

# Sources of Long-Range Dependence

## *Why is the Internet Long-Range Dependent?*

Richard Clegg (richard@manor.york.ac.uk)

Networks and Nonlinear Dynamics Group,

Department of Mathematics,

University of York

Slides prepared using the Prosper package and L<sup>A</sup>T<sub>E</sub>X

# Talk Summary

- (Yet Another) Introduction to Long-Range Dependence (LRD)
  - What is this LRD thing?
  - Why does it matter to the Internet?
  - How can LRD be measured?
- Four sources of LRD in Internet traffic measurements.
  - Inherent at source. Aggregation from heavy tailed sources.
  - TCP feedback mechanisms. Arises in network itself.
- What can experiment tell us?
  - Measuring real world systems.
  - Simulation results.
  - (Future Work) Mathematical models.
- Conclusion: “Where *does* LRD come from in the internet?”

# What is LRD?

- Usually defined in terms of an auto-correlation function (ACF) which sums to infinity.

- For a weakly stationary time-series  $\{X_i : i \in \mathbb{N}\}$  with mean  $\mu$  and variance  $\sigma^2$ , the ACF is defined by

$$\rho(k) = \mathbf{E} [(X_i - \mu)(X_{i+k} - \mu)] / \sigma^2$$

- A time series is LRD if its ACF sums to infinity,

$$\sum_{k=-\infty}^{\infty} |\rho(k)| = \infty.$$

- Often the (stricter) asymptotic form

$$\rho(k) \sim c_\rho k^{(H+1)/2},$$

is given, where  $c_\rho \in \mathbb{R}_+$  and  $H \in (1/2, 1)$  is the Hurst parameter.

# What is Strange about LRD?

- The majority of commonly encountered statistical processes (Poisson processes, Gaussian noise, finite ARIMA models) do not exhibit LRD.
- Usually  $\bar{X}$  converges to  $\mu$  at a rate:  $\text{var}(\bar{X}) \sim C\sigma^2/n$ . [Equality with  $C = 1$  for independent data.] The mean for LRD series converges slower than this:  $\text{var}(\bar{X}) \sim \frac{Cn^{2H-2}}{\sigma^2H(2H-1)}$ .
- The estimator  $S^2$  for sample variance given by

$$S^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n-1},$$

is no longer an unbiased estimator.

- Standard methods for estimating confidence intervals (t-statistic) fail. As sample size  $n \rightarrow \infty$ , any given parameter measurement will, almost surely, fall *outside* a  $(1 - a)$  confidence interval even for arbitrarily small  $a$ .

# Measuring LRD

- Measuring the ACF is not a good way to establish the presence of LRD.
- LRD is detected in the slope at high lags. ACF is only accurate at low lags. (ACF estimator is biased in presence of LRD).
- Some biased estimators with poor convergence performance.
- All are vulnerable to some extent to non-stationarities in the data.
- Periodicity and trends in particular can be a problem.
- While some estimators give confidence intervals, often results from different estimators do not agree even within 95% intervals.
- More information:  
<http://math.bu.edu/people/murad/methods/index.html>

# Methods for Estimating LRD

- R/S plot — the rescaled adjusted range. The oldest method for measuring  $H$ .
- Aggregated Variance — logarithm of variance versus logarithm of aggregation level.
- Periodogram — log periodogram (estimate of spectral density) versus frequency.
- Whittle's Estimator — an approximate MLE.
- Local Whittle — semi-parametric approximate MLE (parametric at frequencies near the origin).
- Wavelet based — a technique based upon a generalisation of Fourier Series.

# Testing Hurst Estimates

| Data   | H     | R/S  | M. R/S | Var  | L.W. | Wlts |
|--------|-------|------|--------|------|------|------|
| FGN    | 0.625 | 0.62 | 0.66   | 0.63 | 0.63 | 0.63 |
| FGN    | 0.75  | 0.71 | 0.74   | 0.73 | 0.77 | 0.76 |
| FGN    | 0.875 | 0.80 | 0.81   | 0.81 | 0.90 | 0.89 |
| Markov | 0.625 | 0.60 | 0.62   | 0.57 | 0.67 | 0.67 |
| Markov | 0.75  | 0.66 | 0.68   | 0.60 | 0.79 | 0.78 |
| Markov | 0.875 | 0.70 | 0.76   | 0.78 | 0.88 | 0.87 |
| Random | 0.5   | 0.54 | 0.55   | 0.47 | 0.49 | 0.49 |

Conclusion: A man with one watch knows the time. A man with two is never sure.

# LRD — Why Should we Care?

- In 1993 LRD was found in a time series of bytes/unit time [3] measured on an Ethernet LAN.
- This finding has been repeated a number of times by a large number of authors (however [1] suggests this may not happen in the core).
- A higher Hurst parameter often increases delays in a network. Packet loss also suffers.
- If buffer provisioning is done using the assumption of Poisson traffic then the network will be underspecified.
- The Hurst parameter is “...a dominant characteristic for a number of packet traffic engineering problems...”. [2]



# Sources of LRD

## (1) Data is LRD at Source

- Claim arises from measurements on VBR video traffic.
- Pictures are updated by sending changes.
- A still scene is few changes, a cut or pan is a lot of changes.

## (2) Data arise from aggregation of heavy tailed ON-OFF sources.

- A random variable  $X$  is heavy-tailed if for all  $\varepsilon > 0$  it satisfies

$$\mathbb{P}[X > x] e^{\varepsilon x} \rightarrow \infty, \quad x \rightarrow \infty.$$

- It can be shown [4] that ON/OFF sources with heavy-tailed train lengths leads to LRD.
- It has been observed that the sizes of files transferred on the internet follow a heavy-tailed distribution [5].

# Sources of LRD (continued)

(3) LRD arises from feedback mechanisms in the TCP protocol.

- This claim comes from Markov models of TCP timeout and retransmission.
- A Markov model is used to show that the doubling of timeouts can cause correlations in timeseries of transmitted data.
- Modelling shows that this can lead to LRD over certain timescales (“local” LRD).

(4) LRD arises from network topology or routing.

- Consider a simulation on a Manhattan network with randomly distributed sources and sinks.
- The sources produce Poisson traffic.
- Packets find their shortest route to the sink (accounting for the traffic on the next hop).
- In this simple situation the aggregated traffic shows LRD.

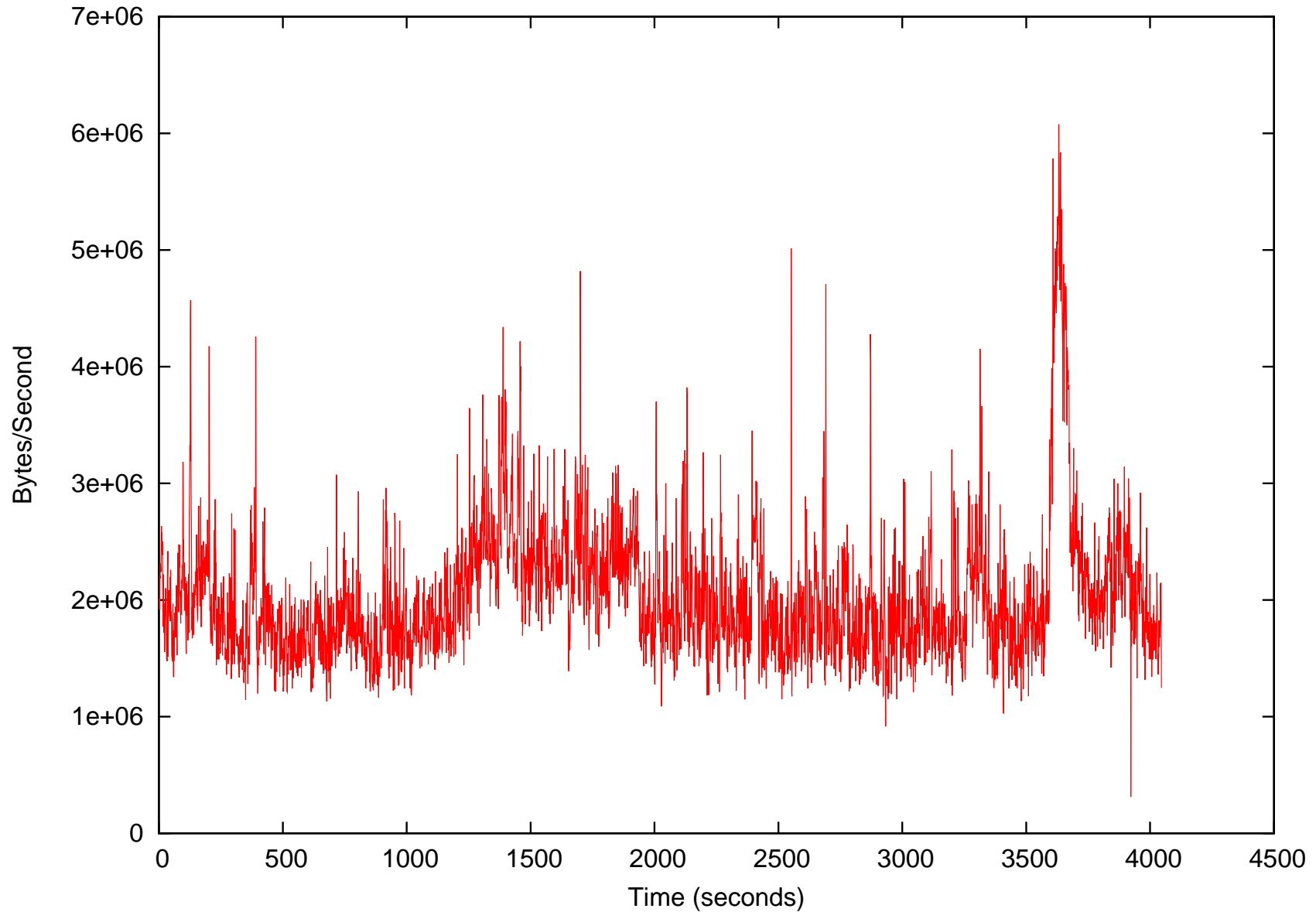
# Data Collection — Data Facts

- Data collected at incoming/outgoing pipe at University of York.
- 8.23 GB of data in 13.6 million packets — 67 minutes of data.
- 7.81 GB of this data is TCP. 0.6MB of data is ICMP. 0.4GB of data UDP.
- Outgoing data: 1.95GB of data in 6.0 million packets (av size: 323 bytes).
- Incoming data: 6.29GB of data in 7.7 million packets (av size: 821 bytes).
- Data has been anonymised and is available for researchers on request (1Gb `tcpdump` format file).

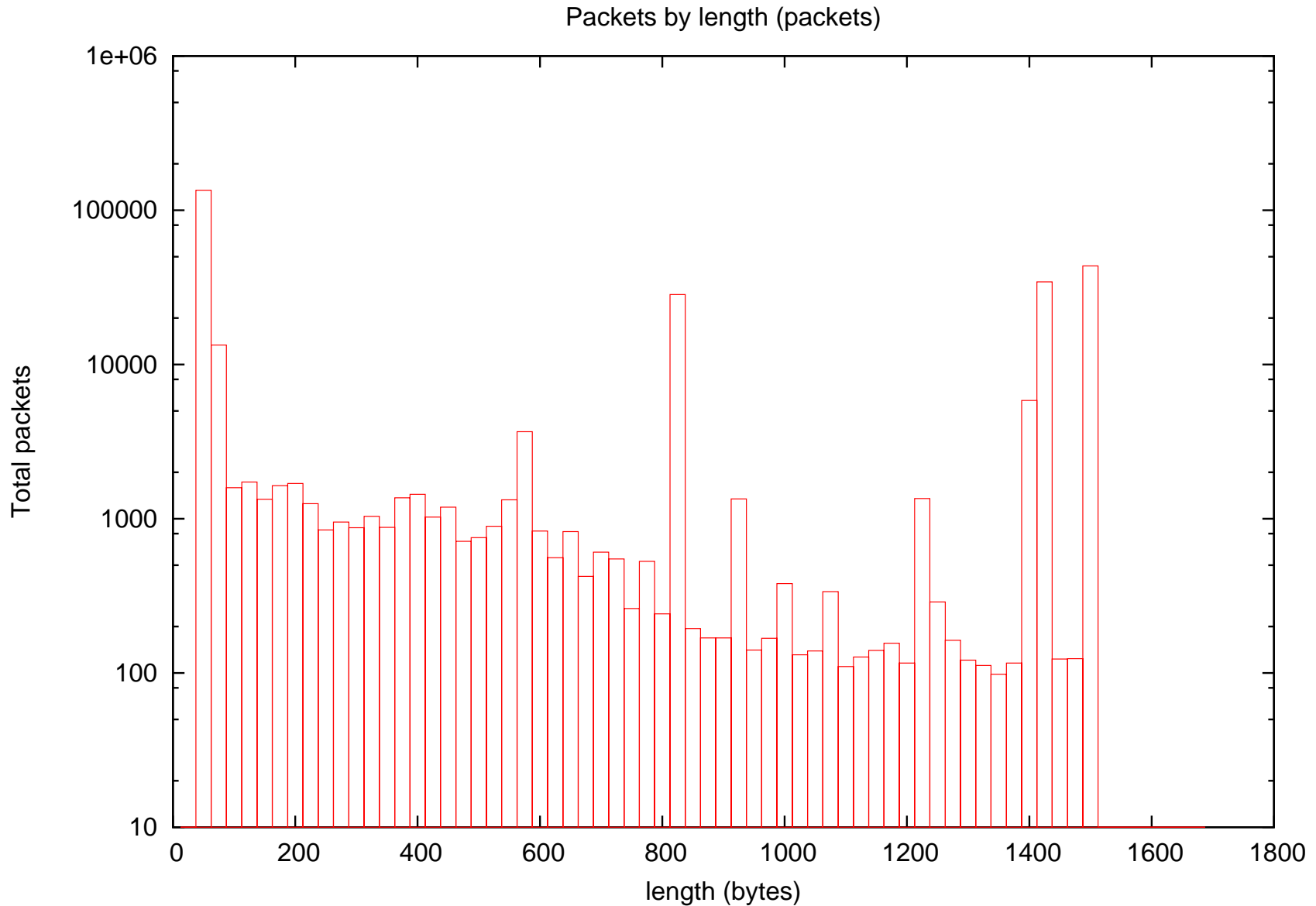
# Disaggregating the data

- In addition to inbound and outbound we can break up traffic by port number.
- Ports are usually associated with particular services.
- Port 80 HTTP (5.78GB) — Web traffic is by far the bulk of the traffic
- Port 25 SMTP (226MB) — Email is large but not by comparison
- Port 21 and 20 FTP (230MB) — FTP is also large (but tends to be more “bursty”).
- Port 53 DNS (33MB) — DNS data doesn’t seem to account for much (but this may be due to where we are looking).

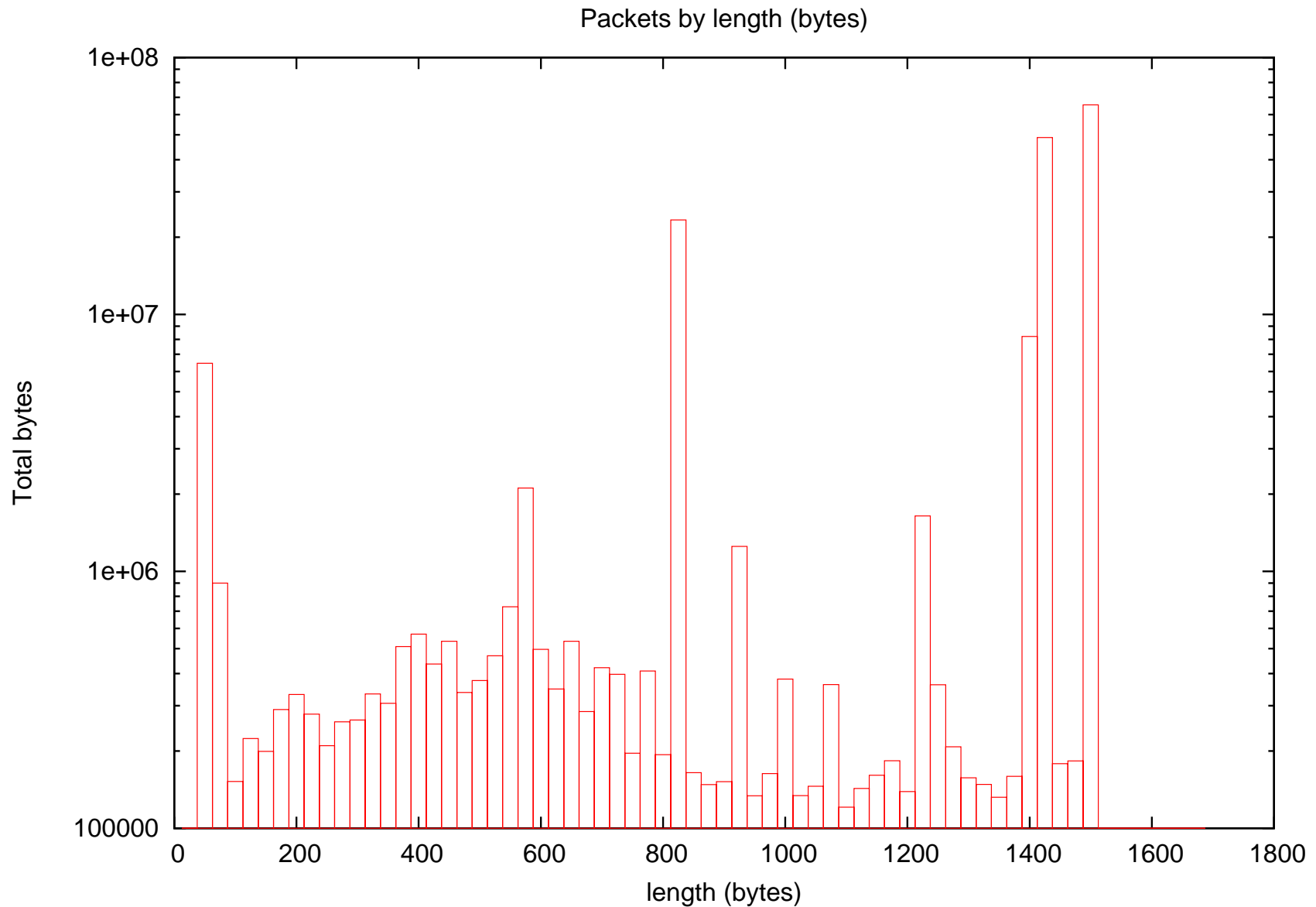
# Raw Data — bytes/second



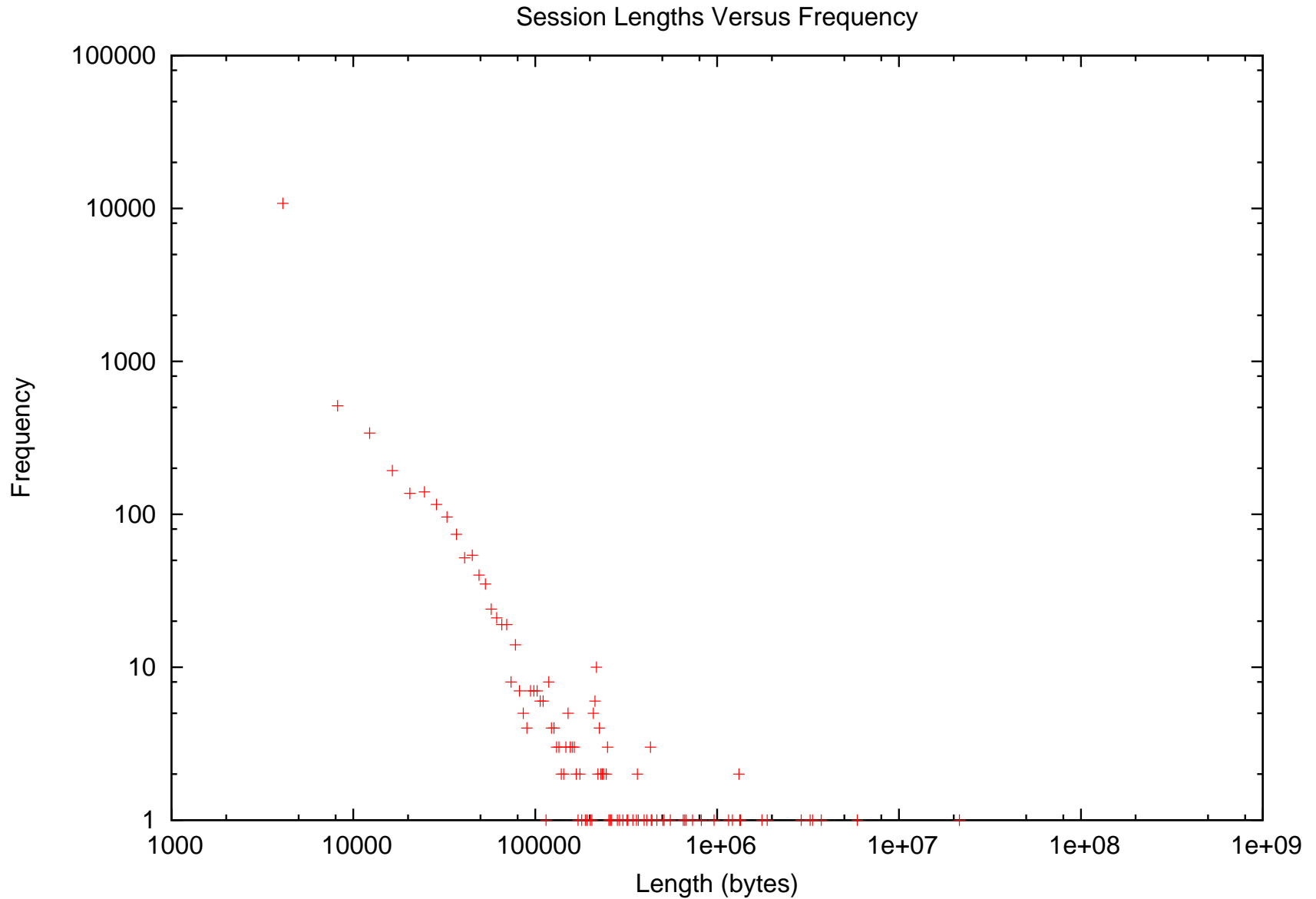
# Packet lengths — by no of packets



# Packet lengths — by no of bytes



# TCP Session Lengths





# Hurst Estimates

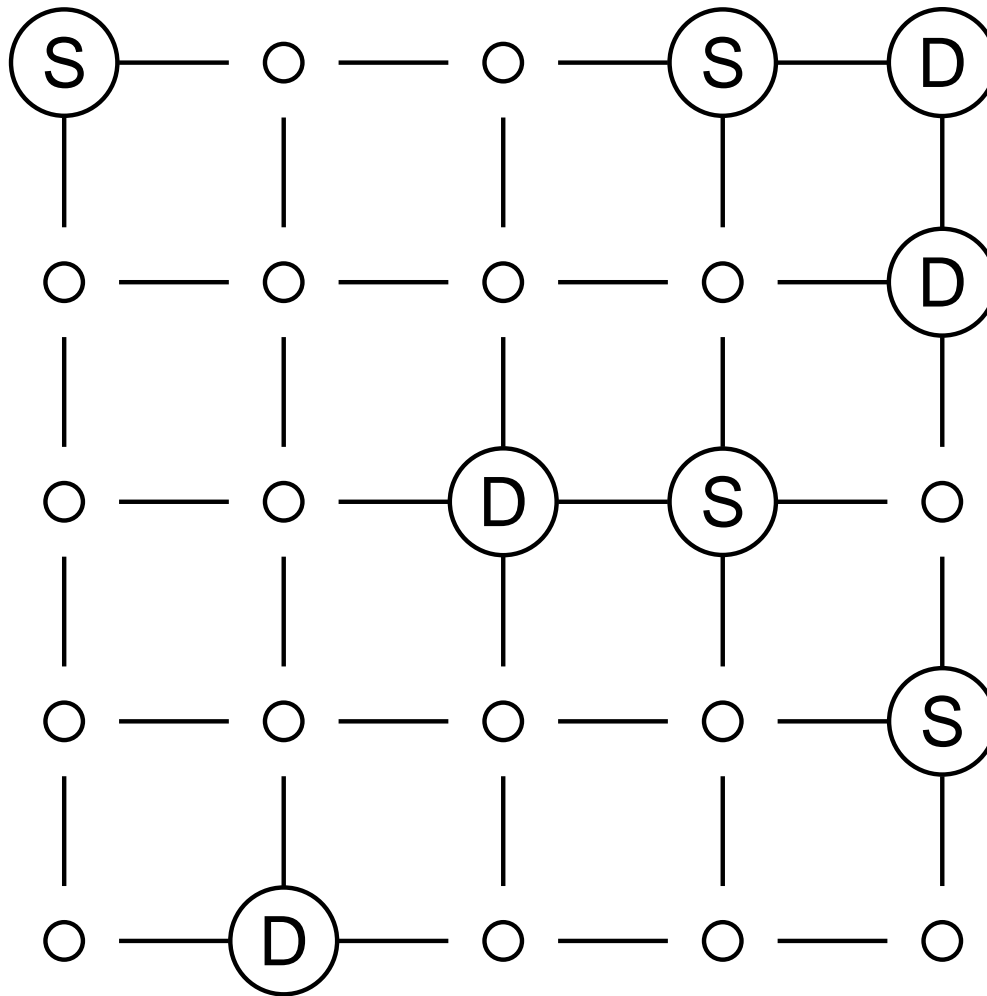
| Data                | R/S  | M. R/S | Var  | L.W.  | Wlts  |
|---------------------|------|--------|------|-------|-------|
| Total Bytes / sec   | 0.75 | 0.87   | 0.88 | 0.98  | 0.95  |
| Total Pkts / sec    | 0.75 | 0.86   | 0.86 | 0.94  | 0.91  |
| Total Bytes / tenth | 0.78 | 0.87   | 0.87 | 0.87  | 0.87  |
| Total Pkts / tenth  | 0.75 | 0.85   | 0.84 | 0.89  | 0.89  |
| http only           | 0.75 | 0.88   | 0.90 | 0.87  | 0.87  |
| smtp only           | 0.83 | 0.78   | 0.76 | 1.06* | 1.06* |
| ftp only            | 0.66 | 0.90   | 0.92 | 0.67  | 0.61  |

Counts are bytes / tenth if not otherwise stated. \* Indicates readings which are outside a reasonable range for H.

# Simulation Model

- Manhattan network with randomly dispersed sources and destinations.
- Sources may produce Poisson or LRD traffic which is sent to a randomly selected destination.
- Packets route based upon a “least hops to destination” algorithm.
- However, when routes are equal hops, the least congested hop is chosen.
- Alternatively, a “fixed route” algorithm may be used.
- Congestion is all at nodes — nodes send one packet per simulation step.
- LRD sources use a double intermittency map to produce traffic with given  $H$ .

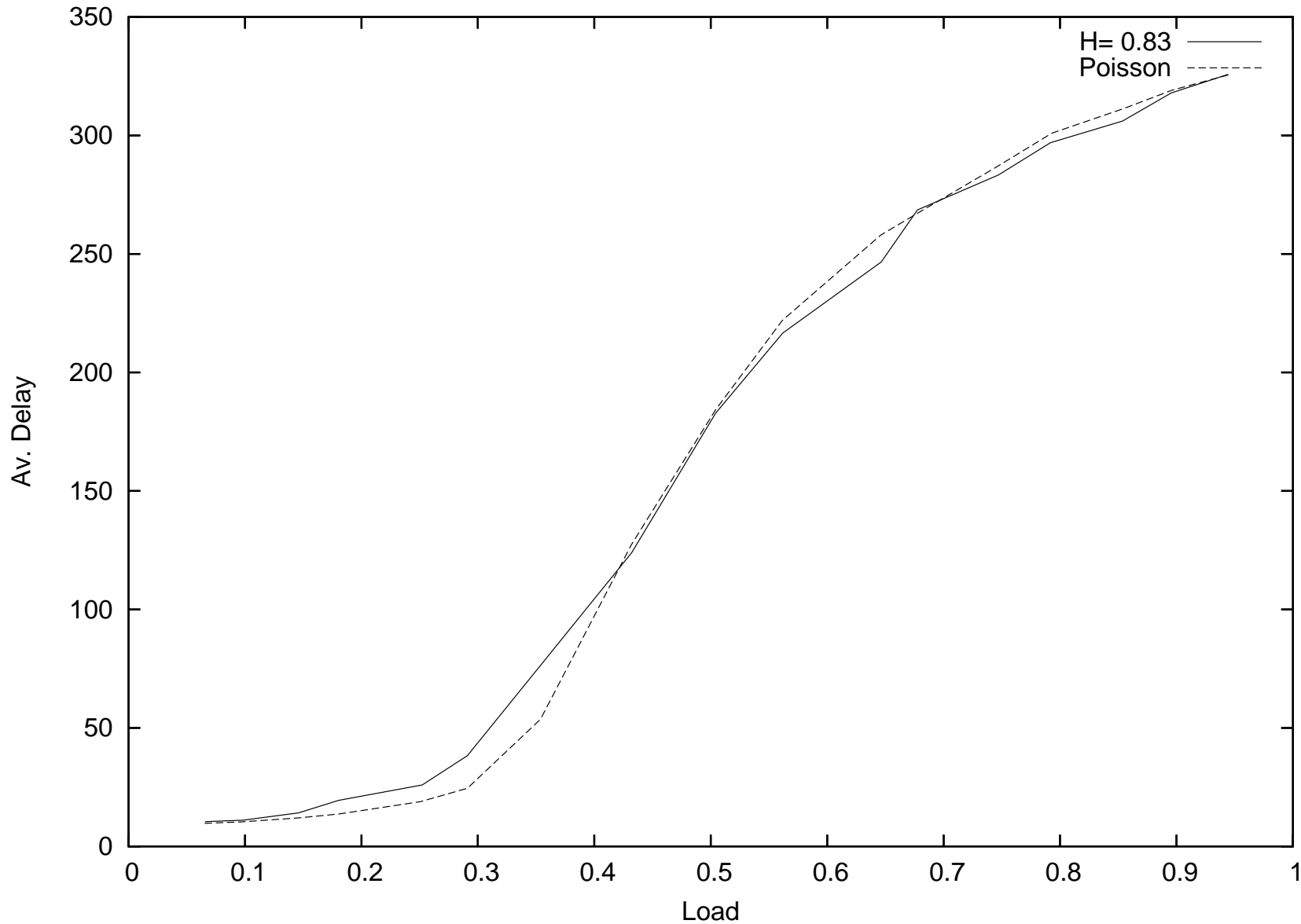
# Typical layout



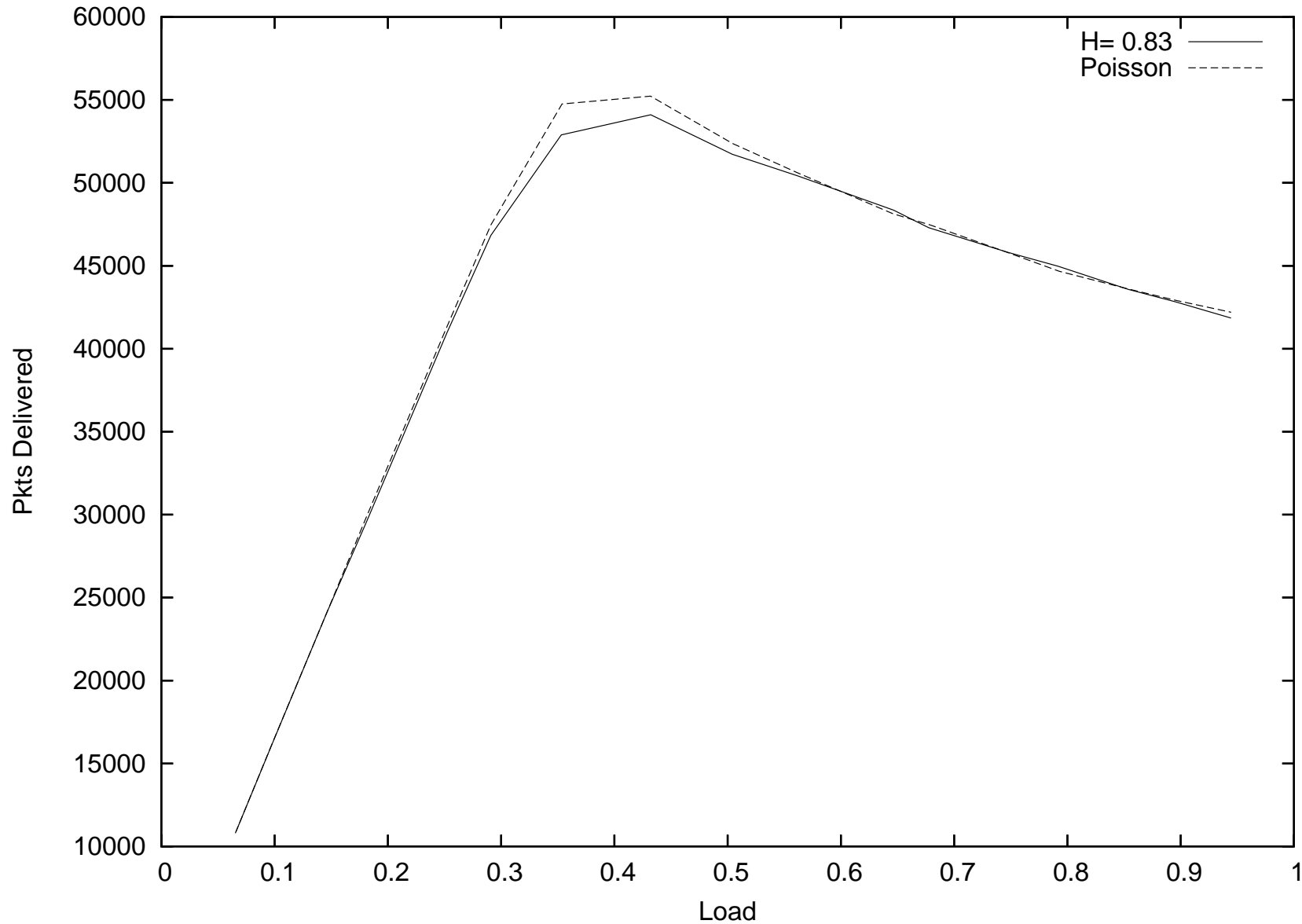
S — Source

D — Dest

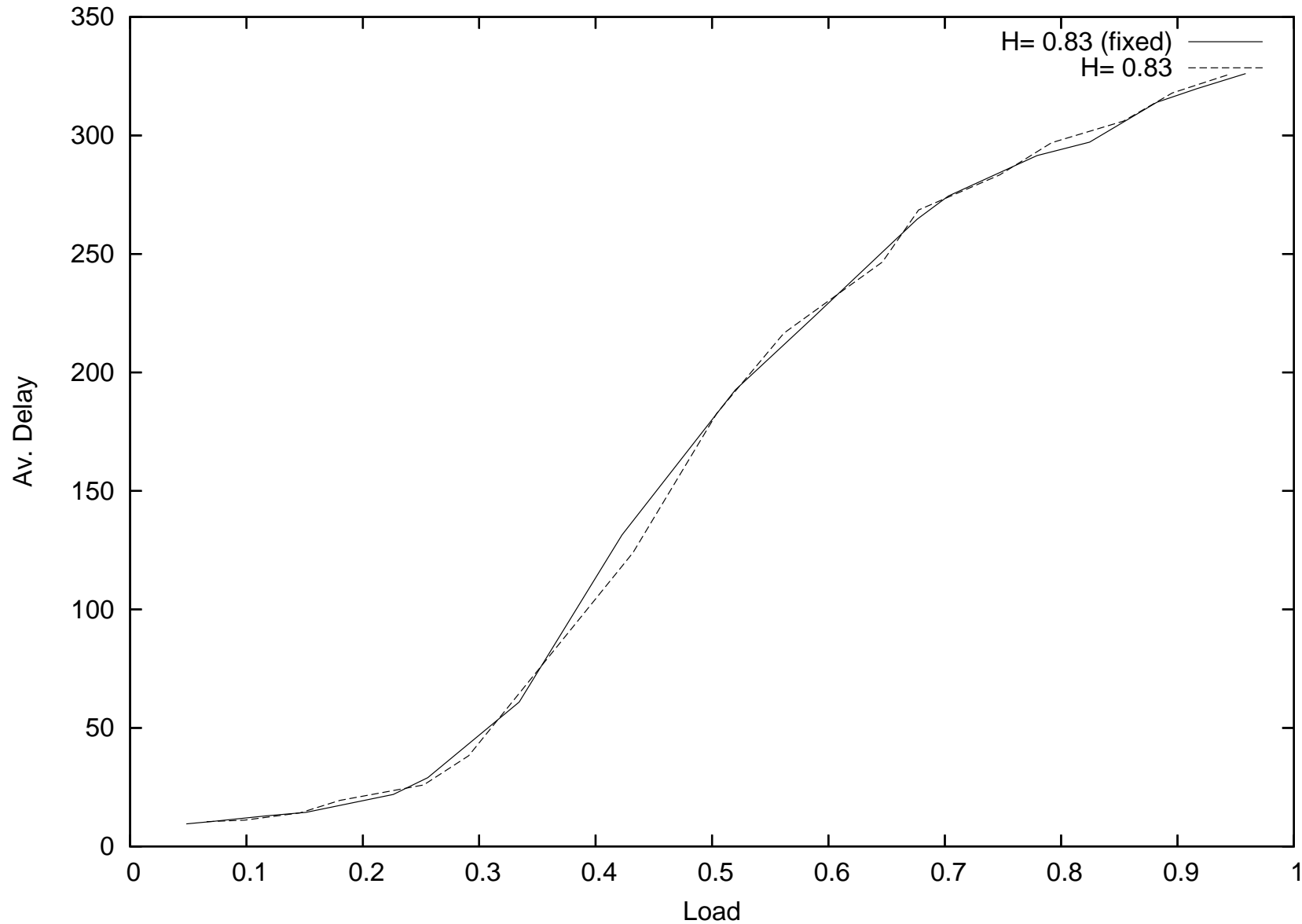
# Manhattan Network – Delays



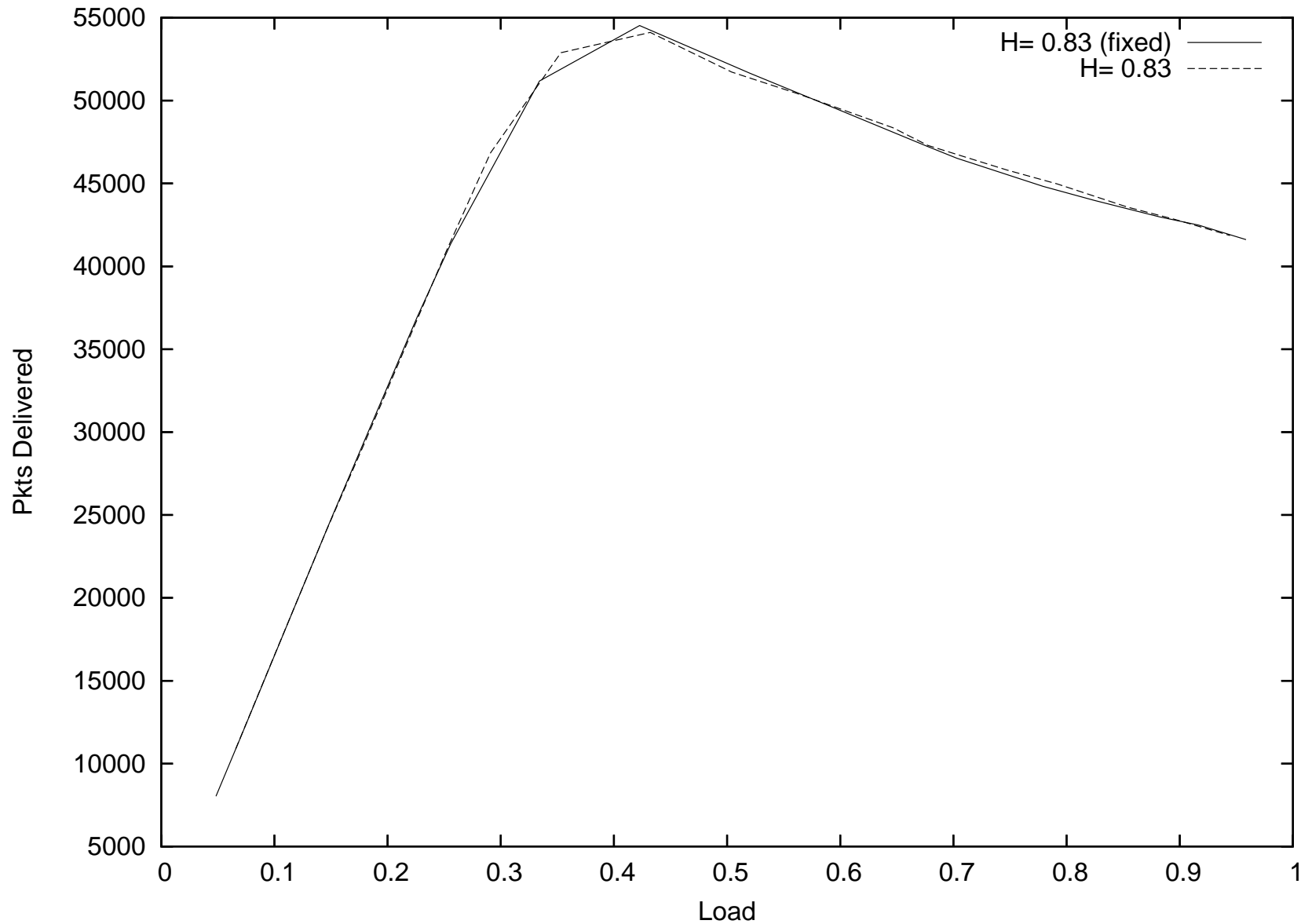
# Manhattan Network – Packets Delivered



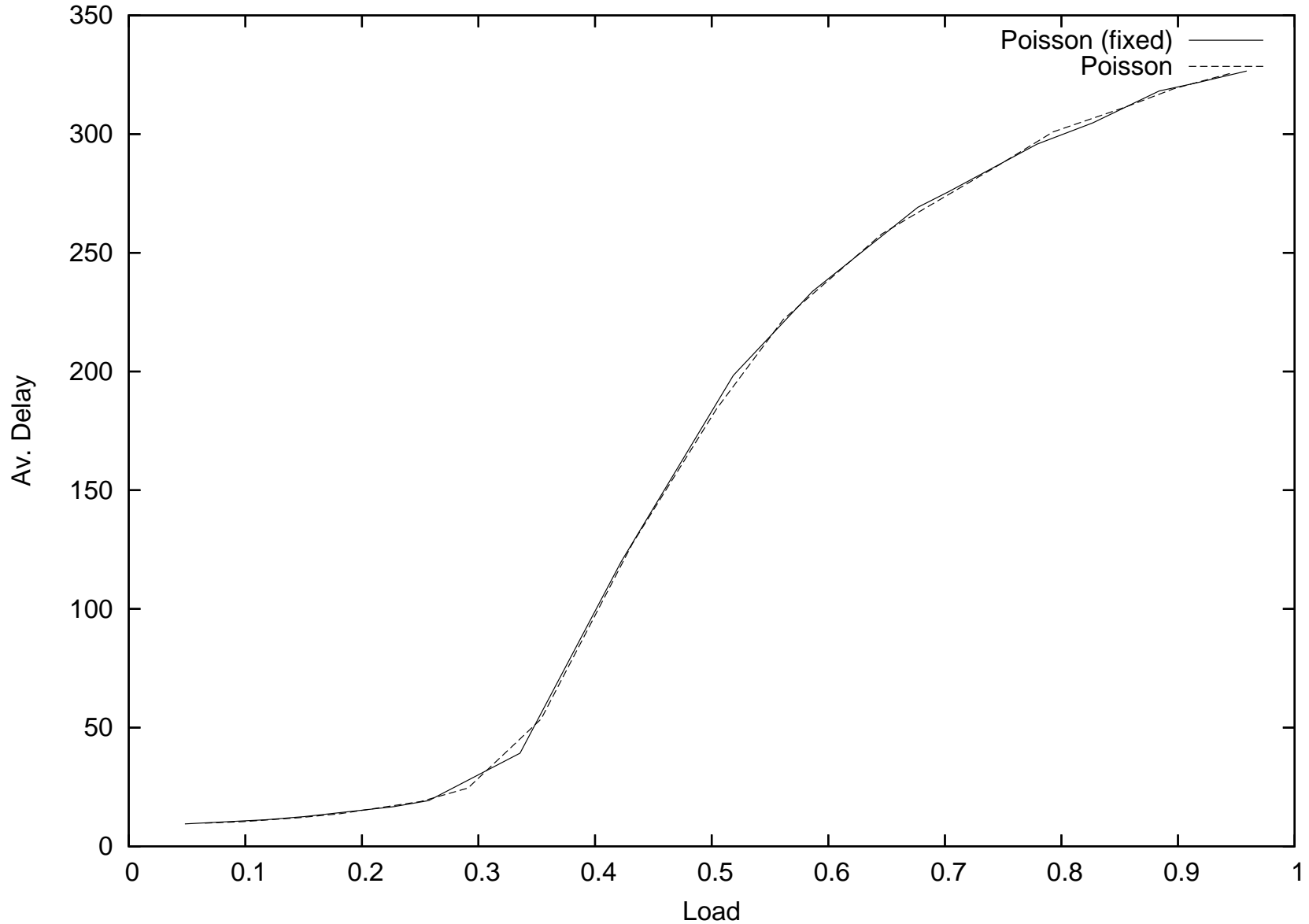
# Delays – Fixed Routing (LRD)



# Packets Delivered – Fixed Routing (LRD)

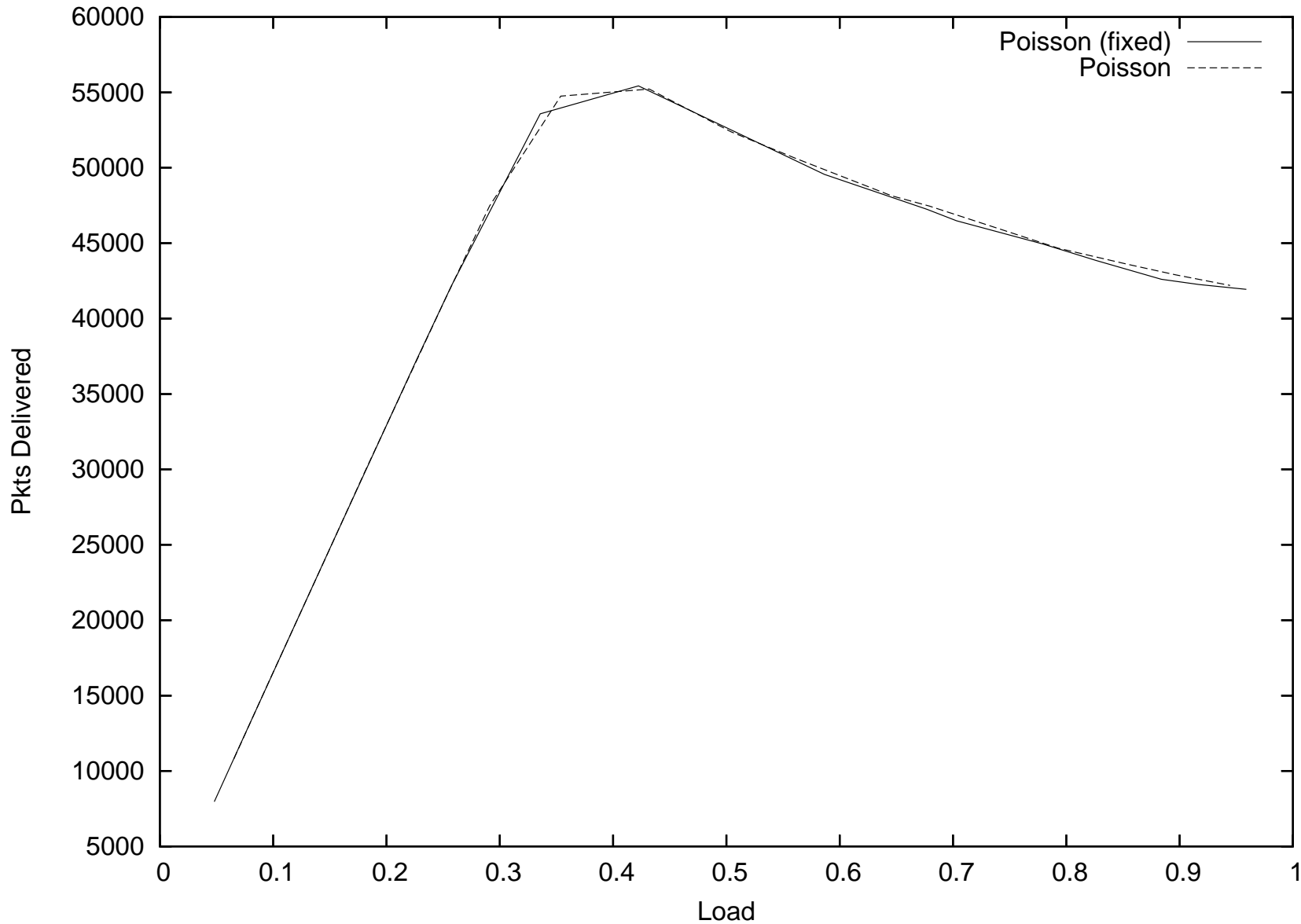


# Delays – Fixed Routing (Poisson)





# Packets Delivered – Fixed Routing (Poisson)



# Hurst Parameters from Simulation

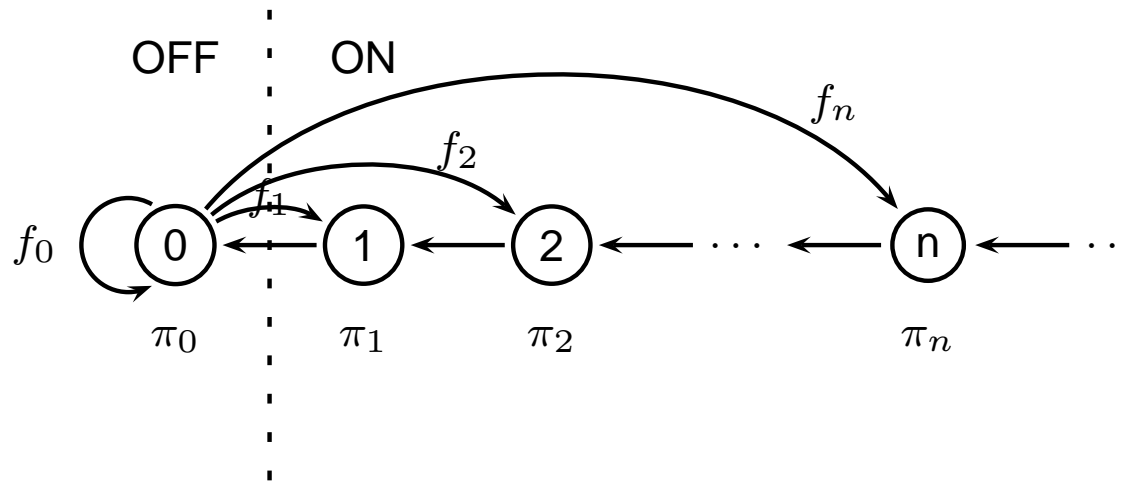
| Route | H     | Load  | R/S  | M. R/S | Var    | L.W.   | Wlts   |
|-------|-------|-------|------|--------|--------|--------|--------|
| Yes   | 0.87  | 0.098 | 0.91 | 0.89   | 0.78   | 1.04 * | 1.01 * |
| Yes   | 0.87  | 0.18  | 0.96 | 0.98   | 1.02 * | 1.35 * | 1.16 * |
| Yes   | 0.87  | 0.29  | 0.95 | 0.98   | 1.02 * | 1.40 * | 1.37 * |
| Yes   | Pois. | 0.098 | 0.73 | 0.72   | 0.51   | 0.76   | 0.75   |
| Yes   | Pois. | 0.18  | 0.78 | 0.78   | 0.59   | 1.00 * | 0.95   |
| Yes   | Pois. | 0.29  | 0.93 | 0.91   | 0.87   | 1.28 * | 1.30 * |
| No    | 0.87  | 0.15  | 0.91 | 0.97   | 0.94   | 1.23 * | 1.12 * |
| No    | Pois. | 0.15  | 0.73 | 0.79   | 0.65   | 0.95   | 0.95   |

Time series is av. queue size at hosts. \* Value out of range.

# Maths Modelling

- To consider the issues with TCP feedback, mathematical models are being developed.
- A model based upon Markov Chains to generate LRD has been created.
- It is hoped to extend existing TCP models.
- By investigation of these models, it is hoped that insight into TCP feedback mechanisms can be provided.
- The Markov model is an interesting and simple model taking two parameters,  $H$  and  $\mu$  and generating traffic with a given Hurst parameter and mean.

# The Infinite Markov Model



$$\mathbf{P} = \begin{bmatrix}
 f_0 & f_1 & f_2 & \dots & f_n & \dots \\
 1 & 0 & 0 & \dots & 0 & \dots \\
 0 & 1 & 0 & \dots & 0 & \dots \\
 0 & 0 & 1 & \dots & 0 & \dots \\
 \vdots & \vdots & \vdots & \ddots & \vdots & \ddots
 \end{bmatrix}$$

# Inducing LRD Correlation Structure

- Unbroken runs of  $k$  0s will clearly decay exponentially with  $k$ . The  $f_i$  values set the decay of unbroken runs of 1s.
- Part of a run of  $k$  or more if  $X_t \geq k$ .
- Control decay of  $\sum_{i=k}^{\infty} \pi_i$ .
- For LRD  $\sum_{i=k}^{\infty} \pi_i \sim Ck^{-\alpha}$ .
- Strict condition  $\sum_{i=k}^{\infty} \pi_i = Ck^{-\alpha}$  for  $k > 0$ .
- It can be proved (details not given here) that if the equilibrium probabilities follow this condition then the time series produced has LRD.
- Since  $\pi_0 = 1 - \sum_{i=1}^{\infty} \pi_i$  then  $C = 1 - \pi_0$ .

# Generating the Correlation Structure

- This system is trivially solved and we can calculate the values of  $f_k$ .
- For  $k > 0$  then (noting problems with some values),

$$f_k = \frac{1 - \pi_0}{\pi_0} [k^{-\alpha} - 2(k+1)^{-\alpha} + (k+2)^{-\alpha}] .$$

- The attractive thing about this series is that it is telescoping. For example.

$$f_0 = 1 - \sum_{i=1}^{\infty} f_i = 1 - \frac{1 - \pi_0}{\pi_0} [1 - 2^{-\alpha}] .$$

- I would be more than happy to discuss this model further with anyone interested.

# Conclusions (1)

- Standard LRD estimators disagree – sometimes wildly. Particularly if the time series used has an “unexpected” nature.
- Because of problems with estimators it is hard to know if different types of traffic have different Hurst parameters and this problem is worse for those types with low traffic levels.
- Considering the small amount of UDP traffic, it is hard to believe the VBR explanation is likely at the York site.
- Would we care in any case? It seems that Poisson traffic aggregates to LRD in practice so does the nature of the source itself matter too much?

# Conclusions(2)

- Routing was tested under quite extreme conditions. Fixed routes were compared with “oscillating” routes which selected the least congested node accurately.
- Routing choices appeared to make no discernable difference to the simulation (these results should be considered preliminary however).
- I still lack insight into whether the topology of the network affects LRD (simulation models on different topologies will follow).
- I still lack insight into how TCP feedback mechanisms might affect LRD (mathematical models may help).



# Bibliography

## References

- [1] J. Cao, W. S. Cleveland, D. Lin, and D. X. Sun. Internet traffic tends toward Poisson and independent as the load increases. In C. Holmes, D. Denison, M. Hansen, B. Yu, and B. Mallick, editors, *Nonlinear Estimation and Classification*. Springer, 2002.
- [2] A. Erramilli, O. Narayan, and W. Willinger. Experimental queueing analysis with long-range dependent packet traffic. *IEEE/ACM Trans. on Networking*, 4(2):209–223, 1996.
- [3] W. E. Leland, M. S. Taqqu, W. Willinger, and D. V. Wilson. On the self-similar nature of Ethernet traffic. In D. P. Sidhu, editor, *Proc. ACM SIGCOMM*, pages 183–193, San Francisco, California, 1993.
- [4] M. S. Taqqu, W. Willinger, and R. Sherman. Proof of a fundamental result in self-similar traffic modelling. *Computer Comm. Rev.*, 27:5–23, 1997.
- [5] W. Willinger, M. S. Taqqu, R. Sherman, and D. V. Wilson. Self-similarity through high-variability: statistical analysis of Ethernet LAN traffic at the source level. *IEEE/ACM Trans. on Networking*, 5(1):71–86, 1997.