

# MALAWI: Aggregated longitudinal analysis of the MAWI dataset

João Taveira Araújo  
University College London  
j.araujo@ee.ucl.ac.uk

Kensuke Fukuda  
National Institute of Informatics  
kensuke@nii.ac.jp

## 1. ABSTRACT

The importance of measurement and analysis of Internet traffic is constantly reasserted as the Internet expands and shifts in often unpredictable ways. The MAWI dataset [1], which provides daily traces across a trans-Pacific link over the past decade, has often been used to analyze traffic from a network perspective. In this paper we focus on information contained at the transport layer and present MALAWI (Measurement and Aggregated Longitudinal Analysis on the WIDE Internet) a new dataset derived from MAWI which extracts information from traced TCP flows and aggregates these statistics by geographical location, AS and network prefix. We briefly illustrate the usefulness of this new dataset by analyzing a month of data to observe the impact of the Tohoku earthquake on delay and loss.

## 2. INTRODUCTION

Measurement and analysis of Internet traffic is critical not only for a deeper understanding of the evolving nature of Internet as a whole but also as an input to designing new elements which are able to act efficiently within the current architecture.

The MAWI dataset contains daily 15-minute traffic traces with transport headers spanning the past decade. While the dataset has been available to the wider community for some time, the short timespan of each trace has lent it to further study in areas where the inexistence of complete flow traces is less significant, such as Internet anomalies [3] or where characterization of traffic is packet-based [2], relying only on the inspection of the IP header and port numbers.

The network layer alone however does not contain much of the information which defines a user's experience. To understand many of the inherent characteristics perceived by an application it is necessary to analyze TCP, extracting values for relevant metrics such as the round trip time (RTT) and loss. This flow analysis requires the partial reconstruction of TCP flows to obtain robust measurements and has been attempted before [4] but limited in scope to complete, bidirectional flows. Additionally, aggregating these statistics in a meaningful manner poses significant challenges due to both the scale of data generated and the availability of external sources to provide context to the original MAWI traces.

MALAWI (Measurement and Aggregated Longitudinal Analysis on the WIDE Internet) builds on the MAWI dataset and will make available flow level statistics aggregated by source and destination prefixes, autonomous system (AS) and geographic location. Both prefix and AS information for each IP is extracted from information contained within the daily BGP routing updates, which are also collected from within WIDE since 2004. Currently geolocation information is obtained from freely available sources as far back as 2008. Prior to this date, country level information is used based on information provided by regional Network Information Centers (NIC).

With the resulting dataset we hope to provide researchers with greater insight into essential metrics without the limitations imposed by IP address anonymization. While the dataset is intended for longitudinal studies of the evolution of TCP behaviour, we illustrate the potential of MALAWI by analyzing the traces for March of 2011 and viewing effect of the Tohoku earthquake on both RTT and loss.

## 3. TOHOKU EARTHQUAKE

While the devastating effects of the Tohoku earthquake on Friday 11th, March 2011 are well known, the impact on network operations within Japan are less than clear. While a significant proportion of both users and infrastructure were in largely unaffected regions, it

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

ACM CoNEXT Student Workshop, December 6, 2011, Tokyo, Japan.

Copyright 2011 ACM 978-1-4503-1042-0/11/0012 ...\$10.00.

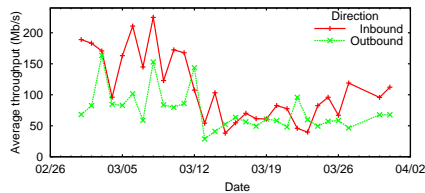


Figure 1: Inbound and outbound average throughput for March 2011 MAWI traces.

is not obvious what the effects of large, if localized, network outages were for the wider network, or even the effect of creeping uncertainty in Internet users themselves as the extent of the disaster became apparent.

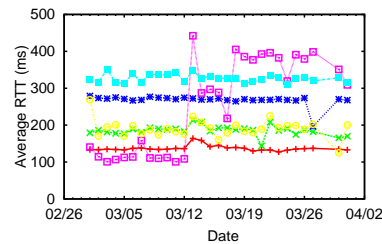
Fig.1 shows the average inbound and outbound throughput for the MAWI traces for the entirety of March. While the total inbound traffic drops from the 11th onwards, outbound traffic rises sporadically on the 12th. Both inbound and outbound traffic would remain below prior levels for the following months, in part due to power saving measures.

While the demand on both sides seemingly fell, the average RTT remained largely unaffected for many destinations as shown in fig.2a. We reduce the dataset to traffic destined to known ports in order to filter out most residential traffic, which is more prone to excessive buffering. The most visible increase in RTT is on the 13th March (Sunday), where it increases for most destinations, likely due to reconfigurations in the routing infrastructure. While the MAWI dataset is often referred to as a trans-Pacific link, in reality the measurement point precedes a peering point with a Japanese commercial provider. Further inspection of the RIB database for this period reveals that many prefixes were shifted away from this provider and towards the US. As a result, the RTT to many hosts in China more than doubled as packets travelled across the Pacific and back.

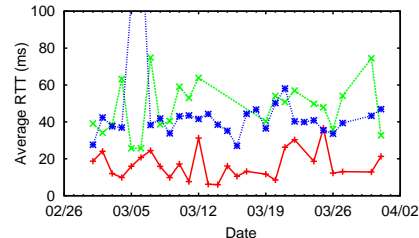
The RTT from the measurement point to the sources within Japan shows an abnormal increase in RTT for a Saturday on the 12th March, and the inexistence of data points for the Fukushima region between the 13th and the 19th of March, most likely due to routing changes which took one or both directions of the flow beyond the measurement point. Additionally the effect on upstream loss (toward destination) is shown in figure 3. Although the figure is hard to analyze, due in part to the decrease in traffic, a network bottleneck is likely to have formed on March 19th, as it affects most destinations. Whether or not this is due to routing reconfigurations requires further analysis.

## 4. CONCLUSIONS

MALAWI will provide a different perspective on the MAWI dataset, offering flow level statistics aggregated according to data provided by BGP routing updates



(a) Average RTT from measurement point to known port numbers in selected destinations.



(b) Average RTT from measurement point to unknown ports in sources for Japan prefectures.

Figure 2: RTT values from measurement point.

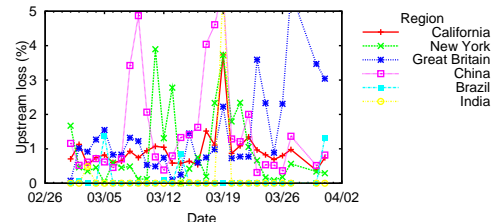


Figure 3: Upstream loss ratio from measurement point to known port numbers in destination regions.

and geolocation sources. The ability to not only observe TCP metrics over time, but also relate these events to changes in the routing topology as experienced by a single vantage point over a decade will hopefully allow further work to be developed in tracing the evolution of the Internet as a whole.

## 5. REFERENCES

- [1] Mawi dataset. <http://tracer.csl.sony.co.jp/mawi/>.
- [2] P. Borgnat, G. Dewaele, K. Fukuda, P. Abry, and K. Cho. Seven years and one day: Sketching the evolution of internet traffic. *INFOCOM 2009, IEEE*, pages 711 – 719, 2009.
- [3] R. Fontugne, P. Borgnat, P. Abry, and K. Fukuda. Mawilab: combining diverse anomaly detectors for automated anomaly labeling and performance benchmarking. *CoNEXT 2010*, page 8, 2010.
- [4] S. Jaiswal, G. Iannaccone, C. Diot, J. Kurose, and D. Towsley. Measurement and classification of out-of-sequence packets in a tier-1 ip backbone. *IEEE/ACM Transactions on Networking*, 2007.