Before proceeding to bad news:

$\longrightarrow$  connection between heavy-tailedness and `google`?

$\longrightarrow$  saga of two lucky kids (aka "grad students")

$\longrightarrow$  lesson to be drawn?

dave                    goliath

Now, to the bad news!

Bad news #1: queueing



- influx rate (`write`) $<$ outflux rate (`read`)

  $\rightarrow$ else buffer will grow out of bound

- during on-time: if `write` rate $<$ `read` rate

  $\rightarrow$ then what?

  $\rightarrow$ economy dictates opposite (suppose 1/2)

  $\rightarrow$ hence: during on-time buffer grows (McDonald's)

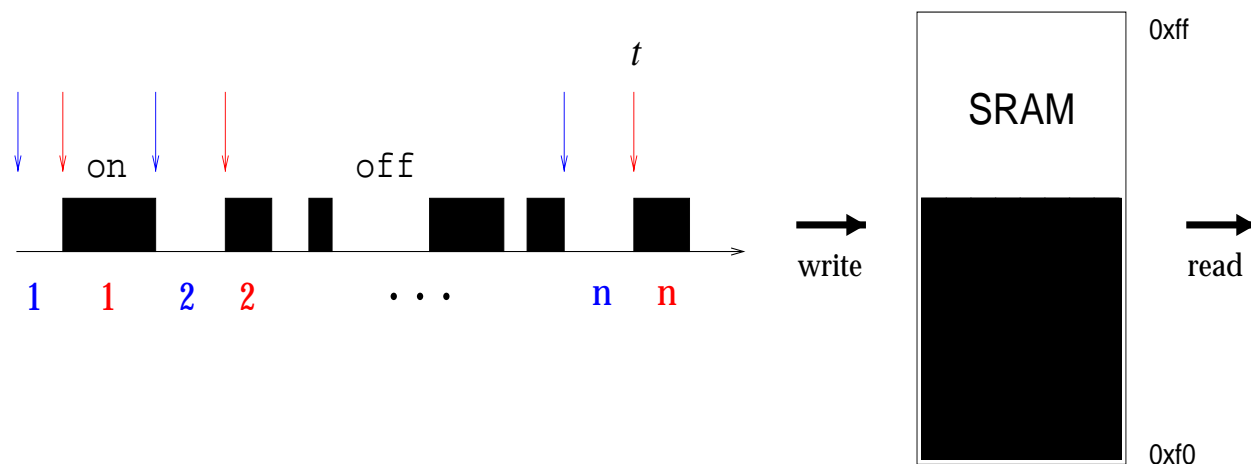Since on/off input is random, so is the buffer/memory occupancy

$\longrightarrow$ at time $t$, could be 10 KB, 120 KB, etc.

$\longrightarrow$ i.e., $\Pr\{Q(t) = 10000\}$ = some value, ...

Want to know: in the long-run $(t \to \infty)$ what is $Q(t)$?

$\longrightarrow$ write as $Q(\infty)$

$\longrightarrow$ practical interest: $\Pr\{Q(\infty) > x\}$ ?

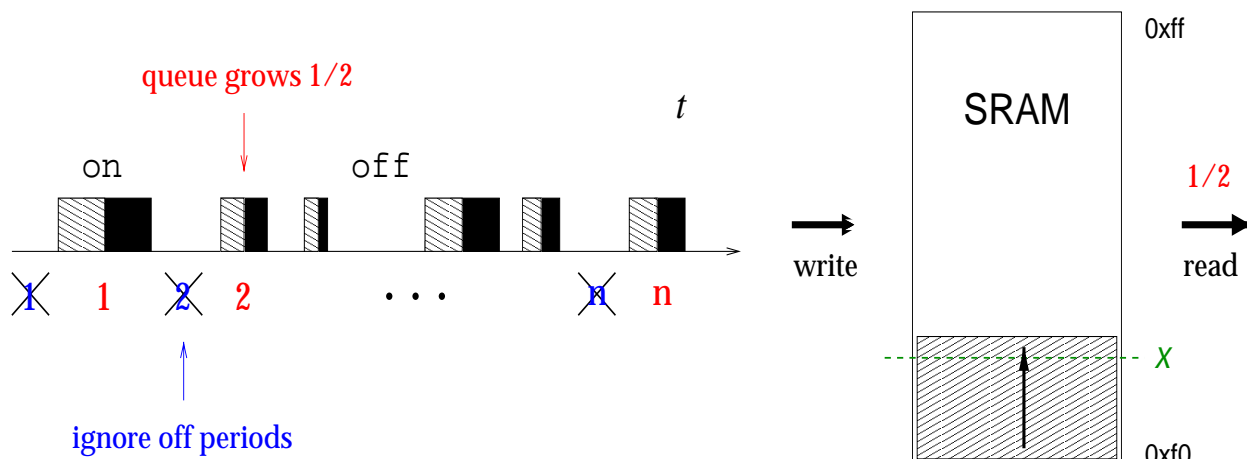$\longrightarrow$ corresponds to excessive delay, buffer loss, etc.

Case I: what shape does $\Pr\{Q(\infty) > x\}$ take when both on and off periods are exponential?

$\longrightarrow$ assume i.i.d. (with $be^{-bt}$)

$\longrightarrow$ first, switch from time unit to count unit

- alternating on/off periods: mutually independent

- how many on and off periods $(n)$ at time $t$?

  $\rightarrow n \approx t/(E[\text{on}] + E[\text{off}]) = t/(\frac{1}{b} + \frac{1}{b}) = bt/2$

  $\rightarrow$ for large $t$

- sum of $n$ on-periods: $S_n$

  $\longrightarrow$ let's upper-bound $\Pr\{Q(t) > x\}$

  $\longrightarrow$ i.e., $\Pr\{Q(t) > x\} < \boxed{?}$

Upper-bounding idea:



- worst-case viewpoint

  → ignore the beneficial effect of off periods

  → McDonald's: groups of people arrive without pause

- thus, to have $Q(t) > x$:

  → $S_n/2 > x$

  → i.e., at the very least

  → hence: $\Pr\{Q(t) > x\} < \Pr\{S_n > 2x\}$

- need to upper bound $\Pr\{S_n > 2x\}$

  $\rightarrow$ for large $x$ (i.e., $2x > nE[\text{on}]$)

  $$\begin{aligned} \Pr\{S_n > 2x\} &= \Pr\{S_n - nE[\text{on}] > 2x - nE[\text{on}]\} \\ &= \Pr\{S_n/n - E[\text{on}] > 2x/n - E[\text{on}]\} \end{aligned}$$

  $\rightarrow$ by LLN $S_n$ is concentrated around its mean!

- we can apply large deviation bound

  $\rightarrow \Pr\{|\frac{S_n(t)}{n} - p| > \varepsilon\} < e^{-an}$

  $\rightarrow$ here: $\varepsilon = 2x/n - E[\text{on}]$

  $\rightarrow$ recall: $a$ depends on $\varepsilon$

- facts: shape of $a(\varepsilon)$

  $\rightarrow$ binary case: $a = \varepsilon \log \frac{\varepsilon}{p} + (1 - \varepsilon) \log \frac{1-\varepsilon}{1-p}$

  $\rightarrow$ exponential case: $a = b\varepsilon - 1 - \log b\varepsilon$

  $\rightarrow$ for large $\varepsilon$ (same as large $x$): $a \approx b\varepsilon$

- apply large deviation bound to $S_n$

$$
\begin{aligned}
\Pr\{S_n > 2x\} &= \Pr\{S_n/n - E[\text{on}] > 2x/n - E[\text{on}]\} \\
&< e^{-an} \\
&\approx e^{-b\varepsilon n} \\
&= e^{-b(2x/n - E[\text{on}])n} \\
&= e^{-2bx + bE[\text{on}]n} \\
&= e^{-2bx + n} \\
&< e^{-bx}
\end{aligned}
$$

$\rightarrow$ for sufficiently large $x$ (used several times)

Thus: $\Pr\{Q(t) > x\} < e^{-bx}$ for large $x$ and $t$

$\longrightarrow \quad \Pr\{Q(\infty) > x\} < e^{-bx}$ for large $x$

$\longrightarrow \quad$ prob. of queue growing large: exponentially small

$\longrightarrow \quad$ for exponential traffic: buffering is effective

$\longrightarrow \quad$ extra buffer/memory $y$ buys a lot:

$$\Pr\{Q(\infty) > x + y\} < e^{-b(x+y)} = e^{-bx}e^{-by}$$

Remarks: analysis method

$\longrightarrow$ con: only holds for large $x$

$\longrightarrow$ pro: very general/powerful

$\longrightarrow$ for exponential case: "excessive force"

$\longrightarrow$ somewhat like "catching fly with a cannon"

$\longrightarrow$ can use more elementary methods

$\longrightarrow$ a course in queueing theory (Markovian input)

$\longrightarrow$ problem: doesn't extend to heavy-tailed input

$\longrightarrow$ but the Internet is heavy-tailed!

Case II: shape of $\Pr\{Q(\infty) > x\}$ when off-period is exponential but on-period is heavy-tailed?

$\longrightarrow$    want to show: $\Pr\{Q(\infty) > x\}$ is heavier as well

$\longrightarrow$    want to contrast with exponential case

$\longrightarrow$    let's lower-bound: $\Pr\{Q(t) > x\} > \boxed{?}$

$\longrightarrow$    why upper-bounding not enough?

Lower-bounding idea:



$\longrightarrow$    sampling viewpoint: ok since i.i.d.

$\longrightarrow$    wait till first long $(> 2x)$ on-period

- how long must one wait?

  $\rightarrow$ on the order of $1/\Pr\{Z > x\} = 1/cx^{-\alpha} \propto x^{\alpha}$

  $\rightarrow$ so, time scale of interest: $t = O(x^{\alpha})$

- number on and off periods at time $t$

  $\rightarrow n \approx t/(E[\text{on}] + E[\text{off}]) = \delta t = O(x^{\alpha})$

  $\rightarrow$ for large $t$ (hence large $x$)

- now: $\Pr\{Q(t) > x\} \approx$ fraction of time during $O(x^{\alpha})$ where queue is bigger than $x$

  $\rightarrow O(x/x^{\alpha}) = O(x^{1-\alpha})$

  $\rightarrow$ where did we apply similar reasoning?

- note: we ignored the contribution of other on periods

  $\rightarrow$ hence: lower-bound

- thus: for large $x$ and $t$

  $\rightarrow \Pr\{Q(t) > x\} > O(x^{1-\alpha})$

  $\rightarrow$ tail $\Pr\{Q(\infty) > x\}$ is at least polynomially heavy

  $\rightarrow$ can also show polynomially upper-bounded

  $\rightarrow$ much more likely to overcrowd

  $\rightarrow$ buffering is not as effective: marginal gain small

  $\rightarrow$ modern view: bandwidth-centric resource provisioning

Remarks:

- heavy-tailed on-times and resultant heavy-tailed queue-ing was a big surprise

  $\rightarrow$ grabbed CS, EE, statistics/probability, OR, some physicists, etc. by surprise!

  $\rightarrow$ huge scientific impact

- one technical aside: for heavy-tailed i.i.d. variables

  $$\Pr\{Z_1 + \cdots + Z_n > x\} = \Pr\{\max\{Z_1, \ldots, Z_n\} > x\}$$

  $\rightarrow$ for large $x$

  $\rightarrow$ when the sum is large, one guy is to blame!

  $\rightarrow$ single long on-period picture: accurate

  $\rightarrow$ yields upper bound

  $\rightarrow$ starkly different from exponential: equal blame

  $\rightarrow$ implication to sampling and simulation: slow convergence

# Sample mean convergenge rate: exponential vs. Pareto

Lastly: characteristic of aggregate traffic

$\longrightarrow$ multiple on/off sources

Recall:



$\longrightarrow$ with many on/off sessions, what does $X(t)$ look like?

$\longrightarrow$ it's fractal, i.e., self-similar!

Some fractal objects:

Menger sponge (picture from `www.ics.uci.edu/~eppstein`):



Fractal fern:



$\longrightarrow$    are fractal objects random?

# Internet traffic: measurement

Host A

Router

Backbone

Host B

LAN X

LAN Y

Traffic Meter

0ms   10ms   20ms   30ms   40ms   t

$\longrightarrow$ traffic time series (at 10ms granularity)

Aggregation (time):

0ms   10ms   20ms   30ms   40ms   50ms   60ms   70ms   80ms   90ms   100ms   110ms   •••

0ms                 100ms                 200ms                 300ms   •••

$\longrightarrow$   analogous to computing sample mean

$\longrightarrow$   aggregation over multiple time scales

$\longrightarrow$   what to expect?

# Internet: self-similar          Telephony: Poisson-like

**100s**

100s aggregation (alpha 1.05)

zoom 10x

100s aggregation (expo)

**10s**

10s aggregation (alpha 1.05)

10s aggregation (expo)

**1s**

1s aggregation (alpha 1.05)

1s aggregation (expo)

**100ms**

0.1s aggregation (alpha 1.05)

0.1s aggregation (expo)

**10ms**

0.01s aggregation (alpha 1.05)

0.01s aggregation (expo)

We observe:

- for Internet traffic burstiness preserved across time scales five orders of magnitude apart

  $\rightarrow$ Poisson traffic: smoothes out quickly

- if traffic were uncorrelated in time, by LLN should smooth out

  $\rightarrow$ how fast should it smooth out?

Self-similar burstiness viewpoint:

Time aggregation of $X(t)$ at level $m$ means

$$X^{(m)}(i) = \frac{1}{m} \sum_{t=m(i-1)+1}^{mi} X(t).$$

Since $X(t)$ are random variables, $X^{(m)}(i)$ in time series is analogous to computing the sample mean.

The visual phenomenon of "burstiness preservation" corresponds to

$$\mathrm{var}(X^{(m)}(i)) \approx \mathrm{var}(X(t))$$

for a range of time scales $m$.

If the $X(t)$'s were independent, then

$$\text{var}(X^{(m)}(i)) = \sigma^2 m^{-1}$$

where $\sigma^2$ is the variance of the $X(t)$'s.

$\longrightarrow$  elementary fact

Consider rewriting expression with parameter $H$ as

$$\sigma^2 m^{-2(1-H)}$$

where $1/2 \leq H < 1$.

If $H = 1/2$, then we have previous expression $\sigma^2/m$.

$\longrightarrow$  $\sigma^2$ decays at rate $m^{-1}$

If $1/2 < H < 1$, then rate of decay is slower.

$$\longrightarrow \quad m^{-\beta} \text{ where } 0 < \beta < 1$$

Thus, if $H \approx 1$, then can expect

$$\mathrm{var}(X^{(m)}(i)) \approx \mathrm{var}(X(t))$$

$\longrightarrow$   burstiness dies out very slowly w.r.t. scaling

$\longrightarrow$   empirically: $H$ is 0.8–0.9 range

$\longrightarrow$   $X(t)$ must be strongly correlated in time

$\longrightarrow$   what causes it?

The principal cause: heavy-tailed file sizes!

$\longrightarrow$  present impacts distant future

$\longrightarrow$  recall predictability discussion

$$\Pr\{Z > x + y \mid Z > y\} = \left(\frac{y}{y + x}\right)^{\alpha}$$

$\longrightarrow$  predictability also leads to long-term correlation

$\longrightarrow$  consequences: heavy-tailed queueing

$\longrightarrow$  periods of over- and under-utilization

$\longrightarrow$  bad for resource provisioning

# Internet: self-similar



# Telephony: Poisson-like