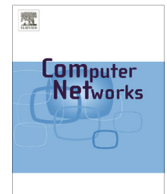




ELSEVIER

Contents lists available at ScienceDirect

Computer Networks

journal homepage: www.elsevier.com/locate/comnet

A flow measurement architecture to preserve application structure [☆]



Myungjin Lee ^{a,*,1}, Mohammad Hajjat ^{b,1}, Ramana Rao Kompella ^c, Sanjay G. Rao ^d

^a School of Informatics, University of Edinburgh, 10 Crichton Street, Edinburgh EH8 9AB, United Kingdom

^b Microsoft Azure, One Microsoft Way, Redmond, WA 98052, United States

^c Department of Computer Science, Purdue University, 305 N. University Street, West Lafayette, IN 47906, United States

^d School of Electrical and Computer Engineering, Purdue University, 465 Northwestern Avenue, West Lafayette, IN 47907, United States

ARTICLE INFO

Article history:

Received 7 April 2014

Received in revised form 14 August 2014

Accepted 10 November 2014

Available online 18 November 2014

Keywords:

NetFlow

Related sampling

Network management

ABSTRACT

The Internet has significantly evolved in the number and variety of applications. Network operators need mechanisms to constantly monitor and study these applications. Modern routers employ passive measurement solution called Sampled NetFlow to collect basic statistics on a per-flow basis (for a small subset of flows), that could provide valuable information for application monitoring. Given modern applications routinely consist of several flows, potentially to many different destinations, only a few flows are sampled per application session using Sampled NetFlow. To address this issue, in this paper, we introduce *related sampling* that allows network operators to give a higher probability to flows that are part of the same application session. Given the lack of application semantics in the middle of the network, our architecture, RelSamp, treats flows that share the same source IP address as related. Our heuristic works well in practice as hosts typically run few applications at any given instant, as observed using a measurement study on real traces. In our evaluation using real traces, we show that RelSamp achieves 5–10× more flows per application session compared to Sampled NetFlow for the same effective number of sampled packets. We also show that behavioral and statistical classification approaches such as BLINC, SVM and C4.5 achieve up to 50% better classification accuracy compared to Sampled NetFlow, while not impairing existing management tasks such as volume estimation too much.

© 2014 Elsevier B.V. All rights reserved.

1. Introduction

The tremendous success of the Internet, and the fertile ground for innovation that it provides has spawned a diverse range of applications, with new applications con-

tinually and rapidly emerging and gaining in prominence. The last decade alone has seen rapid growth in popularity of peer-to-peer (p2p) systems (e.g., BitTorrent, Skype, PPLive), online social networks (e.g., Facebook), and cloud-based applications (e.g., Salesforce [1], Google Apps [2]). While Internet traffic was predominantly dominated by the Web in the 1990s, p2p traffic accounted for over 60% of traffic around 2005, and more recently, video-based applications such as YouTube are gaining in popularity.

Concurrent with the growth of new applications and changes in popularity across applications, we are continually seeing shifts in characteristics and communication

[☆] Portions of this manuscript appeared in IEEE INFOCOM 2011.

* Corresponding author. Tel.: +44 (0)131 650 2713.

E-mail addresses: myungjin.lee@ed.ac.uk (M. Lee), hajjat@purdue.edu (M. Hajjat), kompella@cs.purdue.edu (R.R. Kompella), sanjay@ecn.purdue.edu (S.G. Rao).

¹ The work was performed when the authors were at Purdue University.

patterns of existing applications. For instance, the characteristics of Web traffic have significantly changed with the average size of Web objects increasing from 12 KBytes in 2000 to 68 KBytes in 2007 [3]. Further, p2p applications such as BitTorrent are being redesigned so that communication is localized within ISP networks, rather than crossing ISP boundaries [4,5].

The emergence of new applications, and their rapidly changing characteristics require network operators to continuously measure and monitor traffic characteristics in their networks. These measurements allow operators to, for instance, potentially re-provision their networks, detect any new types of undesirable behavior within applications (e.g., p2p system vulnerabilities [6–9], attacks on Ajax-based web services [10]) and in general, prepare their networks better against any major application trends. Further, such traffic monitoring must ideally occur in a ubiquitous fashion as applications characteristics may differ significantly depending on the location [11].

Router-level measurement solutions such as NetFlow represent the most widely deployed and popular approach used for traffic monitoring today. The widely available nature of NetFlow across many modern routers makes it an ideal candidate for ubiquitous low-cost network monitoring. Unfortunately, however, routers employ *packet sampling* to scale to high line rates, that makes NetFlow ill-suited to monitor the new range of applications evolving in the Internet today. In particular:

- Emerging p2p and cloud-based applications are routinely composed of many different flows to potentially different servers/hosts that are often geographically distributed. With random packet sampling, only a small subset of flows, if any, are sampled for an application session that comprises many different flows. This makes it difficult to accurately characterize application behavior from sampled data.
- Several researchers have pointed out the inadequacies of simple port-based classification for emerging applications such as p2p [12–14]. While several alternate approaches based on statistical techniques, or host behavioral patterns have emerged [15,13,16–18,14,19], much of this work has dealt with unsampled data. The effectiveness of these techniques is likely to degrade with random packet sampling.

Motivated by these limitations of random packet sampling, in this paper, we propose the notion of *related sampling* based on the following key idea: *Once a flow is sampled, all flows that are part of the same application session, are sampled with high probability.* Applying related sampling means that either an application session is (almost) fully sampled, or not sampled at all. Behavioral classifiers benefit from the extra information (of flows that are related) and characterization can be all the more accurate.

We explore the potential of related sampling in the context of the RelSamp architecture. Ideally, flows corresponding to the same application session must be identified as related. However, since determining this is hard, RelSamp considers all flows that contain the same source IP address

created within a given amount of time from each other as related. This heuristic is motivated from a measurement study on a 13-h campus trace. The RelSamp architecture incorporates related sampling with the help of three stages of sampling. First, we use a host selection probability that controls which host (identified using the IP address) gets selected for subsequent packet selection. Once a host gets selected, packets are subject to a flow-selection probability that governs the probability with which a flow that contains the host as the source IP address is created. Finally, the last stage of packet sampling dictates the rate at which flow records are updated. Thus, RelSamp biases packet and flow selection in favor of hosts that are already admitted.

Because RelSamp selects hosts based on source IP address, it can recognize flows from different hosts behind NAT as if they are from a single host. Hence, biasing host selection as proposed can be difficult in an environment where NAT devices are heavily deployed. As such, not all types of networks can rely on our approach. Instead, we identify two key important networks—*enterprise and campus networks*—where accurate application behavior monitoring and classification is crucial. We believe that NAT issue is less concern in the networks; it is relatively easy for operators to locate a right deployment place for RelSamp for mitigating the NAT issue. In that sense, RelSamp can be most suitable for those networks. By the same token, we ignore home and core networks from consideration of deploying RelSamp.

Our study is built upon the networks where NAT boxes are less deployed or operators have a full control over managing them. Under this condition, the paper makes the following contributions:

- We introduce *related sampling* that allows flows that are part of the same application session to be sampled with higher probability. Our architecture allows selecting a large majority of flows from a given application session thus allowing scalable monitoring and characterization of new and emerging Internet applications.
- Using real traces, we extensively evaluate the efficacy of RelSamp. In our results, we observe that RelSamp is capable of obtaining 5–10× more flows per application session compared to sampled NetFlow and flow sampling [20] without significantly compromising the accuracy of aggregate packet count estimates (less than 12% error).
- Using real packet traces with payload, we study the impact of RelSamp on traffic classification. Specifically, we show that the classification accuracy of BLINC, SVM and C4.5 increases by up to 50% in comparison with the flows output by sampled NetFlow for flows that are not easily classifiable using port numbers.

2. Measurement model

Consider a router at the edge of a large-scale campus network or an enterprise network, typically referred to as a *gateway*. Our goal in this paper is to facilitate scalable monitoring of application traffic in order to characterize, study and monitor application behavior in a continuous fashion at such enterprise gateways. Such application

Download English Version:

<https://daneshyari.com/en/article/451792>

Download Persian Version:

<https://daneshyari.com/article/451792>

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