

# Performance Optimization of Distributed-System Models with Unreliable Servers

Ian F. Akyildiz, Senior Member IEEE  
Georgia Institute of Technology, Atlanta  
Wei Liu, Student Member IEEE  
Georgia Institute of Technology, Atlanta

**Key Words** — Distributed system, Performance evaluation, Queuing network, Performance measure, Optimization, Lagrange multiplier technique

**Reader Aids** —

**Purpose:** Widen the state of the art  
**Special math needed for explanations:** Queuing theory, real analysis  
**Special math needed to use results:** Same  
**Results useful to:** Design analysts, network planners

**Abstract** — Models of distributed systems with servers subject to break-down and repair are investigated for optimization of performance measures. The optimization problems are the cost minimization, response time minimization, and throughput maximization. The system is modeled by a preemptive-resume priority queueing network. The mean-value analysis algorithm is applied to derive a relationship between the multiprogramming level and performance measure formulas. Based on this relationship the Lagrange multiplier technique is applied to carry out the optimization of performance measures. Optimal service rates are obtained that reach a target throughput while minimizing the total cost. Servers are also treated individually in order to minimize the mean response time of a particular server in the system in order to find the optimal service rates which minimize the response time of a particular server while reaching a target throughput. Formulas are derived for determining the maximum throughput of the system; unfortunately it has no closed-form solution. However, it can be solved using the binary search and insert value method. Numerical examples illustrate the solutions.

## 1. INTRODUCTION

Advances in communication technology and the reduction in computer hardware costs have made distributed systems feasible. Distributed systems are configured to achieve a higher processing capacity than centralized systems. Cost-evaluation and performance-prediction are an important step in the planning and design of computer systems, distributed systems, computer networks, and flexible manufacturing systems. Queueing network models have received special interest for performance analysis and performance prediction in the past two decades. They can also be used for optimization of performance measures. In the optimization procedure, an objective function (eg, costs, throughput, usage, or response time) is obtained by appropriate selection of system input parameters, eg, service

rates, multiprogramming level, called decision variables. Queueing network models represent the relationship between the objective function and the system input parameters. Using these models it is possible to have an optimal design of the systems mentioned above.

System reliability is an important issue in evaluating the performance of systems [13]. Queueing network models can be applied to capture the break-down, repair, and recovery of computer systems. Vinod & Altioik [15] analyze queueing network models with server break-downs. They capture the failure and repair times by constructing Cox 2-phase (service and break/repair) for each server. From the Cox 2-phase they obtain an average service time for each server which is prolonged by the break-down and repair. They assume that the derived queueing network model has a product-form solution. The disadvantage of their approach is that each job in a server always encounters only one break-down during its service. However, in practice, several jobs could get serviced without the server breaking down, or occasionally, a server breaks-down several times on one job. Another disadvantage of the solution is the approximation itself. Although two examples are given with good accuracy, our tests verified that large deviations can occur. Recently, Ramanjaneyulu & Sarma [9] modeled servers by constructing a queueing network model with preemptive-resume servers and multiple job classes. Their model simulates the essential behavior of unreliable systems. This approach is described in detail in section 3.

The modeling of unreliable servers by a multi-class-job queueing network with preemptive-resume scheduling discipline brings the problem of solving these queueing networks. Since there is no product-form solution for this type of queueing network, several researchers tried to solve them approximately [3, 5]. Bryant et al [3] generalize the solution for an M/M/1 with PR (preemptive resume) or HOL (head-of-line; nonpreemptive) to a priority server at the network. They embed the derived mean response-time to the classical mean-value analysis formula [10] by assuming that the arrival-instant theorem holds [10, 13]. Doremalen et al [5] introduce the Schweitzer/Bard [2, 11] approximation factor into the heuristic mean-value analysis formula [3], thus reducing the recursive computation. They consider the fact that the arrival-instant theorem is not valid in queueing-network models with priorities and adjust the approximation such that it will not overestimate the effect of the lower priority-job class on the higher priority-job class. When it is applied to our model, we use the job-flow pattern and the fact that there is a limited number of jobs in certain classes to simplify it greatly, as described in section 3.

After modeling the unreliable system and solving the model, we derive the relationship among the performance measures and the break-down time and repair rate. This enables us to consider the optimal decision on the systems. Several

authors have discussed the issue of optimization in recent years. Trivedi & Wagner [14] consider a computer configuration design problem where the computer system is modeled by a closed central-server model. The system throughput is the objective function to be maximized by proper choice of device speeds subject to a cost constraint. A nonlinear cost function is considered in the analysis. Chandy, Hogarth, Sauer [4] use a branch and bound algorithm to minimize the mean response-time subject to cost limitations. Von Mayrhauser & Trivedi [8] consider a configuration design problem where the computer system is modeled as a closed queueing-network. The mean response-time to an interactive user request is minimized and the speeds of the devices are the decision variables. Kenevan & von Mayrhauser [6] show that: 1) throughput is a log convex function of the number of items in a closed, single class, network, 2) reciprocal throughput is a convex function of the relative usage of the servers. Kobayashi & Gerla [7] determine the optimal routing in closed queueing network models with multiple classes of jobs. Stecke [2] investigates product-form networks in which she imposes a constraint on the total workload in the system. She shows that throughput is a function of the ratio of the service rate at a server to the sum of the workloads, not purely concave but rather quasi-concave. Akyildiz & Bolch [1] applied mean-value analysis as the optimization basis and derived closed-form solutions for optimal performance measures such as response time, throughput, and cost. The studies all assume that the servers are reliable. However, in practice the resources are prone to failures.

In this paper we study models with unreliable servers. In section 2 the model with unreliable servers is substituted by a model with reliable servers such that the resultant reliable queueing model is a multi-class job-queueing model with the preemptive-resume priority scheduling discipline. We develop a mean-value analysis algorithm for this model and derive the constraint for application of Lagrange multiplier technique in optimization of performance measures in section 4. In section 4.1 we obtain a formula for optimal service rates which provides the minimum cost. Response time minimization of a particular server is given in section 4.2. Section 4.3 contains the throughput maximization. Numerical examples illustrate the solutions.

**Notation**

- MVA mean-value analysis
- FCFS first-come first-served
- $N$  number of servers in the network
- $K$  total number of jobs
- $\kappa$   $(K-1)/K$
- $\mu_i, \theta_i, \psi_i$  service, failure, repair rate of server  $i$
- $p_{ij}$  probability that a job transfers from server  $i$  to server  $j$
- $\mu_{ir}$  service rate for job class  $r$  at server  $i$
- PR preemptive resume
- $\bar{t}_{i,r}, \bar{k}_{i,r}, e_{i,r}, \rho_{i,r}$  mean response-time, mean number, visit ratio, usage of job-class  $r$  at server  $i$
- $\lambda_r$  total network throughput for job-class  $r$
- $K_r$  the total number of class  $r$  jobs

- $C$  total budget
- $c_i$  cost unit for server  $i$

Other, standard notation is given in "Information & Authors" at the rear of each issue.

**2. DISTRIBUTED SYSTEM MODEL**

The distributed system in figure 1 consists of several workstations where each workstation consists of a number of resources (eg, CPUs, disks) which are used by processes that execute in that workstation. Consider, for example, a set of identical workstations that share a set of common resources, such as a file-server and printer, and let the behavior of any of these embedded workstations be the object of the performance analysis. The workstations exchange messages for storing and retrieving data files from the common file-servers. The performance of distributed system can effectively be modeled and predicted by queueing networks. However, performance measures of a distributed system tend to be distorted by irregularities caused by server breakdowns which affect the performance of the system. The presence of unreliable servers prompts the need for performance models that consider the reliability of the system.

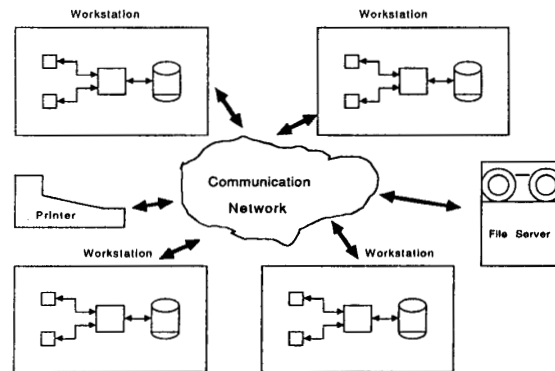


Figure 1. A Distributed System

This distributed system could be represented by a closed queueing-network model employing one server and one job class associated with each workstation. The model consists of  $N$  servers of a distributed system running on a communication network. The servers are prone to time-dependent failures followed by repair times. The failure, repair, and service times are exponentially distributed. The scheduling discipline in each server is FCFS. The application programs are modeled by a closed chain for the entire queueing network. The number of circulating programs (henceforth, jobs) in the network is fixed.

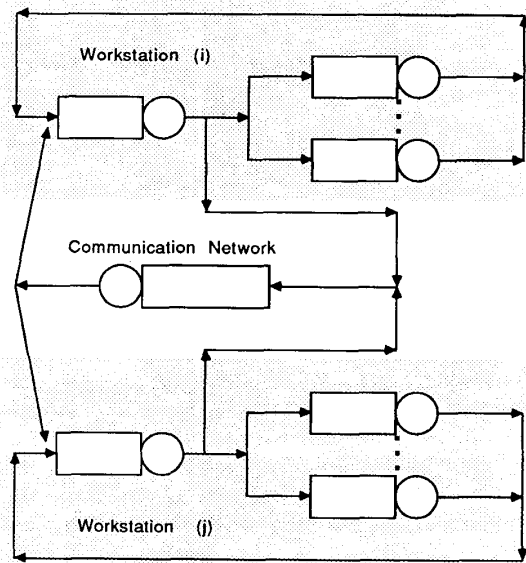


Figure 2. Queuing Network Model of a Distributed System

### 3. MEAN VALUE ANALYSIS FOR UNRELIABLE SERVERS

The unreliable server can be modeled by constructing a virtual server which is connected to and from the server [9] as shown in figure 3.

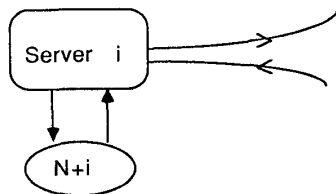


Figure 3. A Virtual Server

For server  $i$ , construct a new virtual server labeled  $(N+i)$ . A virtual job (designated by class  $i$ ) is introduced to circulate only between server  $i$  and server  $(N+i)$ . This virtual job is known as the local job of server  $i$ . When the virtual job is in server  $(N+i)$ , the original jobs are receiving normal service in server  $i$ . As the virtual job jumps to server  $i$ , the server  $i$  breaks down and the original jobs are preempted from the service by the virtual job, which models the repair phase of the server. The virtual job in server  $i$  and server  $(N+i)$  has class  $i$  where the original jobs have class  $(N+1)$ . Then the model becomes a multi-class-job queueing-network with preemptive-resume priority servers.

#### Assumptions [9]

1. The service times of the jobs are exponentially distributed.
2. The mean service times are:

$$1/\mu_{ii} = 1/\Psi_i \tag{3-1}$$

$$1/\mu_{N+i,i} = 1/\Theta_i \tag{3-2}$$

$$1/\mu_{i,N+1} = 1/\mu_i \tag{3-3}$$

for all  $i$ .

This transformation simulates the condition of exponential failures and repairs from the perspective of the original system jobs. The transformed model with all reliable servers is then solved for the performance measures of the original system jobs (of class  $(N+1)$ ) by applying approximation techniques for priority queueing models [3,5].

Our solution of the model is based on the algorithm in [5] which is a mean-value analysis algorithm for multi-class preemptive-resume priority queueing network models. To our knowledge, the algorithm in [5] provides the most accurate results within the existing techniques for the preemptive-resume priority queueing networks.

The mean response-time is:

$$\bar{t}_{i,r} = \left[ \sum_{l=1}^{r-1} \frac{\bar{k}_{i,l}}{\Phi_{i,l} \mu_{i,l}} \right] + \frac{K_r - 1}{K_r} \frac{\bar{k}_{i,r}}{\Phi_{i,r} \mu_{i,r}} + \frac{1}{\Phi_{i,r} \mu_{i,r}}, \tag{3-4}$$

for all  $i$  and  $r$ .

$$\Phi_{i,k} \equiv 1 - \sum_{s=1}^{k-1} \rho_{i,s}$$

The r.h.s. of (3-4) has the following interpretation. Part 1 is the service time for the jobs with higher priority than the tagged arriving job with class  $r$ . Part 2 is the service time for the earlier arrived jobs of the same class, adjusted by the Schweitzer/Bard [2,11] factor  $(K_r - 1)/K_r$ . Part 3 is the service time for the tagged job itself. The  $\Phi_{i,k}$  increases the effective service rate and thus slows down the service times of the servers because of the preemption.

The throughput of the network and the mean number at server  $i$ , all for job class  $r$  is obtained by Little's law:

$$\lambda_r = K_r \sum_{i=1}^N e_{i,r} \bar{t}_{i,r} \tag{3-5}$$

$$\bar{k}_{i,r} = \lambda_r \cdot e_{i,r} \cdot \bar{t}_{i,r} \tag{3-6}$$

$$e_{i,r} = \sum_{j=1}^N \sum_{s=1}^R e_{j,s} \cdot p_{j,s;i,r}$$

for all  $i$  and  $r$ .

Since the jobs with class  $r=i$  is only local to server  $i$  and  $(N+i)$ , and the original job with class  $(N+1)$  is circulating only among servers  $(1, 2, \dots, N)$ , the routing matrix and corresponding job visiting ratio to each server are:

$$P_{i,rj,s} = \begin{cases} 1, & \text{if } i=r=s \text{ \& } j=i+N \\ 1, & \text{if } i=j+N \text{ \& } r=j=s \\ p_{i,j}, & \text{if } r=N+1 \text{ \& } s=N+1 \\ 0, & \text{otherwise} \end{cases} \quad (3-7)$$

$$e_{i,r} = \begin{cases} 1, & \text{if } i=r \text{ or } i=r+N \\ e_i, & \text{if } i=N+1 \\ 0, & \text{otherwise} \end{cases} \quad (3-8)$$

$$e_i = \sum_{j=1}^N e_j p_{ji} \quad (3-9)$$

for all  $i$  and  $r$ .

Eq (3-8) shows that certain jobs never go to certain servers. Thus, only the following response times are interesting:  $\bar{t}_{i,i}$ ,  $\bar{t}_{N+i,i}$ ,  $\bar{t}_{i,N+1}$  (for all  $i$ ). We derive formulas of these response times by analyzing the behavior of the virtual jobs. There is at most one virtual job in server  $i$  and  $(N+i)$ , and that job always has higher priority to preempt other jobs from service, ie, the virtual job gets service immediately after it enters a server. Thus —

$$\bar{t}_{N+i,i} = 1/\Theta_i, \quad (3-10)$$

$$\bar{t}_{i,i} = 1/\Psi_i, \quad (3-11)$$

for all  $i$ .

The response time for the original job can be obtained by simplifying (4) and by considering that if  $e_{i,r}=0$  then  $\bar{k}_{i,r}=0$  for all  $i$  and  $r$ , except  $i=r$ :

$$\bar{t}_{i,N+1} = \frac{\bar{k}_{i,i}}{\bar{\rho}_{i,i} \mu_{i,i}} + \frac{K_{N+1}-1}{K_{N+1}} \frac{\bar{k}_{i,N+1}}{\bar{\rho}_{i,i} \mu_{i,N+1}} + \frac{1}{\bar{\rho}_{i,i} \mu_{i,N+1}} \quad (3-12)$$

$$= \frac{\lambda_{i,i} \bar{t}_{i,i}}{\bar{\rho}_{i,i} \mu_{i,i}} + \frac{K_{N+1}-1}{K_{N+1}} \frac{\bar{k}_{i,N+1}}{\bar{\rho}_{i,i} \mu_{i,N+1}} + \frac{1}{\bar{\rho}_{i,i} \mu_{i,N+1}}$$

$$\bar{\rho}_{i,i} \equiv 1 - \rho_{i,i} \quad (3-13)$$

To further simplify (3-13), we observe the behavior of the arrival rate and usage of the virtual job in server  $i$ . The virtual job circulates between server  $i$  and  $(N+i)$  with the following idle and busy period in server  $i$ :

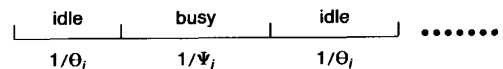


Figure 4. Chart for Usage

Thus —

$$\rho_{i,i} = \frac{1/\Psi_i}{(1/\Theta_i) + (1/\Psi_i)} \quad (3-14)$$

$$\lambda_{i,i} = \rho_{i,i} \mu_{i,i} = \frac{1}{(1/\Theta_i) + (1/\Psi_i)} \quad (3-15)$$

By substituting (3-11), (3-15) into (3-13) we obtain

$$\bar{t}_{i,N+1} = A_i + \frac{K_{N+1}-1}{K_{N+1}} \frac{\bar{k}_{i,N+1}}{B_i \mu_{i,N+1}} + \frac{1}{B_i \mu_{i,N+1}} \quad (3-16)$$

Only jobs of class  $(N+1)$  are in (3-16), thus we can omit the class subscript:

$$\bar{t}_i = A_i + \kappa \frac{\bar{k}_i}{B_i \mu_i} + \frac{1}{B_i \mu_i} \quad (3-17)$$

$$A_i \equiv \Theta_i / \Psi_i^2 \quad (3-18a)$$

$$B_i \equiv \frac{1/\Theta_i}{(1/\Theta_i) + (1/\Psi_i)} \quad (3-18b)$$

Finally we obtain the mean value analysis scheme shown in table 1 for unreliable queuing network models:

TABLE 1  
Mean-Value Analysis for Unreliable Models

Mean Response Time	$\bar{t}_i = \frac{1}{\mu_i B_i} [1 + \kappa \bar{k}_i] + A_i$
Throughput	$\lambda = K / \sum_{i=1}^N e_i \bar{t}_i$
Mean Number of Jobs	$\bar{k}_i = \bar{t}_i e_i \lambda$
Initial Value	$\bar{k}_i = K/N$

#### 4. PERFORMANCE OPTIMIZATION

By substituting  $\bar{k}_i = \lambda e_i \bar{t}_i$  we solve (3-17) for  $\bar{t}_i$ :

$$\bar{t}_i = \frac{A_i B_i \mu_i + 1}{B_i \mu_i - \kappa \lambda e_i} \quad (4-1)$$

Since —

$$\bar{k}_i = \lambda e_i \bar{t}_i = \frac{\lambda e_i (A_i B_i \mu_i + 1)}{B_i \mu_i - \kappa \lambda e_i} \quad (4-2)$$

and in a closed queueing network —

$$\sum_{i=1}^N \bar{k}_i = K \quad (4-3)$$

it follows —

$$\sum_{i=1}^N \frac{\lambda e_i (A_i B_i \mu_i + 1)}{B_i \mu_i - \kappa \lambda e_i} = K \quad (4-4)$$

Eq (4-4) is the constraint used in application of Lagrange multiplier technique for optimization of performance measures. If the throughput is the optimization subject, then (4-4) would be the appropriate constraint. However, if the throughput is given, (4-4) can be simplified to (4-5):

$$\sum_{i=1}^N \frac{D_i}{B_i \mu_i - \kappa \lambda e_i} = \hat{K} \quad (4-5)$$

$$\hat{K} = K - \sum_{i=1}^N \lambda e_i A_i$$

$$D_i \equiv \lambda e_i (A_i \lambda e_i \kappa + 1)$$

In the optimization procedure of performance measures such as response time, throughput, and usage, the objective is to determine the optimal service rates as decision variables subject to certain constraints. The cost constraint is the most common and has the following linear form:

$$\sum_{i=1}^N c_i \mu_i = C. \quad (4-6)$$

The following sub-sections summarize the optimization problem solutions:

1. *Cost Optimization for Fixed Throughput:* The total throughput of the queueing network is assumed to be known. The optimal service rates  $\mu_i^*$  must be determined such that the given total throughput is reached at minimum cost.

2. *Response Time Minimization for Each Server:* The response time of each server for a given queueing network is minimized where a given fixed throughput value is controlled. The service rates are also under cost constraints.

3. *Throughput Maximization under Fixed Costs:* The total cost for the queueing network model is known, viz, the budget available for purchasing a given number of servers with specific service rates. The total throughput is maximized by determining optimal service rates under cost constraints.

#### 4.1 Cost Minimization

This section uses the queueing network model from section 3. The objective is to choose  $N$  servers having service rates  $\mu_i^*$  for minimum costs while the given throughput of the system is maintained. We formulate the optimization problem as:

Minimize the total cost  $C(\mu)$  for a given throughput.

The optimal service rates and total optimal cost are:

$$\mu_i^* = \sqrt{\frac{D_i}{B_i c_i}} \frac{\sum_{j=1}^N \sqrt{\frac{c_j D_j}{B_j}}}{\hat{K}} + \kappa \frac{\lambda e_i}{B_i} \quad (4-7)$$

$$C^*(\mu) = \sum_{i=1}^N \mu_i^* c_i \quad (4-8)$$

The derivation is in the appendix.

*Example 1: Queueing Network Model of an Unreliable System.*

There are  $N=5$  servers and  $K=10$  jobs. Additionally, 5 virtual servers and 5 virtual jobs are introduced to model the unreliability aspects. The mean time-to-failure of each server is  $\Theta_i = 8$  units. The mean repair time is  $\Psi_i = 10$  units. Figure 5 is the queueing network model with virtual servers; it captures the unreliability aspects of the system.

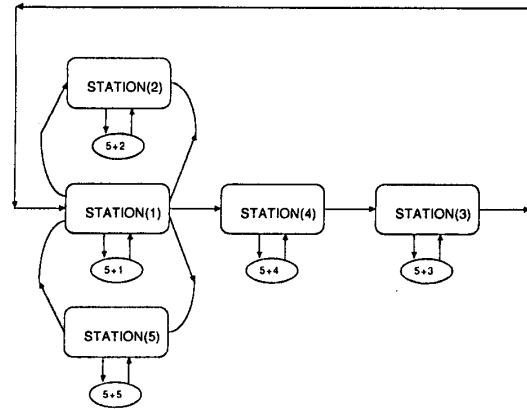


Figure 5. The Queueing Network Model of an Unreliable System

The transition probabilities are:

$$p_{1,2}=0.2; p_{1,4}=0.5; p_{1,5}=0.3; p_{4,3}=p_{2,1}=p_{3,1}=p_{5,1}=1.$$

The costs for each server  $c_i$  are:

$$c_1=10; c_2=20; c_3=15; c_4=8$$

The required throughput is  $\lambda=0.5$ .

The optimal service rates are:

$$\mu^* = [5.997, 1.207, 2.780, 3.657, 2.212]$$

The minimum cost of the system which achieves the throughput requirement is:

$$C^* = 177.204.$$

Table 2 and figure 6 show the dependence of the minimum cost on throughput.

TABLE 2  
Relation between Throughput and Minimum Cost

Throughput	0.1	0.2	0.3	0.4	0.5	0.6	0.62	0.625
Cost	8.43	20.9	40.9	78.8	177	1062	5487	$\infty$

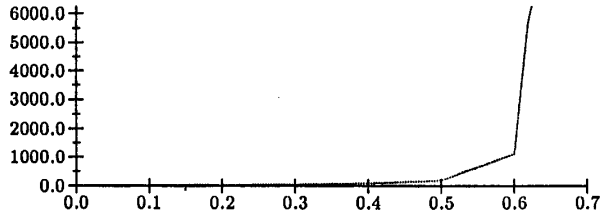


Figure 6. Cost Dependence on Throughput

#### 4.2 Response-Time Minimization

The optimization problem is to minimize the response time  $\bar{t}_i$  for a particular server  $i$  subject to fixed given cost and fixed given throughput constraints. The service rates have a linear cost constraint and the total cost  $C$  is known. The mean response time is given by (4-1).

If we minimize (4-1), as an objective function with conditions, (4-5) and (4-6), the following optimal service rates are obtained:

$$\mu_i^* = \frac{-\beta + \sqrt{\beta^2 - 4\alpha\gamma}}{2\alpha} \quad (4-9a)$$

$$\alpha \equiv \hat{K} B_i c_i$$

$$\beta \equiv \left( \sum_{j \neq i} \kappa \frac{\lambda e_j c_j}{B_j} - C \right) \hat{K} B_i - c_i \left( \kappa \lambda e_i \hat{K} - D_i \right)$$

$$+ \left( \sum_{j \neq i} \sqrt{\frac{c_j D_j}{B_j}} \right) \left( \sum_{j \neq i} \sqrt{\frac{c_j}{B_j}} \right) B_i$$

$$\gamma \equiv \left( e_i - \sum_{j \neq i} \kappa \frac{\lambda e_j e_j}{B_j} \right) \cdot \left( \kappa \lambda e_i \hat{K} - D_i \right) - \left( \sum_{j \neq i} \sqrt{\frac{e_j D_j}{B_j}} \right) \left( \sum_{j \neq i} \sqrt{\frac{e_j}{B_j}} \right) \kappa \lambda B_i$$

The optimal service rates  $\mu_j^*$  for all  $j, j \neq i$ , are:

$$\mu_j^* = \frac{\left( \sqrt{\frac{D_j}{c_j B_j}} \right) \left( \sum_{j \neq i} \sqrt{\frac{c_j D_j}{B_j}} \right)}{\hat{K} - \frac{D_i}{B_i \mu_i^* - \kappa \lambda e_i}} + \kappa \frac{\lambda e_j}{B_j} \quad (4-9b)$$

The minimum response time for server  $i$  is:

$$\bar{t}_i^* = \frac{A_i B_i \mu_i^* + 1}{B_i \mu_i^* - \kappa \lambda e_i} \quad \text{for } i=1, \dots, N \quad (4-10)$$

The derivation is in the appendix.

*Example 2:* Same Model as in Figure 1.

The throughput is  $\lambda=0.5$  and the available budget is  $C=200$ . We obtain the minimum response time for server 5 from (4-10):

$$\bar{t}_5^* = 6.858$$

To achieve this minimum response time for server 5, the service rates for each server are computed from (4-9):

$$\mu^* = [5.184, 1.044, 2.408, 3.147, 7.562].$$

#### 4.3 Throughput Maximization

In the planning phase of distributed systems the budget is limited. The objective is to optimize the performance measures within the budget. Given fixed cost  $C$ , we find an optimum system that achieve the highest throughput. Rewrite (4-4) into:

$$\sum_{i=1}^N \frac{e_i (A_i B_i \mu_i + 1)}{B_i \mu_i - \kappa \lambda e_i} = \frac{K}{\lambda} \quad (4-11)$$

Differentiate both sides of (4-11) by  $\mu_i$ :

$$\frac{\partial \lambda}{\partial \mu_i} = \frac{B_i (e_i A_i B_i \mu_i + e_i) - e_i A_i B_i (B_i \mu_i - \kappa \lambda e_i)}{(B_i \mu_i - \kappa \lambda e_i)^2} \cdot \Gamma \quad (4-12)$$

$$\Gamma \equiv \left[ \sum_{j=1}^N \frac{\kappa e_j^2 (A_j B_j \mu_j + 1)}{(B_j \mu_j - \kappa \lambda e_j)^2} - \frac{K}{\lambda^2} \right]^{-1}$$

The Lagrange function has the following form:

$$L(\mu, y) = \lambda(\mu) + y \left( \sum_{i=1}^N c_i \mu_i - C \right) \quad (4-13)$$

Differentiate  $L$  by  $\mu_i$  and  $y$ , and set derivatives = 0:

$$\frac{\partial L}{\partial \mu_i} = \frac{\partial \lambda}{\partial \mu_i} + c_i y = 0 \quad (4-14a)$$

$$\frac{\partial L}{\partial y} = \sum_{i=1}^N c_i \mu_i - C = 0 \quad (4-14b)$$

Let  $\phi \equiv -y/\Gamma$ ; from (4-12) and (4-13) it follows that:

$$\frac{e_i A_i B_i \mu_i + e_i}{B_i \mu_i - \kappa \lambda e_i} - e_i A_i = \phi \left[ c_i \mu_i - \lambda \kappa \frac{e_i c_i}{B_i} \right]. \quad (4-15)$$

Sum (4-15) over all  $i$  and solve for  $\phi$ .

From (4-15) and the equation for  $\phi -$

$$\mu_i^* = \kappa \frac{\lambda e_i}{B_i} +$$

$$\frac{1}{B_i} \left[ \frac{\left( C - \lambda \kappa \sum_{j=1}^N \frac{e_j c_j}{B_j} \right) (\kappa \lambda e_i^2 A_i B_i + e_i B_i)}{c_i \left( \frac{K}{\lambda} - \sum_{j=1}^N e_j A_j \right)} \right]^{1/2} \quad (4-16)$$

By substituting (4-16) into (4-14b) we derive:  $C =$

$$\sum_{i=1}^N \frac{c_i}{B_i} \cdot \left[ \kappa \lambda e_i + \left[ \frac{\left( C - \lambda \kappa \sum_{j=1}^N \frac{e_j c_j}{B_j} \right) (\kappa \lambda e_i^2 A_i B_i + e_i B_i)}{c_i \left( \frac{K}{\lambda} - \sum_{j=1}^N e_j A_j \right)} \right]^{1/2} \right] \quad (4-17)$$

By solving (4-17) we obtain the optimal throughput value that can be achieved; unfortunately, it has no closed form solution. However, it can be solved by the binary search and insert value method. After solving (4-17) the service rates of each server can be derived from (4-16) by substituting the  $\lambda$ .

*Example 3:* Same Model as in Figure 1.

We need the optimal throughput, which is the reverse of example 1. All parameters are the same as in example 1 except that the throughput is unknown. The total cost ( $C = 177.046$ ) is given as an input. We substitute the given parameter values

into (4-17) and solve for the maximum throughput,  $\lambda^* = 0.5$  by the binary search and insert value method.

Substitute this throughput value into (4-16) and obtain the optimal service rates for each server:

$$\mu^* = [5.997, 1.207, 2.780, 3.657, 2.212]$$

These results match the results in example 1, as they should.

## APPENDIX

### A-1. Derivation of Cost Minimization: (4-7) & (4-8)

We solve this optimization problem using the Lagrange multiplier technique. First we derive the Lagrange function:

$$L(\mu, y) = \sum_{i=1}^N c_i \mu_i + y \left( \sum_{i=1}^N \frac{D_i}{B_i \mu_i - \kappa \lambda e_i} - \hat{K} \right) \quad (A-1)$$

By differentiating  $L(\mu, y)$  by  $\mu_i$  for  $i=1, \dots, N$  and  $y$  we obtain the necessary and sufficient conditions for optimal service rates  $\mu_i^*$ :

$$\frac{\partial L}{\partial \mu_i} = c_i - y \frac{B_i D_i}{(B_i \mu_i - \kappa \lambda e_i)^2} = 0 \quad (A-2)$$

$$\frac{\partial L}{\partial y} = \sum_{i=1}^N \frac{D_i}{B_i \mu_i - \kappa \lambda e_i} - \hat{K} = 0. \quad (A-3)$$

From (A-2) we get:

$$\frac{D_i}{B_i \mu_i - \kappa \lambda e_i} = \sqrt{\frac{c_i D_i}{B_i y}}. \quad (A-4)$$

By substituting (A-4) into (A-3) we derive:

$$\sqrt{\frac{1}{y}} \sum_{i=1}^N \sqrt{\frac{c_i D_i}{B_i}} = \hat{K}. \quad (A-5)$$

Eq (A-5) provides:

$$y = \left( \sum_{i=1}^N \sqrt{\frac{c_i D_i}{B_i}} \right)^2 / \hat{K}^2. \quad (A-6)$$

From (A-2) and (A-6) we obtain the optimal service rates that minimizes the system cost while achieving the required throughput.

### A-2. Derivation of Response-Time Minimization: (4-9) & (4-10)

The Lagrange function is derived as follows:

$$L(\mu, y, z) = \frac{A_i B_i \mu_i + 1}{B_i \mu_i - \kappa \lambda e_i} + z \left( \sum_{i=1}^N c_i \mu_i - C \right) + y \left( \sum_{i=1}^N \frac{D_i}{B_i \mu_i - \kappa \lambda e_i} - \hat{K} \right) \quad (\text{A-7})$$

Differentiate by  $\mu_j$ ,  $\mu_i$  ( $i \neq j$ ),  $y$ ,  $z$ ; set the results = 0.

$$\frac{\partial L}{\partial \mu_j} = z c_j - y \frac{B_j D_j}{(B_j \mu_j - \kappa \lambda e_j)^2} = 0, \quad (\text{A-8})$$

for  $j=1, \dots, N$  &  $j \neq i$ .

$$\frac{\partial L}{\partial y} = \sum_{i=1}^N \frac{D_i}{B_i \mu_i - \kappa \lambda e_i} - \hat{K} = 0 \quad (\text{A-9})$$

$$\frac{\partial L}{\partial \mu_i} = z c_i - \frac{\kappa \lambda A_i B_i e_i + y B_i D_i}{(B_i \mu_i - \kappa \lambda e_i)^2} = 0, \quad (\text{A-10})$$

for  $j=1, \dots, N$  &  $j \neq i$

$$\frac{\partial L}{\partial z} = \sum_{i=1}^N c_i \mu_i - C = 0 \quad (\text{A-11})$$

From (A-8) we get:

$$\frac{D_j}{B_j \mu_j - \kappa \lambda e_j} = \sqrt{\frac{z c_j D_j}{B_j}}. \quad (\text{A-12})$$

By substituting (A-12) into (A-9) we derive:

$$\sqrt{\frac{z}{y}} = \frac{\hat{K} - \frac{D_i}{B_i \mu_i - \kappa \lambda e_i}}{\sum_{j \neq i}^N \sqrt{\frac{c_j D_j}{B_j}}}. \quad (\text{A-13})$$

From (A-13) we obtain (4-9b). By substituting (A-13) and (4-9b) into (A-11) we obtain:  $\alpha \mu_i^2 + \beta \mu_i + \gamma = 0$ . Finally by substituting (A-13) and (4-1) into (4-9b) we obtain the optimal service rates  $\mu_i^*$ .

## REFERENCES

- [1] I. F. Akyildiz, G. Bolch, "Throughput maximization and response time minimization in queueing network models of computer systems", *Proc. Int'l Seminar on Distributed and Parallel Systems*, Hasegawa, Takagi & Takahashi (Eds), North-Holland, 1988 Dec, pp 241-259.
- [2] Y. Bard, "Some extensions to multiclass queueing network analysis", *Proc. Performance '79 Conf.*, North-Holland, 1979 Feb.

- [3] R. Bryant, T. Krzesinski, S. Lakshmi, K. M. Chandy, "The MVA priority approximation", *ACM Trans. Computer Systems*, vol 2, 1984 Nov, pp 335-359.
- [4] K. M. Chandy, J. Hogarth, C. H. Sauer, "Selecting capacities in computer communication systems", *IEEE Trans. Software Engineering*, vol 4, 1977 Jul, pp 290-295.
- [5] J. van Doremalen, J. Wessels, R. Wijbrands, "Approximate analysis of priority queueing networks", *Proc. Conf. Teletraffic Analysis and Computer Performance Evaluation*, Boxma, Cohen, Tijms (Eds), North-Holland, 1986, pp 117-131.
- [6] J. R. Kenevan, A. von Mayrhauser, "Convexity and concavity properties of analytic queueing models for computer systems", *Proc. Performance '84 Conf.*, E. Gelenbe (Ed), North-Holland, 1984, pp 361-375.
- [7] H. Kobayashi, M. Gerla, "Optimal routing in closed queueing networks", *ACM Trans. Computer Systems*, vol 1, 1983 Nov, pp 294-310.
- [8] A. von Mayrhauser, K. Trivedi, "Computer configuration design to minimize response time", *Computer Performance*, Butterworth & Co., vol 3, 1982 Mar, pp 32-39.
- [9] C. S. Ramanjaneyulu, V. V. S. Sarma, "Modeling server-unreliability in closed queueing-networks", *IEEE Trans. Reliability*, vol 38, 1989 Apr, pp 90-95.
- [10] M. Reiser, S. Lavenberg, "Mean value analysis of closed multichain queueing networks", *J. ACM*, 1980 Apr, pp 313-322.
- [11] P. J. Schweitzer, "Approximate analysis of multiclass closed networks of queues", *Proc. Int'l Conf. Stochastic Control and Optimization*, Amsterdam, Netherlands, 1979 May.
- [12] K. Stecke, "On the nonconcavity of throughput in certain closed queueing networks", *Performance Evaluation Journal*, vol 6, 1986, pp 293-305.
- [13] K. S. Trivedi, *Probability and Statistics with Reliability, Queueing, and Computer Science Applications*, Prentice Hall, 1982.
- [14] K. S. Trivedi, R. A. Wagner, "A decision model for closed queueing networks", *IEEE Trans. Software Engineering*, vol SE-5, 1979 Jul.
- [15] B. Vinod, T. Altioik, "Approximating unreliable queueing networks under the assumption of exponentiality", *Operational Research Society*, vol 37, 1986, pp 309-316.

## AUTHORS

Dr. I. F. Akyildiz; School of Information and Computer Science; Georgia Institute of Technology; Atlanta, Georgia 30332-0280 USA.

I. F. Akyildiz (M '86, SM '89) was born in Istanbul, Turkey on 1954. He received the BS, MS, and Doctor of Engineering degrees in Computer Science from the University of Erlangen-Nuernberg, Fed. Rep. GERMANY, in 1978, 1981, and 1984. He is an Associate Professor in the School of Information and Computer Science at Georgia Institute of Technology. He is a co-author of a textbook entitled *Analysis of Computer Systems* published by Teubner Verlag in the Fall of 1982, and is a guest editor of the special issue on Queueing Networks with Finite Capacities in *Performance Evaluation Journal*. He is an associate editor for *Computer Networks and ISDN Journal*. His research interests are performance evaluation, computer networks, distributed systems, and computer security. Dr. Akyildiz is a member of ACM (Sigops and Sigmetrics), and Gesellschaft für Informatik.

Wei Liu; School of Information and Computer Science; Georgia Institute of Technology; Atlanta, Georgia 30332-0280 USA.

Wei Liu (S '89) was born in Guangdong, Peop. Rep. CHINA in 1964. He received the BS degree in Computer Science from Beijing University in 1985 and the MS degree in Information and Computer Science from Georgia Institute of Technology in 1988. He is a PhD student in the School of Information and Computer Science at Georgia Tech. His research interests are computer networks, system performance, queueing theory and distributed database systems.

Manuscript TR88-127 received 1988 June 21; revised 1989 August 10; revised 1989 November 15.

IEEE Log Number 34126

◀TR▶