

The Origins of Network Server Latency & the Myth of Connection Scheduling

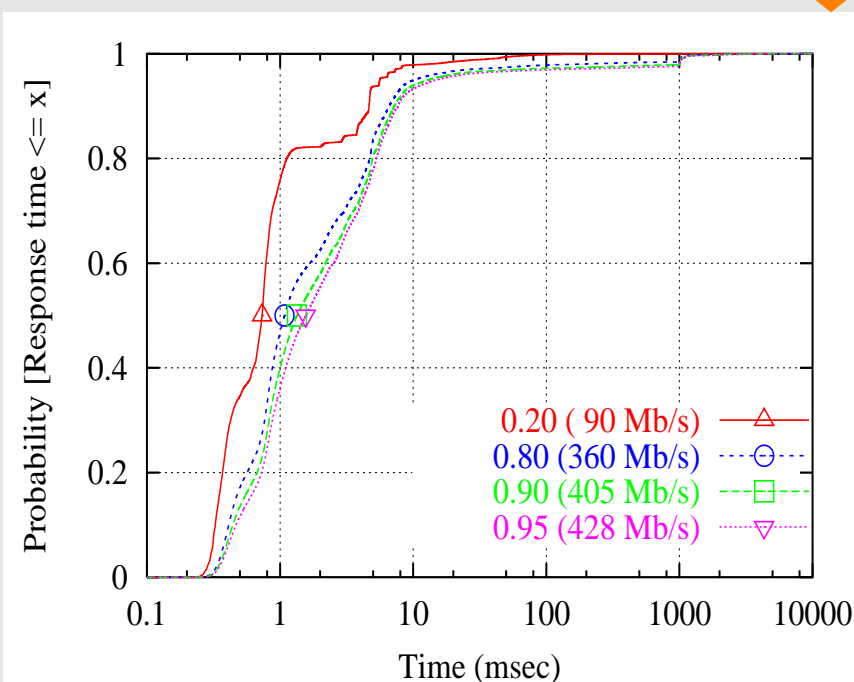
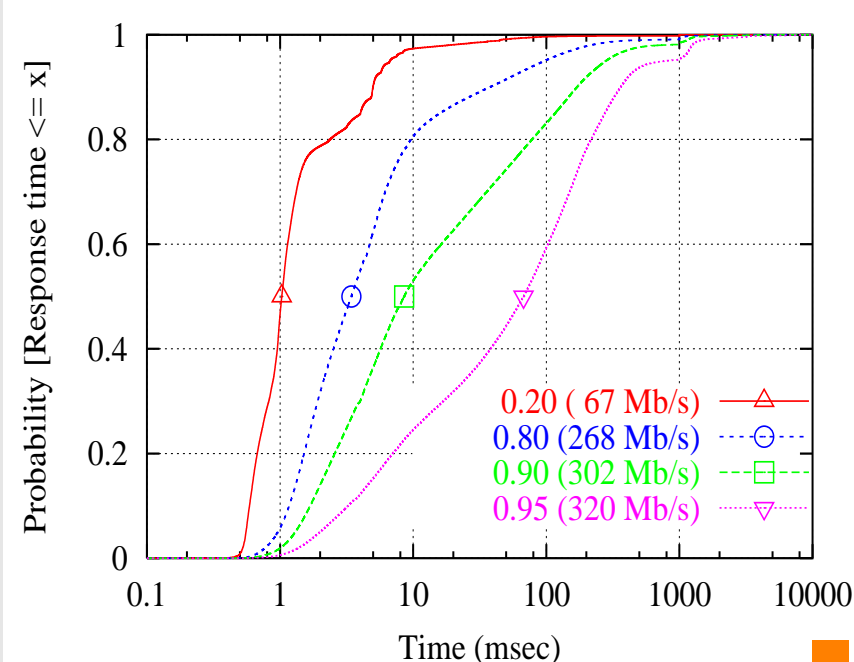


Princeton University

Yaoping Ruan
Vivek S. Pai

{yruan, vivek}
@ cs.princeton.edu

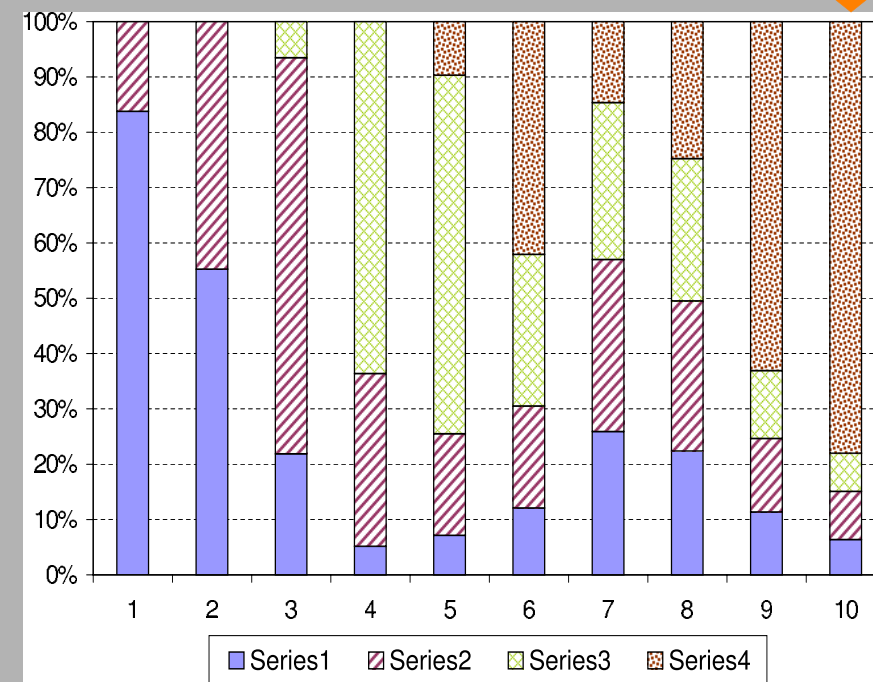
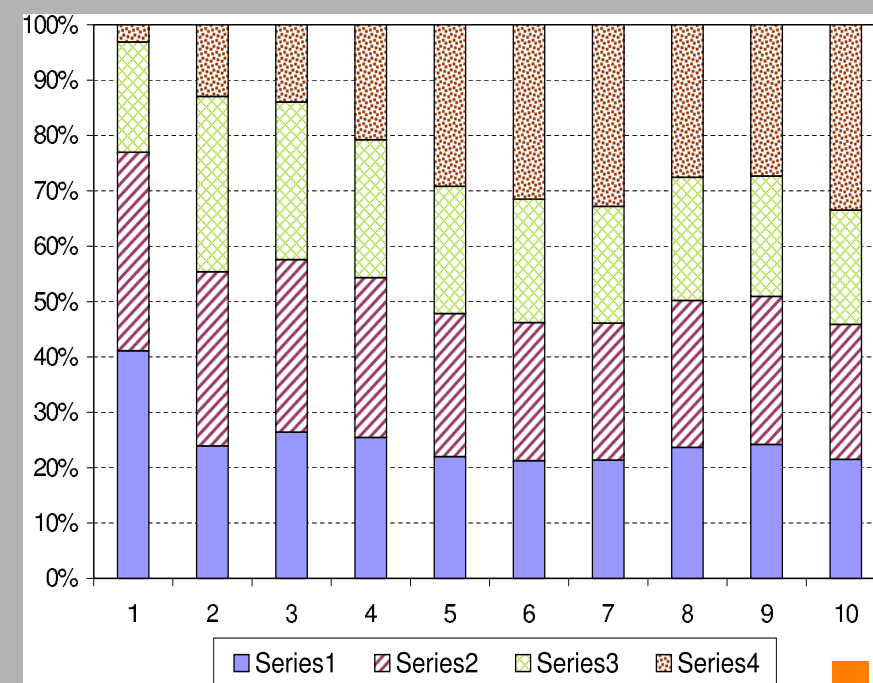
Blocking, Not Queuing, Causes Most Latency



Latency CDFs

- OS head-of-line blocking causes latency growth
- Eliminating blocking greatly reduces the base latency
- Eliminating blocking also curtails the growth with increasing load
- Over 60% of latency caused by blocking

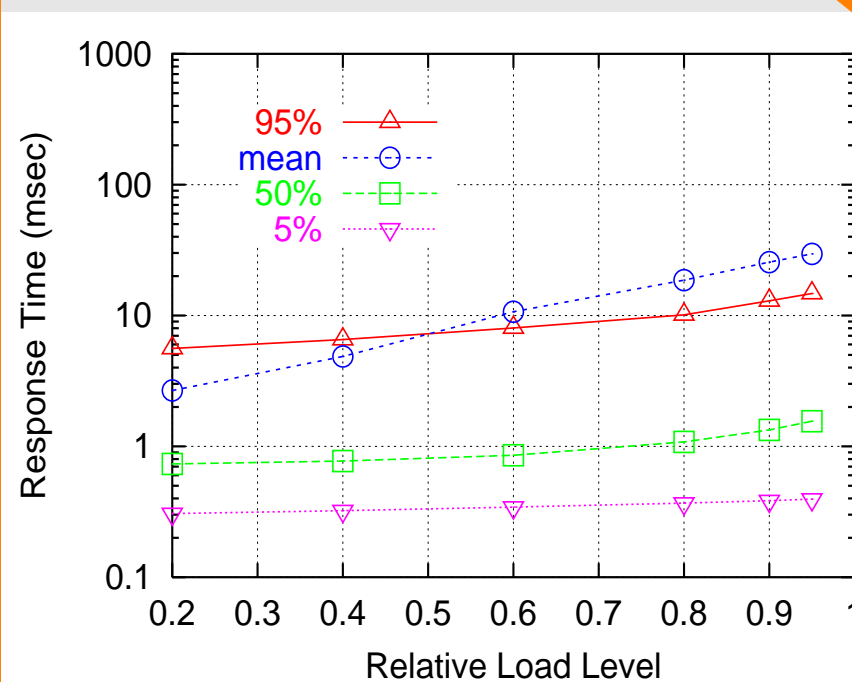
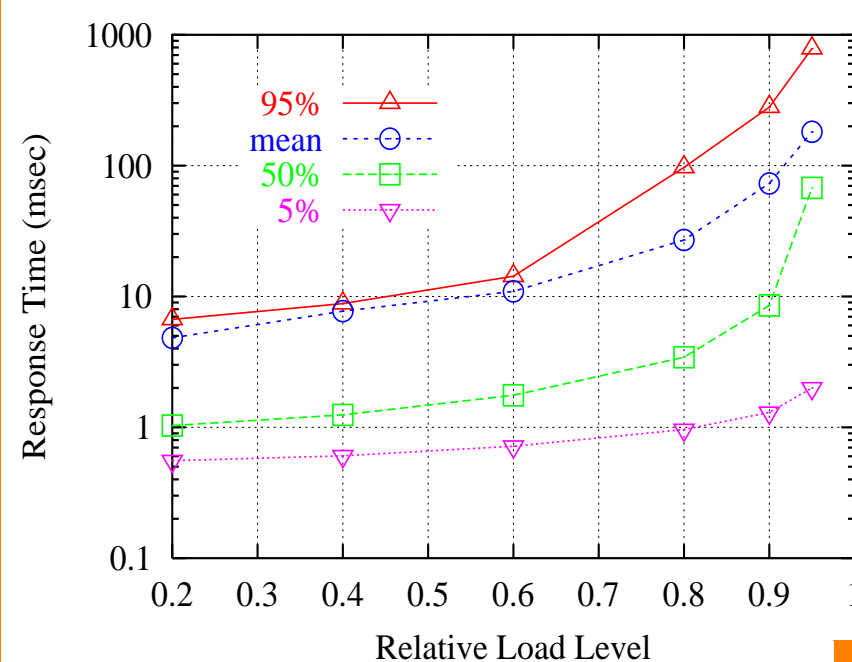
Blocking Degrades Fairness Policies



Request size distribution by latency decile

- Blocking cripples existing fairness schemes
- Unfairness evident as inversion in the service times
- Connection scheduling attempts to address this problem
- Eliminating blocking reduces service inversion

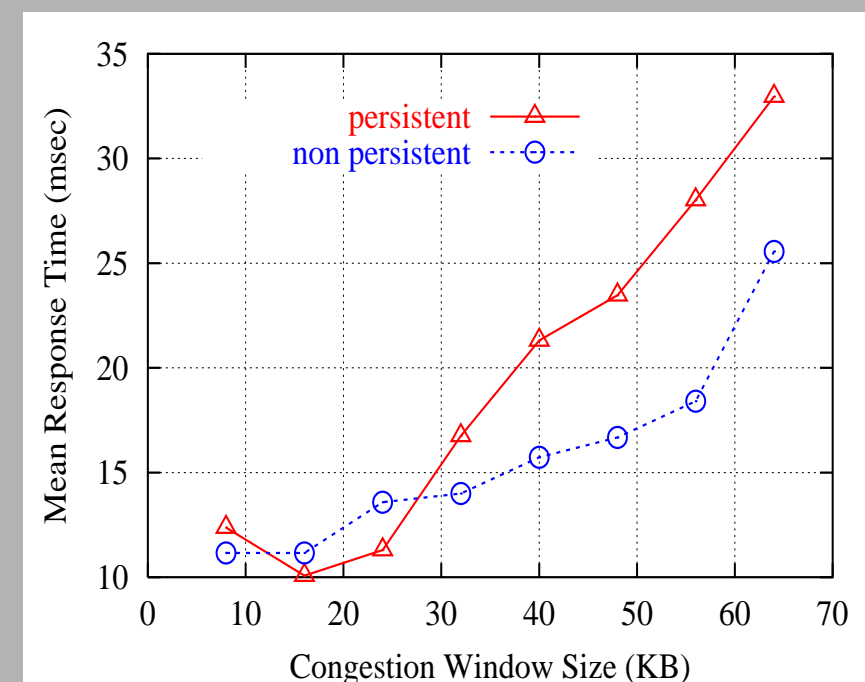
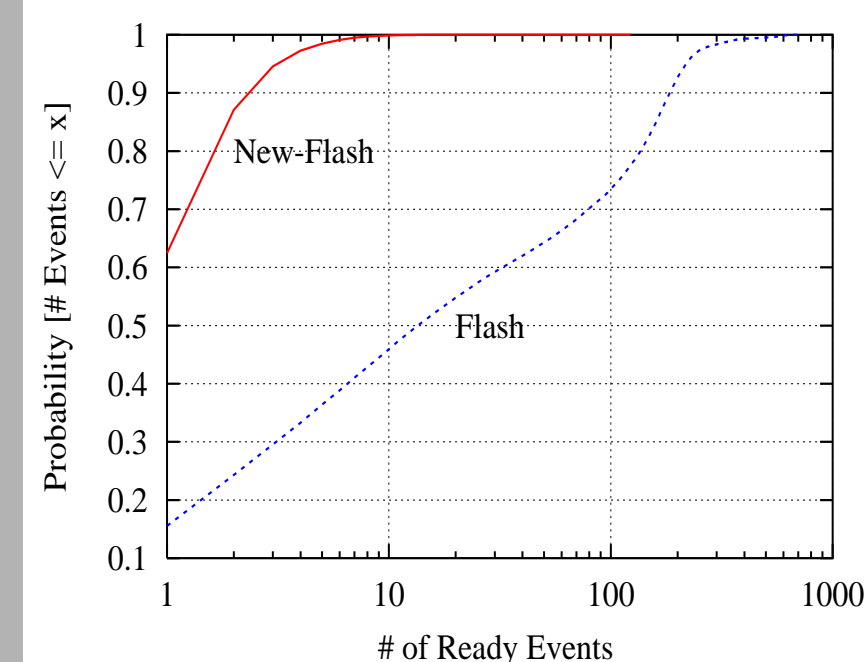
Fixing Blocking Changes Latency Qualitatively



Latency growth vs. load level

- A factor of 6 in mean and 43 in median latency improvement
- Median and 95th percentile are now virtually flat
- Mean latency growth indicates heavy tail
- Most requests served from memory, now unaffected by disk access

Connection Scheduling is Not Necessary



Queue length CDF & congestion window effects

- Blocking causes burstiness and long queues
- Eliminating blocking dramatically reduces event queue lengths
- New server shows no benefit from SRPT scheduling
- Congestion window size has impact on unfairness in network