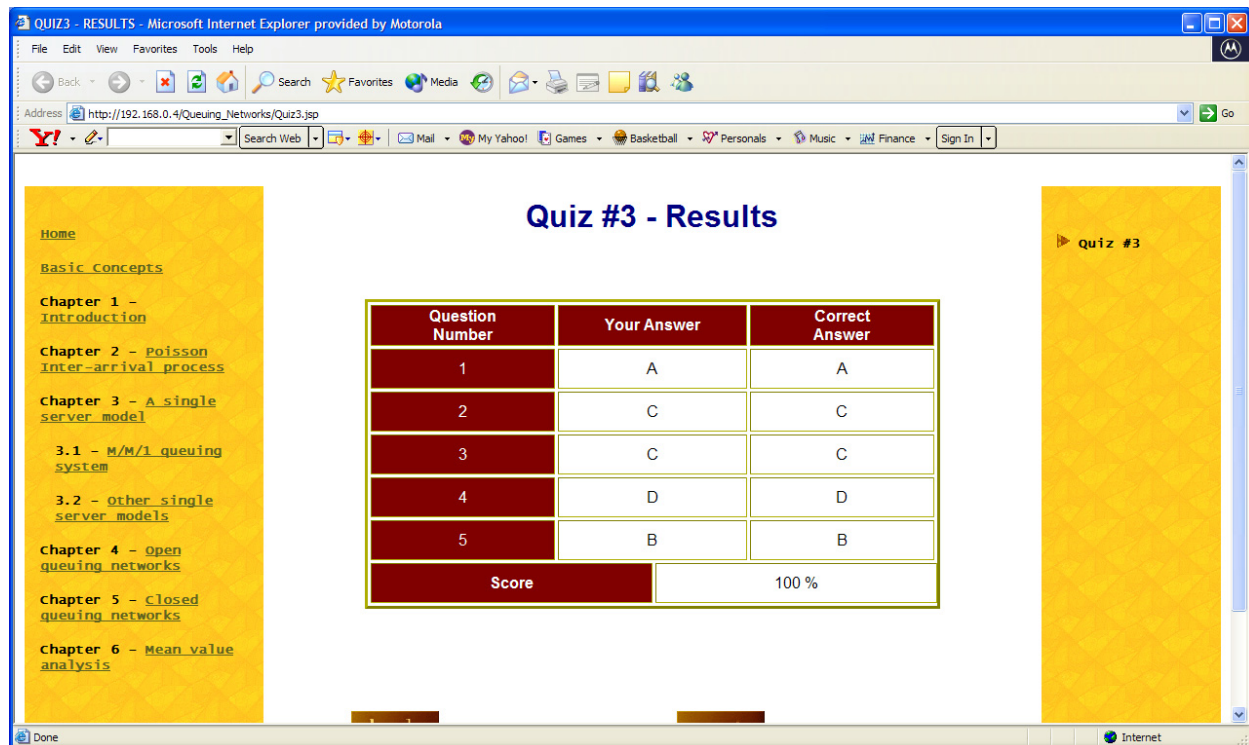Figure 16. An example of the quiz subsystem window

Figure 17. Results for the quiz shown in Fig. 16

A sample quiz and the results for that quiz are shown in Fig. 16 and Fig. 17 respectively.

## 6. System maintenance and evolution

System maintenance of QNS is relatively simple. The Apache Tomcat server should be kept running in order to access the system. It is important to provide support for evolution and updating of QNS. Adding some of the interesting features listed below can further enhance this system:

- This system can be enhanced by adding more queuing network models and by adding more simulators and animators of queuing networks.
- Adding numerical solvers for various additional queuing network models can enhance the system.
- For each presented model it might be useful to add a "derivation button" that user can press to see a complete mathematical derivation of all formulas.
- The quiz subsystem can be expanded furthermore by having more than one quiz per section based on the difficulty level the user chooses. In addition, a final exam (and maybe midterm exams) can be added at the end covering the entire course material. Also, the quizzes can be timed, so that users can practice fast solving of quizzes.

Of course, expansion directions for QNS (and related systems QNAS and RAND) will primarily depend on feedback generated by their users.

## 7. Conclusions

We presented the use and organization of an online educational system for the self-study of queuing models. Our model of online courses is rather general and includes a table of contents, lectures, lab experiments, and quizzes. Consequently, this course model can be used for organizing a variety of similar courses.

The material presented in lectures can be used as a supplement to regular computer performance evaluation courses that present mathematical details and derivations of analytic models. People who have some elementary understanding of queuing phenomena and analytic models can also use it for self-study. QNS can be more effective if used in conjunction with other complementary online systems, RAND, and QNAS. Future work on QNS can be directed in several areas: expansion of the number of presented models, more numerical examples, grading of students, and presentation of derivations of queuing models.

## References

[GUN00] Gunther, N., The Practical Performance Analyst. McGraw Hill, 2000.

[HALL01] Hall, M., More Servlets and JavaServer Pages. Free at http://pdf.coreservlets.com/

[HB04] Hall, M., and L. Brown, Core Servlets and Java Server Pages. Second Edition. Sun Microsystems Press/Prentice Hall, ISBN 0-13-009229-0, 2004.

[JAI91] Jain, R., The Art of Computer Systems Performance Analysis. J. Wiley 1991

[LIU03] Liu, W., An E-learning System for Studying the Quality of Randomness. MS project. Technical Report TR-03.21, Department of Computer Science, San Francisco State University (2003)

[MAD94] Menasce, D., V. Almeida, and L. Dowdy, Capacity Planning and Performance Modeling. Prentice Hall, 1994.

[PAN02a] Pan, F., QNAS – a Queuing Network Animation System. MS project. Department of Computer Science, San Francisco State University (2002)

[PD02] Pan, F. and J.J. Dujmović, QNAS – a Queuing Network Animation System. CMG 2002 Proceedings, Vol. 2 pp. 821-832 (2002)

[QNS] http://violin.sfsu.edu:8080/QNS/index.html

[QNAS] http://violin.sfsu.edu/QNAS

[RAND] http://violin.sfsu.edu:8080/MyWebs/homepage.htm

[SAN05] Sankar, H., An Online System for the Study of Queuing Networks (QNS). MS project. Technical Report TR-05.14, Department of Computer Science, San Francisco State University (2005)

[TOMCAT] Configuring & Using Apache Tomcat, http://www.coreservlets.com/Apache-Tomcat-Tutorial/

[TRI01] Trivedi, K., Probability and Statistics with Reliability, Queuing, and Computer Science Applications. Prentice Hall, 2001.