

Characterizing Internet traffic: When the viewpoint matters

(ongoing work) Part of Msc Thesis

Arsakian Armen

Electrical Engineering & Computer Engineering
Democritus University of Thrace
arsakian at ee.duth.gr

5 November, 2009

Outline

Statistics of Traces

Data sets

Pareto vs Poisson

Asymptotic behavior and data length

Traffic Generators

Network load and estimations

Bibliography

Introduction

Understand the interaction of core internet and endpoint TCP/IP stacks

- ▶ Solution: Statistical analysis of data sets

Available Methods

- ▶ Parametric approach
 - ▶ I estimate the parameters of the distribution
 - ▶ I could use Maximum Likelihood (ML)
- ▶ Non parametric approach
 - ▶ I do not assume the distribution
 - ▶ Estimate empirical density with histograms, kernel based methods

Data Sets

Data set collected inside the network

- ▶ Flow level capturing method gives structured and rich information ([FML⁺03] describe technical issues)
- ▶ Flow is unidirectional set of packets characterized by a five tuple (IP/PORT src/dst, protocol)

From Data sets we can extract various per flow statistics for both network and protocols

- ▶ E.g. Network : RTTs, available bandwidth
- ▶ E.g. Protocols: TCP connection time, file sizes of HTTP FTP.

Landmarks

[WPT],[PF95] investigated WAN traces

- ▶ They defined a Stochastic Process having as Random Variable packet arrival times
- ▶ They adopted parametric application-wise approach

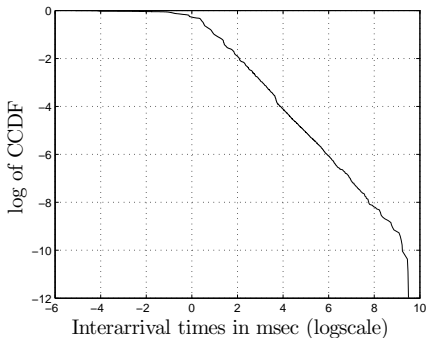
The S.P. exhibits Self Similar behavior

- ▶ A H-SS process is characterized by Hurst exponent $H \in [0, 1]$
- ▶ Reason: The R.V. that generates the packet sequence is heavy tailed
- ▶ Known heavy tail distribution: Pareto

Pareto is characterized by heavy tail index α , The line in a loglog plot is the power law signature of the RV

$$\frac{d \log \overline{F}(x)}{d \log x} = -\alpha$$

- ▶ If the application generates SS Traffic then one or more of the parameters describing the application must be heavy tailed



Application Independent approach

The three assumptions made when using S.P. to describe RVs representing packet arrival times

- ▶ Stationarity
- ▶ Independence
- ▶ Distribution

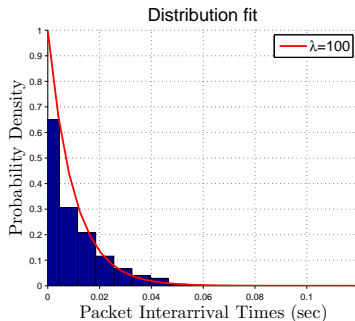
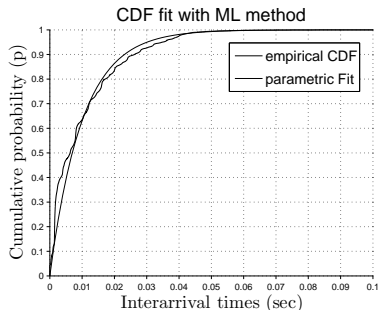
The processes we examine are called process with Independent Stationary Increments (ISI)

- ▶ H-SS is a process with stationary increments
- ▶ stationary increments \rightarrow stationarity
- ▶ Increments could be packet interarrival times X_k as they accumulate and give packet arrival times S_k

$$S_k = X_1 + \dots + X_k = S_{k-1} + X_k$$

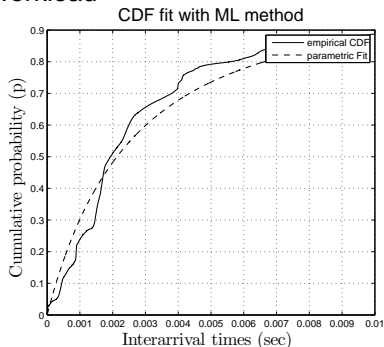
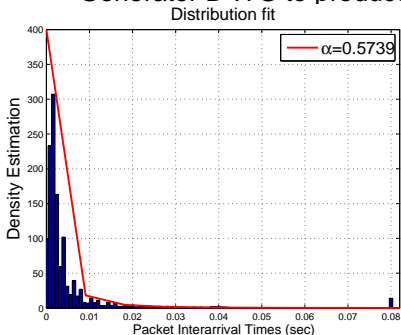
Testing Distribution for Poisson

- ▶ Poisson process $N(t)$ is a point S.P. with i.i.d exponentially distributed interarrival times $N(t) = \max\{k : S_k \leq t\}$
- ▶ Test assumptions with a LAN experiment. I used Traffic Generator D-ITG to produce Poisson workload



Testing Distribution for Pareto

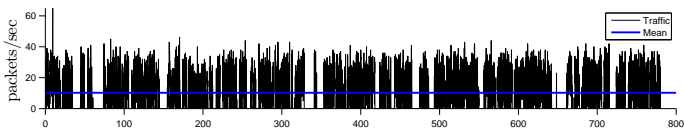
- ▶ Bursty Traffic with i.i.d heavy tail distributed interarrival times
- ▶ Test assumptions with a LAN experiment. I used Traffic Generator D-ITG to produce Bursty workload



Time series of the above LAN experiments

Use structural source models

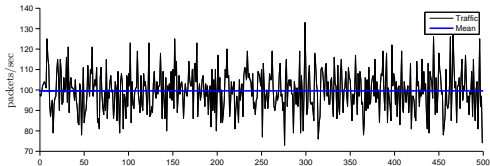
- ▶ Each source sends data according to ON/OFF times which are heavy tailed.
- ▶ Pareto principle: e.g. 20% of the ON times account for 80% coming from that source → Bursty Traffic
- ▶ Observe the ON/OFF periods on the time series
- ▶ Packet arrivals count over 100msecs intervals



Time series of the above LAN experiments pt2

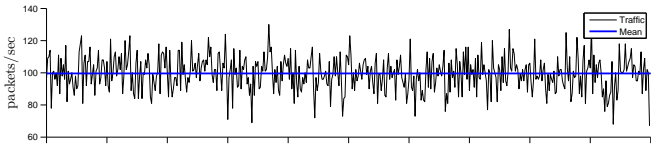
Poisson

- ▶ Observe there is no apparent trend or periodicity
- ▶ Packet arrivals count over 100msecs intervals



Non stationary Poisson

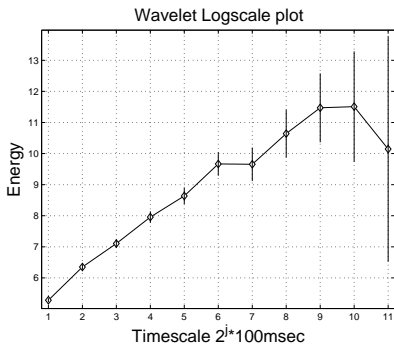
- ▶ generated artificially by adding i.i.d uniformly distributed noise to the mean
- ▶ Observe the trends



LRD or Trend

Long Range Dependence is a statistical property of the RV e.g. packet arrival times

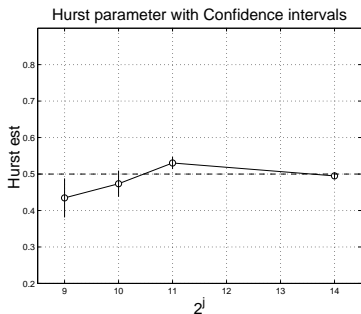
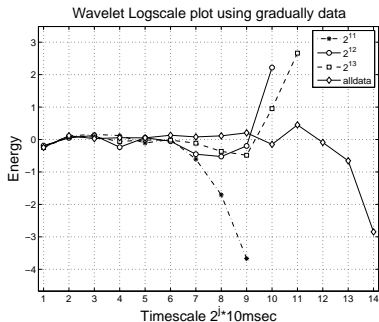
- ▶ When $H > 0.5$ the H-SS process creates LRD
 - ▶ Mathematically: The covariances are non-summable, the process has long memory
 - ▶ there are process that create artificially LRD e.g. FARIMA, fGN (see [KMF04])
-
- ▶ LRD is valid when we observe a straight line in a log-log wavelet plot (for wavelets and LRD see [AV98])



subtleties of LRD

Since LRD is an asymptotic property the measurements should last long enough
 How long?

- ▶ We test Poisson traffic we should see a straight line on the wavelet loglog plot and the Hurst exponent should be $H = 1/2$



Traffic Generators

Traffic generator are used to to perform tests while configuring various parameters.

TGs allow us to investigate the impact of the parameters that are going to be modelled on the traffic

They work using collected traces as inputs (for an overview [HC06])

They are divided into open loop and closed loop

- ▶ Open loop simply replay the collected trace
- ▶ Closed loop take into account the prevalent network conditions
 - ▶ Used frequently with emulators
 - ▶ They are application-oriented or application independent
 - ▶ They use initial trace to configure the TG. This means: find empirical distributions of the parameters that are going to be modelled.
 - ▶ Parameters: Source level e.g. file sizes or HTTP session duration, Network link capacities, link delays

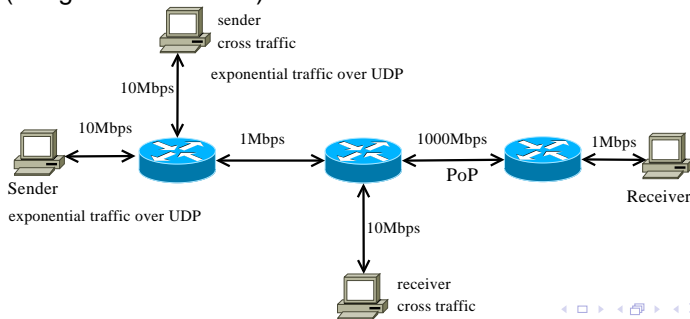
Emulators offer flexibility to administrators when they tune their network characteristics prior to applying them over the real world

Whether we evaluate traces (statistically) or configure TGs, we extract the information needed from backbone measurements.

How reliable and representative of source behavior is the observed flow provided that the instrumentation introduces negligible errors.

- ▶ This depends on where exactly on the path of a flow is the Point of Presence and on network load conditions.

Simulation with non-responsive traffic (UDP) on the following topology (foreground: one flow)



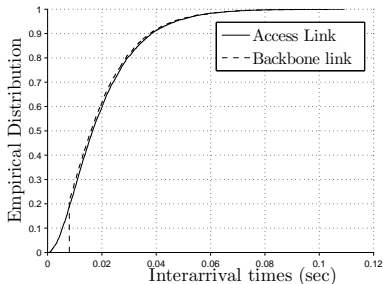
Argumentation

Path has capacities C_i $i = 0 \dots H$. Estimate \hat{F}_X \hat{F}_Y the CDF of interarrival times on the access link l_1 , and backbone link l_j respectively

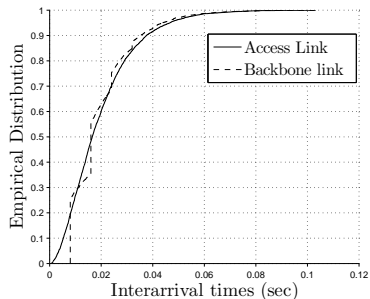
To see how CDF changes we examine the order statistics (non parametric approach)

- ▶ No cross Traffic: Y_k^j can increase if $x_k^1 < \max_{i=2\dots j} \{L/C_i\}$, L : k^{th} packet size (assumed constant)
 - ▶ Cross Traffic with intensity S s.t. $S_i < C_i$: Y_k^j can increase if $x_k^1 < \max_{i=2\dots j} \{L/A_i\}$ where $A_i = C_i - S_i$
 - ▶ Cross Traffic with intensity S s.t. $S_i > C_i$:
 - ▶ Y_k^j increases $x_k^1 < q_k^i - q_{k-1}^i$. Queuing delay q_k^i of k^{th} packet at link i
 - ▶ Y_k^j decreases $x_k^1 > q_k^i - q_{k-1}^i$.
1. No cross traffic see how link capacities C_i affect \hat{F}_Y
 2. Cross traffic with average rate $\bar{S} < C$
 3. Cross traffic with rate \bar{S} such that combined with the foreground flow causes queue build up

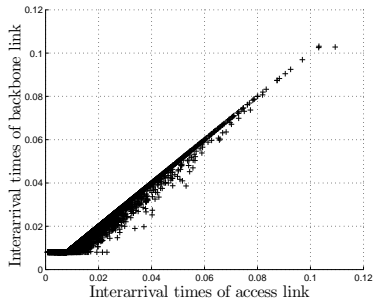
Use exponential T.G. on top of UDP
with average rate 0.4Mbps No
cross traffic



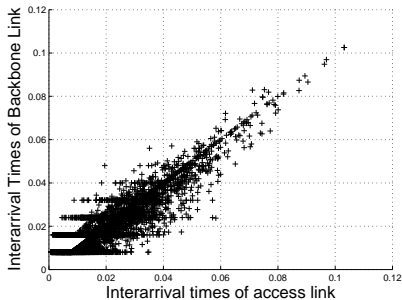
Use exponential T.G. on top of UDP
with average rate 0.4Mbps
exponential cross traffic with
average rate 0.4Mbps



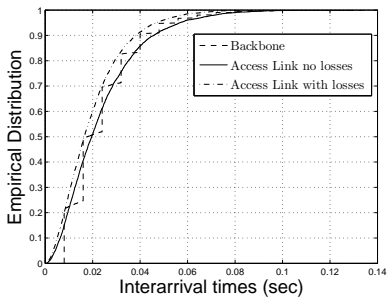
Test the casual relation of X & Y
No cross traffic



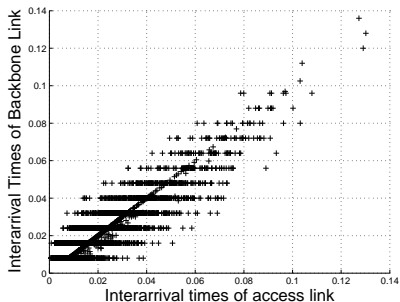
Test the casual relation of X & Y
with cross traffic



Cross traffic that causes queue build up



Test the casual relation of X & Y with queue build up



Order Statistics

Since y_K^j changes this affects order statistics. Examine order statistics for the three cases

1st order statistic equals $L/\min\{C_i\}$ Path of Foreground flow : {10 1 1000 1}

- ▶ No cross traffic

link	1st order statistic	median	IQR
Access	$8 * 10^{-4}$	0.0165	0.017
Backbone	$8 * 10^{-3}$	0.016	0.017

- ▶ Cross Traffic present

link	1st order statistic	median	IQR
Access	$8 * 10^{-4}$	0.0165	0.0171
Backbone	$8 * 10^{-3}$	0.016	0.016
Access with losses	$8 * 10^{-4}$	0.0165	0.0171
Backbone with losses	$8 * 10^{-3}$	0.0177	0.016

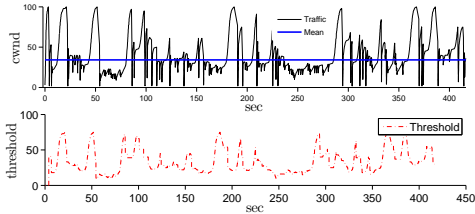
Responsive traffic: an Endpoint view

In contrast with Non Responsive traffic, available BW and RTT affects dynamics and so observed traffic

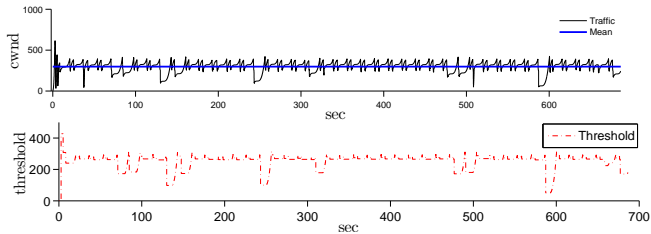
collection method: Tcprobe: kernel built in tool in Linux distros

- ▶ Application Tested: FTP
- ▶ Measurements performed on WLAN Network Characteristics: lossy link, small RTTs 4msecs Available BW: 2.5Mbps
 - ▶ Bursty and severe losses
 - ▶ Large trends
- ▶ Measurements performed in collaboration with S. Lenas on the Internet DUTH → MIT available BW measured by Iperf 15Mbps average RTT 150msecs, 15hops access link: wireless
 - ▶ Periodic less losses
 - ▶ Larger `cwnd` values allow TCP to deal with losses more efficiently

WLAN



Internet



Remarks & Ongoing Work

Remarks

- ▶ TGs used produce stationary and non responsive traffic to varying network load conditions
- ▶ Adopted an application-independent approach
- ▶ In reality during a packet capturing process we are not aware of what is going on the access link, an inference method would not be possible
- ▶ Did not mention RTTs

Ongoing Work

- ▶ Use Responsive traffic TCP
- ▶ test more elaborate scenarios
- ▶ adopt application-oriented approach

Bibliography I



Patrice Abry and Darryl Veitch.

Wavelet analysis of long range dependent traffic.

IEEE TRANSACTIONS ON INFORMATION THEORY, 44:2–15, 1998.



Chuck Fraleigh, Sue Moon, Bryan Lyles, Chase Cotton, Mujahid Khan, Deb Moll, Rob Rockell, Ted Seely, and Christophe Diot.

Packet-level traffic measurements from the sprint ip backbone.

IEEE Network, 17:6–16, 2003.



Félix Hernández-Campos.

Generation and validation of empirically-derived tcp application workloads.

Phd Thesis, University of North Carolina, 2006.



Thomas Karagiannis, Mart Molle, and Michalis Faloutsos.

Long-range dependence: Ten years of internet traffic modeling.

IEEE Internet Computing, 8(5):57–64, 2004.

Bibliography II



Vern Paxson and Sally Floyd.

Wide-area traffic: The failure of poisson modeling.

IEEE/ACM Transactions on Networking, 3:226–244, 1995.



Walter Willinger, Vern Paxson, and Murad S. Taqqu.

Self-similarity and heavy tails: Structural modeling of network traffic.

In *Statistical Techniques and Applications*, pages 27–53. Verlag.