# A fine-grained response time analysis technique in heterogeneous environments

A. Hafsaoui[a], A. Dandoush[b,*], G. Urvoy-Keller[c], M. Siekkinen[d], D. Collange[e]

[a] TENEDIS, 15 avenue du Hoggar - Parc Victoria – Le Vancouver, 91940 LES ULIS, France
[b] ESME-Sudria, Paris sud, 38 rue Moliere, Ivry sur Seine, 94200, France
[c] Laboratoire I3S CNRS/UNS UMR Sophia-Antipolis 7271 06903, France
[d] Aalto University School of Science and Technology, Finland
[e] Orange Labs, Immeuble AGORA 905 rue Albert Einstein, Sophia-Antipolis 06921, France

## ARTICLE INFO

## ABSTRACT

It is crucial for the network operators and Internet service providers (ISPs) to determine the reasons that cause large response time fluctuations. In this paper, we consider passive measurements from heterogeneous environment (ADSL, FTTH and 3G/3G+ access technologies) of an European ISP 'Orange'. Through experimental analysis of real traces, the need of a fine-grained traffic analysis technique is demonstrated. We show that finding the root causes of the observed poor performance using simple metrics such as response time, RTT and packet loss is difficult. In view of this fact, the different factors that play a role in determining the resulting response time are described through examples. Then, a breakdown method that drills down into the passively observed TCP connections is proposed. The method decomposes the end-to-end response time into many time periods and maps each one to a specific parameter or a physical phenomenon. Thus, the impact of not only the network parameters but also the application configuration and user behavior is captured. The resulting time periods are given as input to a clustering algorithm in order to group together transfers with similar performance holding traffic of different application protocols over different access technologies. As a result, the contribution of each participant in the performance bottleneck is identified. The proposed technique is validated through extensive simulations and real passively measured traces and it is compared to other works. Exemplifying the technique on real traces from Internet and enterprise traffic is introduced and discussed to demonstrate the power of the approach and its simplicity. In contrast to some existing tools, ISPs and enterprise administrators do not need to modify their network architecture or to install a new software or a plugin at the client or at the server side in order to use our technique. In addition, data sampling is not used. This is particularly important in order to keep data consistency and to detect metrics peaks. Last, our tool deals with both long and short TCP connections.

## 1. Introduction

Internet Service Providers (ISPs) have a continuous need to measure the offered services and to enhance the performance perceived by their customers. In fact, poor performance does not always mean that the network is to be blamed. This fact particularly makes application-level performance monitoring problematic for the ISPs. Recent works (e.g. [1,2]), have shown that the Quality of Service (QoS) measures like packet loss and throughput do not indicate for many of today's applications anything about the cause of the poor performance perceived by users. Shorter the response time is, better the performance is for most of the nowadays applications. For example, it is shown in [3,4] that the main Web Performance Bottleneck is latency and not the bandwidth capacity. Therefore, studying the reasons of large response time is becoming increasingly important. In general, the poor performance can be due to one or more of the following factors: (i) the applications behavior at the servers and/or the clients side such as throttling sending rates, (ii) the congestion control mechanism of TCP protocol, (iii) and the heterogeneity of access technologies, namely ADSL, FTTH, Cellular and legacy Ethernet.

The purpose of this work, is to present a simple yet efficient visibility solution that can diagnose and troubleshoot the performance of TCP-based services. In fact, a key feature of our time au-

* Corresponding author.
 *E-mail address:* dandoush@esme.fr (A. Dandoush).

dit technique is the ability to divide a TCP connection into certain slots of time using a break-down approach in order to eliminate application response time problems. It is able to deliver for each TCP-based service the classical QoS indicators (cumulative distribution function CDF of packet loss, RTT as well as of throughput). Moreover, it provides a fine grained analysis of response time for capturing the root causes of the poor performance (e.g. application impact, user or server behavior, network problems etc.). Another important feature of our Time Audit technique is the transparency of its use in complex WAN and LAN architectures as it simply uses passively collected traces from any point in the network. The main motivation is that the passive analysis does not have the overhead that active monitoring has. In addition, data sampling is not used as in some active or real-time approaches. This is particularly important in order to keep data consistency and to detect metrics peaks due to the presence of the whole traffic for a given period. The active tools such as NetFlow [5] introduced by CISCO are Router Based Analysis Techniques. These techniques allow granular on-time traffic measurements as well as high-level aggregated traffic collection that can assist in identifying excessive bandwidth utilization or unexpected application traffic. It helps the network admin to understand what is happening in general in his/her network. However, these Techniques are hard-coded into the network devices, i.e. routers. It uses usually a slice of the network capacity for sending continuously statistics or full copy of data for every packets to an external server for doing the analysis. Thus more expensive resources at the routers, core links, and external computing unit are required. Thus, using our approach does not require to modify the networking infrastructure or to install additional software neither at clients nor at servers sides contrary of most network and application performance solutions such as [6,7].

Some related works [8–12] have focused on this problem and developed root cause analysis techniques that can determine the primary cause for the throughput limitation of a TCP flow from a passively captured packet trace. The novelty of our present work is that our technique is completely independent of both the applications and underlying PHY/MAC layer technologies. Our approach enables detailed profiling of short and long flows, unlike some previous work in this area such as [8,9] that only discuss the case of long TCP connections. In addition and contrary to most of the related work, we address the problem for most of the access technologies and not only in a given context like in [13,14]. The approach is general enough to be applicable for studying the impact of any application on the performance and not only for a particular application/service.

We exemplify the above techniques on traces collected on various access networks under the control of the same french ISP (Orange), enterprise traffic and also on a simulated traffic. Moreover, we underscore the limitation of classical techniques to pinpoint actual performance problem.

Last, pieces and preliminary results of this work have been appeared in the following conferences [15–17].

The new contribution in this work consists of (i) the in depth validation of the methodology, (ii) the application of the methodology to different network infrastructures, (iii) the evaluation of interactive services for both Internet and enterprise networks and (v) the comparison of the advantages and the limitations of our methodology with other known works. Also, we provide new materials (Flow charts, Tables, Sections) for the presentation of this complete work.

The remaining of this paper is organized as follows. In Section 2, the limitations of classical approaches to profile the performance of services is highlighted and the new performance metric is introduced. In Section 3, a new approach to address the problem is proposed and explained. In Section 4, an empirical validation of our key algorithms is provided. Section 5 is dedicated to analyze and discuss some results for interactive services. Related works are presented in Section 6 with a comparison between our methodology and other close works from the literature. Section 7 concludes our work.

## 2. Performance metrics

Response time is the performance metric that we focus on in this paper. By response time, we mean the delay between making a request and finishing receiving the response. Such a metric is often used to characterize application performance. That is why it is a particularly interesting metric for ISPs and application service providers to measure.

Given the number of different technologies used today to access the Internet, network characteristics and conditions are clearly among those factors. Also, network engineers observe that the core network is well provisioned and, in most of cases, it is not congested. Therefore, ISPs need to know the causes of 'not acceptable delays' to deal correctly with the problem. For this reason, response time is often complemented with throughput, round-trip time (RTT), and packet loss measurements in order to gain further insights into the observed response time.

However, we will show in the next subsection that the classical QoS indicators (response time with RTT and packet loss) fail to completely capture the underlying root causes of the observed poor performance. Hence, a more systematic approach to such performance analysis is necessary.

### 2.1. RTT estimation

The round trip time corresponds to the spent time between a sender transmitting a segment and the reception of its corresponding acknowledgement. This interval includes propagation, queuing, and other delays at routers and end hosts [18].

Several approaches have been proposed to accurately estimate the RTT from a single measurement point [8,19–21]. To estimate RTT, we adopted two techniques. The first method is based on the observation of the TCP 3-way handshake [20]: one first computes the time interval between the SYN and the SYN-ACK segment, and adds to the latter the time interval between the SYN-ACK and its corresponding ACK. It is important to note that we take losses into account in our analysis. The second method is similar but applied to TCP data and acknowledgement segments transferred in each direction[1]. One then takes the minimum over all samples as an estimate of the RTT. For the present work, RTT estimation will be based on the second approach.

### 2.2. Response time analysis

We exemplify the difficulty of interpreting the root cause of observed performance using simple metrics for the case of the Google search traffic extracted from real traces captured from ADSL, FTTH and 3G/3G+ access technologies of the 'Orange' Internet service provider. The description of the data-sets is presented in Appendix A. To identify the traffic generated by the Google search engine in the traces, we extract the TCP connections that transport HTTP requests containing Google tags in their HTTP header and exclude those to/from other services offered by Google like gmail, map, translate, etc. The used data set comprises of roughly 30 K, 1 K, and 6 K connections for Cellular, FTTH, and ADSL traces, respectively. Given that TCP congestion control has a noticeable impact on transfer times of different sizes we want to

---

[1] Keep in mind that we focus on well-behaved transfers for which there is at least one data packet in each direction.