

Bandwidth Provisioning for Service Overlay Networks

Zhenhai Duan[†], Zhi-Li Zhang[†] and Thomas Yiwei Hou[‡]

[†] Dept. of Computer Science & Engineering
University of Minnesota
Minneapolis, MN 55455
{duan, zhzhang}@cs.umn.edu

[‡] Fujitsu Labs of America
595 Lawrence Expressway
Sunnyvale, CA 94085
thou@fla.fujitsu.com

ABSTRACT

In this paper we study the bandwidth provisioning problem for a service overlay network which buys bandwidth from the underlying network domains to provide end-to-end value-added QoS sensitive services such as VoIP and Video-on-Demand. A key problem in the SON deployment is the problem of *bandwidth provisioning*, which is critical to the cost recovery in deploying and operating value-added services over the SON. In this paper, we mathematically formulate the bandwidth provisioning problem, taking into account various factors such as SLA, service QoS, traffic demand distributions, and bandwidth costs. Analytical models and approximate solutions are developed for long-term static bandwidth provisioning. Numerical studies are also performed to illustrate the properties of the proposed solution and demonstrate the effect of traffic demand distributions and bandwidth costs on the bandwidth provisioning of a SON.

Keywords: Service Overlay Networks, bandwidth provisioning, end-to-end quality of services

1. INTRODUCTION

Today's Internet infrastructure supports primarily *best-effort connectivity* service. Due to historical reasons, the Internet consists of a collection of network domains (i.e., autonomous systems owned by various administrative entities). Traffic from one user to another user typically traverses multiple domains; network domains enter various bilateral business relationships (e.g., provider-customer, or peering) for traffic exchange to achieve global connectivity. Due to the nature of their business relationships, each network domain is only concerned with the network performance of its own domain and responsible for providing service guarantees for its customers. As it is difficult to establish multi-lateral business relationship involving multiple domains, the deployment of end-to-end services beyond the best-effort connectivity that requires support from multiple network domains is still far from reality. Such problems have hindered the transformation of the current Internet into a truly multi-service network infrastructure with end-to-end QoS support.

We propose and advocate the notion of *service overlay network* (SON) as an effective means to address some of the issues, in particular, end-to-end QoS, plaguing the current Internet, and to facilitate the creation and deployment of *value-added Internet services* such as VoIP, Video-on-Demand, and other emerging QoS-sensitive services. The network architecture of a SON relies on well-defined business relationships between the SON, the underlying network domains and users of the SON to provide support for end-to-end QoS: the SON purchases bandwidth with certain QoS guarantees from individual network domains via *bilateral service level agreement* (SLA) to build a logical end-to-end service delivery infrastructure on top of existing data transport networks; via a service contract (e.g., a usage-based or fixed price service plan), users* directly pay a SON provider for using the value-added services provided by the SON.

Figure 1 illustrates the SON architecture. A SON is pieced together via *service gateways* which perform service-specific data forwarding and control functions. The *logical* connection between two service gateways is provided by the underlying network domain with certain bandwidth and other QoS guarantees. These guarantees are specified in a bilateral SLA between the SON and the network domain. This architecture bypasses the peering points among the network domains, and thus avoids potential performance problems associated with them. Relying on the bilateral SLAs a SON can deliver end-to-end QoS sensitive services to its users via appropriate provisioning and service-specific resource management.

Obviously the deployment of a SON is a capital-intensive investment. It is therefore imperative to consider the *cost recovery* issue for a SON. Among many costs incurred in the deployment of a SON (e.g., equipment such as service

*Users may also need to pay (i.e., a monthly fee) the access networks for their right to access the Internet.

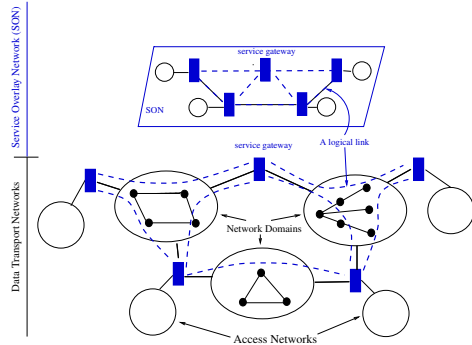


Figure 1. An illustration of a service overlay network.

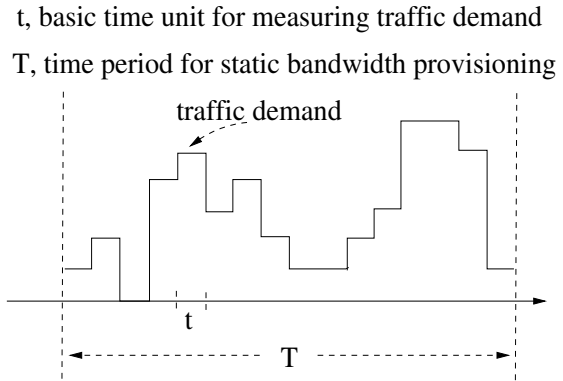


Figure 2. Traffic demands.

gateways), a dominant *recurring* cost is the cost of bandwidth that a SON must purchase from underlying network domains to support its services. A SON must provision adequate bandwidth to support its end-to-end QoS-sensitive services and meet traffic demands while minimizing the bandwidth cost so that it can generate sufficient revenue to recover its service deployment cost and stay profitable. *The bandwidth provisioning problem is therefore a critical issue in the deployment of the SON architecture.* This study is devoted to this issue. The design and implementation of the SON architecture will be left to another paper.

We develop analytical models to study the problem of SON bandwidth provisioning and investigate the impact of various factors on SON bandwidth provisioning: SLAs, service QoS, bandwidth costs and traffic demands. We consider the so-called *pipe* SLA model as an example to illustrate how the SON bandwidth provisioning problem can be formally defined. The analyses and solutions can be adapted to the so-called *hose* SLA model,¹ which due to space limitations we do not consider in this paper. In Section 2 we describe how the SON logical topology can be represented under the pipe SLA model and present the assumptions of our model. Using the pipe SLA model we present a basic static SON bandwidth provisioning solution in Section 3, and study the the more general *long-term static* SON bandwidth provisioning in Section 4. Both analytical model and approximate solution are developed. Numerical studies are also performed to illustrate the properties of the proposed solution and demonstrate the effect of traffic demand distributions and bandwidth costs on SON bandwidth provisioning.

The notion of overlay networks has been used widely in telecommunication and data networks. For example, more recently content distribution networks and application layer multicast networks have been used for multimedia streaming²; *Detour*³ and *Resilient Overlay Network (RON)*⁴ employ the overlay technique to provide better routing support. Moreover, the overlay technique has attracted a lot of attention from industries^{5,6} as a means to deliver diverse QoS-sensitive services over the Internet. The service overlay network we propose here is simply a generalization of these ideas. Perhaps what is particularly interesting is the use of SONs to address the end-to-end QoS deployment issue. The major contribution of our paper however lies in the study of the SON bandwidth provisioning problem. Our approach and formulation also differ from the traditional capacity planning in telephone networks (e.g.^{7,8}) in that we explicitly take into account various factors such as SLAs, QoS, traffic demand distributions.

2. SERVICE OVERLAY NETWORKS: ASSUMPTIONS AND BANDWIDTH PROVISIONING PROBLEMS

In this section we first describe a logical topology representation of a SON under the pipe SLA model and a simplifying assumption on service QoS. We then describe the traffic demand model and a few notations regarding service revenue and bandwidth cost that will be used later in this paper. At the end we present the bandwidth provisioning model that we will study in this paper.

The pipe SLA model is a common SLA model used in today's Internet. Under the pipe model, a SON can request bandwidth guarantees between any two service gateways across a network domain (see Fig. 1); in other words a "pipe"

with certain bandwidth guarantees is provisioned between the two service gateways across the network domain. To emphasize the relationship between service gateways and underlying network domains, we denote the *logical* (uni-directional) connection from a service gateway u to a neighboring service gateway v across a network domain D by $\langle u, v; D \rangle$, and refer to it as a *logical link* (or simply a *link*) between u and v across D . Note that between a SON and the access networks where traffic to the SON originates and terminates, the *hose* SLA model is assumed to be used where certain amount of bandwidth is reserved for traffic *entering* or *exiting* the SON. We can treat each access network A as a *fictitious* service gateway u_A . Then we can talk about “connection” between u_A and a neighboring service gateway v across A and the corresponding “logical link” $\langle u_A, v; A \rangle$.

Given a logical link $l = \langle u, v; D \rangle$, a SON provider will contract with the network domain D to provide a certain amount of bandwidth guarantee c_l between the service gateways u and v across D . The bandwidth provisioning problem of the SON is then to determine how much bandwidth to be provisioned for each link $l = \langle u, v; D \rangle$ so that: 1) the end-to-end QoS required by its services can be supported adequately; and 2) its overall revenue or net income can be maximized.

Although the QoS that a SON must support for its services can be quite diverse (e.g., bandwidth, delay or delay jitter guarantees), in almost all cases a key component in providing such guarantees is to exert some form of control on the link utilization level, i.e., to ensure the overall load on a link does not exceed certain specified condition. Consequently, for the purpose of bandwidth provisioning, we assume that it is possible to map service QoS guarantee requirements to a link utilization threshold[†]. To state this assumption formally, we assume that a link utilization threshold η_l is specified for each link l ; and to ensure service QoS, the bandwidth c_l on link l must be provisioned in such a way that the average link utilization stays below η_l .

We now describe the traffic demand model for a SON. Recall that we assume that traffic always originates from and terminates at access networks. Given a source node s and destination node d , for simplicity we assume that a fixed route r consisting of a series of links connecting s and d is used to forward traffic from s to d . Let R denote the collection of routes between the source and destination nodes. Then the traffic demands over a SON can be represented by the traffic demands over these routes: for each $r \in R$, let ρ_r denote the (average) traffic demand (also referred to as traffic load) along route r measured over some period of time t (see Fig. 2). The period t is relatively short, for example in seconds or a few minutes, compared to the time scale of static bandwidth provisioning, denoted by T , which could be in several hours or days (or longer). The period t is considered as the basic unit of time. The set $\{\rho_r : r \in R\}$ then represents the traffic demands over the SON during the time unit they are measured, and is referred to as the traffic demand matrix of the SON. Note that traffic demands are always measured in units of bandwidth.

To capture traffic demand fluctuations over time, we assume that the traffic demand ρ_r along a route r varies according to some distribution[‡]. We denote the probability density function of the traffic demand distribution of ρ_r by $d\rho_r$. Then the probability that the traffic demand ρ_r exceeds x units of bandwidth is given by $\int_x^\infty d\rho_r$. Let $\bar{\rho}_r = \int_0^\infty \rho_r d\rho_r$, i.e., $\bar{\rho}_r$ is the (long-term) average traffic demand along route r over the time period for static bandwidth provisioning. Furthermore, we assume that traffic demand distributions along different routes are *independent*. In this paper, we will study the bandwidth provisioning problem by considering two different traffic demand models. The first one takes into account the widely observed self-similar property of the Internet traffic by employing the $M/G/\infty$ input process^{12, 13}; the second is based on the measurements of real Internet traffic.

For each route r , we assume that a SON receives e_r amount of revenue for carrying one unit of traffic demand per unit of time along route r . On the other hand, for each logical link or pipe l connecting two service gateways, a SON must pay a cost of $\Phi_l(c_l)$ per unit of time for reserving c_l amount of bandwidth from the underlying network domain. We refer to Φ_l as the bandwidth cost function of link l . Without loss of generality, we assume that Φ_l is a *non-decreasing* function.

In this paper, we consider a *long-term static* bandwidth provisioning mode under the pipe model. In the *static* bandwidth provisioning mode, a SON contracts and purchases a fixed amount of bandwidth *a priori* for each link connecting

[†]This particularly will be the case if the underlying network domain employs aggregate packet scheduling mechanisms such as FIFO or priority queues. For example, it has been shown⁹⁻¹¹ that in order to provide end-to-end delay guarantees, link utilization must be controlled at a certain level. Hence from the bandwidth provisioning perspective we believe that this assumption on service QoS is not unreasonable in practice. In fact it is said that many of today’s network service providers use a similar utilization based rule (e.g., an average utilization threshold of 60% or 70%) to provision their Internet backbones.

[‡]This traffic demand distribution can be obtained, for example, through long-term observation and measurement.

the service gateways from underlying network domains. In other words, the bandwidth is provisioned for a (relatively) long period of time without changing. A key question in bandwidth provisioning for a SON is to determine the appropriate amount of bandwidth to be purchased *a priori* so that the total net income of the SON is maximized while maintaining the service QoS to meet the traffic demands.

3. BASIC STATIC BANDWIDTH PROVISIONING MODEL

In this section, we present a basic static bandwidth provisioning model and analyze its properties. This basic model will serve as the basis for the more generic bandwidth provisioning model that we will consider later. In the basic model, a SON provisions bandwidth on each link based on the long-term average traffic demand matrix $\{\bar{\rho}_r\}$, and attempts to maximize the *expected* net income. To accommodate some degree of fluctuation from the long-term average traffic demands, we introduce an *overprovisioning parameter* ϵ_l on each link l , $\epsilon_l \geq 0$. The meaning of the overprovisioning parameter ϵ_l is given as follows: we will provision c_l amount of bandwidth on link l such that as long as the overall traffic load on link l does not exceed its long-term average load by ϵ_l , the service QoS can be maintained, i.e., the link utilization is kept below the pre-specified threshold η_l . To put it formally, define $\bar{\rho}_l = \sum_{r:l \in r} \bar{\rho}_r$, where $l \in r$ denotes that link l lies on route r . Then scriptsize

$$\bar{\rho}_l(1 + \epsilon_l) = (1 + \epsilon_l) \sum_{r:l \in r} \bar{\rho}_r \leq \eta_l c_l, \forall l \in L \quad (1)$$

where L is the set of all links of the SON.

Given that c_l amount of bandwidth is provisioned on each link l , the expected net income of a SON is $\bar{W} = \sum_{r \in R} e_r \bar{\rho}_r - \sum_{l \in L} \Phi_l(c_l)$. Hence the basic bandwidth provisioning problem can be formulated as the following optimization problem:

$$\max_{c_l: l \in L} \bar{W} \quad \text{subject to (1).}$$

Since Φ_l 's are non-decreasing, it is easy to see that the optimal solution to the optimization problem is given by

$$c_l^* = (1 + \epsilon_l) \bar{\rho}_l / \eta_l, \quad \forall l \in L. \quad (2)$$

Hence under the basic bandwidth provisioning model, once we fix the overprovisioning parameters, the optimal amount of bandwidth to be provisioned for each link can be derived using (2).

Assuming that Φ_l 's are sub-additive[§], we see that a sufficient condition for a SON to have positive expected net income is to ensure that

$$e_r > \frac{\sum_{l \in r} \Phi_l(c_l^*)}{\bar{\rho}_r} = \frac{\sum_{l \in r} \Phi_l(\frac{\bar{\rho}_r(1 + \epsilon_l)}{\eta_l})}{\bar{\rho}_r}. \quad (3)$$

The relationship (3) provides a useful guideline for a SON to determine how it should set its price structure for charging the users to recover its cost of bandwidth provisioning. It has a simple interpretation: we can regard $\frac{\Phi_l(\bar{\rho}_r(1 + \epsilon_l)/\eta_l)}{\bar{\rho}_r}$ as the average cost of carrying one unit of traffic demand per unit of time along route r on link l . Then the right-hand side of (3) is the total cost of carrying one unit of traffic demand per unit of time along route r . To recover its cost, a SON must then charge the users more than this amount. If Φ_l 's are strictly concave (i.e., non-linear), in other words, the per-unit bandwidth cost decreases as the amount of reserved bandwidth increases, the economy of scale will benefit the SON: the higher the average long-term traffic demands, the lower the average cost of providing services, yielding higher net income. In the case Φ_l 's are linear, i.e., $\Phi_l(c_l) = \phi_l c_l$, then (3) becomes $e_r > \sum_{l \in r} \phi_l (1 + \epsilon_l) / \eta_l$ which is independent of traffic demands.

4. STATIC BANDWIDTH PROVISIONING WITH PENALTY

In the basic static bandwidth provisioning model we assume that the overprovisioning parameters are given. We now consider how to obtain the *optimal* overprovisioning parameters under given traffic demand distributions. We study this problem by taking into account the consequence of potential QoS violations when actual traffic demands exceed the target link utilization. For this purpose, we assume that *a SON may suffer a penalty when the target utilization on a link is*

[§] A function $\Phi_l(c)$ is sub-additive if $\Phi_l(c_1 + c_2) \leq \Phi_l(c_1) + \Phi_l(c_2)$.

exceeded, and therefore service QoS may potentially be violated. For example, it is quite likely that the service contract between a SON and its user is such that when the service QoS is poor (e.g., due to network congestion), a lower rate is charged, or the user may demand a refund. In the case that some form of admission control is used by a SON to guide against possible QoS violations, the penalty can be used to reflect the lost revenue due to declined user service requests. We refer to this model as the *static bandwidth provisioning with penalty model*, or in short, *static-penalty model*.

For each route r , let π_r denote the average penalty suffered by per unit of traffic demand per unit of time along route r when the service QoS along the route is potentially violated. Given a traffic demand matrix $\{\rho_r\}$, let $B_r(\{\rho_r\})$ denote the probability that the service QoS along route r is potentially violated, more specifically, *the target utilization on one of its links is exceeded*. Then the total net income of a SON for servicing a given traffic demand matrix $\{\rho_r\}$ can be expressed as follows:

$$W(\{\rho_r\}) = \sum_{r \in \mathcal{R}} e_r \rho_r - \sum_{l \in \mathcal{L}} \Phi_l(c_l) - \sum_{r \in \mathcal{R}} \pi_r \rho_r B_r(\{\rho_r\}), \quad (4)$$

where in the above we use $W(\{\rho_r\})$ to emphasize the dependence of the total net income on the traffic demand matrix $\{\rho_r\}$. When there is no confusion, we may drop $\{\rho_r\}$ from the notation.

Let $d\{\rho_r\}$ denote the joint probability density function of a traffic demand matrix $\{\rho_r\}$, where recall that $d\rho_r$ is the probability density function of a traffic demand ρ_r along route r . Then the expected net income of a SON under the traffic demand distributions $\{d\rho_r\}$ is given by

$$E(W) = \int \cdot \int_{\{\rho_r\}} W(\{\rho_r\}) d\{\rho_r\}, \quad (5)$$

where $\int \cdot \int_{\{\rho_r\}}$ denotes multiple integration under the joint traffic demand distribution $\{d\rho_r\}$.

Now we can state the problem of static bandwidth provisioning with penalty as the following optimization problem: finding the optimal overprovisioning parameters $\{\epsilon_l\}$ to maximize the expected net income, i.e., scripsize

$$\max_{\{\epsilon_l\}} E(W) \text{ subject to (1)}. \quad (6)$$

Unfortunately, the exact solution to this optimization problem is in general difficult to obtain. It depends on both the particular forms of the traffic demand distributions $\{d\rho_r\}$ and the service QoS violation probabilities B_r . To circumvent this difficulty, in the following, we shall derive an approximate solution (a lower bound) based on the so-called *link independence assumption*: the link *overload* events (i.e., exceeding the target utilization threshold) occur on different links *independently*. Clearly this assumption does not hold in reality, but it enables us to express B_r in terms of $B_l(\rho_l, c_l)$, the probability that the target utilization level η_l on link l is exceeded, where $\rho_l = \sum_{r: l \in r} \rho_r$. (Again, we may drop the variables ρ_l and c_l in $B_l(\rho_l, c_l)$ if there is no confusion.) Such link independence assumption has been used extensively in teletraffic analysis and capacity planning in the telephone networks (see e.g.,⁸). Under the link independence assumption, the service QoS violation probability B_r , i.e., at least one of the links on route r is overloaded, is given by

$$B_r = 1 - \prod_{l \in r} (1 - B_l). \quad (7)$$

Before we present the approximate optimal solution, we need to introduce one more set of notations. Define a small real number $\delta > 0$. For each route r , let $\hat{\rho}_r > \bar{\rho}_r$ be such that

$$\int_{\hat{\rho}_r}^{\infty} \rho_r d\rho_r \leq \delta. \quad (8)$$

Since $\int_{\hat{\rho}_r}^{\infty} \rho_r d\rho_r \geq \hat{\rho}_r \int_{\hat{\rho}_r}^{\infty} d\rho_r = \hat{\rho}_r Pr\{\rho_r \geq \hat{\rho}_r\}$, we have $Pr\{\rho_r \geq \hat{\rho}_r\} \leq \delta / \hat{\rho}_r$. In other words, (8) basically says that $\hat{\rho}_r$ is such that the probability the traffic demand along route r exceeds $\hat{\rho}_r$ is very small, and thus negligible.

With these notations in place, we now present a lower bound on $E(W)$ as follows (see Appendix A for the detailed derivation).

$$E(W) \geq \sum_{r \in \mathcal{R}} e_r \bar{\rho}_r - \sum_{l \in \mathcal{L}} \Phi_l(c_l) - \sum_{r \in \mathcal{R}} \pi_r \bar{\rho}_r B_r(\{\hat{\rho}_r\}) - \sum_{r \in \mathcal{R}} \pi_r \delta (1 + \sum_{r' \neq r} \frac{\bar{\rho}_r}{\hat{\rho}_{r'}}).$$

Denote the right-hand side of the above equation by V , then $E(W) \geq V$. Comparing the lower bound V with the expected net income $\bar{W} = \sum_{r \in \mathcal{R}} e_r \bar{\rho}_r - \sum_{l \in \mathcal{L}} \Phi_l(c_l)$ without taking penalty into account, we see that ignoring the

extremal traffic demands (i.e., when $\rho_r \geq \hat{\rho}_r$), we pay at most a penalty of $\pi_r B_r(\{\hat{\rho}_r\})$ per unit of traffic demand on route r for potential service QoS violations. For given $\delta > 0$, the penalty incurred due to extremal traffic demands is upper bounded by $\sum_{r \in \mathcal{R}} \pi_r \delta (1 + \sum_{r' \neq r} \frac{\bar{\rho}_r}{\hat{\rho}_{r'}})$. Note also that $B_r(\{\hat{\rho}_r\})$ is the probability of service QoS violation along route r when the long-term average traffic demands are assumed to be $\hat{\rho}_r$. Thus in using V as an approximation to $E(W)$, we are being conservative by over-estimating the probability of potential QoS violations.

From $E(W) \geq V$, we have $\max_{\{\epsilon_r\}} E(W) \geq \max_{\{\epsilon_r\}} V$. Therefore we can obtain the *best* overprovisioning parameters that maximize V instead of the expected net income $E(W)$ as an approximate solution to the original optimization problem (6). Using the solution to the basic bandwidth provisioning problem (2), we assume $c_l = (1 + \epsilon_l) \bar{\rho}_l / \eta_l$ for a given set of $\{\epsilon_l\}$, i.e., the target utilization constraints (1) hold with equality. Under this assumption, let $\{\epsilon_l^*\}$ be the solution to the optimization problem $\max_{\{\epsilon_r\}} V$, and refer to them as the *approximate optimal overprovisioning parameters*. In the following we demonstrate how $\{\epsilon_l^*\}$ can be derived.

Using (7) we can re-write V as follows:

$$V = \sum_{r \in \mathcal{R}} (e_r - \pi_r) \bar{\rho}_r - \sum_{l \in \mathcal{L}} \Phi_l(c_l) + \sum_{r \in \mathcal{R}} \pi_r \bar{\rho}_r \prod_{l \in r} (1 - B_l(\hat{\rho}_l, c_l)) - \sum_{r \in \mathcal{R}} \pi_r \delta (1 + \sum_{r' \neq r} \frac{\bar{\rho}_r}{\hat{\rho}_{r'}}), \quad (9)$$

where $\hat{\rho}_l = \sum_{r: l \in r} \hat{\rho}_r$.

Assume B_l is a continuous and everywhere differentiable function of c_l . (See the next subsection for a discrete case.) For each link l , define

$$\hat{s}_l = \sum_{r: l \in r} \pi_r \bar{\rho}_r \prod_{k \in r, k \neq l} [1 - B_k(\hat{\rho}_k, c_k)] \zeta_l, \quad (10)$$

where $\zeta_l = -\frac{d}{dc_l} B_l(\hat{\rho}_l, c_l)$.

Through some simple algebraic manipulation, it is not too hard to show that

$$\frac{\partial V}{\partial \epsilon_l} = \frac{\partial V}{\partial c_l} \frac{\partial c_l}{\partial \epsilon_l} = \left(-\frac{\partial \Phi_l(c_l)}{\partial c_l} + \hat{s}_l \right) \frac{\bar{\rho}_l}{\eta_l}. \quad (11)$$

Suppose that $\{\epsilon_l^*\}$ are strictly positive, then a necessary condition for them to be an optimal solution is that the gradient ∇V (with respect to $\{\epsilon_l\}$) must vanish at ϵ_l^* 's. Thus from (11) we must have

$$\frac{\partial \Phi_l(c_l)}{\partial c_l} = \hat{s}_l, \quad \forall l \in \mathcal{L}. \quad (12)$$

Intuitively, \hat{s}_l measures the sensitivity of potential penalty reduction to bandwidth increase on link l , whereas $\frac{\partial \Phi_l(c_l)}{\partial c_l}$ measures the sensitivity of bandwidth cost to bandwidth increase on link l . Hence the “optimal” (or rather, the approximate optimal) overprovisioning parameter ϵ_l^* should be chosen such that the two values coincide. In the following discussion, we will loosely refer to \hat{s}_l as the “per-unit bandwidth gain in potential penalty reduction” and $\frac{\partial \Phi_l(c_l)}{\partial c_l}$ as the “increase in per-unit bandwidth cost.”

In the above derivation of the approximate optimal solution to the static bandwidth provisioning problem, we have simply assumed the existence of B_l but not its form. Its particular form depends on the distribution of (average) traffic demands on link l . In the following subsections, we consider two different traffic demand models—an $M/G/\infty$ traffic demand model and a traffic demand model based on real Internet traffic measurements—to demonstrate the approximate optimal solution to the static bandwidth provisioning problem.

4.1. $M/G/\infty$ traffic demand model

Since the pioneering work by Leland, Taqqu, Willinger and Wilson,¹⁴ the self-similar (or long-range dependent) property has been observed in the Ethernet Local Area Network,¹⁴ Wide Area Network,¹³ and World Wide Web traffic.¹⁵ The observed self-similar property of the Internet traffic has important implications on the dimensioning and provisioning of IP networks. In this section, we consider a traffic demand model, $M/G/\infty$, that captures the (asymptotically) self-similar property of the Internet traffic.^{12, 13}

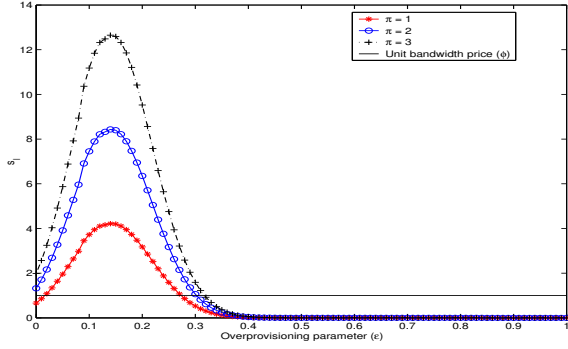


Figure 3. Relationship between \hat{s}_l , ϵ , & ϕ_l .

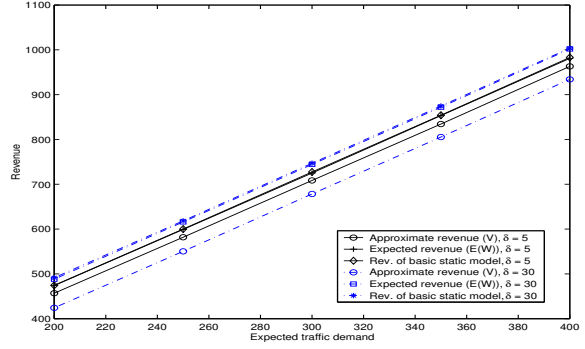


Figure 4. Comparison of V and $E[W]$.

Consider an $M/G/\infty$ input process, where the service time has a heavy-tailed distribution. We assume that the distribution of the service time has a finite mean. Let X_t denote the number of customers in the system at time t , for $t = 0, 1, 2, \dots$. Then the count process $\{X_t\}_{t=0,1,2,\dots}$ is asymptotically self-similar. Let ρ denote the customer arrival rate of the $M/G/\infty$ input process and μ the mean service time, then X_t has a Poisson marginal distribution with mean $\rho\mu$.¹⁶

Now we are ready to present the $M/G/\infty$ traffic demand model. Consider an arbitrary route r . In the $M/G/\infty$ traffic demand model, the (average) traffic demand (i.e., the average traffic arrival rate in each unit of time) on the route is governed by the count process $\{X_t\}_{t=0,1,2,\dots}$ of an $M/G/\infty$ input process. For example, let $\rho_{r,i}$ denote the average traffic demand in the i th unit of time, then we have $\rho_{r,i} = X_i$. Let $\bar{\rho}_r$ denote the long-term average traffic demand on the route. It is easy to see that $\bar{\rho}_r = \rho\mu$, where ρ and μ are the customer arrival rate and the mean service time, respectively, of the $M/G/\infty$ input process. As traffic demands along all the routes are assumed to be independent, the average overall traffic load on a link l is $\bar{\rho}_l = \sum_{r:l \in r} \bar{\rho}_r$.

Given the average overall load $\bar{\rho}_l$ and the link capacity c_l , it can be shown that the probability that the total load on link l exceeds $\bar{c}_l = \eta_l c_l$ during any given unit of time is given by $B_l(\bar{\rho}_l, c_l) = (\sum_{i=(\bar{c}_l+1)}^{\infty} \frac{\bar{\rho}_l^i}{i!}) e^{-\bar{\rho}_l}$. We extend the definition of $B_l(\bar{\rho}_l, c_l)$ to the non-integer values of c_l by linear interpolation. Moreover, at the integer values of c_l we define the derivative of $B_l(\bar{\rho}_l, c_l)$ with respect to c_l to be the left derivative. Then $\frac{d}{dc_l} B_l(\bar{\rho}_l, c_l) = B_l(\bar{\rho}_l, c_l) - B_l(\bar{\rho}_l, c_l - 1)$. Therefore, $\zeta_l = -\frac{d}{dc_l} B_l(\hat{\rho}_l, c_l) = \eta_l \{B_l(\hat{\rho}_l, (\eta_l c_l - 1)) - B_l(\hat{\rho}_l, \eta_l c_l)\} = \eta_l \frac{\hat{\rho}_l^{\lceil \eta_l c_l \rceil}}{\lceil \eta_l c_l \rceil!} e^{-\hat{\rho}_l}$. By this definition of B_l , we can obtain the (approximate) optimal overprovisioning parameters ϵ_l^* 's by solving (12).

We now discuss the effect of the *shapes* of \hat{s}_l ' and Φ_l on (approximate) optimal overprovisioning parameters ϵ_l^* 's as well as their implication in static bandwidth provisioning. Note first that the shape of \hat{s}_l is determined by ζ_l , which has a shape of (skewed) bell-shape with a center approximately at $\hat{\rho}_l$ (it is essentially a Poisson probability density function). Hence \hat{s}_l is a concave function of $\epsilon_l \geq 0$. In particular, there exists $\hat{\epsilon}_l$ such that \hat{s}_l is an increasing function in the range $[0, \hat{\epsilon}_l]$ and a decreasing function in the range $[\hat{\epsilon}_l, \infty)$ (see Fig. 3). Intuitively, this means that as ϵ_l moves from 0 towards $\hat{\epsilon}_l$, there is an increasing benefit in bandwidth overprovisioning in terms of *reducing potential QoS violation penalty*. However, as ϵ_l moves beyond $\hat{\epsilon}_l$, there is a *diminished return* in overprovisioning in terms of reducing potential QoS violation penalty.

Suppose that Φ_l is a linear function, i.e., $\Phi_l(c_l) = \phi_l c_l$. Then $\frac{\partial \Phi_l(c_l)}{\partial c_l} = \phi_l$. Hence (12) becomes $\phi_l = \hat{s}_l$. Suppose $\phi_l = \hat{s}_l$ holds for some $\epsilon_l \geq 0$. Because of the shape of \hat{s}_l , there potentially exist two solutions $\epsilon_{l,1}$ and $\epsilon_{l,2}$, $0 \leq \epsilon_{l,1} \leq \hat{\epsilon}_l \leq \epsilon_{l,2}$ such that $\phi_l = \hat{s}_l$. In particular, as \hat{s}_l is a decreasing function in the range $[\hat{\epsilon}_l, \infty)$, $\epsilon_{l,2}$ always exists. As $\frac{\partial V}{\partial c_l}$ is positive in the range $(\epsilon_{l,1}, \epsilon_{l,2})$, and is negative in the ranges $[0, \epsilon_{l,1})$ and $(\epsilon_{l,2}, \infty)$, we see that with respect to link l , V is maximized at either $\epsilon_l^* = \epsilon_{l,2}$ or at $\epsilon_l^* = 0$ (whereas it is minimized at $\epsilon_{l,1}$). Intuitively, when only a small amount of bandwidth is overprovisioned on link l , the per-unit bandwidth gain in potential penalty reduction is too small to offset the per-unit bandwidth cost, hence V decreases. However, as we increases the amount of bandwidth overprovisioned, the per-unit bandwidth gain in potential penalty reduction becomes sufficiently large and offsets the per-unit bandwidth cost, hence V increases until it reaches a maximum. Due to the diminished return in the per-unit bandwidth gain in potential

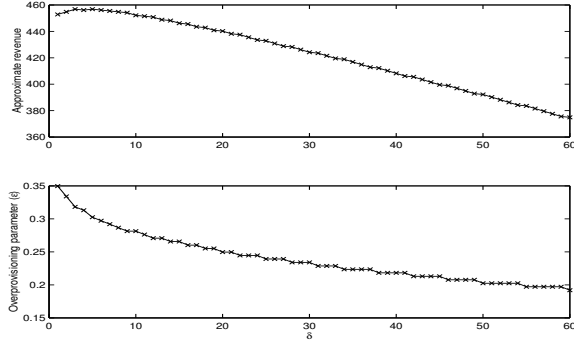


Figure 5. Impact of δ on V and ϵ^* .

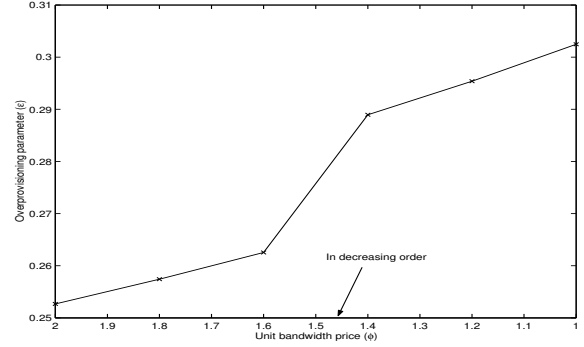


Figure 6. Impact of unit bandwidth price on ϵ^* .

penalty reduction, V decreases again when too much bandwidth is overprovisioned on link l . In the special case that ϕ_l is such that $\phi_l > \hat{s}_l$ for all $\epsilon_l \geq 0$, then as $\frac{\partial V}{\partial c_l} < 0$, V attains its maximum at $\epsilon_l^* = 0$ with respect to link l . Intuitively it says that when the per-unit bandwidth cost on link l is higher than the per-unit bandwidth gain in potential penalty reduction, there is no benefit in overprovisioning any bandwidth on link l to guide against any potential QoS violation penalty. These observations can be extended to other bandwidth cost functions such as concave or convex cost functions. In general we see that the trade-off between the bandwidth cost and overprovisioning bandwidth to guide against service QoS violations is critical to the problem of SON bandwidth provisioning. It is also clear from the above discussion that as the per-unit bandwidth cost decreases, there is more benefit in overprovisioning. Lastly, we comment that from (12) and (10) and the above observations, we can compute the approximate optimal overprovisioning parameters ϵ_l^* 's using *fixed point approximation*.

4.1.1. Numerical examples

We conduct numerical studies to illustrate the properties of the analytic results we obtained and demonstrate the effects of various parameters on static bandwidth provisioning. For this purpose, we consider a simple setting: a single route over a single link. Numerical studies in more complex settings will be performed in a later section.

Unless otherwise stated, the following parameters will be used in the numerical studies: the long-term average traffic demand on the route is 200 (measured in unit of bandwidth per unit of time), i.e., $\bar{\rho}_r (= \bar{\rho}_l) = 200$, and $e_r = 4$, $\phi_l = 1$, $\pi_r = 2$. We set $\delta = 5$ and the target utilization threshold $\eta_l = 0.8$.

Fig. 3 shows \hat{s}_l as a function of ϵ_l with three different values of π_r : $\pi_r = 1, 2, 3$. In the figure we also include a line corresponding to $\phi_l = 1$ to illustrate how ϵ_l^* can be obtained as the solution to $\hat{s}_l = \phi_l$. Recall from Section 4, $\epsilon_l^* = \epsilon_{l,2}$ (the right intersecting point). From Fig. 3 we see that as the penalty π_r increases, ϵ_l^* also increases. Hence for a higher penalty it is necessary to overprovision more bandwidth to guide against potential QoS violations. Likewise, as we increase the per-unit bandwidth cost ϕ_l (i.e., moving up the line of ϕ_l), ϵ_l^* decreases. In other words, as the bandwidth cost increases, it is beneficial to reduce overprovisioned bandwidth so as to maximize the net income.

In Fig. 4 we compare the lower bound V with the actual expected net income $E(W)$ for two given values of δ (5 and 30). For comparison, we also include the expected net income \bar{W} under the basic static model, where the overprovisioning parameter ϵ_l^* is obtained from the static-penalty model. From the figure we see that for both values of δ , the lower bound V provides a reasonable approximation to $E(W)$. Note also that the difference between the actual expected net income $E[W]$ under the static-penalty model and the expected net income \bar{W} under the basic static is almost invisible. This is likely due to the fact that the additional revenue generated when the traffic demand exceeds its long-term average (the first term in $E[W]$) and the potential penalty incurred due to service QoS violations (the third term in $E[W]$) cancel each other out on average. From Fig. 4 it is clear that the lower bound depends on the choice of δ . The smaller the δ is, the closer the approximate revenue V is to the expected revenue $E[W]$. To further explore the relation between δ and V , in Fig. 5 we plot V as a function of δ (upper plot). In the figure, we also include the overprovisioning parameter ϵ_l^* as a function of δ (lower plot). We see that V is a concave function of δ , and thus there is a unique δ that maximizes V . On the other hand, ϵ_l^* is a non-increasing function of δ .

Table 1. Provisioning for the Auckland traffic demands.

	Mean	STD	C.O.V.	ϵ_l^*	V
Day-time	2096	442	0.21	0.67	3446
Night-time	609	240	0.39	0	1672

the histograms of the traffic demands for the day-time (left-hand side) and night-time (right-hand side) separately, where the bin sizes for the day-time traffic demands and the night-time traffic demands are 100 Kb/s and 50 Kb/s , respectively. From the plots we see that the day-time traffic demands are relatively symmetrically centered at its mean arrival rate, while the night-time traffic demands are more skewed. In the following studies, we approximate the day-time traffic demands by a *Normal* distribution, while the night-time traffic demands by a *Lognormal* distribution to retain the different traffic characteristics during the day-time and night-time. Table 1 presents the mean traffic demands and the standard deviations (*STD*) of the day-time and night-time traffic demands, where the basic unit of bandwidth (traffic demand) is 1 Kb/s .

In the following, we conduct numerical studies to illustrate the static bandwidth provisioning using the Auckland data trace. In all these studies, we again consider the simple setting: a single route over a single link. The per-unit bandwidth per-unit time earning $e_r = 4$, and $\phi_l = 1$, $\pi_r = 2$. We set the target utilization threshold $\eta_l = 0.8$.

Similar to the numerical example for the $M/G/\infty$ traffic demand model, in Fig. 9, we show \hat{s}_l as a function of ϵ_l with three different values of π_r : $\pi_r = 1, 2, 3$, for the day-time traffic demands. The value of δ used is 140. In the figure we also include a line corresponding to $\phi_l = 1$ to illustrate how ϵ_l^* can be obtained as the solution to $\hat{s}_l = \phi_l$ (see (12)). Following a similar argument as that in Section 4.1, there potentially exist two solutions $\epsilon_{l,1}$ and $\epsilon_{l,2}$, $0 \leq \epsilon_{l,1} \leq \epsilon_{l,2}$ such that $\phi_l = \hat{s}_l$. Moreover, with respect to link l , V is maximized at either $\epsilon_l^* = \epsilon_{l,2}$ or at $\epsilon_l^* = 0$. From Fig. 9 we can draw similar conclusions as that in the $M/G/\infty$ traffic demand model. In particular, we see that as the penalty π_r increases, ϵ_l^* also increases. Hence for a higher penalty it is necessary to overprovision more bandwidth to guide against potential QoS violations. Likewise, as we increase the per-unit bandwidth cost ϕ_l (i.e., moving up the line of ϕ_l), ϵ_l^* decreases. In other words, as the bandwidth cost increases, it is beneficial to reduce overprovisioned bandwidth so as to maximize the net income. However, compared with the result in Fig. 3, we see that we obtain larger overprovisioning parameters here. This is caused by the high traffic fluctuation in the Auckland data trace. Table 1 gives the *coefficient of variance* for the day-time traffic demands (and the night-time traffic demands) in the column marked as *C.O.V.*. This value (0.21) is much higher than that in Fig. 3, which is 0.07.

To compare the different provisioning behaviors during the day-time and night-time, we present the overprovisioning parameters for both the day-time and night-time traffic demands in Table 1. To obtain these results, we have searched for the best δ 's that yield the maximal V 's, respectively. In the table we also include the approximate revenue V 's (per-unit time) for the day-time and night-time traffic demands. From the table we see that for the day-time traffic demands the overprovisioning parameter $\epsilon_l^* = 0.67$, while for the night-time traffic demands $\epsilon_l^* = 0$. The reason is as follows. Even though the average traffic demands during the night-time are much lower than that during the day-time, we observe a much higher traffic demand fluctuation during the night-time than that during the day-time (see Table 1 for their corresponding coefficients of variance). In this case, it is too expensive to accommodate this high traffic demand variance during the night-time ($\epsilon_{l,2}$ is significantly large). Therefore, it is not beneficial to overprovision any bandwidth and the (approximate) maximal expected revenue is achieved at $\epsilon_l^* = 0$.

4.3. Performance evaluation

We now use two SON topologies—a *tree* (Fig. 4.2(a)) and a *mesh-tree* (Fig. 4.2(b))—to illustrate the effect of traffic load distribution among various routes of a SON on static bandwidth provisioning. In the following $a \rightarrow b$ denotes a route from service gateway a to service gateway b . The path with *minimum* “hop-count” (i.e., service gateways) is used as the route between two service gateways. In case there are two such paths, only one is chosen. In the numerical studies below, we use the $M/G/\infty$ traffic demand model. We set $e_r = 10$, $\pi_r = 2$ for all routes, and $\phi_l = 1$ for all links. The value of δ is chosen in such a way that $\delta_r = \frac{1}{40}\rho_r$.

In the tree topology, four routes are used: $R1 = S3 \rightarrow C1$, $R2 = S1 \rightarrow C1$, $R3 = S4 \rightarrow C2$, and $R4 = S2 \rightarrow C2$. To investigate the effects of different traffic loads on bandwidth provisioning, we consider two types of traffic

Table 2. Tree Topology.

Link ID		1	4	5	7	8	9
Balanced	ρ_l	400	200	800	200	200	400
	ϵ_l^*	0.26	0.3	0.23	0.3	0.3	0.26
	c_l	630	325	1230	325	325	630
Unbalanced	ρ_l	400	250	800	100	150	250
	ϵ_l^*	0.26	0.27	0.23	0.41	0.34	0.33
	c_l	630	397	1230	176	251	416

Table 3. Mesh-Tree Topology.

Link ID		2	6	11	18	19	21
Balanced	ρ_l	1200	800	400	400	400	200
	ϵ_l^*	0.22	0.23	0.26	0.26	0.26	0.3
	c_l	1830	1230	630	630	630	325
Unbalanced	ρ_l	1350	1100	500	400	400	100
	ϵ_l^*	0.21	0.2	0.24	0.26	0.26	0.41
	c_l	2042	1650	775	630	630	176

load distribution among the routes: the *balanced* load where the expected traffic demand for all routes is 200, and the *unbalanced* load where the expected traffic demands on routes $R1, R2, R3, R4$ are 300, 100, 250, and 150, respectively. Table 2 presents the resulting overprovisioning parameter ϵ_l^* and provisioned bandwidth c_l for six representative links: link 1, 4, 5, 7, 8, and 9. The corresponding average traffic loads $\bar{\rho}_l$'s on these links are also given in the table. From the results we see that under the balanced load, links with a higher average traffic load have a smaller overprovisioning parameter. This is due to statistical multiplexing gains for carrying a higher load on a link. In the *unbalanced* case, similar results are observed. Note that even though links 4 and 9 have the same traffic demand load, they are provisioned differently. This is caused by the fact that there are two routes traversing link 9, while there is only one on link 4.

We now consider the mesh-tree topology. In this case there are 10 routes: $R1 = S1 \rightarrow C1, R2 = S2 \rightarrow C2, R3 = S3 \rightarrow C1$ (1), $R4 = S4 \rightarrow C2$ (1), $R5 = S1 \rightarrow C3$ (3), $R6 = S2 \rightarrow C4$ (3), $R7 = S3 \rightarrow C3, R8 = S4 \rightarrow C4, R9 = S5 \rightarrow C5, R10 = S5 \rightarrow C6$. The number in the parentheses following a route shows a link that the route traverses in case there are multiple paths between the source and destination with the same path length. Again for the balanced load case, all the routes have an average traffic demand of 200; while for the unbalanced load case, the average demands for routes $R1$ to $R10$ are 300, 250, 100, 150, 300, 250, 100, 150, 300, and 100 respectively. Table 3 shows the results for six representative links: link 2, 6, 11, 18, 19, and 21. From the table we can see that similar observations also hold for the mesh-tree topology.

5. CONCLUSIONS AND FUTURE WORK

In this paper, we studied the long-term static bandwidth provisioning problem for service overlay networks (SONs). Our formulation of the SON bandwidth provisioning problem took into account various factors such as service QoS, traffic demand distributions, and bandwidth costs. The approximate optimal solution we developed to the static bandwidth provisioning problem is generic in the sense that it applies to different marginal distributions of the traffic demands on the routes in a network, which makes the solution very attractive facing different traffic arrival behaviors. We also investigated the effects of various parameters like bandwidth costs on the revenue that a SON can obtain, which provides useful guidelines on how a SON should be provisioned to stay profitable.

The static bandwidth provisioning model is simple in terms of network resource management but may result in inefficient network resource usage if the traffic demands are highly variable. In this kind of environments, a dynamic bandwidth provisioning model may be called for, which is currently under our investigation. We are also interested in exploring the functionalities of service gateways in support of service-aware (multi-path) routing, which may have great impact on how a SON should be dimensioned and provisioned.

APPENDIX A. A LOWER BOUND ON $E(W)$ OF THE STATIC BANDWIDTH PROVISIONING WITH PENALTY

From (4) and (5), it is easy to see that

$$E[W] = \sum_{r \in R} e_r \bar{\rho}_r - \sum_{l \in L} \Phi_l(c_l) - \sum_{r' \in R} \int \int_{\{\rho_r\}} \pi_{r', \rho_r'} B_{r'}(\{\rho_r\}) d\{\rho_r\}. \quad (13)$$

Moreover,

$$\begin{aligned} \sum_{r' \in \mathcal{R}} \int \cdot \int_{\{\rho_r\}} \pi_{r', \rho_r'} B_{r'}(\{\rho_r\}) d\{\rho_r\} &\leq \sum_{r' \in \mathcal{R}} \int \cdot \int_{\{0\}}^{\{\hat{\rho}_r\}} \pi_{r', \rho_r'} B_{r'}(\{\rho_r\}) d\{\rho_r\} + \sum_{r' \in \mathcal{R}} \sum_{r'' \in \mathcal{R}, r'' \neq r'} \int_{\hat{\rho}_{r''}}^{\infty} \left(\int \cdot \int_{\{0\}}^{\{\infty\}} \right) \pi_{r', \rho_r'} B_{r'}(\{\rho_r\}) d\{\rho_r\} \\ &+ \sum_{r' \in \mathcal{R}} \int_{\hat{\rho}_{r'}}^{\infty} \left(\int \cdot \int_{\{0\}}^{\{\infty\}} \right) \pi_{r', \rho_r'} B_{r'}(\{\rho_r\}) d\{\rho_r\}. \end{aligned} \quad (14)$$

As $B_{r'}(\{\rho_r\}) \leq B_{r'}(\{\hat{\rho}_r\})$ when $\rho_r \leq \hat{\rho}_r, \forall r$,

$$\int \cdot \int_{\{0\}}^{\{\hat{\rho}_r\}} \pi_{r', \rho_r'} B_{r'}(\{\rho_r\}) d\{\rho_r\} \leq \int \cdot \int_{\{0\}}^{\{\hat{\rho}_r\}} \pi_{r', \rho_r'} B_{r'}(\{\hat{\rho}_r\}) d\{\rho_r\} \leq \pi_{r'} B_{r'}(\{\hat{\rho}_r\}) \bar{\rho}_{r'}. \quad (15)$$

Notice $B_{r'}(\{\rho_r\}) \leq 1$ and the definition of δ , we have (note $r'' \neq r'$)

$$\int_{\hat{\rho}_{r''}}^{\infty} \left(\int \cdot \int_{\{0\}}^{\{\infty\}} \right) \pi_{r', \rho_r'} B_{r'}(\{\rho_r\}) d\{\rho_r\} \leq \int_{\hat{\rho}_{r''}}^{\infty} \left(\int_0^{\infty} \pi_{r', \rho_r'} d\rho_{r'} \right) d\rho_{r''} \leq \pi_{r'} \delta \frac{\bar{\rho}_{r'}}{\hat{\rho}_{r''}}. \quad (16)$$

Similarly,

$$\int_{\hat{\rho}_{r'}}^{\infty} \left(\int \cdot \int_{\{0\}}^{\{\infty\}} \right) \pi_{r', \rho_r'} B_{r'}(\{\rho_r\}) d\{\rho_r\} \leq \int_{\hat{\rho}_{r'}}^{\infty} \pi_{r', \rho_r'} d\rho_{r'} \leq \pi_{r'} \delta. \quad (17)$$

Substitute (15), (16), and (17) into (14), and then recursively into (13), we have

$$E(W) \geq \sum_{r \in \mathcal{R}} e_r \bar{\rho}_r - \sum_{l \in \mathcal{L}} \Phi(c_l) - \sum_{r \in \mathcal{R}} \pi_r \bar{\rho}_r B_r(\{\hat{\rho}_r\}) - \sum_{r \in \mathcal{R}} \pi_r \delta \left(1 + \sum_{r' \in \mathcal{R}, r' \neq r} \frac{\bar{\rho}_{r'}}{\hat{\rho}_{r'}} \right). \quad (18)$$

REFERENCES

1. N. Duffield, P. Goyal, A. Greenberg, P. Mishra, K. Ramakrishnan, and J. Merwe, "A flexible model for resource management in virtual private networks," in *Proc. ACM SIGCOMM*, (Cambridge, Massachusetts), Sept. 1999.
2. Y. Chawathe, S. Fink, S. McCanne, and E. Brewer, "A proxy architecture for reliable multicast in heterogeneous environments," in *Proceedings of ACM Multimedia*, (Bristol, U.K.), Sept. 1998.
3. S. Savage, T. Anderson, and et al., "Detour: a case for informed internet routing and transport," *IEEE Micro* **19**, pp. 50–59, Jan. 1999.
4. D. G. Andersen, H. Balakrishnan, M. Kaashoek, and R. Morris, "Resilient overlay networks," in *Proc. 18th ACM SOSP*, (Banff, Canada), Oct. 2001.
5. V. Communications, "http://www.virtel.com."
6. I. N. S. Corporation, "http://www.internap.com."
7. A. Girard, *Routing and Dimensioning in Circuit-Switched Networks*, Addison-Wesley, 1990.
8. F. P. Kelly, "Routing in circuit-switched networks: Optimization, shadow prices and decentralization," *Advances in Applied Probability* **20**, pp. 112–144, 1988.
9. A. Charny and J.-Y. L. Boudec, "Delay bounds in a network with aggregate scheduling," in *Proceedings of QoFIS*, (Berlin, Germany), Oct. 2000.
10. S. Jamin, P. Danzig, S. Shenker, and L. Zhang, "A measurement-based call admission control for integrated services packet networks," in *Proc. ACM SIGCOMM*, pp. 2–13, (Cambridge, MA), Aug. 1995.
11. Z.-L. Zhang, Z. Duan, and Y. T. Hou, "Fundamental trade-offs in aggregate packet scheduling," in *Proceedings of IEEE International Conference on Network Protocols (ICNP)*, (Riverside, CA), Nov. 2001.
12. M. Parulekar and A. M. Markowski, "M/G/∞ input processes: A versatile class of models for network traffic," in *Proc. IEEE INFOCOM*, pp. 419–426, (Kobe, Japan), Apr. 1997.
13. V. Paxson and S. Floyd, "Wide area traffic: The failure of poisson modeling," in *Proc. ACM SIGCOMM*, pp. 257–268, Aug. 1994.
14. W. E. Leland, M. S. Taqqu, W. Willinger, and D. Wilson, "On the self-similar nature of ethernet traffic," *IEEE/ACM Transactions on Networking* **2**(1), 1994.
15. M. Crovella and A. Bestavros, "Self-similarity in world wide web traffic: Evidence and causes," in *Proc. ACM SIGMETRICS*, pp. 160–169, (Philadelphia, PA), May 1996.
16. D. Cox and V. Isham, *Point Processes*, Chapman and Hall, 1980.
17. A. D. Trace, "http://pma.nlanr.net/traces/long/auck2.html."
18. Z.-L. Zhang, Z. Duan, and Y. T. Hou, "On scalable design of bandwidth brokers," *IEICE Transaction on Communications* **E84-B**, Aug. 2001.