



# Efficient elastic bulky traffic transfer with a new pricing scheme based on number of flows



Zhangxiao Feng, Weiqiang Sun\*, Fengqin Li, Weisheng Hu

The State Key Laboratory of Advanced Optical Communication Systems and Networks, Shanghai Jiao Tong University, Shanghai, China

## ARTICLE INFO

### Article history:

Received 9 July 2014

Revised 21 May 2015

Accepted 12 June 2015

Available online 18 June 2015

### Keywords:

Bulk data transfer

Congestion pricing

Utility maximization

## ABSTRACT

Emerging data center and cloud services carry large blocks of data, which leads to intensive bandwidth competition on today's Internet, particularly during peak hours. Congestion pricing is considered to be an effective means to avoid peak-hour competition. Typical congestion pricing schemes proposed in the past require significant enhancements to existing network infrastructures. For instance, complex congestion information collection and price announcement mechanisms must be in place, hindering the deployments of many congestion pricing schemes. In this study, we present one simple yet effective congestion pricing scheme based on the number of flows in single-bottleneck networks carrying elastic traffic. During congested hours, data transmitters share bandwidth on the bottleneck and are also charged according to the number of flows they use. Under this pricing scheme, network capacity is used efficiently, and all users can achieve maximum utility with increasing and strictly concave utility functions. We investigated how a single congestion control parameter, the price index, affects the degree of network congestion. Through network simulation, we proved that under the proposed scheme, the big data transmitters sharply reduce their transfer during congested hours.

© 2015 Elsevier B.V. All rights reserved.

## 1. Introduction

In recent years, both the total number of network users and their access speeds have increased dramatically, injecting an excessive amount of traffic into the network. Emerging cloud and data center services intensify this bandwidth competition. These trends have brought about the “age of big data”. For instance, in applications such as virtual machine migrations and data back-ups/replications, a single flow may consume a significant portion of link bandwidth for a long period of time.

Recently, it has been reported that Internet data traffic is growing by 50–60% every year in North American [1] and that the global volume of digital data will grow to 40,000 EB by 2020 [2]. According to Cisco's White Paper [3], the annual global data center IP traffic will reach 7.7 ZB by the end of 2017. Greenberg et al. [4] reported that the networking cost of running a cloud service data center is approximately 15% of its total cost and is almost as same as the power draw. They believe that networking is the key to reducing the overall costs. Although global data volumes are expanding rapidly, researchers have noted pessimistically that a worldwide fiber infrastructure will not be upgraded in the next decade [5]. Accordingly, how

to make full use of the existing network bandwidth will be an important issue in addressing the growing volume of data.

In a traditional commercial network, internet service providers (ISPs) charge customers under either a usage-based model or a flat fee model. With such pricing strategies, peak-hour congestion is unavoidable, while large amounts of off-peak-hour bandwidth are left unused. The traffic variation within an ISP's network is often regular and predictable and complies with the diurnal variation pattern [6]. In the science domain, large data sets are generated by experimental facilities, such as the large hadron collider (LHC) [7], the large synoptic survey telescope (LSST) [8], and the laser interferometer gravitational wave observatory (LIGO) [9]. Scientific collaborations require the delivery of such data sets to worldwide research institutes. Unlike real-time communications, bulk data delivery is often delay tolerant, e.g., on the order of hours or days. For example, the genomic data produced by DNA sequencers is often too large to be transmitted over the Internet and thus is stored on hard disks and sent via FedEx [10]. By postponing the transmission of delay-tolerant requests to off-peak hours, network congestion can be alleviated, and the overall quantity of the user experience can be improved.

The benefits of taking advantage of the delay tolerance in large data flows have already been demonstrated. Laoutaris et al. [11] proposed taking advantage of the off-peak bandwidth to transmit delay-tolerant bulk (DTB) data within a network deployed with a percentile-based pricing scheme. Bulk data are overwhelming in

\* Corresponding author. Tel.: +86 21 34205359; fax: +86 21 34204597.  
E-mail address: [sunwq@sjtu.edu.cn](mailto:sunwq@sjtu.edu.cn) (W. Sun).

today's networks; a recent report [12] showed that within a typical data center network, 90% of the network traffic is generated by only 1% of the flows whose sizes are larger than 100 MB. The delay tolerance of bulk data may be the key to designing an effective way to transmit data. In [11], researchers designed two transmission policies: an end-to-end scheduling transmission with senders and receivers from the same time zone, and a store-and-forward transmission with an assisting storage node that addresses the problem of non-coincident off-peak valleys. They also designed a system for rescuing unutilized bandwidth across several time zones using a store-and-forward method [13]. Their results showed that it is possible to utilize off-peak capacity to transfer terabyte-sized data.

We propose a new pricing scheme under which customers carrying DTB data are encouraged to avoid peak-hour transmission and fully utilize off-peak capacity. Time-shifting DTB data traffic provides several benefits: (1) eliminating or relieving the network congestion during peak hours, (2) improving the quality of the user experience for light and interactive applications, and (3) reducing both customers' transmission costs and ISPs' network construction expenses.

Recent studies have shown that in single-bottleneck networks with adaptive queue management, the queue length is proportional to the number of active transmission control protocol (TCP) flows [14,15]. Thus, it is natural for ISPs to control the customers' concurrent number of TCP flows in a commercial network as a means of congestion control. According to the bandwidth sharing mechanism in the current Internet, one application that uses more TCP flows is able to obtain more bandwidth in a single-bottleneck network. To alleviate network congestion, incentives must be provided to applications to either reduce the numbers of flows they use on congested networks or avoid peak traffic hours. We have designed a simple but effective pricing scheme based on the number of flows to address this problem. Under this scheme, customers would be charged according to the numbers of flows they use. In this study, we present how a pricing scheme based on the number of flows used may change the way customers utilize bandwidth and how congestion of networks with single bottlenecks might be reduced.

The proposed scheme operates in three stages. First, the ISP announces a price per unit of data that is proportional to each customer's number of flows. Second, customers need to decide how many flows to use, based on the announced unit price and their delay tolerances. Third, the ISP charges each customer based on the customer's number of flows and the volume of transmitted data. By deploying this number-of-flows-based scheme, a stable network-wide price index can be determined by the operation support system (OSS), and microscopic price adjustment can be performed in a distributed manner on the user side. This approach is also an application layer approach that requires no modification to the TCP protocol.

The key contributions of our work are the following:

- We have designed a number-of-flows-based pricing scheme with little central control overhead.
- We prove that under the proposed pricing scheme, all customers' demands are satisfied, and their maximum overall utilities are achieved.
- We show that DTB data transmitters reduce their transfer rates sharply during peak hours, while other transmitters can achieve relatively high transfer rates.
- We demonstrate that the congestion level can be adjusted by changing a single parameter  $\bar{I}$  the price index.

The remainder of this paper is organized as follows. We discuss related work in Section 2. In Section 3, we introduce the system model and list the assumptions that apply throughout this paper. In Section 4, we introduce the design of the pricing scheme, and the details of the scheme and its properties are presented in Section 5. Simulations results are presented and discussed in Section 6. In Section 7,

we discuss the time slot and pricing interval. We present conclusions in Section 8.

## 2. Related work

Many interesting studies on congestion pricing have been carried out since 1990s. In the "Smart Market" approach proposed by MacKie-Mason and Varian in 1995 [16], dynamic prices are set according to price-sensitive and time-sensitive users' responses to the degree of network congestion. The price for a packet is determined by users' "bids" that reflect their willingness to pay, and every packet that is admitted into the network is charged a small market-clearing price that is equal to the lowest "bids" among all of the admitted packets. In 1997, Kelly proposed a method for user utility-based network resource allocation in which users are able to choose the price charged per unit of time and the network allocates rates per unit of charge that are proportionally fair [17,18]. This achieves an equilibrium between the network's resource allocation choices and users' willingness to pay. In 2001, Ganesh et al. [19] introduced a congestion pricing scheme that charges users according to the collective actions of all the users of an elastic traffic network. By applying a rate adaptation mechanism, the network converges to an optimal allocation of bandwidth in terms of users' price predictions and their bandwidth shares. In 2008, Jiang et al. [20] used game theory models to study how users' time-preferences influence network and service providers. In their model, users chose their access times based on their own time-preferences, price and the degree of network congestion. The readers are referred to [21] for a comprehensive review on earlier pricing schemes.

More recently, studies on new congestion pricing schemes are stimulated by the constantly increasing traffic demand and the unfair bandwidth sharing among mice and elephant applications. In 2007, Brisco et al. proposed a congestion pricing scheme based on bandwidth sharing [22]. They suggested that users be charged according to their congestion volume, i.e., the number of packets marked in the network. They also proposed that users choose their numbers of TCP flows to obtain different bandwidth shares. In 2011, Joe-Wong et al. proposed a time-dependent pricing (TDP) scheme that would incentivize users to defer their application sessions to lower-volume periods using time-varying monetary rewards [23]. Under this TDP approach, an ISP achieves a balance between its expenditures on capacity expansion and the costs of offering "rewards" to users who are willing to defer some of their sessions to later times. The authors also presented the deployment and implementation of an end-to-end TDP system called TUBE, and through a trial with real users, the authors showed that users actually did defer their usage to off-peak times and increase their monthly usage [23]. In 2012, Vemu et al. [24] proposed two adaptive multilayered feedback pricing schemes. Prices are announced on the basis of the instantaneous queue length in one of these two schemes and on the basis of the weighted average queue length in the other scheme. Both schemes achieve optimal price levels using the simultaneous perturbation stochastic approximation (SPSA) methodology.

The typical congestion pricing schemes discussed above are summarized in Table 1.

## 3. Network model and assumptions

In this study, we present the pricing problem for a network with a single bottleneck. A group of users with multiple flows share a single bottleneck within an ISP's network, as shown in Fig. 1. It is assumed that each user has unlimited access bandwidth, i.e., that there is no traffic limiting or other bottlenecks in the access or other parts of the network. There is also a business and operations support systems (OSS/BSS) in the network that maintains a billing and accounting system. The OSS/BSS is also responsible for announcing a price index

Download English Version:

<https://daneshyari.com/en/article/447691>

Download Persian Version:

<https://daneshyari.com/article/447691>

[Daneshyari.com](https://daneshyari.com)