

## ABSTRACT

XIONG, KAIQI. Resource Optimization and Security in Distributed Computing. (Under the direction of Professor Harry Perros.)

With the number of e-Business applications dramatically increasing, service level agreements (SLA) will play an important part in distributed service computing. An SLA is a combination of several quality of service (QoS) metrics, such as security, performance, and availability, agreed between a customer and a service provider. Due to the complexity of these metrics, most existing research typically addresses only one of these QoS metrics. In the case of the response time as a performance metric, the average time to process and complete a job is typically used in the literature. However, this may not be of real interest to a customer. A statistically bounded metric, that is, a percentile response time, is more realistic than the average response time. Moreover, in enterprise service computing, customer requests are typically distinguished by different request characteristics and service requirements.

This dissertation includes a study of trustworthiness, percentile response time, service availability, and authentication among service stations or sites that may be owned by different service providers. The first part of this dissertation contains an analysis of percentile response time, which is one of the most important SLA metrics. Effective and accurate numerical solutions for the calculation of the percentile response time in single-class and multi-class queueing networks are obtained. Then, the numerical solution is incorporated in a resource allocation problem. Specifically, we present an approach for the resource optimization that minimizes the total cost of computer resources required while preserving a given percentile of the response time.

In the second part of this dissertation, we extend the approach to consider trustworthiness, service availability, and the percentile of response time in Web services. We clearly define these QoS metrics and provide their quantitative analysis. Then, we take into account these QoS metrics in a trust-based resource allocation problem in which a set of computer resources is used by a service provider to host a typical Web services application for single-class customer services and multiple-class customer services respectively. We formulate the trust-based resource allocation problem as an optimization problem under SLA constraints in which we calculate the number of servers in each service site that minimize a cost function that reflects operational costs for single-class customer services and multiple-class customer services respectively. We solve this problem using an efficient numerical procedure. Experimental results show the applicability of the procedure and validate its accuracy.

Finally, in the third part of this dissertation we first present a thorough performance evaluation of two notable public key cryptography-based authentication techniques, Public-Key Cross Realm Authentication in Kerberos (PKCROSS) and Public Key Utilizing Tickets for Application Servers (PKTAPP, a.k.a. KX.509/KCA), in terms of computational and communication times. We then demonstrate their performance difference using queueing networks. PKTAPP was proposed to address the scalability issue of PKCROSS. However, our in-depth analysis of these two techniques shows that PKTAPP does not perform better than PKCROSS in a large scale system. Thus, we propose a new public key cryptography-based group authentication technique. Our performance analysis demonstrates that the new technique can scale better than PKCORSS and PKTAPP.

©Kaiqi Xiong

2007

Resource Optimization and Security in Distributed Computing

by

Kaiqi Xiong

A dissertation submitted to the Graduate Faculty of  
North Carolina State University  
in partial fulfillment of the  
requirements for the Degree of  
Doctor of Philosophy

Computer Science

Raleigh, NC

2007

Approved By:

---

Dr. David Thuente

---

Dr. Matthias Stallmann

---

Dr. Harry Perros  
Chair of Advisory Committee

---

Dr. Mihail Devetsikiotis

## **Dedication**

To My Family, Parents, Brother and Sister

## **Biography**

Kaiqi Xiong has studied in the Department of Electrical and Computer Engineering at the North Carolina State University where he received the Master of Science in Computer Engineering in December 2000. He has been working in industry for several years before he returned to school for his Ph.D. study in the Department of Computer Science at the same university.

## **Acknowledgements**

I would like to express my sincere appreciation to Dr. Harry Perros for being an adviser and friend. I am very thankful to him for his understanding, guidance, and encouragement during the process of my studies and thesis writing. He was very supportive and always available whenever I wanted to meet him. I very much enjoyed working with him.

I would like to thank Dr. David Thuente, Dr. Matthias Stallmann and Dr. Mihail Devetsikiotis for being in my committee and for giving me valuable comments and suggestions. I am specially thankful to Dr. David Thuente for his great help and support as Director of Graduate Programs. I also want to thank Dr. Matthias Stallmann and Dr. Mihail Devetsikiotis in devoting their time to helping me when needed.

Last, but definitely not the least, I specially want to thank my family for understanding and support.

# Contents

<b>List of Figures</b>	<b>vii</b>
<b>List of Tables</b>	<b>viii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Motivation . . . . .	1
1.2 The Problems . . . . .	3
1.3 Summary of Contributions . . . . .	7
1.4 The Organization of this Dissertation . . . . .	8
<b>2 Related Work</b>	<b>9</b>
2.1 Service Availability . . . . .	9
2.2 Trustworthiness . . . . .	11
2.3 Performance . . . . .	12
2.4 The Resource Optimization Problem Subject to an SLA . . . . .	14
2.5 Public Key Cryptography-based Authentication . . . . .	15
<b>3 Single Class Customers</b>	<b>18</b>
3.1 The Percentile of Response Time . . . . .	18
3.2 A Resource Optimization Problem for Service Models with Single Class Customers	20
3.3 Approaches for the Resource Optimization . . . . .	21
3.4 Numerical Validations . . . . .	26
3.5 The Balanced Condition . . . . .	29
3.6 An Application: Web Services Performance Modeling and Analysis . . . . .	34
3.6.1 A Web Services Performance Model . . . . .	35
3.6.2 Analysis of the Web Services Performance Model . . . . .	36
3.6.3 A Numerical Validation . . . . .	40
3.7 Concluding Remarks . . . . .	42
<b>4 Multiple Class Customers</b>	<b>44</b>
4.1 The SLA Performance Metric in the Case of Multiple Class Customers . . . . .	45
4.2 The Resource Optimization Problem for Multiple Customer Services . . . . .	46

4.3	Approaches for Resource Optimization . . . . .	46
4.3.1	The LSTs of Response Time Distributions for Two Priority Customers . . . . .	46
4.3.2	Algorithms for the Resource Optimization Problem . . . . .	50
4.4	Numerical Validations . . . . .	56
4.5	Concluding Remarks . . . . .	61
<b>5</b>	<b>A Trustworthy Service Model</b>	<b>63</b>
5.1	The Trust-based Resource Optimization Problem . . . . .	64
5.2	A Framework for Solving the Trust-based Resource Provisioning Problem . . . . .	66
5.3	The Calculation of SLA Metrics . . . . .	70
5.3.1	The Trustworthiness of Resource Sites . . . . .	70
5.3.2	The Percentile Response Time . . . . .	72
5.3.3	The Service Availability . . . . .	74
5.4	An Approach for Solving the Trust-based Resource Provisioning Problem . . . . .	75
5.4.1	Single Class Customers . . . . .	75
5.4.2	Multiple Priority Customers . . . . .	81
5.5	Numerical Examples . . . . .	87
5.5.1	Single Class Customers . . . . .	88
5.5.2	Multiple Priority Customers . . . . .	90
5.6	Concluding Remarks . . . . .	94
<b>6</b>	<b>Performance Analysis of Public Key Cryptograph-based Group Authentication</b>	<b>96</b>
6.1	Public Key Cryptography-based Authentication . . . . .	97
6.2	PKCROSS and PKTAPP . . . . .	98
6.2.1	Protocol Analysis . . . . .	99
6.2.2	The Calculation of the Response Time via Queueing Networks . . . . .	102
6.3	A New Group Authentication Technique Using Public-Key Cryptography . . . . .	107
6.3.1	A Single Remote Realm . . . . .	108
6.3.2	Multiple Remote Realms . . . . .	111
6.4	Performance Evaluation of the New Proposed Technique . . . . .	113
6.4.1	The Operations of Encryption and Decryption . . . . .	113
6.4.2	The Calculation of the Response Time via a Queueing Network . . . . .	116
6.4.3	Discussions . . . . .	119
6.5	Concluding Remarks . . . . .	120
<b>7</b>	<b>Summary and Future Work</b>	<b>121</b>
7.1	Research Summary of the Dissertation . . . . .	121
7.2	Future Research Directions . . . . .	123
	<b>Bibliography</b>	<b>126</b>

# List of Figures

1.1	The Execution of Service Requests . . . . .	4
2.1	A Three-node Tandem Network with Overtaking . . . . .	13
3.1	Percentile Response Time vs. Service Rate . . . . .	20
3.2	A Tandem-station Service Model for Single-class Customer Services . . . . .	21
3.3	A Service Model with Feedback for Single-class Customer Services . . . . .	24
3.4	A Two Station Tandem Model . . . . .	30
3.5	A Scenario in Web Services Applications . . . . .	35
3.6	A Web Server Performance Model . . . . .	37
4.1	A Tandem-station Service Model for Multiple-class Customer Services . . . . .	50
4.2	A Service Model with Feedback for Multiple-class Customer Services . . . . .	54
5.1	An SLA-based Web Services Model . . . . .	65
5.2	A Framework for Solving the Trust-based Resource Provisioning Problem . . . . .	67
5.3	A Web Services Computing System for Single-class Customer Services . . . . .	75
5.4	A Web Services Computing System for Multiple-class Customer Services . . . . .	81
6.1	The PKCROSS Message Flow . . . . .	99
6.2	The PKTAPP Message Flow . . . . .	100
6.3	A Queueing Network with $m$ Remote Realms for PKCROSS . . . . .	103
6.4	A Queueing Network with $n$ Application Servers for PKTAPP . . . . .	104
6.5	A Test Scenario with a Remote Realm . . . . .	105
6.6	Response Time vs. Authentication Request Rate . . . . .	106
6.7	Crossover Numbers vs. The Number of Remote Realms . . . . .	107
6.8	The Message Flow of the New Technique in a Single Remote Realm . . . . .	109
6.9	The Message Flow of the New Technique in Multiple Remote Realms . . . . .	112
6.10	A Queueing Network with $n$ Remote Realms for the New Technique . . . . .	117
6.11	Response Times vs. Authentication Request Rates When $m = 1$ . . . . .	117
6.12	Response Times vs. Authentication Request Rates When $m = 2$ . . . . .	118
6.13	Response Times vs. Authentication Request Rates When $m = 8$ . . . . .	118
6.14	Crossover Numbers vs. Remote Realms . . . . .	119

## List of Tables

3.1	The Cumulative Distribution Function (CDF) of the Response Time vs. Service Rate	21
3.2	The Service Rates of the Eight Stations in Model 1	27
3.3	The Cumulative Distribution of the Response Time in Model 1	28
3.4	The Optimal Number of Servers in Model 1	28
3.5	The Service Rates of The Eight Stations in Model 2	28
3.6	The Optimal Number of Servers in Model 2	29
3.7	The Cumulative Distribution of the Response Time in Model 2	30
3.8	The Solutions of Problems I and II with Varied $\gamma$	33
3.9	The Solutions of Problems I and II with Varied $\hat{c}_2$	34
3.10	A Comparison of the Number of Servers by Using Problems I and II	34
3.11	The Cumulative Distribution Functions of The Response Time in the Web Services Performance Model	41
4.1	The Cumulative Distribution of the Low-priority Response Time for $\lambda^{(1)} = \lambda^{(2)} = 50$	58
4.2	The Cumulative Distribution of the Low-priority Response Time for $\lambda^{(1)} = 75$ and $\lambda^{(2)} = 25$	59
4.3	The Cumulative Distribution of The High-priority Response Time in Model 1	59
4.4	The Cumulative Distribution of the Low-priority Response Time in Model 1	60
4.5	The Optimal Number of Servers in Model 1	60
4.6	The Optimal Number of Servers in Model 2	60
4.7	The Cumulative Distribution of The High-priority Response Time in Model 2	61
4.8	The Cumulative Distribution of the Low-priority Response Time in Model 2	62
5.1	The Initial Trust Index of the Ten Stations for Single Class Customers	88
5.2	The Service Rates of the Ten Stations for Single Class Customers	88
5.3	The Server Unavailability Rates of the Ten Stations for Single Class Customers	88
5.4	The Trust Indices of the Ten Stations for Single Class Customers at Times $t = t_2, \dots, t_5, t_{20}, t_{21}$	89
5.5	The Cumulative Distribution of the Response Time for Single Class Customers	90
5.6	The Optimal Number of Servers for Single Class Customers	90
5.7	The Initial Trust Index of the Ten Stations for Priority-class Customers	91
5.8	The Service Rates of the Ten Stations for Priority-class Customers	91

5.9	The Server Unavailability Rates of the Ten Stations for Priority-class Customers . .	91
5.10	The Trust Indices of the Ten Stations for Priority-class Customers at Times $t = t_2, \dots, t_5, \dots, t_{20}, t_{21}, \dots, t_{84}, t_{85}$ . . . . .	92
5.11	The Cumulative Distribution of the High-priority Response Time . . . . .	93
5.12	The Cumulative Distribution of the Low-priority Response Time for Priority-class Customers . . . . .	93
5.13	The Optimal Number of Servers for Priority-class Customers . . . . .	94
6.1	The Operations of Encryption and Decryption When $n$ Application Servers are in $m$ Remote Realms . . . . .	101
6.2	The computational times of Encryption and Decryption Operations . . . . .	105
6.3	A Comparison of Analytic and Simulated Response Times . . . . .	106
6.4	The Notation Used in the Case of a Single Remote Realm . . . . .	108
6.5	The Additional Notation Needed in the Case of Multiple Remote Realms . . . . .	112
6.6	The Operations of Encryption and Decryption When $n$ Application Servers are in $m$ Remote Realms . . . . .	113
6.7	The Minimal Number of Application Servers . . . . .	115
6.8	The Difference of Transaction Times Between Our Technique and PKCROSS (i.e., $g_2(n, m) - g_3(n, m)$ ) . . . . .	116

# Chapter 1

## Introduction

This dissertation is concerned with resource optimization problems subject to various constraints including service availability, performance and security, and the problem of public key cryptography-based group authentication. The motivation of our research is discussed in Section 1.1, the formulation of the optimization problems subject to various constraints is given in Section 1.2, and the performance of public key cryptography-based group authentication is presented in the same section as well. Section 1.3 summarizes the contributions of our research, and Section 1.4 gives the organization of this dissertation.

### 1.1 Motivation

The management of the quality of service is fundamental to a distributed enterprise computing system. The increasing pervasiveness of network connectivity and the proliferation of on demand e-business applications and services in public domains, corporate networks, as well as home environments, give rise to the need for the design of appropriate service solutions. When evaluating how a service provider can support the customer's requirements of an enterprise, several questions often arise:

- How can a service provider more effectively manage its service resources to support the customer's requirements?

- How can a service provider achieve higher return on investment through improved utilization of its existing service resources?
- How can a customer receive its service from a reliable service provider?
- How can a customer receive better quality of service at a lower fee?
- How can multiple service providers authenticate each other such that they can work together to provide the quality of service for a customer?
- How does the process of authentication among service providers affect the overall quality of service?

Accurately predicting e-business application performance based on system statistics and a customer's perceived quality allows a service provider not only to assure the quality of service, but also to avoid over provisioning to meet a service level agreement (SLA).

An SLA sets the expectations between a customer and a service provider, and helps define the relationship between these two parties. It is the cornerstone of how the service provider sets and maintains commitments to the customer. Usually, an SLA is a set of quality of service metrics and a price agreed between a customer and a service provider. It plays an important role in an e-business application. The set of the quality of service metrics generally consists of availability, security, and performance that are defined below:

- *Availability* is the percentage of time that a service provider can offer services. Capabilities in the availability component are determined by the resiliency of the environment of service resources. They are related to recovery mechanisms, and data replications that improve service survivability.
- *Security* can be categorized as *identity security* and *behavior security*. Identity security includes the authentication and authorization between a customer and a service provider, data confidentiality and data integrity. Behavior security includes the trustworthiness among multiple resource sites or stations (both terms are used indistinguishably in this dissertation), and the trustworthiness of these resource sites by customers, including the trustworthiness of computing results provided by these sites. Every IT organization is faced with managing the security of many different resources within its enterprise. Behavior security may also include service survivability and vulnerability.

- *Performance* includes QoS metrics related to the underlying network that interconnects the resource sites, such as, throughput, link utilization, packet loss rate, and end-to-end transfer delay. It also includes similar QoS metrics for the resource sites, such as, response time and throughput. The response time is the time it takes for a service request to be satisfied, and the throughput is the service rate that a service provider can offer.

Customers will receive higher levels of quality of services if they are willing to pay more. Therefore, having different SLAs with different associated costs is a common practice approach.

This dissertation will focus on resource optimization whereby a service provider uses the least resources at each service station but it is still able to achieve the pre-defined SLA. These resources may include servers, storage devices, routers, network links, processors and so on (e.g., see [57] and [91]). It also presents a study of performance of public key cryptography-based group authentication.

## 1.2 The Problems

In this dissertation, we consider a collection of computer resources used by a service provider<sup>1</sup> to host enterprise applications for business customers. An enterprise application running in such a computing environment is associated with an SLA (see [54], [19] and [66]). That is, the service provider is required to execute service requests from a customer within negotiated QoS requirements for a given price. Figure 1.1 depicts a scenario for such an environment. A customer represents a business that generates service requests at a given rate to be processed by the service provider's resource stations according to QoS requirements and for a given fee. As shown in Figure 1.1 a service request is transmitted to the service provider over a network provider. After it is processed at the various resource stations of the service provider the final result is sent back to the customer. For presentation purposes, we assume that each resource station has only one type of server associated with cost  $c_j$ . If they have multiple types of servers, we can decompose each resource station into several individual stations so that each one only contains one type of server with the same cost.

Let  $N_j$  be the number of servers at station  $j$  ( $j = 1, 2, \dots, m$ ). Thus, the resource allo-

---

<sup>1</sup>The computer resources may or may not owned by the service provider. The service provider may only act as a service broker.

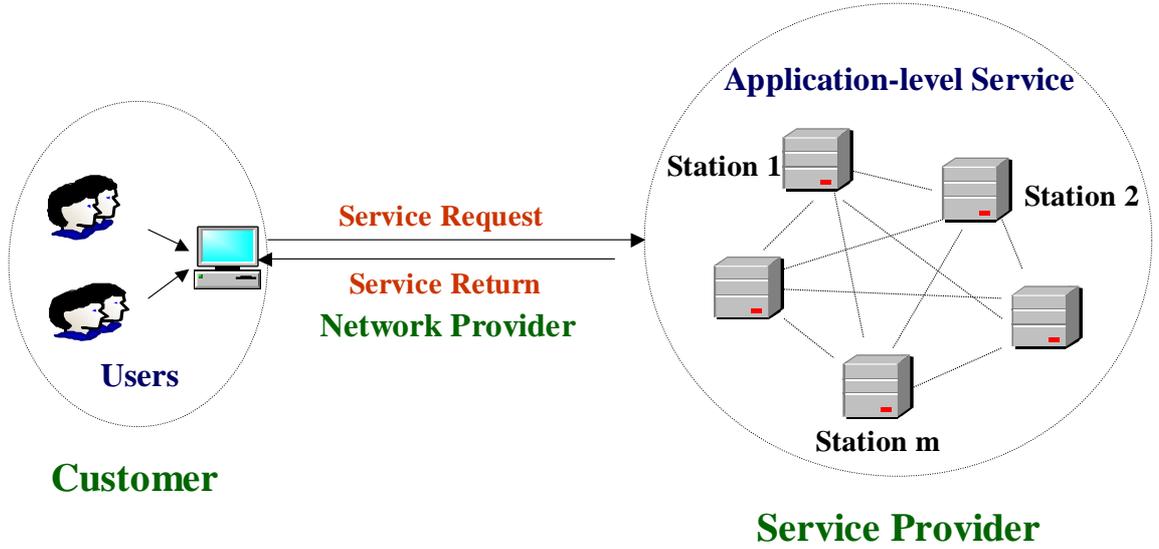


Figure 1.1: The Execution of Service Requests

cation is quantified by solving for  $n_j$  ( $n_j = 1, 2, \dots, N_j$ ) in the following optimization problem:

$$I = \min_{n_1, \dots, n_m} (n_1 c_1 + \dots + n_m c_m) \quad (1.1)$$

subject to SLA constraints. Performance and price are the two most important components for a variety of SLAs for business applications. Hence, in this dissertation, we first consider the minimization of the overall cost of the service provider's computing resources allocated to a multi-class business customer so that  $\gamma\%$  of the time the response time, i.e., the time to execute a service request, is less than a pre-defined value.

In enterprise computing, customer requests often need to be distinguished, with different request characteristics and customer's different service requirements. In this dissertation, we further consider a set of computer resources used by a service provider to host enterprise applications for differentiated customer services subject to an SLA.

Imposing a priority structure with preemption-resume is one way to implement a service for satisfying multiple class customer requests. A priority discipline does not depend on the state of a queue at the arrival of a customer, but is determined by a classification of arriving customers according to some criterion which is independent of the state of the queue. Suppose arriving customers to a single queue are classified into  $R$  different classes, where type  $r_1$  customers ( $r_1 = 1, 2, \dots, R - 1$ ) have always priority for service over those of type  $r_2$  ( $r_2 > r_1$ ;

$r_1, r_2 = 1, 2, \dots, R$ ), while customers of the same type are served in order of arrival (or FIFO). Then, it is said that the single queue operates according to a priority discipline with  $R$  classes.

In this dissertation, we consider a *preemptive-resume* priority discipline, that is, the service of a class  $r_2$  customer can be interrupted if a higher-priority customer of class  $r_1$  ( $r_2 > r_1$ ) arrives during its service. The interrupted customer resumes its service from where it stopped after the higher-priority customer, and any other customer with priority higher than  $r_2$  that may arrive during its service, complete their service. There are many other situations of practical interest (in the fields of modern computer systems, communication networks, computer server maintenance, and computer security checking, for example) in which the order of servicing is determined by preemptive-resume.

Second, we present an approach for the resource optimization that minimizes the total cost of computer resources required while preserving a given percentile of the response time for priority-class customers. We calculate the number of servers in each resource station that minimize a cost function that reflects operational costs.

Web services technology was introduced as a major component of .NET technology by Microsoft in June of 2000, and it is widely considered as an emerging paradigm for the next generation of Internet computing (Zhang et al. [97]). Web services technology has become the most popular computing paradigm for e-business applications and distributed environments. Many companies such as Google and Amazon are boosting their traffic using Web services APIs (Menasce 2004 [63]). By adopting service-oriented architectures (SOAs), services components from several universal service providers can be flexibly integrated into a composite service regardless of their location, platform, and execution speed (Barry [9], and Singh and Huhns [84]). In this type of enterprise service applications, service resources may be located in different places. Thus, security becomes a big concern. Hence, our study further addresses how to choose reliable resource stations from a customer's perspective.

Usually, at the beginning of a service, there is no pre-existing relationship between a customer and a service provider. In this research, "trust" is used to establish the relationship between a customer and a service provider. It is a firm belief in the competence of a resource station that acts as expected. The trustworthiness of resource stations is an indicator of the quality of services provided by these stations based on previous and current job completion experience. It often predicates the future behavior of the quality of services at these stations. Moreover, we only consider the trustworthiness of resource sites from a customer's perspective. Hence, in Chapter 5 of the dissertation, we simply assume that these resource stations trust each other.

In Chapter 5 of the dissertation, we further consider a trust-based resource allocation problem which typically arises in the aforementioned Web services applications. We formulate the trust-based resource allocation problem as an optimization problem under the SLA constraints: trustworthiness, performance and availability.

As discussed above, a trust model is suggested to establish the relationship between a customer and a service provider. However, receiving a reliable service response heavily depends on the security of the chosen service stations. When these chosen service stations are not secure, they can not provide reliable services for the customer. Authentication is the process of reliably verifying the identity of someone or something. It is an essential part in enterprise service computing. Authentication must be done before or with data communication.

Kerberos [52] consists of a client, an application server and a key distribution center (KDC). It is a mature, reliable, secure network authentication protocol that allows a client to prove its identity to a server without sending confidential data across the network. Public key cryptography has been extended to support Kerberos, since it simplifies the distribution of keys in Kerberos. It eliminates a single point of failure. Integrating public key cryptography into Kerberos represents the enhancements of the current Kerberos standard. Several Kerberos-based authentication techniques using public-key cryptography have been proposed in the last decade. Among them include Public-Key Cross Realm Authentication in Kerberos (PKCROSS) [46] and Public-key Cryptography for Initial Authentication in Kerberos (PKINIT) [99]. Moreover, the scalability of network security infrastructures is becoming a serious concern as the explosive growth of collaborative applications such as Web services continues unabated. Public key based Kerberos for Distribution Authentication (PKDA) [85] and Public Key Utilizing Tickets for Application Servers (PKTAPP, a.k.a. KX.509/KCA) [53] and [61] have been proposed to enhance the security and scalability of Kerberos. But, the actual costs (e.g., computational and communication times) associated with these techniques have been poorly understood so far. It remains unknown which technique performs better in a large network where there are multiple KDC remote realms. Both PKCROSS and PKTAPP use variations of PKINIT message types and data structures for integrating public key cryptography with Kerberos in different authentication stages. PKTAPP was originally introduced as PKDA. It implemented PKDA using the message formats and exchanges of PKINIT. Hence, this research only considers these two notable techniques: PKCROSS and PKTAPP.

In this dissertation, we first present a thorough performance evaluation of PKCROSS and PKTAPP in terms of computational and communication times. Then, we demonstrate their performance difference using open queueing networks. An in-depth analysis of these two tech-

niques shows that PKTAPP does not scale better than PKCROSS. Thus, we propose a new public key cryptography-based group authentication technique. Our performance analysis demonstrates that the new technique can achieve better scalability as compared to PKCORSS and PKTAPP.

### 1.3 Summary of Contributions

The contributions of this dissertation are summarized below:

- *End-to-end response time*: The calculation of the response time often becomes critical in solving the resource optimization problem. Existing work addressing resource allocation uses the average response time (or average execution time). Though the average response time is relatively easy to calculate, it does not address the concerns of a customer. Typically, a customer is more inclined to request a statistical bound on its response time than an average response time. To obtain the statistical bound on its response time we need to calculate the probability density function of the end-to-end response time which is not an easy task in a complex computing environment involving many computing nodes. We develop effective and accurate numerical solutions of this probability density function in both a tandem queueing network and a queueing network with feedback. which we then incorporated in this resource allocation problem.
- *Multi-class queueing networks*: In enterprise computing, customer requests often need to be distinguished, with different request characteristics and service requirements. Imposing a priority structure with preemption-resume is one way to implement a service for satisfying multiple class customer requests. However, it is difficult to calculate the probability density function of an end-to-end response time even for a single node in the case of multiple priority customers. In our study, we first develop an efficient and accurate numerical solution for inverting the Laplace-Stieltjes transforms (LSTs) of a multi-class priority customer's response time numerically at a single node, and then extend the result to multi-class queueing networks under our study. A contributing factor is that typically we are interested in values of the cumulative distribution of the response time that correspond to very high percentiles for which our approximate results seem to have a very good accuracy.
- *Trust-based resource optimization in Web services*: Most existing Web services products do not support an SLA that guarantees a level of service delivered to a customer for a given

price. It is not easy to solve resource allocation problem when we consider all the constraints of trustworthiness, an end-to-end response time and service availability. We provide an efficient approach to solving the trust-based resource optimization problem in Web services whereby we are required to minimize the total cost of service providers under the constraints of trustworthiness, an end-to-end response time and service availability.

- *Performance analysis of public key cryptography-based group authentication:* Several authentication techniques have been proposed in the last decade. But, the actual costs (e.g., computational and communication times) associated with these techniques have been poorly understood so far. It remains unknown which technique performs better in a large network where there are multiple KDC remote realms. Based on complexity analysis and queueing theory, the dissertation proposes a performance methodology that is an effective way to analyze the performance of security protocols and suggests a new group authentication that improves the scalability of PKCROSS and PKTAPP. It should be pointed out that the proposed performance method can be extended to analyze other security techniques as well.

## 1.4 The Organization of this Dissertation

The rest of the dissertation is organized as follows. Related work is reviewed in Chapter 2. We give the solution of the resource optimization problem with numerical validations for single class customers in Chapter 3 and for multiple class customers in Chapter 4. Chapter 5 addresses a trust-based resource allocation problem, and provides an approach to solving this problem for single and multiple class customers. Numerical examples are given that demonstrates the validity of this approach for both cases. In Chapter 6, we give an in-depth performance evaluation of authentication techniques: PKCROSS and PKTAPP (a.k.a. KX.509/KCA) by using complexity analysis and queueing networks, and propose a new group authentication technique. We conclude our results and discuss future work in Chapter 7.

## Chapter 2

# Related Work

In this chapter, we discuss existing approaches to addressing SLA metrics including service availability, trustworthiness, performance, and public key cryptography-based group authentication. Then, we provide a review of how a resource optimization problem subject to an SLA is solved in the literature.

The chapter is organized as follows. Section 2.1 gives the definition of service availability and describes its measure methods. The characteristics of trustworthiness and performance metrics is presented in Sections 2.2 and 2.3. Section 2.4 gives a literature review of the resource optimization problem subject to various QoS metrics. The performance of public key cryptography-based group authentication is discussed in Section 2.5.

### 2.1 Service Availability

Availability is a critical metric in today's computer design (Hennessy [45]). It is the percentage of time that a service provider can offer services. A computer system is unavailable due to a variety of causes, such as network failure, hardware failure, software failure, or security attacks. Detecting and preventing these failures and attacks is beyond the scope of our study in this dissertation. Cisco has asserted that the operational failure causes 80% of non-availability [26]. Hence, increasing network availability is becoming a key priority for enterprise and service provider organization, as discussed in [24]. Martin and Nilsson [58] gives an example of how network service

availability is defined in Sprint’s SLA and WorldCom’s SLA.

Availability has been extensively studied for a variety of computer systems in the literature. It has been studied to improve dependability for computer and telephone networks in Gray [40]. The dependability can be defined as the property of a system such that reliance can justifiably be placed on the service it delivers, as defined in Barbacci et al. [8]. Systems with high availability tend to have large quorum<sup>1</sup> sizes and high load. Improving the availability of a system has been discussed in the literature (e.g., see [3], [59], [71], and [86]). Aiyer et al. [3] studied the availability of non-strict quorum systems and proposed  $K$ -quorums that can provide higher availability than the strict quorum systems. Naor and Wool [71] analyzed the load and the availability of traditional quorum systems.

Brown and Patterson [16] defined a new availability metric to capture the variations of the system quality of service over time. It is defined by the number of requests satisfied per second (or the latency of a request service) and the number of server failures that can be tolerated by a system.

In this research, our interest is potential hardware (i.e, servers within each station) failures and their effect on unavailability. To determine this, the service provider who owns these stations needs to understand the *MTBF* (Mean Time to Failure) of all station components and the *MTTR* (Mean Time to Recover) for hardware problems for all devices at each station, where *MTBF* (Mean Time to Failure) is the average time of a server failure, and *MTTR* (Mean Time to Recover) is the average time for recovering a server at each resource station. *MTBF* information can be obtained from a hardware provider. For example, it is mentioned in Cisco [26] that MTBF information is available for all Cisco components and is available upon request to a local account manager. *MTTR* is determined by evaluating how quickly a station owner can repair broken servers. It is a major factor of server availability. To improve service availability, it is necessary to reduce the frequency time of failure, as indicated in Brewer [15].

Network availability data may be found in the Internet. For example, the University of Houston maintains current and historical network availability data on a web site [60]. CSCO [25] presented a formula for the calculation of network availability.

In this dissertation, we consider the percentage of time that a resource is “up” or “down” as a metric, which is the traditional way to define service availability. Then, a two-state Markov chain with the states “up” and “down” is proposed to study the service availability at each service station. The detailed analysis and calculation of the service availability is given in Section 5.3.3 of

---

<sup>1</sup>A quorum system is a collection of sets called quorums such that any two quorums have a non-empty intersection.

Chapter 5.

## 2.2 Trustworthiness

“Trust” is used to deal with the notion of the trustworthiness. In this dissertation, trust is defined as a firm belief in the competence of a resource station to act as expected. Many researchers have studied the trustworthiness of service entities, and suggested several different trust metrics. Zhang et al. [98], and Ziegler and Lausen [100] classified various trust metrics. Vu et al. [89] proposed a new QoS-based semantic Web services selection and ranking solution using a trust and reputation management and assuming known QoS qualities.

The trustworthiness of resource stations is an indicator of the quality of services provided by these stations based on previous and current job completion experience. It often predicates the future behavior of the quality of services at these stations. Most trust models (e.g., Kamvar et al [49], Lee et al [55] and Richardson et al. [81]) assume that the domain of trustworthiness ranges from 0 to 1, which is considered in this research. In this dissertation, we use a rank and a threshold-based approach. That is, according to the trust information of each station maintained by the trust manager, our study addresses how the trust manager selects service stations using a rank and threshold based approach. The rank-based approach is to rank the trustworthiness of all sites who provide the same type of services, such as EigenRep in Kamvar et al. [49]. A threshold-based approach is to choose any of those service sites who can meet trust requirements predefined by customers. For example, if the trustworthiness of a service site is over 0.9 and the predefined trust requirement of a customer is 0.9, then the trust manager can select the site to serve the customer’s service request. The relationship between the rank and the threshold based approaches are discussed in Zhang et al. [98].

Our study introduces a trust manager who represents customers. The trust manager uses the collected trustworthy information of the resource stations to evaluate their security behavior. We consider security behavior by modeling the behavior trusts of all sites, and quantify the trustworthiness of these stations. This approach is based on previous job completion experience assessed by the trust manager and customers. There often exist multiple trust managers in today’s complex computer systems. The assessment also adopts the opinion of those trust managers besides its own customer’s feedback. The domain of feedback is also assumed to be [0, 1].

The feedback is a statement issued by the trust manager’s customers about the quality

of service provided by those chosen service stations for each service request or for a set of service requests during a certain period of time. These customers only provide feedback about their received services rather than those chosen service sites. (Note that the trust manager's customer is usually not aware of which service sites process its service request.) The opinion is defined as the information of trustworthiness provided by those trust managers who are neighbors of the trust manager. These neighbors are a small subset of the trust manager's acquaintances, adaptively selected based on their usefulness. The opinion aggregates the overall impression of those neighborhood trust managers for the service stations. A similar definition is given in Golbeck and Hendler [39], Kamvar et al [49], Lee et al [55], Mui et al [67], and Richardson et al. [81].

In this research, we consider all these factors by combining the feedback of the trust manager's customers and the opinions of neighborhood trust managers with the past trust information of the trust manager. The detailed discussion of updating the trust information for each station will be discussed in Section 5.2 of Chapter 5.

## 2.3 Performance

The SLA performance metric defined in Section 1.1 includes *throughput* and *response time*. As an end user, a customer is in general concerned about response time rather than throughput. So, in the dissertation we only consider the percentile of the response time as the performance metric. This is the time it takes for a service request to be executed on the service provider's multiple resource stations.

To obtain the statistical bound on its response time, i.e., the percentile of the response time, we are required to calculate the probability density function of the response time, whose detailed discussion will be given in Section 3.1 of Chapter 3.

The calculation of the response time often becomes critical in solving the resource optimization problem. The computation of the response time has been extensively studied for a variety of computing systems. However, only the average response time is calculated rather than a percentile of the response time. Menasce and Almeida [64] designed self-managing systems to control QoS. Levy et al. [56] presented a performance management system for cluster-based Web services. In both studies the average response time is used as a metric. Chandra [21] employed an online measurement method, and considered a resource allocation problem based on measured response times.

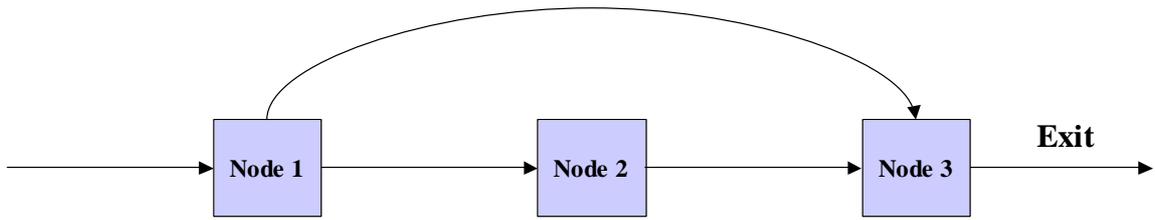


Figure 2.1: A Three-node Tandem Network with Overtaking

Martin and Nilsson [58] measured the average response time of a service request over IP networks. In [74], Osogami and Wierman analyzed an  $M/GI/k$  system with two priority classes and a general phase-type distribution and evaluated the optimal number of servers based on overall *mean* response time. A framework for service management in grid computing was defined in Menasce and Casalicchio [65], but they did not provide a method for calculating the probability distribution of the response time.

In order to compute a percentile of the response time one has to first find the probability distribution function (pdf) of the response time. This is not an easy task in a complex computing environment involving many computing nodes. The calculation of the pdf of the response time is relatively simple for a tandem network in Reich [80] and [79] and overtake-free paths in Jackson and Gordon-Newell networks (Walrand and Varaiya [90] and Daduna [30]). Reich [80] proved that the response times of a customer in each of two  $M/M/1$  queues of a tandem network are independent, and he later extended the result to an arbitrary number of such queues in a tandem network [79]. However, Reich [80] proved that the result cannot be extended to a queueing network with  $E_2/E_2/1$  queues where  $E_2$  is a two-stage Erlangian distributions. Burke [18] further showed that in a tandem network with exponential server queues and customers arriving in a Poisson stream, the first and last nodes of which are multi-server queues, while all other nodes are single-server queues, the response times of nodes are independent. Moreover, Walrand and Varaiya [90] generalized this result and proved that in any single-server open Jackson network, the response times of a customer at the various nodes of an overtake-free path are all mutually independent where the service discipline is first come first served (FCFS). They also proved that when a three-node tandem network has two parallel paths which are not overtaken-free as shown in Figure 2.1, the response time of nodes 1 and 2 are independent, and the response time of nodes 2 and 3 are independent, but the response time of nodes 1 and 3 are not independent. Daduna [30] further showed that the same result is valid for overtake-free paths in Gordon-Newell networks with multiple class customers. Additionally, the

pdf of the response time was derived for a closed queueing network in [69], and the passing time distribution was calculated for large Markov chains in [43].

We develop effective and accurate numerical solutions of the probability density function of the response time for a variety of service models under our study, which we then incorporated in this resource allocation problem in Chapters 3 - 5.

## 2.4 The Resource Optimization Problem Subject to an SLA

Resource optimization problems subject to performance metrics such as response time, throughput, packet loss rate, and link utilization have been extensively studied [1], [14], [13], [36], [58], [65], [73] and [93]. Link utilization is defined as the percentage of the time that the network node is utilized. Packet loss rate is the probability that the packet is lost during transmission through a network since its buffer is full.

The resource optimization problems subject to some of these performance metrics (or QoS metrics) was studied for multiple packet networks in Bouillet, et al. [14], different classes of flows at network nodes in Chassot, et al. [23], IP networks in Martin and Nilsson [58], wireless networks in Aib, et al. [1], and grid computing in Menasce and Casalicchio [65], and optical networks in Nowak, et al. in [73]. Xiong and Perros [93] studied network optimization subject to the percentile of response time, network availability, network utilization and packet loss rate.

These QoS metrics can be estimated by using measurement techniques, see Aikat et al. [2], Allman et al. [4], Gummadi et al. [41] and Sommers et al. [47]. In [23], Chassot et al. dealt with a communication architecture with guaranteed end-to-end QoS in an IPv6 environment. The end-to-end QoS includes an end-to-end delay (i.e., a response time). Chassot et al. only discussed and measured the maximal, minimal and average values of the response time. In [73], Nowak, et al. developed a C++ event driven simulator to measure these three values in optical networks. But, measurement techniques are hard to be incorporated in solving a resource optimization problem with multiple constraints.

Calabrese [20] studied optimal workload allocation at each station in open networks of multi-server queues and formulated it as two nonlinear programming problems. One of their approaches is to maximize average throughput, subject to average number of jobs in the queues. The second approach is to minimize the average number of jobs in the queues, given a fixed average throughput. In [83], Shanthikumar and Yao discussed a server allocation problem in a single class,

closed queueing network. The problem was formulated as an optimization problem in which the average throughput is maximized.

Web services is often contracted through SLAs which typically specify a certain QoS in return for a certain price. Although QoS was not defined in the initial UDDI standard for Web services, many studies have been carried out to extend the initial UDDI, such as WSLA (Keller and Ludwig [50]), WSOL (Tosic et al. [88]), and QML (Dobson [33]). The issue of a reputation-based SLA has been studied by Jurca and Faltings [48] in which the cost is determined by the QoS that was actually delivered.

Resource allocation problems in distributed systems have been widely studied based on one of the metrics: trustworthiness and response time (e.g., see [70]). In Chapter 5 we will consider a resource allocation problem in Web services subject to all three QoS metrics, expressed by trustworthiness, percentile response time and service availability.

When we consider all qualities of services consisting of trustworthiness, percentile response time, and service availability, it becomes extremely difficult to solve a trust-based resource allocation problem whereby a set of computer resources is used by a service provider to host a typical distributed computing application for single-class and multiple-class customer services respectively. To the best of our knowledge, our research is the first attempt toward solving the complex problem. We shall formulate the trust-based resource allocation problem as an optimization problem under all these SLA constraints in which we calculated the number of servers in each service station that minimize a cost function that reflects operational costs for single-class and multiple-class customer services respectively in Chapter 5.

## 2.5 Public Key Cryptography-based Authentication

Kerberos has revolved over the past 15 years and rapidly since 1999. There have been numerous proposals to integrate public key cryptography into Kerberos [30], [46], [61], [72], [99], and [85]. These proposals address various concerns of Kerberos in distributed networks, for instance, security, scalability, and portability. Neuman, et al. [72] proposed PKINIT to enable use of public key cryptography for an initial authentication between the client and its local KDC (Key Distribution Center). PKINIT is an extension of the Kerberos protocol (RFC1510 [52]) and its mechanism hasn't been standardized yet and the specification is still under development directed by the IETF Kerberos WG (see Zhu and Tung [99]). Due to the cost of using public key cryptogra-

phy, Davis [32] suggested that public key cryptography is best suited to securing communications between servers, between sites, and between organizations.

PKCROSS [46] has been proposed to simplify the administrative burden of maintaining cross-realm keys so that it improves the scalability of Kerberos in large multi-realm networks. Public key cryptography takes place only KDC-to-KDC authentication in PKCROSS. PKINIT and PKCROSS are centralized KDC protocols. Sirbu and Chuang [85] extended these two protocols to create PKDA for improving scalability and addressing the single point of failure on the KDC. PKTAPP is only a slight variation on the PKDA specification. Both PKDA and PKTAPP use lengthy public-key message exchanges between the client and the application servers, so they may not be anymore efficient than public key enabled authentication once with a KDC and faster secret-key cryptography for subsequent encryption with application servers [42]. PKTAPP is also known as KX.509/KCA [35] and Windows has its own protocol which is the equivalent to KX.509/KCA (see Altman [6]). Although the evolution of PKTAPP, its protocol structure has not been dramatically changed. We believe that PKCROSS and PKTAPP will revive soon. The Kerberos Consortium at MIT was just formed in 2007 and listed PKCROSS as Project 10 in its proposal [28].

Performance evaluation is a fundamental consideration in the design of security protocols. However, performance associated with most of these protocols has been poorly understood. To analyze the performance of a protocol, we can first implement the protocol and then analyze measurement data taken from the implementation. But, the implementation of a protocol is very time-consuming and constricted to development resources and funding. While the implementation of KX.509 was just released a while ago [53], up to today, PKCROSS has not been implemented yet due to a matter of lack of development resources and funding [7]. Hence, performance modeling has become an attractive approach since it can quickly provide a performance guidance used for a protocol design. Existing studies in Harbitter and Menasce [42] employed the performance modeling approach for examining the impact of PKCROSS and PKTAPP on network throughput based on their skeleton implementations and construction of *closed* queueing networks for a relatively simple case: there is only a single remote realm. In the present research we also employ a performance modeling approach for a study of PKCROSS and PKTAPP but extends Harbitter and Menasce [42] in several essential and important aspects. First, complexity analysis has been often used to evaluate algorithm performance (e.g., [11]). While the performance of PKCROSS and PKTAPP was studied in Harbitter and Menasce [42], it remains unknown which technique performs better in multiple remote realms that are typical in an increasingly large network these days. It is much difficult to analyze the case of multiple remote realms due to the complexity of authentication

message exchanges. The difficulty is in the complexity analysis of these protocols and the building and analysis of queueing network models that reflects the workload of authentication requests for these protocols. Second, we explicitly derive the formulas for calculating the computational and communication times of these protocols so as to easily determine which technique is better. Third, using a closed queueing network in Harbitter and Menasce [42] assumes there exist *constant* authentication requests in the queueing network. Although the assumption can simplify our performance analysis, it is not usually true in a real world. Rather, we adopted an *open* queueing network where the client requests authentication at a given rate, i.e., the number of authentication requests in a computing system under study is *not constant*. Fourth, due to a performance trade off between these two protocols according to our scalability analysis, we propose a new hybrid technique. Our analysis shows that the new technique has better scalability than these two protocols in most cases. We discuss the open queueing network, give the complexity analysis of PKCROSS and PKTAPP, and present the new authentication protocol in Chapter 6.

## Chapter 3

# Single Class Customers

This chapter gives the definition of the percentage of response time with an illustrative example, and discusses a resource optimization problem subject to the percentile of response time and a fee for service models with single class customers. We calculate the number of servers in each resource station that minimize a cost function that reflects operational costs. We first analyze an overtake-free open tandem queueing network, and then we extend our work to an open tandem queueing network with feedback. The chapter is mainly based on the results presented in Xiong and Perros [95], and Xiong in [92].

The rest of the chapter is organized as follows. In Section 3.1 we define the SLA performance metric considered in this chapter. Section 3.2 formulates the resource optimization problem. In Section 3.3, we present two typical real-life models and propose an approach for solving the optimization problem. In Section 3.4, numerical simulations demonstrate the applicability and validity of the proposed approach. Finally, the conclusions are given in Section 3.7.

### 3.1 The Percentile of Response Time

An SLA is a contract between a customer and a service provider that defines all aspects of the service that is to be provided. An SLA generally uses response time as one performance metric.

As discussed in Section 2.3, in this research we are interested in the percentile of the response time. This is the time it takes for a job to be executed in a computing environment consisting

of multiple computing nodes.

Assume that  $f_T(t)$  be the probability distribution function of a response time  $T$ .  $T^D$  is a desired target response time that a customer requests and agrees with its service provider based on a fee paid by the customer. The SLA performance metric that a  $\gamma\%$  SLA service is guaranteed is as follows.

$$\int_0^{T^D} f_T(t) dt \geq \gamma\% \quad (3.1)$$

That is,  $\gamma\%$  of the time a customer will receive its service in less than  $T^D$ .

As an example let us consider an  $M/M/1$  queue with an arrival rate  $\lambda$  and a service rate  $\mu$ . The service discipline is FIFO. The steady-state probability of the system is

$$p_0 = 1 - \rho,$$

and

$$p_k = (1 - \rho)\rho^k, \quad k > 0$$

where  $\rho = \frac{\lambda}{\mu}$ .

The response time  $T$  is exponentially distributed with the parameter  $\mu(1 - \rho)$ , i.e., its probability distribution function is given by (see [12] and Perros [76])

$$f_T(t) = \mu(1 - \rho)e^{-\mu(1-\rho)t}$$

Using the definition given in (3.1), we have that

$$\int_0^{T^D} f_T(t) dt = 1 - e^{-\mu(1-\rho)T^D} \geq \gamma\% \quad (3.2)$$

or

$$\mu \geq \frac{-\ln(1 - \gamma\%)}{T^D} + \lambda \quad (3.3)$$

This means that in order to guarantee higher SLA service levels,  $\mu$  increases when  $T^D$  decreases.

Similarly, for any given arrival rate  $\lambda$  and service rate  $\mu$ , we can use (3.2) to find the percentile of  $\gamma$ .

For example, when  $\lambda = 100$  and  $T^D = 0.05$ , Figure 3.1 gives the cumulative distribution of the response time. As shown in Figure 3.1, the percentile of response time will increase as the service rate increases. Table 3.1 gives the numerical values for the cumulative distribution of the response time. From this table we see that this service rate has to be bigger than 150 in order that 90% of the response time is less than 0.05.

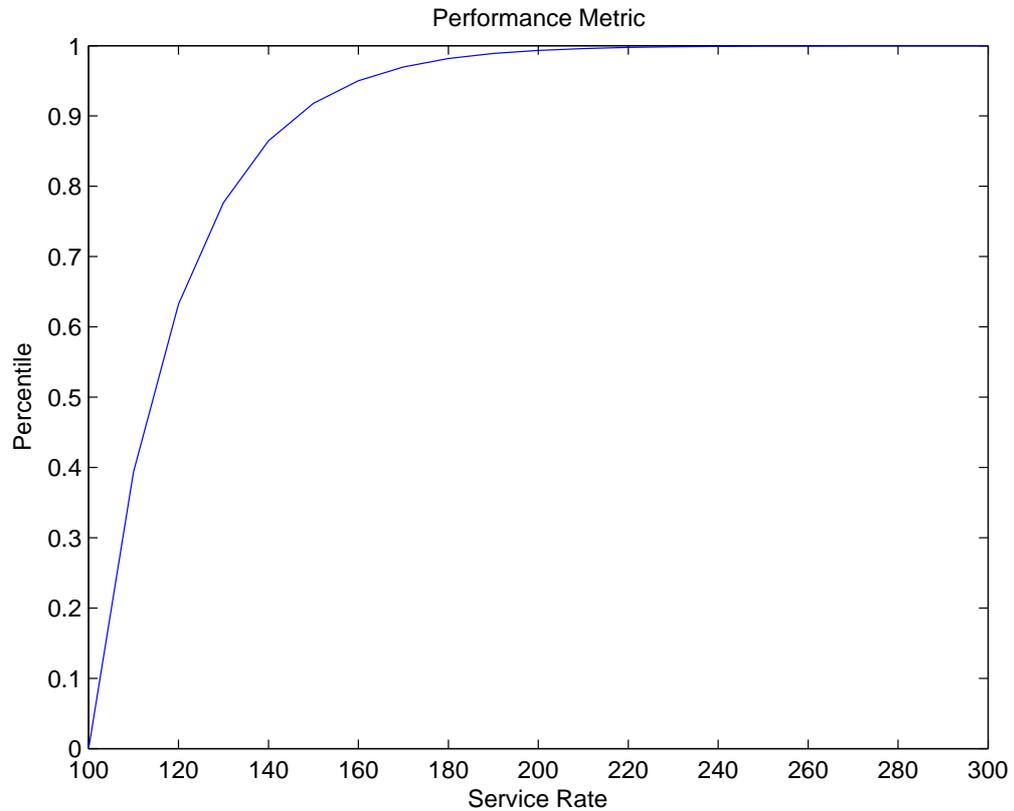


Figure 3.1: Percentile Response Time vs. Service Rate

## 3.2 A Resource Optimization Problem for Service Models with Single Class Customers

In this section, we consider the minimization of the overall cost of the service provider's computing resources allocated to the single class customer so that  $\gamma\%$  of the time the response time, i.e., the time to execute a service request, is less than a pre-defined value. Based on Section 3.1, the resource optimization problem can be formulated as the following optimization problem:

### Resource Optimization Problem:

Find integers  $n_j$  ( $1 \leq n_j \leq N_j; j = 1, 2, \dots, m$ ) in the  $m$ -dimensional integer optimization problem (1.1) under the constraint of a percentile response time as expressed by (3.1), and the constraint:  $I \leq C^D$ , where  $C^D$  is a fee negotiated and agreed upon between a customer and the

Table 3.1: The Cumulative Distribution Function (CDF) of the Response Time vs. Service Rate

Service Rate	100	120	140	150
CDF	0.0000	0.6321	0.8647	0.9179
Service Rate	160	170	180	200
CDF	0.9502	0.9698	0.9817	0.9933
Service Rate	220	240	280	300
CDF	0.9975	0.9991	0.9999	1.0000

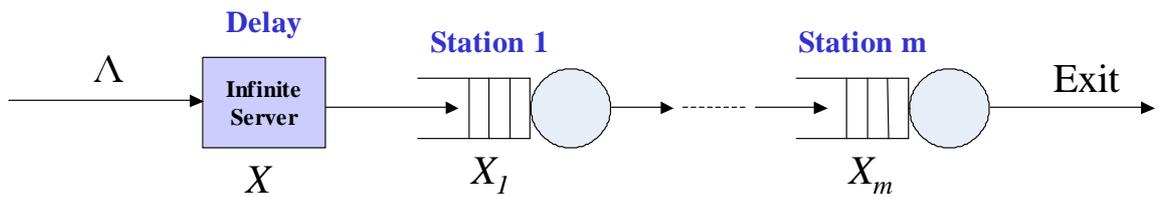


Figure 3.2: A Tandem-station Service Model for Single-class Customer Services

service provider.

### 3.3 Approaches for the Resource Optimization

In this section, we study two queueing network models that depict the path that service requests have to follow through the service provider's resource stations. These two models are shown in Figures 3.2 and 3.3. We refer to these two queueing models as service models since they depict the resources used to provide a service to a customer.

The first service model consists of a single infinite server and  $m$  stations numbered sequentially from 1 to  $m$  as shown in Figure 3.2. Each station  $j$  is modeled as a single FIFO queue served by  $n_j$  identical servers, each providing a service at the rate  $\mu_j$ . Let  $\Lambda$  be the external arrival rate to the infinite server, and let  $\lambda$  and  $\lambda_j$  be the effective arrival rates to the infinite server and station  $j$  ( $j = 1, 2, \dots, m$ ). We assume that all service times are exponentially distributed and the external arrival to the infinite server occurs in a Poisson fashion.

The infinite server represents the total propagation delay from the user to the service provider and back and also from station 1 to  $m$ . Each station carries out a particular function. For instance, it could be a database server, a file server, a web server, a group of CPUs and local disks,

etc. We only consider a single class of customer in this research.

In the following discussion each station is modeled as a single  $M/M/1$  queue with arrival rate  $\lambda_j$  and service rate  $\psi(n_j)\mu_j$ , where  $\psi(n_j)$  is a known function of  $n_j$  and depends on the configuration of servers at each station. It is non-decreasing and can be inverted, i.e.,  $\psi^{-1}$  exists. For instance, suppose that a station represents a group of CPUs. Then,  $\psi(n)$  can be seen as a CPU scaling factor for the number of CPUs from 1 to  $n$ . According to [22],  $\psi(n) = \xi^{\log_2 n}$ , where  $\xi$  is a basic scaling factor from 1 CPU to 2. So,  $\psi^{-1}(n) = \xi^{-\log_2 n}$ .

Since the queueing network is overtake-free, the waiting time of a customer at a station is independent of its waiting times at other stations (see [30] and [90]). Let  $X$  be the service time at the infinite server and  $X_j$  be the time elapsed from the moment a customer arriving at station  $j$  to the moment it departs from the station. Then, the total response time is

$$T = X + X_1 + X_2 + \cdots + X_m,$$

and hence the LST (Laplace-Stieltjes transform) of the response time  $T$  is

$$L_T(s) = L_X(s)L_{X_1}(s) \cdots L_{X_m}(s) \quad (3.4)$$

where  $L_X(s)$  is the LST of the service time  $X$  given by

$$L_X(s) = \frac{\lambda}{s + \lambda} \quad (3.5)$$

and  $L_{X_j}(s)$  is the LST of the response time  $X_j$  at the  $j$ -th station given by

$$L_{X_j}(s) = \frac{\psi(n_j)\mu_j(1 - \rho_j)}{s + \psi(n_j)\mu_j(1 - \rho_j)}, \quad (3.6)$$

where  $\rho_j = \frac{\lambda_j}{\psi(n_j)\mu_j}$  ( $j = 1, 2, \dots, m$ ).

From (3.4), (3.5) and (3.6) we have that

$$L_T(s) = \frac{\lambda}{s + \lambda} \prod_{j=1}^m \frac{\psi(n_j)\mu_j(1 - \rho_j)}{s + \psi(n_j)\mu_j(1 - \rho_j)}$$

We observe that  $f_T(t)$  and  $F_T(t)$  are usually nonlinear functions of  $t$  and  $n_j$ . Hence, the resource optimization problem is an  $m$ -dimensional linear optimization problem subject to nonlinear constraints. In general, it is not easy to solve this problem. However, the complexity of the problem can be significantly reduced by requiring that the service rates of the queues in the service model in Figure 3.2 are all equal. That is, we find the optimum value of  $n_1, \dots, n_m$  such that

$$\psi(n_1)\mu_1 = \cdots = \psi(n_m)\mu_m$$

called *balanced utilization* or *balanced condition*. (We note that in production lines, it is commonly assumed that the service stations are balanced.)

From the traffic equations:

$$\lambda = \lambda_j = \Lambda$$

for  $j = 1, 2, \dots, m$ , we have that the utilization of each station

$$\rho_j = \frac{\lambda_j}{\psi(n_j)\mu_j} = \frac{\Lambda}{\psi(n_j)\mu_j}$$

Thus we have

$$\hat{a}_i = \psi(n_i)\mu_i(1 - \rho_i) = \psi(n_j)\mu_j(1 - \rho_j) = \hat{a}_j \triangleq \hat{a},$$

which implies  $n_j = \psi^{-1}(\frac{\hat{a} + \lambda_j}{\mu_j})$  ( $i, j = 1, 2, \dots, m$ ). Hence, from (3.4) we have

$$f_T(t) = L^{-1}\left\{\frac{\lambda}{s + \lambda} \cdot \frac{\hat{a}^m}{(s + \hat{a})^m}\right\},$$

and subsequently obtain that

$$F_T(t) = L^{-1}\left\{\frac{\lambda}{s(s + \lambda)} \cdot \frac{\hat{a}^m}{(s + \hat{a})^m}\right\} \quad (3.7)$$

Consequently,  $\sum_{j=1}^m n_j c_j$  reduces to a function of variable  $\hat{a}$  due to  $n_j = \lceil \psi^{-1}(\frac{\hat{a} + \lambda_j}{\mu_j}) \rceil$ .

Thus we have the following algorithm for the resource optimization problem.

**Algorithm 3-1:**

1. Find  $\hat{a}$  in the one dimensional optimization problem:

$$\hat{a}^{min} \leftarrow \arg \min_{\hat{a}} F_T(t)|_{t=T^D}$$

subject to the constraint  $F_T(t)|_{t=T^D} \geq \gamma\%$  at  $\hat{a} = \hat{a}^{min}$ , where  $F_T(t)$  is given by (3.7).

2. Compute integers  $n_j$  by using

$$n_j = \lceil \psi^{-1}\left(\frac{\hat{a}^{min}}{\mu_j(1 - \rho_j)}\right) \rceil,$$

and check if  $1 \leq n_j \leq N_j$  ( $j = 1, 2, \dots, m$ ) and  $I \leq C^D$  are satisfied. If yes, the obtained  $n_j$  is the number of servers required at each station. Otherwise, print “the problem cannot be solved.”

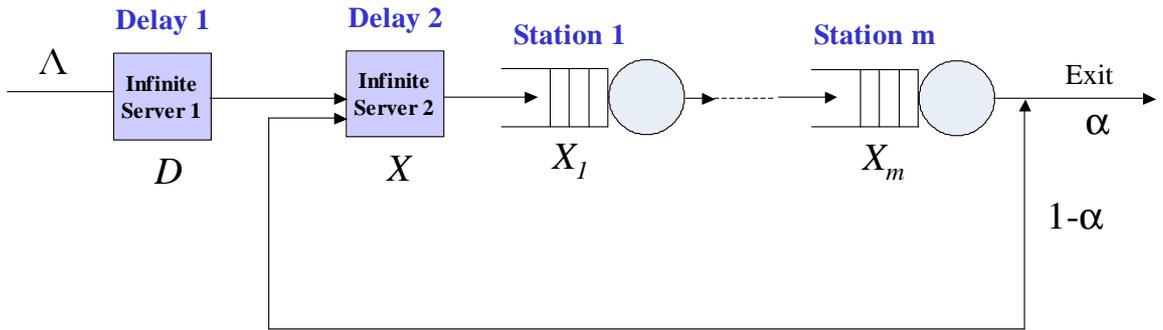


Figure 3.3: A Service Model with Feedback for Single-class Customer Services

In order to allow for a more complex execution path of a service request, we extended the above model to a service model with feedback, as shown in Figure 3.3. Infinite server 1 represents the total propagation delay within a network provider and infinite server 2 represents the propagation delay within the service provider, i.e. among stations 1 to  $m$ . In this figure, a customer upon completion of its service at the  $m$ -th station, exits the system with probability  $\alpha$ , or returns to the beginning of the system with probability  $1 - \alpha$ .

We note that the model shown in Figure 3.3 can be easily extended to a network of queues arbitrarily connected. We reuse the notation in the first model shown in Figure 3.2:  $\Lambda$  as the external arrival rate,  $\lambda_d$ ,  $\lambda$  and  $\lambda_j$  as the effective arrival rates to the second infinite server and station  $j$ , and  $\mu_j$  as the service rate at station  $j$ , where  $j = 1, 2, \dots, m$ .

The main difficulty of this resource optimization problem is to find  $f_T(t)$ , the probability distribution function of  $T$ . We obtain this probability distribution function assuming that the waiting time of a customer at a station is independent of its waiting time at other stations, and each visit at the same station  $j$  is independent of the others. (We note that this assumption of independence does not hold in queueing networks with feedback. However, as will be discussed in the validation section the solution obtained has a good accuracy.) We first have the traffic equations:

$$\lambda_d = \Lambda, \quad \lambda = \Lambda + (1 - \alpha)\lambda_m, \quad \text{and} \quad \lambda_j = \lambda,$$

which implies  $\lambda_j = \lambda = \frac{\Lambda}{\alpha}$ , and the utilization of each station is

$$\rho_j = \frac{\lambda_j}{\psi(n_j)\mu_j} = \frac{\Lambda}{\alpha\mu_j\psi(n_j)} \quad (j = 1, 2, \dots, m)$$

Furthermore, the response time of the  $k$ -th pass at the infinite station and the  $j$ -th station

is considered as the sum of  $m + 2$  random variables

$$T(k) = D + X + X_1 + \cdots + X_m,$$

where we assume that the waiting time of a customer at a station is independent of its waiting times in other visits to the same station.  $D$  and  $X$  are the service times at the first and second infinite servers respectively, and  $X_j$  is the time elapsed from the moment a customer arrives at station  $j$  to the moment it departs from it. Then, the total response time is

$$T = \sum_{k=0}^{\infty} p(k)T(k),$$

where  $p(k)$  is the steady state probability that a request will circulate  $k$  times at the infinite station and the  $j$ -th station through the computing system.  $p(k)$  is determined by

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

Thus the LST of the response time  $T$  is

$$L_T(s) = L_D(s) \sum_{k=0}^{\infty} p(k) L_X^k(s) L_{X_1}^k(s) \cdots L_{X_m}^k(s),$$

which can be re-written as follows:

$$L_T(s) = \frac{\alpha L_D(s) L_X(s) \prod_{j=1}^m L_{X_j}(s)}{1 - (1 - \alpha) L_X(s) \prod_{j=1}^m L_{X_j}(s)} \quad (3.8)$$

where  $L_D(s)$  is the LST of the service time  $D$  given by

$$L_D(s) = \frac{\Lambda}{s + \Lambda},$$

and replacing  $L_X(s)$  and  $L_{X_j}(s)$  ( $j = 1, 2, \dots, m$ ) by (3.5) and (3.6) in (3.8), we have that

$$L_T(s) = \frac{\Lambda^2 \prod_{j=1}^m \hat{a}_j}{(s + \Lambda)[(s + \lambda) \prod_{j=1}^m (s + \hat{a}_j) - (1 - \alpha) \lambda \prod_{j=1}^m \hat{a}_j]} \quad (3.9)$$

To find the response time distribution  $f_T(t)$ , we are required to invert the above LST using partial fraction decomposition of a rational function. However, the partial fraction decomposition of the rational function requires searching for roots of a high-order polynomial. It is usually not an easy task when the order of the polynomial is more than 5. Instead, in this research the LST is inverted numerically.

Similarly, we want to find  $n_1, \dots, n_m$  such that the best utilization of these stations is achieved, which implies that each station has the same maximal service capacity. That is,

$$\hat{a}_i = \hat{a}_j = \hat{a} \quad (i, j = 1, \dots, m)$$

Then, from equation (3.9) and  $F_T(t) = L^{-1}\{L_T(s)/s\}$  we have

$$F_T(t) = L^{-1}\left\{\frac{\Lambda^2 \hat{a}^m}{s(s + \Lambda)[(s + \lambda)(s + \hat{a})^m - (1 - \alpha)\lambda \hat{a}^m]}\right\} \quad (3.10)$$

Thus we have the following algorithm for the resource optimization problem in the model shown in Figure 3.3.

**Algorithm 3-2:**

Steps 1 and 2 are the same as Steps 1 and 2 in Algorithm 3-1 except  $F_T(t)$  given by (3.10).

Note that if we cannot get a solution for the resource optimization problem using Algorithm 3-1 (or 3-2), then the service provider cannot execute the service request for the service model 1 (or 3-2) due to at least one of the following reasons:

1. The service provider has an insufficient resource (i.e.,  $N_j$  is too small).
2. A pre-specific fee is too low (i.e.,  $I > C^D$ ).
3. A network connection is either too slow or has a problem so that (3.1) cannot be satisfied.

Using these information, we may detect and debug either a network problem or a service provider's capacity problem, or the SLA needs to be re-negotiated.

In Algorithms 3-1 and 3-2, the run-time for Step 2 is  $O(m)$ . Thus, the run-time of either Algorithm 3-1 or 3-2 is a sum of  $O(m)$ , the run-time for inverting the LST of the response time, and the run-time for finding the maxima of the resulting function (or  $F(t)$ ). While an efficient approach to finding the maxima of  $F(t)$  can be found in [78], inverting the LST of the response time can be easily done by using the methods presented in [40].

### 3.4 Numerical Validations

In this section we demonstrate the accuracy and applicability of our proposed approximation method.

Table 3.2: The Service Rates of the Eight Stations in Model 1

Service Rates	$\mu_1$	$\mu_2$	$\mu_3$	$\mu_4$	$\mu_5$	$\mu_6$	$\mu_7$	$\mu_8$
Values	52	18	80	35	41	15	25	35

Two types of errors are introduced in our proposed approximation method. The first, hereafter referred to as Class I error, comes from numerically inverting the Laplace transform. The other, hereafter referred to as Class II error, is due to the assumptions that the waiting time of a customer at each station is independent of the waiting times at the other stations, and it is also independent of its waiting times in other visits to the same station.

The relative error % is used to measure the accuracy of the approximate results compared to model simulation results, and it is defined as follows

$$\text{Relative error \%} = \frac{\text{Approximate Result} - \text{Simulation Result}}{\text{Simulation Result}} \times 100 \quad (3.11)$$

We study the accuracy of our proposed approximation method using two examples below.

We first verify the accuracy of our approach for the first service model shown in Figure 3.2. Let  $m = 8$ ,  $\lambda = 100$ ,  $N_j = 100$ ,  $c_j = 2$ ,  $\psi(n_j) = 1.5^{\log_2 n_j}$  and  $C^D = 400$  ( $j = 1, \dots, 8$ ). The service rates of these eight stations are listed in Table 3.2.

We simulated the queueing network using Arena and the analytical method was implemented in Mathematica. The simulation results are considered as “exact” since the simulation model is an exact representation of the queueing network under study.

Table 3.3 shows the simulated and approximate cumulative distribution of the response time. In the table, the column labeled “Simul” gives the simulation result, the column labeled “Approx” gives the approximate result, and the column labeled “R-Err %” gives their relative errors. The same abbreviations are also used in Table 3.7. It appears that the results obtained by Algorithm 1 are very accurate. The optimal number of servers required for 97.5% of the response time to be less than  $T^D = 0.16$  is shown in Table 3.4. The exact optimal number of servers, obtained by exhaustive search using the simulation model, and assuming that each station has the same utilization, or balanced utilization, is consistent with the ones shown in Table 3.4. Thus,  $I = 382 < C^D$ . We point out that the relative errors shown in Table 3.3 are only due to the Class I error since the Class II error is not present for this service model.

Let us now consider an example of the service model 2 shown in Figure 3.3. We choose  $m = 8$ ,  $\Lambda = 100$ ,  $\alpha = 0.67$ ,  $N_j = 250$ ,  $c_j = 1$ ,  $\psi(n_j) = 1.5^{\log_2 n_j}$ , and  $C^D = 580$  ( $j = 1, \dots, 8$ ).

Table 3.3: The Cumulative Distribution of the Response Time in Model 1

Response Time	Simul	Approx	R-Err %
0.04	0.0213	0.0214	0.4393
0.06	0.1517	0.1528	0.7004
0.08	0.4070	0.4075	0.1112
0.10	0.6681	0.6672	-0.1377
0.12	0.8468	0.8450	-0.2158
0.14	0.9398	0.9379	-0.1974
0.16	0.9785	0.9780	-0.0498
0.18	0.9931	0.9929	-0.0157
0.20	0.9979	0.9979	0.0000
0.22	0.9995	0.9994	-0.0077
0.24	0.9999	0.9998	-0.0051
0.26	1.0000	1.0000	0.0000

Table 3.4: The Optimal Number of Servers in Model 1

Station	1	2	3	4	5	6	7	8
#Servers	11	62	5	20	16	84	35	20

The service rates of these eight stations are listed in Table 3.5. Thus it follows from equation  $\lambda = \frac{\Lambda}{\alpha}$  that  $\lambda = 149.25$ .

We obtained the cumulative distribution of the response time by solving (3.10) using the package of the inverse Laplace transforms in Graf [40]. Table 3.6 shows the number of servers in the eight stations necessary to ensure the 95% SLA guarantee for  $T^D \leq 0.6$ . We also simulated the tandem queueing network and validated using the brute-force approach that these numbers of servers obtained by our approximate method are in fact optimal, provided that each station has balanced utilization. The optimal number of servers is given in Table 3.6. It derives that  $I = 560 < C^D$ , i.e., Step 2) in Algorithm 3-2 is met. Table 3.7 gives the cumulative distribution of the response time

Table 3.5: The Service Rates of The Eight Stations in Model 2

Service Rates	$\mu_1$	$\mu_2$	$\mu_3$	$\mu_4$	$\mu_5$	$\mu_6$	$\mu_7$	$\mu_8$
Values	10	45	100	20	32	18	8	85

Table 3.6: The Optimal Number of Servers in Model 2

Station	1	2	3	4	5	6	7	8
#Servers	168	13	4	52	23	62	246	5

obtained using the approximate method and the simulation method, and the relative error %. The relative error comes from both Classes I and II errors. We note that our approximate method has a very good accuracy.

In both two examples as above, the run-time for the approximate method is less than 1 second when Algorithm 3-1 or 3-2 was implemented in Mathematica, and the run-time for the simulation result is less than 1 minute when Arena was used.

Extensive numerical results point to the fact that the independence assumption has little impact on the accuracy of the results when the number of nodes is large. A contributing factor is that typically we are interested in values of the cumulative distribution of the response time that correspond to very high percentiles for which the approximate results seem to have a very good accuracy through a comparison with simulation results that are considered as “exact.”

Additionally, to the best of our knowledge, this is the first work that provides an analytical solution of the resource optimization problem subject to the constraints of a percentile response time and a price. Hence, we do not give a comparison of our proposed method with other methods in this research.

### 3.5 The Balanced Condition

The aforementioned method to solving the resource optimization problem requires a balanced condition on all server stations. From a practical point of view, this condition is a fairly reasonable assumption in the resource optimization problem. This is because each server station may be owned by different entities, such as different organizations or different service providers, who have their individual objectives. These entities would price their services in such a way that they might maximize their profits. When a customer pays a certain fee for the service, these entities have to collaborate in allocating enough resources and determining the prices that they charge the customer. Otherwise, as studied in [44], a self-centered competition for more revenues through either the insufficient allocation of resources or the high pricing may inflate the price charged to

Table 3.7: The Cumulative Distribution of the Response Time in Model 2

Response Time	Simul	Approx	R-Err %
0.20	0.4865	0.4781	-1.7284
0.30	0.7418	0.7267	-2.0336
0.40	0.8551	0.8541	-0.1201
0.50	0.9189	0.9226	0.4067
0.60	0.9538	0.9589	0.5327
0.70	0.9733	0.9782	0.5000
0.80	0.9845	0.9884	0.3967
0.90	0.9908	0.9938	0.3070
1.00	0.9946	0.9967	0.2136
1.10	0.9967	0.9983	0.1563
1.20	0.9980	0.9991	0.1079
1.30	0.9988	0.9995	0.0714
1.40	0.9997	0.9997	0.0000
1.60	0.9999	0.9999	0.0026
1.80	0.9999	1.0000	0.0079
2.00	1.0000	1.0000	0.0000

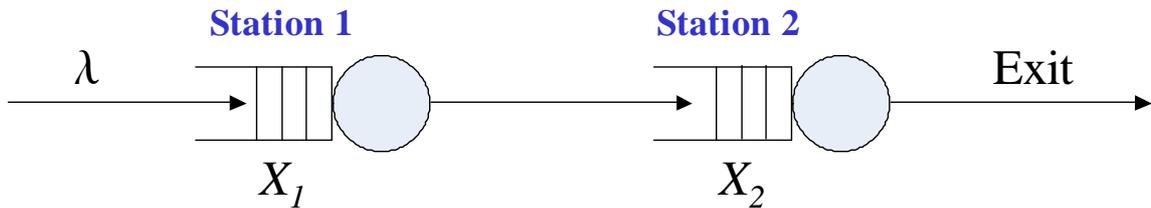


Figure 3.4: A Two Station Tandem Model

a customer and reduce the potential demand for the service. This means that these entities should collaborate with each other in allocating sufficient resources. Therefore, imposing the balanced condition on all server stations is one of good ways such that these entities are evenly utilized.

In this section, we specifically study the balanced condition purely from a mathematical point of view. We want to investigate the change of the obtained solution, i.e., the number of servers in each station, when the condition is not imposed. This section restricts our discussion to the case of two stations in tandem without an infinite server in order to reduce the amount of unnecessary formalization which would lead to tedious computations giving no new insight into the study. The tandem is shown in Figure 3.4. A similar extension to the case of  $m$  stations in tandem can be

obtained for  $m > 2$ .

The LST of the response time is computed by

$$L_T(s) = \prod_{j=1}^2 \frac{\hat{\mu}_j(1 - \rho_j)}{s + \hat{\mu}_j(1 - \rho_j)}$$

where  $\hat{\mu}_j = \psi(n_j)\mu_j$  and  $\rho_j = \frac{\lambda_j}{\hat{\mu}_j}$  for  $j = 1, 2$ . It thus derives the PDF of the response time given by

$$f_T(t) = \frac{\alpha_1 \alpha_2}{\alpha_2 - \alpha_1} (e^{-\alpha_1 t} - e^{-\alpha_2 t})$$

where  $\alpha_j = \hat{\mu}_j(1 - \rho_j) = \hat{\mu}_j - \lambda_j$  for  $j = 1, 2$ .

Furthermore, the percentile response time given in (3.1) can be rewritten as

$$G(T^D) = \int_{T^D}^{+\infty} f_T(t) dt = \frac{1}{\alpha_2 - \alpha_1} (\alpha_2 e^{-\alpha_1 T^D} - \alpha_1 e^{-\alpha_2 T^D}) \leq 1 - \gamma\%.$$

Again, in order to reduce the amount of unnecessary formalization which would lead to tedious computations giving no new insight into the study, we reformulate the resource optimization problem given in Section 3.2 as

$$\min_{\hat{\mu}_1, \hat{\mu}_2 \geq 0} (\hat{c}_1 \hat{\mu}_1 + \hat{c}_2 \hat{\mu}_2) \quad (3.12)$$

where  $\hat{c}_j$  is the cost at a per unit of service rate in station  $j$  for  $j = 1, 2$ , under the constraint:

$$G(T^D) = \frac{1}{\alpha_2 - \alpha_1} (\alpha_2 e^{-\alpha_1 T^D} - \alpha_1 e^{-\alpha_2 T^D}) \leq 1 - \gamma\%$$

which is equivalent to the constraint:

$$G(T^D) = \frac{1}{\hat{\mu}_2 - \hat{\mu}_1} \left[ (\hat{\mu}_2 - \lambda_2) e^{-(\hat{\mu}_1 - \lambda_1) T^D} - (\hat{\mu}_1 - \lambda_1) e^{-(\hat{\mu}_2 - \lambda_2) T^D} \right] \leq 1 - \gamma\% \quad (3.13)$$

due to  $\alpha_j = \hat{\mu}_j - \lambda_j$ , where  $\alpha_j \geq 0$ . Note that by using the traffic equation in queueing theory we have  $\lambda_1 = \lambda_2$ .

Next, when imposing the balanced condition, i.e.,  $\alpha_1 = \alpha_2 \triangleq \alpha$ , we have that

$$L_T(s) = \frac{\alpha^2}{(s + \alpha)^2}$$

Thus, the PDF of response time is given by

$$p(t) = \alpha^2 t e^{-\alpha t}$$

Since  $\alpha_j = \hat{\mu}_j - \lambda_j$  and  $\lambda_1 = \lambda_2 \triangleq \lambda$ , we obtain  $\hat{\mu}_1 = \hat{\mu}_2 \triangleq \hat{\mu}$ . Subsequently, the objective function in the resource optimization problem is written as

$$c_1 \hat{\mu}_1 + c_2 \hat{\mu}_2 = (c_1 + c_2) \hat{\mu} \quad (3.14)$$

and the percentile response time constraint can be computed by

$$H(T^D) = \int_{T^D}^{+\infty} p(t) dt = [(\hat{\mu} - \lambda) + 1] e^{-(\hat{\mu} - \lambda)T^D} \leq 1 - \gamma\% \quad (3.15)$$

Note that  $G(T^D)$  is a function of variables  $\hat{\mu}_1$  and  $\hat{\mu}_2$ , and  $H(T^D)$  is a function of variable  $\hat{\mu}$ . To simplify notation, we reuse the notation  $G$  and  $H$ , and write them as  $G(\hat{\mu}_1, \hat{\mu}_2)$  and  $H(\hat{\mu})$ , i.e.,  $G(\hat{\mu}_1, \hat{\mu}_2) = G(T^D)$  and  $H(\hat{\mu}) = H(T^D)$ . Moreover, it is easy to see that

$$\frac{d(H(\hat{\mu}))}{d\hat{\mu}} = -(T^D)^2(\hat{\mu} - \lambda) e^{-(\hat{\mu} - \lambda)T^D} \leq 0 \quad (3.16)$$

due to  $\hat{\mu} \geq \lambda$ , which implies that  $H(\hat{\mu})$  is a decreasing function of variable  $\hat{\mu}$ . By considering (3.12)-(3.16) we can rewrite the resource optimization problem as follows.

When the balanced condition is not imposed, the resource optimization problem is to find  $\hat{\mu}_1$  and  $\hat{\mu}_2$  in the minimization problem:

$$\begin{aligned} & \min_{\hat{\mu}_1, \hat{\mu}_2} (\hat{c}_1 \hat{\mu}_1 + \hat{c}_2 \hat{\mu}_2) \\ \textbf{(Problem I)} \quad & \text{subject to} \\ & G(\hat{\mu}_1, \hat{\mu}_2) \leq 1 - \gamma\%, \text{ and } \hat{\mu}_j \geq \lambda_j, \quad j = 1, 2 \end{aligned}$$

where  $G(\hat{\mu}_1, \hat{\mu}_2) = G(T^D)$  is given in (3.13).

When the balanced condition is imposed, the resource optimization problem is to find  $\hat{\mu} = \hat{\mu}_1 = \hat{\mu}_2$  in the maximization problem:

$$\begin{aligned} & \arg \max_{\hat{\mu}} H(\hat{\mu}) \\ \textbf{(Problem II)} \quad & \text{subject to} \\ & H(\hat{\mu}) \leq 1 - \gamma\%, \text{ and } \hat{\mu} \geq \lambda \end{aligned}$$

where  $H(\hat{\mu}) = H(T^D)$  is given in (3.15).

We are interested in investigating the difference between the solutions of Problems I and II. Let first study the property of these solutions. Assume that  $\hat{\mu}_1 = \hat{\mu}_1^*$  and  $\hat{\mu}_2 = \hat{\mu}_2^*$  be the solution of Problem I, denoted by

$$(\hat{\mu}_1^*, \hat{\mu}_2^*) \leftarrow \textbf{(Problem I)}$$

Clearly, the constraint function  $G(\hat{\mu}_1, \hat{\mu}_2) - (1 - \gamma\%)$  is a symmetric function with respect to the line  $\hat{\mu}_1 = \hat{\mu}_2$ . Furthermore, if  $\hat{c}_1 = \hat{c}_2$ , then the objective function  $\hat{c}_1 * \hat{\mu}_1 + \hat{c}_2 * \hat{\mu}_2$  is a symmetric function with respect to the line  $\hat{\mu}_1 = \hat{\mu}_2$  as well, which implies that  $\hat{\mu}_1 = \hat{\mu}_2^*$  and  $\hat{\mu}_2 = \hat{\mu}_1^*$  is the

solution of Problem I as well. Hence, when there exists a unique solution of Problem I,  $\hat{\mu}_1^*$  should be equal to  $\hat{\mu}_2^*$ . This derives that Problem I has the same solution as Problem II. We summarize the conclusion below.

**Proposition 3.5.1** *Suppose that  $\hat{c}_1 = \hat{c}_2$  and  $(\hat{\mu}_1^*, \hat{\mu}_2^*)$  is a unique solution of Problem I. Then,  $\hat{\mu}_1^* = \hat{\mu}_2^*$ , that is, Problem I has the same solution as Problem II.*

This means that the balanced condition is satisfied when the two entities representing two resource stations charge the same rate at a per unit of service rate. Hence, the balanced condition is a fairly reasonable assumption as well from the mathematical point of view.

Note that  $H(\hat{\mu})$  is a continuous and non-increasing function. Hence, Problem II is equivalent to the problem in finding a solution of  $H(\hat{\mu}) = 1 - \gamma\%$ . Thus, in this case, the solution of Problem II can be easily obtained. In general, when the assumption of Proposition 3.5.1 does not hold, Problem II can be solved numerically.

Next, we present numerical implementations of these problems. Let us first assume that  $\hat{c}_1 = \hat{c}_2$ . We choose  $\lambda = 100$  and  $T^D = 3.5$  with varied  $\gamma$ . Table 3.8 gives the solutions of Problems I and II. Clearly, both optimization problems have the same solution, i.e.,  $\hat{\mu}_1 = \hat{\mu}_2 = \hat{\mu}$ ,

Table 3.8: The Solutions of Problems I and II with Varied  $\gamma$

$\gamma$	80	85	90	95	97
$\hat{\mu}_1$ by Problem I	100.856	100.964	101.111	101.355	101.530
$\hat{\mu}_2$ by Problem I	100.856	100.964	101.111	101.355	101.530
$\hat{\mu}$ by Problem II	100.856	100.964	101.111	101.355	101.530

which confirms our analysis as presented in Proposition 3.5.1.

We further consider the case of  $\hat{c}_1 \neq \hat{c}_2$ . Let us choose  $\lambda = 100$ ,  $T^D = 3.5$  and  $\gamma = 95$  with a fixed  $\hat{c}_1 = 1$  and varied  $\hat{c}_2$ . We also choose  $\mu_1 = 10$ ,  $\mu_2 = 20$  and  $\psi_j(n_j) = 1.5^{\log_2 n_j}$  where  $j = 1, 2$ . Table 3.9 showed the solutions of Problems I and II. Note that the solution of Problem II does not change with varied  $\hat{c}_2$  because Problem II is not related to  $\hat{c}_2$ . Table 3.9 demonstrates that the differences among  $\hat{\mu}$ ,  $\hat{\mu}_1$  and  $\hat{\mu}_2$  are small. Furthermore, by using  $\mu_1 = 10$ ,  $\mu_2 = 20$  and  $\psi_j(n_j) = 1.5^{\log_2 n_j}$  where  $j = 1, 2$ , we found the number of servers required to ensure that 95% of the service requests from a customer can be received in less than  $T^D = 3.5$ , as shown in Table 3.10. As is seen in Table 3.10 there is only at most one server difference in the number of servers

Table 3.9: The Solutions of Problems I and II with Varied  $\hat{c}_2$ 

$\hat{c}_2$	1	2	4	8	16
$\hat{\mu}_1$ by Problem I	101.355	101.572	101.868	102.27	102.825
$\hat{\mu}_2$ by Problem I	101.355	101.202	101.097	101.026	100.977
$\hat{\mu}_1 - \hat{\mu}_2$	0.000	0.370	0.771	1.244	1.848
$\hat{\mu}$ by Problem II	101.355	101.355	101.355	101.355	101.355
$\hat{\mu} - \hat{\mu}_1$	0.000	-0.217	-0.513	-0.915	-1.470
$\hat{\mu} - \hat{\mu}_2$	0.000	0.153	0.258	0.329	0.378

Table 3.10: A Comparison of the Number of Servers by Using Problems I and II

$\hat{c}_2$	1	2	4	8	16
$n_1$ by Problem I	53	53	53	54	54
$n_1$ by Problem II	53	53	53	53	53
$n_2$ by Problem I	17	16	16	16	16
$n_2$ by Problem II	17	17	17	17	17

obtained by Problems I and II. Again, according to Proposition 3.5.1 and the above numerical implementations, we see that the balanced condition does not introduce a significant inaccuracy. We have also done a lot of simulations with varied  $\lambda$  that are not reported here. These simulation results have shown that solving Problem II can provide an accurate solution for Problem I.

In addition, the balanced condition seems to be a fairly reasonable assumption from the practical point of view.

**Remark.** We can similarly show that Proposition 3.5.1 holds in the case of  $m$  stations in tandem with or without feedback when the cost of each station at a per unit of service rate is the same.

### 3.6 An Application: Web Services Performance Modeling and Analysis

As described early, the proposed approach can be extended to a general queueing network in which its nodes are arbitrarily linked as long as they can be quantified. In this section we study

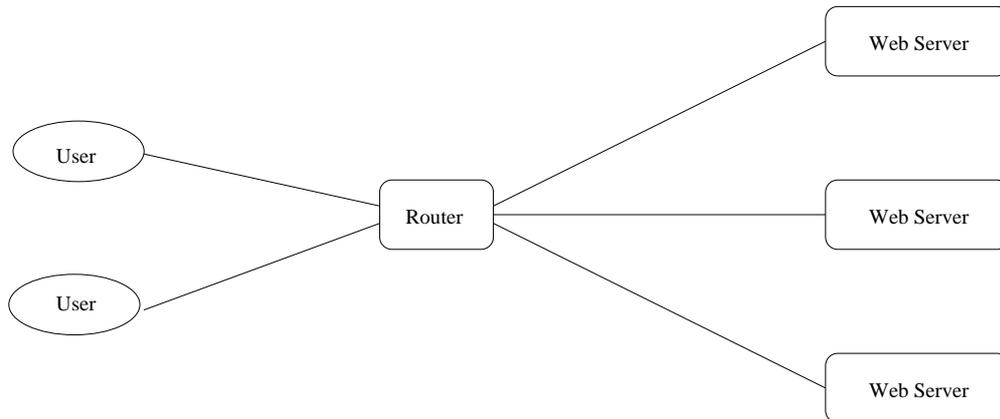


Figure 3.5: A Scenario in Web Services Applications

a two-station queueing network that describes a Web services application model. We demonstrate how to apply our proposed approach to analyzing the Web services performance. This section is based on the results in Xiong [92].

In this section, we first introduce a typical Web services model and then study the Web services performance by using our aforementioned approach for the calculation of the percentile response time. Moreover, a numerical validation is given in Section 3.6.3 that demonstrates the accuracy of this method. The obtained result can be used in the design of a computing system for Web services applications subject to the percentile response time guarantee.

### 3.6.1 A Web Services Performance Model

Understanding the characteristics of Web service performance has become critical for e-business applications. For the commercial success of these e-business services, the ability to deliver QoS guarantees is crucial. One often considers the percentile of response time as the QoS metric. In the section we consider a typical Web services performance model. The model consists of a collection of Web services resources - Web servers - used by a service provider to host e-business applications for business customers, as shown in Figure 3.5. A service request sent by a user is transmitted to Web servers that are owned by the service provider over network [58]. An e-business application running in such a computing environment is associated with a SLA.

A customer represents a business that may consist of multiple users and requests services at a given rate to be processed by the service provider's resources: the Web servers, according to

QoS requirements and for a given fee. One of our goals is to analyze Web services performance in term of a percentile response time, which is different from most existing work addressing QoS requirements in Web services performance in which the average response time (or average execution time) are usually used, for example, si.e., the time to execute a service request, is less than a pre-defined value.

Our second goal is to study the independent assumption required for the queueing network with feedback given in Figure 3.3. Similarly, the two-station queueing network shown in Figure 3.5 is not overtake-free. Hence, in order to obtain the probability distribution function of the response time we have to assume that the waiting time of a customer at a station is independent of its waiting time at another station, and each visit at the same station  $j$  is independent of the others. The two-station queueing network is selected. It is because a closed-form solution of the probability distribution function can be obtained under the independent assumption. In this case, the numerical computation of an inverse Laplace transform is not needed so that there is no computational error coming from numerically inverting the Laplace transform.

In next section, we derive an approximation method for the calculation of the probability and cumulative distributions of the response time, and then show the accuracy of the proposed approximation method.

### 3.6.2 Analysis of the Web Services Performance Model

In this section, we derive the percentile of the response time in a Web services performance model consisting of a router and a Web server station for Web services applications.

#### The Response Time Distribution

The Web services performance model, shown in Figure 3.6, consists of a router and a Web server station. Both the router and the Web server station are modeled as a single  $M/M/1$  queue. The service discipline is FIFO. External arrivals to the computing station are Poisson distributed with a rate  $\lambda$  and service times are exponentially distributed. Upon completion of a service at the Web server station, the customer exits the system with probability  $1 - \beta$ , or continues to be served at the router with probability  $\beta$ . Furthermore, after a customer is processed at the router, it returns to the beginning of the system with probability  $1 - \alpha$ , or it exits the computing system with probability

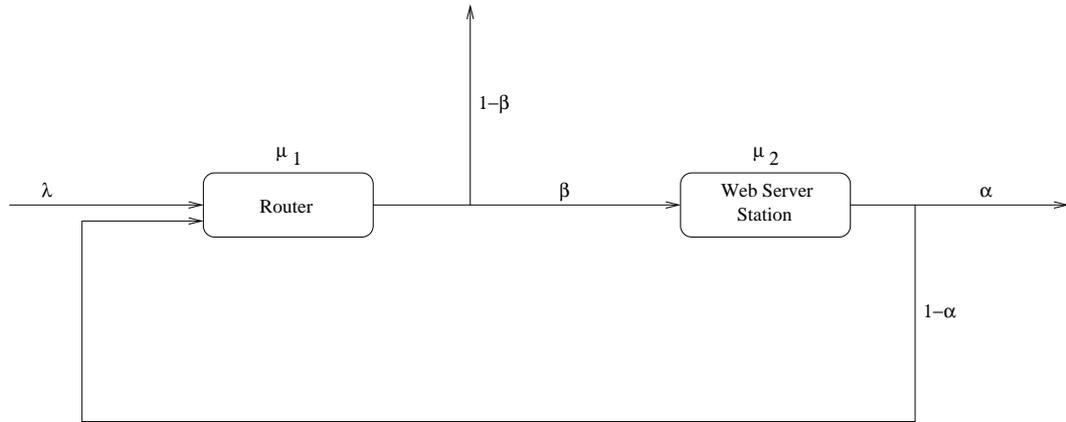


Figure 3.6: A Web Server Performance Model

$\alpha$ .

Let  $(i, j)$  be the number of visits in the router and the Web server station where  $i$  and  $j$  are the number of visits in the router and the number of visits in the Web server station respectively. Let  $p(i, j)$  be the probability of  $i$  visits to the router and  $j$  visits to the Web server station. There may be one time visit difference between the router and the Web server station. This means that either  $j = i$ , or  $j = i - 1$ . Let  $T_1$  (or  $T_2$ ) be the waiting time, that is, the time from the moment a customer arrives at the router (or the Web server station) to the moment it leaves the router (or the Web server station). We assume that  $T_1$  and  $T_2$  are the same for each visit. Then, we have

# Visits	Response Time	Probability
(1, 0)	$T_1$	$p(1, 0) = 1 - \beta$
(1, 1)	$T_1 + T_2$	$p(1, 1) = \beta\alpha$
(2, 1)	$(T_1 + T_2) + T_1$	$p(2, 1) = \beta(1 - \alpha)(1 - \beta)$
(2, 2)	$2(T_1 + T_2)$	$p(2, 2) = \beta^2(1 - \alpha)\alpha$
...	...	...

Therefore, the average response time  $E(T)$  is a weighted sum of individual response times

$$\begin{aligned}
 E[T] = & \sum_{j=1}^{\infty} \beta^{j-1} (1 - \alpha)^{j-1} (1 - \beta) [(j - 1)(T_1 + T_2) \\
 & + T_1] + \sum_{j=1}^{\infty} \beta^j (1 - \alpha)^{j-1} \alpha [j(T_1 + T_2)] \quad (3.17)
 \end{aligned}$$

Note that in equation (3.17), its coefficients maintain the following identity:

$$\sum_{j=1}^{\infty} [\beta^{j-1} (1 - \alpha)^{j-1} (1 - \beta) + \beta^j (1 - \alpha)^{j-1} \alpha] = 1$$

Now, let  $T_1^i$  and  $T_2^i$  be the response time at the router and the Web server station for the  $i$ -th visit respectively. Then we have the following expression for  $E(T)$

$$\begin{aligned}
E[T] &= \sum_{j=1}^{\infty} \beta^{j-1} (1-\alpha)^{j-1} (1-\beta) \left[ \sum_{i=1}^{j-1} (T_1^i + T_2^i) \right. \\
&\quad \left. + T_1^j \right] + \sum_{j=1}^{\infty} \beta^j (1-\alpha)^{j-1} \alpha \left[ \sum_{i=1}^j (T_1^i + T_2^i) \right]
\end{aligned} \tag{3.18}$$

Let  $R(j, j-1)$  be the response time for  $j$  visits to  $T_1$  and  $j-1$  visits to  $T_2$ , and  $R(j, j)$  be the response time for  $j$  visits to  $T_1$  and  $T_2$ . Then

$$\begin{aligned}
R(j, j-1) &= \sum_{i=1}^{j-1} (T_1^i + T_2^i) + T_1^j \\
R(j, j) &= \sum_{i=1}^j (T_1^i + T_2^i)
\end{aligned} \tag{3.19}$$

Assume that  $T_1^i$  and  $T_1^k$ ,  $T_2^i$  and  $T_2^k$ , and  $T_1^i$  and  $T_2^i$  are mutually independent ( $i \neq k$ ). Under these assumptions, the conditional Laplace-Stieltjes transforms (LSTs) of the response times  $R(j, j-1)$  and  $R(j, j)$  can be expressed as follows.

$$\begin{aligned}
L_{R(j, j-1)}(s) &= (L_{T_1}(s))^j \times (L_{T_2}(s))^{j-1} \\
L_{R(j, j)}(s) &= (L_{T_1}(s))^j \times (L_{T_2}(s))^j
\end{aligned} \tag{3.20}$$

Combining equations (3.18) and (3.20), the LST of the total response time is

$$\begin{aligned}
L_T(s) &= \sum_{j=1}^{\infty} \beta^{j-1} (1-\alpha)^{j-1} (1-\beta) L_{T_1}^j(s) L_{T_2}^{j-1}(s) \\
&\quad + \sum_{j=1}^{\infty} \beta^j (1-\alpha)^{j-1} \alpha L_{T_1}^j(s) L_{T_2}^j(s)
\end{aligned}$$

That is,

$$L_T(s) = \frac{(1-\beta)L_{T_1}(s) + \alpha\beta L_{T_1}(s)L_{T_2}(s)}{1 - \beta(1-\alpha)L_{T_1}(s)L_{T_2}(s)} \tag{3.21}$$

Let us now remind that the router and the Web server station are each modeled as an  $M/M/1$  queue. Let  $\lambda_i$  and  $\mu_i$  ( $i = 1, 2$ ) be the arrival rates and service rates at the first and second  $M/M/1$  queues respectively. Then, traffic equations are given by

$$\begin{cases} \lambda_1 &= \lambda + (1-\alpha)\lambda_2 \\ \lambda_2 &= \beta\lambda_1 \end{cases}$$

Thus, the following expressions of  $\lambda_1$  and  $\lambda_2$  can be derived from these local balance equations:

$$\begin{cases} \lambda_1 = \frac{\lambda}{1-(1-\alpha)\beta} \\ \lambda_2 = \frac{\lambda\beta}{1-(1-\alpha)\beta} \end{cases}$$

Moreover, the LSTs of the response time at each queue is

$$L_{T_1}(s) = \frac{a_1}{s + a_1},$$

where

$$a_1 = \mu_1(1 - \rho_1), \quad \rho_1 = \frac{\lambda_1}{\mu_1},$$

and

$$L_{T_2}(s) = \frac{a_2}{s + a_2},$$

where

$$a_2 = \mu_2(1 - \rho_2), \quad \rho_2 = \frac{\lambda_2}{\mu_2}.$$

It follows from (3.21) that

$$\begin{aligned} L_T(s) &= \frac{(1 - \beta)\frac{1}{L_{T_2}(s)} + \alpha\beta}{\frac{1}{L_{T_1}(s)L_{T_2}(s)} - \beta(1 - \alpha)} \\ &= \frac{a_1(1 - \beta)(s + a_2) + a_1a_2\alpha\beta}{(s + a_1)(s + a_2) - a_1a_2\beta(1 - \alpha)} \end{aligned} \quad (3.22)$$

We observe that the de-numerator is a quadratic polynomial in  $s$ :

$$\begin{aligned} (s + a_1)(s + a_2) - a_1a_2\beta(1 - \alpha) &= \\ s^2 + (a_1 + a_2)s + a_1a_2(1 - \beta + \alpha\beta) & \end{aligned}$$

This polynomial has the roots

$$s_{1,2} = \frac{-(a_1 + a_2) \pm \sqrt{(a_1 + a_2)^2 - 4a_1a_2(1 - \beta + \alpha\beta)}}{2} \quad (3.23)$$

Notice that both  $\alpha$  and  $\beta$  range from 0 to 1. Hence,  $1 - \beta + \alpha\beta$  is non-negative. This means that  $s_1$  and  $s_2$  must be non-positive, and either  $s_1$  or  $s_2$  is zero if and only if

$$1 - \beta + \alpha\beta = 1 - \beta(1 - \alpha) = 0$$

i.e.,  $\alpha = 0$  and  $\beta = 1$ , which is a less interesting case.

Moreover, from (3.22) we can have that

$$\begin{aligned} L_T(s) &= \frac{a_1(1-\beta)s + a_1a_2(1-\beta+\alpha\beta)}{(s-s_1)(s-s_2)} \\ &\triangleq \frac{B_1}{s-s_1} + \frac{B_2}{s-s_2} \end{aligned} \quad (3.24)$$

where constants  $B_1$  and  $B_2$  are determined by

$$\begin{cases} B_1 + B_2 = a_1(1-\beta) \\ S_1B_2 + s_2B_1 = -a_1a_2(1-\beta+\alpha\beta) \end{cases}$$

That is, constants  $B_1$  and  $B_2$  are given by

$$\begin{cases} B_1 = \frac{a_1s_1(1-\beta)+a_1a_2(1-\beta+\alpha\beta)}{s_1-s_2} \\ B_2 = -\frac{a_1s_2(1-\beta)+a_1a_2(1-\beta+\alpha\beta)}{s_1-s_2} \end{cases} \quad (3.25)$$

Therefore it follows from (3.24) that the probability distribution function of response time  $T$  is

$$f_T(t) = B_1 e^{s_1 t} + B_2 e^{s_2 t} \quad (3.26)$$

To ensure that  $\gamma\%$  of the response times for customer service requests are not more than a desired target response time  $T^D$ , we require that

$$\begin{aligned} G(T^D) &= \int_{T^D}^{\infty} f_T(t) dt \\ &= -\left( \frac{B_1}{s_1} e^{s_1 T^D} + \frac{B_2}{s_2} e^{s_2 T^D} \right) \\ &\leq 1 - \gamma\% \end{aligned} \quad (3.27)$$

From equation (3.27), we determine the service level ( $= \gamma\%$ ) for a given arrival rate and a given service rate, or find a service rate necessary to ensure a certain service level as in (3.2) and (3.3).

In this section, we consider a Web services performance model consisting of a router and a Web server station. Our study for the Web services performance as above can be easily extended to complex Web services scenarios where the number of stations is more than 2.

### 3.6.3 A Numerical Validation

In this section we demonstrate the accuracy and applicability of our proposed approximation method. We study the accuracy of our proposed approximation method using an example below.

Table 3.11: The Cumulative Distribution Functions of The Response Time in the Web Services Performance Model

$T^D$	$\alpha=0.65, \beta=0.4$			$\alpha=0.65, \beta=0.2$			$\alpha=0.5, \beta=0.2$		
	Simul	Approx	R-Err%	Simul	Approx	R-Err%	Simul	Approx	R-Err%
0.005	0.5179	0.5166	-0.26	0.6395	0.6407	0.19	0.6291	0.6303	0.20
0.010	0.7478	0.7457	-0.28	0.8566	0.8556	-0.11	0.8436	0.8441	0.06
0.015	0.8624	0.8610	-0.17	0.9379	0.9375	-0.04	0.9278	0.9287	0.09
0.020	0.9232	0.9227	-0.05	0.9714	0.9718	0.04	0.9652	0.9659	0.08
0.025	0.9569	0.9568	-0.01	0.9863	0.9870	0.07	0.9824	0.9834	0.10
0.030	0.9752	0.9758	0.06	0.9936	0.9939	0.03	0.9910	0.9918	0.09
0.035	0.9854	0.9864	0.10	0.9968	0.9972	0.04	0.9953	0.9960	0.07
0.040	0.9914	0.9924	0.10	0.9983	0.9987	0.04	0.9975	0.9980	0.05
0.045	0.9950	0.9957	0.07	0.9991	0.9994	0.03	0.9985	0.9990	0.05
0.050	0.9969	0.9976	0.07	0.9994	0.9997	0.03	0.9992	0.9995	0.03
0.055	0.9981	0.9986	0.05	0.9996	0.9999	0.03	0.9994	0.9998	0.04
0.060	0.9989	0.9992	0.03	0.9998	0.9999	0.01	0.9996	0.9999	0.03
0.065	1.0000	0.9996	-0.04	0.9999	1.0000	0.01	0.9997	0.9999	0.02
0.070	1.0000	0.9998	-0.02	0.9999	1.0000	0.01	0.9998	1.0000	0.02
0.075	1.0000	0.9999	-0.01	1.0000	1.0000	0.00	0.9999	1.0000	0.01
0.080	1.0000	0.9999	-0.01	1.0000	1.0000	0.00	0.9999	1.0000	0.01
0.085	1.0000	1.0000	0.00	1.0000	1.0000	0.00	1.0000	1.0000	0.00

We shall verify the accuracy of the approximate method for the two-station system analyzed in Section 3.6.2. We let  $\lambda = 100$ ,  $\mu_1 = 380$ ,  $\mu_2 = 200$ , and  $\alpha$  and  $\beta$  were varied.

We simulated the queueing network using Arena; see [5], and the analytical method was implemented in Matlab and also in Mathematica. The simulation model in Arena exactly represents the Web services performance model under study, so the simulation results in Arena are considered as “exact”.

Table 3.11 presented on next page shows the simulated and approximate cumulative distribution function of the response time for different values of  $\alpha$  and  $\beta$ . In the table, columns labeled “Simul” gives the simulation result, columns labeled “Approx” gives the approximate result, and columns labeled “R-Err %” gives their relative errors. Those abbreviations are also used in other tables of this section. It appears that the results obtained by our approximate method are fairly good. For example, when  $\alpha = 0.5$ ,  $\beta = 0.2$ , both results reflects that 98% of time the customer requests will be responded in less than  $T^D = 0.025$  as shown in Table 3.11.

Additionally, in Web services applications subject to a SLA, customers are typically inter-

ested in a high percentile of the response time for which our approximation method seems to have a very good accuracy as shown in Table 3.11.

In this approximation method, we assume that the waiting time of a customer at the router (or the Web server station) is independent of the waiting times at the Web server station (or the router), and it is also independent of its waiting times in other visits to the same router (or the same Web server station). As discussed early, there is no error coming from numerically inverting the Laplace transform. Hence, the relative error as shown in Table 3.11 is only due to the above assumption. This means that our approximation method can provide a very good accuracy for the two-station queueing network even if the independent assumption is imposed.

The section has provided an example of solving a non-overtaken free queueing network using our approximate method.

In summary, we have considered Web services performance modeling and analysis in this section. The main difficulty for understanding the characteristics of Web services performance is in the computation of the probability distribution function of a response time. In this section, we have first proposed a typical Web services performance model, and then developed an approximation method that used the Laplace transform of the response time distribution at the router (or the Web server station). Under the assumptions that the studied computing system is analyzed by considering the router (the Web server station) as being separate from the other the Web server station (or the router), and the waiting of a customer at the router (the Web server station) is independent of its waiting times in other visits to the same router (the same Web server station), the Laplace transform of the total response time was computed.

Furthermore, we have conducted numerical experiments to investigate the above assumptions. Numerical results showed the the proposed approximation method provided a good accuracy for the calculation of probability and cumulative distributions of the response time. Hence, the proposed method provides an efficient and accurate solution for the calculation of probability and cumulative distributions of a customer's response time. It can be used in Web services performance prediction.

### **3.7 Concluding Remarks**

We proposed an approach for resource optimization in a service provider's computing environment, whereby we minimize the total cost of computer resources allocated to a customer so

that satisfies a given percentile of the response time. We have formulated the resource optimization problem as an optimization subject to SLA constraints for a service model with or without feedback. In the model without feedback, the obtained LST of a customer's response time is exact, and in the case of the model with feedback, it is approximate. We also developed an efficient and accurate numerical solution for inverting the LST of a customer's response time numerically. Validation tests showed that our approach has a very good accuracy.

In this research we assumed that service requests are served in a queue in a FIFO manner. Priority service disciplines will be discussed in Chapter 4.

## Chapter 4

# Multiple Class Customers

This chapter considers a set of computer resources used by a service provider to host enterprise applications subject to a service level agreement for differentiated customer services under a *preemptive-resume* priority. We present an approach for the resource optimization that minimizes the total cost of computer resources required while preserving a given percentile of the response time for priority-class customers. We first analyze an open tandem queueing network, and then we extend our work to an open tandem queueing network with feedback. We only consider two priority customers in this chapter. High-priority class customers are indexed 1 and low-priority class customers 2. The obtained results in this chapter can be easily extended to the case of multiple priority customers by using a class-based decomposition. (This extension is not considered in this dissertation.) The chapter is mainly based on the results obtained in Xiong and Perros [94].

The chapter is organized as follows. In Section 4.1 we define the SLA performance metric considered in this dissertation. Section 4.2 formulates the resource optimization problem for differential customer services. In Section 4.3, we first give the probability distribution of the response time distribution for a single priority queue with preemptive-resume. Then, we present two typical real-life models and propose an approach for solving the optimization problem for two priority-class customers. In Section 4.4, numerical simulations demonstrate the applicability and validity of the proposed approach. Finally, the conclusions are given in Section 4.5.

## 4.1 The SLA Performance Metric in the Case of Multiple Class Customers

Let us recall that an SLA is a contract between a customer and a service provider that defines all aspects of the service that is to be provided. In this chapter the SLA consists of service performance and a fee under multiple customer services. We consider the percentile of the response time as the performance metric. This is the time it takes for a service request to be executed on the service provider's multiple resource stations.

Assume that  $f_T(t)$  is the probability distribution function of a response time  $T$  of a certain priority class. For example, in the case of two priority classes, for the high-priority class  $T = T^{(1)}$  and for the low-priority class  $T = T^{(2)}$ .  $T_D^{(r)}$  is a desired target response time for priority class  $r$  ( $r = 1, 2$ ) that a customer requests and agrees upon with its service provider based on a fee paid by the customer. The SLA performance metric used in this chapter can be expressed as follow:

$$\int_0^{T_D^{(r)}} f_{T^{(r)}}(t) dt \geq \gamma^{(r)}\%, \quad (r = 1, 2) \quad (4.1)$$

That is,  $\gamma^{(r)}\%$  of the time a customer will receive its service in less than  $T_D^{(r)}$  ( $r = 1, 2$ ).

As an example let us consider an  $M/M/1$  queue with an arrival rate  $\lambda^{(r)}$  and a service rate  $\mu^{(r)}$  ( $r = 1, 2$ ). The service discipline is preemptive-resume. We want to describe the SLA performance metric (4.1) for the high-priority class. As discussed in section 4.1, in this case, the steady-state probability of the system as far the high-priority class is concerned is  $p_0 = 1 - \rho^{(1)}$ , and  $p_k = (1 - \rho^{(1)})(\rho^{(1)})^k$ ,  $k > 0$ , where  $\rho^{(1)} = \frac{\lambda^{(1)}}{\mu^{(1)}}$  (see [76]). The response time  $T^{(1)}$  is exponentially distributed with the parameter  $\mu^{(1)}(1 - \rho^{(1)})$ , i.e., the probability distribution of the high-priority response time is given by

$$f_{T^{(1)}}(t) = \mu^{(1)}(1 - \rho^{(1)})e^{-\mu^{(1)}(1 - \rho^{(1)})t} \quad (4.2)$$

Using the definition given in (4.1), we have that

$$\int_0^{T_D^{(1)}} f_{T^{(1)}}(t) dt = 1 - e^{-\mu^{(1)}(1 - \rho^{(1)})T_D} \geq \gamma^{(1)}\% \quad (4.3)$$

or  $\mu^{(1)} \geq \frac{-\ln(1 - \gamma^{(1)}\%)}{T_D^{(1)}} + \lambda^{(1)}$ . This means that in order to guarantee a higher SLA service level,  $\mu^{(1)}$  increases when  $T_D^{(1)}$  decreases. Similarly, for any given arrival rate  $\lambda^{(1)}$  and service rate  $\mu^{(1)}$ , we can use (4.3) to find the percentile of  $\gamma^{(1)}$ .

## 4.2 The Resource Optimization Problem for Multiple Customer Services

The computer resource optimization problem for multiple customer services can be formulated as the following optimization problem.

### Resource Optimization Problem for Multiple Class Customers:

Find integers  $n_j$  ( $1 \leq n_j \leq N_j$ ;  $j = 1, 2, \dots, m$ ) in the  $n$ -dimensional integer optimization problem (1.1) under the constraints of percentile response times for differential customer services as expressed by (4.1), and the constraint:  $I \leq C_D$ , where  $C_D$  is a fee negotiated and agreed upon between a customer and the service provider.

## 4.3 Approaches for Resource Optimization

In this section, we propose an approach to solving the resource optimization problem for two typical service models that depict the path that service requests have to follow through the service provider's resource stations. Before presenting the approach, we need to derive the Laplace-Stieltjes transforms (LST) of the response time distributions for priority-class customers.

### 4.3.1 The LSTs of Response Time Distributions for Two Priority Customers

Let us recall that high-priority class customers are indexed 1 and low-priority class customers 2. In this section, we assume that the arrival processes of the two classes are independent of each other. Let  $B^{(r)}(t)$  be the service time cumulative distribution of class  $r$  with mean  $1/v^{(r)}$  and second moment  $1/v_2^{(r)}$  ( $r = 1, 2$ ). The total service time distribution  $B(t)$  is given by  $B(t) = \frac{\lambda^{(1)}}{\lambda^{(1)} + \lambda^{(2)}} B^{(1)}(t) + \frac{\lambda^{(2)}}{\lambda^{(1)} + \lambda^{(2)}} B^{(2)}(t)$ , and the arrival rate into the queue is  $\lambda = \lambda^{(1)} + \lambda^{(2)}$ . It follows that the total utilization  $\rho = \rho^{(1)} + \rho^{(2)}$ , where  $\rho^{(r)} = \frac{\lambda^{(r)}}{v^{(r)}}$ , equals to the occupation time given by  $\lambda \int_0^\infty t dB(t) = \frac{\lambda^{(1)}}{v^{(1)}} + \frac{\lambda^{(2)}}{v^{(2)}} = \rho^{(1)} + \rho^{(2)}$ . Assume that the stability condition of the queueing system holds, i.e.,  $\rho = \rho^{(1)} + \rho^{(2)} < 1$ .

The LST of the service time distribution of class  $r$  is

$$g^{(r)}(s) = \int_0^\infty e^{-st} dB^{(r)}(t), \quad r = 1, 2 \quad (4.4)$$

and the LST of the residual service time distribution for class  $r$ , ( $r=1, 2$ ) is:

$$g_e^{(r)}(s) = \int_0^\infty e^{-st} (1 - B^{(r)}(t)) v^{(r)} dt = \frac{(1 - g^{(r)}(s)) v^{(r)}}{s} \quad (4.5)$$

The busy period time is the time between an interruption moment at which the server becomes busy due to an arriving customer and the first moment at which the server becomes available again. The LST of the busy time distribution of the high priority class, denoted by  $\delta^{(1)}(s)$ , is the smallest root of the *Kendall functional equation* (Cohen [27])

$$\delta^{(1)}(s) = g^{(1)}(s + \lambda^{(1)}(1 - \delta^{(1)}(s))) \quad (4.6)$$

where  $g^{(1)}$  is defined by (4.4).

Let  $c$  be the completion time of a customer, that is, the time elapsed from the moment when the service of the customer begins to the moment where the customer is completely served. As derived in (3.51) of [27], the LST,  $L_c(s)$ , of the complete time  $c$  for both preemptive-resume and non-preemptive resume, is:

$$L_{c(s)} = g^{(2)}(z(s)) \quad (4.7)$$

where

$$z(s) = s + \lambda^{(1)}(1 - \delta^{(1)}(s)) \quad (4.8)$$

For the preemptive-resume discipline, class 2 does not exist as far as class 1 is concerned. Thus in this case the waiting time distribution of the high-priority class is the same as in the FIFO queue without priorities. The probability distribution of the high-priority class was given in (4.2). Hence, in the following discussion of this section we only need to find the LSTs of the waiting time distribution of the low-priority class.

Denote  $W^{(2)}$  the cumulative distribution of the steady-state waiting time of the low-priority customers. According to (3.63) and (3.67) in [27], the LST of the low-priority waiting time  $W^{(2)}(t)$  is given by

$$L_{W^{(2)}}(s) = \frac{1 - \rho}{1 - \rho f(s)}, \quad s > 0 \quad (4.9)$$

where

$$f(s) = \frac{\rho^{(1)}}{\rho} h_0^{(1)}(s) + \frac{\rho^{(2)}}{\rho} g_e^{(2)}(z(s)) \quad (4.10)$$

$$h_0^{(1)}(s) = \frac{1 - \delta^{(1)}(s)}{\frac{1}{v^{(1)}} s + \rho^{(1)}(1 - \delta^{(1)}(s))} \quad (4.11)$$

Note that  $g_e^{(2)}(z(s))$  and  $\delta^{(1)}(s)$  are determined by (4.4) and (4.6) respectively.

Therefore, the low-priority response time denoted by  $T^{(2)}$  equals to the sum of the low-priority waiting time and the low-priority completion time whose LST is given by (4.7). The low-priority response time distribution is a convolution of the distributions of the low-priority waiting time and the low-priority completion time. Thus the LST of the low-priority response time distribution denoted by  $L_{T^{(2)}}(s)$  is:  $L_{T^{(2)}}(s) = L_{W^{(2)}}(s) \times L_{c(s)}$ . That is,

$$L_{T^{(2)}}(s) = \frac{(1 - \rho) \times g^{(2)}(z(s))}{1 - \rho f(s)}, \quad s > 0 \quad (4.12)$$

The distributions of the low-priority response time can be obtained by numerically inverting the LST of (4.12) respectively. This procedure can usually be done via the software package provided in [40].

For presentation purposes, we only consider an  $M/M/1$  queue in this dissertation since the distributions of the low-priority response time given by (4.12) can be specified and simplified, as discussed in the following subsection.

### The LST of The Low-priority Response Time Distribution In An $M/M/1$ Queue

In an  $M/M/1$  priority queue, assume that the service rates are  $\mu^{(r)}$  ( $r = 1, 2$ ). This means that  $v^{(r)} = \mu^{(r)}$ . In this case, from (4.4) and (4.5) we have that for  $r = 1, 2$ ,

$$g^{(r)}(s) = \frac{\mu^{(r)}}{s + \mu^{(r)}} \quad (4.13)$$

$$g_e^{(r)}(s) = \frac{\mu^{(r)}}{s + \mu^{(r)}} \quad (4.14)$$

Therefore, from (4.6) and (4.13) we derive the relation:

$$\delta^{(1)} = \frac{1 - \delta^{(1)}}{\frac{1}{u^{(1)}} s + \rho^{(1)}(1 - \delta^{(1)})} \quad (4.15)$$

and by solving for  $\delta^{(1)}$  we obtain

$$\delta^{(1)} = \frac{\eta - \sqrt{\eta^2 - 4\lambda^{(1)}\mu^{(1)}}}{2\lambda^{(1)}} \quad (4.16)$$

where  $\eta$  is determined by

$$\eta = s + \lambda^{(1)} + \mu^{(1)} \quad (4.17)$$

Furthermore, it follows from (4.10), (4.11) and (4.15) that  $h_0^{(1)}(s) = \delta^{(1)}$ , and

$$f(s) = \frac{\rho^{(1)}}{\rho} \frac{\eta - \sqrt{\eta^2 - 4\lambda^{(1)}\mu^{(1)}}}{2\lambda^{(1)}} + \frac{\rho^{(2)}}{\rho} \frac{\mu^{(2)}}{z(s) + \mu^{(2)}} \quad (4.18)$$

where  $z(s)$  is given by (4.8).

Moreover, from (4.6), (4.8) and (4.13) we have that

$$\delta^{(1)} = \frac{\mu^{(1)}}{z(s) + \mu^{(1)}}, \quad \text{or} \quad z(s) = \mu^{(1)}[(\delta^{(1)})^{-1} - 1] \quad (4.19)$$

Thus due to (4.19), the LST of the low-priority waiting time distribution given in (4.9) reduces to

$$\begin{aligned} L_{W^{(2)}}(s) &= \frac{1-\rho}{1-\rho^{(1)} \delta^{(1)} - \rho^{(2)} \frac{\mu^{(2)}}{z(s)+\mu^{(2)}}} \\ &= \frac{(1-\rho)\{\mu^{(1)}[(\delta^{(1)})^{-1}-1]+\mu^{(2)}\}}{(\mu^{(1)}-\mu^{(2)})\rho^{(1)}\delta^{(1)}+\mu^{(1)}(\delta^{(1)})^{-1}-\xi} \end{aligned} \quad (4.20)$$

where  $\xi$  is given by

$$\begin{aligned} \xi &= \mu^{(1)}(\rho^{(1)} + \rho^{(2)}) + (\mu^{(1)} - \mu^{(2)})(1 - \rho^{(2)}) \\ &= \mu^{(1)}\rho + (\mu^{(1)} - \mu^{(2)})(1 - \rho^{(2)}) \end{aligned} \quad (4.21)$$

Hence, from (4.16) and (4.20) we have the LST of the low-priority waiting distribution for an  $M/M/1$  queue

$$L_{W^{(2)}}(s) = \frac{(1-\rho)(\frac{\eta}{2} + \frac{1}{2}\sqrt{\eta^2 - 4\lambda^{(1)}\mu^{(1)}} - (\mu^{(1)} - \mu^{(2)}))}{(1 - \frac{\mu^{(2)}}{2\mu^{(1)}})\eta - \xi + \frac{\mu^{(2)}}{2\mu^{(1)}}\sqrt{\eta^2 - 4\lambda^{(1)}\mu^{(1)}}} \quad (4.22)$$

Consequently, replacing (4.22) in (4.12) gives the LST of the low-priority response time distribution

$$\begin{aligned} L_{T^{(2)}}(s) &= \mu^{(2)}(1-\rho)[\frac{\eta}{2} + \frac{1}{2}\sqrt{\eta^2 - 4\lambda^{(1)}\mu^{(1)}} \\ &\quad - (\mu^{(1)} - \mu^{(2)})]\{(z(s) + \mu^{(2)}) \times [(1 - \frac{\mu^{(2)}}{2\mu^{(1)}})\eta \\ &\quad - \xi + \frac{\mu^{(2)}}{2\mu^{(1)}}\sqrt{\eta^2 - 4\lambda^{(1)}\mu^{(1)}}]\}^{-1}, \quad s > 0 \end{aligned} \quad (4.23)$$

where  $z(s)$ ,  $\eta$  and  $\xi$  are given by (4.8), (4.17) and (4.21) respectively.

Additionally, when the service rates of the two priority classes are the same, i.e.,  $\mu^{(1)} = \mu^{(2)}$ , (4.18) reduces to the following expression:

$$\begin{aligned} f(s) &= \frac{\rho^{(1)}}{\rho} \delta^{(1)} + \frac{\rho^{(2)}}{\rho} \frac{\mu^{(2)}}{z(s) + \mu^{(2)}} \\ &= \frac{\rho^{(1)}}{\rho} \delta^{(1)} + \frac{\rho^{(2)}}{\rho} \frac{\mu^{(1)}}{z(s) + \mu^{(1)}} = \delta^{(1)} \end{aligned}$$

due to (4.19). Hence, equation (4.9) becomes  $L_{W^{(2)}}(s) = \frac{1-\rho}{1-\rho\delta^{(1)}}$ . Thus, the LST of the low-priority response time distribution is

$$L_{T^{(2)}}(s) = \frac{(1-\rho)\delta^{(1)}}{1-\rho\delta^{(1)}} \quad (4.24)$$

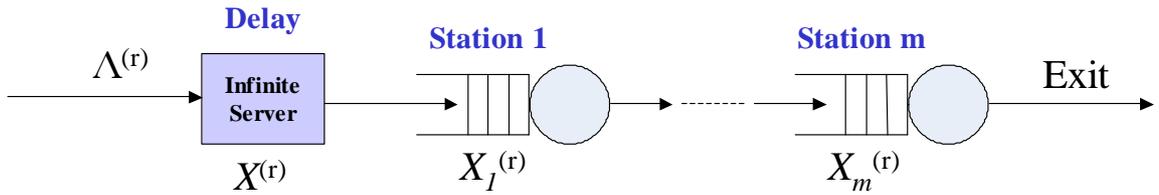


Figure 4.1: A Tandem-station Service Model for Multiple-class Customer Services

based on (4.12), (4.19) and  $\mu^{(1)} = \mu^{(2)}$ . It should be pointed out that when the utilization  $\rho$  is fixed, the LST of low-priority response time distribution does not depend on the arrival rate  $\lambda^{(2)}$  of the low-priority class.

### 4.3.2 Algorithms for the Resource Optimization Problem

In this subsection, we study two queueing network models that depict the path that service requests have to follow through the service provider's resource stations. These two models are shown in Figures 4.1 and 4.2. We refer to these two queueing models as service models since they depict the resources used to provide a service to a customer.

The first service model consists of a single infinite server and  $m$  stations numbered sequentially from 1 to  $m$  as shown in Figure 4.1. Each station  $j$  is modeled as a priority queue served by  $n_j$  identical servers, each providing a service at the rate  $\mu_j^{(r)}$ , where  $r = 1, 2$ . Let  $\Lambda^{(r)}$  be the external arrival rate to the infinite server, and let  $\lambda^{(r)}$  and  $\lambda_j^{(r)}$  be the effective arrival rates to the infinite server and station  $j$ ,  $j = 1, 2, \dots, m$ . We assume that all service times are exponentially distributed and the external arrival to the infinite server occurs in a Poisson fashion.

The infinite server represents the total propagation delay from the user to the service provider and back, and also from station 1 to  $m$ . Each station carries out a particular function. For instance, it could be a database server, a file server, a web server, a group of CPUs and local disks, etc. We consider two priority classes of customers in this dissertation. In the following discussion each station is modeled as a single  $M/M/1$  priority queue with arrival rate  $\lambda_j^{(r)}$  and service rate  $\psi^{(r)}(n_j)\mu_j^{(r)}$ , where  $\psi^{(r)}(n_j)$  is a known function of  $n_j$  ( $r = 1, 2$ ), and depends on the configuration of servers and the type of customers at each station. It is non-decreasing and can be inverted, i.e.,  $(\psi^{(r)})^{-1}$  exists ( $r = 1, 2$ ). For instance, suppose that a station represents a group of CPUs. Then,  $\psi^{(r)}(n)$  can be seen as a CPU scaling factor for the number of CPUs from 1 to  $n$ .

According to [22],  $\psi^{(r)}(n) = (\xi^{(r)})^{\log n}$ , where  $\xi^{(r)}$  is a basic scaling factor from 1 CPU to 2, and it ranges from 1 to 2 ( $r = 1, 2$ ). So,  $(\psi^{(r)})^{-1}(n) = (\xi^{(r)})^{-\log n}$ . Additionally, each station that is modeled as a *single M/M/1* priority queue with service rate  $\psi^{(r)}(n_j)\mu_j^{(r)}$  can only serve either high-priority or low priority customers at one time. Hence,  $\psi^{(1)}(n_j)\mu_j^{(1)}$  is considered as the same as  $\psi^{(2)}(n_j)\mu_j^{(2)}$ .

Since the queueing network is overtake-free, the waiting time of a customer at a station is independent of its waiting times at other stations (see [30] and [90]). Let  $X^{(r)}$  be the service time at the infinite server and  $X_j^{(r)}$  be the time elapsed from the moment a customer arriving at station  $j$  to the moment it departs from the station ( $r = 1, 2$ ). Then, the total response time is  $T^{(r)} = X^{(r)} + X_1^{(r)} + X_2^{(r)} + \dots + X_m^{(r)}$ , and hence the LST (Laplace-Stieltjes transform) of the response time  $T$  is

$$L_{T^{(r)}}(s) = L_{X^{(r)}}(s)L_{X_1^{(r)}}(s) \cdots L_{X_m^{(r)}}(s) \quad (4.25)$$

where  $L_{X^{(r)}}(s)$  is the LST of the service time  $X^{(r)}$  given by

$$L_{X^{(r)}}(s) = \frac{\lambda^{(r)}}{s + \lambda^{(r)}} \quad (4.26)$$

and  $L_{X_j^{(r)}}(s)$  is the LST of the response time  $X_j^{(r)}$  at the  $j$ -th station, where  $r = 1, 2$ .

Due to the preemptive-resume priority, the high-priority response time is the same as in the single class FIFO *M/M/1* whose probability distribution can be expressed as (4.2) in section 4.1. Thus  $L_{X_j^{(1)}}(s)$  is determined by

$$L_{X_j^{(1)}}(s) = \frac{\psi^{(1)}(n_j)\mu_j(1 - \rho_j^{(1)})}{s + \psi^{(1)}(n_j)\mu_j(1 - \rho_j^{(1)})}, \quad (j = 1, 2, \dots, m) \quad (4.27)$$

and  $L_{X_j^{(2)}}(s)$  is the LST of the low-priority response time given by

$$L_{X_j^{(2)}}(s) = \frac{(1 - \rho_j)\delta_j^{(1)}}{1 - \rho_j\delta_j^{(1)}}, \quad (j = 1, 2, \dots, m) \quad (4.28)$$

due to (4.24), where  $\rho_j^{(r)} = \frac{\lambda_j^{(r)}}{\psi^{(r)}(n_j)\mu_j}$ ,  $\rho_j = \frac{\lambda_j^{(1)} + \lambda_j^{(2)}}{\psi^{(r)}(n_j)\mu_j}$ ,  $\delta_j^{(1)}$  is given by

$$\delta_j^{(1)} = \frac{\eta_j - \sqrt{\eta_j^2 - 4\psi^{(1)}(n_j)\lambda_j^{(1)}\mu_j^{(1)}}}{2\psi^{(1)}(n_j)\lambda_j^{(1)}}$$

and  $\eta_j = s + \lambda_j^{(1)} + \psi^{(1)}(n_j)\mu_j^{(1)}$  for  $j = 1, 2, \dots, m$ .

From (4.25), (4.26), (4.27) and (4.28) we have that

$$L_{T^{(1)}}(s) = \frac{\lambda^{(1)}}{s + \lambda^{(1)}} \prod_{j=1}^m \frac{\psi^{(1)}(n_j)\mu_j(1 - \rho_j^{(1)})}{s + \psi^{(1)}(n_j)\mu_j(1 - \rho_j^{(1)})}$$

and

$$L_{T^{(2)}}(s) = \frac{\lambda^{(2)}}{s + \lambda^{(2)}} \prod_{j=1}^m \frac{(1 - \rho_j)\delta_j^{(1)}}{1 - \rho_j\delta_j^{(1)}}$$

We observe that  $f_{T^{(r)}}(t)$  and  $F_{T^{(r)}}(t)$  ( $r = 1, 2$ ) are usually nonlinear functions of  $t$  and  $n_j$ . Hence, the resource optimization problem is an  $n$ -dimensional linear optimization problem subject to nonlinear constraints. In general, it is not easy to solve this problem. However, the complexity of the problem can be significantly reduced by requiring that the service rates of the queues in the service model in Figure 4.1 are all equal. That is, we find the optimum value of  $n_1, \dots, n_m$  such that  $\psi^{(r)}(n_1)\mu_1^{(r)} = \dots = \psi^{(r)}(n_m)\mu_m^{(r)}$  ( $r = 1, 2$ ), called *balanced utilization* or *balanced condition* as similarly defined in Chapter 3.

From the traffic equations:  $\lambda^{(r)} = \lambda_j^{(r)} = \Lambda^{(r)}$  ( $j = 1, 2, \dots, m$ ), we have that the utilization of each station  $\rho_j^{(r)} = \frac{\lambda_j^{(r)}}{\psi^{(r)}(n_j)\mu_j^{(1)}} = \frac{\Lambda^{(r)}}{\psi^{(r)}(n_j)\mu_j^{(1)}}$ . Thus we have that for the high-priority queue,  $\hat{a}_i = \psi^{(1)}(n_i)\mu_i^{(1)}(1 - \rho_i^{(1)}) = \psi^{(1)}(n_j)\mu_j(1 - \rho_j^{(1)}) = \hat{a}_j \triangleq \hat{a}$  that implies  $n_j = (\psi^{(1)})^{-1}\left(\frac{\hat{a} + \lambda_j^{(1)}}{\mu_j}\right)$  ( $i, j = 1, 2, \dots, m$ ). Hence, from (4.25) we have  $f_{T^{(1)}}(t) = L^{-1}\left\{\frac{\lambda^{(1)}}{s + \lambda^{(1)}} \cdot \frac{\hat{a}^m}{(s + \hat{a})^m}\right\}$ , and subsequently obtain that

$$F_{T^{(1)}}(t) = L^{-1}\left\{\frac{\lambda^{(1)}}{s(s + \lambda^{(1)})} \cdot \frac{\hat{a}^m}{(s + \hat{a})^m}\right\} \quad (4.29)$$

Moreover, for the low-priority queue, we have that  $\rho_{j_1} = \rho_{j_2} \triangleq \hat{\rho}$ , and  $\delta_{j_1}^{(1)} = \delta_{j_2}^{(1)} \triangleq \hat{\delta}^{(1)}$ , for  $j_1, j_2 = 1, 2, \dots, m$ , due to  $\psi^{(2)}(n_1)\mu_1^{(2)} = \dots = \psi^{(2)}(n_m)\mu_m^{(2)} \triangleq \hat{b}$ . Thus  $L_{T^{(2)}}(s)$  reduces to

$$L_{T^{(2)}}(s) = \frac{\lambda^{(2)}}{s + \lambda^{(2)}} \frac{(1 - \hat{\rho})^m (\hat{\delta}^{(1)})^m}{(1 - \hat{\rho}\hat{\delta}^{(1)})^m} \quad (4.30)$$

which is a function of only one variable  $\hat{b}$ , since  $\hat{\rho}$  and  $\hat{\delta}^{(1)}$  are considered as functions of only one variable  $\hat{b}$ . Consequently,  $\sum_{j=1}^m n_j c_j$  reduces to a function of variable  $\hat{a}$  due to

$$n_j^{(1)} = \lceil (\psi^{(1)})^{-1}\left(\frac{\hat{a} + \lambda_j^{(1)}}{\mu_j}\right) \rceil$$

for a high-priority customer and  $n_j^{(2)} = \lceil (\psi^{(2)})^{-1}\left(\frac{\hat{b}}{\mu_j^{(2)}}\right) \rceil$  for a low-priority customer. Thus the resource optimization problem can be decomposed into the two one-dimensional resource optimiza-

tion problems for both high-priority and low-priority customers respectively. We have the following algorithm for the resource optimization problem.

**Algorithm 4-1:**

1. The minimization problem for high-priority customers: Find  $\hat{a}$  in the one dimensional optimization problem:

$$\hat{a}^{min} \leftarrow \arg \min_{\hat{a}} F_{T^{(1)}}(t)|_{t=T_D^{(1)}}$$

subject to the constraint  $F_{T^{(1)}}(t)|_{t=T_D^{(1)}} \geq \gamma^{(1)}\%$  at  $\hat{a} = \hat{a}^{min}$ , where  $F_{T^{(1)}}(t)$  is given by (4.29).

2. The minimization problem for low-priority customers: Find  $\hat{b}$  in the one-dimensional optimization problem:

$$\hat{b}^{min} \leftarrow \arg \min_{\hat{b}} F_{T^{(2)}}(t)|_{t=T_D^{(2)}}$$

subject to the constraint  $F_{T^{(2)}}(t)|_{t=T_D^{(2)}} \geq \gamma^{(2)}\%$  at  $\hat{b} = \hat{b}^{min}$ , where  $F_{T^{(2)}}(t)$  is given by (4.30).

3. Compute integers  $n_j^{(1)}$  and  $n_j^{(2)}$  by using the expressions:

$$n_j^{(1)} = \lceil (\psi^{(1)})^{-1} \left( \frac{\hat{a}^{min}}{\mu_j^{(1)}(1 - \rho_j^{(1)})} \right) \rceil \text{ and } n_j^{(2)} = \lceil (\psi^{(2)})^{-1} \left( \frac{\hat{b}^{min}}{\mu_j^{(2)}} \right) \rceil$$

Then, calculate  $n_j = \max\{n_j^{(1)}, n_j^{(2)}\}$  ( $j = 1, 2, \dots, m$ ).

4. Check if  $1 \leq n_j \leq N_j$  ( $j = 1, 2, \dots, m$ ) and  $I \leq C^D$  are satisfied. If yes, the obtained  $n_j$  is the number of servers required at each station. Otherwise, print “the problem cannot be solved.”

Note that when the balanced utilization as above is satisfied, the suboptimal solution obtained by Algorithm 4-1 is optimal for the resource optimization problem. This statement is also true for Algorithm 4-2 that will be given later.

In order to allow for a more complex execution path of a service request, we extended the above model to a service model with feedback, as shown in Figure 4.2. Infinite server 1 represents the total propagation delay within a network provider and infinite server 2 represents the propagation delay within the service provider, i.e. among stations 1 to  $m$ . In this figure, a customer upon completion of its service at the  $m$ -th station, exits the system with probability  $\alpha$ , or returns to the beginning of the system with probability  $1 - \alpha$ . We note that the model shown in Figure 4.2 can

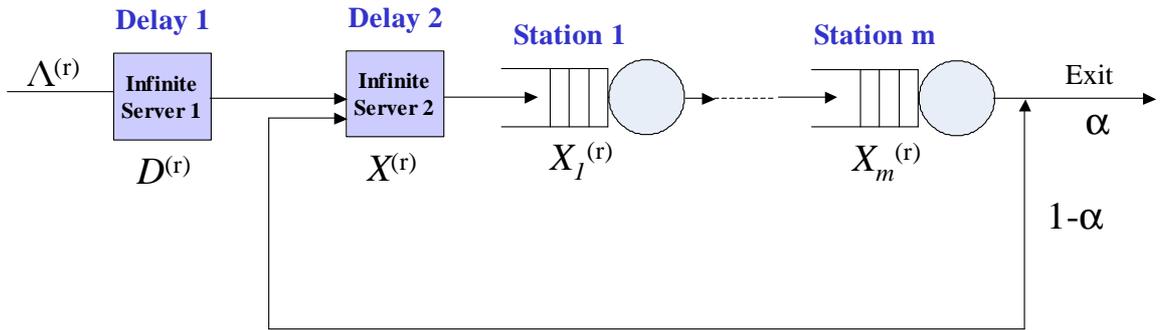


Figure 4.2: A Service Model with Feedback for Multiple-class Customer Services

be easily extended to a network of queues arbitrarily connected. We reuse the notation in the first model shown in Figure 4.1:  $\Lambda^{(r)}$  as the external arrival rate,  $\lambda_d^{(r)}$ ,  $\lambda^{(r)}$  and  $\lambda_j^{(r)}$  as the effective arrival rates to the second infinite server and station  $j$ , and  $\mu_j^{(r)}$  as the service rate at station  $j$ , where  $j = 1, 2, \dots, m$  and  $r = 1, 2$ .

The main difficulty of this resource optimization problem is to find  $f_{T^{(r)}}(t)$ , the probability distribution function of  $T^{(r)}$  ( $r = 1, 2$ ). We obtain this probability distribution function assuming that the waiting time of a customer at a station is independent of its waiting time at other stations, and each visit at the same station  $j$  is independent of the others. (We note that this assumption of independence does not hold in queueing networks with feedback. However, as will be discussed in the validation section the solution obtained has a good accuracy.) We first have the traffic equations:  $\lambda_d^{(r)} = \Lambda^{(r)}$ ,  $\lambda^{(r)} = \Lambda^{(r)} + (1 - \alpha)\lambda_m^{(r)}$  and  $\lambda_j^{(r)} = \lambda^{(r)}$  that implies  $\lambda_j^{(r)} = \lambda^{(r)} = \frac{\Lambda^{(r)}}{\alpha}$ . and the utilization of each station is  $\rho_j^{(1)} = \frac{\lambda_j^{(r)}}{\psi^{(r)}(n_j)\mu_j} = \frac{\Lambda^{(r)}}{\alpha\mu_j\psi^{(r)}(n_j)}$  ( $j = 1, 2, \dots, m$ ).

Then, the high-priority and low-priority response times of the  $k$ -th pass at the infinite station and the  $j$ -th station are considered as the sum of  $m + 2$  random variables  $T^{(r)}(k) = D^{(r)} + X^{(r)} + X_1^{(r)} + \dots + X_m^{(r)}$ , where we assume that the high-priority (or low-priority) waiting time of a customer at a station is independent of its high-priority (or low-priority) waiting times in other visits to the same station.  $D^{(r)}$  and  $X^{(r)}$  are the service times of class  $r$  at the first and second infinite servers respectively, and  $X_j^{(r)}$  is the time elapsed from the moment a class  $r$  customer arrives at station  $j$  to the moment it departs from it. Then, the total response time is  $T^{(r)} = \sum_{k=0}^{\infty} p(k)T^{(r)}(k)$ , where  $p(k)$  is the steady state probability that a request will circulate  $k$  times at the infinite station and the  $j$ -th station through the computing system.  $p(k)$  is determined by  $p(k) = \alpha(1-\alpha)^{k-1}$ . Thus the LST of the response time  $T^{(r)}$  is  $L_{T^{(r)}}(s) = L_D(s) \sum_{k=0}^{\infty} p(k) L_X^k(s) L_{X_1^{(r)}}^k(s) \dots L_{X_m^{(r)}}^k(s)$ ,

which can be re-written as follows:

$$L_{T^{(r)}}(s) = \frac{\alpha L_{D^{(r)}}(s) L_X(s) \prod_{j=1}^m L_{X_j^{(r)}}(s)}{1 - (1 - \alpha) L_X(s) \prod_{j=1}^m L_{X_j^{(r)}}(s)} \quad (4.31)$$

where  $L_{D^{(r)}}(s)$  is the LST of the service time  $D^{(r)}$  given by  $L_{D^{(r)}}(s) = \frac{\Lambda^{(r)}}{s + \Lambda^{(r)}}$ , and replacing  $L_X(s)$  and  $L_{X_j^{(r)}}(s)$  ( $j = 1, 2, \dots, m$ ) by (4.26), (4.27) and (4.28) in (4.31), we have that

$$L_{T^{(1)}}(s) = \frac{(\Lambda^{(1)})^2 \prod_{j=1}^m \hat{a}_j}{(s + \Lambda^{(1)}) [(s + \lambda^{(1)}) \prod_{j=1}^m (s + \hat{a}_j) - (1 - \alpha) \lambda^{(1)} \prod_{j=1}^m \hat{a}_j]} \quad (4.32)$$

where  $\hat{a}_j = \psi^{(1)}(n_j) \mu_j^{(1)} (1 - \rho_j^{(1)})$ , and

$$\begin{aligned} L_{T^{(2)}}(s) &= (\Lambda^{(2)})^2 \prod_{j=1}^m [(1 - \rho_j) \delta_j^{(1)}] (s + \Lambda^{(2)})^{-1} \\ &\quad \times \{ (s + \lambda^{(2)}) \prod_{j=1}^m (1 - \rho_j \delta_j^{(1)}) \\ &\quad - (1 - \alpha) \lambda^{(2)} \prod_{j=1}^m [(1 - \rho_j) \delta_j^{(1)}] \}^{-1} \end{aligned} \quad (4.33)$$

To find the response time distribution  $f_{T^{(r)}}(t)$ , we are required to invert the LST given by (4.32) using partial fraction decomposition of a rational function. However, the partial fraction decomposition of the rational function requires searching for roots of a high-order polynomial. It is usually not an easy task when the order of the polynomial is more than 5. Instead, in this dissertation the LST given by (4.32) is inverted numerically, as (4.33).

Similarly, we want to find  $n_1, \dots, n_m$  such that the best utilization of these stations is achieved, which implies that each station has the same maximal service capacity. That is,  $\hat{a}_i = \hat{a}_j = \hat{a}$ ,  $\hat{\rho}_i = \hat{\rho}_j = \hat{\rho}$  and  $\hat{\delta}_i^{(1)} = \hat{\delta}_j^{(1)} = \hat{\delta}^{(1)}$  ( $i, j = 1, \dots, m$ ). Then, from equations (4.32) and (4.33), and  $F_{T^{(r)}}(t) = L^{-1}\{L_{T^{(r)}}(s)/s\}$  ( $r = 1, 2$ ) we have

$$F_{T^{(1)}}(t) = L^{-1}\left\{ \frac{(\Lambda^{(1)})^2 \hat{a}^m}{s(s + \Lambda^{(1)}) [(s + \lambda^{(1)}) (s + \hat{a})^m - (1 - \alpha) \lambda^{(1)} \hat{a}^m]} \right\} \quad (4.34)$$

and

$$\begin{aligned} F_{T^{(2)}}(t) &= L^{-1}\left\{ (\Lambda^{(2)})^2 [(1 - \hat{\rho}) \hat{\delta}^{(1)}]^m s^{-1} (s + \Lambda^{(2)})^{-1} \right. \\ &\quad \times \{ (s + \lambda^{(2)}) (1 - \hat{\rho} \hat{\delta}^{(1)})^m \\ &\quad \left. - (1 - \alpha) \lambda^{(2)} [(1 - \hat{\rho}) \hat{\delta}^{(1)}]^m \}^{-1} \right\} \end{aligned} \quad (4.35)$$

Thus we have the following algorithm for the resource optimization problem in the model shown in Figure 4.2.

**Algorithm 4-2:**

Steps 1-4 are the same as Steps 1-4 in Algorithm 4-1 except  $F_{T(2)}(t)$  and  $F_{T(2)}(t)$  given by (4.34) and (4.35) respectively.

Note that if we cannot get a solution for the resource optimization problem using Algorithm 4-1 (or 4-2), then the service provider cannot execute the service request for the service model 1 (or 2) due to at least one of the following reasons: (i) the service provider has an insufficient resource (i.e.,  $N_j$  is too small), (ii) a pre-specific fee is too low (i.e.,  $I > C_D$ ), and (iii) a network connection is either too slow or has a problem so that (4.1) cannot be satisfied. Using these information, we may detect and debug either a network problem or a service provider's capacity problem, or the SLA needs to be re-negotiated.

**4.4 Numerical Validations**

In this section we demonstrate the accuracy and applicability of our proposed approximation method.

Two types of errors are introduced in our proposed approximation method. The first, hereafter referred to as Class I error, comes from numerically inverting the Laplace transform. The other, hereafter referred to as Class II error, is due to the assumptions that the waiting time of a customer at each station is independent of the waiting times at the other stations, and it is also independent of its waiting times in other visits to the same station.

Let us recall that the relative error % is used to measure the accuracy of the approximate results compared to model simulation results, and it is defined as follows

$$\text{Relative error \%} = \frac{\text{Approximate Result} - \text{Simulation Result}}{\text{Simulation Result}} \times 100$$

as given in (3.11).

As is seen, our proposed approximation method heavily depends on the computation of the inverse Laplace transform

$$f(t) = \frac{1}{2\pi i} \int_{c-i\infty}^{c+i\infty} e^{ts} \hat{f}(s) ds$$

where  $\hat{f}(s)$  is the image function of the inverse Laplace function  $f(t)$ , and it is defined by

$$\hat{f}(t) = \int_0^{\infty} e^{-st} f(t) dt$$

where  $f(t)$  is called the original function of  $\hat{f}(s)$  and it is a real or complex-valued function defined on the positive part  $\mathcal{R}_+$ .

The numerical inversion of the Laplace transform has been widely studied and several efficient methods have been proposed in the past a few decades (see Graf [40]). However, as is described in [40], the numerical computation of an inverse Laplace transform is an ill-posed problem. The inverse Laplace transform is determined by the singularities of the image function  $\hat{f}(s)$ . This means that the behavior of the image function near the singularities determines the inverse Laplace transform. Hence, we need to consider the singularities of the image function in our numerical validations. Additionally, whereas some numerical methods work for certain image functions well, they may provide poor results for other image functions. No single method works for every image functions. Thus, in our numerical validations we employed several different numerical methods for a given image function. If two or more methods can reach about the same results, then we are confident that the derived numerical inverse Laplace transform is correct. These numerical methods include the inversion methods using Laguerre functions and Fourier functions (see Graf [40]), Gaussian quadrature formulas (see Piessens [77]), and the method by Gaver [38] and Stehfest [87].

We study the accuracy of our proposed approximation method using several examples below.

Contrary to the distribution of the high-priority response time whereby it is the single class FIFO  $M/M/1$  case, as given in section 4.1, the distribution of the low-priority response time in an  $M/M/1$  queue does not have a closed-form solution. The distribution whose Laplace transform is given by (4.24) have to be evaluated numerically by using one of the numerical methods for inverting a Laplace transform (e.g., [40]). Hence, we first verify the accuracy of the approximate method for the single queue given by (4.24). Let  $\mu^{(1)} = \mu^{(2)} = \frac{1000}{9} = 111.11$  and  $\lambda^{(r)}$  is varied ( $r = 1, 2$ ).

We simulated the queueing network using Arena (see [5]) and the analytical method was implemented in Mathematica using the package of the inverse Laplace transform in Graf [40]. The simulation results are considered as “exact” since the simulation model is an exact representation of the queueing network under study.

Tables 4.1 and 4.2 show the simulated and approximate cumulative distributions of the low-priority response times for the two different cases: (i)  $\lambda^{(1)} = \lambda^{(2)} = 50$ ; (ii)  $\lambda^{(1)} = 75$  and  $\lambda^{(2)} = 25$ . In these two tables, the column labeled “Simul” gives the simulation result, the column labeled “Approx” gives the approximate result, and the column labeled “R-Err %” gives their relative errors. The same abbreviations are also used in other tables in the rest of this section.

Table 4.1: The Cumulative Distribution of the Low-priority Response Time for  $\lambda^{(1)} = \lambda^{(2)} = 50$ 

Response Time	Simul	Approx	R-Err %
0.1	0.5183	0.4821	-6.9844
0.3	0.8691	0.8318	-4.2918
0.5	0.9691	0.9447	-2.5178
0.7	0.9911	0.9818	-0.9384
0.9	0.9997	0.9940	-0.5702
1.1	1.0000	0.9980	-0.2000
1.3	1.0000	0.9994	-0.0600
1.5	1.0000	0.9998	-0.0200
1.7	1.0000	0.9999	-0.0100
1.9	1.0000	1.0000	0.0000

It appears that using the package of the inverse Laplace transform in Graf [40] we can get a good accuracy for the numerical inversion of the low-priority response time distribution given by (4.24) respectively.

Then, we verify the accuracy of our approach for the first service model shown in Figure 4.1. Let  $m = 8$ ,  $\lambda^{(1)} = 100$ ,  $\lambda^{(2)} = 50$ ,  $\mu_1^{(1)} = 48$ ,  $\mu_2^{(1)} = 18$ ,  $\mu_3^{(1)} = 85$ ,  $\mu_4^{(1)} = 32$ ,  $\mu_5^{(1)} = 49$ ,  $\mu_6^{(1)} = 24$ ,  $\mu_7^{(1)} = 28$ ,  $\mu_8^{(1)} = 38$ ,  $\mu_1^{(2)} = 42$ ,  $\mu_2^{(2)} = 15$ ,  $\mu_3^{(2)} = 60$ ,  $\mu_4^{(2)} = 25$ ,  $\mu_5^{(2)} = 41$ ,  $\mu_6^{(2)} = 18$ ,  $\mu_7^{(2)} = 26$ , and  $\mu_8^{(2)} = 35$ . We also choose  $N_j = 100$ ,  $c_j = 1$ ,  $\psi^{(1)}(n_j) = 1.5^{\log n_j}$ ,  $\psi^{(2)}(n_j) = 1.55^{\log n_j}$  and  $C_D = 300$  ( $j = 1, \dots, 8$ ).

Tables 4.3 and 4.4 show the simulated and approximate cumulative distributions of the high-priority and low-priority response times respectively. It appears that the results obtained by Algorithm 4-1 are very accurate. It is shown in Table 4.5 that the optimal number of servers is required for 97.5% of the high-priority response time to be less than  $T_D^{(1)} = 0.16$  and for 97.5% of the low-priority response time to be less than  $T_D^{(2)} = 0.8$ . Moreover, the optimal number of servers is required for satisfying both the high-priority and the low-priority response times is shown in Table 4.5. The exact optimal number of servers, obtained by exhaustive search using the simulation model, and assuming that each station has balanced utilization, is consistent with the ones shown in Table 4.5. So,  $I = 214 < C_D$ . We point out that the relative errors shown in Tables 4.3 and 4.4 are only due to the Class I error since the Class II error is not present for this model.

Let us now consider an example of the service model 2 shown in Figure 4.2. We choose  $m = 8$ ,  $\Lambda^{(1)} = 55$ ,  $\Lambda^{(2)} = 42$ ,  $\mu_1^{(1)} = 12$ ,  $\mu_2^{(1)} = 46$ ,  $\mu_3^{(1)} = 95$ ,  $\mu_4^{(1)} = 25$ ,  $\mu_5^{(1)} = 35$ ,  $\mu_6^{(1)} = 20$ ,  $\mu_7^{(1)} = 10$ , and  $\mu_8^{(1)} = 98$ ,  $\mu_1^{(2)} = 15$ ,  $\mu_2^{(2)} = 42$ ,  $\mu_3^{(2)} = 90$ ,  $\mu_4^{(2)} = 18$ ,  $\mu_5^{(2)} = 28$ ,  $\mu_6^{(2)} = 15$ ,

Table 4.2: The Cumulative Distribution of the Low-priority Response Time for  $\lambda^{(1)} = 75$  and  $\lambda^{(2)} = 25$

Response Time	Simul	Approx	R-Err %
0.1	0.4175	0.3837	-8.0958
0.3	0.7115	0.6783	-4.6662
0.5	0.8617	0.8225	-4.5491
0.7	0.9338	0.9003	-3.5875
0.9	0.9669	0.9435	-2.4201
1.1	0.9835	0.9679	-1.5862
1.3	0.9922	0.9817	-1.0583
1.5	0.9970	0.9895	-0.7523
1.7	0.9982	0.9940	-0.4208
1.9	0.9996	0.9966	-0.3001
2.1	1.0000	0.9980	-0.2000
2.3	1.0000	0.9999	-0.0100
2.5	1.0000	1.0000	0.0000

Table 4.3: The Cumulative Distribution of The High-priority Response Time in Model 1

Response Time	Simul	Approx	R-Err %
0.02	0.0002	0.0002	0.0000
0.04	0.0214	0.0214	0.0000
0.06	0.1542	0.1528	-0.9079
0.08	0.4085	0.4075	-0.2448
0.10	0.6691	0.6672	-0.2840
0.12	0.8459	0.8450	-0.1064
0.14	0.9385	0.9379	-0.0639
0.16	0.9781	0.9780	-0.0102
0.18	0.9931	0.9929	-0.0201
0.20	0.9980	0.9979	-0.0100
0.22	0.9995	0.9994	-0.0100
0.24	0.9998	0.9998	0.0000
0.26	0.9999	1.0000	0.0100
0.28	1.0000	1.0000	0.0000

Table 4.4: The Cumulative Distribution of the Low-priority Response Time in Model 1

Response Time	Simul	Approx	R-Err %
0.2	0.1767	0.1473	-16.6384
0.3	0.4538	0.4396	-3.1291
0.4	0.6982	0.7082	1.4323
0.5	0.8541	0.8719	2.0841
0.6	0.9389	0.9503	1.2142
0.7	0.9774	0.9824	0.5116
0.8	0.9925	0.9942	0.1713
0.9	0.9978	0.9982	0.0401
1	0.9993	0.9995	0.0200
1.1	0.9998	0.9999	0.0100
1.2	1.0000	1.0000	0.0000

Table 4.5: The Optimal Number of Servers in Model 1

Station	1	2	3	4	5	6	7	8
High-priority Customer	12	62	5	23	12	38	29	18
Low-priority Customer	12	61	7	27	13	46	26	16
All the Customers	12	62	7	27	13	46	29	18

$\mu_7^{(2)} = 5$ , and  $\mu_8^{(2)} = 82$ , and  $\alpha = 0.67$ . Let us also select  $N_j = 250$ ,  $c_j = 1$ ,  $\psi^{(1)}(n_j) = 1.5^{\log n_j}$ ,  $\psi^{(2)}(n_j) = 1.55^{\log n_j}$ , and  $C_D = 800$  ( $j = 1, \dots, 8$ ). Thus it follows from equation:  $\lambda^{(r)} = \frac{\Lambda^{(r)}}{\alpha}$  ( $r = 1, 2$ ) that  $\lambda^{(1)} = 82.09$  and  $\lambda^{(2)} = 62.69$ .

We obtained the cumulative distribution of the response time by solving (4.34) and (4.35) using the package of the inverse Laplace transforms in Graf [40]. Table 4.6 shows the number of servers in the eight stations necessary to ensure the 95% SLA guarantee for  $T_D^{(1)} = 0.3$  and  $T_D^{(2)} = 1.0$  respectively, and the number of servers in the eight stations necessary to ensure all the customers. We also simulated the tandem queueing network and validated using the brute-force

Table 4.6: The Optimal Number of Servers in Model 2

Station	1	2	3	4	5	6	7	8
High-priority Customer	123	13	4	35	20	52	168	4
Low-priority Customer	61	12	4	46	23	61	342	5
All the Customers	123	13	4	46	23	61	342	5

Table 4.7: The Cumulative Distribution of The High-priority Response Time in Model 2

Response Time	Simul	Approx	R-Err %
0.1	0.3879	0.3850	-0.7476
0.2	0.8343	0.8332	-0.1318
0.3	0.9528	0.9547	0.1994
0.4	0.9867	0.9877	0.1013
0.5	0.9961	0.9966	0.0502
0.6	0.9989	0.9991	0.0200
0.7	0.9997	0.9998	0.0100
0.8	0.9999	0.9999	0.0000
0.9	1.0000	1.0000	0.0000

approach that these numbers of servers obtained by our approximate method are in fact optimal, provided that each station has balanced utilization. The optimal number of servers is given in Table 4.6. It derives that  $I = 617 < C_D$ , i.e., Step 4 in Algorithm 4-2 is met. Tables 4.7 and 4.8 gives the cumulative distributions of the high-priority and low-priority response times obtained using the approximate method and the simulation method, and their relative error %. The relative error comes from both Classes I and II error. We note that our approximate method has a very good accuracy when percentiles are high. Extensive numerical results (not reported here due to lack of space) point to the fact that the independence assumption has little impact on the accuracy of the results when the number of nodes is large. A contributing factor is that typically we are interested in values of the cumulative distribution of the response time that correspond to very high percentiles for which the approximate results seem to have a very good accuracy.

## 4.5 Concluding Remarks

We proposed an approach for resource optimization in a service provider's computing environment, whereby we minimize the total cost of computer resources allocated to a priority-class customer so that satisfies a given percentile of the response time for each class customer. We have formulated the resource optimization problem as an optimization subject to SLA constraints for a service model with or without feedback. We have derived the LSTs of a customer's high-priority and low-priority response times. We further developed an efficient and accurate numerical solution for inverting the LSTs of a high-priority and low-priority customer's response time numerically. In the

Table 4.8: The Cumulative Distribution of the Low-priority Response Time in Model 2

Response Time	Simul	Approx	R-Err %
0.2	0.1679	0.1467	-12.6266
0.4	0.6144	0.6116	-0.4557
0.6	0.8219	0.8207	-0.1460
0.8	0.9126	0.9177	0.5588
1.0	0.9567	0.9622	0.5749
1.2	0.9784	0.9826	0.4293
1.4	0.9889	0.9920	0.3135
1.6	0.9943	0.9963	0.2011
1.8	0.9971	0.9978	0.0702
2.0	0.9985	0.9992	0.0701
2.2	0.9993	0.9996	0.0300
2.4	0.9996	0.9998	0.0200
2.6	0.9998	0.9999	0.0100
2.8	0.9999	1.0000	0.0100
3.0	0.9999	1.0000	0.0100
3.2	1.0000	1.0000	0.0000

resource optimization problem we are typically interested in values of the cumulative distribution of the response time that correspond to very high percentiles. Validation tests showed that our approach has a very good accuracy in this case.

## Chapter 5

# A Trustworthy Service Model

This chapter considers a set of computer resources used by a service provider to host enterprise applications for customer services subject to an SLA. The SLA defines three QoS metrics, namely, trustworthiness, percentile response time and availability. We present an approach for resource optimization in such an environment that minimizes the total cost of computer resources used by a service provider for such an application while satisfying all these three QoS metrics in a trust-based resource provisioning problem which typically arises in Web services. We formulate the trust-based resource provisioning problem as an optimization problem under SLA constraints, and we solve it using an efficient numerical procedure. Part of materials in this chapter has been published in Xiong and Perros [96].

The rest of the chapter is organized as follows. In Section 5.2 we present a framework for solving the trust-based resource provisioning problem, and give the calculation of SLA metrics in Section 5.3. In Section 5.4 we propose an approach for solving the trust-based resource provisioning problem for single-class and multiple-class customer services respectively. A numerical example is given in Section 5.5 that demonstrates the validity of this approach. We conclude our results in Section 5.6.

## 5.1 The Trust-based Resource Optimization Problem

The main standard in Web services is Extensible Markup Language, XML. XML provides a foundation for many core standards including Web Services Description Language (WSDL), Universal Description, Discovery and Integration specification (UDDI), and Simple Object Access Protocol (SOAP) (Curbera et al. [29]). WSDL allows developers to describe what services are offered, and it helps Web services of e-business to be accessed in public. UDDI defines XML-based registries in which businesses can upload information about themselves and the services they offer. SOAP gives us a way to move XML messages between locations, and it enables programs on separate computers to interact across any network. Thus, Web services allows us to exchange data with other applications on different computers by using Internet protocols.

However, most existing Web services products do not support service level agreements (SLAs) that guarantee a level of service delivered to a customer for a given price. An SLA is a formal contract between a customer and a service provider that defines all aspects of the service being provided. It generally consists of security, performance and availability. *Security* can be categorized as *identity security* and *behavior security*. Identity security includes the authentication and authorization between a customer and a service provider, data confidentiality and data integrity. Behavior security describes the trustworthiness among multiple resource sites, and the trustworthiness of these resource sites by customers, including the trustworthiness of computing results provided by these sites. Performance includes the two following aspects (Menasce [62]).

1. *Response time* is the time for a service request to be satisfied. That is, this is the time it takes for a service request to be executed on the service provider's multiple resource sites.
2. *Throughput* is the service rate that a service provider can offer. It is defined by the maximum throughput or by the undergoing change of throughput with service intensity.

Finally, *availability* is the percentage of time that a service provider can offer services.

In the chapter we consider a resource management problem for Web services under SLA guarantees. Specifically, we define and solve a trust-based resource provisioning problem that occurs in Web services applications, subject to the constraints of trustworthiness, percentile response time and availability. Figure 5.1 depicts a typical scenario for these applications. A customer represents a business that submits a service request consisting of a stream of service jobs at a given rate with a certain price for a given quality of services. After the trust manager who represents the customer negotiates an SLA with a service broker who represents resource sites (alternatively called

service sites)  $S_1, S_2, \dots, S_M$  where  $M > 0$ , it checks the trustworthy information of the resource sites and selects  $m$  ( $0 < m \leq M$ ) of those sites which meet pre-define trustworthy requirements for serving the service request, simply say sites  $1, 2, \dots, m$ . All service managers who manage the selected sites need to work together to achieve the SLA's requirement, and they overlook all resources within their sites. Upon the completion of a service job, the completed result is sent back to the trust manager. The trust manager forwards the completed job to the customer after checking and updating trustworthiness information in its database. Meanwhile, the customer periodically sends feedbacks to the trust manager who uses them to update its trustworthiness information as well.

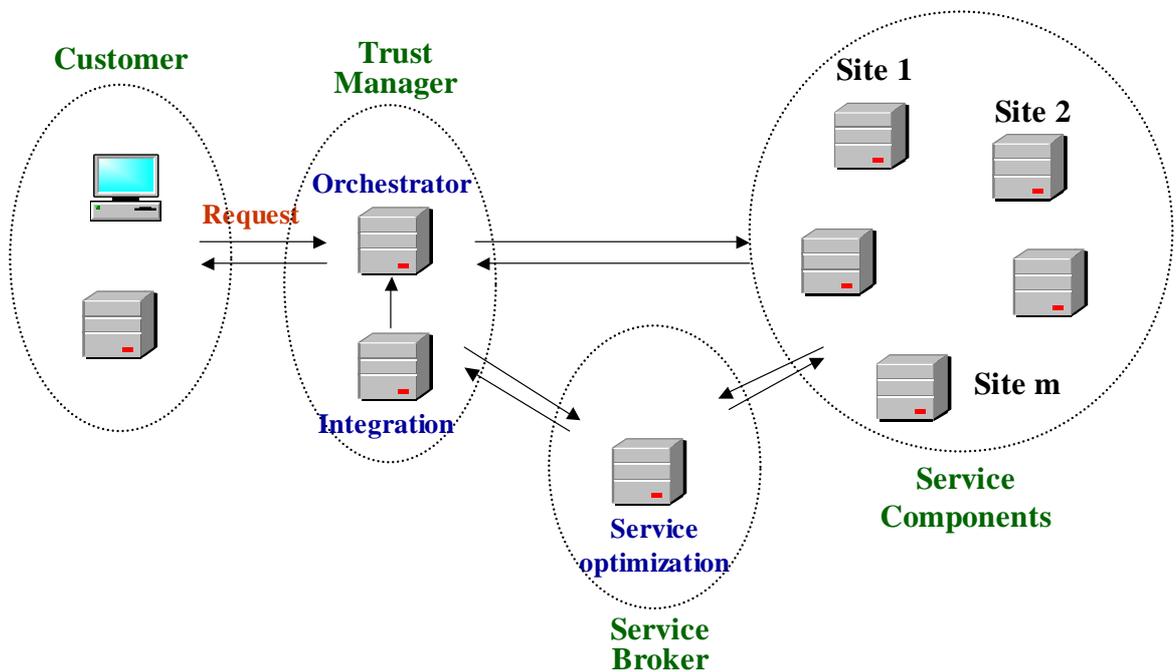


Figure 5.1: An SLA-based Web Services Model

The trust-based resource provisioning problem is to minimize the overall cost of the trusted computing resources required while satisfying SLA requirements. For presentation purpose, we assume that each resource site has only one type of server, each with cost  $c_j$ . Otherwise, if they have multiple types of servers, we can decompose each resource site into several sites so that each one only contains one type of server with the same cost. Let  $N_j$  be the number of servers at site  $j$  ( $j = 1, 2, \dots, M$ ). Thus, the trust-based resource provisioning is quantified by solving for

$n_j$  ( $1 \leq n_j \leq N_j$ ) in the following optimization problem:

$$\min_{n_1, \dots, n_m} (n_1 c_1 + \dots + n_m c_m) \quad (5.1)$$

subject to constraints by an SLA. We discuss these constraints in Section 5.3 and present a framework for the detailed workload of the above Web services scenario in Section 5.2.

In this chapter, we present an approach for the resource optimization that minimizes the total cost of computer resources required while preserving given three QoS metrics for single class and priority-class customers respectively. We calculate the number of servers in each resource station that minimize a cost function that reflects operational costs in the trust-based resource provisioning problem. We analyze an open tandem queueing network. We note that the proposed approach can be also applied to queueing networks consisting of nodes arbitrarily linked.

## 5.2 A Framework for Solving the Trust-based Resource Provisioning

### Problem

Services components from several universal service providers can be flexibly integrated into a composite service with cross-language and cross-platform regardless of their location, platform, execution speed, and process. Delivering quality services to meet customer's requirement under SLA is very important and challenging due to the dynamical and unpredictable nature of Web services applications. In the chapter, we propose a Web services framework for solving the trust-based resource provisioning problem as shown in Figure 5.2. The framework consists of a customer, a trust manager, a service broker and a service processor. The functions of these service entities are explained as follows.

- The customer represents a business that negotiates a contract for particular Web services with a service broker and submits a service request to be processed by a number of resource sites. The customer consists of a number of business end-users (simply called users), and the service request defines service jobs generated by the users at a given rate, and QoS requirements with a fee.
- The trust manager is an entity that plays in the following roles:
  - service integration. According to the customer's request, it integrates services chosen from a number of resource sites.

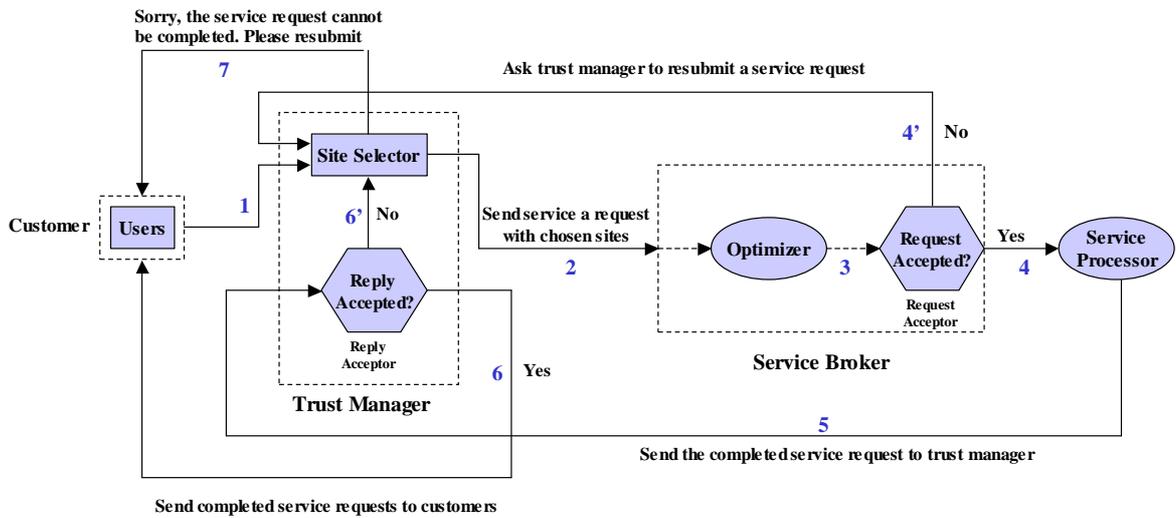


Figure 5.2: A Framework for Solving the Trust-based Resource Provisioning Problem

- trust monitoring and determination. It monitors and determines the trustworthiness of service sites.
- service selection. It chooses service sites that meet trust level requirements predefined by the customer.
- service orchestration. It defines the sequence and conditions in which one Web service invokes other Web services in order to realize some useful function. That is, it defines the pattern of interactions that service components must follow in order to process the customer's requests.

Based the customer request, the trust manager generates a new service request that contains a set of service sites selected by the trust manager.

- The service broker is an entity that represents service sites. After it receives the trust manager's service request, and then the service broker optimizes the number of service resources in each service site to ensure the customer's QoS requirements. The service resource could be a piece of software, hardware, or both and it plays a key role in completing the customer's request. In this research, we consider it as a hardware device, such as a blade, a CPU, a storage device, a disk, a router and so on.
- The service processor is an entity that manages and instructs service sites to process the

customer request based on the trust manager's service request. The service sites are entities that provide specific services. They may be owned by different service providers. Each site consists of a number of service resources managed by a service manager.

Furthermore, we outline the workflow for an approach to solving the trust-based resource provisioning problem in the framework as shown in Figure 5.2. The workflow can be presented in the following steps.

- Step 1.** A customer submit a service request to the trust manager who represents the customer. Let us recall that the customer represents a business consisting a number of users within this business, and the service request consists of a number of service jobs generated by users.
- Step 2.** The site selector who is part of the trust manager selects service sites with the trust indices that meet the customer's trust requirement. Then, the trust manager submits the customer request to the service broker.
- Step 3.** An optimizer who is part of the service broker runs an optimization algorithm to find the number of servers required at each resource site to ensure the customer's SLA guarantee. Then, it will decide whether to accept the service request based on the profitability of the service broker.
- Step 4.** If the service broker can make an acceptable profit with a service provider who owns these chosen service sites, then the service request is accepted. Subsequently, the service request is submitted to a service processor whereby these chosen service sites process the service request based on a rule defined in the trust manager's request. If the service request is not accepted, go to Step 4'. The service processor is an entity that manages service sites.
- Step 5.** After these chosen service sites finish the service request, the service processor sends the completed service request back to the trust manager. Then, the reply acceptor who is part of the trust manager will decide whether to accept the service reply based on the trust indices of these chosen sites at the time when it receives the response from the service processor.
- Step 6.** If the trust indices of these chosen sites meet the customer's requirement, then the service reply is accepted. Subsequently, the completed service request is forwarded to the customer.
- Step 6'.** If the trust indices of these chosen sites do not meet the customer's requirement, then the service reply cannot be accepted. Subsequently, the service acceptor forwards the service request to the site selector and asks it to resubmit the service request.

**Step 4’.** If the service broker cannot make an acceptable profit with a service provider who owns these chosen service sites, then the service request cannot be accepted. Thus, the service broker notifies the trust manager to resubmit the service request by re-selecting service sites.

**Step 7.** After the number of resubmissions is more than a threshold predefined by the trust manager or the service broker, then the trust manager will notify the customer that its request cannot be completed. Then, the customer needs to either abort the service request or modify the QoS requirement with a fee accordingly.

**Remarks:**

- The trust manager and the service broker represent the customer and the chosen service sites respectively. Hence, we assume that they only receive service commissions from the customer and the service processor respectively based on the number of service completions.

In general, the trust manager may be a business entity who not only receive a service commission from the customer, but also make money since it may work with the service broker to find service sites who meet predefined trust requirements but also the chosen service sites have an overall lower fee.

- The trust manager selects service sites on behalf of its customer. Hence, it does not care about how much cost is needed to process the customer’s request. But, if the number of service resubmissions surpasses a threshold predefined by the trust manager or the service broker, then the trust manager will notify the customer that her/his request cannot be completed, as mentioned in Step 7.
- If the trust manager constantly resubmits a service request to the service broker, then the service request submissions constantly consume the service broker’s resource. Thus, a denial of service attack may occur when an attacker impersonates the trust manager to keep submitting bogus service requests to the service broker. We do not analyze such an attack here since it is beyond the scope of our study in this research.
- The above steps can be regarded as an SLA negotiation process between the trust manager and the service broker. The process is often finished before service delivery. That is, the SLA is static. However, the SLA can be dynamically changed as well. In this case, the SLA will be periodically verified and/or negotiated. The length of the period of time for the verification and/or negotiation depends on individual service types. It may be a week, or a month.

### 5.3 The Calculation of SLA Metrics

Before presenting an algorithm to solving the above trust-based resource provisioning problem, we need to qualitatively analyze three SLA metrics: trustworthiness, percentile response time and service availability. In this section, we provide the calculation of three SLA metrics.

#### 5.3.1 The Trustworthiness of Resource Sites

In Section 5.1 we classified security into *identity security* and *behavior security*, similar to identity trust and behavior trust defined in a grid computing system (Azzedin et al. [5]). Identity security includes the authentication and authorization between a customer and a service provider, data confidentiality and data integrity. It has been widely studied in the literature (e.g., see Bishop [10], Kim et al. [51], Perrig et al. [75] and J. Rosenberg and D. Remy [82]). The identity security is beyond the scope of our study in this chapter. We shall discuss group authentication in Chapter 6.

In this study, “trust” is used to deal with the notion of the trustworthiness in behavior security. Its definition is varied in the literature. In this research, trust is a firm belief in the competence of a resource site to act as expected. The trustworthiness of resource sites is an indicator of the quality of services provided by these sites based on previous and current job completion experience. It often predicates the future behavior of the quality of services at these sites. Moreover, we are only interested in the trustworthiness of resource sites from a customer’s perspective. Hence, in this study we simply assume that these resource sites trust each other.

Furthermore, assume that the trust manager is a trusted agent who represents customers. The trust manager uses the collected trustworthy information regarding the resource sites to evaluate their security behavior. We consider security behavior by modeling the behavior trusts of all sites, and quantify the trustworthiness of these sites using a rank and a threshold-based approach. This approach is based on previous job completion experience assessed by the trust manager and customers. The assessment may also include the opinion of other trust managers besides its own customer’s feedback. The domain of feedback is assumed to be  $[0, 1]$ .

The feedback is a statement issued by the trust manager’s customers about the quality of service provided by those chosen service sites in each service request or for a set of service requests during a certain period of time. These customers only provide feedbacks about their received services rather than those chosen service sites. (Note that the trust manager’s customer is usually not aware of which service sites process its service request.) Then, the trust manager scores the trust

indices of those chosen service sites that serve these customers' requests. In the meanwhile, the trust manager may have its own feedback. In this case, the trust indices of those chosen service sites will aggregate both the trust manager's feedback and its customers' feedbacks.

The opinion is defined as the information of trustworthiness provided by those trust managers who are *neighbors* of the trust manager. These neighbors are a small subset of the trust manager's acquaintances, adaptively selected based on their usefulness. The opinion aggregates the overall impression of those neighborhood trust managers for the service sites.

We first discuss the case in which the trust indices of service sites are only determined by the trust manager's and its customers' feedbacks. Let us consider the discrete times  $t_1, t_2, \dots, t_k, \dots$  in an increasing order ( $k = 1, 2, \dots$ ). Let  $I_j^{t_k}$  be the trust index of site  $j$  at time  $t_k$ . Then, a *trust function* is defined by

$$I_j^{t_{k+1}} = \xi \frac{R_j^s(k+1)}{R_j^a(k+1)} + (1 - \xi)I_j^{t_k} \quad (5.2)$$

for each site  $j$  ( $j = 1, 2, \dots, M$ ), where  $R_j^s(k+1)$  is the number of service jobs completed at site  $j$  that satisfied customers provided that the trust manager itself does not assess each completed job, or both the trust manager and its customers provided that both of them assess each completed job.  $R_j^a(k+1)$  is the total number of service jobs submitted to site  $j$  during the time period  $[t_k, t_{k+1}]$ . The trust manager's satisfaction is assessed by the validation of a customer's SLA requirement after a service is completed, and a customer's satisfaction is based on the customer's feedback assessed by the customer through a comparison of the job completion experience with its expectation.

For presentation purpose, we simply assume that the trust manager has the same feedback as its customers. (Otherwise, the following discussion and notation needs to be adjusted accordingly which can be easily done.) Denote  $r_j(k+1)$  by

$$r_j(k+1) = \frac{R_j^s(k+1)}{R_j^a(k+1)} \quad (5.3)$$

Moreover,  $r_j(k+1)$  is the satisfactory rate at site  $j$  from time  $t_k$  to  $t_{k+1}$  and it obviously ranges from 0 to 1. Clearly, when a set of the chosen sites is unchanged during this period of time, the number of completion jobs is the same for all chosen sites. Thereby,  $r_j(k+1)$  is the same for these chosen sites as well.  $\xi$  is a parameter determined by the trust manager. It is chosen depending on the type of resource sites. If a critical service job is processed at site  $j$  and site  $j$ 's security is sensitive with the change of time, then  $\xi$  should be close to 1 (e.g., bigger than 0.7). Otherwise, it should be close to 0 (e.g., smaller than 0.3). Thus,  $I_j^{t_k}$  ranges from 0 to 1, and it is called a *trust index*. It gives

the percent of a time that a resource site completes jobs to the satisfaction of the trust manager and its customers.

Then, we discuss the case in which the trust indices of service sites are determined by not only the trust manager's and its customers' feedbacks, but also the opinions of other trust managers. In today's large scale complex computer system, a central trust manager would not be able to represent all customers, and would not know or would not be able to keep up with the trustworthiness information of all service sites in the system. This implies that there would exist multiple or many trust managers in the large scale complex computer system. In this case, the trust index function in (5.2) should be modified by considering the opinions of the trust manager's neighbors besides its own customer's feedback. Thus, the satisfactory rate  $r_j(k+1)$  consists of the following two components: the first component is based on the feedbacks of the trust manager and its own customers that is denoted by  $r_j^b(k+1)$ . Thus,  $r_j^b(k+1)$  is given by (5.3), i.e.,

$$r_j^b(k+1) = \frac{R_j^s(k+1)}{R_j^a(k+1)}$$

The second component is determined by aggregating all the opinions of the trust manager's neighbors and it is denoted by  $r_j^o(k+1)$ . Moreover, the satisfactory rate  $r_j(k+1)$  is given by

$$r_j(k+1) = \xi^b \times r_j^b(k+1) + \xi^o \times r_j^o(k+1)$$

where  $\xi^b + \xi^o = 1$ . The two parameters are used to adjust the balance between the feedbacks and the opinion, and they are determined by the trust manager. Subsequently, the trust index function in (5.2) is rewritten as

$$\begin{aligned} I_j^{t_{k+1}} &= \xi r_j(k+1) + (1-\xi)I_j^{t_k} \\ &= \xi [\xi^b \times r_j^b(k+1) \\ &\quad + \xi^o \times r_j^o(k+1)] + (1-\xi)I_j^{t_k} \end{aligned}$$

The trust function describes the trustworthiness of resource sites monitored and updated by the trust manager. Therefore, the trust function reflects a probabilistic security behavior of the resource sites from a customer's perspective.

### 5.3.2 The Percentile Response Time

The SLA performance metric defined in Section 5.1 includes *throughput* and *response time*. As an end user, a customer is in general concerned about response time rather than throughput. So, in the study we only consider *the response time*. In the trust-based resource provisioning

problem for Web services, a response time is the time it takes for a job to be executed and completed in a distributed computing environment consisting of a trust manager and multiple resource sites.

In the literature, typically the average response time (or an average execution time) is used (e.g., see Menasce [62] and [63]). The average response time is heavily influenced by “outliers,” which occur in almost all measurements. Therefore, although the average response time is relatively easy to calculate, it may not address the concerns of a customer. Typically, a customer is more inclined to request a statistical bound on its response time than an average response time. For instance, a customer can request that at least 95% of the time its response time should be less than a desired value. Hence, in this research we are concerned with finding an optimal resource provisioning from trusted resource sites that meets a desired percentile response time.

Let us recall that  $f_T(t)$  is the probability distribution function of a response time  $T$  of a customer as defined in Section 3.1 for the case of single class customers, and a certain priority class customer as defined in Section 4.1 for the case of multiple class customers. For example, in the case of two priority classes, for the high-priority class  $T = T^{(1)}$  and for the low-priority class  $T = T^{(2)}$ .  $T_D^{(r)}$  is a desired target response time for priority class  $r$  ( $r = 1, 2$ ) that a customer requests and agrees upon with its service provider based on a fee paid by the customer. Then, the SLA performance metric used in this chapter is:

$$\int_0^{T_D} f_{T(t)}(t) dt \geq \gamma\% \quad (5.4)$$

for the case of single class customers. That is,  $\gamma\%$  of the time a customer will receive its service in less than  $T_D$ , where  $T_D$  is a desired target response time that a customer requests and agrees upon with its service provider based on a fee paid by the customer. And, the SLA performance metric used in this chapter is:

$$\int_0^{T_D^{(r)}} f_{T^{(r)}}(t) dt \geq \gamma^{(r)}\%, \quad (r = 1, 2) \quad (5.5)$$

for the case of multiple class customers. That is,  $\gamma^{(r)}\%$  of the time a customer will receive its service in less than  $T_D^{(r)}$  ( $r = 1, 2$ ).

As an example given in Section 4.1, when considering an  $M/M/1$  queue with an arrival rate  $\lambda^{(r)}$  and a service rate  $\mu^{(r)}$  ( $r = 1, 2$ ). The service discipline is preemptive-resume. we have that the probability distribution of the high-priority response time is given by

$$f_{T^{(1)}}(t) = \mu^{(1)}(1 - \rho^{(1)})e^{-\mu^{(1)}(1-\rho^{(1)})t}. \quad (5.6)$$

### 5.3.3 The Service Availability

Availability is a critical metric in today's computer design (Hennessy [45]). Brown and Patterson [16] used an availability metric to describe the variations in system quality of service over time. It is defined by the latency of a request service and the number of failures that can be tolerated by a system. The former was discussed in Section 5.3.2. So, we only need to study the latter in this section. We consider the percentage of time that a resource is "up" or "down" as a metric, which is the traditional way to define service availability.

Let  $MMTF_j$  (Mean Time to Failure) be the average time of a server failure, and  $MTTR_j$  (Mean Time to Recover) be the average time for recovering a server at the resource site  $j$ . Then,  $MTBF_j$  (Mean Time Between Failures) is a sum of  $MMTF_j$  and  $MTTR_j$ .  $MTBF$  information can be obtained from a hardware provider. For example, it is mentioned in Cisco [26] that  $MTBF$  information is available for all Cisco components and is available upon request to a local account manager.  $MTTR$  is determined by evaluating how quickly a site owner can repair broken servers. It is a major factor of server availability. To improve service availability, it is necessary to reduce the frequency time of failure, as indicated in Brewer [15].

Network availability data may be also found over the Internet. For example, the University of Houston maintains current and historical network availability, see [60].

Furthermore, each server fails at a rate of  $\frac{1}{MMTF_j}$ , denoted by  $a_j$ , and recovers (i.e., it is put back into operation) at a rate of  $\frac{1}{MMTR_j}$ , denoted by  $b_j$  ( $j = 1, \dots, m$ ). Thus, a two-state Markov chain with the states "up" and "down" can be used to study the service availability at site  $j$ . The failure rate  $a_j$  is the state of transition from "up" to "down," and the recovery rate  $b_j$  represents the rate of transition from "down" to "up." Then, the probability  $p_j^i$  that  $i$  servers are down is given by (see Bloch et al. [12])

$$p_j^i = \frac{N_j!}{i!(N_j - i)!} \eta_j^i p_j^0, \quad \text{for } i = 1, \dots, N_j \quad (5.7)$$

where  $\eta_j = \frac{a_j}{a_j + b_j}$  is the server unavailability rate, and  $p_j^0$  is given by

$$p_j^0 = [N_j! \sum_{i=0}^{N_j} \frac{\eta_j^i}{i!(N_j - i)!}]^{-1} \quad (5.8)$$

The probability that no more than  $N_j - n_j$  servers at site  $j$  are down is:

$$P_j(n_j, N_j) = p_j^0 \sum_{i=0}^{N_j - n_j} \frac{N_j!}{i!(N_j - i)!} \eta_j^i \quad (5.9)$$

(5.9) can also be seen as the probability that at least  $n_j$  servers in site  $j$  are available.

## 5.4 An Approach for Solving the Trust-based Resource Provisioning

### Problem

As we saw in Section 5.2, a Web services framework consists of a customer, a trust manager, a service broker, and a service processor. In this framework, an optimizer within the service broker is an important key components. It calculates the number of service resources required to ensure that the response time of a service request meets the requirement of a predefined percentile response time under a given fee. The calculation of the response time is a sum of the processing time of the trust manager and the execution time of chosen service sites for the customer request.

Without any confusion, we reuse  $m$  ( $0 < m \leq M$ ) as the number of resource sites necessary for processing a customer's service job. Assume that the trust manager first selects  $m$  resource sites from the  $M$  ones. We consider the following two cases: single customer services and multiple priority customer services respectively.

#### 5.4.1 Single Class Customers

We can model those  $m$  resource sites as a queueing network as shown in Figure 5.3. Without loss of generality, we assume that the first  $m$  sites are chosen. In Figure 5.3, the tandem queueing network consists of a trust server and  $m$  stations numbered sequentially from 1 to  $m$ . The trust server represents the trust manager, and each station represents a resource site. Each station carries out a particular function, such as a database server, a computing server, a file server, or a web server. The execution procedure of a service request may be complicated in the real world, but the main idea of our proposed modelling approach can be extended to describe any collaborative relationship among resource sites as long as the relationship can be quantified.

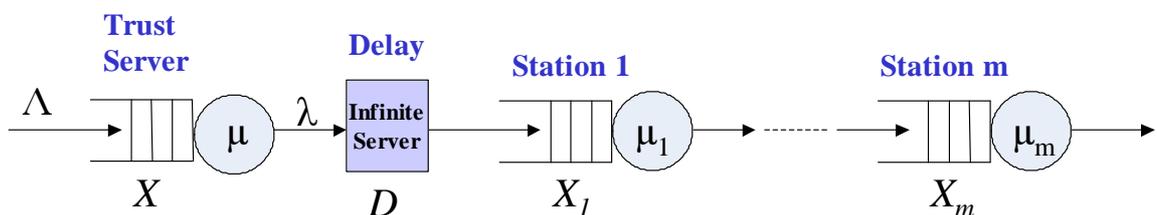


Figure 5.3: A Web Services Computing System for Single-class Customer Services

In Figure 5.3, each station  $j$  is modeled as a single FIFO queue served by  $n_j$  identical servers, each providing a service at the rate  $\mu_j$ . Let  $\Lambda$  be the external arrival rate to the infinite server, and let  $\lambda$  and  $\lambda_j$  be the effective arrival rates to the infinite server and station  $j$  ( $j = 1, 2, \dots, m$ ). The notion of server here is defined as a service resource at each site that processes users' jobs. For example, as we mentioned early, it could be a blade, a CPU, a disk, a storage device, and so on. We assume that all service times are exponentially distributed and the external arrival to the trust server occurs in a Poisson fashion. The trust server provides a service at the rate  $\mu$ .

In Figure 5.3, the infinite server represents the total propagation delay from the user to the service provider and back and also from station 1 to  $m$ . We only consider a single class of customer in this research.

In this research we are interested in minimizing the overall cost of the above Web services computing system so that the desired SLA is guaranteed. Before presenting an algorithm for solving the minimization problem, we first need to calculate the response time of a customer's service request. Recall that the response time refers to the time elapsed from the moment a customer's service job joins the computing system to the time it departs. It is the most important QoS metric. The response time also reflects security behavior and the availability of services in some degree.

Let  $T$  be a random variable that represents the response time of a service job. Also, assume that  $f_T(t)$  is the probability distribution of  $T$ . Denote the overall service cost (1.1) by

$$g(n_1, n_2, \dots, n_m) = \sum_{j=1}^m n_j c_j \quad (5.10)$$

where  $n_1, n_2, \dots, n_m$  are the number of servers allocated to a specific stream of jobs with arrival rate  $\lambda$ .

Then, the trust-based resource provisioning problem can be formulated as the following three sub-problems:

1. Select  $m$  resource sites within a predefined trust index at time  $t = t_k$ .
2. Solve for  $n_j$  in the  $m$ -dimensional integer optimization problem:

$$\min_{n_1, \dots, n_m} g(n_1, n_2, \dots, n_m) \quad (5.11)$$

Under the constraint of a percentile response time of (3.1), and the constraint of service availability:

$$P_j(n_j, N_j) \geq \delta_j \% \quad (5.12)$$

where  $T^D$  is a desired response time defined by a customer, and  $\delta_j$  is a desired percentage of service availability at site  $j$ .  $P_j(n_j, N_j)$  is given by (5.9).

3. Update the trust indices of all  $M$  sites based on the activity during the time interval  $[t_k, t_k + T^D]$ . Then, the trust manager decides if a completed job is accepted. If each trust index at those selected sites who complete the service job meet a predefined index value, then the completed job is accepted. Otherwise, then the completed job is discarded and the trust manager needs to resubmit the job.

Note that in the model the verifying time of a trust index is ignored since we are interested in the processing time of a job at resource sites. While Sub-problems 1 and 3 are solved at the customer side to ensure a customer service received from reliable resource sites, Sub-problem 2 is solved by the service broker who represents these resource sites. Hence, in the current model the trustworthiness of resource sites in Sub-problems 1 and 3 are not considered as a constraint for the optimization problem in Sub-problem 2.

The above resource optimization problem as described in Sub-problem 2 of Section 3.3 is to find a minimal cost presented in (5.11) such that  $\gamma\%$  of the time a customer's response time is less than a predefined value  $T^D$ , and  $\delta_j\%$  of time the selected resource sites have  $n_j$  servers available for the job.

In the following discussion each server station is modelled as a single  $M/M/1$  queue with arrival rate  $\lambda_j$  and service rate  $\psi_j\mu_j$ , where  $\psi_j \in [1, n_j]$  is a function of  $n_j$  to be determined by the configuration of servers at each station ( $j = 1, 2, \dots, m$ ). It is non-decreasing and can be inverted, i.e.,  $\psi^{-1}$  exists. As discussed in Section 3.3, (1) suppose that a domain represents a group of CPUs in a model as discussed in [57]. Then,  $\psi(n)$  can be seen as a CPU scaling factor for the number of CPUs from 1 to  $n$ , and then  $\psi(n)$  can be expressed as  $\psi(n) = \xi^{\log_2 n}$ , where  $\xi$  is a basic scaling factor from 1 CPU to 2. So,  $\psi^{-1}(n) = \xi^{-\log_2 n}$ . (2) Suppose that a domain represents a group of routers in a network model. Then,  $\psi(n) = 1$  when these  $n$  routers are serially placed, and  $\psi(n) = n$  when they are in parallel.

Since the queueing network is overtake-free, the waiting time of a customer at a station is independent of its waiting times at other stations (see [30] and [90]). Let  $D$  be the service time at the infinite server,  $X$  be the time elapsed from the moment a customer arriving at the trust manager server to the moment it departs from the server, and  $X_j$  be the time elapsed from the moment a customer arriving at station  $j$  to the moment it departs from the station. Then, similar to Chapter 3,

the total response time can be calculated by

$$T = X + D + X_1 + X_2 + \cdots + X_m,$$

and hence the LST (Laplace-Stieltjes transform) of the response time  $T$  is

$$L_T(s) = L_X(s)L_D(s)L_{X_1}(s) \cdots L_{X_m}(s) \quad (5.13)$$

where  $L_X(s)$  and  $L_D(s)$  are the LST of the service time  $D$  and  $X$  respectively given by

$$L_X(s) = \frac{\mu(1-\rho)}{s + \mu(1-\rho)}, \quad L_D(s) = \frac{\lambda}{s + \lambda}, \quad (5.14)$$

and  $L_{X_j}(s)$  is the LST of the response time  $X_j$  at the  $j$ -th station given by

$$L_{X_j}(s) = \frac{\psi(n_j)\mu_j(1-\rho_j)}{s + \psi(n_j)\mu_j(1-\rho_j)}, \quad (5.15)$$

where  $\rho = \frac{\lambda}{\mu}$  and  $\rho_j = \frac{\lambda_j}{\psi(n_j)\mu_j}$  ( $j = 1, 2, \dots, m$ ).

From (5.13), (5.14) and (5.15) we have that

$$L_T(s) = \frac{\mu(1-\rho)}{s + \mu(1-\rho)} \cdot \frac{\lambda}{s + \lambda} \times \prod_{j=1}^m \frac{\psi(n_j)\mu_j(1-\rho_j)}{s + \psi(n_j)\mu_j(1-\rho_j)}$$

Then, we can obtain the probability and cumulative distributions of the response time by solving the following:

$$f_T(t) = L^{-1}\{L_T(s)\} \quad \text{and} \quad F_T(t) = L^{-1}\left\{\frac{L_T(s)}{s}\right\}$$

We observe that  $f_T(t)$  and  $F_T(t)$  are usually nonlinear functions of  $t$  and  $n_j$ . Hence, the resource optimization problem is an  $m$ -dimensional linear optimization problem subject to nonlinear constraints. In general, it is not easy to solve this problem. However, the complexity of the problem can be significantly reduced by requiring that the service rates of the queues in the service model in Figure 5.3 are all equal. That is, we find the optimum value of  $n_1, \dots, n_m$  such that

$$\psi(n_1)\mu_1 = \cdots = \psi(n_m)\mu_m$$

From the traffic equations:

$$\lambda = \lambda_j = \Lambda$$

for  $j = 1, 2, \dots, m$ , we have that the utilization of each station

$$\rho_j = \frac{\lambda_j}{\psi(n_j)\mu_j} = \frac{\Lambda}{\psi(n_j)\mu_j}$$

Thus we have

$$\hat{a}_i = \psi(n_i)\mu_i(1 - \rho_i) = \psi(n_j)\mu_j(1 - \rho_j) = \hat{a}_j \triangleq \hat{a},$$

which implies  $n_j = \psi^{-1}\left(\frac{\hat{a} + \lambda_j}{\mu_j}\right)$  ( $i, j = 1, 2, \dots, m$ ). Hence, from (5.13) we have

$$f_T(t) = L^{-1}\left\{\frac{\mu(1 - \rho)}{s + \mu(1 - \rho)} \cdot \frac{\lambda}{s + \lambda} \cdot \frac{\hat{a}^m}{(s + \hat{a})^m}\right\},$$

and subsequently obtain that

$$F_T(t) = L^{-1}\left\{\frac{\mu(1 - \rho)}{s + \mu(1 - \rho)} \cdot \frac{\lambda}{s(s + \lambda)} \cdot \frac{\hat{a}^m}{(s + \hat{a})^m}\right\} \quad (5.16)$$

Consequently,  $\sum_{j=1}^m n_j c_j$  reduces to a function of variable  $\hat{a}$  due to  $n_j = \lceil \psi^{-1}\left(\frac{\hat{a} + \lambda_j}{\mu_j}\right) \rceil$ .

Thus we have the following algorithm for the resource optimization problem.

Moreover, the constraint of service availability at each resource site in (5.12) is rewritten as

$$G_j(\hat{a}) \stackrel{def}{=} P_j(\lceil \psi^{-1}\left(\frac{\hat{a} + \lambda_j}{\mu_j}\right) \rceil, N_j) \quad (5.17)$$

At this time,  $g(n_1, \dots, n_m)$  in (5.10) reduces to a function of one variable. Thus the trust-based resource provisioning problem is solved using the following iterative algorithm.

#### Algorithm 5.4.1-1

1. At time  $t_k$ , select  $m$  resource sites within a predefined trust index  $\hat{I}_j$ , and the highest  $m$  trust indices. If such a selection is impossible, then the trust manager informs the customer “We cannot process the service request at this moment.” After the customer request is waited for more than a given threshold time, the trust manager still cannot find those service sites that meet the predefined trust index  $\hat{I}_j$ . Then, the trust manager informs the customer “You need to either revise the trust requirement or abort the request and resubmit it later.” Otherwise, continue to do Step 2.
2. Find  $\hat{a}$  in the following two one-dimensional optimization problems.
  - (a) The minimization problem of a percentile response time:

$$\hat{a}^{(1)} \leftarrow \arg \min_{\hat{a}} F_T(t)|_{t=T^D}$$

subject to the constraint  $F_T(t)|_{t=T^D} \geq \gamma\%$  at  $\hat{a} = \hat{a}^{(1)}$ , where  $F_T(t)$  is given by (5.16).

(b) The minimization problem of service availability:

$$\hat{a}_j^{(2)} \leftarrow \arg \min_{\hat{a}} G_j(\hat{a})$$

subject to the constraint  $G_j(\hat{a}^{(2)}) \geq \delta_j\%$ , where  $G(\hat{a})$  is given by (5.17).

3. Compute integers  $n_j$  by using  $n_j = \lceil \frac{\hat{a}_j^M}{\mu_j(1-\rho_j)} \rceil$ , where  $\hat{a}_j^M = \max(\hat{a}^{(1)}, \hat{a}_j^{(2)})$  for  $1 \leq n_j \leq N_j$  and  $j = 1, 2, \dots, m$ .
4. Update  $I_j^{t_k}$  to  $I_j^{t_k+T^D}$  based on the trust function (5.2). If  $|I_j^{t_k+T^D}| \geq \hat{I}_j$ , then the completed service job is accepted. Otherwise, repeat Steps 1-3.
5. When the repeated number is more than a predefined threshold, the trust manager notifies the customer “We need to re-negotiate an SLA with the service broker.”

As mentioned early, Algorithm 5.4.1-1 is an iteration algorithm. The number of iterations will be determined by the trust manager and the service broker. When the number of iterations is more than a predefined threshold, the trust manager will notify the customer to modify the SLA. The customer may abort the service request or re-submit the service request with a new QoS requirement with a new fee. Additionally, the trust manager selects those service sites with the highest trust indices who also meet predefined trust requirements in Step 1 in Algorithm 5.4.1-1. But, for its own benefit, the trust manager can choose those service sites who only meet trust requirements predefined by a customer.

Moreover, Algorithm 5.4.1-1 shows that the  $m$ -dimensional optimization problem in Subproblem 2 significantly reduces to an one-dimensional optimization problem. Surprisingly, the optimal number of servers obtained by Algorithm 5-1 is independent of server costs  $c_j$ , due to  $\hat{a}_i = \hat{a}_j$  ( $i, j = 1, 2, \dots, m$ ). In general, each of the above two one-dimensional optimization problems can be solved numerically.

Note that the service broker cannot provide the quality of service at the moment in which an SLA negotiation is requested. In this case, the service broker may get a service penalty. For presentation purpose, it will be not further discussed here since the above approach can be easily adjusted to deal with this case. Additionally, we have only considered a single-class customer. The above results can be extended to the case of two-class customers in next section.

### 5.4.2 Multiple Priority Customers

We can model those  $m$  resource sites as a queueing network as shown in Figure 5.4. Without loss of generality, we assume that the first  $m$  sites are chosen. In Figure 5.4, the tandem queueing network consists of a trust server and  $m$  stations numbered sequentially from 1 to  $m$ . The trust server represents the trust manager, and each station represents a resource site. Each station carries out a particular function, such as, it could be a database server, a computing server, a file server, and a web server. The execution procedure of a service request may be complicated in the real world, but the main idea of our proposed modeling approach can be still extended to describe any collaborative relationship among resource sites as long as the relationship can be quantified.

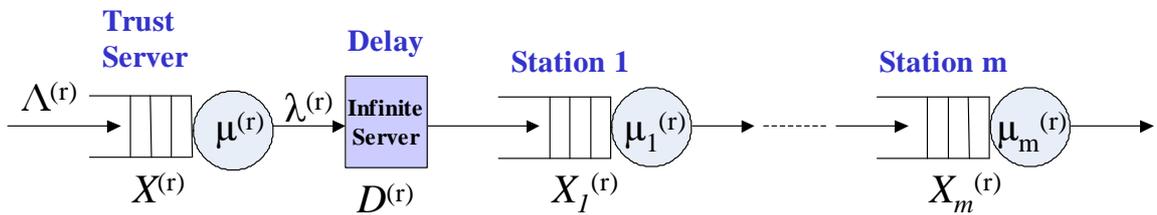


Figure 5.4: A Web Services Computing System for Multiple-class Customer Services

In Figure 5.4, each station  $j$  is modeled as a single FIFO queue served by  $n_j$  identical servers, each providing a service at the rate  $\mu_j^{(r)}$ , where  $r = 1, \dots, R$ . Let  $\Lambda^{(r)}$  be the external arrival rate to the infinite server, and let  $\lambda^{(r)}$  and  $\lambda_j^{(r)}$  be the effective arrival rates to the infinite server and station  $j$ ,  $j = 1, 2, \dots, m$ . We assume that all service times are exponentially distributed and the external arrival to the infinite server occurs in a Poisson fashion. The trust server provides a service at the rate  $\mu^{(r)}$ .

In Figure 5.4, the infinite server represents the total propagation delay from the user to the service provider and back and also from station 1 to  $m$ . We only consider two classes of customers in this research. That is,  $R = 2$ .

In this research we are interested in minimizing the overall cost of the above Web services computing system so that the desired SLA is guaranteed. Before presenting an algorithm for solving the minimization problem, we first need to calculate the response time of a customer's service request. Recall that the response time refers to the time elapsed from the moment a customer's service job joins the computing system to the time it departs. It is the most important QoS metric. The response time also reflects security behavior and the availability of services in some degree.

Let  $T^{(r)}$  be a random variable that represents the response time of a service job with class  $r$  customers. Also, assume that  $f_T^{(r)}(t)$  is the probability distribution of  $T^{(r)}$ . Denote the overall service cost (5.1) by

$$g(n_1^{(r)}, n_2^{(r)}, \dots, n_m^{(r)}) = \sum_{j=1}^m n_j^{(r)} c_j \quad (5.18)$$

where  $n_1^{(r)}, n_2^{(r)}, \dots, n_m^{(r)}$  are the number of servers allocated to a specific stream of jobs with arrival rate  $\Lambda^{(r)}$  ( $r = 1, 2$ ).

Then, the trust-based resource provisioning problem can be formulated as the following three sub-problems:

1. Select  $m$  resource sites within the predefined trust indices of all class customers at time  $t = t_k$ .
2. Solve for  $n_j^{(r)}$  in the  $m$ -dimensional integer optimization problem:

$$\min_{n_1^{(r)}, \dots, n_m^{(r)}} g(n_1^{(r)}, n_2^{(r)}, \dots, n_m^{(r)}) \quad (5.19)$$

Under the constraint of a percentile response time of (5.5), and the constraint of service availability:

$$P_j(n_j^{(r)}, N_j) \geq \zeta_j^{(r)}\% \quad (5.20)$$

where  $T_D^{(r)}$  is desired response time defined by class  $r$  customers, and  $\zeta_j^{(r)}\%$  is a desired percentage of service availability at site  $j$ .  $P_j(n_j^{(r)}, N_j)$  is given by (5.9). Then, calculate that

$$n_j = \max\{n_j^{(1)}, n_j^{(2)}\} \quad (j = 1, 2, \dots, m)$$

3. Update the trust indices of all  $M$  sites based on the activity during the time interval  $[t_k, t_k + T^D]$ . Then, the trust manager decides if a completed job is accepted. If each trust index at those selected sites who complete the service job meet a predefined index value, then the completed job is accepted. Otherwise, then the completed job is discarded and the trust manager needs to resubmit the job.

Note that in the model the verifying time of a trust index is ignored since we are interested in the processing time of a job at resource sites. While Sub-problems 1 and 3 are solved at the customer side to ensure a customer service received from reliable resource sites, Sub-problem 2 is solved by the service broker who represents these resource sites. Hence, in the current model the trustworthiness of resource sites in Sub-problems 1 and 3 are not considered as a constraint for the optimization problem in Sub-problem 2.

The above resource optimization problem as described in Sub-problem 2 of Section 5.4 is to find a minimal cost presented in (5.19) such that  $\gamma^{(r)}\%$  of the time a customer's response time is less than a predefined value  $T_D^{(r)}$ , and  $\zeta_j^{(r)}\%$  of time the selected resource sites have  $n_j^{(r)}$  servers available for the job of class  $r$  ( $r = 1, 2$ ).

In the following discussion each station is modeled as a single  $M/M/1$  priority queue with arrival rate  $\lambda_j^{(r)}$  and service rate  $\psi^{(r)}(n_j)\mu_j^{(r)}$ , where  $\psi^{(r)}(n_j)$  is a known function of  $n_j$  ( $r = 1, 2$ ), and depends on the configuration of servers and the type of customers at each station. It is non-decreasing and can be inverted, i.e.,  $(\psi^{(r)})^{-1}$  exists ( $r = 1, 2$ ). The detailed explanation of  $(\psi^{(r)})^{-1}$  is referred to Section 4.3.2.

Since the queueing network is overtake-free, the waiting time of a customer at a station is independent of its waiting times at other stations (see [30] and [90]). Let  $D^{(r)}$  be the service time at the infinite server,  $X^{(r)}$  be the time elapsed from the moment a customer of class  $r$  arriving at the trust manager server to the moment it departs from the server, and  $X_j^{(r)}$  be the time elapsed from the moment a customer of class  $r$  arriving at station  $j$  to the moment it departs from the station. Then, the total response time is

$$T^{(r)} = X^{(r)} + D^{(r)} + X_1^{(r)} + X_2^{(r)} + \dots + X_m^{(r)},$$

for the customer of class  $r$  and hence the LST (Laplace-Stieltjes transform) of the response time  $T^{(r)}$  is

$$L_{T^{(r)}}(s) = L_{X^{(r)}}(s)L_{D^{(r)}}(s)L_{X_1^{(r)}}(s) \cdots L_{X_m^{(r)}}(s) \quad (5.21)$$

where  $L_{D^{(r)}}(s)$  is the LST of the service time  $D^{(r)}$  given by

$$L_{D^{(r)}}(s) = \frac{\lambda^{(r)}}{s + \lambda^{(r)}} \quad (5.22)$$

Also,  $L_{X^{(r)}}(s)$  is the LST of the response times  $X^{(r)}$  at the trust server, and  $L_{X_j^{(r)}}(s)$  is the LST of the response time  $X_j^{(r)}$  at the  $j$ -th station, where  $r = 1, 2$ .

Due to the preemptive-resume priority, the high-priority response time is that same as in the single class FIFO  $M/M/1$  whose probability distribution can be expressed as (5.6) in section 5.2. Thus  $L_{X^{(1)}}(s)$  is determined by

$$L_{X^{(1)}}(s) = \frac{\mu(1 - \rho^{(1)})}{s + \mu(1 - \rho^{(1)})}, \quad (j = 1, 2, \dots, m) \quad (5.23)$$

and  $L_{X^{(2)}}(s)$  is the LST of the low-priority response time given by

$$L_{X^{(2)}}(s) = \frac{(1 - \rho)\delta^{(1)}}{1 - \rho\delta^{(1)}}, \quad (j = 1, 2, \dots, m) \quad (5.24)$$

due to (25) in [94], where  $\mu \triangleq \mu^{(1)} = \mu^{(2)}$ ,  $\rho^{(r)} = \frac{\lambda^{(r)}}{\mu}$ ,  $\rho = \frac{\Lambda^{(1)} + \Lambda^{(2)}}{\mu}$ , and  $\delta^{(1)}$  is given by

$$\delta^{(1)} = \frac{\eta - \sqrt{\eta^2 - 4\Lambda^{(1)}\mu^{(1)}}}{2\Lambda^{(1)}}, \quad \text{and} \quad \eta = s + \Lambda^{(1)} + \mu^{(1)}$$

Furthermore,  $L_{X_j^{(1)}}(s)$  is given by

$$L_{X_j^{(1)}}(s) = \frac{\psi^{(1)}(n_j)\mu_j(1 - \rho_j^{(1)})}{s + \psi^{(1)}(n_j)\mu_j(1 - \rho_j^{(1)})}, \quad (j = 1, 2, \dots, m) \quad (5.25)$$

and  $L_{X_j^{(2)}}(s)$  is the LST of the low-priority response time given by

$$L_{X_j^{(2)}}(s) = \frac{(1 - \rho_j)\delta_j^{(1)}}{1 - \rho_j\delta_j^{(1)}}, \quad (j = 1, 2, \dots, m) \quad (5.26)$$

due to (25) in [94], where  $\rho_j^{(r)} = \frac{\lambda_j^{(r)}}{\psi^{(r)}(n_j)\mu_j}$ ,  $\rho_j = \frac{\lambda_j^{(1)} + \lambda_j^{(2)}}{\psi^{(r)}(n_j)\mu_j}$ , and  $\delta_j^{(1)}$  is given by

$$\delta_j^{(1)} = \frac{\eta_j - \sqrt{\eta_j^2 - 4\psi^{(1)}(n_j)\lambda_j^{(1)}\mu_j^{(1)}}}{2\psi^{(1)}(n_j)\lambda_j^{(1)}}, \quad \text{and} \quad \eta_j = s + \lambda_j^{(1)} + \psi^{(1)}(n_j)\mu_j^{(1)}$$

for  $j = 1, 2, \dots, m$ .

From (5.21)-(5.26) we have that

$$L_{T^{(1)}}(s) = \frac{\lambda^{(1)}}{s + \lambda^{(1)}} \frac{\mu(1 - \rho^{(1)})}{s + \mu(1 - \rho^{(1)})} \prod_{j=1}^m \frac{\psi^{(1)}(n_j)\mu_j(1 - \rho_j^{(1)})}{s + \psi^{(1)}(n_j)\mu_j(1 - \rho_j^{(1)})}$$

and

$$L_{T^{(2)}}(s) = \frac{\lambda^{(2)}}{s + \lambda^{(2)}} \frac{(1 - \rho)\delta^{(1)}}{1 - \rho\delta^{(1)}} \prod_{j=1}^m \frac{(1 - \rho_j)\delta_j^{(1)}}{1 - \rho_j\delta_j^{(1)}}$$

We observe that  $f_{T^{(r)}}(t)$  and  $F_{T^{(r)}}(t)$  ( $r = 1, 2$ ) are usually nonlinear functions of  $t$  and  $n_j$ . Hence, the resource optimization problem is an  $n$ -dimensional linear optimization problem subject to nonlinear constraints. In general, it is not easy to solve this problem. However, the complexity of the problem can be significantly reduced by requiring that the service rates of the queues in the Web services model in Figure 5.4 are all equal. That is, we find the optimum value of  $n_1, \dots, n_m$  such that  $\psi^{(r)}(n_1)\mu_1^{(r)} = \dots = \psi^{(r)}(n_m)\mu_m^{(r)}$  ( $r = 1, 2$ ), i.e., *balanced utilization*.

From the traffic equations:  $\lambda^{(r)} = \lambda_j^{(r)} = \Lambda^{(r)}$  ( $j = 1, 2, \dots, m$ ), we have that the utilization of each station  $\rho_j^{(r)} = \frac{\lambda_j^{(r)}}{\psi^{(r)}(n_j)\mu_j^{(1)}} = \frac{\Lambda^{(r)}}{\psi^{(r)}(n_j)\mu_j^{(1)}}$ . Thus we have that for the high-priority queue,  $\hat{a}_i = \psi^{(1)}(n_i)\mu_i^{(1)}(1 - \rho_i^{(1)}) = \psi^{(1)}(n_j)\mu_j(1 - \rho_j^{(1)}) = \hat{a}_j \triangleq \hat{a}$  that implies

$n_j = (\psi^{(1)})^{-1}(\frac{\hat{a} + \lambda_j^{(1)}}{\mu_j})$  ( $i, j = 1, 2, \dots, m$ ). Hence, from (5.21) we have

$$f_{T^{(1)}}(t) = L^{-1}\left\{\frac{\lambda^{(1)}}{s + \lambda^{(1)}} \cdot \frac{\mu(1 - \rho^{(1)})}{s + \mu(1 - \rho^{(1)})} \cdot \frac{\hat{a}^m}{(s + \hat{a})^m}\right\}$$

and subsequently obtain that

$$F_{T^{(1)}}(t) = L^{-1}\left\{\frac{\lambda^{(1)}}{s(s + \lambda^{(1)})} \cdot \frac{\mu(1 - \rho^{(1)})}{s + \mu(1 - \rho^{(1)})} \cdot \frac{\hat{a}^m}{(s + \hat{a})^m}\right\} \quad (5.27)$$

Moreover, since  $\hat{a}_i = \hat{a}_j = \hat{a}$  and  $\psi^{(1)}(n_i)\mu_i^{(1)} = \psi^{(2)}(n_i)\mu_i^{(2)}$  ( $i, j = 1, 2, \dots, m$ ), we have that  $\psi^{(2)}(n_1)\mu_1^{(2)} = \dots = \psi^{(2)}(n_m)\mu_m^{(2)} \triangleq \hat{b}$ . It is easy to see that  $\hat{b} = \hat{a} + \lambda_j^{(1)}$ . Thus for the low-priority queue, we obtain that  $\rho_{j_1} = \rho_{j_2} \triangleq \hat{\rho}$ , and  $\delta_{j_1}^{(1)} = \delta_{j_2}^{(1)} \triangleq \hat{\delta}^{(1)}$ , for  $j_1, j_2 = 1, 2, \dots, m$ . Furthermore,  $L_{T^{(2)}}(s)$  reduces to

$$L_{T^{(2)}}(s) = \frac{\lambda^{(2)}}{s + \lambda^{(2)}} \cdot \frac{(1 - \rho)\delta^{(1)}}{1 - \rho\delta^{(1)}} \cdot \frac{(1 - \hat{\rho})^m(\hat{\delta}^{(1)})^m}{(1 - \hat{\rho}\hat{\delta}^{(1)})^m} \quad (5.28)$$

which is a function of only one variable  $\hat{b}$ , i.e., variable  $\hat{a}$ . This is because  $\hat{\rho}$  and  $\hat{\delta}^{(1)}$  are considered as functions of only one variable  $\hat{b}$ , i.e., variable  $\hat{a}$ . Hence, we obtain the cumulative distribution of the response time for a low-priority customer

$$F_{T^{(2)}}(s) = L^{-1}\left\{\frac{\lambda^{(2)}}{s(s + \lambda^{(2)})} \cdot \frac{(1 - \rho)\delta^{(1)}}{1 - \rho\delta^{(1)}} \cdot \frac{(1 - \hat{\rho})^m(\hat{\delta}^{(1)})^m}{(1 - \hat{\rho}\hat{\delta}^{(1)})^m}\right\} \quad (5.29)$$

Since  $n_j^{(1)} = \lceil (\psi^{(1)})^{-1}(\frac{\hat{a} + \lambda_j^{(1)}}{\mu_j^{(1)}}) \rceil$  for a high-priority customer and  $n_j^{(2)} = \lceil (\psi^{(2)})^{-1}(\frac{\hat{b}}{\mu_j^{(2)}}) \rceil$  for a low-priority customer, the constraint of service availability at each resource site in (5.20) for a high-priority customer rewritten as

$$G_j^{(1)}(\hat{a}) \stackrel{def}{=} P_j^{(1)}(\lceil (\psi^{(1)})^{-1}(\frac{\hat{a} + \lambda_j^{(1)}}{\mu_j^{(1)}}) \rceil, N_j) \quad (5.30)$$

and the constraint of service availability at each resource site in (5.20) for a low-priority customer rewritten as

$$G_j^{(2)}(\hat{a}) \stackrel{def}{=} P_j^{(2)}(\lceil (\psi^{(2)})^{-1}(\frac{\hat{b}}{\mu_j^{(2)}}) \rceil, N_j) = P_j^{(2)}(\lceil (\psi^{(2)})^{-1}(\frac{\hat{a} + \lambda_j^{(1)}}{\mu_j^{(2)}}) \rceil, N_j) \quad (5.31)$$

Consequently,  $\sum_{j=1}^m n_j c_j$  reduces to a function of variable  $\hat{a}$  due to  $n_j^{(1)} = \lceil (\psi^{(1)})^{-1}(\frac{\hat{a} + \lambda_j^{(1)}}{\mu_j^{(1)}}) \rceil$  for a high-priority customer and to a function of variable  $\hat{b}$  due to  $n_j^{(2)} = \lceil (\psi^{(2)})^{-1}(\frac{\hat{b}}{\mu_j^{(2)}}) \rceil$  for a

low-priority customer. That is,  $g(n_1, \dots, n_m)$  in (5.18) reduces to a function of one variable. Thus the resource optimization problem can be decomposed into the one-dimensional resource optimization problems for both high-priority and low-priority customers respectively. Thus the trust-based resource provisioning problem is solved using the following iterative algorithm.

**Algorithm 5.4.2-1**

1. At time  $t_k$ , select  $m$  resource sites within predefined trust indices  $\hat{I}_j^{(r)}$ , and the highest  $m$  trust indices. If such a selection is impossible, then the trust manager informs the customer “We cannot process the service request at this moment.” After the customer request is waited for more than a given threshold time, the trust manager still cannot find those service sites that meet the predefined trust indices  $\hat{I}_j^{(r)}$ . Then, the trust manager informs the customer “You need to either revise the trust requirement or abort the request and resubmit it later.” Otherwise, continue to do Step 2.
2. Find  $\hat{a}$  in the following two one-dimensional optimization problems.

- (a) The minimization problem of a percentile response time for class  $r$  customers:

$$\hat{a}^{(r)}[1] \leftarrow \arg \min_{\hat{a}} F_{T^{(r)}}(t)|_{t=T_D^{(r)}}$$

subject to the constraint  $F_{T^{(r)}}(t)|_{t=T_D^{(r)}} \geq \gamma^{(r)}\%$  at  $\hat{a} = \hat{a}^{(r)}[1]$ , where  $F_{T^{(r)}}(t)$  are given by (5.27) and (5.29) respectively. Then, calculate the number of servers  $n_j^{(r)}[1]$  required for availability guarantee by using

$$n_j^{(r)}[1] = \lceil (\psi^{(r)})^{-1} \left( \frac{\hat{a}_j^{(r)}[1] + \lambda_j^{(1)}}{\mu_j^{(r)}} \right) \rceil$$

for  $r = 1, 2$ , and  $j = 1, 2, \dots, m$ .

- (b) The minimization problem of service availability:

$$\hat{a}_j^{(r)}[2] \leftarrow \arg \min_{\hat{a}} G_j^{(r)}(\hat{a})$$

subject to the constraint  $G_j^{(r)}(\hat{a}_j^{(r)}[2]) \geq \zeta_j^{(r)}\%$ , where  $G_j^{(r)}(\hat{a})$  are given by (5.30) and (5.31) respectively. Then, calculate the number of servers  $n_j^{(r)}[2]$  required for availability guarantee by using

$$n_j^{(r)}[2] = \lceil (\psi^{(r)})^{-1} \left( \frac{\hat{a}_j^{(r)}[2] + \lambda_j^{(1)}}{\mu_j^{(r)}} \right) \rceil$$

for  $r = 1, 2$ , and  $j = 1, 2, \dots, m$ .

3. Compute integers  $n_j^{(r)}$  by using

$$n_j^{(r)} = \max\{n_j^{(r)}[1], n_j^{(r)}[2]\}$$

Thus the number of servers required is  $n_j = \max\{n_j^{(1)}, n_j^{(2)}\}$  for station  $j$  ( $j = 1, 2, \dots, m$ ).

4. Update  $I_j^{t_k}$  to  $I_j^{t_k + T_D^{(r)}}$  based on the trust function (5.2). If  $|I_j^{t_k + T_D^{(r)}}| \geq \hat{I}_j^{(r)}$ , then the completed service job for class  $r$  customers is accepted. Otherwise, repeat Steps 1-3.
5. When the repeated number is more than a predefined threshold, the trust manager notifies the customer “We need to re-negotiate an SLA with the service broker.”

As mentioned early, Algorithm 5.4.2-1 is an iteration algorithm. The number of iterations will be determined by the trust manager and the service broker. When the number of iterations is more than a predefined threshold, the trust manager will notify the customer to modify the SLA. The customer may abort the service request or re-submit the service request with a new QoS requirement with a new fee. Additionally, the trust manager selects those service sites with the highest trust indices who also meet predefined trust requirements in Step 1 in Algorithm 5.4.2-1. But, for its own benefit, the trust manager can choose those service sites who only meet trust requirements predefined by a customer.

Moreover, Algorithm 5.4.2-1 shows that the  $m$ -dimensional optimization problem in Subproblem 2 significantly reduces to an one-dimensional optimization problem. Surprisingly, the optimal number of servers obtained by Algorithm 5.4.2-1 is independent of server costs  $c_j$ , due to  $\hat{a}_i = \hat{a}_j$  ( $i, j = 1, 2, \dots, m$ ). In general, each of the above two one-dimensional optimization problems can be solved numerically.

Note that the service broker cannot provide the quality of service at the moment in which an SLA negotiation is requested. In this case, the service broker may get a service penalty. For presentation purpose, it will be not further discussed here since the above approach can be easily adjusted to deal with this case. Additionally, we have only considered two-class customers. The above results can be extended to the case of customers with more than two classes.

## 5.5 Numerical Examples

In this section we demonstrate how to apply our algorithm to solve the trust-based resource provisioning problem. We consider the cases of single customer services and multiple prior-

Table 5.1: The Initial Trust Index of the Ten Stations for Single Class Customers

Station	1	2	3	4	5	6	7	8	9	10
$I^{t_1}$	0.8	0.5	0.2	0.9	0.6	0.4	0.8	0.7	0.3	0.6

Table 5.2: The Service Rates of the Ten Stations for Single Class Customers

Service Rates	$\mu_1$	$\mu_2$	$\mu_3$	$\mu_4$	$\mu_5$	$\mu_6$	$\mu_7$	$\mu_8$	$\mu_9$	$\mu_{10}$
Values	28	18	80	35	39	41	15	25	82	35

ity customer services respectively below.

### 5.5.1 Single Class Customers

Let us first consider a ten resource site example modelled by the tandem queueing network presented in Section 3.3. We choose  $m = 7$ ,  $\xi = 0.6$ , and the trust index at time  $t_1$  is listed in Table 5.1.

Assume that  $r_j(k)$  is uniformly distributed in  $[0.75, 1]$ .  $r_i$  and  $r_j$  are independent for any  $i \neq j$  ( $i, j = 1, \dots, 10$ ). Furthermore, we choose  $\lambda = 100$ ,  $T^D = 0.16$ ,  $\gamma = 97.5$ ,  $\mu = 200$ , and the service rates of these ten stations are given in Table 5.2. We also choose  $c_j = 2$ ,  $N_j = 200$ ,  $\psi(n_j) = 1.5^{\log_2 n_j}$ ,  $\delta_j = 99.999$ ,  $\hat{I}_j = 0.9$  ( $j = 1, \dots, 10$ ), and the server unavailable rates of these ten stations are listed in Table 5.3.

A customer submits a service job at time  $t_5$  with  $t_{k+1} - t_k = 0.01$  ( $k = 1, 2, \dots$ ) and it requires two of these 5 resource sites satisfying the predefined  $\hat{I}_j$  for processing the job.

First, we generated  $I^{t_k}$  ( $k = 2, \dots, 5, \dots, 21$ ) in Matlab as shown in Table 5.4. As we see, sites 1, 2, 3, 4, 6, 7, 8, and 10 meet the trust requirement at  $t = t_5$ . Thus sites 2, 3, 4, 6, 7, 8, and 10 are selected because they have the highest seven trust indices.

Then, we simulated the model in Arena 10.01. The simulation results are considered as

Table 5.3: The Server Unavailability Rates of the Ten Stations for Single Class Customers

Unavailable Rates	$\eta_1$	$\eta_2$	$\eta_3$	$\eta_4$	$\eta_5$	$\eta_6$	$\eta_7$	$\eta_8$	$\eta_9$	$\eta_{10}$
Values	0.012	0.01	0.005	0.04	0.015	0.03	0.02	0.045	0.008	0.01

Table 5.4: The Trust Indices of the Ten Stations for Single Class Customers at Times  $t = t_2, \dots, t_5, t_{20}, t_{21}$

	Station									
	1	2	3	4	5	6	7	8	9	10
$I^{t_2}$	0.8625	0.6412	0.5709	0.9156	0.7908	0.6651	0.8317	0.7892	0.3205	0.7428
$I^{t_3}$	0.8761	0.8190	0.7925	0.8501	0.8335	0.7698	0.8918	0.8593	0.4846	0.8865
$I^{t_4}$	0.8420	0.8532	0.8209	0.9014	0.7961	0.8725	0.9016	0.8946	0.4592	0.9132
$I^{t_5}$	0.9005	0.9129	0.9314	0.9248	0.8552	0.9153	0.9421	0.9189	0.5828	0.9588
...	...	...	...	...	...	...	...	...	...	...
$I^{t_{20}}$	0.8728	0.9215	0.9338	0.9182	0.8736	0.9326	0.9419	0.9546	0.6223	0.9235
$I^{t_{21}}$	0.8812	0.9513	0.9602	0.9278	0.9046	0.9462	0.9392	0.9698	0.8588	0.9351

“exact” since the simulation model is an exact representation of the queueing network under study. We further implemented (5.16) by using the inverse Laplace transform method in Graf [40] that is an approximate solution. The obtained cumulative distributions are shown in Table 5.5. It is seen that our approximation result has a good accuracy, where the relative error in Table 5.5 was calculated by

$$\text{Relative error } \% = \frac{\text{Approximation Result} - \text{Simulation Result}}{\text{Simulation Result}} \times 100$$

as given in (3.11). It is used to measure the accuracy of the approximate results compared to model simulation results.

Moreover, we obtained that  $\hat{a}^{(1)} = 100$ , and the optimal number of servers required for 97.5% of the response time to be less than  $T^D = 0.16$  is shown in Table 5.6. The exact optimal number of servers, obtained by exhaustive search using the simulation model, and assuming that each station has the same utilization, or balanced utilization, is consistent with the ones shown in Table 5.6. We validated that they are consistent with the result obtained by a brute force search using the simulation model in Arena, and assuming that each station has the same utilization, or balanced utilization.

Furthermore, we obtain the number of servers required for 99.999% service availability in these seven stations using Step 2 (b) of Algorithm 5.4.1-1, as shown in Table 5.6. By doing the calculation in Step 3 of Algorithm 3, we obtained the numbers of servers required for the response time and service availability guarantees at these seven stations respectively.

Finally, the trust manager needs to determine whether to accept the job completed by sites 2, 3, 4, 6, 7, 8, and 10 when it receives the completed job at  $t = t_5 + T^D = t_5 + 0.16 = t_{21}$ . As is

Table 5.5: The Cumulative Distribution of the Response Time for Single Class Customers

Response Time	Simul	Approx	R-Err %
0.04	0.0213	0.0214	0.4393
0.06	0.1517	0.1528	0.7004
0.08	0.4070	0.4075	0.1112
0.10	0.6681	0.6672	-0.1377
0.12	0.8468	0.8450	-0.2158
0.14	0.9398	0.9379	-0.1974
0.16	0.9785	0.9780	-0.0498
0.18	0.9931	0.9929	-0.0157
0.20	0.9979	0.9979	0.0000
0.22	0.9995	0.9994	-0.0077
0.24	0.9999	0.9998	-0.0051
0.26	1.0000	1.0000	0.0000

Table 5.6: The Optimal Number of Servers for Single Class Customers

Station	2	3	4	6	7	8	10
#Servers necessary to ensure 97.5% response time	62	5	20	16	84	35	20
#Servers necessary to ensure 99.999% availability	10	7	22	18	15	23	10
#Servers necessary to ensure these two SLA metrics	62	7	22	18	84	35	20

seen, these stations meet the trust requirement at  $t = t_{21}$ , and consequently the job is accepted.

The above ten-sites example has demonstrated how to apply our approach to solving the trust-based resource provisioning problem in the case of single class customers.

Next, we are going to demonstrate how to use our proposed algorithm to solve the trust-based resource provisioning problem in the case of multiple class customers.

## 5.5.2 Multiple Priority Customers

In the case of multiple class customers, we still consider a ten resource site example modelled by the tandem queueing network presented in Section 5.4. We choose  $m = 7$ ,  $\xi = 0.6$ , and the trust index at time  $t_1$  is unchanged, i.e., it is listed in Table 5.1.

Again, assume that  $r_j(k)$  is uniformly distributed in  $[0.75, 1]$ .  $r_i$  and  $r_j$  are independent for any  $i \neq j$  ( $i, j = 1, \dots, 10$ ). Moreover, we choose  $\Lambda^{(1)} = 100$ ,  $\Lambda^{(2)} = 50$ ,  $T_D^{(1)} = 0.16$ ,  $T_D^{(2)} = 0.8$ ,  $\gamma^{(1)} = 97.5$ ,  $\gamma^{(2)} = 99$ ,  $\mu = 200$ , and the service rates of these seven stations are given

Table 5.7: The Initial Trust Index of the Ten Stations for Priority-class Customers

Station	1	2	3	4	5	6	7	8	9	10
$I^{t_1}$	0.8	0.5	0.2	0.9	0.6	0.4	0.8	0.7	0.3	0.6

Table 5.8: The Service Rates of the Ten Stations for Priority-class Customers

Service Rates	$\mu_1$	$\mu_2$	$\mu_3$	$\mu_4$	$\mu_5$	$\mu_6$	$\mu_7$	$\mu_8$	$\mu_9$	$\mu_{10}$
$\mu_j^{(1)}$	88	18	85	32	31	49	24	28	21	38
$\mu_j^{(2)}$	96	15	60	25	55	41	18	26	31	35

in Table 5.8. We also choose  $c_j = 1$ ,  $N_j = 200$ ,  $\psi(n_j) = 1.55^{\log_2 n_j}$ ,  $\xi_j^{(1)} = 99.999$ ,  $\xi_j^{(2)} = 99.9$ ,  $\hat{I}_j^{(1)} = 0.9$ ,  $\hat{I}_j^{(2)} = 0.91$  ( $j = 1, \dots, 10$ ), and the server unavailable rates of these ten stations are listed in Table 5.9.

Again, a customer submits a service job at time  $t_5$  with  $t_{k+1} - t_k = 0.01$  ( $k = 1, 2, \dots$ ) and it requires two of these 5 resource sites satisfying the predefined  $\hat{I}_j$  for processing the job.

We first extended the generation of trust indices at the ten station. That is, we generated  $I^{t_k}$  ( $k = 2, \dots, 5, \dots, 21, \dots, 84, 85$ ) in Matlab shown in Table 5.10, as an extension of Table 5.4. As is seen, while sites 1, 2, 3, 4, 6, 7, 8, and 10 meet the trust requirement for high-priority customers at  $t = t_5$ , only sites 2, 3, 4, 6, 7, 8, and 10 meet the trust requirement for low-priority customers at  $t = t_5$ . Thus sites 2, 3, 4, 6, 7, 8, and 10 are selected because they meet both trust requirements.

Then, we simulated the model in Arena 10.01. The simulation results are considered as “exact” since the simulation model is an exact representation of the queueing network under study. We further implemented (5.27) by using the inverse Laplace transform method in Graf [40] that is an approximate solution. Tables 5.11 and 5.12 show the simulated and approximate cumulative distributions of the high-priority and low-priority response times respectively. It appears that the

Table 5.9: The Server Unavailability Rates of the Ten Stations for Priority-class Customers

Unavailable Rates	$\eta_1$	$\eta_2$	$\eta_3$	$\eta_4$	$\eta_5$	$\eta_6$	$\eta_7$	$\eta_8$	$\eta_9$	$\eta_{10}$
Values	0.012	0.01	0.005	0.04	0.015	0.03	0.02	0.045	0.008	0.01

Table 5.10: The Trust Indices of the Ten Stations for Priority-class Customers at Times  $t = t_2, \dots, t_5, \dots, t_{20}, t_{21}, \dots, t_{84}, t_{85}$

	Station									
	1	2	3	4	5	6	7	8	9	10
$I^{t_2}$	0.8625	0.6412	0.5709	0.9156	0.7908	0.6651	0.8317	0.7892	0.3205	0.7428
$I^{t_3}$	0.8761	0.8190	0.7925	0.8501	0.8335	0.7698	0.8918	0.8593	0.4846	0.8865
$I^{t_4}$	0.8420	0.8532	0.8209	0.9014	0.7961	0.8725	0.9016	0.8946	0.4592	0.9132
$I^{t_5}$	0.9005	0.9129	0.9314	0.9248	0.8552	0.9153	0.9421	0.9189	0.5828	0.9588
...	...	...	...	...	...	...	...	...	...	...
$I^{t_{20}}$	0.8728	0.9215	0.9338	0.9182	0.8736	0.9326	0.9419	0.9546	0.6223	0.9235
$I^{t_{21}}$	0.8812	0.9513	0.9602	0.9278	0.9046	0.9462	0.9392	0.9698	0.8588	0.9351
...	...	...	...	...	...	...	...	...	...	...
$I^{t_{84}}$	0.9126	0.9318	0.9458	0.9288	0.8538	0.9225	0.9518	0.9449	0.7256	0.9338
$I^{t_{85}}$	0.8918	0.9815	0.9731	0.9388	0.9126	0.9565	0.9291	0.9588	0.8528	0.9588

results obtained by Algorithm 5.4.2-1 are very accurate. It is shown in Table 5.13 that the optimal number of servers is required for 97.5% of the high-priority response time to be less than  $T_D^{(1)} = 0.16$  and for 99% of the low-priority response time to be less than  $T_D^{(2)} = 0.8$ . As is seen in Table 5.11, our approximation result has a good accuracy, where the relative error was calculated by (3.11). Again, the relative error is used to measure the accuracy of the approximate results compared to model simulation results.

Furthermore, the optimal number of servers is required for satisfying both the high-priority and the low-priority response times is shown in Table 5.13. The exact optimal number of servers, obtained by exhaustive search using the simulation model, and assuming that each station has balanced utilization, is consistent with the ones shown in Table 5.13. We validated that they are consistent with the result obtained by a brute force search using the simulation model in Arena, and assuming that each station has the same utilization, or balanced utilization.

Additionally, we obtain the number of servers required for 99.999% service availability guarantee for high-priority customers and for 99.9% service availability guarantee for low-priority customers in these seven stations using Step 2 (b) of Algorithm 5.4.2-1, as shown in Table 5.13. By doing the calculation in Step 3 of Algorithm 3, we obtained the numbers of servers required for the response time and service availability guarantees at these seven stations respectively.

To the end, the trust manager needs to determine whether to accept the job completed by sites 2, 3, 4, 6, 7, 8, and 10 when it receives the completed job at  $t = t_5 + T_D^{(1)} = t_5 + 0.16 = t_{21}$  for

Table 5.11: The Cumulative Distribution of the High-priority Response Time

Response Time	Simul	Approx	R-Err %
0.02	0.0002	0.0002	0.0000
0.04	0.0214	0.0214	0.0000
0.06	0.1542	0.1528	-0.9079
0.08	0.4085	0.4075	-0.2448
0.10	0.6691	0.6672	-0.2840
0.12	0.8459	0.8450	-0.1064
0.14	0.9385	0.9379	-0.0639
0.16	0.9781	0.9780	-0.0102
0.18	0.9931	0.9929	-0.0201
0.20	0.9980	0.9979	-0.0100
0.22	0.9995	0.9994	-0.0100
0.24	0.9998	0.9998	0.0000
0.26	0.9999	1.0000	0.0100
0.28	1.0000	1.0000	0.0000

Table 5.12: The Cumulative Distribution of the Low-priority Response Time for Priority-class Customers

Response Time	Simul	Approx	R-Err %
0.2	0.1767	0.1473	-16.6384
0.3	0.4538	0.4396	-3.1291
0.4	0.6982	0.7082	1.4323
0.5	0.8541	0.8719	2.0841
0.6	0.9389	0.9503	1.2142
0.7	0.9774	0.9824	0.5116
0.8	0.9925	0.9942	0.1713
0.9	0.9978	0.9982	0.0401
1.0	0.9993	0.9995	0.0200
1.1	0.9998	0.9999	0.0100
1.2	1.0000	1.0000	0.0000

Table 5.13: The Optimal Number of Servers for Priority-class Customers

The number of Servers	Stations						
Metrics	2	3	4	6	7	8	10
97.5% response time guarantee for high-priority customers	62	5	23	12	38	29	18
99% response time guarantee for low-priority customers	61	7	27	13	46	26	16
99.999% service availability guarantee for high-priority customers	10	7	22	18	15	23	10
99.9% service availability guarantee for low-priority customers	7	5	14	15	11	19	7
#Servers necessary to ensure all these SLA metrics	62	7	27	18	46	29	18

high-priority customers and at  $t = t_5 + T_D^{(2)} = t_5 + 0.8 = t_{85}$  for low-priority customers respectively. As is seen, these stations meet the trust requirements at  $t = t_{21}$  for high-priority customers and  $t = t_{85}$  for low-priority customers, and consequently the completed jobs are accepted.

We have demonstrated how to apply our approach to solving the trust-based resource provisioning problem in the above ten-site example for the case of multiple class customers.

As we know, most existing Web services products do not support an SLA that guarantees a level of service delivered to a customer for a given price. It is not easy to solve a resource provisioning problem when we consider all these constraints: trustworthiness, an end-to-end response time and service availability. The last two sections demonstrated how to apply our efficient algorithm to solving the trust-based resource provisioning problem by using the above two illustrative examples in the cases of single and multiple class customers respectively.

## 5.6 Concluding Remarks

In this chapter, we have proposed a trust-based Web services model, and provided a framework for studying the model. We have discussed a trust-based resource provisioning problem that arises in these typical Web services applications. The problem has been constructed by minimizing the total cost of service providers while satisfying SLA guarantees. We have further formulated it as an optimization problem under the constraints of percentile response time and service availability in the cases of single class customers and multiple class customers respectively.

In our approach we considered the percentile response time that may better address a customer's concern as opposed to the average response time that is commonly used in the literature.

We obtained an efficient and accurate numerical solution for calculating the percentile response time. Then, we proposed efficient approaches for solving the trust-based resource provisioning problem in the cases of single class customers and multiple class customers respectively. We used two numerical examples to illustrate the use of our proposed algorithms in these two cases. Our numerical validations showed that our algorithms have provided a good accuracy.

## Chapter 6

# Performance Analysis of Public Key

## Cryptograph-based Group

### Authentication

Several Kerberos-based authentication techniques using public-key cryptography have been proposed in the last decade. The use of public-key cryptography is to eliminate the problem of a single point failure in a Key Distribution Center (KDC), and achieve better scalability. Among them, the two notable techniques are Public Key Cryptography for Cross Realm Authentication in Kerberos (PKCROSS) and Public Key Utilizing Tickets for Application Servers (PKTAPP, a.k.a. KX.509/KCA). The latter was proposed to improve the scalability of the former. But, the actual costs (e.g., computational and communication times) associated with these techniques have been poorly understood so far. It remains unknown which technique performs better in a large network where there are multiple KDC remote realms.

This chapter is organized as follows. An overview of public key cryptography-based authentication is given in Section 6.1. In Section 6.2, we first give an in-depth discussion of PKCROSS and PKTAPP and then present their performance evaluation using queueing theory. Section 6.3 gives a definition of the new public key cryptography-based group authentication technique along with

the details of messages. We present a performance analysis of the proposed technique and compare it to PKCROSS and PKTAPP in Section 6.4. We finally conclude our discussion in Section 6.5.

## 6.1 Public Key Cryptography-based Authentication

In the last few years we have witnessed an explosive growth in the usage of Internet collaborative applications, such as video and audio conferencing, replicated servers and databases, and in particular Web services whereby service components from several universal service providers can be flexibly integrated into a composite service regardless of their location, platform, and execution speed [9]. To ensure quality of service, the service providers are required to work together. The rapid growth in collaborative applications has heightened the need for a reliable group communication system. The system, in turn, is impossible to be reliable without group authentication.

Kerberos [52] consisting of a client, an application server and a key distribution center (KDC) is a mature, reliable, secure network authentication protocol that allows a client to prove its identity to a server without sending confidential data across the network. Public key cryptography has been extended to support Kerberos since it simplifies the distribution of keys in Kerberos. It eliminates a single point of failure. Integrating public key cryptography into Kerberos represents the enhancements of the current Kerberos standard. Several Kerberos-based authentication techniques using public-key cryptography have been proposed in the last decade. Among them are Public-Key Cross Realm Authentication in Kerberos (PKCROSS) [46] and Public-key Cryptography for Initial Authentication in Kerberos (PKINIT) [99]. Moreover, the scalability of network security infrastructures is becoming a serious concern as the explosive growth of collaborative applications such as Web services continues unabated. Public key based Kerberos for Distribution Authentication (PKDA) [85] and Public Key Utilizing Tickets for Application Servers (PKTAPP, a.k.a. KX.509/KCA) [53] and [61] have been proposed to enhance the security and scalability of Kerberos.

PKINIT is a core specification among these Internet drafts. Both PKCROSS and PKTAPP use variations of PKINIT message types and data structures for integrating public key cryptography with Kerberos in different authentication stages. PKTAPP was originally introduced as PKDA. It implemented PKDA using the message formats and exchanges of PKINIT. Microsoft has adopted the Internet draft specification of PKINIT for the support of public key cryptography in the Windows 2000 and 2003 implementations of Kerberos [37]. It has its own protocol that is the equivalent of

KX509/KCA. There were preliminary discussions between the Kerberos WG and Microsoft about using a common protocol. The MIT Kerberos consortium will drive these discussions [7]. According to Altman [6], the standardization of PKCROSS and PKTAPP will be the next for Kerberos and PKI integration. The Kerberos Consortium at MIT was just formed in September 2007 and listed PKCROSS as Project 10 in its proposal [28]. We believe that PKCROSS and PKTAPP will revive soon. Hence, this research only considers these two notable techniques: PKCROSS and PKTAPP.

It has been argued that PKCROSS would not scale well in large networks in [85]. PKTAPP was proposed to improve the scalability of PKCROSS. However, the actual costs associated with these techniques have been *poorly understood* so far. PKTAPP has been shown in [42] to poorly perform when there are two or more application servers in one remote realm. In addition, it remains *unknown* as to which technique performs better in a large network where as is typical in many applications, application servers are within *multiple* KDC remote realms.

In next section, we first present a thorough performance evaluation of PKCROSS and PKTAPP in terms of computational and communication costs. Then, we demonstrate their performance difference using open queueing networks.

## 6.2 PKCROSS and PKTAPP

Kerberos [52] was developed at the Massachusetts Institute of Technology in 1988. It is a network authentication protocol for providing a secure communication between a user workstation (or called a client) and application servers. The latest version of Kerberos is Version 5. It divides the world into realms, each with user workstations, a single primary KDC, back-up KDCs, and application servers in which the KDC is a trusted intermediary. In the Kerberos protocol, the client engages in a multiple-step authentication to obtain access to the application server whereby the client first obtains a relatively short lived credential: a Ticket-Granting Ticket (TGT) from the Authentication Service running on a KDC, and then obtain a session ticket for a particular application server by presenting the TGT to a centralized Ticket-Granting Service (TGS) running on the KDC. The client presents the session ticket to the application server for authenticating herself/himself by showing knowledge of a secret session key. The secret session key was securely passed to the client by the KDC. Kerberos is stateless, and this is extremely valuable from the scalability point of view. Cross-realm authentication is necessary when a client and an application server with different network domains fall into different Kerberos realms.

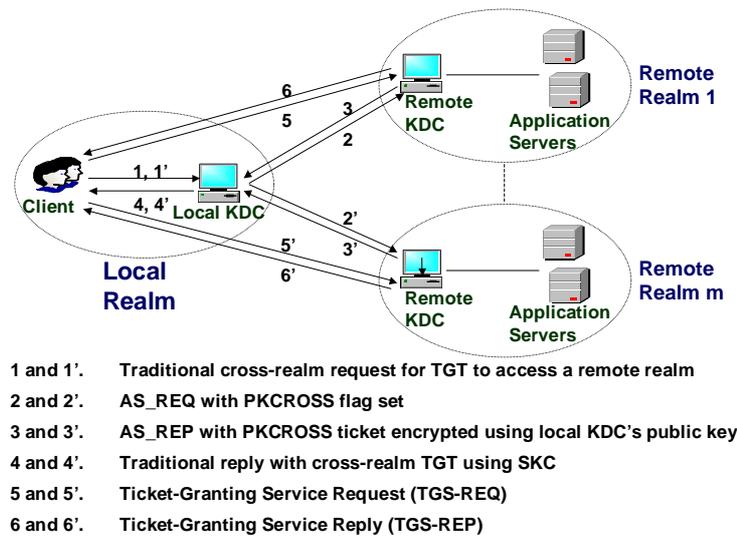


Figure 6.1: The PKCROSS Message Flow

It is well-known that a public key security system is easier to administer, more secure, less trustful, and more scalable than a symmetric key security system. Public key security does not use a trusted key management infrastructure since the burden of key management falls to public key clients [32]. In a public key infrastructure, public key clients need to constantly and vigorously check the validity of the public keys that they use. Public key cryptography shifts the burden on key management from the KDC/TGS in Kerberos to its Certificate Authority (CA) who may be considered as a trusted intermediary. CA issues a public key certificate that is relatively long lived credential. The burden is determined by the number of times clients want to authenticate to application servers in Kerberos. It might not be affordable in time-sensitive applications if one large-scale PKI deployment is needed for group authentication in a large network.

### 6.2.1 Protocol Analysis

PKCROSS [46] is one notable protocol of integrating public key cryptography with Kerberos to address the problem of network authentication among the client and the application servers in a large number of realms. Figure 6.1 illustrates the authentication steps of PKCROSS. The cross realm KDC to KDC authentication is achieved by using public key cryptography. First, just like Kerberos the KDC in PKCROSS is burdened by the need to constantly renew short-lived TGTs,

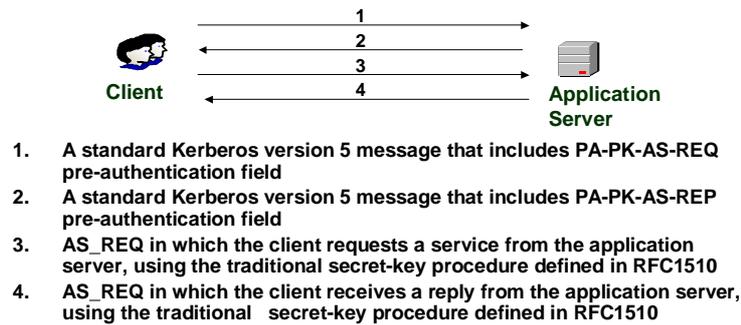


Figure 6.2: The PKTAPP Message Flow

and TGS must be involved whenever a client wants to request a service from an application server. Thus, if a large number of users in a single realm request services which may be a typical case in Web services, the KDC in the realm becomes a serious single point of failure due to a KDC compromise and possibly a performance bottleneck. Recovering from such a compromise requires the re-establishment of secret keys for all users within this realm. When a number of users is large, such a recovery is very costly. PKINIT facilitates public-key based authentication between users and their local KDC. Hence, one possible solution to eliminating the single point of failure is to combine PKCROSS with PKINIT, and thus public-key based authentication becomes integrated throughout the *entire* Kerberos environment. However, it is easy to see that this solution is not feasible for time-sensitive applications, since the users might suffer from a long delay or even denial of services due to a high computational cost for the calculation of a public key used in the entire environment.

Second, in PKCROSS the local KDC of the client issues all short-lived TGTs and all session tickets in its realm, and communicates with a remote KDC. Hence, it can easily become a performance bottleneck since all these authentication transactions have to transit the KDC. Thus, PKTAPP [61] has been proposed to address the issue. Figure 6.2 shows the message exchange of PKTAPP. The PKTAPP technique allows the client to communicate directly with the application servers so that the number of messages between them is reduced in an authentication process. But, as is seen in [42], even though PKTAPP requires fewer message exchanges than PKCROSS during client authentication with remote application servers, PKCROSS outperforms PKTAPP when two or more application servers are in a single remote realm. This is because PKTAPP requires more public-key message exchanges than PKCROSS which only requires one pair of public-key message

Table 6.1: The Operations of Encryption and Decryption When  $n$  Application Servers are in  $m$  Remote Realms

Protocols	Entities	# Secret Key	# Private Key	# Public Key
PKTAPP	Client	$2n+1$	$n+1$	$3n$
	Application server	$3n+1$	$n$	$4n$
	Total	$5n+2$	$2n+1$	$7n$
PKCROSS	Client	$m+6$	$0$	$0$
	Local KDC	$5$	$m+1$	$3m$
	Remote KDC	$3m+n$	$m$	$4m$
	Application server	$3n$	$0$	$0$
	Total	$4(n+m)+11$	$2m+1$	$7m$

exchanges. Below, we study their performance in multiple remote realms.

Table 6.1 summarizes the encryption and decryption operations performed in PKTAPP and PKCROSS, where  $m$  and  $n$  are the number of remote realms and the number of applications servers within these realms respectively. Denote by  $c_1$ ,  $c_2$ , and  $c_3$  the computational times per secret, private, or public key operation respectively. Then,  $c_1 < c_2 < c_3$ . Let  $f_j(n, m)$ ,  $j = 1, 2$ , be the total computational times of encryption and decryption operations in PKTAPP and PKCROSS. Then, from Table 6.1 we have that

$$f_1(n, m) = (5c_1 + 2c_2 + 7c_3)n + 2c_1 + c_2$$

$$f_2(n, m) = 4c_1n + (4c_1 + 3c_2 + 7c_3)m + 11c_1 + c_2$$

Note that  $f_1(n, m)$  does not depend on  $m$ , which can be hence abbreviated as  $f_1(n)$ . Then, we have the following proposition.

**Proposition 6.2.1** *For each authentication, PKCROSS requires less computational time than PKTAPP if and only if the number of application servers  $n$  is more than  $\lceil m + \frac{3(m+3)c_1}{c_1+2c_2+7c_3} \rceil$ , or the number of remote realms  $m$  is less than  $\lfloor n - \frac{3(n+3)c_1}{4c_1+2c_2+7c_3} \rfloor$ .*

*Proof.* According to  $f_1(n, m)$  and  $f_2(n, m)$ , their difference is given by  $f_1(n, m) - f_2(n, m) = -9c_1 + (c_1 + 2c_2 + 7c_3)n - (4c_1 + 2c_2 + 7c_3)m$ . Hence,  $f_1(n, m) - f_2(n, m) \geq 0$  if and only if  $-9c_1 + (c_1 + 2c_2 + 7c_3)n - (4c_1 + 2c_2 + 7c_3)m \geq 0$ , which implies Proposition 6.2.1. The proof is complete.

It is anticipated that the client is connected to the local KDC by a LAN, the client and the local KDC are connected to a remote KDC by a WAN, and a remote KDC and its application

servers are connected by a WAN. Assume that all WANs have an identical communication time and a local area network (LAN) has a negligible communication time compared to a WAN. From Figures 6.1 and 6.2, we note that the number of WAN communications are  $4n$  for PKTAPP and  $4m + 2n$  for PKCROSS. Let  $g_j(n, m)$ ,  $j = 1, 2$ , be the transaction times of PKTAPP and PKCROSS. *The transaction time* is defined as the computational time of the total encryption and decryption operations plus the communication time per authentication in a technique. Also, denote by  $d$  the time spent in a WAN communication. Then, the following conclusion holds.

**Proposition 6.2.2** *For each authentication, PKCROSS uses less transaction time than PKTAPP if and only if the number of application servers  $n$  is more than  $m + \lceil \frac{(3c_1+2d)m+9c_1}{c_1+2c_2+7c_3+2d} \rceil$ .*

*Proof.* For  $j = 1, 2$ ,  $g_j(n, m)$  can be computed by  $g_1(n, m) = f_1(n) + 4n$  and  $g_2(n, m) = f_2(n, m) + 4m + 2n$  and thus their difference is given by  $g_1(n, m) - g_2(n, m) = (c_1 + 2c_2 + 7c_3 + 2d)n - (4c_1 + 2c_2 + 7c_3 + 4d)m - 9c_1$  which easily derives Proposition 6.2.2. The proof is complete.

Note that if  $n = 1$  (so,  $m = 1$ ),  $m + \lceil \frac{(3c_1+2d)m+9c_1}{c_1+2c_2+7c_3+2d} \rceil > 1 = n$ , and if  $m = 1$ ,  $m + \lceil \frac{(3c_1+2d)m+9c_1}{c_1+2c_2+7c_3+2d} \rceil = 1 + \frac{12c_1+2d}{c_1+2c_2+7c_3+2d} < 2$  since  $c_1$  is significantly smaller than  $c_3$ . This means that

**Corollary 1** *When  $m = 1$ , we have that PKTAPP requires less transaction time than PKCROSS if  $n = 1$  but more transaction time than PKCROSS if  $n \geq 2$ .*

In Proposition 6.2.2 we have shown that the number of application servers should be  $\lceil \frac{(3c_1+2d)m+9c_1}{c_1+2c_2+7c_3+2d} \rceil$  more than the number of remote KDC realms so as to ensure that PKCROSS uses less transaction time than PKTAPP. But, the transaction time does not take into account the time required to wait in a queue for a service (i.e., waiting time) when multiple authentication requests present in any one of the authentication entities. Response time is used when such a case is considered. That is, *the response time* is the transaction time plus the waiting time per authentication request. A further discussion of the response time is given in next section.

## 6.2.2 The Calculation of the Response Time via Queueing Networks

A queueing network is an efficient tool to analyze the scalability of a distributed system. To investigate the scalability of PKTAPP and PKCROSS, we first model the entities, the client, the

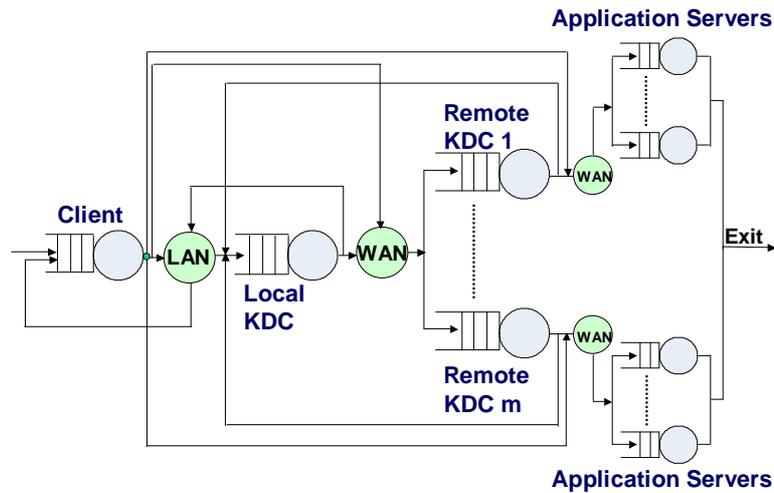


Figure 6.3: A Queueing Network with  $m$  Remote Realms for PKCROSS

local KDC, the remote KDCs, the application servers and communication networks, as a queueing network. The client may represent a single or multiple users that request group authentication at a given rate and the authentication request is processed in these entities according to these two techniques. Figures 6.3 and 6.4 give two queueing networks which depict the message flows of PKCROSS and PKTAPP where each system resource is modeled as a queue associated with a queueing discipline. Since a public key consumes a significantly higher computation cost than a private key, it is not reasonable to assume that all client requests are served at the same average service time. Instead, a class-switching technique as in [17] and [68] is employed to model the class transaction with switching from low- to high-priority class, as different types of encryption/decryption operations are required with different service times. Preemption-resume is one good way to implement service for satisfying multiple class client requests. In the Chapter, we use a *preemptive-resume* priority discipline, i.e., the service of a class  $r_2$  request can be interrupted if a higher-priority request of class  $r_1$  ( $r_2 > r_1$ ) arrives during its service. The interrupted request resumes its service from where it stopped after the higher-priority request, and any other request with priority higher than  $r_2$  that may arrive during its service, complete their service.

Our goal is to use the queueing networks to calculate the response time of an authentication request. The calculation is sketched below. Assume that the client requests group authentication at a rate  $\lambda$ . Based on the forced law, the throughput of an entity (the client station, the local KDC, the remote KDCs, or the application servers) is  $X^{(j)} = \lambda v^{(j)}$  for class  $j$  job where  $v^{(j)}$  is the num-

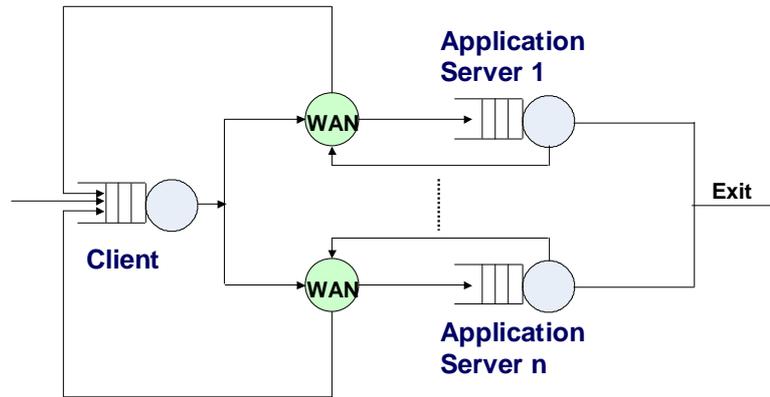


Figure 6.4: A Queueing Network with  $n$  Application Servers for PKTAPP

ber of visits to the entity by class  $j$  jobs. Then, the total throughput  $X$  is a sum of  $X^{(j)}$  over all job classes at the entity. Thus, according to the utilization law, the utilization of the entity by class  $j$  jobs is  $\rho^{(j)} = X^{(j)}\mu^{(j)} = \lambda v^{(j)}\mu^{(j)}$  where  $\mu^{(j)}$  is the service time per visit to the entity by class  $j$  jobs. Hence, the total utilization  $\rho$  is a sum of  $\rho^{(j)}$  over all classes at the entity. Denote by  $v$  a total of the number of visits at the entity. Thereby, the response time of the entity is  $R = \frac{\mu}{1-\rho}$  where  $\mu$  is an average service time of all classes at the entity, and the total response time per authentication request is a sum of  $vR$  over all entities.

To validate the accuracy of the queueing networks, we first simulated the scenario of these entities as shown in Figure 6.5 by doing the reference implementations of these two techniques under Windows XP. Moreover, a public-key cipher is usually associated with the calculations of 1024-bit precision numbers, so a public-key operation is computationally expensive and it costs as much as a factor of 1000 than an equivalent secret-key operations [34]. In the reference implementations we adopted the results from the Crypto++ ran on an Intel Core 2 1.83 GHz processor under Windows XP SP 2 in 32-bit mode [31]. Table 6.2 gives the time required to encrypt and decrypt a 64-byte block of data. Without loss of generality, we chose  $c_1 = 0.000248$  msec,  $c_2 = 0.07$  msec and  $c_3 = 1.52$  msec. (The performance evaluation based on an ECC key will be discussed in my future research.) Table 6.3 shows the accuracy of analytical response times obtained by the queueing methods compared to simulated results based on the scenario shown in Figure 6.5, where R-Err% is the relative error used to measure the accuracy of the analytic results compared to model simulation results, and it is defined by  $(\text{analytic result} - \text{simulated result})/\text{simulated result} \times 100$ . As seen in the table, the analytic response times match the simulated results very well.

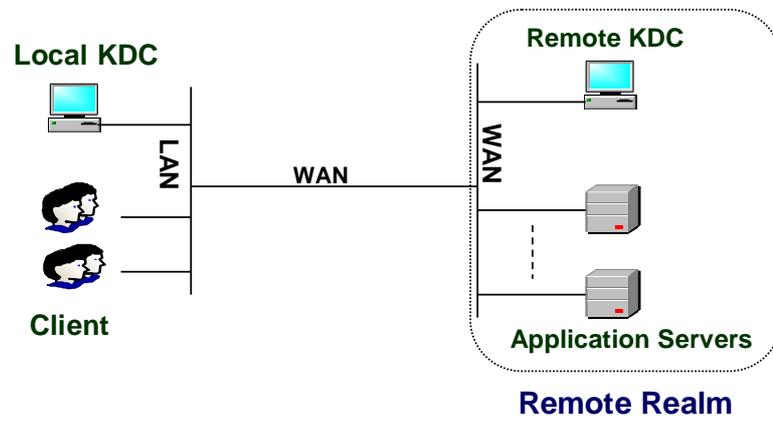


Figure 6.5: A Test Scenario with a Remote Realm

Table 6.2: The computational times of Encryption and Decryption Operations

Protocols and Operation	Key Length	Computational times (msec)
AES/ECB	128	0.000248
AES/ECB	256	0.000312
RSA encryption	1024	0.07
RSA decryption	1024	1.52
RSA encryption	2048	0.15
RSA decryption	2048	5.95

Table 6.3: A Comparison of Analytic and Simulated Response Times

$m=1$		The Number of Application Servers				
Protocols	Methods	1	2	4	8	16
PKCROSS	Analytic	102.1	122.1	162.1	242.1	402.1
	Simulated	102.3	122.5	162.9	244.2	408.8
	R-Err%	-0.22	-0.30	-0.47	-0.85	-1.63
PKTAPP	Analytic	82.10	164.05	327.96	655.76	1311.38
	Simulated	82.23	164.56	329.97	663.87	1344.23
	R-Err%	-0.15	-0.30	-0.61	-1.22	-2.44

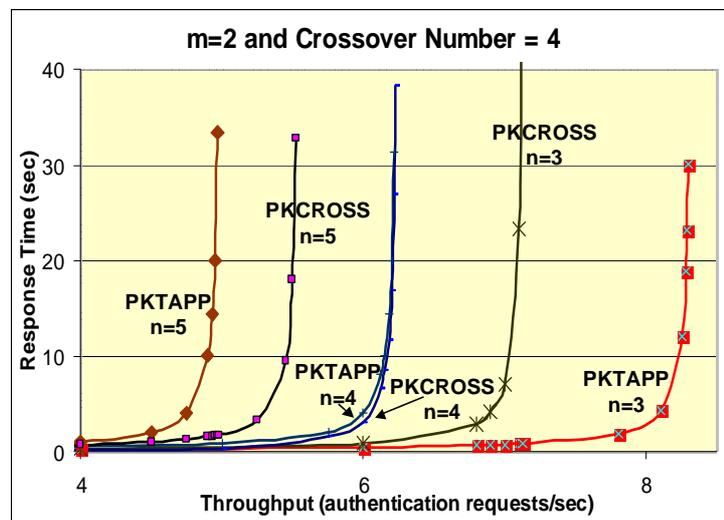


Figure 6.6: Response Time vs. Authentication Request Rate

To investigate the performance with an increased number of application servers and remote realms, we further present response times as a function of authentication request rates, i.e., throughput, when there are two remote realms, as shown in Figure 6.6. As is seen, PKCROSS has a slightly lower response time than PKTAPP when  $n = 4$ , called *the crossover number*. That is, PKCROSS is more efficient than PKTAPP if  $n \geq 4$ , but less efficient than PKTAPP if  $n < 4$ . Clearly, it follows from Table 6.3 that the crossover number is equal to 2 when  $m = 1$ . In general, Figure 6.7 shows crossover numbers with varying number of remote realms. The crossover numbers can be 99.8% perfectly fitted by the straight line:  $n = 1.8916m + 0.2879$ . This means that PKCROSS is more efficient than PKTAPP when  $n \geq \lceil 1.8916m + 0.2879 \rceil$ . In other word, PKTAPP

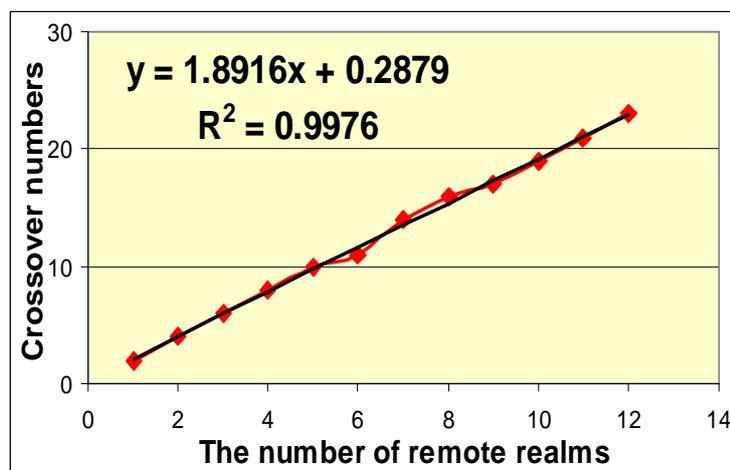


Figure 6.7: Crossover Numbers vs. The Number of Remote Realms

does not scale better than PKCROSS. That is different from what PKTAPP's designers expected. In the next section we propose a hybrid approach that has better scalability than PKCROSS.

### 6.3 A New Group Authentication Technique Using Public-Key Cryptography

As seen above, in PKCROSS the client is first required to communicate with the KDC before talking to the application server. PKTAPP uses a direct communication between the client and the application servers for the reduction of message exchanges and the relieve of a single point of failure on the client's KDC. While PKTAPP is more efficient in the case of only a single application server and thus only a single remote realm too, PKCROSS significantly performs better if the number of application servers is more than a crossover number as shown in Table 6.3 and Figure 6.7. Our proposed technique below takes into consideration the advantages of both techniques. While it allows the client to deal directly with the application servers so that the number of messages is reduced in the authentication process like PKTAPP, the new technique still relies on the KDC for the authentication of the client and the application servers like PKCROSS.

Table 6.4: The Notation Used in the Case of a Single Remote Realm

$C$	Client
$S_j$	Application Server $j$ ( $j = 1, 2, \dots, n$ )
$KDC_L$	Local KDC, i.e., Client's KDC
$KDC_R$	Remote KDC, i.e., Application servers' KDC
$K_C$	Secret key of $C$
$K_{S_j}$	Secret key of $S_j$
$K_{C,S}$	Group key shared by $C$ and all $S_j$
$\{\text{Message}\}_{KDC_L}$	Message encrypted with $KDC_L$ 's public key
$\{\text{Message}\}_{KDC_R}$	Message encrypted with $KDC_R$ 's public key
$[\text{Message}]_{KDC_R}$	Message signed with $KDC_R$ 's private key
$KB_{KDC_R}$	Public key of $KDC_R$
$Cert_{KDC_L}$	Certificate of $KDC_L$
$Cert_{KDC_R}$	Certificate of $KDC_R$
$N_C$	Nonce generated by $C$
$N_{KDC_R}$	Nonce generated by $KDC_R$
$TMS_C$	Timestamp generated by $C$
$TMS1_R, TMS2_R$	Timestamps generated by $KDC_R$

### 6.3.1 A Single Remote Realm

In this section we define the group authentication technique among a client (denoted by  $C$ ) and application servers (denoted by  $S_j$ ) that are in a single remote realm (denoted by  $KDC_R$ ), where  $j = 1, \dots, n$ . The client  $C$  is a group key initiator, for example, representing either a group member in a video conference or a service provider in Web services applications who wants to initiate secure communication among groups or service providers. (How  $C$  is chosen is beyond the scope of our study in this dissertation.) The application servers  $S_j$  are the rest of either group members or service providers. Table 6.4 lists the notation used in this section. Figure 6.8 shows the message flow of the new group authentication technique whose message exchanges is given in detail as follows.

**Message 1,  $C \rightarrow KDC_R$ :**

$N_C, "C", "KDC_R: S_1, \dots, S_n", K_C\{N_C, K_{C,S}, "KDC_R: S_1, \dots, S_n", TMS_C\}$

The client  $C$  initiates a group key request for a secure communication with the application servers  $S_j$ . To do this, it will encrypt the request message with a key,  $K_{C,S}$ , that the client  $C$  invents. To make sure that only  $S_j$  can get the key (of course,  $KDC_L$  and  $KDC_R$  should be able to get it as well), the key will be encrypted using the client  $C$ 's key along with the destination,  $KDC_R: S_1,$

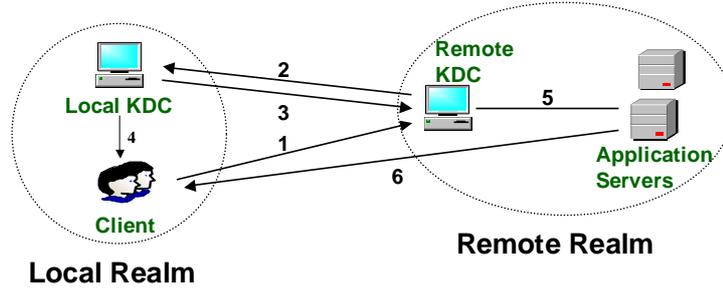


Figure 6.8: The Message Flow of the New Technique in a Single Remote Realm

...,  $S_n$ . For  $S_j$  to get  $K_{C,S}$ , the  $KDC_R$  will first need to be authenticated by the KDC of the client  $C$ , and then  $S_j$  will need to be authenticated by their KDC:  $KDC_R$  and have  $KDC_R$  give  $S_j$  the key,  $K_{C,S}$ .  $KDC_R$  will be responsible for making sure that  $S_j$  are valid recipients and the client  $C$  is a valid sender validated by the client  $C$ 's KDC,  $KDC_L$ . A common nonce,  $N_C$ , will also be invented to track the transaction and give the recipient something to authenticate with to the remote KDC server,  $KDC_R$ . A timestamp  $TMS_C$  is included to prevent replay attacks or prevent the key,  $K_{C,S}$ , from being given out more than once.

**Message 2,  $KDC_R \rightarrow KDC_L$ :**

$N_C$ , "C",  $K_C\{N_C, K_{C,S}, "KDC_R : S_1, \dots, S_n", TMS_C\}$ ,  $\{N_C, N_{KDC_R}, TMS1_R, \text{AuthInfo}\}_{KDC_L}$ , where AuthInfo consists of "KDC\_R",  $Cert_{KDC_R}$ ,  $[N_C, N_{KDC_R}, TMS1_R, "KDC_L",  $KB_{KDC_R}]_{KDC_R}$ .$

After receiving the message from the client  $C$ ,  $KDC_R$  invents a nonce,  $N_{KDC_R}$ , and generates AuthInfo necessary to authenticate  $KDC_R$ . Then, it constructs a message to send to the client  $C$ 's local KDC that contains  $N_{KDC_R}$ ,  $N_C$ ,  $TMS1_R$ , and AuthInfo, encrypted with  $KDC_L$ 's public key.  $KDC_R$  will also forward the part of the original message encrypted with  $C$ 's key. This is enough information for  $KDC_L$  to authenticate the client  $C$  and  $KDC_R$  to make sure only if  $KDC_L$  can decrypt this message. Note that the message uses "KDC\_R" instead of "KDC\_R:  $S_1, \dots, S_n$ " for the sake of a privacy consideration. This is because the client  $C$ 's local KDC is unnecessary to know who will be part of the secure communication except the client  $C$ . Additionally, this message does not explicitly contain the identity  $KDC_R$  since it is included in AuthInfo. "KDC\_R" and its public key  $KB_{KDC_R}$  are uniquely determined by  $Cert_{KDC_R}$  in AuthInfo. Thus,  $Cert_{KDC_R}$  is used to prevent man-in-the-middle attacks.  $TMS1_R$  serves to avert replay attacks.

**Message 3,  $KDC_L \rightarrow KDC_R$ :**

$\{"C", "KDC_R", N_{KDC_R}, K_{C,S}, Cert_{KDC_L}, TMS_C, TMS1_R\}_{KDC_R}$

Using the key of the name given,  $KDC_L$  will decrypt the message encrypted with its public key and the original message encrypted by  $K_C$ , and verify  $KDC_R$ 's signature. Then, it will check if the  $N_C$  encrypted by key  $K_C$  matches the  $N_C$  encrypted by its public key. If they match, then the client  $C$  is really the client  $C$  and its message is not altered.  $KDC_L$  will then read the destination, in this case, " $KDC_R$ ." (Again, it does not include " $KDC_R: S_1, \dots, S_n$ " in the message for the sake of a privacy consideration.) Furthermore,  $KDC_L$  will make sure if  $KDC_R$  is a valid member. If so, it continues the processing. Accordingly,  $KDC_L$  will make sure that the timestamp is within the allowed clock-skew and that the key for this request has not already been given out.  $KDC_L$  will then give out the key,  $K_{C,S}$ .  $KDC_L$  will also encrypt  $N_{KDC_R}$  with  $KDC_R$ 's public key to authenticate  $KDC_L$  to  $KDC_R$ .  $KDC_R$  will verify that the  $N_{KDC_R}$  received from  $KDC_L$  matches the  $N_{KDC_R}$  sent to  $KDC_L$ , and it will distribute the returned  $K_{C,S}$  to  $S_j$ .

**Message 4,  $KDC_L \rightarrow C$ :**

$K_C\{K_{C,S}, "KDC_L", "FromKDC_R", N_C, TMS_C, TMS_L\}$

The client decrypts the share key  $K_{C,S}$ , " $FromKDC_R$ ",  $N_C$ ,  $TMS_C$ , and  $TMS_L$ . The client needs to make sure that the  $N_C$  matches the nonce of a pending request, the timestamp is within the allowed clock-skew, and  $K_{C,S}$  matches the group key that was previously created. Using timestamps and nonce numbers makes the encrypted message random to some extent, and thus prevents a brute-force cryptographic attack.

**Message 5,  $KDC_R \rightarrow S_j$ :**

$K_{S_j}\{ "C", K_{C,S}, K_{C,S}\{N_C, TMS_C, TMS_{2R}\}$

The application servers  $S_j$  will decrypt the first message to get  $K_{C,S}$ , and then use the shared key  $K_{C,S}$  to decrypt the second message to get  $TMS_C$  and  $TMS_{2R}$ .  $TMS_C$  is included so that  $S_j$  knows when the group key was inverted by  $C$ . To make sure that it is not a replay,  $S_j$  should keep all timestamps that were recently received, say in the last five minutes that is a parameter set approximately for the maximum allowable time skew. Then,  $S_j$  should check that each received timestamp from a given request initiator is different from any of the stored values. Any authentication request older than five minutes (or whenever the value of the maximum allowable time skew) would be rejected anyway, so  $S_j$  would not remember values older than 5 minutes.

**Messages 6,  $S_j \rightarrow C$ :**

$K_{C,S}\{N_C, TMS_C\}$ , where  $j = 1, 2, \dots, n$ .

The application servers send replies to the client. Then, by using pre-generated  $K_{C,S}$ , the client decrypts the message to make sure that the group key  $K_{C,S}$  is not altered. It also verifies nonce

$N_C$  and makes sure the received timestamp of the pending request is within the allowed clock-skew of the timestamp  $TMS_C$ .

In the new technique, note that  $K_{C,S}$  can be replaced by  $K_{C,S_j}$  when  $C$  only wants to securely communicate with any one of application servers  $S_j$ . The aforementioned technique can be also modified so that a reply is issued to  $C$  by either  $KDC_L$  or  $KDC_R$  whenever needed after  $C$  is authenticated.

### 6.3.2 Multiple Remote Realms

In a large network, application servers are often within different network domains. For instance, in Web services service providers may be partitioned into many different realms due to their location flexibility. To work together, a service provider may wish to gain access to other service providers' application servers in remote realms. To support "cross-realm" authentication, the service provider's KDC needs to establish an either direct or indirect trust relationship with the other service providers' KDCs. Here, we briefly describe how the proposed technique is extended in multiple remote realms.

Assume that the  $n$  application servers  $S_j$  are distributed in  $m$  realms, each with KDC servers denoted by  $KDC_i$  ( $i = 1, \dots, m$ ). Let  $n_i$  be the number of application servers within  $KDC_i$ , where  $0 \leq n_i \leq m$  and  $\sum_{i=1}^m n_i = n$ . Clearly, the case of a single remote realm is associated with  $n_1 = n$  and  $n_i = 0$  ( $i = 2, \dots, m$ ). We further let  $S_{i,k}$  be those application servers which belong to the set  $\{S_j | j = 1, \dots, n\}$  within realm  $i$ , where  $k = 1, \dots, n_i$ . Then, the set  $\cup_{i=1}^m \{S_{i,k} | k = 1, \dots, n_i\}$  is identical to the set  $\{S_j | j = 1, 2, \dots, n\}$ . Table 6.5 lists additional notation needed in this section. Figure 6.9 shows how the client authenticates to application servers. Authentication messages are given below.

**Messages 1 and 1'**,  $C \rightarrow KDC_i$  for each fixed  $i$ :

$N_C, "C", "KDC_i: S_{i,1}, \dots, S_{i,n_i}", K_C\{N_C, K_{C,S}, "KDC_i: S_{i,1}, \dots, S_{i,n_i}", TMS_C\}$

**Messages 2 and 2'**,  $KDC_i \rightarrow KDC_L$ :

$N_C, "C", K_C\{N_C, K_{C,S}, "KDC_i: S_{i,1}, \dots, S_{i,n_i}", TMS_C\}, \{N_C, N_{KDC_i}, TMS1_i, \text{AuthInfo}_i\}_{KDC_L}$ , where  $\text{AuthInfo}_i$  consists of " $KDC_i$ ",  $Cert_{KDC_i}$ , and  $[N_C, N_{KDC_i}, TMS1_i, "KDC_L", KB_{KDC_i}]_{KDC_i}$

**Messages 3 and 3'**,  $KDC_L \rightarrow KDC_i$  where  $i = 1, \dots, n$ :

$\{ "C", "KDC_i", N_{KDC_i}, K_{C,S}, Cert_{KDC_L}, TMS_C, \}_{KDC_i}$

**Messages 4 and 4'**,  $KDC_L \rightarrow C$ :

Table 6.5: The Additional Notation Needed in the Case of Multiple Remote Realms

$S_{i,k}$	Application Server $k$ in realm $i$
$KDC_i$	Remote KDC server in realm $i$
$K_{S_{i,k}}$	Secret key of $S_{i,k}$
$KB_{KDC_i}$	Public key of $KDC_i$
$\{\text{Message}\}_{KDC_i}$	Message encrypted with $KDC_i$ 's public key
$[\text{Message}]_{KDC_i}$	Message signed with $KDC_i$ 's private key
$N_{KDC_i}$	Nonce generated by $KDC_i$
$Cert_{KDC_i}$	Certificate of $KDC_i$
$TMS1_i$ and $TMS2_i$	Timestamps generated by $KDC_i$

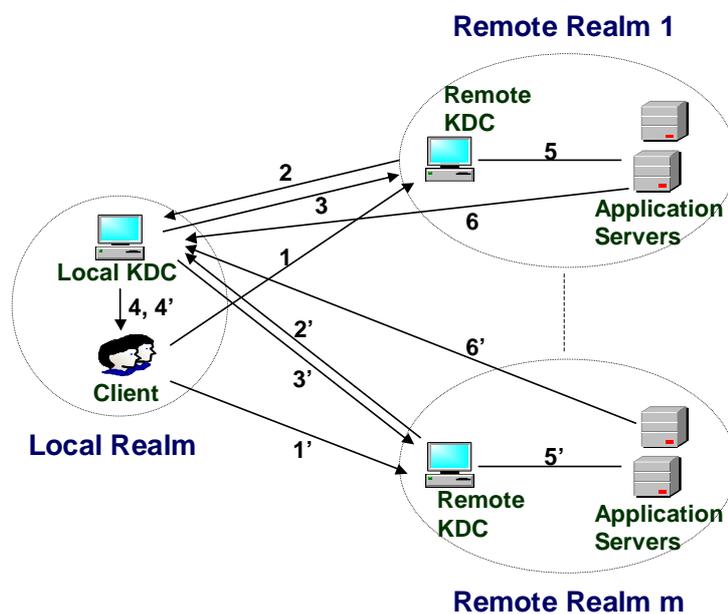


Figure 6.9: The Message Flow of the New Technique in Multiple Remote Realms

Table 6.6: The Operations of Encryption and Decryption When  $n$  Application Servers are in  $m$  Remote Realms

Entities	# Secret Key	# Private Key	# Public Key
Client	3	0	0
Local KDC	2	1	$4m$
Remote KDC	$n+1$	$2m$	$3m$
Application server	$3n$	0	0
Total	$4n+6$	$2m+1$	$7m$

$K_C\{K_{C,S}, "KDC_L", "FromKDC_i", N_C, TMS_C, TMS_L\}$

**Messages 5 and 5'**,  $KDC_i \rightarrow S_{i,k}$ :

$K_{S_{i,k}}\{"C", K_{C,S}\}, K_{C,S}\{N_C, TMS_C, TMS_{2i}\}$ , where  $i = 1, \dots, m$  and  $k = 1, \dots, n_i$ .

**Messages 6 and 6'**,  $S_{i,k} \rightarrow C$ :

$K_{C,S}\{N_C, TMS_C\}$ , where  $i = 1, \dots, m$  and  $k = 1, \dots, n_i$ .

In the first three messages, the client  $C$  authenticates to her local KDC through remote  $KDC_i$ . Message 4 or 4' checks if the group key  $K_{C,S}$  is valid, and Message 5 or 5' distributes the group key  $K_{C,S}$  and authenticates designated application servers within their individual KDC realms. An explanation of these messages are similar to the one in Section 6.3.1.

## 6.4 Performance Evaluation of the New Proposed Technique

In this section we focus on studying the efficiency of the proposed group authentication technique. We first compute its computational and communication costs. Then, we give a thorough performance evaluation of the new technique using the queueing method proposed in Section 6.2.

### 6.4.1 The Operations of Encryption and Decryption

The baseline transactions are constructed with one or more application servers in a remote realm as shown in Figure 6.8. Without any confusion, we directly consider the case of  $m$  remote realms, each with  $n_i$  application servers where  $0 < n_i \leq n$  and  $i = 1, 2, \dots, m$ . Table 6.6 summarizes the number of encryption and decryption operations performed in the proposed technique. As is shown in Tables 6.1 and 6.6, in our technique the computational burden on the client is mitigated to its local KDC compared to PKTAPP. Specifically, in term of the calculation of public key opera-

tions the burden on the client's local KDC is  $O(m)$  in both the proposed technique and PKCROSS but the burden on the client is  $O(n)$  in PKTAPP. Let us recall that  $m$  is the number of remote KDC servers and  $n$  is a total number of application servers where  $1 \leq m \leq n$ . Hence, their computational burdens may be significantly different when  $m \ll n$ . Note that the proposed technique uses public key cryptography to authenticate the client's KDC and the remote KDCs of application servers. So, it reduces the risk of a single failure on these KDCs.

Next, let us consider the operation costs of the proposed technique. Denote by  $f_3(n, m)$  the total computational time of its encryption and decryption operations. Then, it follows from Table 6.6 that  $f_3(n, m)$  is computed by

$$f_3(n, m) = 4c_1n + (2c_2 + 7c_3)m + 6c_1 + c_2$$

Thus, we have the following proposition.

**Proposition 6.4.1** (a) *The proposed technique requires less computational time than PKCROSS.*

(b) *For  $n \geq 4$ , the proposed technique requires less computational time than PKTAPP. But, when  $1 \leq n < 4$ , the proposed technique requires less computational time PKTAPP if and only if the number of remote realms  $m$  is less than  $\lfloor n + \frac{(n-4)c_1}{2c_2+7c_3} \rfloor$ .*

*Proof.* The differences of computational times among the three techniques are given by  $f_1(n, m) - f_3(n, m) = -4c_1 + (c_1 + 2c_2 + 7c_3)n - (2c_2 + 7c_3)m$  and  $f_2(n, m) - f_3(n, m) = 4c_1m + 5c_1$ . Obviously, our technique requires less computational time than PKCROSS due to  $f_2(n, m) - f_3(n, m) > 0$ . Moreover,  $f_1(n, m) - f_3(n, m) \geq 0$  if and only if  $m \leq n + \frac{(n-4)c_1}{2c_2+7c_3}$ , which implies (b) due to  $n \geq m$ . This proves Proposition 6.4.1.

Similarly, we can easily get that

**Proposition 6.4.2** *For  $m \geq 4$ , the proposed technique requires less computational time than PKTAPP. But, when  $1 \leq m < 4$ , the proposed technique requires less computational time PKTAPP if and only if the number of application servers  $n$  should be more than  $\lceil m - \frac{(m-4)c_1}{c_1+2c_2+7c_3} \rceil$ .*

Furthermore, let  $g_3(n, m)$  be the transaction time of the proposed technique, i.e., the computational time of its encryption and decryption operations plus its communication time. Note that the number of WAN communications required in the technique is  $3m+2n$ , and recall that  $d$  is the time spent in a WAN communication. Then, the following statements hold.

Table 6.7: The Minimal Number of Application Servers

# Servers	Number of Remote Realms									
	1	2	3	4	5	6	7	8	9	10
d=0.12	2	3	4	5	6	7	8	9	10	11
d=4.8	2	3	4	5	7	8	9	10	12	13
d=10	2	3	4	6	7	8	10	11	12	14

**Proposition 6.4.3** (a) *The proposed technique has less transaction time than PKCROSS. (b) The proposed technique uses less transaction time than PKTAPP if and only if the number of application servers  $n$  should be more than  $\lceil m + \frac{(d-c_1)m+4c_1}{c_1+2c_2+7c_3+2d} \rceil$ .*

*Proof* given by  $g_1(n, m) - g_3(n, m) = [f_1(n, m) - f_3(n, m)] + (2n - 3m)d = (c_1 + 2c_2 + 7c_3 + 2d)n - (2c_2 + 7c_3 + 3d)m - 4c_1$ , and  $g_2(n, m) - g_3(n, m) = [f_2(n, m) - f_3(n, m)]$ . These easily derive (a)-(b) in this proposition. The proof is complete.

Note that when  $n = 1$ , we get that  $m + \frac{(d-c_1)m+4c_1}{c_1+2c_2+7c_3+2d} > 1 = n$  due to usually  $d > c_1$ , which imply  $g_1(n, m) - g_3(n, m) < 0$ . Also, when  $m = 1$  and  $n \geq 2$ ,  $m + \frac{(d-c_1)m+4c_1}{c_1+2c_2+7c_3+2d} = 1 + \frac{d+3c_1}{c_1+2c_2+7c_3+2d} < 2 \leq n$  due to  $c_1 < c_3$ , which imply  $g_1(n, m) - g_3(n, m) > 0$ . Thus, we have that

**Corollary 2** *PKTAPP is more efficient than the proposed technique when  $n = 1$ , but less efficient when  $m = 1$  and  $n \geq 2$ , in term of transaction time.*

Using Proposition 6.4.3 we calculated the minimal number of application servers so as to ensure that our technique requires a lower transaction time than PKTAPP with varied  $d = 0.12$  msec, 4.8 msec and 10 msec when  $c_1 = 0.000248$  msec,  $c_2 = 0.07$  msec and  $c_3 = 1.52$  msec. The results are shown in Table 6.7. We observe that the minimal number of application servers is sensitive to  $d$  rather than  $c_j$  ( $j = 1, 2, 3$ ). Table 6.8 further presents the difference of transaction times between our technique and PKCROSS. We also noted that our proposed technique requires significantly less transaction time than PKCROSS, and the difference of transaction times between the two techniques are independent of the number of application servers. That is, the proposed technique is more efficient than PKCORSS.

Table 6.8: The Difference of Transaction Times Between Our Technique and PKCROSS (i.e.,  $g_2(n, m) - g_3(n, m)$ )

(a)  $m=1, 2, 3, 4, 5$

	Number of Remote Realms				
Difference (msec)	1	2	3	4	5
d=0.12	0.1222	0.2432	0.3642	0.4852	0.6062
d=4.8	4.8022	9.6032	14.404	19.205	24.006
d=10	10.002	20.003	30.004	40.005	50.006

(b)  $m=8, 12, 16, 20, 24$

	Number of Remote Realms				
Difference (msec)	8	12	16	20	24
d=0.12	0.9691	1.45314	1.93711	2.42108	2.90504
d=4.8	38.409	57.6131	76.8171	96.0210	115.225
d= 10	80.009	120.013	160.017	200.021	240.025

#### 6.4.2 The Calculation of the Response Time via a Queueing Network

Similar to Section 6.2, we can use a queueing network to characterize the message flow of the proposed technique where each system resource is modeled as a queue associated with a queueing discipline shown in Figure 6.10. Figures 6.11, 6.12 and 6.13 show response times as a function of authentication request rates, i.e., throughput, in the case of one, two or eight remote realms with varying number of application servers. As can be seen, our technique performs better than PKCROSS in all cases. It has also been demonstrated that our technique is more efficient than PKTAPP when  $n \geq 2$  for one remote realm in Figure 6.11, when  $n \geq 4$  for two remote realms in Figure 6.12, and when  $n \geq 12$  for eight remote realms in Figure 6.13 (roughly  $n \geq 1.5m$  if  $m > 1$ ). These crossover numbers from  $m = 1$  to 12 are further depicted in Figure 6.14 where the crossover numbers are 99.6% perfectly fitted by the straight line:  $n = 1.4406m + 0.6364$ , which predicts the number of application servers required to ensure that our technique performs better than PKTAPP. Moreover, the line  $n = 1.4406m + 0.6364$  is below the straight line  $n = 1.8916m + 0.2879$  given in Section 6.2 due to  $1.8916m + 0.2879 > 1.4406m + 0.6364$  when  $m \geq 1$ . This again confirms that our technique is more efficient than PKCROSS. We also noted that the numbers  $n$  in Table 6.7 are smaller than the crossover numbers in Figure 6.14. This means that it is not sufficient *only if* transaction time is employed to analyze performance that researchers have often used.

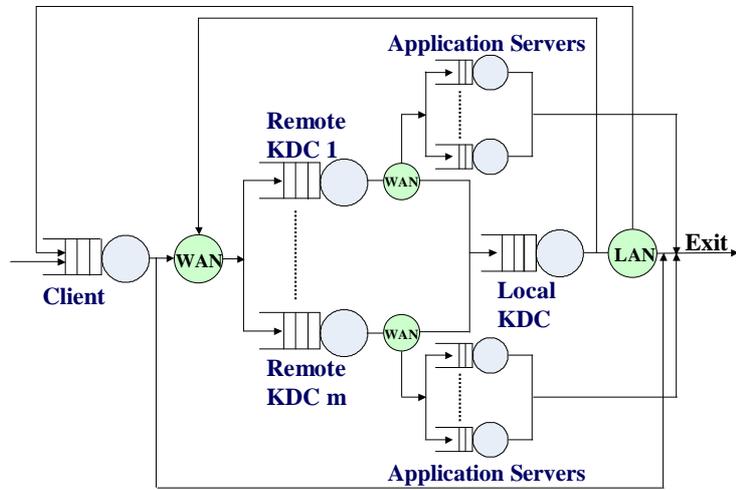


Figure 6.10: A Queuing Network with  $n$  Remote Realms for the New Technique

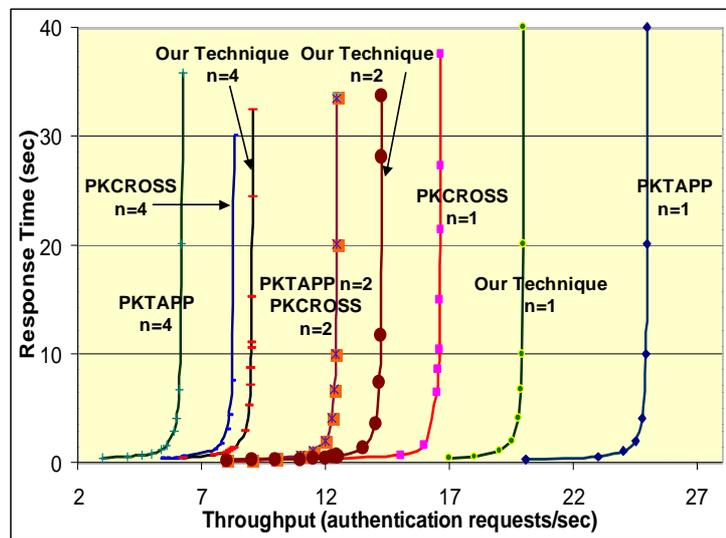


Figure 6.11: Response Times vs. Authentication Request Rates When  $m = 1$

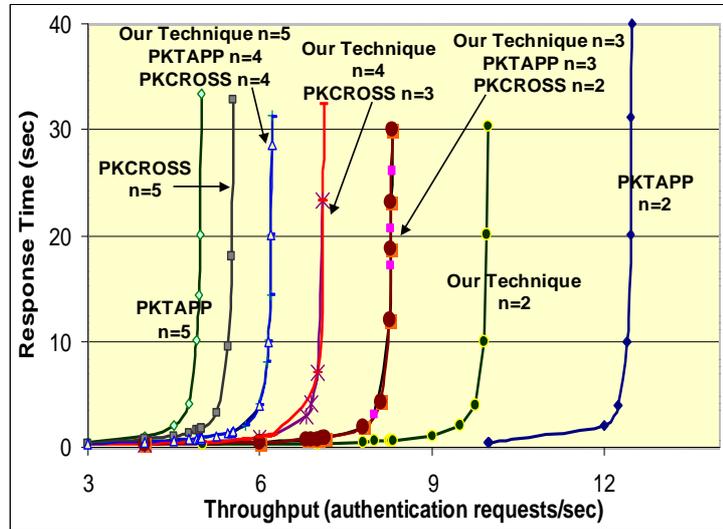


Figure 6.12: Response Times vs. Authentication Request Rates When  $m = 2$

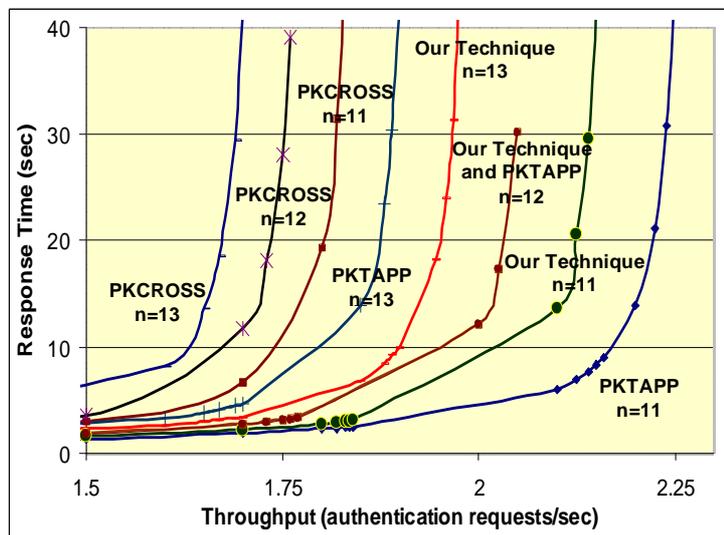


Figure 6.13: Response Times vs. Authentication Request Rates When  $m = 8$

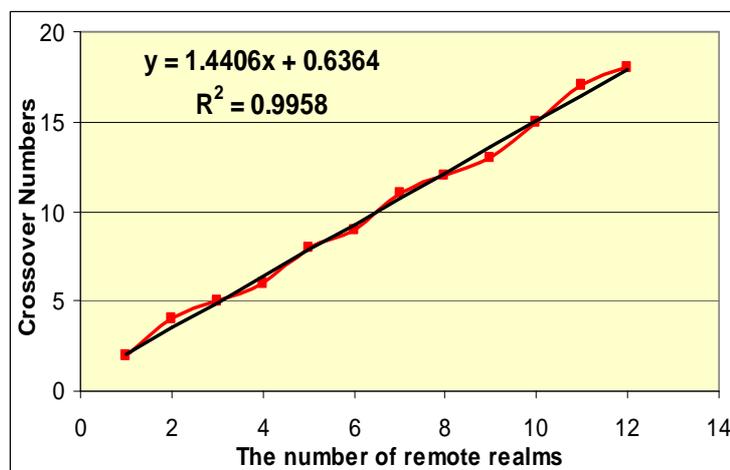


Figure 6.14: Crossover Numbers vs. Remote Realms

### 6.4.3 Discussions

The proposed technique is conceptually simple as it involves the client that generates a key and the KDCs that authenticate and distributes the key to group members via a pair-wise secure channel established with each group member. It works well in one-to-many multicast scenarios, for example, secure communications among service providers in Web services.

While our technique has several strengths as compared to PKCROSS and PKTAPP, we have also realized that the client's request message passes through the remote KDC to the client's local KDC with no verification. It requires to securely migrate and maintain an authentication state at every remote KDC, and might result in serious DOS attacks. But, our technique is motivated by secure communications in Web services applications and it is mainly toward to e-business applications in which a trust-but-verify framework for Web services authorization has been proven an efficient method [86]. So, under the same framework we trust a client in a certain degree and allow the client to contact the KDC server of the application servers directly. Then, the authentication procedures proposed in this chapter will follow. However, if the remote KDC server is badly under DOS attacks due to unauthenticated messages from the client, then we leave the option to switch back to the way in which a client needs to be first authenticated by its local KDC, e.g., using PKCROSS.

Additionally, in the proposed technique, the client is allowed to create the group key  $K_{C,S}$ . This is particularly useful when the client does not need to get a reply from its local KDC for a further verification. If the reply is required, then we can easily mitigate the creation of  $K_{C,S}$  from

the client to its local KDC by slightly modifying the proposed technique accordingly.

## 6.5 Concluding Remarks

Public-key enabled Kerberos-based techniques such as PKINIT, PKCROSS, PKDA and PKTAPP (a.k.a KX.509/KCA) give a potential effective way for cross-realm authentication. However, the authentication problem is simple to describe but hard to solve, especially for multiple realms. Their performance has been poorly understood. The chapter presented a throughout performance evaluation of PKCROSS and PKTAPP in terms of computational and communication times, and through the use of validated analytical queueing models. Our analysis revealed that PKTAPP does not perform better than PKCROSS in a large network where there are many application servers within either a single or multiple remote realms. Thus, we proposed a new public key cryptography-based group authentication technique. Our performance analysis has showed that the proposed technique outperforms PKCROSS in terms of computational and communication times in Propositions 6.4.1 and 6.4.3, and response time in Figures 6.11, 6.12 and 6.13. The chapter also gave the predicted minimal number of application servers so as to ensure that the proposed approach is more efficient than PKTAPP in multiple remote realms. Roughly speaking, the proposed technique performs better than PKTAPP when the number of application servers is about 50% more than the number of remote realms (i.e.,  $n \geq 1.5m$ ) in Section 6.4.2. As shown, our performance methodology based on complexity analysis and queueing theory is an effective way to analyze the performance of security protocols.

## Chapter 7

# Summary and Future Work

This chapter gives the research summary of the dissertation and provides research directions for future study.

### 7.1 Research Summary of the Dissertation

In this dissertation we considered a collection of computer resources used by a service provider to host enterprise applications for business customers. An enterprise application running in such a computing environment is associated with a service level agreement (SLA). That is, the service provider is required to execute service requests from a customer within negotiated quality of service (QoS) requirements for a given price. The QoS requirements may include service availability, performance and security, and the service requests may come from either single class customers or multiple class customers. It is assumed that the service provider owns and controls a number of resource stations or sites. In Chapter 1 we gave an overview of service availability, performance and security, and described a resource optimization problem subject to these QoS metrics, and a related group authentication problem in the distributed computing. A literature review was given in Chapter 2.

Chapter 3 presented an approach for service optimization in such an enterprise application environment that minimizes the total cost of computer resources required while satisfying an SLA for single class customers. We considered a response time in the QoS. Typically, in the literature,

the response time is taken into account through its mean. However, this may not be of real interest to a customer. In view of this, we studied percentiles of the response time, and calculated the number of servers in each resource station that minimize a cost function that reflects operational costs. Numerical results showed the applicability of the proposed approach and validated its accuracy in the case of single class customers. We found that there is a relative error less than 0.71% in the cumulative distribution function of response time calculated by between our approximate algorithm and a simulation method for the service models under our study. We also found that the number of servers in each resource station obtained by our approximate algorithm is in fact optimal, provided that each station has balanced utilization.

However, in enterprise computing customer requests often need to be distinguished, with different request characteristics and customer's service requirements. In this dissertation, we further considered a set of computer resources used by a service provider to host enterprise applications for differentiated customer services subject to a percentile response time and a fee, as given in Chapter 4. We proposed algorithms to solving the resource optimization problem in the case of multiple class customers. The accuracy of our algorithms was verified by numerical simulation. A contributing factor to our algorithms' accuracy is that typically we are interested in values of the cumulative distribution of the response time that correspond to very high percentiles for which the approximate results seem to have a very good accuracy.

With the number of e-Business applications dramatically increasing, the SLA will play an important part in Web services. It is often a combination of several quality of services (QoS), such as security, performance, and availability, agreed between a customer and a service provider. Most existing research addresses only one of these QoS metrics.

In Chapter 5, we presented a study of three QoS metrics, namely, trustworthiness, percentile response time and availability. Then, we considered all these three QoS metrics in a trust-based resource provisioning problem which typically arises in Web services, and formulated the trust-based resource provisioning problem as an optimization problem under SLA constraints in the cases of single class customers and multiple class customers respectively. By the use of an effective and accurate numerical solution for the calculation of the percentile response time presented in Chapters 3 and 4 we solved the optimization problem using an efficient numerical procedure in the both cases.

Chapter 6 presented a thorough performance evaluation of PKCROSS and PKTAPP (a.k.a. KX.509/KCA) in terms of computational and communication times. Then, we demonstrated their performance difference using open queueing networks. An in-depth analysis of these two techniques

showed that PKTAPP does not scale better than PKCROSS. Thus, we proposed a new public key cryptography-based group authentication technique. Our performance analysis demonstrated that the new technique can achieve better scalability as compared to PKCORSS and PKTAPP.

## 7.2 Future Research Directions

In this section we give an extensive discussion of our approach for solving the resource optimization problem and the problem of public key cryptography-based group authentication for our future study.

**The Selection of Service Sites:** We introduce a trust manager who is a trusted agent representing customers. We assume that the trust manager is aware of the service sites that process its service request. Moreover, the trust manager is able to maintain the trustworthiness information of the service sites. A customer, however, does not have this information because the trust manager is its representative who deals with all service matters including an SLA negotiation, and service requests' submission and processing.

As shown in Figure 5.2, the service processor uses a composition of services selected from these service sites to complete the customer's request. These service sites are chosen by the trust manager based on the trust information of each site. In our proposed algorithms, the trust manager keeps the ranking of all service sites. When the trust manager receives a service request from a customer, it will check the trustworthiness information of service sites in its database, and then select those service sites with the highest indices who meet trust requirements predefined by its customers as well. This selection maximizes the customer's benefit in term of trustworthiness.

However, the higher reliable service site usually requires more security measures so as to provide better the quality of service. Thus, the higher trust index the service site, the higher cost the service. The trust manager on the other hand may only choose those service sites who meet predefined trust requirements with less costs. In this case, the trust manager may require the service broker to find those service sites who meet predefined trust requirements but have the lowest total cost for the customer's request. The service broker can find those

sites by modifying (1.1) into the following optimization

$$\begin{aligned} & \min_{\substack{\forall s_j \in S; I_{s_j} \leq \hat{I}_j \\ j = 1, \dots, m}} \min_{n_1, \dots, n_m} (n_1 c_1 + \dots + n_m c_m) \end{aligned} \quad (7.1)$$

and solving the problem subject to constraints by an SLA, where  $S$  is a set of service sites consisting of  $S_1, S_2, \dots, S_M$ ,  $s_j$  ( $j = 1, \dots, m$ ) are any  $m$  of those service sites whose trust indices  $I_{s_j}$  are less than  $\hat{I}_j$ , as predefined by the customer. This approach maximizes the profitability of the trust manager. It may be a typical case in the real world, for example, when the trust manager and the service broker may be integrated into only one entity who represents both the customer and the service sites.

In this dissertation, we only considered the case in which the trust manager select those service sites with the highest indices, as used in Sections 3.3 for the best benefit of the customer. An extension to study the above alternative case will be of interest.

**The Relationship Among Trustworthiness, Response Time and Availability:** The trust indices of resource sites is an indicator of the quality of service provided by these sites, and it is based on previous and current job completion experience. Clearly, when the response time does not meet a predefined requirement, or resource sites are not able to serve a customer request, the trust manager will give these resource sites low trust indices.

Similarly, the longer time the service at a resource site is unavailable, the longer the end-to-end response time becomes. Particularly, when resource sites are under denial of service attacks, the end-to-end response time of a service request may become infinite. Hence, when a service request suffers from a much longer response time, it may mean that some of resource sites cannot meet the requirement of service availability. In other words, these three SLA metrics, trustworthiness, response time, and service availability, interact each other, and their relationship is very complicated. A throughout study of these three metrics will be conducted in our future research.

**The Trust Updating Approach and the Server Repairing Mechanism:** In our approach to solving the trust-based resource provisioning problem, we used a rank and a threshold-based approach to deal with the trust indices of resource sites, and assumed that a service discipline was FIFO and each resource site could intermediately repair a server. In our future study

we will extend our discussion by evaluating different trust approaches and using other service disciplines with a general repairing mechanism for secure resource management in Web service applications and other distributed computing environments.

**The Study of Other Security Mechanisms:** As is shown, our performance methodology, based on complexity analysis and queueing theory presented in Chapter 6, is an effective way to analyze the performance of security protocols. In the future, we plan to apply the methodology to other security protocols. Particularly, it will be interesting to apply our method to those protocols which are used in either time-sensitive or resource-limited applications, e.g., those in wireless sensor networks.

## Bibliography

- [1] Issam Aib, Nazim Agoulmine, and Guy Pujolle. The generalized service level agreement model and its application to the SLA driven management of wireless environments. In *International Arab Conference on Information Technology*, ACIT, December 2004.
- [2] J. Aikat, J. Kaur, FD Smith, and K. Jeffay. Variability in tcp round-trip times. In *Proceedings of the ACM SIGCOMM Internet Measurement Conference*, October 2003.
- [3] A. Aiyer, L. Alvisi, and R. Bazzi. On the availability of non-strict quorum systems. In *Proceedings of the 19th International Symposium on Distributed Computing (ISDC 2005)*, pages 48–62, September 2005.
- [4] M. Allman, W. Eddy, and S. Ostermann. Estimating loss rates with TCP. *ACM Performance Evaluation Review*, 31(3), December 2003.
- [5] Tayfur Altioek and Benjamin Melamed. *Simulation Modeling and Analysis with Arena*. Cyber Research, Inc. and Enterprise Technology Solutions, Inc., 2001.
- [6] J. Altman. Nist pki'06: Integrating pki and kerberos. In [www.secure-endpoints.com/talks/nist-pki-06-kerberos.pdf](http://www.secure-endpoints.com/talks/nist-pki-06-kerberos.pdf), 2007.
- [7] J. Altman. Personal communication. In [www.secure-endpoints.com/talks/nist-pki-06-kerberos.pdf](http://www.secure-endpoints.com/talks/nist-pki-06-kerberos.pdf), 2007.
- [8] M. Barbacci, M. H. Klein, T. H. Longstaff, and C. B. Weinstock. Quality attributes. In *CMU Technical Report: CMU/SEI-95-TR-021, 1995*. Software Engineering Institute, Carnegie Mellon University, 1995.

- [9] Douglas K Barry. *Web Services and Service-Oriented Architecture: Your Road Map to Emerging IT*. Morgan Kaufmann, 2003.
- [10] M. Bishop. *Computer Security*. Addison Wesley, Boston, MA, 2002.
- [11] M. Blaze. *High-Bandwidth Encryption with Low-Bandwidth Smartcards*. Lecture Notes in Computer Science, Number 1039, 1996.
- [12] Gunter Bolch, Stefan Greiner, Hermann Meer, and Kishor Shridharbhai Trivedi. *Queueing Networks and Markov Chains: Modeling and Performance Evaluation with Computer Science Application*. John Wiley & Sons, New York, 1998.
- [13] J. Bolot. End-to-end packet delay and loss behavior in the Internet,. In *Proceedings of the ACM SIGCOMM*, 1993.
- [14] E. Bouillet, D. Mitra, and K. Ramakrishnan. The structure and management of service level agreements in networks. *IEEE Journal on Selected Areas in Communications*, 20(4):691–699, 2002.
- [15] E. A. Brewer. Lessons from giant-scale service. *IEEE Internet computing*, pages 46–55, July-August 2001.
- [16] A. Brown and D. Patterson. Towards availability benchmarks: A case study of software RAID systems. In *Proceedings of 2000 USENIX Annual Technical Conference*. USENIX, June 2000.
- [17] S. Bruell and G. Balbo. *Computational Algorithms for Closed Queueing Networks*. Edited by P. J. Denning, Science Library, Elsevier North Holland, Inc., New York, 1980.
- [18] P. Burke. Output processes and tandem queues. In *Proceedings of the Symposium on Computer Communication Networks and Teletraffic*. Brooklyn, April 1972.
- [19] E. Lieurain D. Salomon L. Winkelbauer B. Jacob S. Mui J. Pannu S. Park H. Raguette J. Schneider C. Matthys, P. Bari and L. Vanel. *On Demand Operating Environment: Managing the Infrastructure (Virtualization Engine Update)*. IBM Redbooks, June, 2005.
- [20] Joel Calabrese. Optimal workload allocation in open queueing networks in multiserver queues. *Management Science*, 38(12):1792–1802, December 1992.

- [21] A. Chandra, W. Gong, and P. Shenoy. Dynamic resource allocation for shared data centers using online measurements. In *Proceedings of Eleventh International Conference on Quality of Service (IWQoS 2003)*, June 2003.
- [22] J. Chang. Processor performance, Update 1. In <http://www.sql-server-performance.com>, 2005.
- [23] C Chassot, F Garcia, G Auriol, A Lozes, E Lochin, and P Anelli. Performance analysis for an IP differentiated services network. In *Proceedings of IEEE International Conference on Communication (ICC'02)*, pages 976–980, 2002.
- [24] CISCO. Increasing network availability. In <http://www.cisco.com/warp/public/779/largeent/learn/technologies/ina/IncreasingNetworkAvailability-Abstract.pdf>.
- [25] CISCO. Network availability: How much do you need? how do you get it? In <http://www.cisco.com/en/US/netsol/ns206>.
- [26] CISCO. Service level management: Best practice. In <http://www.cisco.com/en/US/tech/tk869/tk769>.
- [27] J. Cohen. *The Single Server Queue*. North-Holland, Amsterdam, New York, Oxford, 1982.
- [28] The MIT Kerberos Consortium. Proposal for corporate sponsors. In <http://www.kerberos.org/join/proposal.pdf>, 2007.
- [29] F. Curbera, M. Duftler, R. Khalaf, W. Nagy, N. Mukhi, and S. Weerawarana. Unraveling the Web services Web: An introduction to SOAP, WSDL, and UDDI. *IEEE Internet Computing*, 6(2):86–93, March-April 2002.
- [30] H. Daduna. Burke's theorem on passage times in Gordon-Newell networks. *Adv. Appl. Prob.*, 16, 1984.
- [31] W. Dai. Crypto++ 3.1 benchmarks. In <http://www.eskimo.com/~weidai/benchmark.html>, 2007.
- [32] D. Davis. Compliance defects in public-key cryptography. In *Proceedings of the Sixth USENIX UNIX Security Symposium (USENIX Security'96)*, San Jose, California, July 1996.

- [33] G. Dobson. Quality of service in service-oriented architectures. <http://digs.sourceforge.net/papers/qos.html>, 2002.
- [34] P. Dongara and T. N. Vijaykumar. Accelerating private-key cryptography via multithreading on symmetric multiprocessors. In *Proc. IEEE Int'l Symp. Performance Analysis of Systems and Software (ISPASS 03)*, pages 58–69, IEEE Press, 2003.
- [35] W. Doster, M. Watts, and D. Hyde. The kx.509 protocol. <http://www.citi.umich.edu/techreports/reports/citi-tr-01-2.pdf>, 2001.
- [36] S. Floyd, M. Handley, J. Padhye, and J. Widmer. Equation based congestion control for unicast applications. In *Proceedings of the ACM SIGCOMM*, 2000.
- [37] J. Garman. *Kerberos: The Definitive Guide*. O'Reilly, 2003.
- [38] D. Gaver. Observing stochastic processes, and approximate transform inversion. *Operation Research*, 14(3), 1966.
- [39] J. Golbeck and J. Hendler. Accuracy of metrics for inferring trust and reputation in semantic Web-based social networks. In *Proceedings of International Conference on knowledge engineering and knowlege management (EKAW)*, October 2004.
- [40] J. Gray. Dependability in the Internet era. In *Presentation at the High dependability computing consortium meeting*, <http://research.microsoft.com/~gray/talks>, May 2001.
- [41] K. Gummadi, S. Saroiu, and S. Gribble. King: Estimating latency between arbitrary Internet end hosts. In *Proceedings of the ACM SIGCOMM Internet Measurement Workshop*, November 2002.
- [42] A. Harbitter and D. Menasce. Performance of public-key-enabled Kerberos authentication in large networks. In *Proceedings of 2001 IEEE Symposium on Security and Privacy*, Oakland, California, 2001.
- [43] P. Harrison and W. Knottenbelt. Passage time distributions in large Markov chains. In *Proceedings of the ACM SIGMETRICS*, 2002.
- [44] L. He and J. Walrand. Pricing and revenue sharing strategies for internet service providers. In *Proceedings of the IEEE INFOCOM*, 2005.

- [45] J. Hennessy. The future of systems research. *IEEE Computer*, 32(8):27–33, August 1999.
- [46] M. Hur, B. Tung, T. Ryutov, C. Neuman, A. Medvinsky, G. Tsudik, and B. Sommerfeld. Public key cryptography for cross-realm authentication in Kerberos (pkcross). <http://tools.ietf.org/html/draft-ietf-cat-kerberos-pk-cross-07>, May 2001.
- [47] N. Duffield J. Sommers, P. Barford and A. Ron. Improving accuracy in end-to-end packet loss measurement. In *Proceedings of the ACM SIGCOMM*, August 2005.
- [48] R. Jurca and B. Faltings. Reputation-based service level agreements for Web services. In *Third International Conference on Service Oriented Computing (ICSOC 2005)*. Amsterdam, The Netherlands, December 2005.
- [49] S. D. Kamvar, M. T. Schlosser, and H. Garcia-Molina. EignRep: Reputation management in P2P networks. In *Proceedings of the World-Wide Web Conference*, 2003.
- [50] A. Keller and H. Ludwig. The WSLA framework: Specifying and monitoring service level agreements for Web services. *Journal of Network and Systems Management, Special Issue on "E-Business Management"*, 11(1):27–33, March 2003.
- [51] Y. Kim, A. Perrig, and G. Tsudik. Simple and fault-tolerant key agreement for dynamic collaborative groups. In *Proceedings of the 7th ACM Conference on Computer and Communications Security (ACM CCS 2000)*, pages 235 – 244. ACM, 2000.
- [52] J. Kohl and C. Neuman. RFC 1510: The Kerberos network authentication service (v5). <http://rfc.net/rfc1510.html>, 1993.
- [53] KX.509. Kx.509 source. <http://kx509.cvs.sourceforge.net/kx509/>, 2007.
- [54] J. Lee and R. Ben-Natan. *Integrating Service Level Agreements : Optimizing Your OSS for SLA Delivery*. Addison-Wesley Publishing Company, Reading, Massachusetts, 2002.
- [55] S. Lee, R. Sherwood, and B. Bhattacharjee. Cooperative peer groups in NICE. In *Proceedings of the IEEE INFOCOM*, 2003.
- [56] R. Levy, J. Nagarajao, G. Pacifici, M. Spreitzer, A. Tantawi, and A. Youssef. Performance management for cluster based Web services. In *The 8th IFIP/IEEE International Symposium on Integrated Network Management (IM2003)*. IEEE, 2003.

- [57] S. Chong M. Shin and I. Rhee. Dual-resource TCP/AQM for processing-constrained networks. In *Proceedings of the IEEE INFOCOM*, April 2006.
- [58] J. Martin and A. Nilsson. On service level agreements for IP networks. In *Proceedings of the IEEE INFOCOM*, June 2002.
- [59] D. S. Amr El Abbadi and F. Cristian. An efficient and fault-tolerant protocol for replicated data management. In *Proceedings of PODS*, 1985.
- [60] University of Houston. Network availability. In <http://www.telecomm.uh.edu/stats/NetStatus.html>.
- [61] A. Medvinsky, M. Hur, and C. Neuman. Public key utilizing tickets for application servers (PKTAPP). <http://tools.ietf.org/html/draft-ietf-cat-pktapp-00>, January 1997.
- [62] D. Menasce. QoS issues in Web services. *IEEE Internet Computing*, 6(4):72–75, November-December 2002.
- [63] D. Menasce. Response-time analysis of composite Web services. *IEEE Internet Computing*, 8(1):90–92, January-February 2004.
- [64] D. Menasce and M. Bennani. On the use of performance models to design self-managing computer systems. In *Proceedings of the 2003 Computer measurement group conference*. IEEE, 2003.
- [65] Daniel Menasce and E. Casalicchio. A framework for resource allocation in grid computing. In *Proceedings of the 12th IEEE International Symposium on Modeling, Analysis, and Simulation of Computer and Telecommunications Systems (MASCOTS'04)*, pages 259–267. IEEE, October 2004.
- [66] Sun Microsystems. Service level agreement in the data center. In <http://www.sun.com/blueprints/0402/sla.pdf>.
- [67] L. Mui, M. Mohtashemi, and A. Halberstadt. A computational model of trust and reputation. In *Proceedings of the 35th Hawaii International Conference on System Science*, 2002.
- [68] R. Muntz, K. Chandy, F. Baskett, and F. Palacios. Open, closed, and mixed networks of queues with different classes of customers. *Journal of the ACM*, April 1975.

- [69] J. Muppula, K. Trivedi, V. Mainkar, and V. Kulkarni. Numerical computation of response time distributions using stochastic reward nets. *Annals of Oper. Res.*, 48:155–184, 1994.
- [70] J. Nabrzysbi, J. Schopf, and J. Weglarz. *Grid Resource Management*. Kluwer Academic Publication, Boston, MA, 2004.
- [71] M. Naor and A. Wool. The load, capacity, and availability of quorum systems. *SIAM Journal on Computing*, 27(2):423–447, 1998.
- [72] B. Neuman, B. Tung, J. Way, and J. Trostle. Public key cryptography for initial authentication in Kerberos servers (PKINIT-02). <http://ietf.org/internet-drafts/draft-ietf-cat-\mbox{K}erberos-pk-init-02.txt>, October 1996.
- [73] D. Nowak, P. Perry, and J. Murphy. Bandwidth allocation for service level agreement aware Ethernet passive optical networks. In *IEEE Globecom 2004*, pages 1953–1957, 2002.
- [74] T. Osogami, A. Wierman, Mor Harchol-Balter, and Alan Scheller-Wolf. How many servers are best in a dual-priority FCFS system? In *CMU Technical Report: CMU-CS-03-201, November, 2003*. Carnegie Mellon University, November 2003.
- [75] A. Perrig, R. Canetti, D. Song, and D. Tygar. Efficient and secure source authentication for multicast. In *Proceedings of Network and Distributed System Security Symposium*, February 2001.
- [76] H. Perros. *Queueing Network with Blocking, Exact and Approximate Solutions*. Oxford University Press, 1994.
- [77] R. Piessens. Gaussian quadrature formulas for the numerical integration of Bromwich’s integral and the inversion of the Laplace transform. *Journal of Engineering Mathematics*, 5(1), 1971.
- [78] W. Press, B. Flannery, S. Teukolsky, and W. Vetterling. *Numerical Recipes in Fortran*. Cambridge University Press, 1997.
- [79] E. Reich. Notes on queues are in tandem. *Annals of Mathematical Statistics*, 34:338–341.
- [80] E. Reich. Waiting times when queues are in tandem. *Annals of Mathematical Statistics*, 28(3):768–773.

- [81] M. Richardson, R. Agrawal, and P. Domingos. Trust management for the semantic Web. In *Proceedings of the Second International Semantic Web Conference*, 2003.
- [82] J. Rosenberg and D. Remy. *Securing Web services with WS-Security: demystifying WS-Security, WS-Policy, SAML, XML Signature, and XML Encryption*. SAMS, Indianapolis, IN, 2004.
- [83] J. George Shanthikumar and David D. Yao. On server allocation in multiple center manufacturing systems. *Operations Research*, 36(2):333–342, November 1988.
- [84] Munindar P. Singh and Michael N. Huhns. *Service-Oriented Computing*. John Wiley & Sons, Ltd, 2003.
- [85] M. Sirbu and J. Chuang. Distributed authentication in Kerberos using public key cryptography. In *IEEE Symposium On Network and Distributed System Security (NDSS'97)*, 1997.
- [86] D. Skeen and D. D. Wright. Increasing availability in partitioned database systems. In *Proceedings of PODS*, pages 290–299, 1984.
- [87] H. Stehfest. Algorithm 386, numerical inversion of Laplace transforms. *Communications of the ACM*, 13(1), January 1970.
- [88] V. Tasic, K. Patel, and B. Pagurek. WSOL - Web service offerings language. In *Lecture Notes In Computer Science, Vol. 2512, Revised Papers from the International Workshop on Web services, E-Business, and the Semantic Web*, pages 57–67, 2002.
- [89] L. Vu, M. Hauswirth, and K. Aberer. QoS-based service selection and ranking with trust and reputation management. In *Infoscience's Technical Report*. <http://infoscience.epfl.ch/search.py?recid=52732>, 2005.
- [90] J. Walrand and P. Varaiya. Sojourn times and the overtaking condition in Jacksonian networks. *Adv. Appl. Prob.*, 12(4):1000–1018, 1980.
- [91] X. Xiao and L. M. Ni. Internet QoS: a big picture. *IEEE Network*, March-April 1999.
- [92] K. Xiong. Web services performance modeling and analysis. In *HONET and IEEE Xploer*, September 2006.
- [93] K. Xiong and H. Perros. Network optimization subject to service level agreement. *to submit*.

- [94] K. Xiong and H. Perros. Computer resource optimization for differentiated customer services. In *Proceedings of Conference on Measurement and Simulation of Computer and Telecommunication Systems (MASCOTS)*, September 2006.
- [95] K. Xiong and H. Perros. Resource optimization subject to a percentile response time sla for enterprise computing. In *Proceedings of IEEE Globecom: Quality Reliability and Performance Modeling for Emerging Network Services (Globecom-QRPM)*, November-December 2006.
- [96] K. Xiong and H. Perros. Trust-based resource allocation in web services. In *Proceedings of IEEE International Conference on Web Services (ICWS)*, September 2006.
- [97] J. Zhang and L. J. Zhang. Editorial preface: A framework to ensure trustworthy Web services. *International Journal of Web Services Research*, 2(3):1–7, July-September 2005.
- [98] Q. Zhang, T. Yu, and K. Irwin. A classification scheme for trust functions in reputation-based trust management. In *International Workshop on Trust, Security, and Reputation on the Semantic Web*. Hiroshima,, November 2004.
- [99] L. Zhu and B. Tung. RFC 4556: Public key cryptography for initial authentication in Kerberos (PKINIT). <http://www.ietf.org/rfc/rfc4556.txt>, June 2006.
- [100] C. Ziegler and G. Lausen. Spreading activation models for trust propagation. In *IEEE International Conference on e-Technology, e-commerce, and e-Service (EEE '04)*, April 2002.