# Fair bandwidth sharing and delay differentiation: Joint packet scheduling with buffer management

Xiaobo Zhou [a,*], Dennis Ippoliti [a], Liqiang Zhang [b]

[a] Department of Computer Science, University of Colorado at Colorado Springs, CO 80918, USA
[b] Department of Computer and Information Sciences, Indiana University South Bend, IN 46615, USA

ABSTRACT

Packet delay and bandwidth are two important metrics for measuring quality of service (QoS) of Internet services. Traditionally, packet delay differentiation and fair bandwidth sharing are studied separately. In this paper, we first propose a generalized model for providing fair bandwidth sharing with delay differentiation, namely FBS-DD, at the same time. It essentially aims to provide multi-dimensional proportional differentiation with respect to both QoS metrics. We design size-based packet scheduling schemes that take both packet delay and packet size into scheduling considerations, without assuming admission control or policing. Furthermore, we propose a PID control-theoretic buffer management scheme. The packet scheduling with buffer management approach provides delay and bandwidth differentiation in an integrated way, while existing approaches consider delay and loss rate differentiation as orthogonal issues. It enhances the flexibility of network resource management and multi-dimensional QoS provisioning. It is capable of self-adapting to varying workloads from different classes, which automatically builds a firewall around aggressive clients and hence protects network resources from saturation. Extensive simulation results by the use of trace files demonstrate that the packet scheduling schemes can provide predictable fair bandwidth sharing with delay differentiation at various situations. The control-theoretic buffer management scheme further improves the controllability.

© 2008 Published by Elsevier B.V.

## 1. Introduction

Differentiated Services (DiffServ) is one of the major efforts to meet the demand of provisioning different levels of quality of service (QoS) on the Internet so as to support different types of network applications and various user requirements. It aims to provide differentiated services between classes of aggregated traffic flows within a router, rather than offer QoS guarantees to individual flows [1]. To receive different levels of QoS, packets are assigned with different service types or traffic classes at the network edges. DiffServ-compatible routers in the network core perform stateless prioritized packet forwarding or dropping, called "per-hop behaviors" (PHBs), to the classified packets. Due to its per-class stateless processing, the DiffServ architecture exhibits good scalability. Its provisioning is an active research topic [5,7,9,10,17,21–25].

There are two basic schemes to DiffServ provisioning. Absolute DiffServ aims to provide statistical assurances for a class's received performance measures, such as a minimum service rate or maximum delay. Relative DiffServ is to quantify the quality spacings between different classes. The proportional differentiation model,

proposed by Dovrolis, *et al.* [4], is a popular relative DiffServ model. It aims to provide per-class QoS level in proportion to the pre-specified differentiation parameters of the classes, independent of those class workloads. Delay and bandwidth are two important QoS metrics considered in the model. The algorithms for proportional delay differentiation (PDD) consider lossless and work-conserving packet scheduling [4–6,11,13,14,17,18]. When the overall workload of classes is close to or exceeds the link bandwidth capacity, the algorithms for proportional bandwidth differentiation aim to enforce that the ratio of the loss rates of two classes be proportional to the ratios of their differentiation parameters [3,7,21,23]. However, most of those algorithms consider delay differentiation and bandwidth differentiation as orthogonal issues.

While the PDD model is excellent due to its delay proportionality fairness to clients, it is insufficient and might be unfair from the perspective of the network resource providers. It is because the model does not consider another important issue, fair bandwidth sharing. Fair bandwidth sharing is a classic issue. Its short-term behaviors were originally studied as fair queueing [2]. While those PDD algorithms can ensure that experienced delay of different classes be proportional, there is no assumption nor guarantee on the fair bandwidth sharing, be in short term or in long term. Consider two traffic classes (Class-1 and Class-2) with the pre-specified differentiation parameters 2 and 1, respectively. Consider the

* Corresponding author. Tel.: +1 719 262 3493.
  *E-mail addresses:* zbo@cs.uccs.edu (X. Zhou), liqzhang@iusb.edu (L. Zhang).

scenario that Class-1's workload is 80% of the link capacity and Class-2's workload is 5% of the link capacity. According to the proportional delay differentiation model, the ratio of the average packet delay of Class-1 to that of Class-2 would be 1 to 2. However, the workload of Class-1 is 16 times of that of Class-2 while their differentiation parameter ratio is only 2 to 1. The scenario illustrates that the current workload-independent proportional differentiation model can be very unfair to some network traffic. Even worse, some aggressive or malicious clients can utilize this unfairness and weak controllability to attack the network resources.

Note that we do not intend to deny the merit of the PDD model. Essentially, it considers the single-dimensional QoS provisioning with respect to delay. It needs the support of admission control schemes that shape the traffic according to the service level agreements or some adaptive schemes that promote the differentiation parameters dynamically according to the workload conditions. Generally, the pre-specified differentiation parameters are used by the network operators to control the quality spacings between the multiple classes. They are often associated to the differentiated pricing, say proportionally. But the model is insufficient when multiple QoS metrics exist and multi-dimensional QoS provisioning should be considered.

Given that both bandwidth and delay are important metrics for measuring QoS of Internet services, we need to consider fair bandwidth sharing and delay differentiation at the same time. The primary contributions of our work are:

1. We propose a generalized model, FBS-DD, for providing fair bandwidth sharing with delay differentiation at the same time. It is to ensure that the ratio of the average delay of two classes normalized by their achieved bandwidth ratio be proportional to the pre-specified differentiation parameters. It essentially aims to provide multi-dimensional proportional differentiation with respect to both QoS metrics, packet delay and bandwidth. One uniqueness is that the delay differentiation and loss rate differentiation are integrated with traffic policing capabilities for providing better controllability to network operators and fairness to clients.
2. We design size-based packet scheduling algorithms for FBS-DD provisioning, modified from the waiting-time priority (WTP) algorithms which are excellent schedulers for performing proportional delay differentiation. Two VPS (various packet size) algorithms take both packet size and packet delay into consideration in packet scheduling. For packets with the uniform size, the VPS schemes are reduced to UPS (uniform packet size) schemes.
3. We further study the performance controllability with control-based buffer management. When the overall workload of the classes is below the link capacity, the FBS-DD model actually is to achieve the proportional delay differentiation weighted by the workloads of the classes in the long term. When the overall workload of classes is beyond the link capacity so that there will be packet loss, the FBS-DD model is to achieve the proportional delay differentiation weighted by the experienced bandwidth ratio of the classes. This is however a non-trivial issue. We propose a PID control-theoretic buffer management scheme to further provide proportional loss rate differentiation along with the FBS-DD provisioning. The controller enhances the controllability of network resource management.
4. We conduct extensive performance evaluation based on the simulation by the use of Bell Labs-I IP trace files. Results show that the proposed scheduling and buffer management schemes are capable of self-adapting to varying workloads of different classes. They automatically build a firewall around aggressive clients and protect network resources from saturation.

Our work is to address the integration of traffic policing with proportional differentiation. The study provides insights to the multi-dimensional differentiated services provisioning. The structure of the paper is as follows. In Section 2, we review existing packet scheduling and dropping algorithms for proportional differentiation provisioning. Section 3 presents the FBS-DD model with packet scheduling and buffer management schemes.Section 4 focuses on the performance evaluation. Section 5 concludes the paper.

## 2. Related work

Fair bandwidth sharing was initially studied as fair queueing [2], which aims to allow each flow passing through a network device to have a fair share of network resources. There are classic mechanisms for achieving the short-term per-flow fairing sharing, see PGPS [20] and a random scheme in [16] for examples. There is also recent study on fair load sharing in multipath communication networks [12]. In the context of DiffServ, the QoS provisioning is concerned with per-class behaviors. The FBS-DD model considers the long-term fair bandwidth sharing with delay differentiation.

### 2.1. Packet scheduling for proportional delay differentiation

Delay differentiation in packet networks is an active research topic. The PDD model is to provide differentiated delay services among traffic classes [4,5]. A class is assigned a delay differentiation parameter. The packet scheduler of a router aims to keep the ratio of average delay of a higher priority class to that of a lower priority class equal to the pre-specified value. The existing PDD algorithms can be classified into three categories [25].

*Rate-based packet scheduling algorithms* adjust service rate allocations of classes dynamically to meet the proportional delay differentiation constraints [4,13,14]. BPR [4] adjusts the service rate of a class according to its backlogged queue length so that the class service rates are proportional to the corresponding ratios of class loads. JoBS [14] allocates the service rate of a class based on delay predictions of its backlogged traffic. It forms the service rate allocation into an optimization problem when the system is heavy-loaded. The objective of JoBS is to enforce absolute delay and loss constraints. The accuracy of the rate-based algorithms over the delay ratio is unfortunately dependent of class load conditions [4]. This is because they rely on the relationship between queueing delay and service rate for a backlogged queue. However, the class load distribution on a router tends to change quickly. This limits the the applicability of the algorithms.

*Time-dependent priority packet scheduling algorithms* adjust the priority of a backlogged class according to the experienced delay of its head-of-line packet. In WTP [5], the priority of a backlogged class is set equal to the waiting time of the head packet normalized by its differentiation parameter dynamically on departure of each packet. A packet of a backlogged class with the highest priority will be forwarded next. Albeit simple, WTP implements the PDD model only when the system utilization approaches 100%. In AWTP [11,6], the control parameter of class is adjusted according to its class load dynamically. Moreover, a necessary condition was derived in [11], with respect to the class load conditions, for feasible WTP control parameters to achieve desired class delay ratios. It has better accuracy and adaptivity, in comparison with WTP, in both short and long timescales.

*Little's law-based packet scheduling algorithms* correlate the average queue length to the average arrival rate and the average queueing delay of packets. They control the actual delay ratio between two different classes by equalizing their normalized queue lengths with the pre-specified delay differentiation parameters.