

Stochastic Traffic Engineering for Demand Uncertainty and Risk-Aware Network Revenue Management

Debasis Mitra, *Fellow, IEEE*, and Qiong Wang, *Member, IEEE*

Abstract—We present a stochastic traffic engineering framework for optimizing bandwidth provisioning and route selection in networks. The objective is to maximize revenue from serving demands, which are uncertain and specified by probability distributions. We consider heterogeneous demands with different unit revenues and uncertainties. Based on mean-risk analysis, the optimization model enables a carrier to maximize mean revenue and contain the risk that the revenue falls below an acceptable level. Our framework is intended for off-line traffic engineering design, which takes a centralized view of network topology, link capacity, and demand. We obtain conditions under which the optimization problem is an instance of convex programming and therefore efficiently solvable. We also study the properties of the solution and show that it asymptotically meets the stochastic efficiency criterion. We derive properties of the optimal solution for the special case of Gaussian distributions of demands. We focus on the impact of demand uncertainty on various aspects of traffic engineering, such as link utilization, bandwidth provisioning and total revenue. The carrier's tolerance to risk is shown to have a strong influence on traffic engineering and revenue management decisions. We develop the *efficient frontier*, which is the entire set of Pareto optimal pairs of mean revenue and revenue risk, to aid the carrier in selecting an appropriate operating point.

Index Terms—Demand uncertainty, economics, mathematical programming, risk, traffic engineering.

I. INTRODUCTION

TRAFFIC engineering (TE) is a mechanism for traffic and revenue management in networks [3], [4], [7], [9]. The TE mechanism takes two complementary forms, on-line and off-line [4], [30]. On-line TE responds to changes of the network state in real-time. See [14] and [28] for works where the focus is on distributed on-line traffic engineering and provisioning. Off-line TE applies on a longer time-scale. Instead of focusing on instantaneous network states and individual connections, it takes as input statistics of aggregated traffic demands between node pairs. Combining this demand information with a centralized view of network topology and link capacities, off-line TE selects the topology of routes and provisions resources on the selected routes for carrying the demands. These decisions are optimized globally for demands of various service types and origin-destination pairs [2], [23], [24], [29]. The solution of the off-line optimization has been proposed as a reference point

for on-line operations. For example, in [29] capacity preallocated by the off-line TE process is used as the threshold for on-line admission control. Similarly, in [8] the off-line provisioning process is used to guide the real-time routing and admission control.

This paper focuses on the optimization of off-line traffic engineering. The optimization in our model is with respect to the topology of the paths serving end-to-end demands and the amount of provisioned bandwidth on these paths, among other decision variables. The problem has previously been formulated as a deterministic multicommodity flow (MCF) model, where demand for each service and node pair is given as a fixed quantity, such as the expected value of forecasted traffic load [23], [29]. The goal is to find an appropriate amount of traffic to admit for each demand, and capacitated route(s) to carry the traffic, so that the carrier's objective, usually formulated as revenue earned by serving demands, is optimized under capacity constraints.

This paper develops a stochastic traffic engineering framework for decision making that takes into account demand uncertainty and its consequence, namely, risk. The overall objective is to employ traffic engineering and revenue management techniques to achieve robustness in network performance and resulting revenue. This framework uses probabilistic distributions of demands as inputs for off-line optimization. Such distributional information is typically a byproduct of statistical procedures for forecasting network traffic from measurements [5], [6], [22], [31]. However, such information has not been used extensively in the past for lack of a suitable modeling framework.

The framework of this paper remedies several shortfalls of the deterministic approach to traffic engineering. For instance, in the deterministic MCF model, revenue derived from carrying demand is assumed to increase linearly with the amount of provisioned capacity up to the point where all traffic demand is satisfied. When demand is random, if more capacity is provisioned then the probability that the capacity is fully utilized decreases. Consequently, mean revenue increases with a decreasing rate with the amount of provisioned capacity. This nonlinear effect is captured in the stochastic traffic engineering framework. Incidentally, this also illustrates one of the reasons why the latter framework is intrinsically nonlinear. Importantly, the framework also allows the variability of uncertainty in the traffic demand to be taken into account, and its impact assessed. Even when mean demands stay fixed, increasing uncertainty significantly affects the optimal traffic engineering solution. Typically the consequences include lower link utilizations and

Manuscript received September 9, 2003; revised January 26, 2004; approved by IEEE/ACM TRANSACTIONS ON NETWORKING Editor A. Orda.

The authors are with Bell Laboratories, Lucent Technologies, Murray Hill, NJ 07974 USA (e-mail: mitra@lucent.com; chiwang@lucent.com).

Digital Object Identifier 10.1109/TNET.2005.845527

higher link shadow costs. Note that optimal selections of path topologies and provisioning of resources are based on link shadow costs [24], hence uncertainty impacts route selection and provisioning.

Uncertainty breeds risk. A carrier must concern itself not only with mean revenue and the strategies for its maximization, but also with the risk of revenue falling below acceptable levels. This risk is not calculable in the deterministic traffic engineering framework. Calculating risk and managing it are important components of network revenue management, and also of this work. The framework here allows the carrier to trade off different objectives, such as the maximization of mean revenue against minimization of the risk of revenue shortfall.

The objective function in the optimization model of this paper incorporates a risk index. The rationale for its selection draws extensively from the mean-risk analysis that was originally developed in the finance community to address the needs of balancing growth and risk in portfolio management [12], [20]. Risk in networking has its own characteristics. At the most basic level, demands of various service types between node pairs are associated with varying degrees of uncertainty; these demands compete for fixed network resources. The allocation of network resources to demands is also subject to service providers' attitude toward risk, some carriers being more risk-averse than others. Taking the carrier's risk aversion into account in decision-making for network revenue management is clearly important, and this is done in the framework of this paper.

Considerable attention is given in this paper to finding a measure of risk that is appropriate specifically for network revenue management in the mean-risk approach. There are several candidates for such a risk index. Some of these candidates, such as the variance of revenue, lead to inferior solutions that are stochastically dominated by other feasible solutions in the context of our problem. On the other hand, there are candidates, such as the Tail Value at Risk (TVaR), which are stochastically efficient but rather difficult to handle analytically and in the optimization. We make the case that the standard deviation of network revenue as a measure of risk is an attractive compromise between stochastic efficiency and tractability.

The framework in the paper is sufficiently general to include demands with different degrees of uncertainty, even for the same (source, destination) pair. Carriers routinely complement retail services by making wholesale contracts with large customers for guaranteed demands at discounted prices; for example, see [16] and [17] regarding Level 3's wholesale agreement with Microsoft. Our framework is used to model the carrier's decision-making in allocating bandwidth to demands that range from higher-priced, but uncertain, to lower-priced, but with little or no uncertainty. A proper mix of these types of demands allows the carrier to mitigate revenue risk while maximizing expected revenue. This capability of our model is demonstrated numerically in Section VI.

As another aid to decision-making, we develop the efficient frontier, which is the complete set of Pareto optimal pairs of mean revenue and revenue risk, and thus is the totality of all rational solutions that the carrier needs to consider. As we will see in Section VI, the shape of the efficient frontier provides

valuable information to the carrier in selecting the appropriate operating point.

With natural extensions, the framework of this paper has the potential of contributing to studies of future developments in the telecommunications industry. For instance, if the industry structure stabilizes to one with several competing carriers with distinct and possibly overlapping footprints, then resource sharing among all or a subset of these carriers will benefit not only the participating carriers but society and consumers also. Modeling frameworks will be needed to quantify the value proposition from such bandwidth sharing agreements. It is possible to envisage carriers being engaged in dynamic, game-theoretic iterations of actions and reactions, in which basic moves of each player are computed by an extension of the present model.

This paper is organized as follows. In Section II we formulate the stochastic traffic engineering problem and present the optimization model. In Section III we discuss modeling of the risk term in the objective function and in Section IV we analyze the complexity of obtaining the global optimal solution. In Section V we derive properties of the optimal solution that provide important insights on bandwidth provisioning and route selection. In Section VI we give numerical results that show the impact of demand uncertainty and carriers' risk tolerance on network traffic engineering and revenue management. We also develop the efficient frontier. We present our concluding remarks in Section VII.

II. MODEL

A. Conception

1) *Demand, Provisioning, and Revenue Management:* In this paper, we define demand as the aggregated amount of traffic volume that customers are willing to pay the carrier to deliver from one network node to another. The price per unit of delivered traffic is assumed to be given. It is sometimes useful to define multiple demands for the same node pair, with the demands distinguished by the unit price, quality of service requirements, and volume statistics. For example, in [25] and [26], for the same end nodes, there is a distinction between retail demand, which generates higher unit revenue but for which the volume is subject to uncertainty, and wholesale demand for which the volume is known with certainty.

In the case of deterministic traffic engineering, the value of a demand is characterized by its unit price. In the presence of demand uncertainty, unit price alone is not sufficient to characterize the value of a demand. It is possible that a demand of high unit price is very volatile, in which case it is necessary to take into account the probability that demand volume will fall below the provisioned bandwidth. Consequently, the carrier should anticipate the possibility of a smaller revenue from such demands than from demands with lower unit prices but guaranteed volume. To optimally provision bandwidth in our framework, unit price should be used in conjunction with distributional information of demands to characterize the mean revenue return and revenue risk from bandwidth provisioning.

Revenue management is further complicated by the carrier's attitude toward risk. In the presence of demand uncertainty, the carrier may want to give up a certain amount of mean revenue

in exchange for reduced level of revenue risk. The difference in carrier's willingness to trade mean revenue with risk should be reflected in a revenue management model, as is the case in this paper.

We digress to observe that in the model here prices are not influenced by capacity provisioning decisions. It is possible to couple the two through price-demand relationships, as in [18] and [24], for instance. The value in these models is in decision-making over time scales that are longer than is of interest in this paper.

2) *Admissible Route Sets*: QoS and policy considerations are major constraints on provisioning decisions. The notion of *admissible route sets* allows these constraints to be taken into account in the optimization. For example, routes may be required to have lengths not exceeding specified thresholds, on account of propagation delay, and there may also be restrictions on the number of hops, since each hop is associated with a switching node and consequent incremental delay. The admissibility of a route may also depend on policy, which typically reflects diverse considerations, such as security, the capability of switching nodes in the routes to handle certain services, etc. Generating the admissible route sets is a substantial task in itself. It is assumed in this paper that these sets are given.

B. Model Formulation

We formulate the network as a collection of nodes and links $(\mathcal{N}, \mathcal{L})$, where link $l \in \mathcal{L}$ has bandwidth C_l . Let \mathcal{V} be the set of all demands and denote demand volume by $T_v (v \in \mathcal{V})$. In a restricted case where there is only one demand between each pair of network nodes, \mathcal{V} corresponds to a subset of all node pairs. However, in this paper, each source-destination pair is allowed to have multiple demands, in which case $v \in \mathcal{V}$ specifies both node pair and demand type. Different demand types for the same node pair may be characterized by the unit price, admissible route set, GoS requirement, as well as traffic volume distribution that reflects the degree of uncertainty.

Let demand volume T_v be a random variable characterized by its probability density function (PDF) and cumulative distribution function (CDF), denoted by $f_v(x)$ and $F_v(x)$, respectively. Let $d_v (v \in \mathcal{V})$ be the amount of capacity provisioned to demand v . The provisioned quantity, d_v , may be routed on one or more admissible routes. Denote the admissible route set for $v \in \mathcal{V}$ by $\mathcal{R}(v)$ and let $\xi_r (r \in \mathcal{R}(v))$ be the amount of capacity provisioned on route r . Then

$$d_v = \sum_{r \in \mathcal{R}(v)} \xi_r. \quad (1)$$

It follows that the amount of carried demand

$$x_v(d_v) = \min(T_v, d_v). \quad (2)$$

Let $m_v(d_v)$ and $s_v^2(d_v)$ be the mean and variance of x_v , respectively. Then

$$m_v(d_v) = \int_0^{d_v} x f_v(x) dx + d_v \bar{F}_v(d_v) = \int_0^{d_v} \bar{F}_v(x) dx \quad (3)$$

$$\begin{aligned} s_v^2(d_v) &= \int_0^{d_v} x^2 f_v(x) dx + d_v^2 \bar{F}_v(d_v) - m_v^2(d_v) \\ &= 2 \int_0^{d_v} x \bar{F}_v(x) dx - m_v^2(d_v) \end{aligned} \quad (4)$$

where $\bar{F}_v(x) \equiv 1 - F_v(x)$. Notice that

$$\begin{aligned} \frac{\partial m_v}{\partial d_v} &= \bar{F}_v(d_v) \geq 0 \\ \frac{\partial s_v^2}{\partial d_v} &= 2[d_v - m_v(d_v)] \bar{F}_v(d_v) \geq 0 \end{aligned} \quad (5)$$

i.e., both the mean and variance of carried demand increase with the amount of bandwidth provisioned. Their maximum values, denoted by $m_v(\infty)$ and $s_v^2(\infty)$, are the mean and variance of the demand v , respectively.

It is worth noting that the model for deterministic traffic engineering is a special case of the model above. In the former case, network demand is given by the "demand matrix" of constants \bar{T}_v . We formulate the demand matrix as an instance of our model by using the following degenerate distribution function:

$$F_v(x) = \begin{cases} 0, & x < \bar{T}_v \\ 1, & x \geq \bar{T}_v. \end{cases}$$

In this case, from (3) and (4)

$$m_v(d_v) = \min(\bar{T}_v, d_v) \quad s_v^2(d_v) = 0. \quad (6)$$

Our model also accommodates mixed cases where both deterministic and random demands exist [25], [26].

Given fixed unit revenue, π_v , the revenue from serving demand v equals $\pi_v x_v(d_v)$, for which the mean is $\pi_v m_v(d_v)$ and the variance is $\pi_v^2 s_v^2(d_v)$. The total revenue is

$$W = \sum_{v \in \mathcal{V}} \pi_v x_v(d_v). \quad (7)$$

In the general case, $x_v(d_v)$ is a random variable for at least some $v \in \mathcal{V}$, and so is W . The objective function, denoted by Θ , is itself a function of W . The composition of Θ should reflect the carrier's desire to increase mean revenue and reduce the risk of revenue shortfall due to the randomness of W . In the next section, we discuss the specification of Θ that satisfies this objective.

The overall optimization model is as follows:

$$\max_{d_v, \xi_r} \Theta(W) \quad (8)$$

subject to

$$\sum_{r \in \mathcal{R}(v)} \xi_r = d_v \quad (v \in \mathcal{V}) \quad (9)$$

$$\sum_{v \in \mathcal{V}} \sum_{r \in \mathcal{R}(v): l \in r} \xi_r \leq C_l, \quad l \in \mathcal{L} \quad (10)$$

$$0 < \underline{d}_v \leq d_v (v \in \mathcal{V}) \quad 0 \leq \xi_r (r \in \mathcal{R}(v)) \quad (11)$$

The constraints in (9) specify bandwidth provisioning for demands on admissible routes, and (10) limit the total amount of provisioned bandwidth from exceeding link capacities. The lower bounds (11) on the provisioned bandwidth deserve additional comment. The parameter \underline{d}_v is defined as the minimum bandwidth that must be provisioned for demand v . The parameter value is determined by factors such as regulation, service level agreements, and customer expectations on the grade of service.

Notice that with fixed link capacities, it may be infeasible to satisfy the minimum bandwidth requirement of some demands. Should this situation arise, the carrier has to supplement link capacities by reconfigurations or by buying bandwidth from other carriers. Our model can be generalized to incorporate buying decisions to cover this scenario [26]. However, the generalized model is not the focus of this paper. We will analyze the original model in (8)–(11), and assume that a feasible solution exists for the existing link capacities.

III. MODELING RISK

In this section, we take uncertainty in demand into consideration in the composition of an objective function that reflects both the maximization of mean revenue and the containment of the risk of revenue shortfall. We review two relevant risk modeling frameworks in Section III-A and discuss our characterization of revenue risk in Section III-B.

A. Relevant Frameworks for Risk Modeling

1) *Mean-Risk Model*: Mean-risk analysis addresses the issue of risk averseness by offering a broader optimization objective. The approach starts by developing a risk index, which quantifies the risk of revenue shortfall based on the revenue distribution. It then maximizes the weighted combination of the mean revenue and the risk index, i.e., mean $- \delta$ (risk index), where $\delta \geq 0$ is a parameter. Different levels of risk averseness can be reflected by choosing different values for δ . A higher value of δ indicates greater willingness to sacrifice the mean revenue to avoid risk.

2) *Stochastic Dominance*: Stochastic dominance theory defines a partial ordering of random variables based on their probability distributions [15]. Let W_1 and W_2 be two random variables, which represent revenues under two different bandwidth management decisions. Then W_1 stochastically dominates W_2 to the first degree iff the former renders the carrier a better chance to exceed *any* revenue target w , i.e.,

$$\forall w \quad \Pr(W_1 \geq w) \geq \Pr(W_2 \geq w). \quad (12)$$

Furthermore, W_1 stochastically dominates W_2 to the second degree iff

$$\int_w^\infty \Pr(W_1 \geq \zeta) d\zeta \geq \int_w^\infty \Pr(W_2 \geq \zeta) d\zeta \quad \forall w. \quad (13)$$

We consider a solution to our problem to be stochastically efficient if the corresponding revenue distribution is not dominated in either degree. It suffices to prove efficiency by showing that the revenue distribution is not subject to second-degree

dominance, which, by definition, is a necessary condition for dominance in the first degree.

B. Formulation of the Objective Function

We formulate the objective function as an instance of the mean-risk model and require that the solution that optimizes it be stochastically efficient. Whether these conditions can be met depends on the characteristics of the revenue distribution and the choice of the risk index. In portfolio management, variance is commonly used as the risk index and the objective is to maximize (mean $- \delta$ variance) [15]. However, as we show in Section III-B.1, the mean-variance model in network revenue management often leads to solutions that are stochastically dominated. On the other hand, applying other risk measures that guarantee stochastic efficiency, such as Tail Value at Risk defined in Section III-B.2, lead to models that are too complex to solve. As a compromise, in Section III-B.3 we propose standard deviation as the risk index.

B.1 Variance: Variance of network revenue is a natural candidate for the risk index, especially since it has been widely used in many mean-risk models. The objective function that uses the variance as the risk index is

$$\Theta(W) = E(W) - \delta \text{Var}(W) = \sum_{v \in \mathcal{V}} \pi_v [m_v(d_v) - \delta \pi_v s_v^2(d_v)] \quad (14)$$

Notice that (14) reflects the aforementioned assumption that demands between different node pairs are independent. Applying (3) and (4)

$$\begin{aligned} \frac{\partial \Theta}{\partial d_v} &= \pi_v \bar{F}_v(d_v) [1 - 2\delta \pi_v (d_v - m_v)] \\ &= \pi_v \bar{F}_v(d_v) \left[1 - 2\delta \pi_v \int_0^{d_v} F_v(x) dx \right]. \end{aligned} \quad (15)$$

The equation $2\delta \pi_v \int_0^{d_v} F_v(x) dx = 1$ has a unique solution. Therefore, there exists a unique \hat{d}_v such that

$$\frac{\partial \Theta}{\partial d_v} \geq 0 \quad \frac{\partial^2 \Theta}{\partial d_v^2} \leq 0, \quad \text{if } d_v \leq \hat{d}_v$$

and

$$\frac{\partial \Theta}{\partial d_v} < 0, \quad \text{if } d_v > \hat{d}_v. \quad (16)$$

Based on this observation, we impose $\max(\underline{d}_v, \hat{d}_v)$ as an upper bound on d_v . If $\hat{d}_v \leq \underline{d}_v$, then $d_v = \underline{d}_v$ by (11). Otherwise, $d_v \leq \hat{d}_v$, which is a new constraint that reduces the original feasible region [a polyhedron defined by (9)–(11)] to a convex set on which the objective function (14) is concave, which makes our model a concave maximization problem that can be solved efficiently. This additional restriction on d_v does not affect the optimal solution since by (16), increasing d_v beyond \hat{d}_v only reduces the value of Θ .

We now show that the optimal solution of the mean-variance model can be stochastically dominated. Note that the optimal value of d_v is bounded by $\max(\underline{d}_v, \hat{d}_v)$. When the network has as much capacity as this bound for every demand, and has more for at least one demand, a new feasible solution can be devised that will carry as much traffic for each demand and more in total

than the optimal solution to the mean-variance model. Consequently, the new solution generates more revenue, and thus dominates the optimal solution to the first degree.

B.2 Tail Value at Risk: Besides variance, other distributional parameters have also been proposed as candidates for the risk index. Of particular interest are indexes that guarantee stochastic efficiency of the optimal solution. One example is the Tail Value at Risk (TVaR), defined as

$$\text{TVaR}_W(p) = \int_0^p \frac{q_W(\eta) d\eta}{p} \quad (17)$$

where $q_W(\eta) = \inf\{w | \Pr(W \leq w) \geq \eta\}$ is the η -quantile of revenue W . As in [15], the dominance condition (13) is equivalent to

$$\int_0^p q_{W_1}(\eta) d\eta \geq \int_0^p q_{W_2}(\eta) d\eta \quad \forall 0 < p \leq 1 \quad (18)$$

indicating that if W_1 dominates W_2 , then $\text{TVaR}_{W_1}(p) \geq \text{TVaR}_{W_2}(p)$ for all p . For instance, in the above discussion on variance, we illustrated an inferior solution that is optimal in the mean-variance model but does not use all available bandwidth. We showed that the resulting revenue is dominated by another solution which uses all capacity. It is easy to verify that the latter solution has higher TVaR for every p . Note that condition (18) implies that revenue W_0 is stochastically efficient if for some p , $\text{TVaR}_W(p)$ has a unique maximum at W_0 .

As in [27], the mean-risk model that uses $-\text{TVaR}_W(p)$ as the risk index

$$\Theta(W) = E(W) + \delta \text{TVaR}_W(p) \quad (19)$$

maximizes the weighted sum of $E(W)$ (which equals $\text{TVaR}(1)$) and $\text{TVaR}(p)$. Consequently, any other solution must be smaller either in $\text{TVaR}(1)$ or $\text{TVaR}(p)$, and thus cannot stochastically dominate the optimal solution.

However, optimizing the above requires specifying quantiles of the revenue distribution (i.e., the distribution of W) as functions of the demand distributions (i.e., distributions of T_v) and the provisioned bandwidth (d_v). This approach introduces a significant amount of computational complexity, and is therefore not adopted in this paper.

B.3 Standard Deviation as Risk Index: Our choice of risk index is the standard deviation of network revenue, $s(W)$, i.e.,

$$\Theta(W) = E(W) - \delta s(W). \quad (20)$$

It is computationally more tractable for optimization than TVaR. For networks with few node-pair demands, an optimal solution to the mean-standard deviation model may still be stochastically dominated. Now consider the difference between the optimal solution that uses the standard deviation as the risk index and the optimal solution that maximizes TVaR for some p . This difference approaches zero as the number of node pairs increases. This property serves our purpose well, since bandwidth transport networks usually have tens or even hundreds of nodes, and thus a large number of node pairs.

As in (7), W is a summation of many *independent* random revenues, $\pi_v x_v(d_v)$. From (5), the variance of each revenue,

$\pi_v^2 s_v^2(d_v)$, monotonically increases with d_v , and thus has a lower bound $\pi_v^2 s_v^2(\underline{d}_v)$ from (11), and an upper bound $\pi_v^2 s_v^2(+\infty) = \pi_v^2 \text{var}(T_v)$. Therefore, for every feasible solution

$$\lim_{|\mathcal{V}| \rightarrow +\infty} \frac{\pi_{v'}^2 s_{v'}^2(d_{v'})}{\sum_{v \in \mathcal{V}} \pi_v^2 s_v^2(d_v)} = 0 \quad \forall v' \in \mathcal{V}$$

i.e., Lindeberg's condition is satisfied, and the Central Limit Theorem can be applied as follows:

$$\frac{W - E(W)}{s(W)} = \frac{\sum_{v \in \mathcal{V}} \pi_v [x_v(d_v) - m_v(d_v)]}{\sqrt{\sum_{v \in \mathcal{V}} \pi_v^2 s_v^2(d_v)}} \rightarrow N(0, 1) \quad \text{as } |\mathcal{V}| \rightarrow +\infty. \quad (21)$$

where $N(0, 1)$ is the standard normal distribution. It follows that if $|\mathcal{V}|$ is sufficiently large

$$\begin{aligned} \text{TVaR}_W(p) &= \int_0^p \frac{q_W(\eta) d\eta}{p} = \int_{-\infty}^{q_W(p)} \frac{x dF_W(x)}{p} \\ &\sim \int_{-\infty}^{q_N(p)} \frac{[E(W) + s(W)\theta] dF_N(\theta)}{p} \\ &= E(W) - \frac{e^{-\frac{q_N^2(p)}{2}}}{p\sqrt{2\pi}} s(W) \end{aligned} \quad (22)$$

where $q_N(p)$ is the p -quantile and $F_N(\theta)$ is the CDF of $N(0,1)$. Since $e^{-q_N^2(p)/2}/(p\sqrt{2\pi})$ changes monotonically from 0 to $+\infty$ as p varies from 1 to 0, $\delta = e^{-q_N^2(p)/2}/(p\sqrt{2\pi})$ has a unique inverse $p(\delta)$. By (22), the solution that maximizes $E(W) - \delta s(W)$ is stochastically efficient since it also maximizes $\text{TVaR}[p(\delta)]$.

IV. MODEL ANALYSIS AND PROPERTIES

The preceding section has established that the standard deviation is an appropriate choice of the risk index in large networks. In this section we demonstrate that the use of standard deviation also leads to a tractable optimization model. We also discuss the implications of the necessary conditions for optimality.

Using the standard deviation as the risk index, the objective function in (8) becomes

$$\max \Theta(W(d_v, \xi_r)) = \sum_{v \in \mathcal{V}} \pi_v m_v(d_v) - \delta \sqrt{\sum_{v \in \mathcal{V}} \pi_v^2 s_v^2(d_v)}. \quad (23)$$

Given that all constraints of the above model are linear, the complexity of finding the global optimum depends on the shape of the objective function Θ . If Θ is concave in all nonlinear variables $d_v (v \in \mathcal{V})$, the model falls into the class of concave maximization problems. In this case, the global optimum can be found efficiently with existing standard algorithms.

In general, Θ is not concave everywhere if $\delta > 0$, as can be verified by considering a restricted case with only one nonlinear variable. Despite this inconvenience, we show in this section that in many circumstances, the model can still be solved as a concave maximization problem. Our approach is based on two theorems developed in Section IV-A and the subsequent analysis in Section IV-B.

A. Shape of the Objective Function

In the following, we use m_v and s_v^2 to represent $m_v(d_v)$ and $s_v^2(d_v)$, as defined in (3) and (4), respectively. We also omit the

argument d_v in the distribution functions, and use F_v , \bar{F}_v , and f_v to represent $F_v(d_v)$, $\bar{F}_v(d_v)$, and $f_v(d_v)$, respectively. Lemma 1 is the basis for Theorems 1 and 2.

Lemma 1: For any $v \in \mathcal{V}$ and $d_v \geq 0$

$$F_v s_v^2 \geq (d_v - m_v)^2 \bar{F}_v. \quad (24)$$

Proof: Because $f_v \geq 0$, from (3) and (4)

$$\begin{aligned} \frac{\partial (F_v s_v^2)}{\partial d_v} &= 2F_v \bar{F}_v (d_v - m_v) + f_v s_v^2 \\ &\geq 2F_v \bar{F}_v (d_v - m_v) - f_v (d_v - m_v)^2 \\ &= \frac{\partial [(d_v - m_v)^2 \bar{F}_v]}{\partial d_v}. \end{aligned} \quad (25)$$

Since $F_v s_v^2 = (d_v - m_v)^2 \bar{F}_v = 0$ at $d_v = 0$

$$F_v s_v^2 \geq (d_v - m_v)^2 \bar{F}_v \quad \forall d_v \geq 0. \quad (26)$$

□

We show in Theorem 1 that Θ is unimodal in every d_v .

Theorem 1: For $v \in \mathcal{V}$, given fixed values for $d_{v'} (v' \neq v)$

$$\frac{\partial \Theta}{\partial d_v} = \begin{cases} \geq 0 & \text{if } 0 \leq d_v \leq \hat{d}_v \\ < 0 & \text{if } d_v > \hat{d}_v \end{cases} \quad (27)$$

$$\text{where } \hat{d}_v = \sup \left\{ d_v : \frac{(d_v - m_v)^2}{s_v^2 + \Psi_v} \leq \frac{1}{\delta^2} \right\}. \quad (28)$$

Here $\Psi_v = \sum_{v' \neq v} (\pi_{v'} / \pi_v)^2 s_{v'}^2$, and \hat{d}_v is the unique solution to (28), i.e., \hat{d}_v is calculated by solving the equation obtained by replacing \leq by $=$.

Proof: Let $S^2 = \text{Var}(W) = \sum_{v \in \mathcal{V}} \pi_v^2 s_v^2$

$$\begin{aligned} \frac{\partial \Theta}{\partial d_v} &= \pi_v \frac{\partial m_v}{\partial d_v} - \delta \frac{\partial S}{\partial d_v} \\ &= \pi_v \bar{F}_v - \delta \pi_v \bar{F}_v \frac{(d_v - m_v)}{\frac{S}{\pi_v}} \\ &= \pi_v \bar{F}_v \left[1 - \delta \frac{(d_v - m_v)}{\sqrt{s_v^2 + \Psi_v}} \right]. \end{aligned} \quad (29)$$

Notice that by Lemma 1

$$\frac{\partial \left[\frac{(d_v - m_v)}{\sqrt{s_v^2 + \Psi_v}} \right]}{\partial d_v} = \frac{F_v (s_v^2 + \Psi_v) - (d_v - m_v)^2 \bar{F}_v}{\sqrt{(s_v^2 + \Psi_v)^3}} \geq 0. \quad (30)$$

Therefore $(d_v - m_v) / \sqrt{s_v^2 + \Psi_v}$ monotonically increases from 0 to $+\infty$ as d_v goes from 0 to $+\infty$. It follows that \hat{d}_v as defined by (28) is unique, and $\partial \Theta / \partial d_v >, =, \text{ or } < 0$, depending on whether $d_v <, =, \text{ or } > \hat{d}_v$. □

Holding (28) at equality and applying the Implicit Function Theorem

$$\frac{\partial \hat{d}_v}{\partial \Psi_v} = \frac{d_v - m_v}{F_v (s_v^2 + \Psi_v) - (d_v - m_v)^2 \bar{F}_v} \geq 0 \quad (31)$$

indicating that \hat{d}_v increases with Ψ_v . This result will be used in Section IV-B.

Clearly, any maximum point of Θ can be reached only in areas where $d_v \leq \hat{d}_v$ for all $v \in \mathcal{V}$. Theorem 2 shows that Θ is concave in this region. Before presenting the theorem, we first give the second-order derivatives

$$\begin{aligned} \frac{\partial^2 \Theta}{\partial d_v^2} &= -\pi_v f_v \left[1 - \delta \frac{(d_v - m_v)}{\frac{S}{\pi_v}} \right] \\ &\quad - \delta \frac{\pi_v \bar{F}_v}{\frac{S}{\pi_v}} \left[F_v - \left(\frac{(d_v - m_v)}{\frac{S}{\pi_v}} \right)^2 \bar{F}_v \right] \quad v \in \mathcal{V} \end{aligned} \quad (32)$$

and for $v \neq v'$

$$\frac{\partial^2 \Theta}{\partial d_v \partial d_{v'}} = \frac{\delta}{S^3} \pi_v^2 \pi_{v'}^2 (d_v - m_v) \bar{F}_v (d_{v'} - m_{v'}) \bar{F}_{v'}. \quad (33)$$

Theorem 2: Let $H(\Theta)$ be the Hessian matrix of Θ . If

$$\frac{\partial \Theta}{\partial d_v} = \pi_v \bar{F}_v \left[1 - \delta \frac{(d_v - m_v)}{\frac{S}{\pi_v}} \right] \geq 0 \quad \forall v \in \mathcal{V} \quad (34)$$

then $H(\Theta)$ is negative semi-definite.

Proof:

$$H(\Theta) = \begin{pmatrix} H_1 & 0 \\ 0 & 0 \end{pmatrix} \quad (35)$$

where H_1 is a square matrix of dimension $n = |\mathcal{V}|$, and its entries are $\partial^2 \Theta / \partial d_{v_i} \partial d_{v_j} (v_i, v_j \in \mathcal{V})$. Other elements that take the zero value correspond to second-order derivatives with respect to ξ_r , which are linear variables in our model. $H(\Theta)$ is negative semi-definite if and only if H_1 is negative semi-definite.

Apply (32) and (33), and note that $S^2 = \pi_v^2 s_v^2 + \sum_{v' \neq v} \pi_{v'}^2 s_{v'}^2$ for any $v \in \mathcal{V}$, we have

$$H_1 = -H_{1a} - \frac{\delta}{S^3} H_{1b} - \frac{\delta}{S^3} H_{1c} \quad (36)$$

where elements of matrix H_{1a} , H_{1b} , and H_{1c} , denoted correspondingly as $h_a(i, j)$, $h_b(i, j)$, and $h_c(i, j)$. H_{1a} and H_{1b} are diagonal matrices with

$$h_a(i, i) = \pi_{v_i} f_{v_i} \left[1 - \delta \frac{(d_{v_i} - m_{v_i})}{\frac{S}{\pi_{v_i}}} \right]$$

and

$$h_b(i, i) = \pi_{v_i}^4 \bar{F}_{v_i} \left[s_{v_i}^2 F_{v_i} - (d_{v_i} - m_{v_i})^2 \bar{F}_{v_i} \right]. \quad (37)$$

Also,

$$h_c(i, i) = \pi_{v_i}^2 \bar{F}_{v_i} F_{v_i} \left(\sum_{v \neq v_i} \pi_v^2 s_v^2 \right)$$

and for $i \neq j$

$$h_c(i, j) = -\pi_{v_i}^2 \pi_{v_j}^2 (d_{v_i} - m_{v_i}) \bar{F}_{v_i} (d_{v_j} - m_{v_j}) \bar{F}_{v_j}. \quad (38)$$

H_{1a} is positive semi-definite by assumption, and H_{1b} is positive semi-definite by Lemma 1. To prove that $H(\Theta)$ is negative semi-definite, it suffices to show that H_{1c} is positive semi-definite, i.e., for any real vector \vec{X} of dimension $n = |\mathcal{V}|$

$$\begin{aligned} \vec{X}^T H_{1c} \vec{X} &= \sum_{v \in \mathcal{V}} x_v^2 \pi_v^2 \bar{F}_v F_v \left(\sum_{v' \neq v} \pi_{v'}^2 s_{v'}^2 \right) \\ &\quad - \sum_{v' \neq v} x_v x_{v'} \pi_v^2 \pi_{v'}^2 (d_v - m_v) \\ &\quad \times \bar{F}_v (d_{v'} - m_{v'}) \bar{F}_{v'} \geq 0. \end{aligned} \quad (39)$$

This is true because by Lemma 1

$$\sqrt{\bar{F}_v F_v \bar{F}_{v'} F_{v'} s_v^2 s_{v'}^2} \geq (d_v - m_v) \bar{F}_v (d_{v'} - m_{v'}) \bar{F}_{v'} \quad (40)$$

and by the A-G Mean Inequality (i.e., $a^2 + b^2 \geq 2ab$)

$$\begin{aligned} \sum_v x_v^2 \pi_v^2 \bar{F}_v F_v \left(\sum_{v' \neq v} \pi_{v'}^2 s_{v'}^2 \right) \\ \geq \sum_{v' \neq v} \pi_v^2 \pi_{v'}^2 x_v x_{v'} \sqrt{\bar{F}_v F_v \bar{F}_{v'} F_{v'} s_v^2 s_{v'}^2}. \end{aligned} \quad (41)$$

□

B. Discussion of Properties

The two theorems in the last section define a set

$$\Omega \equiv \left\{ d_v(v \in \mathcal{V}) : \underline{d}_v \leq d_v \leq \hat{d}_v \right\} \quad (42)$$

over which Θ is concave (if $\underline{d}_v \geq \hat{d}_v$, then d_v is fixed at \underline{d}_v). The set Ω also contains all local maximum point(s) of Θ . The optimal value of the objective function is obtained by $d_v(v \in \mathcal{V})$ from the set Ω and $\xi_r(r \in \mathcal{R}(v), \xi_r \geq 0)$. The model is a concave maximization problem if either Ω or its intersection with the feasible region is a convex set. In this case, we can find the global optimum efficiently.

The set Ω may or may not be convex, depending on the demand distributions, parameters \underline{d}_v , and δ . To verify that Ω may not be convex, consider a single-link network that serves two demands. Assume that volumes of both demands are uniformly distributed over $[0,1]$, and the ratio of their unit prices $\pi_2/\pi_1 = 2$. Let $\delta = 1$ and $\underline{d}_1 = \underline{d}_2 = 0.1$. From (28) and (42), Ω is the set of (d_1, d_2) that satisfies

$$d_1 \geq 0.1, \quad d_2 \geq 0.1$$

and

$$3d_1^4 - 2d_1^3 \leq 8d_2^3 - 6d_2^4, \quad 6d_2^4 - 4d_2^3 \leq d_1^3 - \frac{3d_1^4}{4}$$

The set is not convex. For instance, it contains $(0.67, 0.11)$ and $(0.79, 0.33)$ but not their average $(0.73, 0.22)$.

In case Ω is not convex, we can still solve the model as a concave maximization problem under certain conditions. Consider $\hat{d}_v(v \in \mathcal{V})$ in Theorem 1, which are obtained by solving

$$\frac{[\hat{d}_v - m_v(\hat{d}_v)]^2}{s_v^2(\hat{d}_v) + \Psi_v} = \frac{1}{\delta^2} \quad \Psi_v = \sum_{v' \neq v} \left(\frac{\pi_{v'}}{\pi_v} \right)^2 s_{v'}^2(d_{v'}). \quad (43)$$

Since $\underline{d}_v > 0$ is the minimum amount of provisioned bandwidth for v , and $s_v^2(d_{v'})$ increases with $d_{v'}$

$$\bar{\Psi}_v = \sum_{v' \neq v} \left(\frac{\pi_{v'}}{\pi_v} \right)^2 s_{v'}^2(\underline{d}_{v'})$$

is the lower bound of Ψ_v . For each v , we solve (43) in which $\Psi_v = \bar{\Psi}_v$ and denote the solution by \hat{d}_v^e , i.e.,

$$\frac{[\hat{d}_v^e - m_v(\hat{d}_v^e)]^2}{s_v^2(\hat{d}_v^e) + \bar{\Psi}_v} = \frac{1}{\delta^2}. \quad (44)$$

Then $\hat{d}_v^e \leq \hat{d}_v$, because as mentioned in IV-A, \hat{d}_v that solves (43) increases with Ψ_v . Consequently

$$\Omega^e \equiv \left\{ d_v(v \in \mathcal{V}) : \underline{d}_v \leq d_v \leq \hat{d}_v^e \right\}$$

is a subset of Ω . If Ω^e is large enough to contain the global optimum, the model can be solved by using a concave maximization algorithm on the polyhedron.

We now discuss conditions for the global optimum to be “trapped” inside Ω^e . Define $k_v = \bar{\Psi}_v/s_{v,\infty}^2$, where $1/k_v$ is the ratio of the maximum variance of the revenue from node pair v to the minimum variance of the revenue from all other node pairs. Because $s_{v,\infty}^2 \geq s^2(\hat{d}_v^e)$

$$\frac{1}{\delta^2} = \frac{[\hat{d}_v^e - m_v(\hat{d}_v^e)]^2}{s_v^2(\hat{d}_v^e) + \bar{\Psi}_v} \leq \frac{[\hat{d}_v^e - m_v(\hat{d}_v^e)]^2}{(k_v + 1)s_v^2(\hat{d}_v^e)}.$$

By Lemma 1, $(1/\delta^2) = ([\hat{d}_v^e - m_v(\hat{d}_v^e)]^2/s_v^2(\hat{d}_v^e) + \bar{\Psi}_v) \leq (F_v(\hat{d}_v^e)/(k_v + 1)\bar{F}_v(\hat{d}_v^e))$, so $F_v(\hat{d}_v^e) \geq k_v + 1/k_v + 1 + \delta^2$. If $k_v \gg \delta^2$ for each v , then $F_v(\hat{d}_v^e) \approx 1$, and it becomes certain that Ω^e contains the global optimum. For example, when $k_v = 20$, $\delta = 1$, then $F_v(\hat{d}_v^e) \geq 0.95$. This means that if it is uneconomical to provision bandwidth for demand in excess of the 95% percentile, then the optimal solution can be obtained by concave maximization over Ω^e .

For the condition $k_v \gg \delta^2$ to hold, either δ is small (in an extreme case when $\delta = 0$, $\hat{d}_v^e = \infty$), or k_v is large. The latter corresponds to situations where the contribution of each node pair to the variance of total revenue is insignificant. This happens when the network has many node pairs, and the total revenue is not dominated by the revenue from an individual pair. This is typically the real-world case for large networks. For smaller networks, straightforward enumeration of extreme and boundary points is a viable option.

C. Necessary Conditions for Optimality and Implications

The Lagrangian of our model takes the following form:

$$\Lambda = \sum_{v \in \mathcal{V}} \pi_v m_v - \delta \sqrt{\sum_{v \in \mathcal{V}} \pi_v^2 s_v^2} + \sum_{v \in \mathcal{V}} \chi_v \left(\sum_{r \in \mathcal{R}(v)} \xi_r - d_v \right) + \sum_{l \in \mathcal{L}} \lambda_l \left(C_l - \sum_{v \in \mathcal{V}} \sum_{r \in \mathcal{R}(v): l \in r} \xi_r \right). \quad (45)$$

It follows that the first-order necessary conditions are [using (29)]

$$\left. \begin{aligned} \frac{\partial \Lambda}{\partial d_v} = \pi_v \bar{F}_v \left[1 - \delta \frac{\pi_v (d_v - m_v)}{S} \right] - \chi_v &\leq 0 \\ d_v \left(\frac{\partial \Lambda}{\partial d_v} \right) = 0, \quad d_v \geq 0, \quad \chi_v \geq 0 &(v \in \mathcal{V}) \end{aligned} \right\} \quad (46)$$

$$\left. \begin{aligned} \frac{\partial \Lambda}{\partial \xi_r} = \chi_v - \sum_{l \in r} \lambda_l &\leq 0 \\ \xi_r \left(\frac{\partial \Lambda}{\partial \xi_r} \right) = 0, \quad \xi_r \geq 0, \quad \lambda_l \geq 0 &\forall r, l. \end{aligned} \right\} \quad (47)$$

The quantity λ_l is interpreted as the *link shadow cost*, which reflects the marginal value of capacity on link l and corresponds to dual variables in [1]. It is a critical quantity for route selection [1], for real-time routing and rate control [13], [33], and for joint off-line optimization of pricing and provisioning [24]. Specifically

1. By (47), for any route $r_0 \in \mathcal{R}(v)$, $\xi_{r_0} > 0$ only when

$$\sum_{l \in r_0} \lambda_l = \min_{r \in \mathcal{R}(v)} \sum_{l \in r} \lambda_l$$

indicating that traffic is carried solely on the *shortest* path(s) of all the admissible routes, where link shadow cost, λ_l , is the distance metric.

2. We define $\chi_v = \min_{r \in \mathcal{R}(v)} \sum_{l \in r} \lambda_l$ as the *demand shadow cost*, i.e., the opportunity cost of carrying demand between $v \in \mathcal{V}$. By (46), the optimal quantity to be provisioned for demand v is determined at the point where the marginal increase of mean revenue, $\pi_v \bar{F}_v$, compensated by the marginal change of risk, $\delta \pi_v (d_v - m_v) / \sqrt{\sum_{v \in \mathcal{V}} \pi_v^2 s_v^2}$, equals the demand shadow cost, χ_v .

V. SPECIAL CASE: TRUNCATED GAUSSIAN DISTRIBUTION

When demand between a node pair comes from many independent individual sources, the total demand can be approximated by the Gaussian distribution. That the aggregated traffic demand between network nodes follows the Gaussian distribution in real networks has been extensively observed and documented, specially in studies that extract traffic models from measurements and inference techniques, such as network tomography. See for instance [5], [6], [21], [22]. We will make this assumption in what follows. Of course, the distribution needs to be restricted to nonnegative values, and the PDF should also be normalized properly, so that the total probability over the restricted sample space is unity. As a result, we will consider the

Truncated Gaussian Distribution characterized by the following PDF:

$$f_v(x) = \frac{1}{\sqrt{2\pi}\sigma_v G_v} e^{-\frac{(x-\mu_v)^2}{2\sigma_v^2}}, \quad x \geq 0 \quad (48)$$

where the normalizing parameter is

$$G_v = \frac{Erfc(-\tau_v)}{2} \quad \text{and} \quad \tau_v = \frac{\mu_v}{\sqrt{2}\sigma_v} \quad (49)$$

where $Erfc() = 1 - Erf()$.

In the rest of this paper, we will assume that μ_v is sufficiently larger than σ_v ($\sigma_v \leq 0.35 \mu_v$, in our numerical investigations) so that the truncation effect is negligible. In this case, μ_v and σ_v approximately equal the mean and standard deviation of the demand distribution, and are used as proxies for these parameters.

Let d_v be the amount of bandwidth provisioned for demand v . Then the mean and standard deviation of carried demand are

$$\begin{aligned} m_v(d_v) &= \frac{1}{\sqrt{2\pi}\sigma_v G_v} \int_0^\infty \min(x, d_v) e^{-(x-\mu_v)^2/2\sigma_v^2} dx \\ &= \mu_v + \sigma_v \gamma(\tilde{d}_v) \end{aligned} \quad (50)$$

$$\begin{aligned} s_v^2(d_v) &= \frac{\int_0^\infty \min(x^2, d_v^2) e^{-(x-\mu_v)^2/2\sigma_v^2} dx}{\sqrt{2\pi}\sigma_v G_v} - m_v^2 \\ &= \sigma_v^2 \left[1 - \frac{Erfc(\tilde{d}_v)}{2G_v} - \gamma^2(\tilde{d}_v) + \sqrt{2}\tilde{d}_v \gamma(\tilde{d}_v) \right. \\ &\quad \left. - \frac{\tilde{d}_v + \tilde{\mu}_v}{\sqrt{\pi}G_v} e^{-\tilde{\mu}_v^2} \right] \end{aligned} \quad (51)$$

where $\tilde{d}_v = (d_v - \mu_v) / \sqrt{2}\sigma_v$, $\tilde{\mu}_v = \mu_v / \sqrt{2}\sigma_v$, and $\gamma(\tilde{d}_v) = [e^{-\tilde{\mu}_v^2} - e^{-\tilde{d}_v^2} + \sqrt{\pi}\tilde{d}_v Erfc(\tilde{d}_v)] / \sqrt{2\pi}G_v$.

The following inequalities provide important insights of the optimal solution.

Theorem 3: Let $F_v(x)$ be the CDF of the distribution specified by (48). Then

$$a) \quad \frac{\partial (m_v/d_v)}{\partial d_v} \leq 0 \quad (52)$$

$$b) \quad \frac{\partial m_v}{\partial \sigma_v} \leq 0, \quad \text{if } d_v \leq 2\mu_v \quad (53)$$

$$c) \quad \frac{\partial \bar{F}_v}{\partial \sigma_v} \geq 0, \quad \text{if } d_v \geq \left(1 + \frac{e^{-\mu_v^2/2\sigma_v^2}}{G_v} \right) \mu_v. \quad (54)$$

Equation (52) shows that the ratio of the mean carried traffic to the provisioned bandwidth decreases with the increase of provisioned bandwidth. We note that this trend of declining return from bandwidth provisioning is not specific to the case of the truncated Gaussian distribution.

Theorems 3-b) and 3-c) reveal an interesting implication of demand uncertainty, which is parameterized by the standard deviation σ_v . In (53), the condition $d_v \leq 2\mu_v$ is imposed to accommodate the effect of truncating the demand distribution, and can be ignored if the distribution is close to Gaussian. The inequality shows that increasing the standard deviation of a demand reduces the mean carried traffic, indicating that demand uncertainty is detrimental to revenue. Does that mean that less bandwidth should be provisioned to a demand if its uncertainty

increases? Not necessarily. It is implied by Theorem 3-c) that under certain conditions, it is optimal to allocate more bandwidth to a demand as its uncertainty increases. For ease of exposition, we elaborate below on this point for the special case of $\delta = 0$, while noting that it holds in general for $\delta > 0$.

The condition on d_v in (54) requires the bandwidth provisioned for the demand to be more than its mean, amplified slightly by the term $e^{-\mu_v^2/2\sigma_v^2}/G_v$ to accommodate the truncation of the Gaussian distribution. Suppose the requirement is satisfied, as is typical in networks that are not overloaded, so that $\partial\bar{F}_v/\partial\sigma_v \geq 0$. When $\delta = 0$, the necessary condition (46) becomes $\pi_v\bar{F}_v = \chi_v$. To maintain the above equality as σ_v increases, either d_v has to increase to cancel out the impact on $\pi_v\bar{F}_v$ by the increase of σ_v (notice that $\partial\bar{F}_v/\partial d_v \leq 0$), or χ_v has to increase. The increase in χ_v leads to an increase in the shadow costs of some links on route $r \in \mathcal{R}(v)$. The latter effectively makes routes of other demands that share these links with v more expensive. As a result, the amount of bandwidth provisioned to these demands will be reduced, and hence more bandwidth will be left for demand v .

To summarize, while a demand of higher uncertainty usually gives lower mean revenue, it may also have a higher mean *marginal* revenue (represented by the term $\pi_v\bar{F}_v$), which is the incremental revenue obtained from an additional amount of provisioned bandwidth. This upward revenue potential may lead to more bandwidth being dedicated to demands with higher uncertainty, despite the fact that the average revenue from these demands are lower.

The carrier's attitude to risk influences the provisioning of bandwidth to demands with higher uncertainty. Referring to the necessary condition (46)

$$\pi_v\bar{F}_v \left[1 - \delta \frac{\pi_v(d_v - m_v)}{S} \right] = \chi_v$$

in the above discussion of the special case of $\delta = 0$ we showed that when σ_v increases, d_v also needs to be increased to rebalance the equation by reducing its left-hand side. In case $\delta \neq 0$, raising d_v not only reduces \bar{F}_v but also increases $\pi_v(d_v - m_v)/S$. Therefore, the equation can be rebalanced with a smaller increase in d_v . This suggests that a carrier with lower risk tolerance provides less bandwidth for demands with higher uncertainty.

VI. NUMERICAL STUDIES

In this section, we discuss implications of demand uncertainty and the carrier's risk tolerance on traffic engineering and revenue management through numerical examples. We first describe the network topology and base case scenario in Section VI-A, and present results in Sections VI-B–E.

A. Network and Base Case

We consider a sample network which has 12 nodes and 14 bidirectional links. The network topology, as well as the labels of nodes and links are shown in Fig. 1. All links have 150 units of capacity, except links 6, 7, 12, 13, and 14, which have 200 units. The latter links are given higher capacities since they are likely

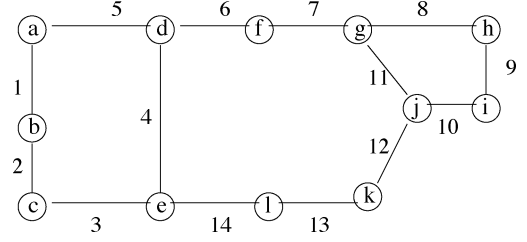


Fig. 1. Network topology.

to carry more traffic on account of their central locations. Between each node pair in the network there is some random demand which is symmetric in both directions. We make the following assumptions about these demands.

- 1) Let h_v be the minimum number of hops between a node pair v . The admissible route set of demand v , $\mathcal{R}(v)$, contains all paths that have no more than $(h_v + 2)$ links.
- 2) The price per unit of carried bandwidth for demand v , π_v , is proportional to h_v , i.e., $\pi_v = \kappa h_v$, where κ is set at 50 in the numerical studies. Hence, the price for a given demand is determined by the length of the shortest path and is independent of its selected routes.
- 3) The volume of each demand follows the Truncated Gaussian distribution, see (48). The ratio of the two defining parameters of the distribution, $CV_v = \sigma_v/\mu_v$, is kept below 0.35 for all v . Consequently, the effect of truncating the original Gaussian distribution is negligible and μ_v, σ_v and CV_v are close estimates of the mean, standard deviation and coefficient of variation, respectively.
- 4) To satisfy the GoS requirement, the amount of capacity provisioned to serve each demand is required to be no less than its mean, i.e., $\underline{d}_v = \mu_v$.
- 5) All demands v are assumed to have the same mean, i.e., $\mu_v = \bar{\mu}$ for all v . Then $\sum_{v \in \mathcal{V}} \mu_v h_v = \bar{\mu} \sum_{v \in \mathcal{V}} h_v$ is an estimate of average network demand. The ratio of this quantity to the sum of all link capacities, $\rho = (\bar{\mu} \sum_{v \in \mathcal{V}} h_v) / \sum_{l \in \mathcal{L}} C_l$, is defined to be the network *load factor*. In numerical examples we vary the load factor ρ by appropriately setting the mean demand. In the sample network (Fig. 1), $\sum_{l \in \mathcal{L}} C_l = 2350$ and $\sum_{v \in \mathcal{V}} h_v = 176$, so $\mu_v = \bar{\mu} = 13.35\rho$. For given μ_v , to vary the coefficient of variation, CV_v , we appropriately set σ_v .

In addition to the *random* demands specified above, the network has another set of *guaranteed* demands between all network node pairs. The latter demands have no uncertainty and constitute a special case of random demands, as indicated in (6). While we assume here that the guaranteed demand is abundant, the unit price of carrying it is a (small) fixed fraction, ϕ ($0 < \phi < 1$), of the unit price for carrying random demand for the same node pair. This fraction is uniform across all node pairs and its value is set at either 0.1 or 0.2. Route selection of guaranteed demands is subject to the same constraints that apply to random demands. However, from the above assumptions on guaranteed demands, specifically, prices are proportional to minimum distances (measured by hops) and volumes are unlimited, it should be clear that it is optimal to carry such demands only on minimum-hop routes.

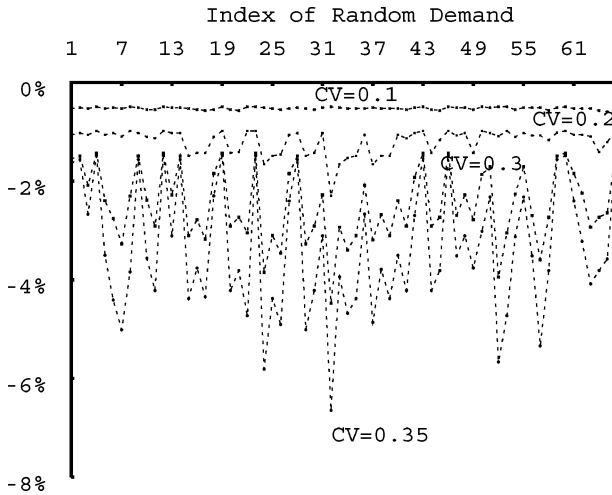


Fig. 2. Mean revenue from random demand drops with variability.

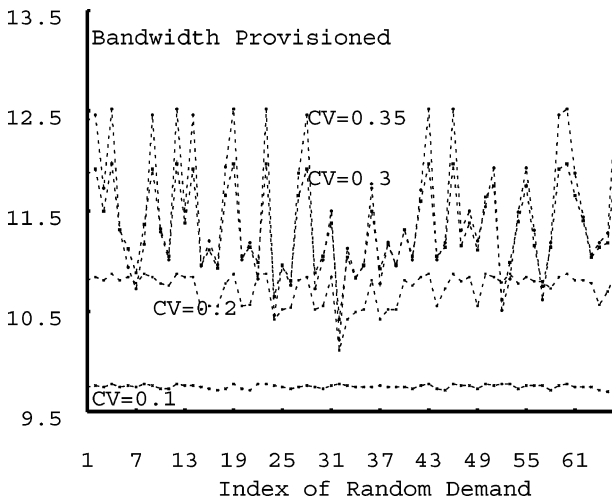


Fig. 3. Bandwidth provisioned to random demand increases with variability.

Observe that the carrier decides on bandwidth allocation to random and guaranteed demands, where the former has potentially higher but uncertain revenue, e.g., from retail demands, while the latter has lower but secure revenue, e.g., from wholesale contracts. The optimal allocation is influenced by both demand uncertainty and the carrier's risk-averseness.

B. Impact of Demand Uncertainty

We start with the base case and set the risk parameter $\delta = 0.5$, and the ratio of unit prices for guaranteed and random demands for the same node pair, $\phi = 0.1$. Let the coefficient of variation (CV_v) increase from the base case value of 0.1 to 0.2, 0.3, and 0.35, and keep all other parameters unchanged. Figs. 2 and 3 demonstrate vividly the revenue implications of demand uncertainty, which are consistent with Theorem 3 in the last section. In both figures, the horizontal axis corresponds to the index of random demand, which is associated with a node pair. Given various uncertainty levels, Fig. 2 plots the percentage difference of mean revenue from the benchmark, which is the mean revenue obtained at zero uncertainty ($CV_v = 0$). The mean revenue for each random demand is an implication of the optimum network-wide bandwidth provisioning computed by the proce-

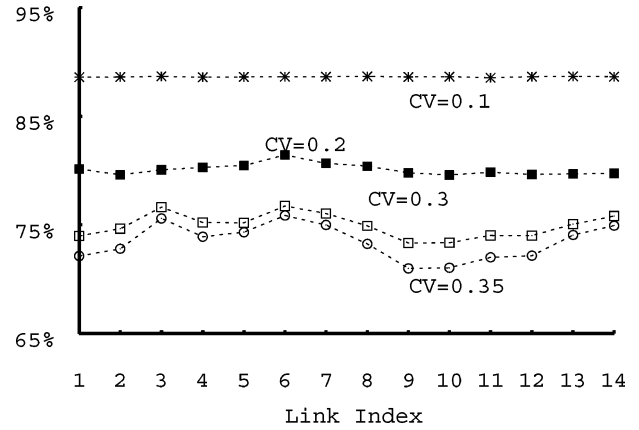


Fig. 4. Link utilization decreases with variability.

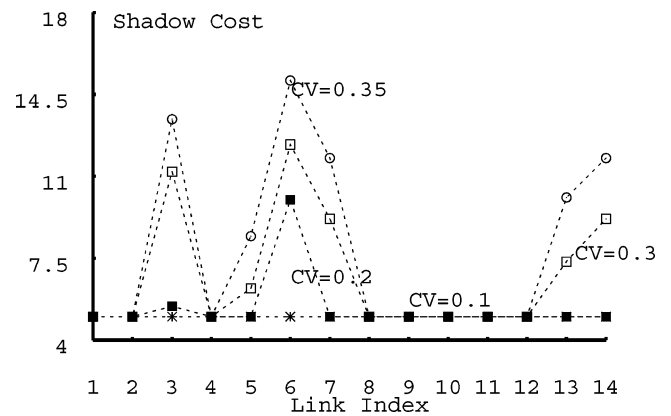


Fig. 5. Shadow cost increases with variability.

cedure discussed in the preceding sections. The figure reveals the detrimental impact of demand uncertainty as the differences are always negative and higher values of CV_v always lead to larger percentage drops. The amounts of provisioned bandwidth to these demands are shown in Fig. 3, which are increasing with demand uncertainty, indicating that more network capacities are diverted from serving guaranteed demands to serving random demands. This is exactly what Theorem 3 postulated, i.e., despite the fact that demands of higher uncertainty result in lower mean revenues, more bandwidth is provisioned to random demands.

For random demands, as uncertainty increases, the utilization of provisioned bandwidth decreases. Define link usage rate as the ratio of mean carried traffic to the amount of provisioned bandwidth on a link, both calculated after aggregation over all random demands using the link. Fig. 4 plots the usage rate for every link in the network at several uncertainty values. Despite the decreased efficiency of bandwidth usage, the marginal value of link capacity, characterized by the link shadow cost, increases with uncertainty, as shown in Fig. 5. This fact, surprising at first, is due to the need to provision more bandwidth to random demands when uncertainty is higher, which makes capacity a more constrained resource.

C. Risk Mitigation by Bandwidth Allocation

We highlight a risk mitigation mechanism derived from adjusting bandwidth allocations to the guaranteed and random de-

TABLE I
EFFECT OF RISK AVERSENESS ON ADJUSTED BANDWIDTH ALLOCATION AND MEAN REVENUE COMPOSITION

δ	bandwidth allocation		mean revenue composition	
	random	guarant.	random	guarant.
0.0	49.0%	51.0%	62.5%	37.5%
0.4	46.5%	53.5%	60.6%	39.4%
0.8	43.8%	56.2%	58.6%	41.4%
1.2	40.8%	59.2%	56.2%	43.8%
1.6	37.2%	62.8%	53.2%	46.8%
2.0	35.0%	65.0%	51.2%	48.8%
2.4	32.8%	67.2%	49.1%	50.9%

mands. The former has no uncertainty, its associated revenue is risk-free, albeit lower, so allocating more bandwidth to it mitigates the revenue risk. We vary the risk-averseness parameter, δ , and investigate the implied allocation of bandwidth to the two types of demands in the optimal solution. The results show that as the carrier's attitude varies from risk-neutral to highly risk-averse, the proportion of bandwidth allocated to guaranteed demands increases, as does the proportion of revenue from guaranteed demands in the total expected revenue.

Table I gives the results for δ ranging from 0 to 2.4. For random demands, we set $\rho = 0.65$ and $CV_v = 0.35$ for all v . The ratio of prices for carrying guaranteed and random demands for the same node pair, ϕ is fixed at 0.2. In all cases, the parameter $\underline{d}_v = \mu_v$ for all v in the GoS constraint, see (11). This implies that 1531 units of capacity, out of the total of 2350, is set aside for random services to meet the GoS constraint. The carrier allocates the remaining 819 units of bandwidth to random and guaranteed demands. Table I gives *adjusted* numerical values of bandwidth and expected revenue corresponding to the excess over the respective fixed quantities committed to satisfy the GoS constraint.

Table I shows that a risk-neutral carrier ($\delta = 0$) should provision 49% of uncommitted capacity to serve random demands. As the carrier becomes more risk-averse, the percentage of bandwidth allocated to random demands declines. The last row of the table shows that a conservative carrier ($\delta = 2.4$) should allocate no more than 32.8% of the uncommitted capacity to random demands. Differences in provisioning lead to differences in revenue composition. When $\delta = 0$, 62.5% of the revenue is generated from random demands. When $\delta = 2.4$, the corresponding number is 49.1%.

D. Impact of Risk-Averseness on Route Selection

Risk-averseness not only affects bandwidth allocation, as we demonstrated above, but also route selection, as we will illustrate next. We associate with any routing solution its *distance-bandwidth product* (DBP)

$$DBP = \sum_{v \in \mathcal{V}} \sum_{r \in \mathcal{R}(v)} \xi_r |r|$$

where ξ_r is the amount of bandwidth provisioned on route $r \in \mathcal{R}(v)$ as defined in Section II and $|r|$ is the number of hops in route r . A reference value is

$$DBP_0 = \sum_{v \in \mathcal{V}} h_v \sum_{r \in \mathcal{R}(v)} \xi_r$$

Deviation of DBP from Min-hop Routing

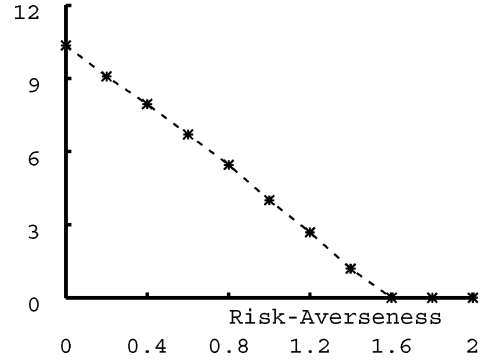


Fig. 6. Change of route selection with risk-averseness.

TABLE II
EFFECT OF RISK AVERSENESS ON ROUTE SELECTION

δ	min-hop	non min-hop	
	rte. 1	rte. 2	rte. 3
0.0	67.8%	5.4%	26.8%
0.4	76.1%	7.2%	16.7%
0.8	86.4%	6.4%	7.2%
1.2	93.0%	3.5%	3.5%
1.6	100.0%	0.0%	0.0%

which is the distance-bandwidth product when bandwidth is provisioned only on routes with the minimum number of hops. We let $\Delta(DBP) = (DBP - DBP_0)$ measure the departure of the routing solution from minimum-distance routing. In Fig. 6, we show the dependence of $\Delta(DBP)$ on the risk-averseness parameter δ . The noteworthy feature of the figure is the strictly monotonic decrease of $\Delta(DBP)$ with increasing δ until it reaches 0, where the optimal solution coincides with the minimum-distance routing. Note that $\Delta(DBP)$ sticks to 0 for higher values of δ .

Table II provides additional insight to corroborate the results in Fig. 6. We consider routing of demand from node h to node e in Fig. 1 on three candidate routes, route 1 ($\{h-g-f-d-e\}$), which has the minimum distance, routes 2 ($\{h-g-j-k-l-e\}$) and 3 ($\{h-i-j-k-l-e\}$), which have an additional hop. Table II shows that with risk-neutrality ($\delta = 0$), roughly one-third of the demand is carried over the two nonminimum-distance routes. For $\delta \geq 1.6$, the routing is exclusively minimum distance.

E. Efficient Frontier

The effect of the carrier's tolerance to risk is summarized by the *efficient frontier*, which is obtained by solving the optimization model for various values of δ . In Fig. 7 we show the frontier for the base case network. Each point on the curve gives the optimum combination of mean and standard deviation of revenue at a given value of δ , and represents the maximum expected revenue (reward) obtainable at a given level of risk, or the minimum risk the carrier has to take for a given reward.

A key feature to observe in the displayed efficient frontier in Fig. 7 is the "knee" of the curve, which corresponds approximately to mean revenue of 11 000. To the left of this point, the tradeoff between risk and mean revenue is roughly linear and benign, and to the right the curve is increasingly nonlinear and

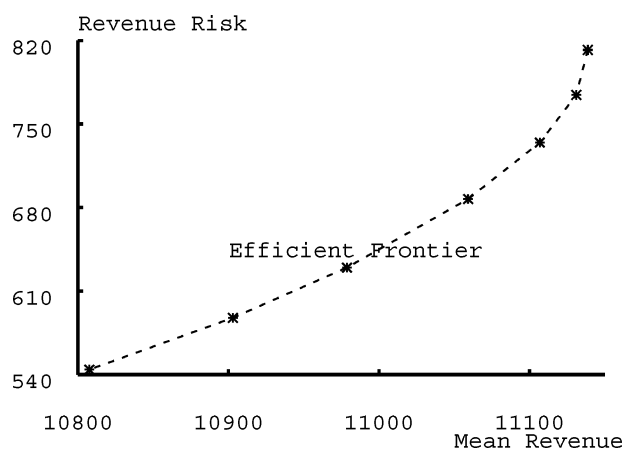


Fig. 7. Efficient frontier.

any small gain in the expected revenue is accompanied by a substantial cost in increased risk. Moreover, the figure highlights the importance of precision in locating the knee since a relatively small variation of about 1%–2% in mean revenue around the knee can lead to substantial changes in the slope of the efficient frontier.

VII. CONCLUSION

We have presented and analyzed a stochastic traffic engineering framework for off-line planning of bandwidth provisioning and routing. The framework is based on an optimization model that uses probability distributions of demands as inputs and maximizes the weighted combination of the mean revenue and the risk of revenue shortfall. We discuss properties of the objective function, and strategies for solving the model as a concave maximization problem. In our numerical studies, we analyze the impacts of demand uncertainty on various aspects of traffic engineering design. We observe significant changes in mean revenue, bandwidth provisioning, link shadow costs and utilization with demand uncertainty. We demonstrate that changes in bandwidth provisioning with uncertainty are strongly influenced by the carrier's tolerance to risk, and give the efficient frontier, which characterizes the tradeoff between expected revenue and revenue risk.

Our analysis can be extended in several directions. For instance, an extension of the model can give carriers an optimization tool for use in interactions with peers in bandwidth sharing and network interconnections. These extensions are currently being investigated, and will be reported in future publications.

ACKNOWLEDGMENT

The authors acknowledge K. G. Ramakrishnan for his contributions to the early development of this work, and A. Weiss for the benefit of discussions.

REFERENCES

- [1] R. K. Ahuja, T. L. Magnanti, and J. B. Orlin, *Network Flows*. Upper Saddle River, NJ: Prentice-Hall, 1993.
- [2] D. Applegate and M. Thorup, "Load optimal MPLS routing with N+M labels," presented at the IEEE INFOCOM, San Francisco, CA, Apr. 2003.
- [3] P. Aukia, M. Kodialam, P. V. N. Koppol, T. V. Lakshman, H. Sarin, and B. Suter, "RATES: a server for MPLS traffic engineering," *IEEE Network*, vol. 14, no. 2, pp. 34–41, Mar./Apr. 2000.
- [4] D. O. Awduche, A. Chiu, A. Elwalid, I. Widjaja, and X. Xiao, "Overview and principles of Internet traffic engineering," IETF, RFC 3272, May 2002.
- [5] J. Cao, D. Davis, S. V. Wiel, and B. Yu, "Time-varying network tomography: router link data," *J. Amer. Statistical Assoc.*, vol. 95, pp. 1063–1075, 2000.
- [6] J. Cao, S. V. Wiel, B. Yu, and Z. Zhu, "A scalable method for estimating network traffic matrices," Bell Labs Tech. Memo., 2000.
- [7] A. Elwalid, C. Jin, S. Low, and I. Widjaja, "MATE: MPLS traffic engineering," in *Proc. IEEE INFOCOM*, 2001, pp. 1300–1309.
- [8] A. Elwalid, D. Mitra, I. Sanjeev, and I. Widjaja, "Routing and protection in GMPLS networks: from shortest paths to optimized designs," *J. Lightw. Technol.*, vol. 21, no. 11, pp. 2828–2838, Nov. 2003.
- [9] A. Feldmann, A. Greenberg, C. Lund, N. Reingold, and J. Rexford, "NetScope: traffic engineering for IP networks," *IEEE Network*, vol. 14, no. 2, pp. 11–19, Mar./Apr. 2001.
- [10] B. Fortz and M. Thorup, "Internet traffic engineering by optimizing OSPF weights," in *Proc. IEEE INFOCOM*, 2000, pp. 519–528.
- [11] M. Grossglauser and J. Rexford, "Passive traffic measurement for IP operations," presented at the INFORMS Telecom Meeting, Ft. Lauderdale, FL, Mar. 2002.
- [12] J. Hirshleifer and J. G. Riley, *The Analytics of Uncertainty and Information*. Cambridge, U.K.: Cambridge Univ. Press, 1992.
- [13] F. Kelly, A. Maulloo, and D. Tan, "Rate control in communication networks: shadow prices, proportional fairness and stability," *J. Oper. Res. Soc.*, vol. 49, pp. 237–252, 1998.
- [14] C. Lagoa and H. Che, "Decentralized optimal traffic engineering in the Internet," *Comput. Commun. Rev.*, vol. 30, no. 5, pp. 39–47, Oct. 2000.
- [15] H. Levy, *Stochastic Dominance: Investment Decision Making Under Uncertainty*. Boston, MA: Kluwer, 1998.
- [16] Level 3. (2001, Jun.) Level 3 Signs New Agreement with Microsoft. [Online]. Available: <http://www.level3.com/press/2053.html>
- [17] Level 3. (3)Crossroads Wholesale Internet Access. [Online]. Available: <http://www.level3.com/560.html>
- [18] S. G. Lanning, D. Mitra, Q. Wang, and M. H. Wright, "Optimal planning for optical transport networks," *Philos. Trans. Royal Soc.*, pp. 2183–2196, 2000.
- [19] D. Luenberger, *Introduction to Linear and Nonlinear Programming*. Reading, MA: Addison-Wesley, 1973.
- [20] H. Markowitz, *Mean Variance Analysis in Portfolio Choice and Capital Markets*. New York: Basil Blackwell, 1987.
- [21] R. Morris and D. Lin, "Variance of aggregated web traffic," presented at the IEEE INFOCOM, Tel Aviv, Israel, Apr. 2000.
- [22] A. Medina, N. Taft, K. Salamatiyan, S. Bhattacharyya, and C. Diot, "Traffic matrices estimation: existing techniques and new directions," presented at the ACM SIGCOMM, Pittsburgh, PA, Aug. 2002.
- [23] D. Mitra and K. G. Ramakrishnan, "A case study of multiservice multipriority traffic engineering design for data networks," in *Proc. IEEE GLOBECOM*, Dec. 1999, pp. 1077–1083.
- [24] D. Mitra, K. G. Ramakrishnan, and Q. Wang, "Combined economic modeling and traffic engineering: joint optimization of pricing and routing in multiservice networks," in *Traffic Engineering in the Internet Era: Proc. 17th Int. Teletraffic Congr. (ITC-17)*, J. M. de Souza, N. L. S. da Fonseca, and E. A. de Souza e Silva, Eds. Amsterdam, The Netherlands: Elsevier, 2001, pp. 73–85.
- [25] D. Mitra and Q. Wang, "Stochastic traffic engineering, with applications to network revenue management," presented at the IEEE INFOCOM, San Francisco, CA, Apr. 2003.
- [26] —, "Risk-aware network profit management in a two-tier market," presented at the 18th Int. Teletraffic Congr. (ITC-18), Berlin, Germany, Aug.–Sep. 2003.
- [27] W. Ogryczak and A. Ruszczyński, "Dual stochastic dominance and related mean-risk models," *SIAM J. Optim.*, vol. 13, pp. 60–78, 2002.
- [28] N. Semret, R. R.-F. Liao, A. T. Campell, and A. A. Lazar, "Peering and provisioning of differentiated Internet services," in *Proc. IEEE INFOCOM*, Tel Aviv, Israel, Apr. 2000, pp. 100–107.
- [29] S. Suri, M. Waldvogel, and P. R. Warkhede, "Profile-based routing: a new framework for MPLS traffic engineering," presented at the Quality of Future Internet Services Conf. (QofIS'01), Coimbra, Portugal, 2001.
- [30] P. Trimintzios, I. Andrikopoulos, G. Pavlou, P. Flegkas, D. Griffin, P. Georgatsos, D. Goderis, Y. T' Joens, L. Georgiadis, C. Jacquenet, and R. Egan, "A management and control architecture for providing IP differentiated services in MPLS-based networks," *IEEE Commun. Mag.*, vol. 39, no. 5, pp. 80–87, May 2001.

- [31] Y. Vardi, "Network tomography: estimating source-destination traffic intensities from link data," *J. Amer. Statist. Assoc. Theory Meth.*, vol. 91, pp. 365–377, 1996.
- [32] J. von Neuman and O. Morgenstern, *Theory of Games and Economic Behavior*. Princeton, NJ: Princeton Univ. Press, 1953.
- [33] W. Wang, M. Palaniswami, and S. H. Low, "Optimal flow control and routing in multipath networks," *Perform. Eval.*, vol. 52, no. 2–3, pp. 119–132, Apr. 2003.



Debasis Mitra (F'89) received the Ph.D. degree in electrical engineering from London University, London, U.K.

He joined Bell Laboratories as a Member of Technical Staff in 1968. During the fall semester of 1984, he was Visiting McKay Professor at the University of California, Berkeley. As Vice President of Mathematical Sciences Research at Bell Labs, Murray Hill, NJ, he directs research in fundamental mathematics, mathematics of networks and systems, statistics and data mining, information and communications

theory, and industrial mathematics. He serves on the Air Force Science and Technology Board of the National Academies. In 2003 he served as the Chair of the Telecom review panel of the New Jersey Commission on Jobs Growth and Economic Development and as the Albert Winsemius Professor at the Nanyang Technical University in Singapore. His personal research interests are currently in optical networking, IP/optical convergence, stochastic traffic engineering, network economics and network revenue management.

Dr. Mitra is a member of the National Academy of Engineering and a Bell Labs Fellow. He is a co-recipient of the 1998 IEEE Eric E. Sumner Award with the citation, "For the conception and development of voice echo cancelers." He is the recipient of the 1993 Steven O. Rice Prize Paper Award and the 1982 Guillemin–Cauer Prize Paper Award of the IEEE. He is also the recipient of awards from the 1995 ACM Sigmetrics/Performance Conference, the Institution of Electrical Engineers (U.K.) and the *Bell Systems Technical Journal*. He has been a member of the editorial boards of the IEEE/ACM TRANSACTIONS ON NETWORKING, the IEEE TRANSACTIONS ON COMMUNICATIONS, the IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS, and *Queueing Systems (QUESTA)*. He is currently the Area Editor of *Operations Research* for Telecommunications and Networking.



Qiong Wang (M'98) received the Ph.D. degree in engineering and public policy from Carnegie-Mellon University, Pittsburgh, PA, in 1998.

He is a Member of Technical Staff at Bell Labs, Murray Hill, NJ, where he does research on network economics and operations management in the telecommunications industry.