

Analysis of Internet Data

Tony Field, Peter Harrison and Uli Harder
Imperial College London

Network layers

Network connection between user application and the transmission media are divided into layers (OSI model)

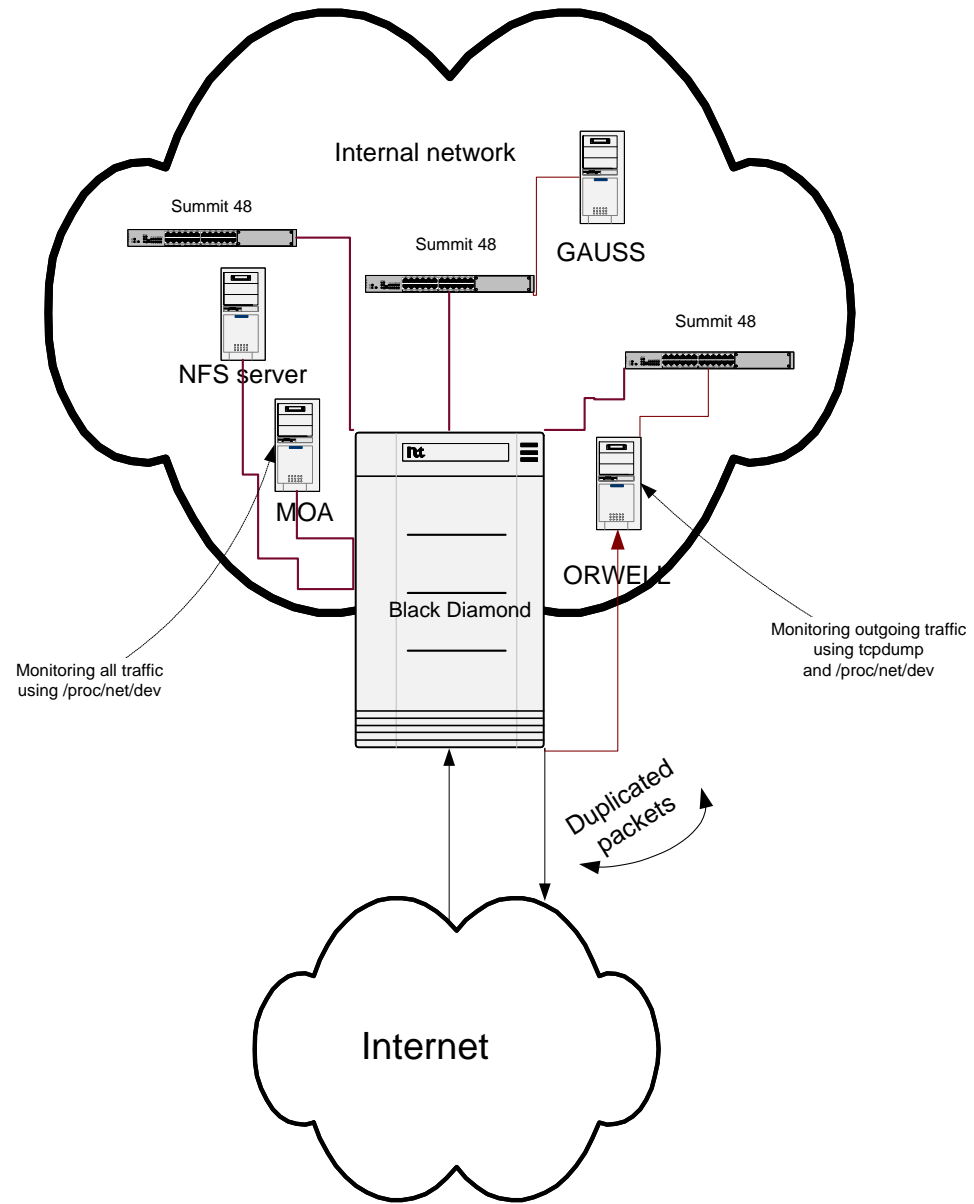
1. Physical (Ethernet cable)
2. Data link (frame format, MAC, single link)
3. Network (multiple links)
4. Transport (IP)
5. Session (TCP/UDP)
6. Presentation
7. Application (e-mail, http, ...)

1 Monitoring

- tcpdump on single machine and on dedicated machine using mirrored ports and MAC level replication
- GILK
- /proc/net/dev
- application logs (web server, and NFS home directories)

Other options

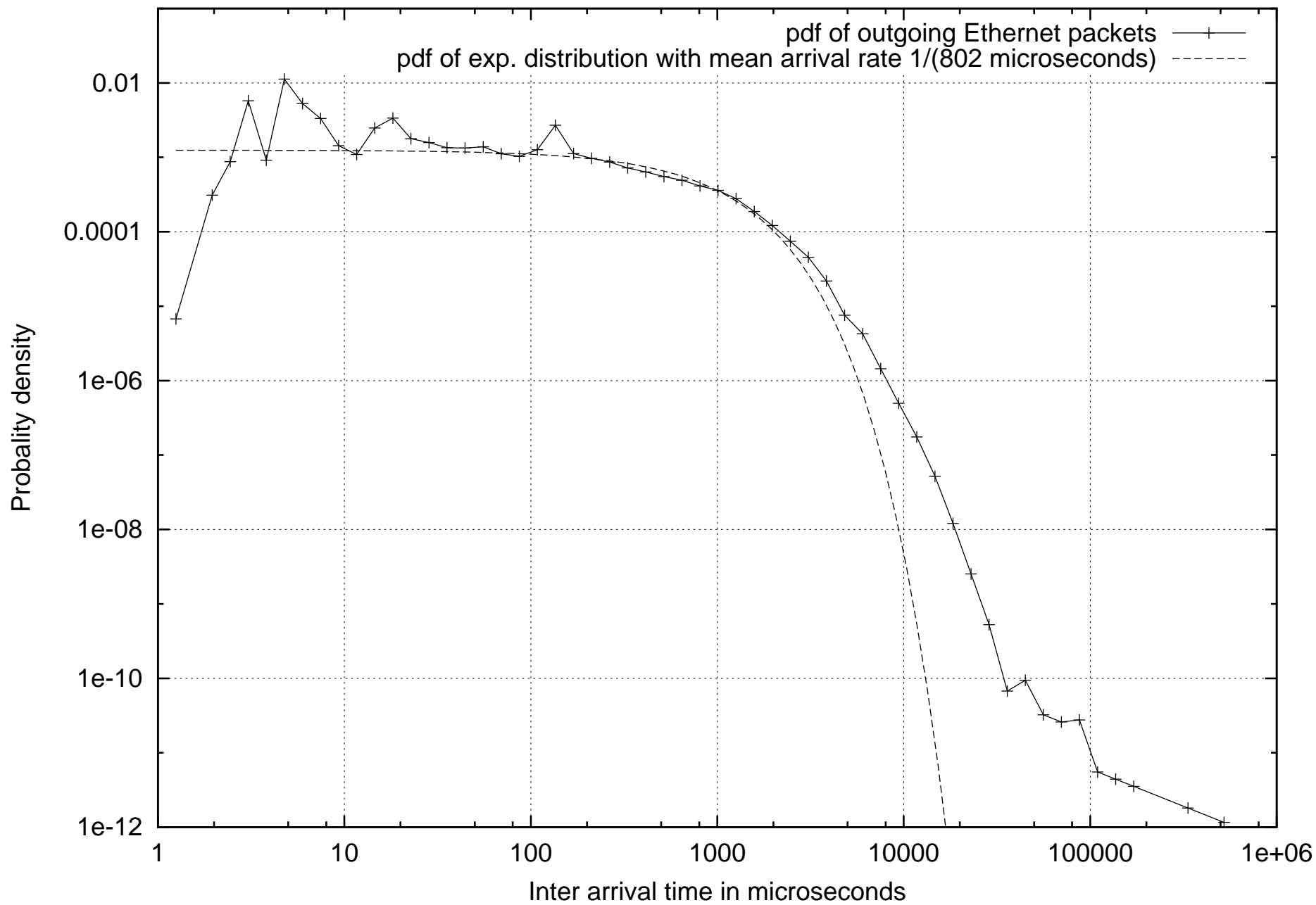
- switch summary data
- ping times
- router information

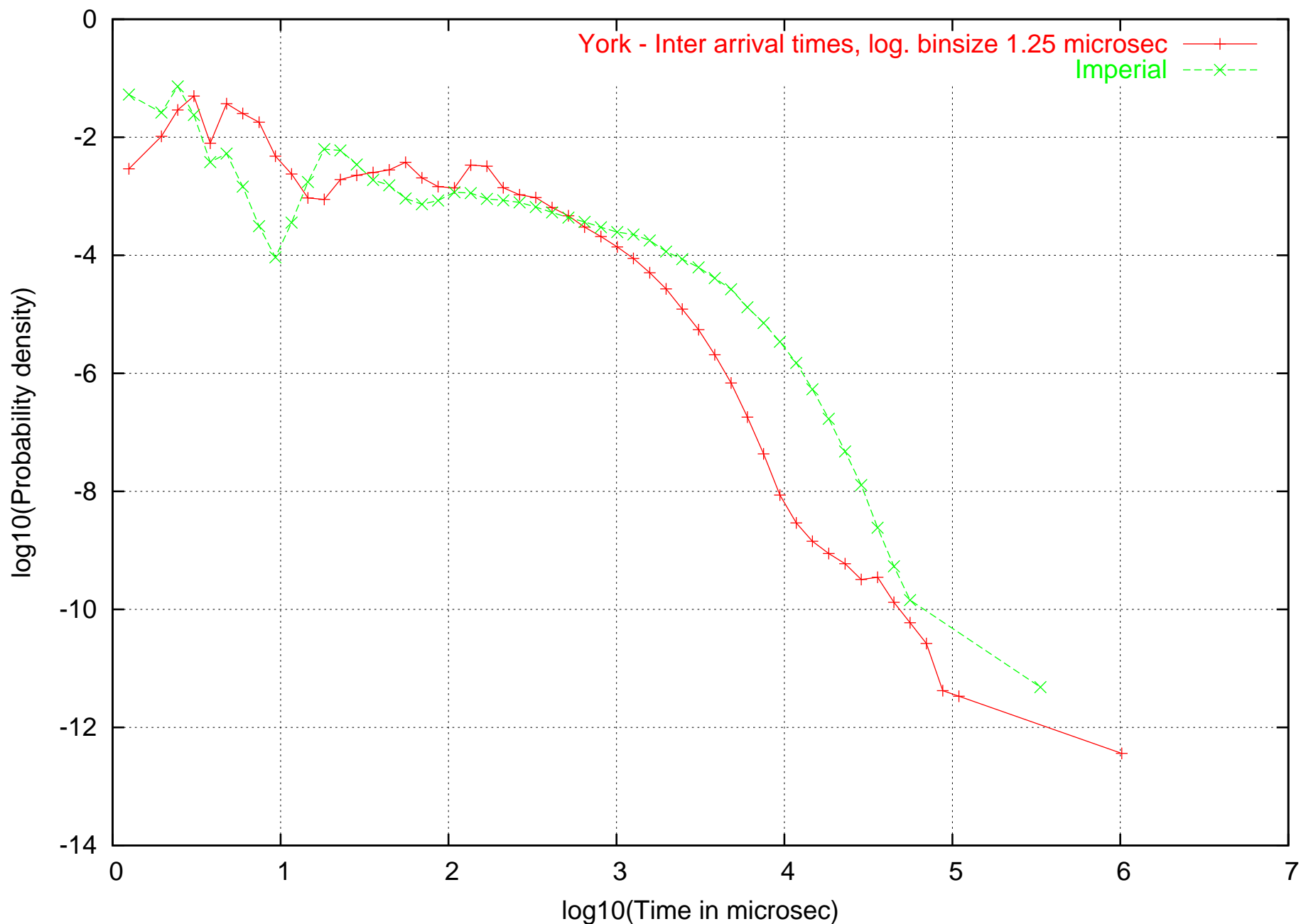


Inter arrival time histograms

- `tcpdump` provides detailed data on each packet seen by the data
- This is essentially information about a point process.
- What is the distribution of inter arrival times?
- If events happen at times $t_i, i \in I \subseteq \mathbb{N}$. The event that occurs at time t_i is called E_{t_i} .
- There are $n \in \mathbb{N}$ events, the first happening at t_1 and the last one at t_n . The observation period may begin before the first event and end after the last, so we define it to be $T = [t_0, t_{n+1}] \subset \mathbb{R}$ with $t_0 \leq t_1 \leq \dots \leq t_n \leq t_{n+1}$, for arbitrary t_0, t_{n+1} bounding the set of event-instants.
- The inter-event times, $\Delta t_i, 1 \leq i \leq n - 1$, are defined as

$$\Delta t_i = t_{i+1} - t_i$$





Power laws

- Suppose a time series $X(t)$ exhibits the scaling law $X(t\alpha) = g(\alpha)X(t)$ for some function $g(\alpha)$. Then $X(t) = bg(t)$ and $g(\alpha) = \alpha^c$, for real constants b, c , is the only non-trivial
- This property is related to what is known as self-similarity.
- a pdf exhibits a power law if

$$p(x) \propto \beta x^\gamma$$

as $x \rightarrow \infty$, for $\beta > 0, \gamma < -1$.

- Gutenberg-Richter, mass extinction events, sand piles, rice piles

Aggregation

Observation period T divided into N contiguous intervals of size $T_N = T/N$. In each interval count the number (or property) of events, so the time series consists

of N values

$$a_i = \left| \{E_t | t_0 + iT_N \leq t < t_0 + (i+1)T_N\} \right|.$$

for $i = 1, 2, \dots, N$. Sometimes it is preferred to use the quantity $A_i = a_i/T_N$. For data gathered with `/proc/net/dev` this is not the necessary.

Power spectrum and auto correlation function

- For a time series $X(t)$ with zero mean the auto-correlation function at lag τ is defined as

$$C(\tau) = \lim_{T \rightarrow \infty} \frac{1}{2T} \int_{-T}^T dt X(t + \tau) X(t).$$

- The power spectrum $S(f)$ is the Fourier transform of the ACF, indicating how much signal/noise is created by what frequency.
- The power spectrum is related to the time original time series by the

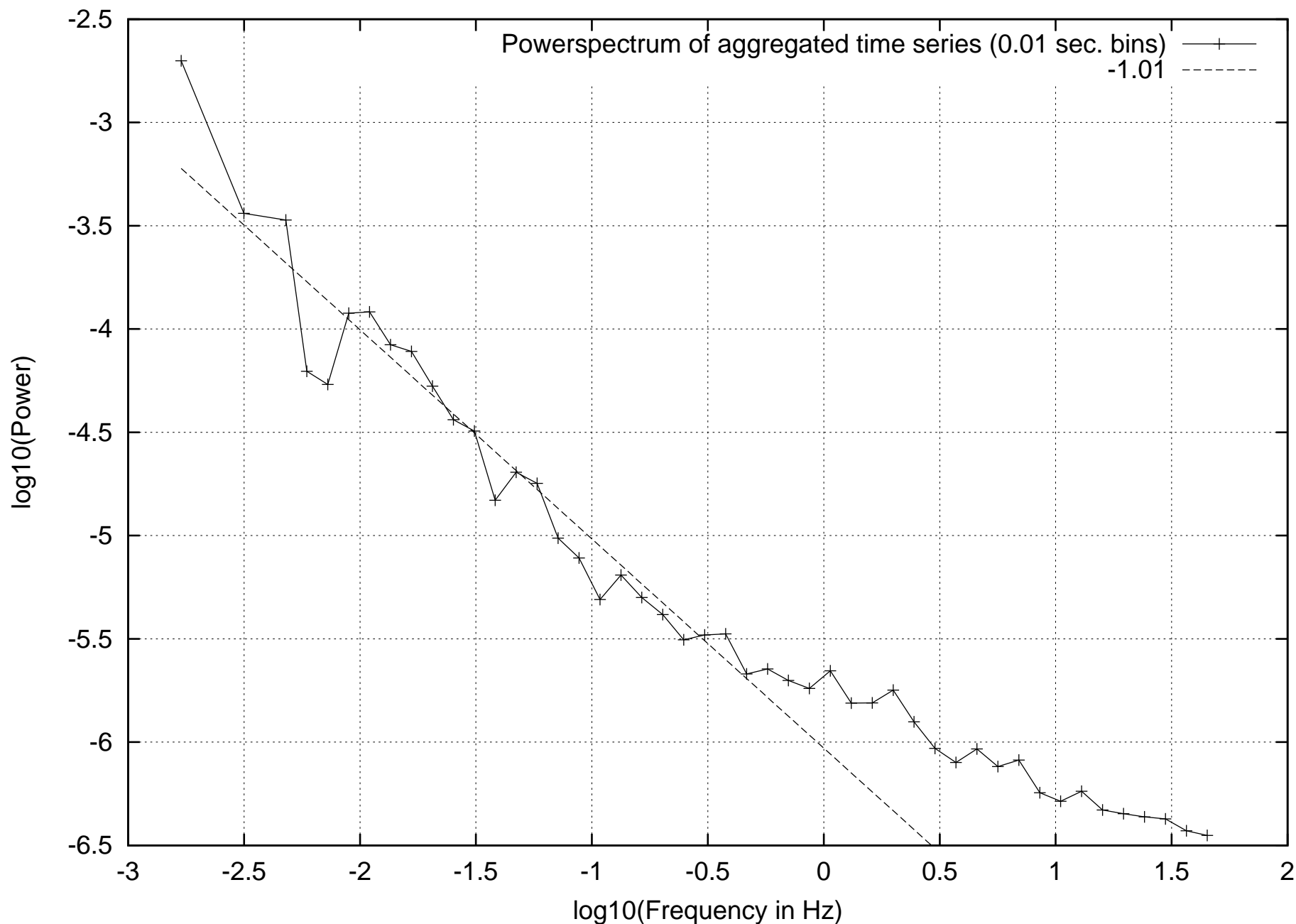
Wiener-Khinchine theorem:

$$S(f) = \lim_{T \rightarrow \infty} \frac{1}{4\pi T} \left| \int_{-T}^T dt X(t) e^{-i2\pi f t} \right|^2.$$

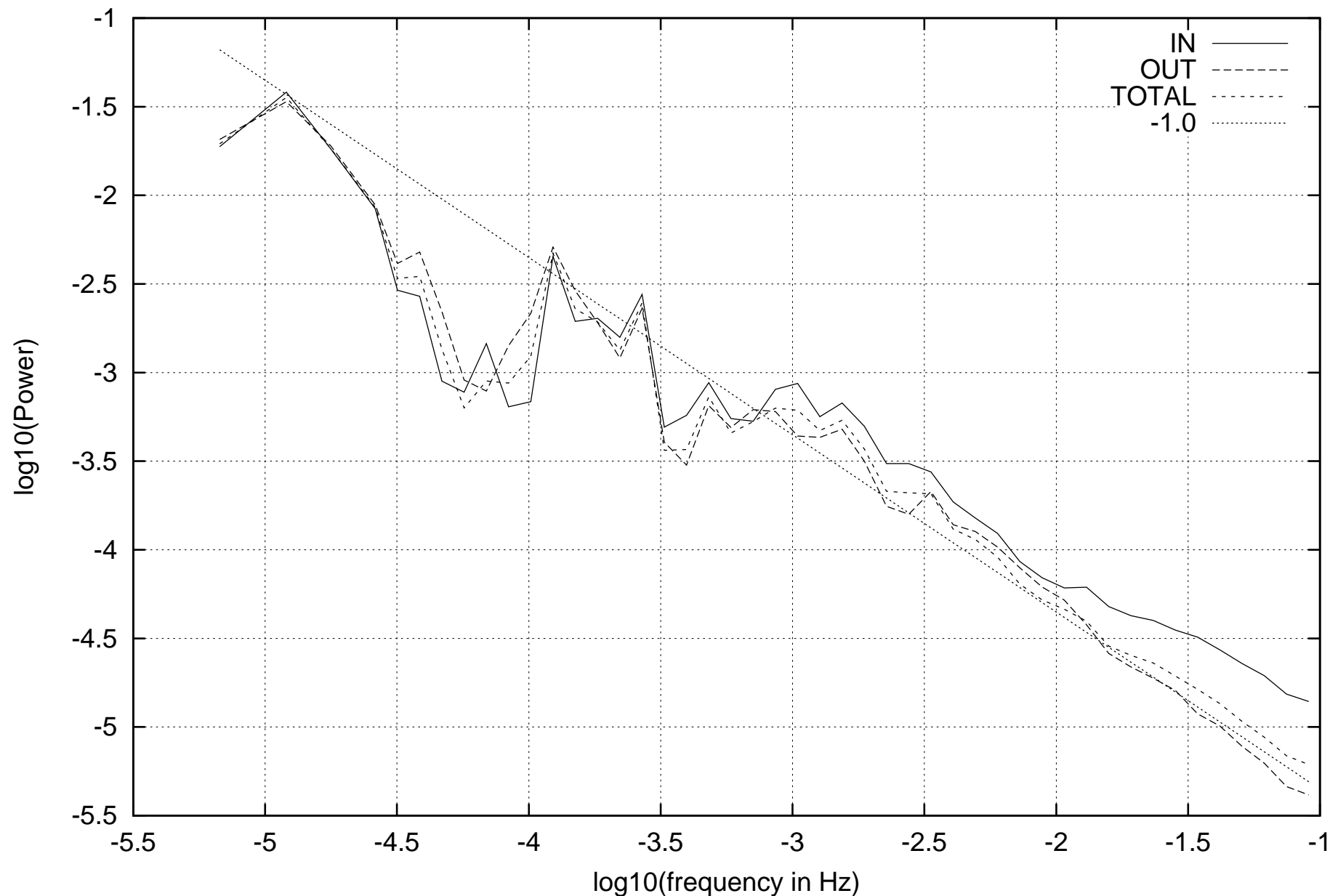
- In fact $S(f)$ is not a consistent estimator. To make it consistent we use filters or windows.
- The power spectrum exhibits a power law if $S(f)$ behaves like $S(f) \propto 1/f^\alpha$, where f is the frequency. (1/f noise)
- The exponent α turns out to be 0 for white noise and 2 for a Brownian motion.
- From the relation of the power spectrum to the auto-correlation function

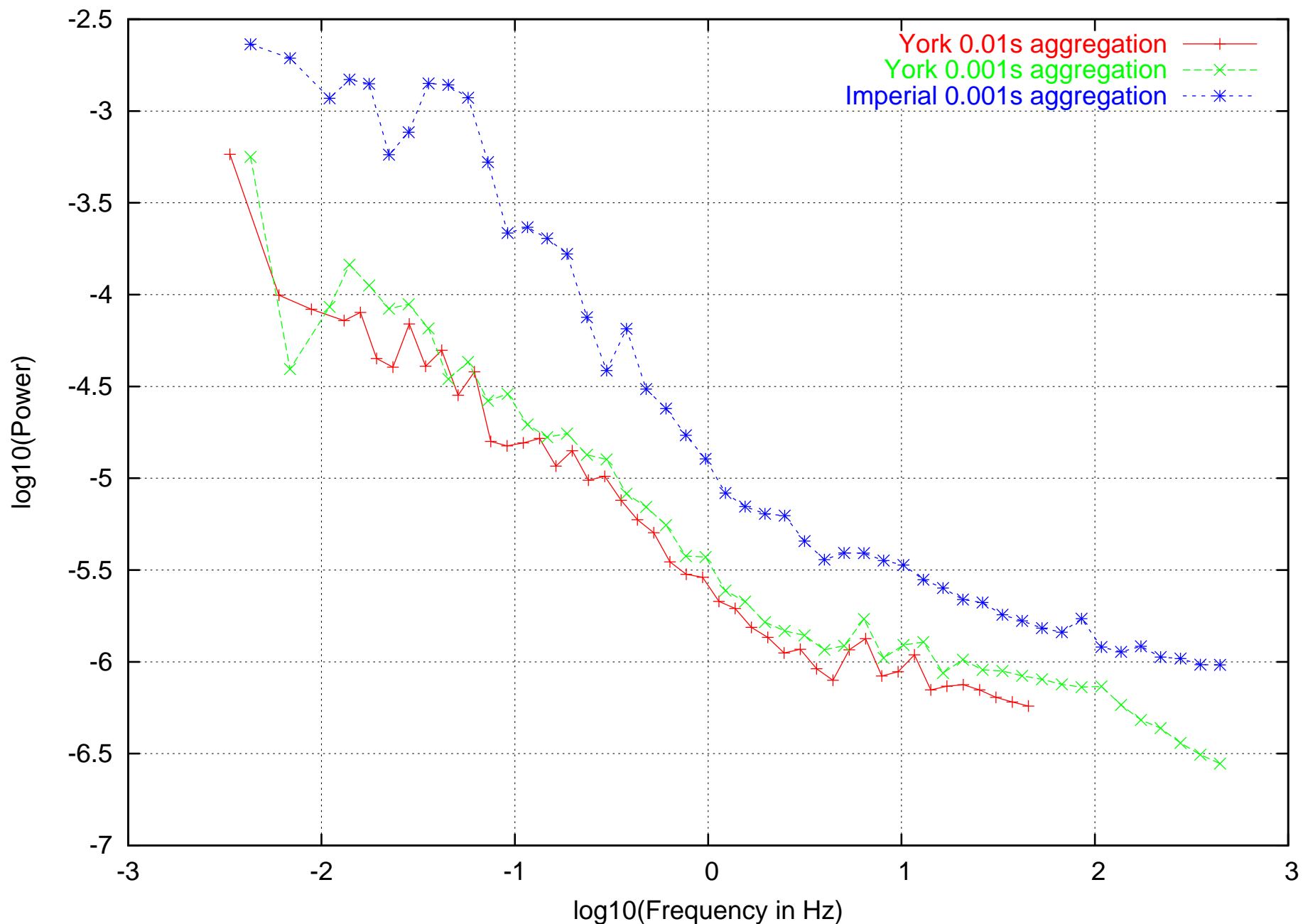
$$C(\tau) \propto |\tau|^{\alpha-1} \quad \text{for } 0 < \alpha < 1$$

- it also follows that an exponent α close to but smaller than 1 corresponds to long term correlations.



Power spectra of the network traffic





Some heavy tailed distributions

- Zipf's law

$$P(\text{file or request size} > x) \approx \frac{1}{x}.$$

- One pdf that can exhibit this behaviour is the Pareto distribution

$$p(x) = \alpha k^\alpha x^{-\alpha-1},$$

where $\alpha, k > 0$ and $x \geq k$. If $\alpha = 1$ the Pareto distribution shows the behaviour of the Zipf law for large x . In a double logarithmic plot, this distribution is a straight line with gradient $-(1 + \alpha)$.

- The symmetric Cauchy (aka Lorentz or Breit-Wigner) distribution has a pdf given by

$$p(x) = \frac{1}{\pi} \frac{s}{s^2 + x^2} \tag{1}$$

- The truncated Cauchy distribution has a pdf $c(x)$ defined by:

$$c(x) = \begin{cases} p(x)/C & 0 \leq x \leq x_{\max} \\ 0 & \text{else} \end{cases}$$

where $p(x)$ is given by eq. 1 and C is a normalisation constant

$$C = \int_0^{x_{\max}} p(x)dx.$$

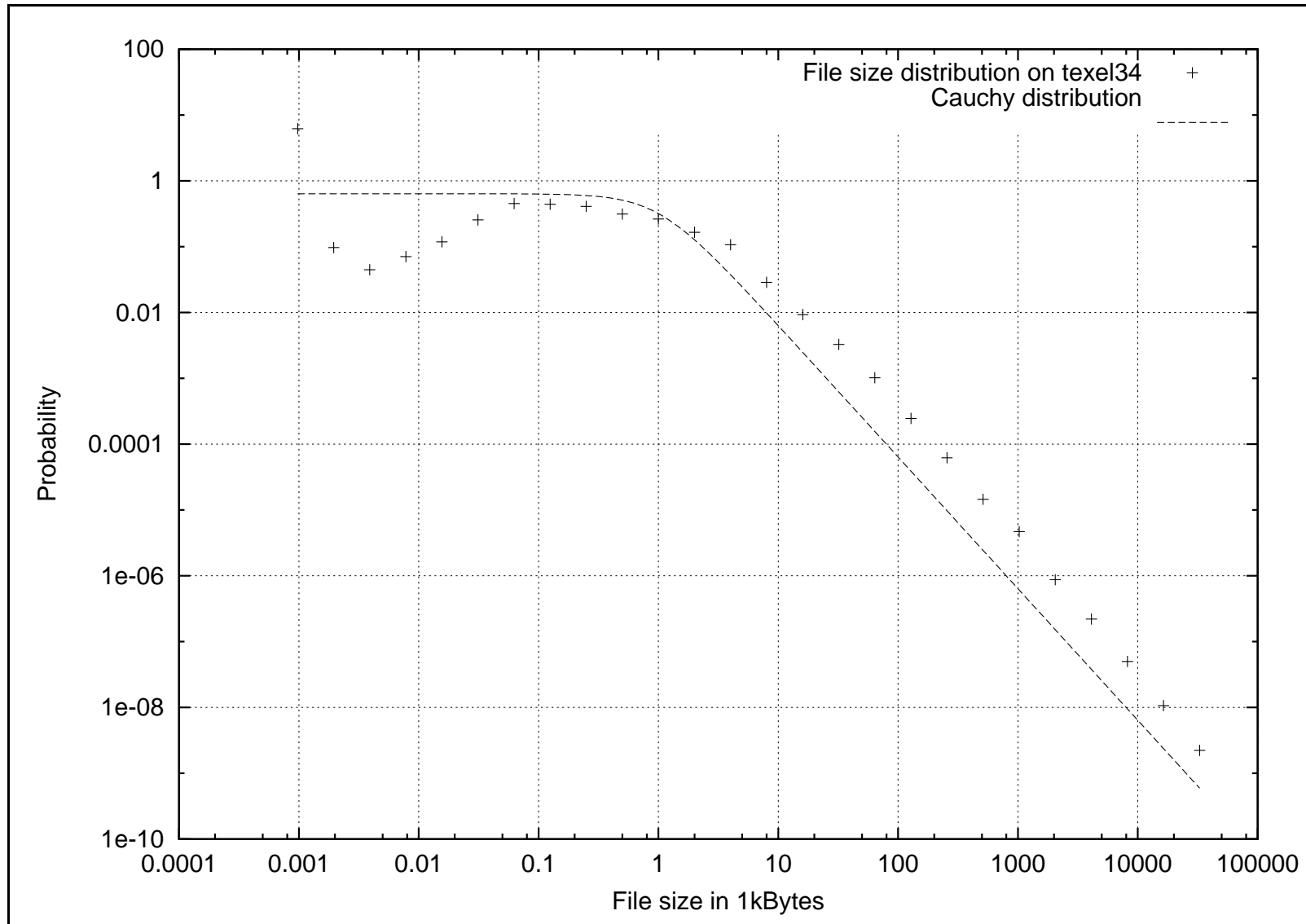
The truncation of the Cauchy distribution gets rid of its usually prohibiting features like infinite moments.

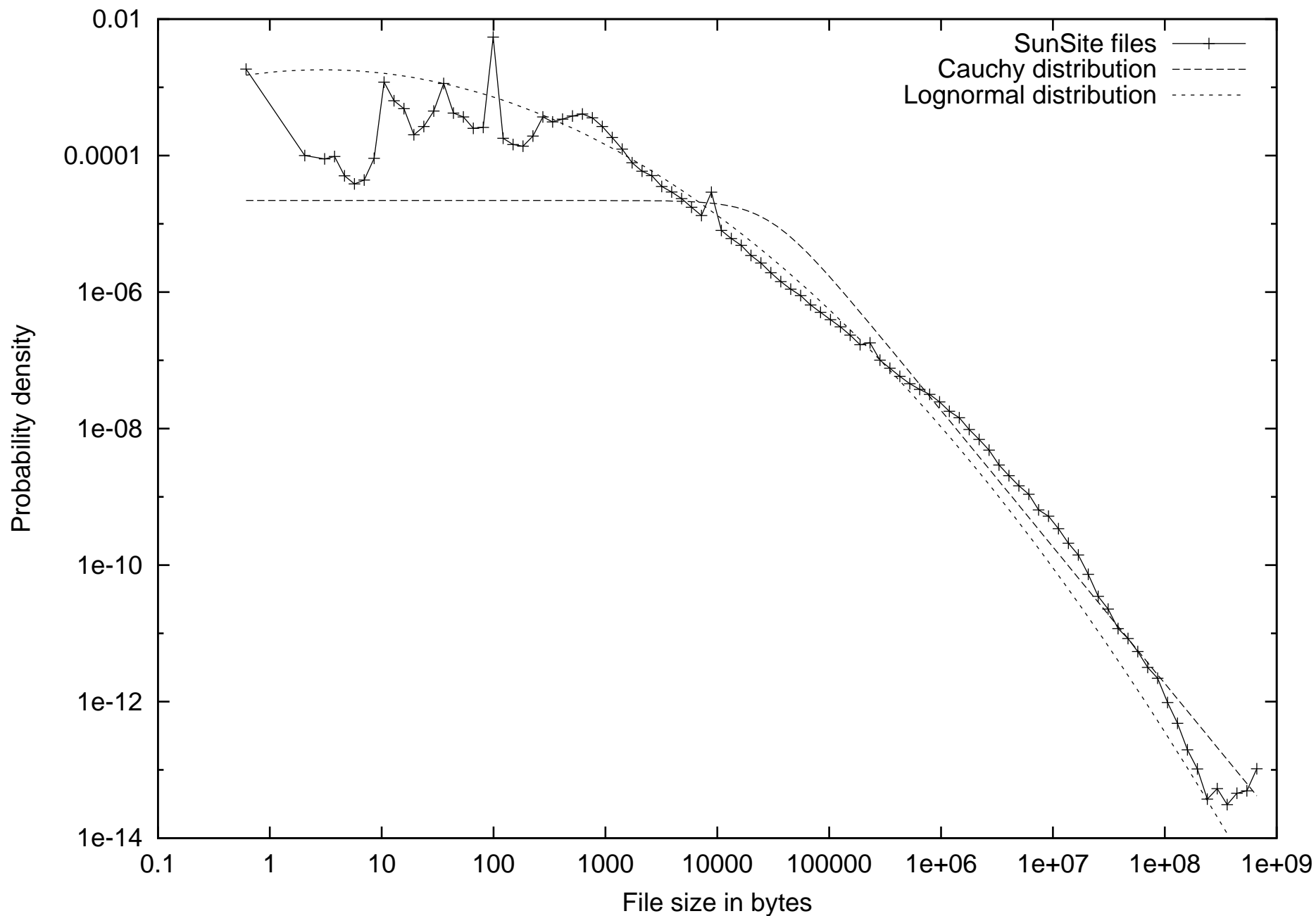
- A strictly stable Lévy distribution with characteristic exponent α is defined by its characteristic function [?]

$$G(k) = e^{-|k|^\alpha \gamma}. \quad (2)$$

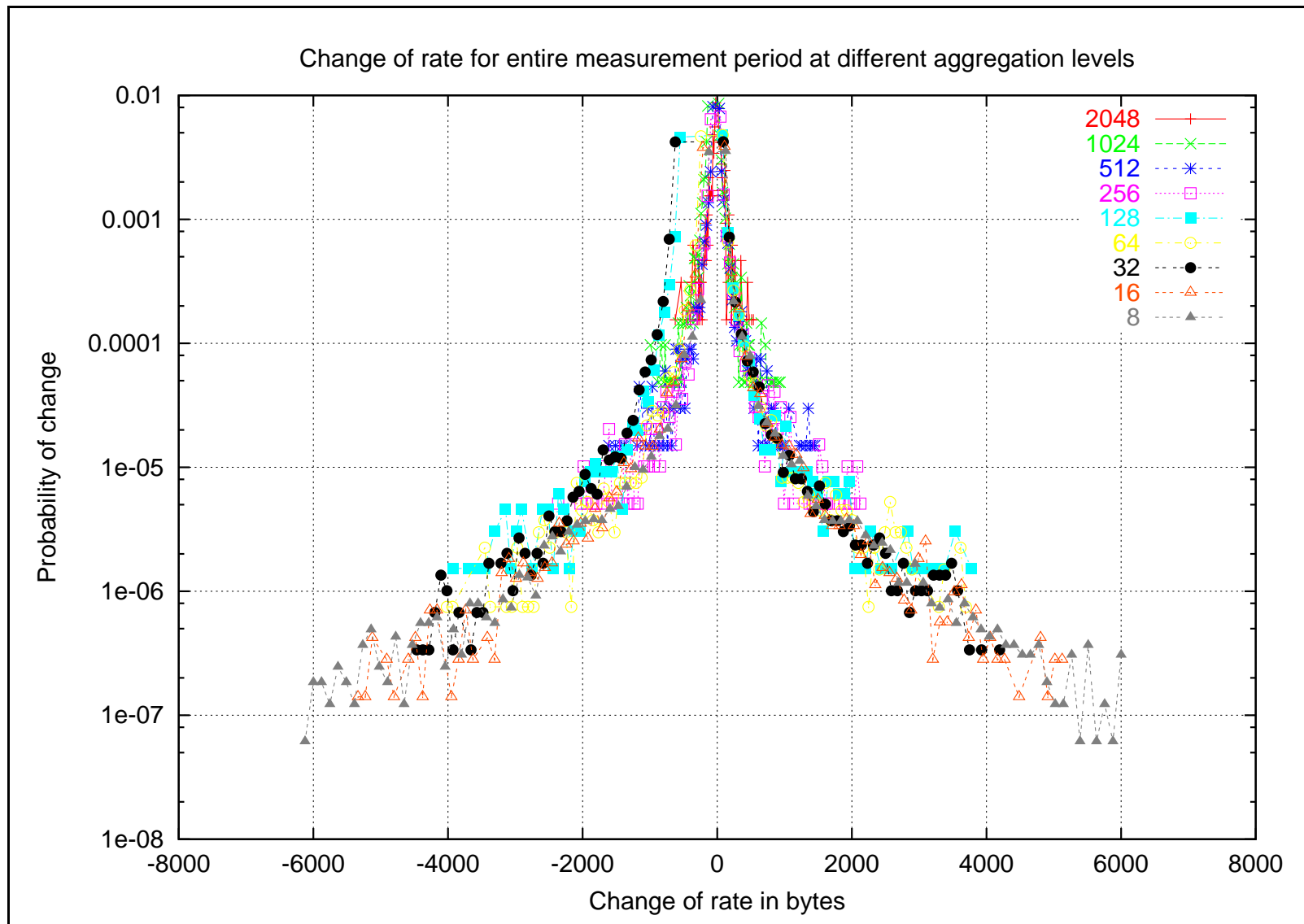
The scale parameter γ has to satisfy $\gamma > 0$. The distributions are stable for $0 \leq \alpha \leq 2$. Ranging from normal to Cauchy distribution

File and request size distribution





Changes in the packet rate



Models

- ON/OFF
- TCP
- SOC
- anything fractal
- lognormal file sizes
- M/G/1, G/G/1, M/P/ ∞ , MMPPs

Our model

- Lévy requests
- chopped to Ethernet packet sizes
- Going into (infinte) buffer

- Released from network “server”
- We measure the departure events. Similar to outgoing network traffic on webserver.

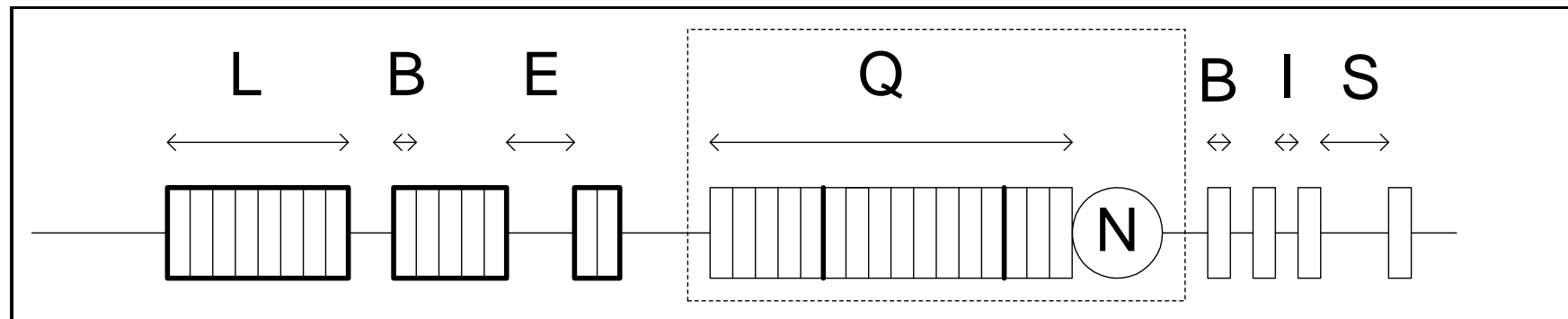
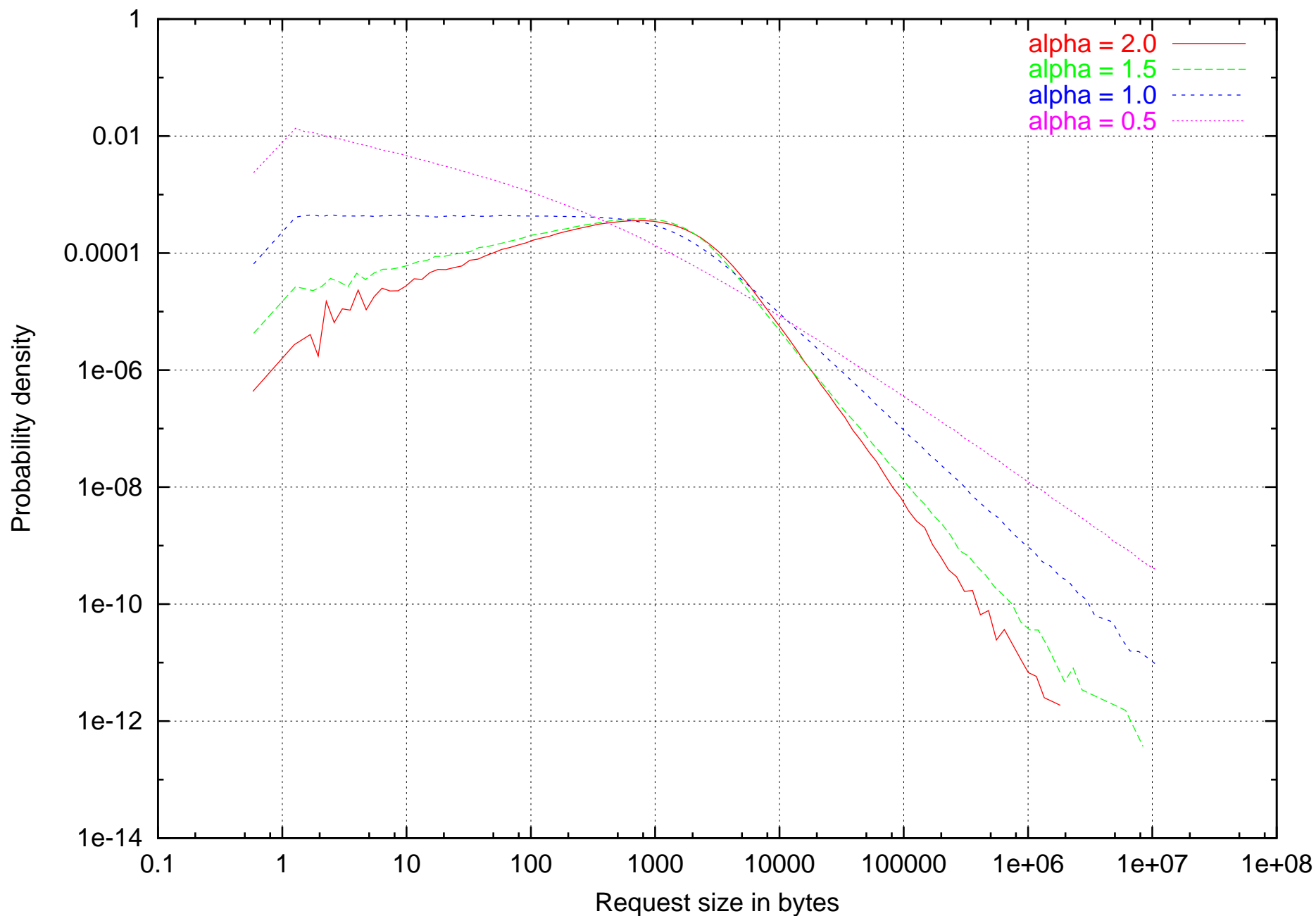
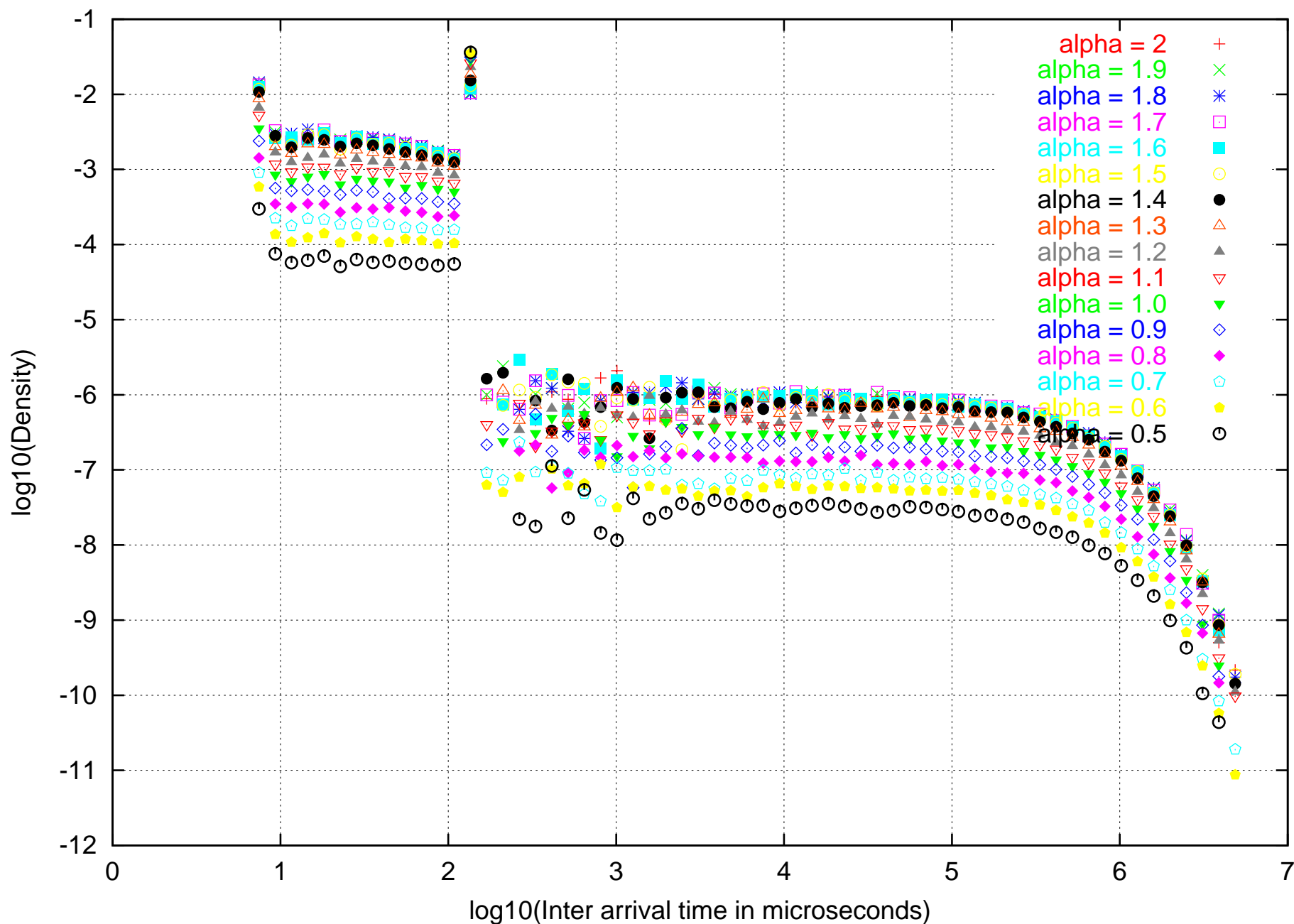
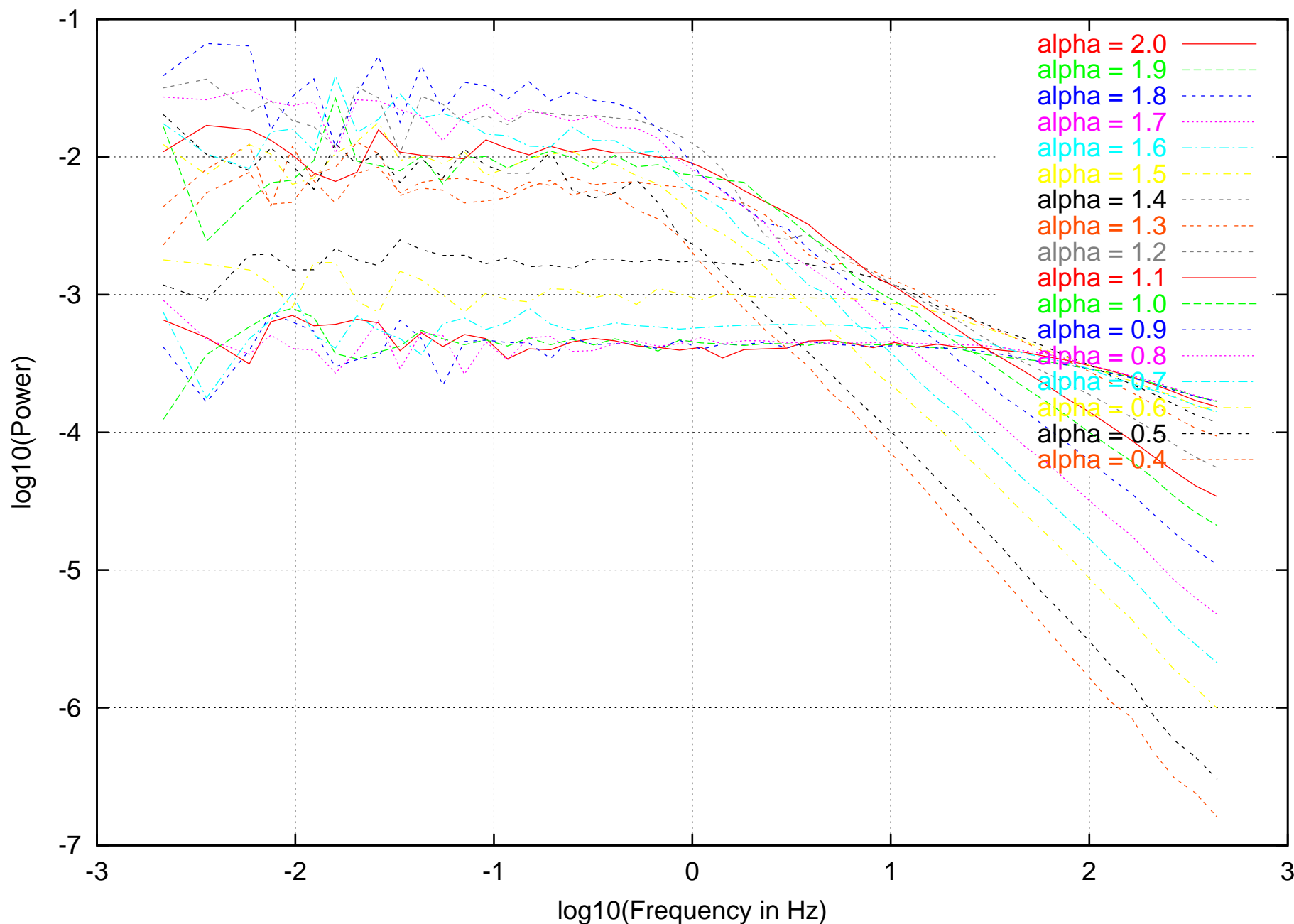


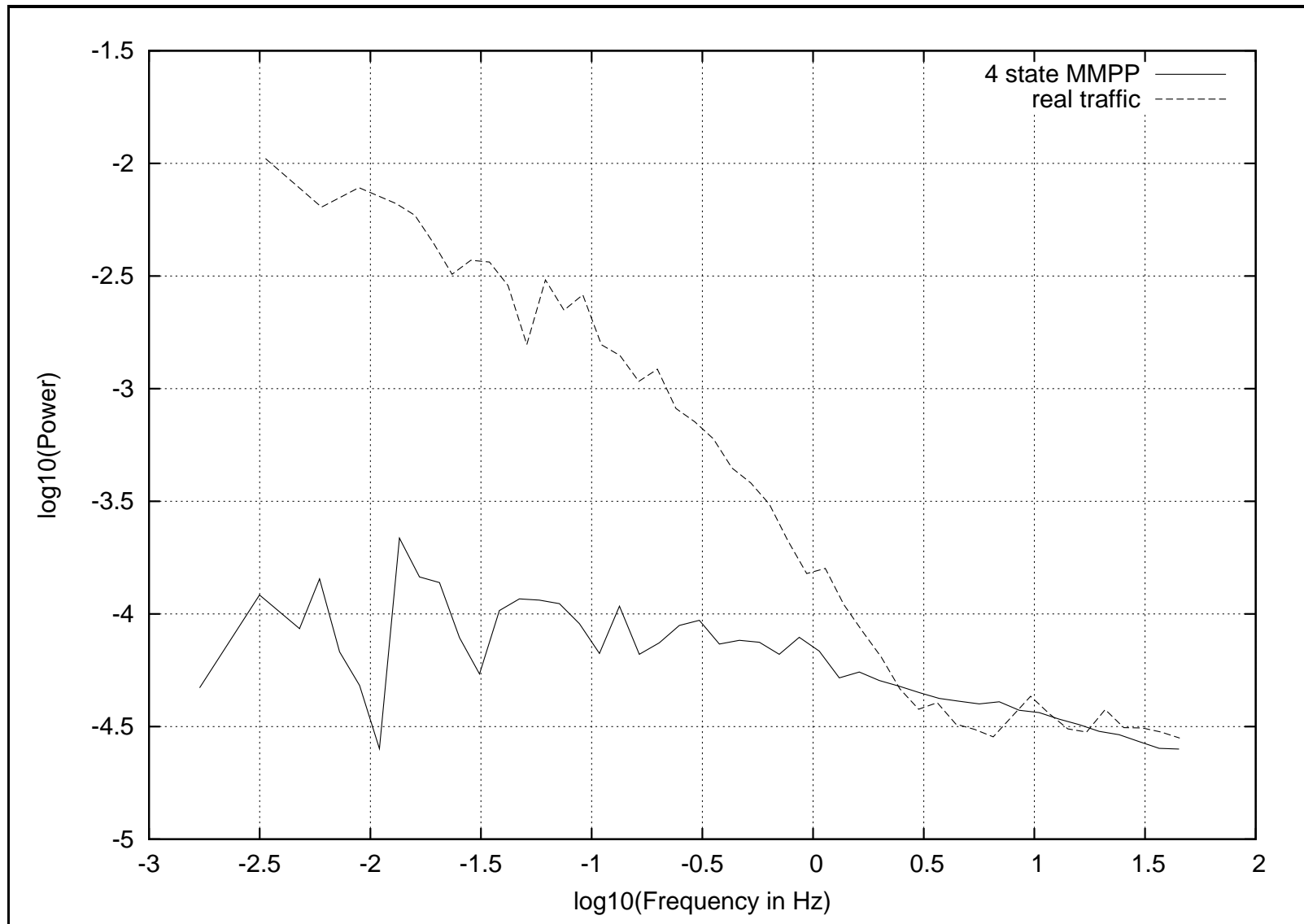
Figure 3: The model







4 state MMPP



Acknowledgements

The research was funded by EPSRC (research grant QUANT, GR/M80826).