

QoS Architectures for the IP Network

N

Harry G. Perros

North Carolina State University, USA

INTRODUCTION

When we call someone over the Internet using a service such as Skype or Google talk, we may experience certain undesirable problems. For instance, we may not be able to hear the other person very well, or even worse, the call may be dropped. This is in contrast to using the regular telephone system where the quality of the voice is always very good. Similarly, during a conversational video call, the picture may freeze, or there may be pixilation, or the call may be dropped. The reason for these problems is that the IP packets that carry the contents of our call are not delivered on time at the destination so that they can be played out at the right time. Also, some of the packets may be lost while they are traversing the Internet. In order to eliminate these problems, the underlying IP network has to be able to provide Quality of Service (QoS) guarantees. Several schemes have been developed that enable the IP network to provide such guarantees. Of these schemes, the Multi-Protocol Label Switching (MPLS) and the Differentiated Services (DiffServ) are the most widely used. In this article, some of the salient features of MPLS and DiffServ are reviewed.

BACKGROUND

QoS is a well-understood and studied topic within the networking community. It is typically expressed in term of the following three metrics: the end-to-end delay, the jitter, and the packet loss rate. The end-to-end delay is the amount of time it takes to transfer a packet from the transmitter to the receiver, and it consists of a) the end-to-end propagation delay, b) delays induced by transmission systems and processing times inside the routers, and c) delays a packet encounters due to queueing in the buffers of the routers. Jitter refers to the variability of the inter-arrival times of the packets at the destination, and the packet loss rate is the percent of packets that are lost.

Different applications have different tolerance to these QoS metrics. Table 1 relates various common networking applications to the end-to-end delay and packet loss rate. For instance, for conversational voice and video it is important that packets should be delivered to the destination in less than 150 msec in order to maintain user satisfaction. (Studies have showed that in fact an end-to-end delay of up to 220 msec can be tolerated.) On the other hand, a packet loss rate of about

Table 1. QoS metrics for common networking services

Tolerance for packet loss	<i>Tolerant</i>	Conversational voice and video	Voicemail	Streaming audio and video	Fax
	<i>Intolerant</i>	Remote app., command and control games	e-commerce web browsing	Texting, file transfer (foreground)	File transfer (background), email
		<i>Interactive delay</i> < 1 s	<i>Responsive delay</i> ~ 1 s	<i>Timely delay</i> ~ 10 s	<i>Background delay</i> >> 10 s
		Tolerance for delay			

DOI: 10.4018/978-1-5225-2255-3.ch573

1 in 100 can be tolerated. That is, conversational voice and video type of applications are packet-loss tolerant but they have a strict end-to-end delay constraint, i.e. they are delay intolerant. On the other hand, a file transfer service is delay tolerant but packet-loss intolerant. This is because we do not expect a file to be delivered immediately, but the integrity of the file is important, and any lost packets have to be re-transmitted.

In view of the above, the question that arises is how can the network provide different QoS to different applications. In order to answer this question, let us first take a look at how the IP network routes packets. Each IP packet consists of a header and a payload, and the header contains different fields one of which is the destination IP address. When a packet arrives at an IP router, the header is examined and the destination address is used in a forwarding routing table in order to find out the next IP router to which the IP packet should be sent. This forwarding operation is carried out at each router along the path followed by the packet, until the packet reaches its destination. The forwarding routing table in each IP router is constructed using a routing protocol, such as the Open Shortest Path First (OSPF). The path that a packet follows is the shortest path in terms of the number of hops, i.e., routers. The advantage of this type of routing is that it is simple. However, since it minimizes the number of hops, it is difficult to guarantee any QoS metrics, such as end-to-end delay, jitter, and packet loss rate. For instance, the fact that the path that a packet follows has the smallest number of hops, does not necessarily mean that it has the shortest end-to-end delay. On the other hand, if all routers have approximately the same packet loss rate, then the shortest path will result to the lowest end-to-end packet loss rate.

Another problem is that a router cannot distinguish packets without an additional mechanism, such as, packet inspection or packet classification. Therefore, it cannot give packets from delay-intolerant applications a higher priority for transmission out of an output port over packets from delay-tolerant applications. (This is necessary, if we want to minimize the end-to-end delay of packets from

a delay-intolerant application.) In view of this, the only way that delay sensitive applications can be served satisfactorily is to under-utilize the entire network so that the queues of packets waiting for transmission in the output ports of the routers are never too long, and as a result, the delay to transmit a packet out of an output port is negligible. This is known as *over-engineering* the network. This solution is expensive since the links are under-utilized, and also it does not prevent the occurrence of transient traffic congestion. An advantage of over-engineering the network, is that when a link failure occurs, traffic can be redirected over other links without saturating them.

In order to provide QoS guarantees in an IP network without having to operate it at very low utilization, we need a scheme that can carry out *call admission control* and *packet classification*:

- **Call Admission Control:** Before a user or an application starts transmitting packets over the network, we have to make sure that the network has the necessary bandwidth to carry the new flow of packets that will be generated at the expected QoS, without affecting the QoS of the flows of packets that are currently being transmitted over the network. This procedure is known as call admission control;
- **Packet Classification:** Packets should be classified to different classes of service with different priorities, so that they are transmitted out of the output port of a router according to the priority of their class of service.

Several schemes have been developed that enable the IP network to provide QoS guarantees. Of these schemes, the Multi-Protocol Label Switching (MPLS), see IETF RFC 3031 (2001), and the Differentiated Services (DiffServ), see IETF RFC 2474 (1998), are the most widely used. Below, some of the salient features of MPLS and DiffServ are reviewed. For further details the reader is referred to Perros (2014).

7 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage:

www.igi-global.com/chapter/qos-architectures-for-the-ip-network/184355

Related Content

Unmanned Bicycle Balance Control Based on Tunicate Swarm Algorithm Optimized BP Neural Network PID

Yun Li, Yufei Wu, Xiaohui Zhang, Xinglin Tan and Wei Zhou (2023). *International Journal of Information Technologies and Systems Approach* (pp. 1-16).

www.irma-international.org/article/unmanned-bicycle-balance-control-based-on-tunicate-swarm-algorithm-optimized-bp-neural-network-pid/324718

Metadata Diversity in the Cultural Heritage Repositories

Sumeer Gul, Shahkar Riyaz Tramboo and Humma Ahangar (2015). *Encyclopedia of Information Science and Technology, Third Edition* (pp. 1843-1854).

www.irma-international.org/chapter/metadata-diversity-in-the-cultural-heritage-repositories/112590

Experiment Study and Industrial Application of Slotted Bluff-Body Burner Applied to Deep Peak Regulation

Tianlong Wang, Chaoyang Wang, Zhiqiang Liu, Shuai Ma and Huibo Yan (2024). *International Journal of Information Technologies and Systems Approach* (pp. 1-15).

www.irma-international.org/article/experiment-study-and-industrial-application-of-slotted-bluff-body-burner-applied-to-deep-peak-regulation/332411

Digital Literacy in Theory and Practice

Heidi Julien (2018). *Encyclopedia of Information Science and Technology, Fourth Edition* (pp. 2243-2252).

www.irma-international.org/chapter/digital-literacy-in-theory-and-practice/183937

Sociotechnical Change Perspective for Enterprise Resource Planning System Implementation

Jessy Nair, D. Bhanu Sree Reddy and Anand A. Samuel (2015). *Encyclopedia of Information Science and Technology, Third Edition* (pp. 923-933).

www.irma-international.org/chapter/sociotechnical-change-perspective-for-enterprise-resource-planning-system-implementation/112485