

# Performance analysis of relative service using TCP-aware marking and dynamic WRED

Christos Bouras<sup>1,2,\*</sup>,† and Afrodite Sevasti<sup>2,3</sup>

<sup>1</sup>*RA Computer Technology Institute—RACTI, 61 Riga Feraiou Str, 26221 Patras, Greece*

<sup>2</sup>*Department of Computer Engineering and Informatics, University of Patras, 26500 Rion, Patras, Greece*

<sup>3</sup>*Greek Research and Technology Network-GRNET, 56 Mesogion Av., 11574, Athens, Greece*

## SUMMARY

The implementation of successful assured forwarding (AF) services according to the DiffServ framework remains a challenging problem today, despite the numerous proposals for assured forwarding per-hop-behaviour (AF PHB) mechanisms and AF-based service implementations. The interaction of the TCP and UDP traffic under an AF-based service and a number of relative issues such as fairness among flows, achievable bandwidth guarantees and qualitative performance have been taken into consideration in this work in order to address the existing limitations. We propose two modules, the TCP-window aware marker (TWAM) and the dynamic WRED (WRED) mechanism for implementing the differentiated services (DiffServ) AF PHB. We provide analytical models and an experimental evaluation in order to demonstrate how they succeed in enhancing the quality, improving the performance and easing the deployment of a production level AF-based service. Copyright © 2008 John Wiley & Sons, Ltd.

Received 6 August 2007; Revised 28 July 2008; Accepted 25 August 2008

KEY WORDS: assured forwarding; marking; active queue management; DiffServ; QoS provision; bounded delay

## 1. INTRODUCTION

The introduction of the Differentiated Services (DiffServ) framework has been a quite recent development in the direction of providing differential treatment to IP packets in backbone networks. DiffServ was introduced in order to solve the scalability problems emerging when deterministic, resource reservation mechanisms (e.g. RSVP and IntServ) were applied to networks serving large aggregates of individual microflows.

The DiffServ framework deals with traffic aggregates and keeps complexity at the edge of a network. It operates on the basis of policing and aggregating at the ingress of a DiffServ-enabled

---

\*Correspondence to: Christos Bouras, RA Computer Technology Institute—RACTI, 61 Riga Feraiou Str, 26221 Patras, Greece.

†E-mail: bouras@cti.gr, bouras@ceid.upatras.gr

domain, marking the packets of individual flows that belong to a certain QoS class with a single differentiated services code point (DSCP) value (or a group of DSCPs). In the interior of such a domain, queuing and scheduling of packets are performed according to their DSCP value and the per-hop-behaviour (PHB) configured for their DSCP, rather than the particular microflow to which they belong. DiffServ compromises the quality guarantees provided to IP traffic in order to introduce mechanisms that ensure differentiation while not requiring per-packet or per-flow management. In this way, scalability is achieved.

A lot of research is currently in progress on the issue of building successful services based on the DiffServ framework and on the mechanisms used in building corresponding PHBs. A DiffServ can be considered successful only if it provides reliable and quantifiable quality guarantees and at the same time can be exploited by a significant number of applications. Relative work has mainly focused on defining and provisioning end-to-end virtual-leased-line-like services with quality guarantees for reliable transmission of time-critical data according to the expedited forwarding (EF) PHB. In this work, we deal with another area of DiffServ, attempting to provide to IP traffic a service that is qualitatively better than that of the traditional 'best-effort' model, without deterministic, high-assurance quality guarantees. Such services are built on the assured forwarding per-hop-behaviour (AF PHB) of the DiffServ framework, as defined in [1].

An AF-based service ensures with very high probability the delivery of packets for flows or aggregates that conform to a predetermined profile. Packets belonging to AF class flows that fall outside the corresponding profile are served (or delivered to their destination) with a much smaller probability under congestion conditions. AF-based services do not provide strict delay/jitter guarantees, since they allow the use of non-reserved resources for the transmission of non-conformant packets when such resources are available. What they do provide is bandwidth guarantees for in-profile traffic under any conditions. In our previous work [2], we have proposed and outlined the basic principles of a service based on the AF PHB, denoted as Relative service.

Like all AF-based services, the Relative service provides a distinguishable from best-effort performance during congestion, while in un-congested periods the preferential treatment of eligible Relative traffic is not visible. The principal aim of the Relative service is to enable profile-conformant traffic to be served without realizing the negative effects of congestion, when congestion occurs. Among possible uses-applications of the Relative service would be:

- The provision of an assured transmission rate with minimized packet losses due to congestion between two end points over a backbone network.
- The controlled resources' allocation to flows over a congested transmission link.

AF-based services have not been very successful so far and are rarely deployed in production networks. There are quite a few reasons for this. The definition of the AF PHB leaves a lot of issues open for service designers, varying from service dimensioning principles to configuration of individual AF PHB mechanisms, ensuring fairness to different types of traffic, quantifiable means to assess performance and others. A great amount of research work, presented in the following section, has dealt with the individual open issues in AF-based services provisioning and the limitations of AF PHB mechanisms, but until today practical realization of such services has not been easy and a number of issues remain to be solved.

In this work we attempt to address the limitations of the Relative service provisioning. We aim to enhance the quality, improve the performance and ease the deployment of the Relative service, by proposing two distinct mechanisms for implementing the AF PHB and providing configuration guidelines for them. The proposed TCP-window-aware marker (TWAM), applied at

the network ingress, addresses the unfairness issues occurring when serving TCP flows under the Relative service with much less processing overhead than other proposed marking schemes. The proposed dynamic weighted random early detection (DWRED) mechanism, applied at the core routers, achieves much higher utilization than other active queue management (AQM) mechanisms, while it successfully adapts to changing network conditions. When combined for Relative service provisioning, with the aid of MPLS explicit routing, TWAM and DWRED are shown to outperform current AF mechanisms, achieving adaptability under transient load, fair differentiation and high performance. Early results of the work presented here are available in [3]. Together with the generic framework of our earlier work [2], the proposed mechanisms and specific configuration guidelines provided here contribute to the realization and practical implementation of AF-based services in production DiffServ-enabled networks. The Relative service and the proposed mechanisms adhere to AF PHB principles and, therefore, for the rest of this paper, the terms Relative and AF service will be used interchangeably.

In Section 2 we provide a thorough overview of the related work both on AF-based services and individual mechanisms. In Section 3 we give an overview of the Relative service framework and provisioning principles. In Sections 4 and 5 the proposed mechanisms are presented and approached theoretically, while in Section 6 a thorough experimental evaluation of the two mechanisms under the Relative services principles is conducted. This paper concludes with our future work and conclusions.

## 2. RELATED WORK

The effects of AF-based services provisioning to the TCP and UDP traffic and a number of relative issues such as fairness among flows, achievable bandwidth guarantees and qualitative performance improvement achieved by AF traffic have been a research topic for some years now. The main areas of work are:

- Refinements/variations of the assured service definition, as it was first introduced in [1].
- Fairness of distribution of resources between TCP and UDP flows under an AF service. Although related work on this topic is extensive [4–7], it is not analysed here as we focus on AF treatment to TCP traffic.
- General enhancements for TCP performance improvement under the AF service model.
- Focus on bandwidth fairness for TCP flows served under the AF service model.

Our work in this paper focuses on the latter two issues. However, in the following sections we provide an overview of the related work in all of the areas identified above.

### *2.1. Refinements/variations of the assured service definition*

In [1], where the idea for a service qualitatively better than that of the Internet best-effort was first introduced, a number of guidelines for the AF-based service provisioning are provided. However, due to the non-deterministic nature of the service itself and the PHB defined, a number of issues are left open for research.

In [8], it is stressed out that an ‘assured service’ can be implemented using a profile metre at the edge and a buffer management (or AQM) scheme at the core routers. This is the approach we follow in this work, with the introduction of the TWAM and DWRED correspondingly. An estimation and an assurance of the transmission rate for a flow can be imposed with the use of

markers such as the token bucket marker or the time-sliding window (TSW) marker [9] for marking packets as in- or out-of-profile. Different buffer management algorithms are evaluated and simple analytical models for the effective throughput of packets as well as the expected queuing delay are derived, based on Poisson arrivals of packets.

In [10], the possibility to provide assured rate (AR) guarantees over DiffServ-enabled MPLS networks, by using E-LSPs or L-LSPs, is investigated. In this way, the admission control per LSP can provide AR services with deterministic guarantees.

Azeem *et al.* [11] propose a number of improvements to the dropping, marking, shaping algorithms used by contemporary DiffServ-based services, so that they will become more user friendly. A TCP-friendly packet marker for TCP aggregates is proposed, which operates with a per-packet granularity, based on tokens per TCP flow. In [12], random early management (REM) is introduced in order to overcome the problems caused by RED in AF service differentiation. In REM, the target is to stabilize input rate at routers around the existing output link capacity and the router queue size around a target value.

In [13], it is emphasized how important it is for all flows participating in an AF-based differentiation schema to encounter a fair proportional loss rate. On this basis, the authors propose FRED, an AQM mechanism that maintains statistics for average queue occupancy per flow as well as the current number of packets buffered per flow. However, keeping per-flow state at core routers is against the principles of DiffServ. In [14], the author proposes a new service called alternative best effort (ABE), which achieves lower queuing delay for specifically marked traffic, while protecting best-effort traffic from perceiving the effects of this enhancement. However, as the author claims, ABE can be viewed as being positioned between the flat best-effort service and AF services.

## 2.2. TCP performance improvement under the AF service model

In [6], an analysis for the end-to-end guarantees provided by AF-based services is made. According to this analysis, the parameters that affect/obstruct the provision of quantitative guarantees to TCP traffic are the effects of round trip time (RTT) and the self-clocked sliding window mechanism, the number of concurrent flows in an AF class, the variety in packet sizes, the interference between TCP and UDP traffic and the contracted AF capacity. The latter parameters affect mainly the provision of guarantees for sharing unused resources. As an improvement to the findings of [6], our approach succeeds in addressing the limitations imposed by almost all of these parameters. As far as fairness in the distribution of excess bandwidth resources between flows of an AF class is concerned, [6] finds it quite difficult to achieve. In the experimental approach of our proposal, the results are more than encouraging.

Hollot *et al.* [15] and a number of related publications are based on classical linear feedback control theory and stochastic differential equation analysis to examine how the closed-loop system of TCP traffic and RED behaves. Both in [15] and [16], during an evaluation of the performance for aggregate flows under variations of AQM mechanisms (RED, gentle-RED, drop tail and instantaneous gentle RED), it is found that the mechanisms based on an instantaneous queue size give better results in terms of achieved queue sizes and thus packet delay. In [16] it is further shown that a small change in RED parameters can have a large impact on aggregate performance. Indeed, choosing RED parameters is a real challenge in an operational router, where the traffic fluctuates due to time of day effects. Our proposed DWRED addresses this issue successfully.

In [17] the authors propose separate configuration rules for TCP and UDP traffic under an AF-based service. These rules provide some guidelines for the configuration of the token bucket

policers and the WRED parameters of the AF queues serving TCP and UDP traffic. In [18] the authors are comparing the QoS perceived by two differentiation mechanisms: priority queuing and threshold dropping (a general definition of AQM mechanisms like RED and WRED). A result relative to our proposed work is that in order to fully utilize the benefit of a differential marking scheme, the choice of the maximum TCP-window size should take several parameters into account, such as the loss probabilities of in and out packets and the sender service profile. A threshold  $W_{\alpha}$  in the size of the congestion window is used, above which packets are marked as out. In and Out packets are dropped with different probabilities. It is shown that the incremental gain of achieved throughput by increasing  $W_{\alpha}$  decreases rapidly once exceeding a certain value. Our proposed TWAM is based on the TCP congestion window and its operation defines a specific marking policy, shown to achieve fairness and maximum throughput.

### 2.3. Bandwidth fairness achieved by the AF service model

In [19] deficiencies of classical RED are identified with respect to achieving fair AF-based service differentiation. In the AF framework, fair service provisioning means that flows or aggregates sharing the same resources obtain bandwidth in proportion to their contracted guaranteed capacity. In [19] it is shown how classic RED congestion notification is not dependent upon the number of active flows, a conclusion confirmed by our findings in [2, 3] where we have also shown that the feedback provided by classic WRED is not very successful in service differentiation. The authors of [19] propose adaptive RED according to which the average queue length is constantly observed and updated.

In [20] an effort to study the effects of different RTTs of individual flows in the capacity obtained by the AF aggregate to which they belong is made. The results show that differences in RTT for TCP flows affect the capacity observed by these flows individually. In order to balance the unfairness in throughput perceived by TCP flows with different RTTs, the authors of [21] propose to set the ECN bit of individual flows with a different probability according to their RTT. In our proposed scheme, the unfairness due to RTT differences between TCP flows is taken into consideration, for a more efficient AF service provisioning.

In [5] it is emphasized that the regulating factor in capacity distribution under an AF-based service model is the quantity of in-profile traffic. In [9], the TSW packet marker is proposed in an effort to achieve fairness in capacity distribution among flows based on the reserved capacity for each flow. The marker estimates the current rate of each flow on a per-packet basis and marks packets accordingly in two levels, as in and out, depending on whether the current rate of a flow exceeds or falls short of the flow's contracted capacity. It is based on a TSW estimation of a flow's current rate and uses two levels (or 'colours') of marking for packets. Therefore, it can be referred to as TSW two colour marker (TSW2CM). The mechanism itself and certain variations (like the TSW three colour marker) appear often in the literature. Taking this into consideration, we have chosen to compare the performance of our proposed TWAM with the TSW2CM.

In [22], the problem of unfairness of bandwidth allocation among TCP flows under an AF-based service is solved with the adoption of core-stateless fair queueing. Edge routers classify the incoming packets into flows keeping their state in the form of flow arrival rate and upon packet reception, the core router re-computes the fair share for all flows in transit and probabilistically punishes unfair flows. This approach requires state at a flow granularity to be kept at the core routers and, thus, does not fully comply with the DiffServ principles.

Another approach for improving the performance of TCP flows using the RED with in/out bit (RIO) AQM mechanism at the routers is presented in [23]. The objective of the proposed marking is to selectively mark packets, based on TCP flow state, in order to reduce the possibility of TCP entering the slow start phase due to retransmission timeouts. In the simulation study that follows, the flows with smaller RTT obtain better gains with the proposed mechanisms. Our proposed mechanisms also succeed in preventing TCP flows from entering slow start, achieving an improved performance in terms of resources' utilization.

In [24] the authors use a Markov process to model the window size evolution of TCP connections sharing a bottleneck router, the queue of which operates either under the tail drop principle or the RED AQM. It is shown that the RED router results in smaller dispersion of the window sizes, which implies that RED can improve fairness among many TCP connections. In [25] a variation of RED is proposed where AQM thresholds apply only to out packets. An interesting proposal here is to consider the TCP congestion window (CWND) as consisting of two parts, a reserved part (RWND) equal to the product of the reserved rate and the estimated RTT and a variable part (CWND-RWND) that tries to estimate the residual capacity. At the beginning of fast recovery, only the variable part of the CWND is reduced by half, and at the beginning of slow start, CWND is set to RWND+1 instead of 1. Ssthresh is updated accordingly. Our approach also takes into consideration the TCP congestion window and improves the throughput performance of the AF-based service, without requiring changes to the TCP stack.

In [26] a TCP-friendly packet marker is proposed that allocates only in tokens to flows with small TCP window, and makes a 'max-min' fair allocation of the rest of the tokens to the remaining flows. It also maintains optimum spacing between in and out tokens allocated for a flow. Packet marking is made according to the per-flow allocation of tokens. Less TCP timeouts and packet losses are observed, while the TCP throughput is increased and less variable than when using a plain token bucket marker.

After reviewing the related work in fairness achieved by AF treatment, it is important to stress that our proposed approach is different from the ones presented here. It does not only concern efficient marking and AQM but also introduces MPLS-enabled explicit routing for flow aggregates, thus, adding slightly to the complexity but achieving an improved performance.

Several issues continue to be open for investigation concerning AF-based services provisioning and dimensioning in the cases of realistic topologies. Such issues include:

- Marking mechanisms that provide proportional capacity distribution.
- Provision of AF-based QoS to isolated TCP sessions (effects of RTT etc.).
- Better than best-effort service quality (e.g. bounded average delay) for the AF traffic.

Our approach as presented in the following sections outlined in the sequel aims to address these issues and provide an evaluation of the proposed mechanisms against existent proposals and implementations.

### 3. SERVICE DEFINITION

This section attempts to provide an overview of the Relative service framework as an AF-based or AF PHB compliant service to which the proposed mechanisms of TWAM and DWRED are applicable. A more detailed specification of the Relative service can be found in [2].

By definition, the service is provided to individual TCP flows or aggregates of TCP flows as they cross a DiffServ-enabled backbone network. According to the related research work already

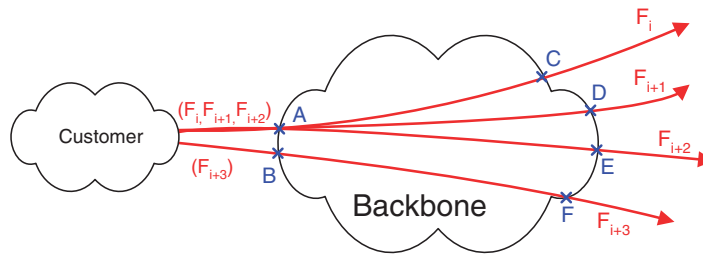


Figure 1. A destination-aware AF service model.

presented, isolation of TCP and UDP traffic in the case of AF-based services is recommended for fairness and guarantees' provisioning purposes. Owing to the fact that provisioning of an AF-based service to TCP traffic presents, unlike UDP traffic, a number of challenges, we will focus our approach to provisioning principles for TCP traffic. Thus, any references to Relative traffic from now on will imply TCP traffic. Furthermore, for simplicity purposes of the upcoming analysis, we will consider only individual TCP flows rather than flow aggregates.

According to the proposed Relative service model, the customers agree with the network on a guaranteed average service rate and bounded average end-to-end delay for in-profile Relative packets. The guarantees provided together with a number of provisioning specifications and a profile for the customer traffic are included in the service level agreement (SLA) signed between each customer and the network service provider.

The Relative service and more specifically the proposed mechanisms ensure:

- That each customer's AF flow or aggregate obtains at least an average throughput equal to that specified by the corresponding SLA regardless of congestion.
- That the average size of router queues serving Relative traffic is bounded, so that served (non-dropped) packets face a bounded delay.
- That AF traffic exploits additional resources whenever those are available, while preserving bounded average queuing delay.
- That fairness in capacity distribution between TCP flows is ensured regardless of the network operating conditions.

### 3.1. Service provisioning

We propose the adoption of a destination-aware service, according to which the SLA signed with each customer identifies as many AF flows served by the network as the number of destinations or egress routers (ER) that the customer wishes to address his Relative traffic to. In the example of Figure 1, an individual customer signs an SLA for Relative service provisioning with the backbone network provider such that for each ingress to ER pair between which customer's traffic is served, a different Relative service provisioning instance is identified. Thus, TCP flows  $F_i, F_{i+1}, F_{i+2}$  from the upstream customer domain of Figure 1, although entering the backbone at the same edge router A, are distinguished to different provisioning instances based on their corresponding egress routers, thus C, D, E.

Moreover, the DiffServ-enabled backbone itself adopts traffic engineering to pin the entire path for the set of customers' AF aggregates between a certain ingress and ER pair. This can be achieved with the use of multi-protocol label switching (MPLS) as a traffic engineering technique. In MPLS,

labels are assigned to packets and used in the forwarding process in each router to explicitly define the next hop in the path from the source (ingress router (IR)) to the destination (ER). All AF traffic flows from different customers that enter the backbone network in the same ingress (edge) router and exit the backbone network from the same ER can be assigned to the same MPLS labelled switched path (LSP) and thus, comprise a forwarding equivalence class (FEC) according to the MPLS terminology.

For the purposes of the proposed Relative service, we define such an FEC comprised of AF flows as an AF FEC. Each Relative service customer can participate in an AF FEC with a single flow or aggregate. It is compulsory for all the TCP traffic from a single customer between an egress and an IR of the backbone network to be assigned to the same AF FEC and be treated as a single entity. However, traffic from different customers can use different LSPs and thus, AF FECs between the same edge routers. Still, the AF FEC LSP notion preserves scalability as it remains independent of the number of transit flows.

Although TCP traffic is unidirectional, TCP ACKs from the egress back to the IR also need to be explicitly routed. Hence, each AF FEC LSP needs to be configured bidirectionally.

Explicit routing, in the context of the Relative service, is essential as it will be shown in the sequel. It ensures that all the TCP flows participating in a Relative service provisioning instance use the same path and thus, experience the same RTT. By using a destination-aware model and MPLS traffic engineering, unfairness due to unequal RTTs is no longer a concern in our provisioning model.

### 3.2. Per-hop mechanisms

In order for the proposed service definition to be complete, after defining the destination aware, explicit routed paths and excluding UDP traffic, the PHB principles remain to be described. Here is where the core of our contribution in this work lies.

The proposed Relative service is based on the AF DiffServ PHB. According to the AF PHB specification, AF-based services do not require policing of AF flows. Instead, AF traffic is examined against certain predetermined profiles at the network edges and packets are marked as in- and out-of-profile. Marking of packets at the edge routers when combined with buffer management techniques in the core is a means to exercise conformance and confinement of AF traffic to the resources available. This is particularly the case under the assumption of TCP traffic, due to the TCP adaptation and congestion avoidance mechanisms.

According to the proposed model, within an AF FEC, one AF profile with a guaranteed average service rate is defined for each participating customer. This profile is used by a marker at the AF FEC LSP IR to distinguish between AF packets of a customer's flow that fall within and out of the profile. In-profile packets are thus marked as 'in' and the remaining packets are marked as 'out' ones. In the following sections, we propose a TCP-window-aware marker to implement the marking block of the AF-based service.

In the core of the AF-enabled network, a dedicated queue  $Q_{AF}$  is configured at each core router along the AF FEC LSP for serving the aggregated AF FEC traffic with a minimum guaranteed service rate equal to  $C_{AF}$ .  $C_{AF}$  can vary in the individual core routers used by a single AF FEC LSP; however, in the analysis that follows, we make the simplifying assumption that  $C_{AF}$  is constant along the AF FEC LSP. Moreover,  $C_{AF}$  has to be configured in such a way that under congestion or heavy load in the router due to traffic of other classes, AF traffic is provided with a minimum assured service rate. A queuing mechanism to achieve this is weighted fair queuing (WFQ), which



allows one  $Q_j^{\text{AF}}$  with a guaranteed service rate of  $C_j^{\text{AF}}$  for each  $\text{LSP}_j^{\text{AF}}$  implemented over router's outbound link. This service rate has to be such that in all queues along the AF FEC LSP that the average service rate guarantees included in SLAs with individual customers participating in the specific LSP are not violated. As shown in [2] and in the experimental analysis that follows, the recommended value for  $C_{\text{AF}}$  is slightly higher than the sum of contracted capacity for the flows belonging to a certain AF FEC LSP.

Apart from the provisioned minimum  $C_{\text{AF}}$  service rate for AF packets at the AF FEC LSP core routers, the AF PHB includes also a buffer management or AQM mechanism to manage the size of  $Q_{\text{AF}}$ . AQM mechanisms have been used as traditional congestion avoidance mechanisms in routers well before DiffServ came into the picture. A large majority (such as RED, WRED, RIO etc.) use packet marking or dropping as a means for congestion notification towards the traffic sources and are relying on the adaptiveness of TCP flows.

According to our approach and in the framework of AF service provisioning, an AQM mechanism operates upon in and out packets in an AF service queue, so as to control the number of accumulated packets in  $Q_{\text{AF}}$  and drop packets whenever AF resources become scarce. AF traffic sources are thus, notified of congestion in order to adapt their sending rates. In the following sections, we adopt a widely used AQM mechanism in our service model and propose a methodology for dynamic configuration of its parameters so as to achieve:

- Dropping of out packets during congestion, as a means of congestion notification to the TCP flows in a way proportional to their contracted capacity.
- Bounded average queue size in  $Q_{\text{AF}}$  and thus, bounded average queuing delay for the AF packets, despite of the network conditions.
- Utilization of excess resources when available.
- Adaptation to preserve the fairness and high utilization properties when new aggregates are added to the AF FEC LSP or when the load in other traffic classes fluctuates.

In the following sections, we provide a detailed specification of the proposed mechanisms for the case of a single AF FEC LSP. We propose the TCP-window-aware marker applied to each customer's AF flow participating in the LSP. We also propose specific configuration guidelines for the weighted random early detection AQM, used to manage the size of  $Q_{\text{AF}}$ . We are assuming equal  $Q_{\text{AF}}$  maximum sizes and  $C_{\text{AF}}$  configured along the routers of the LSP. However, our analysis can be extended for generic cases, provided that the WRED configuration principles are enforced at least on the 'bottleneck' of the LSP, thus, on the router with the minimum configured value for  $C_{\text{AF}}$ . The proposed TWAM in conjunction with the proposed use of WRED is verified experimentally to achieve the goals of the Relative service as already defined.

#### 4. A TCP-WINDOW-AWARE MARKER

In this section, we present a TCP-window-aware marker for use in the IR of each AF FEC LSP. We model the TWAM operation principles for the case of a single LSP, but they can be easily generalized for the case of multiple AF LSPs crossing a network domain. We also make the assumption that a single TCP flow rather than an aggregate participates to the AF LSP on behalf of each customer. Experimental data verify that this assumption does not compromise the performance advantages of the TWAM.

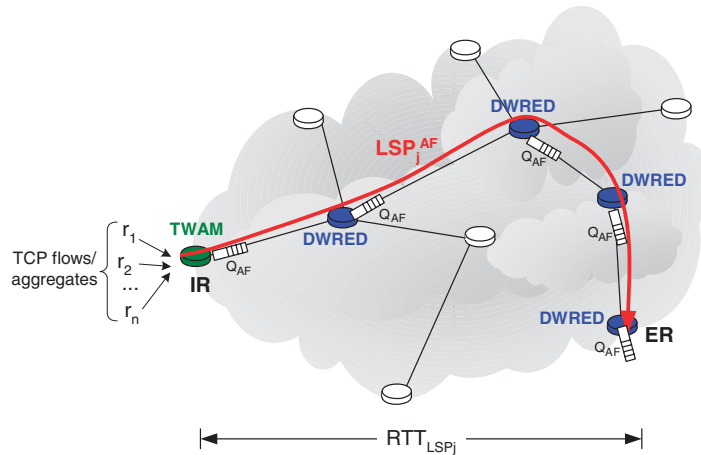


Figure 2. Deployment of the TWAM and DWRED mechanisms over an AF FEC LSP.

A number of TCP flows, forming an AF FEC, are injected to IR and routed to the LSP ER of an AF LSP denoted by  $LSP_j^{AF}$ , as depicted in Figure 2. The TCP flows are assumed to always have packets to send. According to the SLA of the corresponding customer, each TCP flow  $i$  has to be served with an average rate of  $r_i$ . Owing to the fact that all flows of  $LSP_j^{AF}$  follow the same route, they experience the same RTT (an assumption verified by the experimental results later on, see Figure 15). The TWAM operates on each flow or customer aggregate separately at IR. It marks the packets of each flow or aggregate as in and/or out according to the TWAM parameters for the specific flow.

The proposed TWAM is self-tuned in intervals of  $\alpha * RTT_{LSP_j}$  seconds (with  $\alpha \in \{1, 2, 3, \dots\}$ ) for the specific  $LSP_j^{AF}$ . The  $\alpha$  parameter determines how often in consecutive  $RTT_{LSP_j}$  intervals does the proposed marker update its configuration parameters for each of the TCP flows. In other words, it determines how prompt will the TWAM be to adapt its operation according to changing network conditions. For simplicity purposes, the analysis that follows adopts a value of  $\alpha = 1$ , assuming that the TWAM configuration parameters are reconsidered, based on the current conditions in the network, every  $RTT_{LSP_j}$  seconds. However, the experimental study that follows shows that for small values of  $\alpha$ , a good approximation is also achieved and the results still confirm the analytical approach.

We also make the simplifying assumption of an equal packet size among the TCP flows in  $LSP_j^{AF}$ , denoted by  $S$ . However, the analysis that follows applies in the case of different packet sizes if all TCP-window figures are expressed in terms of bytes rather than segments, and of course the benefits of the TWAM are preserved.

We define for each flow, the reserved congestion window  $W_i^r$  as

$$W_i^r = \frac{r_i \times RTT_{LSP_j}}{8 \times S} \quad (1)$$

to denote the minimum amount of in segments that the flow can inject to the  $LSP_j^{AF}$  in an interval equal to  $RTT_{LSP_j}$ . We denote by  $W_i^{cwnd}$  the current congestion window of each flow and we assume that the TCP receiver-advertised window is always larger than  $W_i^{cwnd}$ , so that  $W_i^{cwnd}$  is the window used by the TCP sources to transmit packets. Using the assumptions and definitions

already provided, the following proposition can be formed for the proposed TWAM:

*Proposition*

The TWAM distributes the resources available for the total traffic of an AF FEC among the individual AF FEC flows in a fair manner, analogous to the  $W_i^r$  of each flow.

In order for this proposition to be honoured, the TWAM should mark arriving packets of flow  $i$  as in with probability  $p_i$  and as out with probability  $1-p_i$  so that during an  $RTT_{LSP_j}$  interval, for each flow  $i$ , the number of packets marked as in and the number of packets marked as out are correspondingly:

$$N_i^{\text{In}} = p_i * W_i^{\text{cwnd}} = \frac{W_i^r}{\sum_i W_i^r} * \sum_i W_i^{\text{cwnd}} \quad \text{if } \frac{W_i^r * \sum_i W_i^{\text{cwnd}}}{\sum_i W_i^r * W_i^{\text{cwnd}}} < 1 \quad (2a)$$

$$N_i^{\text{Out}} = (1 - p_i) * W_i^{\text{cwnd}} = W_i^{\text{cwnd}} - \frac{W_i^r}{\sum_i W_i^r} * \sum_i W_i^{\text{cwnd}}$$

$$N_i^{\text{In}} = W_i^{\text{cwnd}} \quad \text{otherwise} \quad (2b)$$

$$N_i^{\text{Out}} = 0$$

The analysis that supports this proposition and provides the value of  $p_i$  follows.

Based on the slow start and congestion avoidance mechanisms of TCP, a flow at a certain moment in time, can be in one of two conditions:

*Case 1:*  $W_i^{\text{cwnd}} \leq W_i^r$ . All segments transmitted in the current  $RTT_{LSP_j}$  interval fall within the reserved congestion window of the flow.

*Case 2:*  $W_i^{\text{cwnd}} > W_i^r$ . A part of the segments transmitted in the current  $RTT_{LSP_j}$  interval fall out of the reserved congestion window of the flow. The amount of these packets is equal to  $W_i^{\text{cwnd}} - W_i^r$ .

We also define the bandwidth-delay product for the aggregated AF traffic of  $LSP_j^{\text{AF}}$  or the  $LSP_j^{\text{AF}}$  congestion window (thus, the number of segments that can be sent over  $LSP_j^{\text{AF}}$  in the upcoming RTT interval) as

$$W_{LSP_j}^{\text{cwnd}} = \frac{R_{\text{AF}} \times RTT_{LSP_j}}{8 \times S} \quad (3)$$

where  $R_{\text{AF}}$  is the current aggregate rate achieved by all of the flows participating in  $LSP_j^{\text{AF}}$ . The exact value of  $R_{\text{AF}}$  and thus,  $W_{LSP_j}^{\text{cwnd}}$  depends on the congestion conditions along the  $LSP_j^{\text{AF}}$  routers, e.g. the load of traffic belonging to higher priority QoS classes. Under severe congestion in other traffic classes, the maximum amount of AF segments that can be sent through  $LSP_j^{\text{AF}}$  in an interval equal to  $RTT_{LSP_j}$  is given by

$$W_{LSP_j}^{\text{cwnd}}(\text{min}) = \frac{C_{\text{AF}} \times RTT_{LSP_j}}{8 \times S} \quad (4)$$

This result is self-evident for a single TCP flow and also  $W_{LSP_j}^{\text{cwnd}}(\text{min})$  is the upper bound for an aggregate of TCP flows:

$$W_{LSP_j}^{\text{cwnd}}(\text{min}) \geq \sum_i W_i^r \stackrel{(1)}{\leftarrow} \quad (5)$$

$$C_{\text{AF}} \geq \sum_i r_i \quad (6)$$

As already specified in the Relative service definition in Section 3.2, (6) must hold at all times along  $LSP_j^{AF}$ .

The aim of the TWAM is to mark the  $W_i^{cwnd}$  packets of each TCP flow  $i$  arriving within an interval  $RTT_{LSP_j}$  in proportion to  $W_i^r$ . This achieves a fair distribution of  $W_{LSP_j}^{cwnd}$  to the TCP flows, according to their subscribed  $r_i$  and thus, their reserved congestion window  $W_i^r$ . This is anticipated to achieve fairness in the distribution of the current  $R_{AF}$  among the individual flows of  $LSP_j^{AF}$ . Similar approaches for achieving fairness among TCP flows, such as the TSW2CM of Clark and Feng [9], have so far only relied on rate approximations per TCP flow, trying to mark packets so that each TCP source would adjust the sending rate to a value proportional to its contracted  $r_i$ . However, real-time per-packet rate estimation is error prone, estimation parameters differ under different scenarios and the various methodologies are susceptible to deviations under transient network load. Moreover, many of the proposed schemes of rate estimation impose excessive overhead, as they impose rate estimation at each packet arrival. Using the TCP congestion window and a time granularity of one or a few RTTs according to our proposal, provide a more accurate and efficient estimation improving the achieved fairness. Moreover, TWAM can be successfully used with the proposed in the following section DWRED configuration scheme to achieve bounded average queuing delay under different operating conditions of the Relative service.

For the operation of TWAM, in order to fairly mark packets, and since it is not straightforward to estimate the exact value of  $R_{AF}$  and thus of  $W_{LSP_j}^{cwnd}$  for each RTT interval, we make the following approximation:

$$W_{LSP_j}^{cwnd} \approx \sum_i W_i^{cwnd} \quad (7)$$

Another alternative would be to use the TCP flows' ssthresh values emerging from the previous RTT interval, thus,

$$W_{LSP_j}^{cwnd} \approx \sum_i ssthresh_i \quad (8)$$

since ssthresh reflects what TCP perceives to be the optimal operating point under the current network conditions. Both (7) and (8) approximate the value of  $W_{LSP_j}^{cwnd}$ , which is actually an estimator of the current conditions on the network and more specifically of the resources available for the traffic of  $LSP_j^{AF}$ . We will use (7) in this work and leave the study of (8) for our future work.

Once  $W_{LSP_j}^{cwnd}$  is estimated, TWAM has to enforce that it is redistributed among the existent flows in a fair manner. In order to achieve this, TWAM marks arriving packets of flow  $i$  as in with probability  $p_i$  and as out with probability  $1 - p_i$  throughout the current RTT interval, where

$$p_i = \begin{cases} \frac{W_i^r * \sum_i W_i^{cwnd}}{\sum_i W_i^r * W_i^{cwnd}} & \text{if } \left( W_i^r * \sum_i W_i^{cwnd} \right) < \left( \sum_i W_i^r * W_i^{cwnd} \right) \\ 1 & \text{otherwise} \end{cases} \quad (9)$$

Hence, during the  $RTT_{LSP_j}$  interval, for each flow  $i$ , the number of packets marked as in and the number of packets marked as out are provided by (2a) and (2b). In this way the available  $LSP_j^{AF}$  congestion window (see (4) and (7)) is 'coloured' in a way that reflects the individual TCP flows' reservations, rather than the current achieved congestion windows by each flow individually. Owing to (5) each flow is allowed to send in packets in proportion to its reserved congestion window

Table I. An example of the TWAM marking principles.

	$r_i$ (Mbps)	$W_i^r$	$p_i$	$W_i^{cwnd}$	In packets	Out packets
1	2	50	0.73	60	$0.73 * 60 = 44$	16
2	3	75	0.55	120	$0.55 * 120 = 66$	54
3	10	250	1	150	150	—

subscription (provided by (1)). In addition, flows are allowed to send in packets in excess of  $W_i^r$  in a controlled and fair manner.

As can be seen in the example of Table I, flow 1 and flow 2 exceed their reserved congestion window  $W_i^r$  according to (1) (with  $RTT = 0.2$  s and  $S = 1000$  bytes). Thus, TWAM allows for only a percentage of their packets in the current interval to be marked as in and the rest of the packets in their current congestion window  $W_i^{cwnd}$  are marked as Out. Flow 3 comes short of its  $W_i^r$  and thus, can mark all of its packets in the current interval as in.

A set of experiments will be provided in the following section in order to verify that the proposed TWAM achieves fairness in terms of sharing resources among the TCP flows that share the same AF LSP over an AF-enabled backbone, without preventing adaptations towards the use of excess resources whenever they are available. One of the most valuable achievements of TWAM is that it succeeds in driving TCP flows to continuously operate in congestion avoidance at different levels under any conditions (both congested and un-congested periods) and eliminates occurrences of the slow start phase. This important property ensures high utilization and eliminates oscillations for TCP flows.

## 5. DYNAMIC WRED

Recall that WRED is an AQM mechanism mainly used in core backbone routers. RED (on which WRED is based) uses the average queue occupancy as a parameter in a random function that determines whether the congestion avoidance mechanisms (e.g. packet drops) should be triggered. As the average queue occupancy increases so does the packet drop probability. According to RED, when the average queue occupancy ( $Q_{avg}$ ) is below a minimum threshold,  $min_{th}$ , then no packet is dropped, while when  $Q_{avg}$  exceeds  $min_{th}$  the packet drop probability increases linearly up to a maximum value  $max_p$ . When  $Q_{avg}$  exceeds  $max_{th}$ , then the packet drop probability is equal to 1, thus, all packets are dropped. In WRED, packet drop probability for each packet arriving to the queue is calculated according to a single value of  $Q_{avg}$ , but the values of  $min_{th}$ ,  $max_{th}$ ,  $max_p$  differ according to the colour of each packet. The WRED adopted for the Relative service model is used to drop only out packets. Although it is based on the occupancy of  $Q_{AF}$  by both in and out packets to estimate the value of  $Q_{avg}$ , it defines only one set of parameters  $min_{th}$ ,  $max_{th}$  and  $max_p$ . According to the WRED principles, when the size of  $Q_{AF}$  is between  $min_{th}$  and  $max_{th}$ , an arriving out packet is dropped with probability  $max_p$  equal to

$$max_p \times \frac{Q_{avg} - min_{th}}{max_{th} - min_{th}} \quad (10)$$

A lot of research work conducted for investigating the principles of operation of AQM mechanisms such as RED and WRED. Generic recommendations exist for the configuration of such AQM

mechanisms both in best-effort networks and in the framework of AF services' provisioning. In almost all cases, RED or WRED parameters are fixed during the operation of a router queue (see the Appendix for a number of configurations proposed in the literature). As we will also experimentally verify, this can lead to inefficient utilization and unfair differentiation of resources.

We propose the use of a dynamically re-configurable WRED AQM mechanism to be used by the Relative service model defined in Section 3. Dynamic configuration of the WRED parameters ( $\min_{th}$ ,  $\max_{th}$  and  $\max_p$ ) in the AF queues of the  $LSP_j^{AF}$  routers takes place (if necessary) at intervals equal or multiple of  $RTT_{LSP_j}$ , thus, in accordance with the operation intervals of the TWAM. The idea here is to 'co-ordinate' the AQM with the TWAM operation. The proposed WRED configuration mechanism responds to fluctuations in available resources, allowing the use of excessive resources whenever they are available, in a way that achieves a bounded average queuing delay for packets in each  $Q_{AF}$  along  $LSP_j^{AF}$ . At the same time, fairness of service among TCP flows sharing each  $Q_{AF}$  is ensured.

In our previous work [2, 3], we have demonstrated how the adaptation of the  $\min_{th}$ ,  $\max_{th}$  parameters of WRED according to the bandwidth-delay ( $BW * D$ ) product of a link achieves a bounded average queue size of approximately  $\{0.2 \dots 0.25\} * (BW * D)$ , while at the same time ensures a high utilization of available resources by the TCP traffic. This is the generic principle of operation for the proposed DWRED.

In this work, we use the TWAM operation principles to guide the DWRED parameters configuration. We propose that for each operating interval  $RTT_{LSP_j}$ , the  $\max_{th}$  parameter of WRED is set equal to  $W_{LSP_j}^{cwnd}$ , derived either from (7) or (8):

$$\max_{th} = \sum_i W_i^{cwnd} = W_{LSP_j}^{cwnd} \quad (11)$$

In this way, DWRED incorporates fluctuations in available resources by exploiting the implicit indication provided by the congestion windows of TCP flows and staying in accordance with the result presented in [2, 3]. DWRED uses the value of  $W_{LSP_j}^{cwnd}$  as the indication of the end-to-end throughput of  $LSP_j^{AF}$ , with the proposed value for  $\max_{th}$  embodying all related parameters otherwise addressed by the bandwidth-delay product, including the number of hops on the particular LSP.

In a static WRED configuration,  $\min_{th}$  and  $\max_{th}$  have fixed values. If these values are set too low, this results in too many packet drops, which do not allow TCP flows to open up their congestion windows when excess bandwidth resources are available. On the other hand, if they are set too high, fewer out packets fall within the  $\min_{th}$ ,  $\max_{th}$  region and thus, dropping is not as effective in providing congestion notification back to the sources. Unlike static WRED configurations, DWRED allows for adaptation of the AQM mechanism to the current operating conditions and ensures that both of the aforementioned undesirable situations are avoided.

DWRED is particularly important for ensuring fairness among TCP flows, as in combination with TWAM, signals to the appropriate set of TCP sources that they should reduce their sending rate in order for fairness to be achieved.

Under the framework of DWRED, the use of values

$$\min_{th} = b * \max_{th}, \quad \max_p = 0.1 \quad (12)$$

is also proposed, thus,  $\min_{th}$  is proposed to also be updated for each operating interval in order to 'follow' the fluctuations of  $\max_{th}$ . For the value of  $b$  in (12), our previous work and related research

appoint values in the interval (0.1, 0.2) as preferable. We have adopted the value of  $b=0.15$  for the experimental evaluation that follows.

The proposed DWRED attempts to improve the performance and better adapt to available resources without compromising bounded average queuing delay guarantees.

## 6. EXPERIMENTAL APPROACH

For the evaluation of the proposed AF service mechanisms, we have conducted a series of experiments, through which we have compared the performance of our approach with well-known mechanisms for AF-based services provisioning. For our experimentation, we have used the ns-2 simulator [27] and more specifically the DiffServ functionality built into this simulation environment. We have implemented our proposed PHB mechanisms and integrated them into the simulator. The experiments presented in the following sections were carried out using the TCP Reno implementation of ns-2, while WFQ was used to configure the queue service rates for Relative traffic.

### 6.1. DWRED evaluation

In order to demonstrate the inefficiency of a static WRED configuration in a more generic setting than that of Relative service provisioning, we have conducted a simple experiment of 10 TCP flows competing over a 30 Mbps link (**scenario #1**). By parameterizing the propagation delay on the link, we ran experiments with RTTs of 60, 240 and 480 ms. The TCP segment size is set to 1500 bytes and the TCP sources are conducting FTP transfers and thus, are assumed to always have data to send. All experiments were conducted for 1000 s and in the middle of their duration (at 500 s), the UDP traffic of a rate exceeding 10 Mbps was injected to the same link traversed by TCP flow. UDP traffic was served with absolute priority from a separate queue with rate limitation at 10 Mbps. A static WRED configuration of ( $\min_{th}=10$ ,  $\max_{th}=30$ ,  $\max_p=0.02$ ) is compared with a DWRED configuration at the router queue serving TCP traffic.

As can be seen from Figure 3(I), static WRED allows TCP flows to achieve a smaller average congestion window. Thus, with the exception of  $RTT=60$  ms (Figure 3(II)), where the results are comparable, the throughput of TCP flows (total number of successfully transmitted packets) is quite higher in the DWRED case (Figure 3(III) and (IV)). For this particular experimental setting, the static WRED configuration seems to approximate the DWRED performance only when  $RTT=60$  ms, as for this case the equivalent DWRED parameters ( $\min_{th}$  and  $\max_{th}$ ) obtain values closer to those of the static WRED parameters.

A different set of WRED parameters ( $\min_{th}=40$ ,  $\max_{th}=160$ ,  $\min_p=0.1$ ) (**scenario #2**) seem to better approximate the DWRED performance for the particular experimental set-up as shown in Figure 4. It is obvious that the optimum WRED configuration for a router queue differs according to a number of parameters such as the topology, RTT values, traffic load, service principles, etc. This is one of the main reasons why research has focused on studying RED and WRED dynamics and yet no specific WRED tuning guidelines with universal validity exist today.

DWRED provides a simple but efficient WRED configuration methodology, applicable to different environments. In its generic form, as initially evaluated in [2, 3], DWRED is nothing but a means to dynamically configure an AQM mechanism in backbone routers of the best-effort Internet. In the following sections, it will be experimentally verified how the use of DWRED in

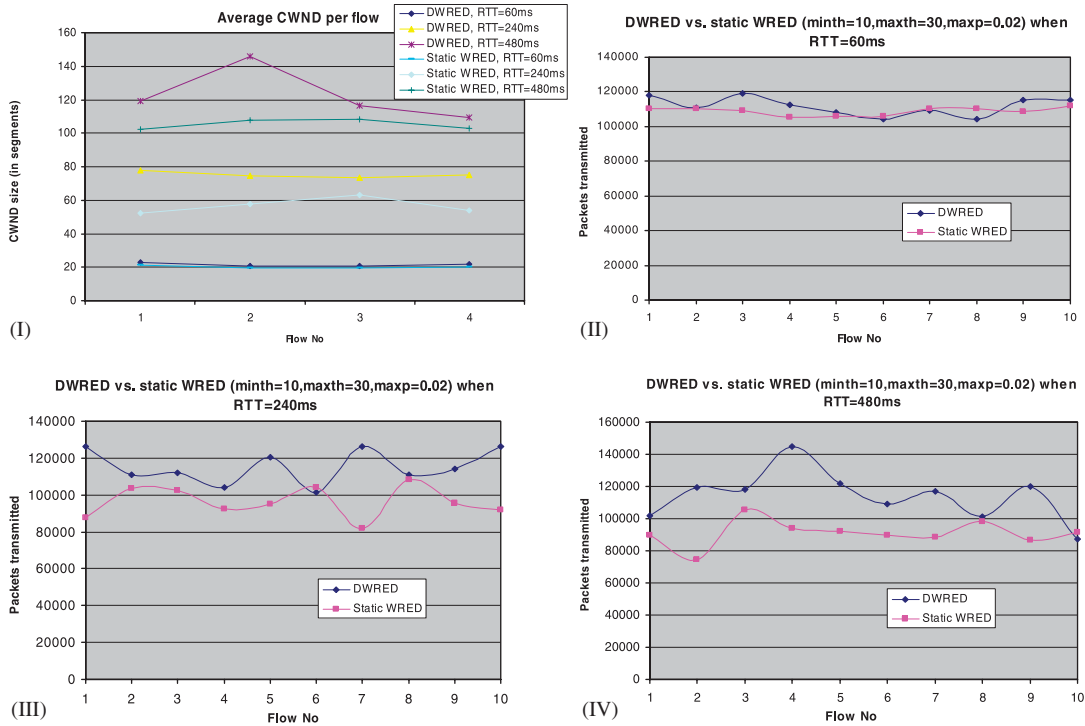


Figure 3. Comparison of static WRED and DWRED configurations under scenario #1.

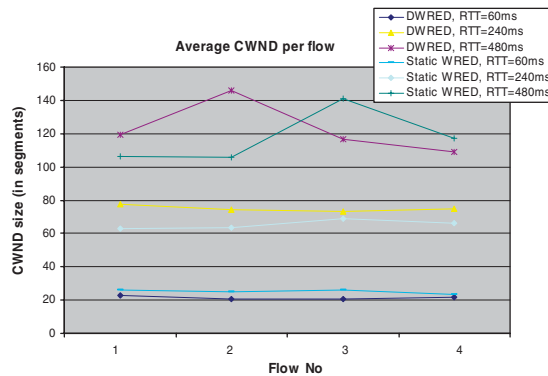


Figure 4. Comparison of static WRED and DWRED configurations under scenario #2.

the framework of the Relative service (as specified in expressions (11) and (12)) combined with the TWAM can achieve fair allocation as well as very high utilization of available resources (up to 100%).



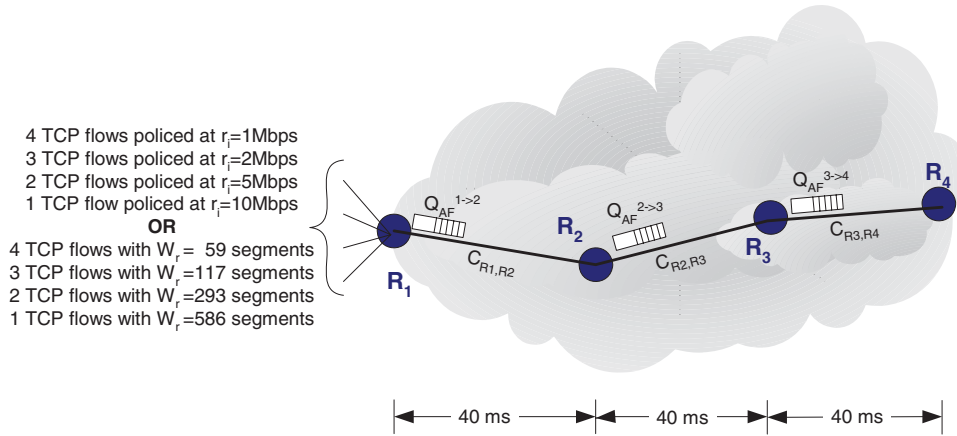


Figure 5. The experimental set-up of scenario #3.

Table II. TCP flows participating in scenario #3.

Flow no.	Contracted capacity $r_i$ (or CIR $_i$ ) (Mbps)	Reserved congestion window $W_i^f$
1	1	59
2	1	59
3	1	59
4	1	59
5	2	117
6	2	117
7	2	117
8	5	293
9	5	293
10	10	586

### 6.2. Basic TWAM and DWRED evaluation

For the comparison of well-known mechanisms for AF-based services provisioning (such as the TSW2CM and static WRED) with the TWAM and DWRED proposed in this work, the experimental set-up of Figure 5 has been used. A three-hop AF FEC LSP is used for serving Relative traffic entering a DiffServ-enabled domain in router  $R_1$  and exiting from router  $R_4$ . For this initial set of experiments (**scenario #3**), the backbone links have a capacity of 30 Mbps and a dedicated router queue ( $Q_{AF}$ ) is configured to serve Relative traffic at each router. The aggregate of TCP flows injected in  $R_1$  is provided in Table II.

Two different configurations are compared:

- Use of the TSW2CM in  $R_1$  to examine upon each packet arrival the current rate of the TCP flow (marking packets as out when the flow  $r_i$  is exceeded) and use of static WRED with  $\{\min_{th}=40, \max_{th}=160, \min_p=0.1\}$  in each of the  $Q_{AF}$ s depicted.

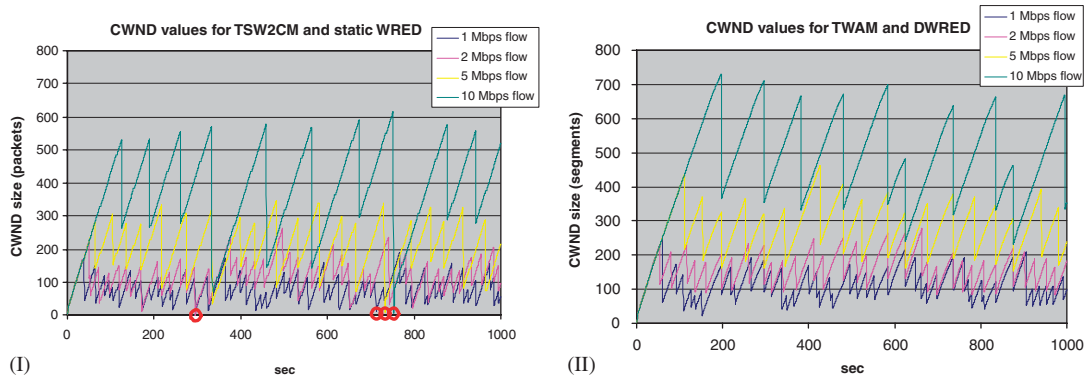


Figure 6. CWND measured for the two configurations of scenario #3.

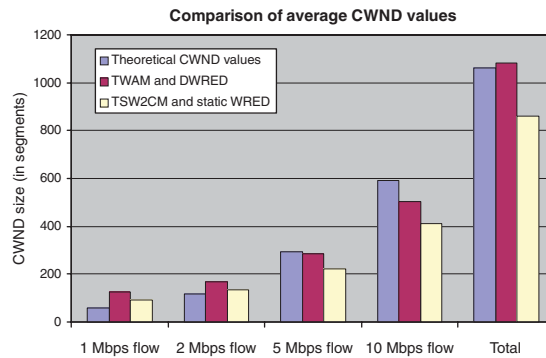


Figure 7. Average CWND achieved by TCP flows with different Relative class reserved capacities under scenario #3.

- Use of the TWAM marker in  $R_1$  to calculate the probability of marking incoming packets of a flow as in (for intervals according to the RTT measured at the specific LSP) and use of DWRED in each of the  $Q_{AF}$ s depicted.

In Figure 5 and Table II, the  $r_i$  values used by TSW2CM and the corresponding  $W_i^r$  values (according to (1)) used by the TWAM are presented for the group of the TCP flows entering  $R_1$ .

For this experimental set-up, Relative traffic is initially the only traffic over the path depicted. Each  $Q_{AF}$  is configured to serve Relative traffic with a rate equal to 30Mbps and thus, the  $Q_{AF}$  configuration provides the minimum resources required for in-profile traffic (see also (6)).

$$\sum_i r_i = 30\text{Mbps} \quad (13)$$

As can be seen from Figure 6, the TWAM+DWRED combination achieves the desired differentiation at the level of average congestion windows, while at the same time allowing all TCP flows to achieve larger average CWND sizes than those achieved by TSW2CM+static WRED. At the same time, no TCP flow experiences a slow start phase when TWAM+DWRED is used and CWND

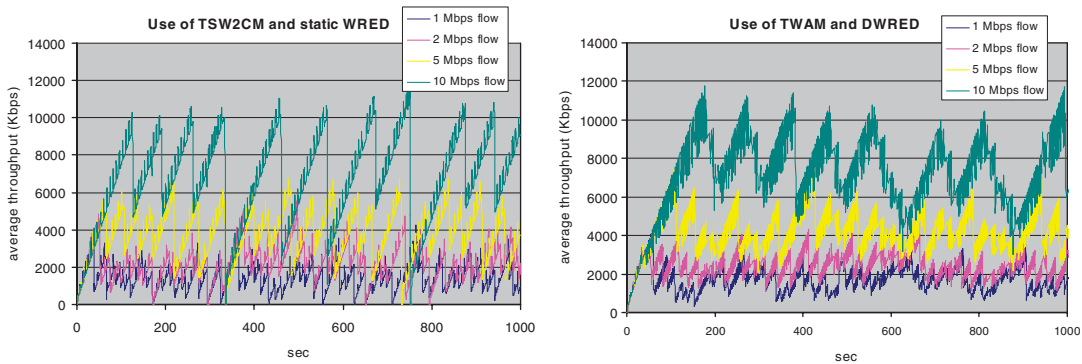


Figure 8. Instantaneous throughput achieved by TCP flows with different Relative class reserved capacities under scenario #3.

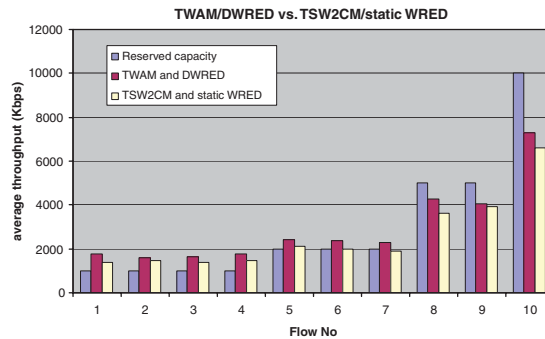


Figure 9. Average throughput achieved by TCP flows of Table II under scenario #3.

values for TCP flows are ‘stabilized’ within intervals proportional to each flows reserved  $W_i^r$  and thus  $r_i$ . In the case of TSW2CM+static WRED, CWND values achieved are less differentiated and flows are prone to slow start as the experiment evolves (see marked points in Figure 6(I)).

In Figure 7, it is shown how the proposed combination of TWAM+DWRED achieves average CWND values that better approximate the theoretical ones. When TWAM+DWRED are used, the sum of CWNDs achieved matches the bandwidth-delay product of the path and thus 100% utilization is achieved, unlike the TSW2CM+static WRED case. Besides better utilization, TWAM+DWRED achieve fairness by ensuring differentiation among flows that the Relative service requires.

In terms of average throughput, as shown in Figure 9 (while instantaneous throughput values are shown in Figure 8), the use of TWAM and DWRED outperforms TSW2CM and static WRED for each individual TCP flow. However, although Relative service is provisioned in a way that the sum of contracted service rates for all TCP flows equals the configured service rate  $C_{AF}$  along the AF FEC LSP (see (13)), both mechanisms provide excess bandwidth to low  $r_i$  TCP flows at the expense of high  $r_i$  TCP flows (e.g. the flow with  $r_i = 10$ Mbps). It is a topic of our future work to improve the Relative service provisioning mechanisms so as to achieve more strict differentiation.

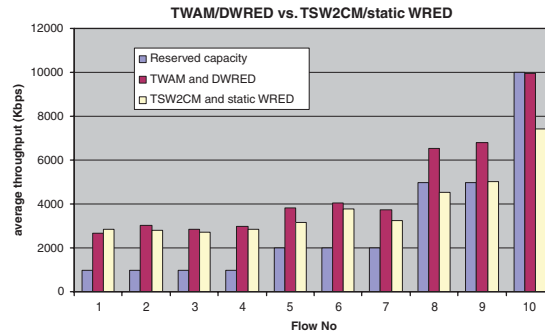


Figure 10. Average throughput of TCP flows for scenario #4: service rate in  $Q_{AF}$  is 50 Mbps.

For the next set of experiments, we have used over-provisioning for the set-up of Figure 5 with  $C_{R1,R2} = C_{R2,R3} = C_{R3,R4} = 50$  Mbps and a configuration of  $Q_{AF}$  to serve Relative traffic with a rate of 50 Mbps, which is higher than  $\sum_i r_i$  (**scenario #4**). Over-provisioning has ensured that all Relative flows achieve the contracted capacity under TWAM and DWRED (Figure 10). The limitations of TSW2CM and static WRED are once again apparent here since, despite of the over-provisioning, the larger  $r_i$  flows (flows 8, 9 and 10) fail to reach their contracted  $r_i$ . Thus, our proposed model when co-existing with some level of over-provisioning clearly outperforms TSW2CM and static WRED, while at the same time better approximates the optimal performance. It is important to mention here, the zero packet loss of in packets observed in the simulations, as an indisputable advantage of the Relative service model.

### 6.3. Evaluation under realistic operating conditions

After the initial evaluation experiments for TWAM and DWRED, the effectiveness of these mechanisms in Relative service performance will be demonstrated under more realistic conditions.

More specifically, two different scenarios are presented:

- In **scenario #5**, the set-up of Figure 5 is used with  $C_{R1,R2} = C_{R2,R3} = C_{R3,R4} = 50$  Mbps and  $C_{AF} = 50$  Mbps in all Relative traffic queues. However, flows no. 3,4,7 and 9 (Table II) start transmitting 200 s after the simulation starts. The purpose here is to compare the efficiency of TWAM and DWRED with that of TSW2CM and static WRED in fluctuations of Relative traffic. Ideally, flows starting to transmit later, when all other flows have fully opened up their TCP window to occupy available resources, should be able to quickly obtain their fair share of resources according to contracted Relative capacity. In equilibrium, flows with the same contracted capacity should be provided with the same share of resources.
- In **scenario #6**, the set-up of Figure 5 is used with  $C_{R1,R2} = C_{R2,R3} = C_{R3,R4} = 50$  Mbps. However, in router  $R_2$  a separate queue  $Q_{EF}$  serving traffic with absolute non-pre-emptive priority and a rate up to 10 Mbps is configured. During the simulation duration, UDP traffic at a rate of 11 Mbps (to ensure that  $Q_{EF}$  reaches its maximum serving capacity) is injected in  $R_2$  and exits the AF FEC LSP from  $R_3$  at regular intervals. This set-up aims to reproduce normal operating conditions in a backbone router and investigate how transient load in other traffic classes affects the Relative service fairness and the capacity guarantees provided.

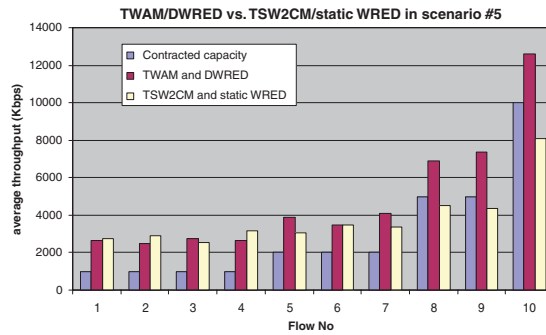


Figure 11. Average throughput of TCP flows for scenario #5: flows 3, 4, 7 and 9 start transmitting 200 s later, but obtain their share of resources when TWAM and DWRED are used.

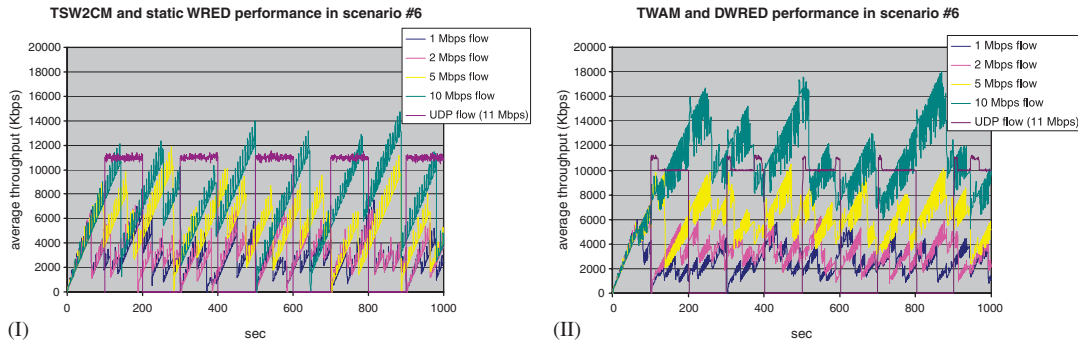


Figure 12. Instantaneous throughput achieved by TCP flows with different Relative class reserved capacities under scenario #6.

Figure 11 shows the average throughput of all TCP flows measured in the interval from  $t_1=400$  to 1000 s for scenario #5. Use of TWAM and DWRED ensures that the differentiation between TCP flows is preserved in a fair way. Measured average throughput is higher than the contracted capacity for all TCP flows due to over-provisioning ( $C_{AF}=50$  Mbps). However, when TSW2CM and static WRED are used flows no. 8, 9 and 10 fail to obtain their contracted rate despite of the over-provisioning conditions. Moreover, differentiation among TCP flows is not fair.

In Figure 12, the performance advantages of TWAM and DWRED become clearly evident in the framework of scenario #6. High-priority UDP traffic is only transmitted in intervals  $\{[100, 200], [300, 400], [500, 600], [700, 800], [900, 1000]\}$  s. Despite the periodic load of 110% in the high-priority class, TCP flows succeed in obtaining and preserving a throughput proportional to their reserved capacity.

Throughput is distributed in a fair way among TCP flows and the average throughput measured for each flow is above the guaranteed value despite the fluctuations in high-priority traffic load. Again packet losses for in packets are zero. On the contrary, the use of the TSW2CM and static WRED (Figure 12(II)) demonstrates a poor performance.

Despite the minimum of 40 Mbps service rate available to Relative traffic along the AF FEC LSP (since the high-priority traffic is limited to 10 Mbps service rate), the TCP flows are unable to preserve their contracted capacity. Differentiation among them is also very poor, demonstrating thus, the inefficiency of the TSW2CM and static WRED configuration to support an AF-based service under realistic operation conditions.

#### 6.4. Update interval

In Sections 4 and 5, both TWAM and DWRED are proposed to operate at intervals of  $\alpha \cdot \text{RTT}_{\text{LSP}_j}$  seconds (with  $\alpha \in \{1, 2, 3, \dots\}$ ). The optimal value of  $\alpha$  has been the subject of an experimental study, the results of which are shown in Figure 13.

For small values of  $\alpha$ , the highest utilization is achieved (highest sum of average throughput) as well as better differentiation among TCP flows. As the operating intervals increase ( $\alpha > 2$ ), performance deteriorates in both perspectives. For the experiments presented in this work, the value of  $\alpha = 2$  has been used. More detailed study on this will be part of our future work; however, at this stage it is obvious that a trade-off exists between more brief operating intervals (and thus higher overhead) and less optimal results.

#### 6.5. Queuing behaviour

In this section, further insight to the implications of using TWAM and DWRED for Relative service provisioning is provided. Owing to the adaptive nature of DWRED, which allows for TCP flows to expand their congestion windows up to the level of available resources, the use of DWRED is associated with a higher level of router queue occupancy than static WRED. As can be seen from Figure 14, for scenario #5 as presented in Section 6.3, when static WRED is used, the Relative queue  $Q_{\text{AF}}$  has a very low average size of 3 packets and for 99% of the time its size varies from 1 to 44 packets (close to  $\text{min}_{\text{th}} = 40$  packets). On the contrary, during the use of DWRED, the Relative queue size varies from 1 to 490 for 99% of the time (close to  $\text{min}_{\text{th}} = 0.15 \cdot \text{max}_{\text{th}} = 0.15 \cdot 2930 = 440$  packets). This is reasonable since  $\text{min}_{\text{th}}$  and  $\text{max}_{\text{th}}$  in the case of DWRED are tuned according to the  $W_{\text{LSP}_j}^{\text{wnd}}$  of scenario #5 (expressions (11) and (12)).

Thus, the advantages of DWRED come with some cost in queuing delay. The resulting average queue size for the generic application of DWRED was shown in [2] not to exceed  $\{0.2 \dots 0.25\} \cdot (\text{BW} \cdot D)$ . In Figure 14(II), the average queue size is measured at 360 packets, thus a 12.3% of  $W_{\text{LSP}_j}^{\text{wnd}}$ . Experiments performed with a single TCP flow over the setting of scenario #5 have resulted in an average queue size of 580 packets, thus a 20% of  $W_{\text{LSP}_j}^{\text{wnd}}$ .

As a conclusion, the application of DWRED results in increased queuing delays when compared with static WRED configurations due to higher average queue sizes. However, the average queue size under DWRED seems to be bounded by a percentage around 20% of the value of  $\text{max}_{\text{th}}$ . As a result, the average queuing delay perceived by TCP traffic under DWRED is also bounded, although higher than that achieved with a static WRED configuration. This is confirmed by Figure 15, depicting the measured RTT for a TCP flow in scenario #5, when DWRED and static WRED are used. Instantaneous RTT for all TCP flows was measured to be identical to the values depicted here, in both configurations.

RTT in the case of DWRED is quite higher but stays below 290 ms with a few exceptions. The higher utilization and the improved differentiation achieved when DWRED is combined with TWAM (see Figures 11 and 12) together with the fact that the excess queuing delay imposed by

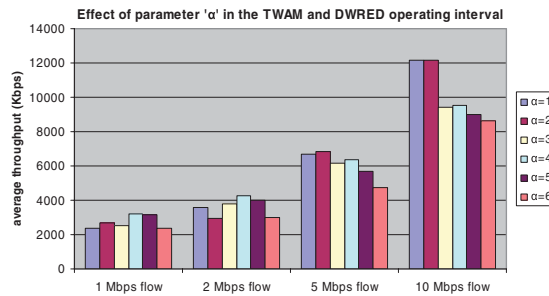


Figure 13. Average throughput achieved by TCP flows for different operating intervals of TWAM and DWRED applied on scenario #5.

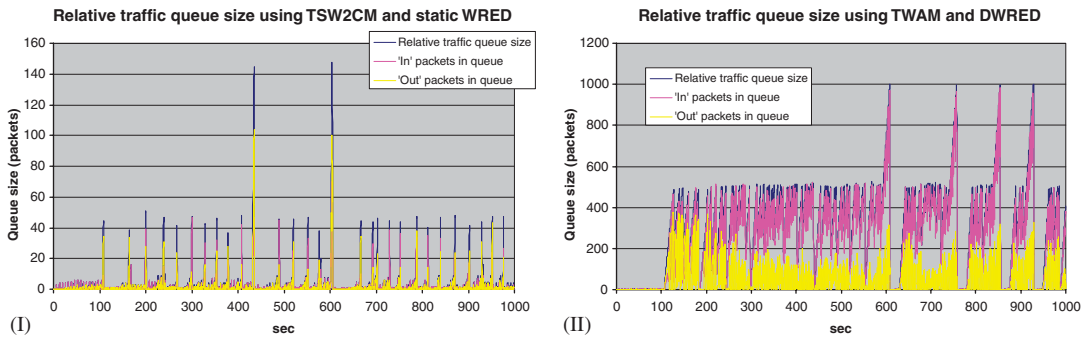


Figure 14. Relative traffic queue size for scenario #1: comparison of TSW2CM/static WRED and TWAM/DWRED configurations.

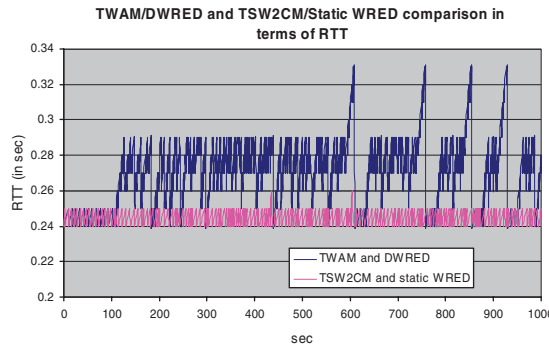


Figure 15. RTT measured for TCP flows participating in scenario #5 under (i) the TWAM/DWRED and (ii) the TSW2CM/static WRED configurations.

DWRED is bounded and predictable, make the proposed mechanisms quite valuable tools for the provisioning of the Relative service.

### 6.6. Processing requirements of the proposed mechanisms

In the case of TWAM and DWRED, the following set of computations are required:

- For each operating interval, thus every  $\alpha \cdot \text{RTT}_{\text{LSP}_j}$  seconds (with  $a \in \{1, 2, 3, \dots\}$ ), for each TCP flow, the TWAM in the IR computes (9).
- Similarly, every  $\alpha \cdot \text{RTT}_{\text{LSP}_j}$  seconds (with  $a \in \{1, 2, 3, \dots\}$ ), in each core router, the DWRED parameters for the Relative traffic queue  $Q_{\text{AF}}$  are calculated according to (11) and (12).

Without the need for detailed analysis, it is obvious that the computational complexity of TWAM is much lower than marking schemes such as that of TSW2CM requiring a set of computations to be made upon each packet arrival at the IR. DWRED imposes some overhead when compared with static WRED configurations; however, this overhead only depends on the update intervals chosen and scales with the number of TCP flows, the topology, etc.

## 7. IMPLEMENTATION ISSUES

After the analytical and experimental evaluation of the proposed mechanisms, a number of issues that concern the implementation and provisioning of an AF service based on TWAM and WRED have to be considered.

The principles of operation of the proposed TWAM and DWRED mechanisms rely on the following strong assumptions:

- All the TCP flows belonging to the same  $\text{LSP}_{\text{AF}}$  across an AF-enabled backbone update their CWND values in the same intervals.
- All the TCP flows served by the same AF service router queue  $Q_{\text{AF}}$  in a core router update their CWND values in the same intervals.

The first assumption comes as a natural consequence of the traffic engineering properties of MPLS, forcing each packet belonging to the same AF FEC LSP to follow the same route across a network and thus, experience the same end-to-end delay. By using AF FEC LSPs for TCP traffic and the corresponding acknowledgments between two network edges, we can ensure that all TCP flows belonging to the same  $\text{LSP}_{\text{AF}}$  experience the same RTT and thus, update their CWND values in the same intervals. However, the use of MPLS LSPs for AF traffic when the Relative service co-exists with other services, does not impose any requirements on the use of LSPs for these services.

The second assumption is more limiting in terms of implementation. It does not allow to serve AF LSPs associated with different RTT values through a single router queue of a core router. Intuitively, for flows participating in LSPs with smaller RTT values, this would mean that within an operating interval of the TWAM and DWRED they are allowed to increase their CWND more times than flows in LSPs with larger RTTs can. Or in other words, aggregates with high RTTs take longer to ramp up after a packet drop occurs. This can lead to unfairness in throughput distribution, shown by Figure 16 where two flows with the same contracted  $r_i$  but different RTTs



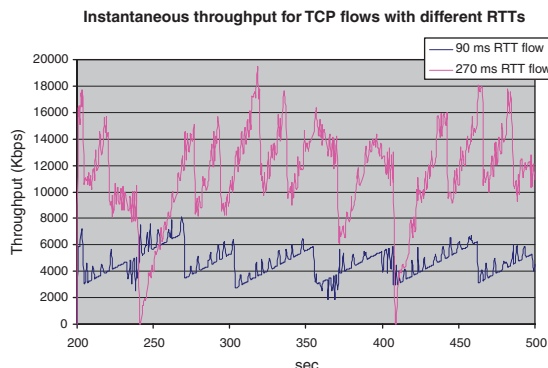


Figure 16. Using TWAM and DWRED to serve TCP flows experiencing different RTTs by the same  $Q_{AF}$  in a core router.

obtain different throughput. It is obvious from Figure 16 how the smaller RTT flow is able to increase its CWND in a higher rate than the larger RTT flow.

It has to be noted here that unfairness in throughput distribution between TCP flows experiencing different RTTs is not a limitation of the TWAM and DWRED as proposed here. No scalable solutions to the problem of unfairness between TCP flows experiencing different RTTs under AF service provisioning exist. However, under the Relative service provisioning model proposed in this work, this unfairness can be avoided if a suitable queuing scheme, such as WFQ, is used for AF traffic in queue routers. Using WFQ on each core router, a guaranteed minimum service rate  $C_j^{AF}$  can be obtained for each AF FEC LSP  $LSP_j^{AF}$ . Ensuring a minimum  $C_j^{AF}$  service rate for each  $LSP_j^{AF}$  avoids the case where TCP flows with differing RTT value compete for the same resources. Thus, the theoretical model of Section 4, based on a single  $LSP_j^{AF}$ , can be applied for multiple  $LSP_j^{AF}$  in a real network scenario.

The Relative service implementation requires that the traffic carried by AF FEC LSPs configured to carry TCP ACKs from ER back to IR should be treated with topmost priority in the core routers (e.g. using an EF PHB).

A last but not least implementation issue is that of obtaining the input parameters required for the operation of TWAM and DWRED mainly of the values of each flow  $W_i^r$  and  $W_i^{cwnd}$  over sequential operation intervals, required in (9), (11) and (12) thus, for TWAM and DWRED tuning. Maintaining per-flow or TCP aggregate information such as  $W_i^r$  and  $W_i^{cwnd}$  values for the calculation of  $p_i$  at the IR, where the TWAM operation is completely in accordance with the DiffServ principles.  $W_i^r$  values occur according to the contracted  $r_i$  for each flow and the path data of each AF LSP.  $W_i^{cwnd}$  values can be signalled by TCP flows to the IR either constantly or every time the TCP flow CWND is updated. TWAM has thus all the data it needs to operate locally at the IR. The only additional requirement of the TWAM and DWRED schema is to signal core routers with an  $\sum_i W_i^{cwnd}$  value per AF LSP in every operating interval, so that DWRED can tune itself according to (11) and (12). Again here in-band signalling or the TCP receiver-advertised window can be used and scalability is not endangered since a single value is required for each AF LSP aggregating multiple flows. This makes the amount of state information kept in the core routers constant and independent of the number of individual flows.

## 8. FUTURE WORK

Despite the several advantages of the Relative service model and of the individual AF PHB mechanisms proposed in this work, there are a number of issues that require further study and will be part of our future work on this topic.

We intend to investigate refinements of the TWAM principles of operation, such as optimization of the length of the operating intervals and of the estimation of the  $W_{LSP_j}^{cwnd}$  for each operating interval. Although TWAM and DWRED were shown to improve the differentiation among TCP flows, better differentiations are always a topic of interest. More experimentation is also required on how do the proposed mechanisms performs under diverse scenarios and different flow aggregation levels. There are indications that the utilization achieved decreases as the number of aggregated flows increases and part of our future work will be to look into this issue.

Alternative implementations of DWRED will also be evaluated, such as calculating the values of DWRED based on average queue sizes for in packets only and/or for the total of packets in  $Q_{AF}$ .

It is anticipated that further refinements in the proposed model will succeed in lowering the level of over-provisioning required for the Relative service rate guarantees to be honored under diverse scenarios. It is part of our future work to investigate this possibility also.

## 9. CONCLUSIONS

In this paper we have proposed and evaluated two mechanisms for the AF service PHB implementation that can be used in the framework of the Relative service model, as this was specified in our previous work. Our main target was to overcome a number of limitations in AF services' implementation and provide a deployable, predictable and successful service for TCP traffic.

According to related research work, among the factors that affect/obstruct the provision of quantitative guarantees to TCP traffic are the effects of RTT, the self-clocked sliding window mechanism, the number of concurrent flows in an AF class, the variety in packet sizes and the interference between the TCP and UDP traffic. Our Relative service model and the proposed TWAM and DWRED mechanisms address most of these issues, while at the same time providing specific configuration guidelines for implementing such a service.

The TWAM ensures the distribution of available resources among TCP flows within an AF class in a fair manner, according to their contracted average service rates. Its efficiency is reinforced by the DWRED and in particular by the self-tuning capability of DWRED, since that WRED parameters are dynamically adjusted to the current perceived load at the TWAM. Both mechanisms require much less overhead than existing equivalent mechanisms and achieve improved differentiation, fairness, adaptation to transient network conditions and high utilization of available resources. TCP flows under the TWAM and DWRED operation demonstrate a controlled behaviour, with smooth adaptation to network conditions, remaining at the congestion avoidance phase and avoiding the performance handicaps that slow start entails. The assured average service rate of the AF service provisioning model is thus achieved, even over small time scales.

Provided that the proposed mechanisms are implemented by routing equipment, the results of our work provide a number of promising indications for the successful deployment of AF-based services in IP networks.

## APPENDIX

A set of AQM configurations found in the literature is provided in the table below:

---

[20]	RIO with $\min_{th}^{In}=40$ , $\max_{th}^{In}=100$ , $\max_p^{In}=0.02$ $\min_{th}^{Out}=20$ , $\max_{th}^{Out}=40$ , $\max_p^{Out}=0.5$
[8]	RIO with $\min_{th}^{In}=40$ , $\max_{th}^{In}=70$ , $\max_p^{In}=0.2$ $\min_{th}^{Out}=10$ , $\max_{th}^{Out}=30$ , $\max_p^{Out}=0.2$ Or $\min_{th}^{In}=\min_{th}^{Out}=50$ , $\max_{th}^{In}=\max_{th}^{Out}=100$ , $\max_p^{In}=0.1$ , $\max_p^{Out}=0.9$
[9]	RIO with $\min_{th}^{In}=40$ , $\max_{th}^{In}=70$ , $\max_p^{In}=0.02$ $\min_{th}^{Out}=10$ , $\max_{th}^{Out}=30$ , $\max_p^{Out}=0.5$ RED with $\min_{th}=15$ , $\max_{th}=40$ , $\max_p=0.02$ RED with $\min_{th}=10$ , $\max_{th}=30$ , $\max_p=0.02$
[13]	RED with $\min_{th}=BS/4$ , $\max_{th}=BS/2$ , $\max_p=0.02$ , where BS is the queue buffer size
[19]	RED with $\min_{th}=20\text{Kb}$ , $\max_{th}=80\text{Kb}$ , $\max_p=\{variable\}$
[12]	RIO with $\min_{th}^{In}=20$ , $\max_{th}^{In}=80$ , $\max_p^{In}=0.1$ $\min_{th}^{Out}=10$ , $\max_{th}^{Out}=30$ , $\max_p^{Out}=0.1$
[16]	RED with $\min_{th}=30$ , $\max_{th}=130$ , $\max_p=0.1$
[25]	RED with $\min_{th}=20\text{Kb}$ , $\max_{th}=80\text{Kb}$ , $\max_p=0.02$
[26]	RIO with $\min_{th}^{In}=150$ , $\max_{th}^{In}=200$ , $\max_p^{In}=0.05$ $\min_{th}^{Out}=1$ , $\max_{th}^{Out}=20$ , $\max_p^{Out}=0.2$

---

## REFERENCES

1. Heinanen J, Baker F, Weiss W, Wroclawski J. Assured Forwarding PHB Group. *RFC 2597*, 1999.
2. Bouras C, Sevasti A. Performance analysis for a DiffServ-enabled network: the case of relative service. *Proceedings of the 2nd IEEE International Symposium on Network Computing and Applications (NCA-03)*, Cambridge, MA, U.S.A., 16–18 April 2003; 381–388.
3. Bouras C, Sevasti A. Performance enhancement of an AF service using TCP-aware marking and dynamic WRED. *10th IEEE Symposium on Computers and Communications, LaManga del Mar Menor*, Cartagena, Spain, 27–30 June 2005; 624–647.
4. Medina O, Orozco J, Rosaf D. Bandwidth sharing under the assured forwarding PHB. *Internal Report, IRISA No. 1478*, September 2002.
5. Goyal M, Durresi A, Liu C, Jain R. Performance analysis of assured forwarding. *Internet Draft (draft-goyal-diffserv-afstdy-00)*, Internet Engineering Task Force, February 2000.
6. Nandy B, Seddigh N, Piedad P. Diffserv's assured forwarding PHB: what assurance does the customer have? *Proceedings of the 9th International Workshop on Network and Operating Systems Support for Digital Audio and Video (NOSSDAV'99)*, New Jersey, NJ, U.S.A., June 1999.
7. Seddigh N, Nandy B, Piedad P. Study of TCP and UDP Interaction for AF PHB. *Internet Draft (draft-nsbnpp-diffserv-udptcpaf-01.pdf)*, Internet Engineering Task Force, August 1999.
8. May M, Bolot J-C, Jean-Marie A, Diot C. Simple performance models of tagging schemes for service differentiation in the Internet. *Proceedings of IEEE INFOCOM 1999*, New York, NY, U.S.A., March 1999; 1385–1394.
9. Clark D, Feng W. Explicit allocation of best effort packet delivery service. *IEEE/ACM Transactions on Networking* 1998; **6**(4):362–374.

10. Seddigh N, Nandy B, Heinanen J. An assured rate per-domain behaviour for differentiated services. *Internet Draft (draft-ietf-diffserv-pdb-ar-00.txt)*, Internet Engineering Task Force, 2000.
11. Azeem F, Rao A, Lu X, Kalyanaraman S. TCP-friendly traffic conditioners for differentiated services. *Internet Draft (draftazeem-tcpfriendly-diffserv-00.txt)*, Internet Engineering Task Force, February 1999.
12. Athuraliya S, Li VH, Low SH, Yin Q. REM: active queue management. *IEEE Network*, vol. 15, May/June 2001; 48–53.
13. Lin D, Morris R. Dynamics of random early detection. *Proceedings of ACM SIGCOMM '97*, Cannes, French Riviera, France, October 1997; 127–137.
14. Hurley P. The provision of a low delay service within the best-effort Internet. *Ph.D. Thesis, Number 2516, EPFL-I&C-LCA*, Lausanne, Switzerland, December 2001.
15. Hollot CV, Misra V, Towsley D, Gong W. A control theoretic analysis of RED. *Proceedings of IEEE INFOCOM 2001*, April 2001, Anchorage, AK, 2001; 1510–1519.
16. Iannaccone G, May M, Diot C. Aggregate traffic performance with active queue management and drop from tail. *Computer Communication Review* 2001; 31(3):4–13.
17. Banchs A, Tartarelli S, Orlandi F, Sato S, Kobayashi K, Pan H. Configuration of DiffServ routers for high-speed links. *Proceedings of Workshop on High Performance Switching and Routing (HPSR 2002)*, Kobe, Japan, 26–29 May 2002.
18. Sahu S, Towsley D, Kurose J. A quantitative study of differentiated services for the Internet. *Proceedings of IEEE Global Internet, Globecom'99*, Rio de Janeiro, Brazil, December 1999; 1808–1817.
19. Feng W, Kandlur D, Saha D, Shin K. A self-configuring RED gateway. *Proceedings of IEEE INFOCOM 1999*, New York, NY, U.S.A., 1999; 1320–1328.
20. Yeom I, Reddy N. Impact of marking strategy on aggregated flows in a differentiated services network. *Proceedings of IWQoS Workshop*, June 1999, London, U.K., 1999; 156–158.
21. Matsuda T, Nagata A, Yamamoto M. TCP rate control using active ECN mechanism with RTT-based marking probability. *Proceedings of the 16th International Workshop on Communications Quality and Reliability (CQR 2002)*, Okinawa, May 2002; 112–116.
22. Kapoor R, Casetti C, Gerla M. Core-stateless fair bandwidth allocation for TCP flows. *Proceedings of IEEE ICC2001*, Helsinki, Finland, June 2001; 146–150.
23. Mellia M, Stoica I, Zhang H. TCP-aware packet marking in networks with DiffServ support. *The International Journal of Computer and Telecommunications Networking* 2003; 42(1):81–100.
24. Hasegawa G, Murata M. Analysis of dynamic behaviors of many TCP connections sharing tail-drop/RED routers. *Proceedings of Global Telecommunications Conference '01 (GLOBECOM '01)*, vol. 3(25–29), 2001; 1811–1815.
25. Feng W, Kandlur D, Saha D, Shin K. Understanding TCP dynamics in a differentiated services Internet. *IEEE/ACM Transactions on Networking* 1999; 7:173–187.
26. Azeem F, Rao A, Kalyanaraman S. TCP-friendly traffic marker for IP differentiated services. *Proceedings IWQoS'2000*, Pittsburgh, PA, June 2000; 35–48.
27. McCanne S, Floyd S. ns network simulator. Available from: <http://www.isi.edu/nsnam/ns/>.

#### AUTHORS' BIOGRAPHIES



**Christos Bouras** obtained his Diploma and PhD from the Computer Science and Engineering Department of Patras University (Greece). He is currently Professor in the above department. Also he is a scientific advisor of Research Unit 6 in Research Academic Computer Technology Institute (CTI), Patras, Greece. His research interests include Analysis of Performance of Networking and Computer Systems, Computer Networks and Protocols, Telematics and New Services, QoS and Pricing for Networks and Services, e learning, Networked Virtual Environments and WWW Issues. He has extended professional experience in Design and Analysis of Networks, Protocols, Telematics and New Services. He has published 300 papers in various well-known refereed conferences and journals. He is a co-author of 8 books in Greek. He has been a PC member and referee in various international journals and conferences. He has participated in R&D projects such as RACE, ESPRIT, TELEMATICS, EDUCATIONAL MULTIMEDIA, ISPO, EMPLOYMENT, ADAPT, STRIDE, EUROFORM, IST, GROWTH and others. Also he is member of experts in the

Greek Research and Technology Network (GRNET), Advisory Committee Member to the World Wide Web Consortium (W3C), IEEE-CS Technical Committee on Learning Technologies, IEEE ComSoc Radio Communications Committee, IASTED Technical Committee on Education W 6.4 Internet Applications Engineering of IFIP, ACM, IEEE, EDEN, AACE, New York Academy of Sciences and Technical Chamber of Greece.



**Dr Afrodite Sevasti** holds the position of Services and Infrastructure Development Coordinator at the Greek Research and Technology Network (GRNET) S.A. She is also a member of the GEANT2 project Technical Committee, coordinating the research activity on Bandwidth on Demand. She has participated in a number of R&D projects and currently leads the Standardization and Liaisons activity of the FEDERICA project (EU FP7). Afrodite received her PhD in the field of QoS enhancements on IP-based networks from the Computer Engineering and Informatics Department, Engineering School of Patras University (Greece). She also holds an MSc degree in Information Networking from Carnegie Mellon University. She has published 30 papers in refereed conferences and journals.