
Fundamental Bounds on the Accuracy of Network Performance Measurements.

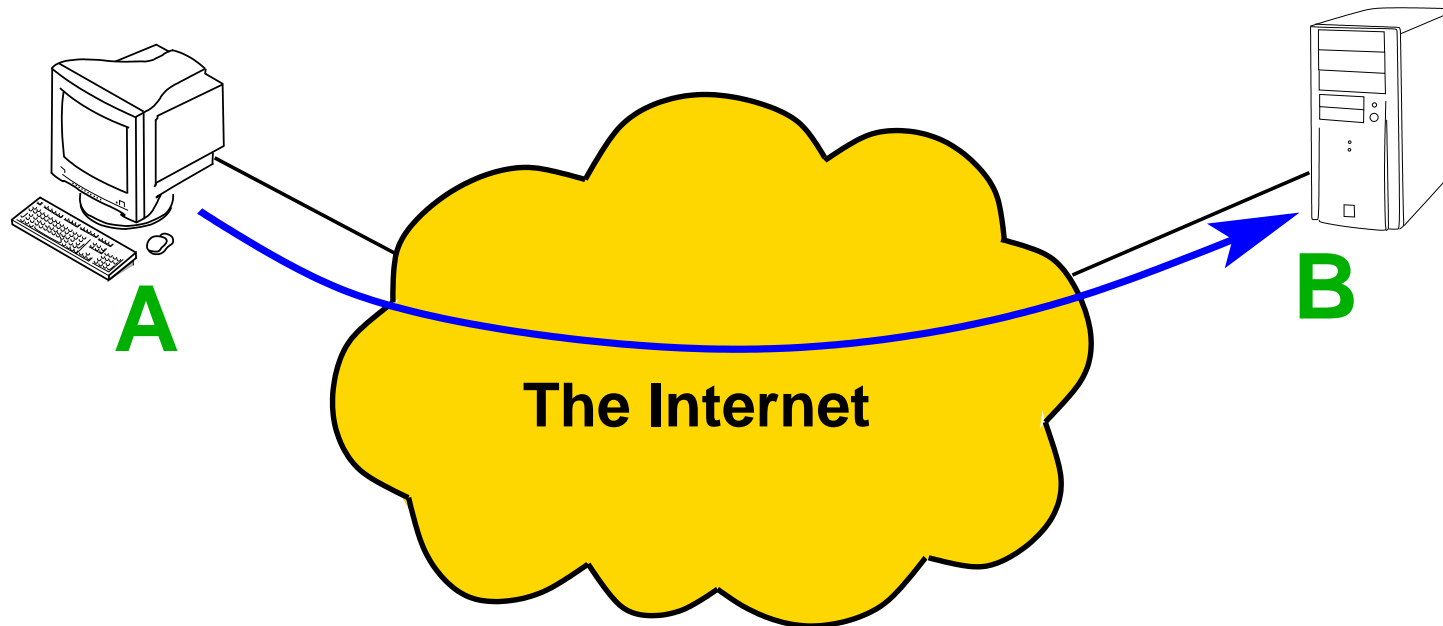
Matthew Roughan

[<matthew.roughan@adelaide.edu.au>](mailto:matthew.roughan@adelaide.edu.au)

Discipline of Applied Mathematics
School of Mathematical Sciences
University of Adelaide

The problem

- Active performance measurements
- Send probe packets from $A \rightarrow B$ across the network
- e.g. measure the delays experienced by packets



- How many probe packets should we send?
 - really we need to be a little more specific

Motivation

Another way to state the problem is how accurate will a set of N measurements be?

- What do I mean by accurate?
 - not equipment accuracy!
 - assume perfect infrastructure
 - we mean statistical accuracy
- Can I achieve arbitrary accuracy?
 - naively you might say yes: take $N \rightarrow \infty$
- In reality there are **fundamental bounds**

Related problems

Applications

- **network quality control**
- anomaly detection
- streaming playout buffer size estimation
- load balancing & TE
- TCP RTO est.
- Vegas congestion meas.
- tomography (topology)
- location mapping

Measurements

- **packet delay**
- packet loss rate
- packet jitter
- packet reordering
- throughput

Statistical Accuracy

What do we mean by accuracy

- often individual measurements are inaccurate.
- implicit assumption of stationary ergodic process
⇒ a time average converges to an ensemble average
- measurements over time can be averaged to give a better estimate of the mean delay
- variance can be directly quantified by the **Central Limit Theorem**
- assume Gaussian limit, quantify accuracy by confidence bounds for estimates.

Accuracy of **estimates** not individual measurements

Central Limit Theorem

Set of independent, identically distributed RVs X_i
with sample mean

$$\hat{X} = \frac{1}{N} \sum_{i=1}^N X_i,$$

then $E[\hat{X}] = E[X_0]$, and

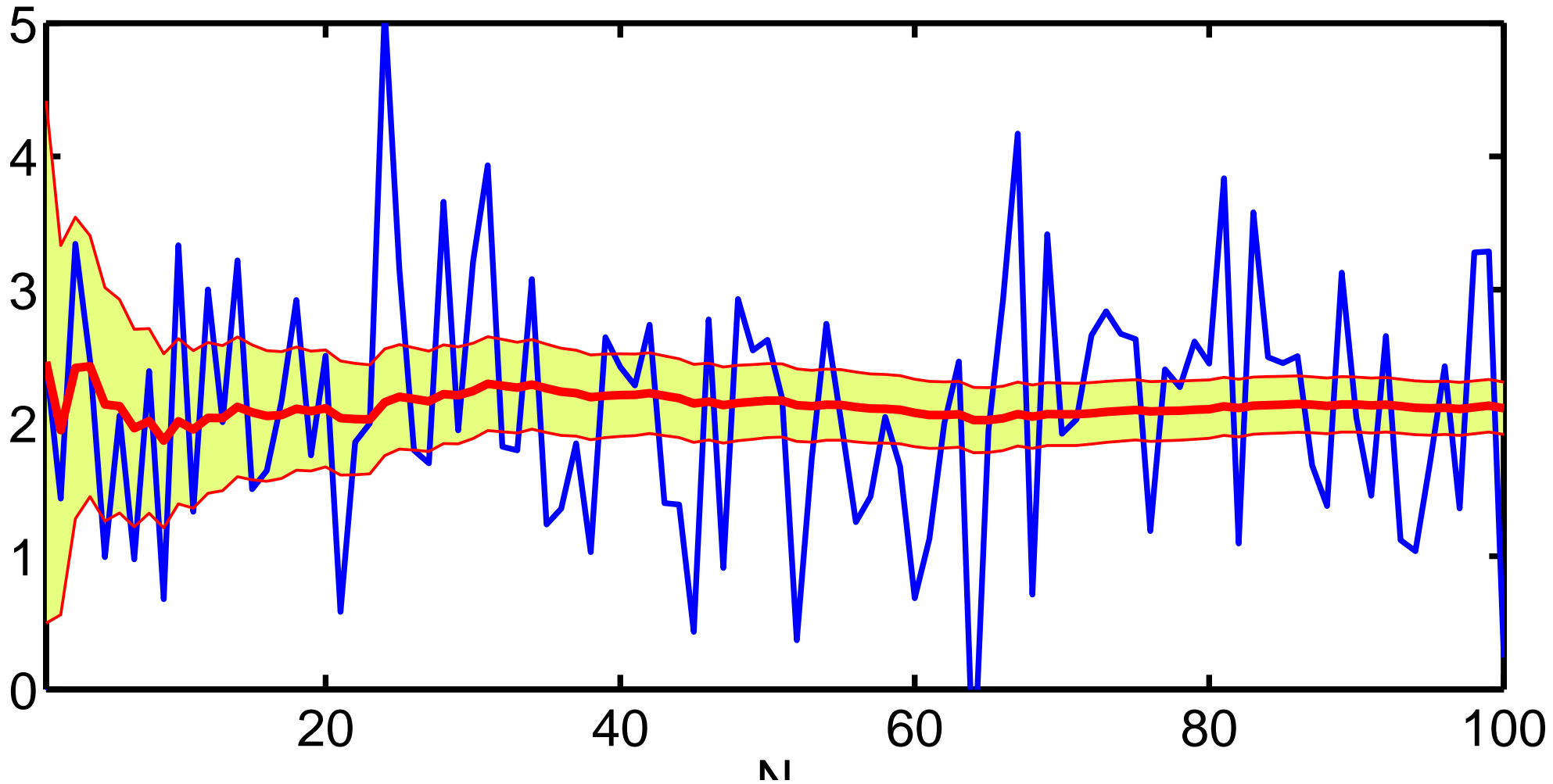
$$\sqrt{N} (\hat{X} - E[X_0]) \rightarrow N(0, \sigma^2)$$

in distribution as $N \rightarrow \infty$, where $\sigma^2 = \text{Var}[X_0]$

So the 95th% CIs of estimate \hat{X} are $\pm 1.96\sigma/\sqrt{N}$

Example

White noise



Time is short

- Stationarity is at best an approximation
 - approx. on short (e.g. < 1 min.) intervals
 - not true for long (e.g. > 24 hour) intervals
- We need to detect problems quickly
 - problems may be transient
 - diagnose problems within minutes to fix
- Some applications aren't around long enough
 - TCP RTT measurements
 - Streaming playout buffer needs to be determined at start of stream.

Constrained time interval

- constrained measurement interval
- perfect measurements (no artifacts)
- passive measurements

How accurate can we be?

- To increase N , measure more frequently.
- Optimal is continuous measurements, $N \rightarrow \infty$.
- Does estimate variance go to zero?

Need a continuous-time version of the CLT

Central Limit Theorem: cont. time

Continuous time process $X(t)$ where the sample mean

$$\hat{X} = \frac{1}{T} \int_0^T X(u) du$$

converges to the true mean $\hat{X} \rightarrow E[X]$, and

$$\sqrt{T} (\hat{X} - E[X]) \rightarrow N(0, s^2)$$

in distribution as $T \rightarrow \infty$, where

$$s^2 = 2\sigma^2 \int_0^\infty r(u) du$$

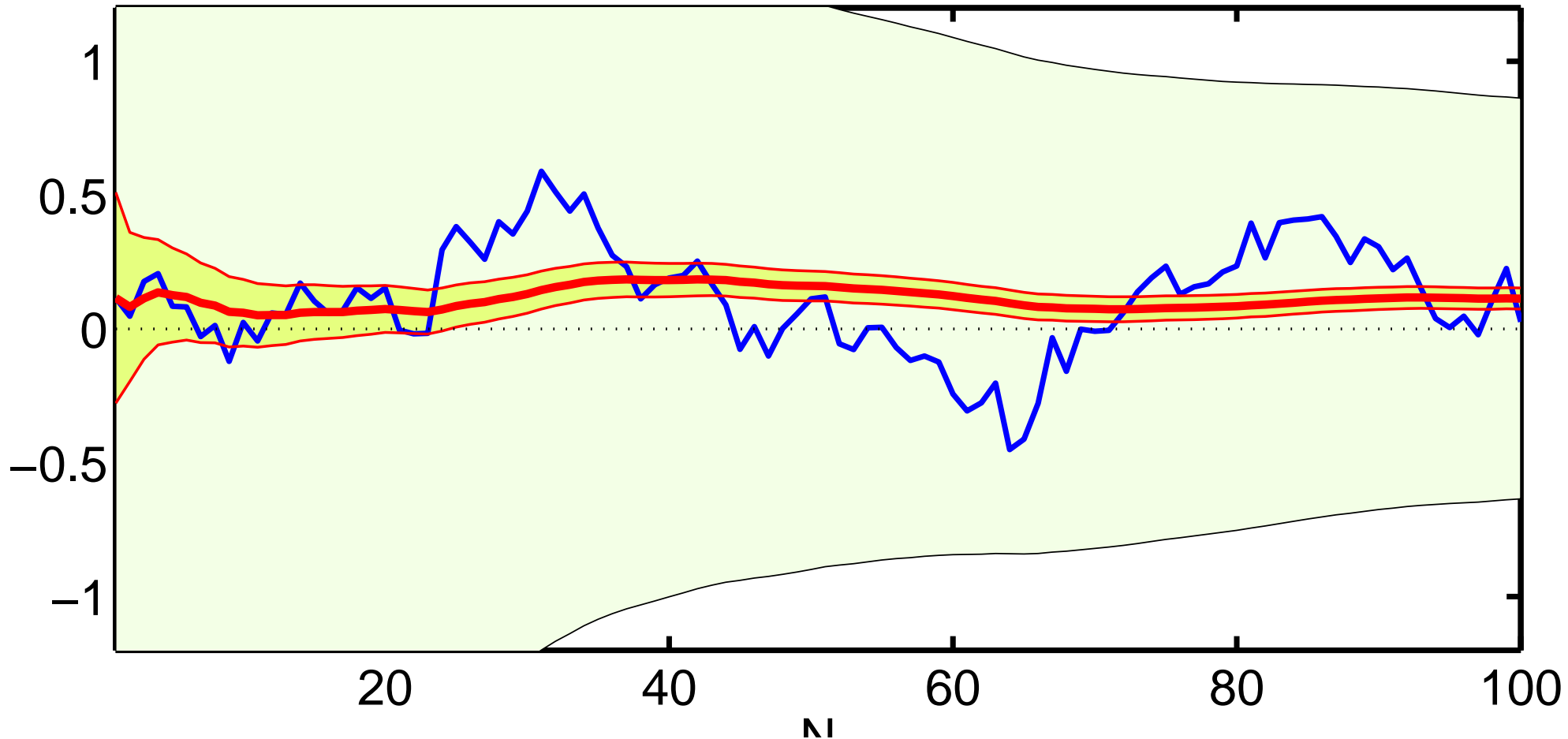
where $\sigma^2 = \text{Var}[X]$, and $r(s)$ is the autocorrelation of $X(t)$.

What does it mean

- closer samples are more correlated
- less information gained per sample
- There is a **limit** as $N \rightarrow \infty$
- Captured in the **asymptotic variance s^2**
- Asymptotic results, but similar impact on short term measurements.
- Accuracy determined by T , σ and $r(s)$.

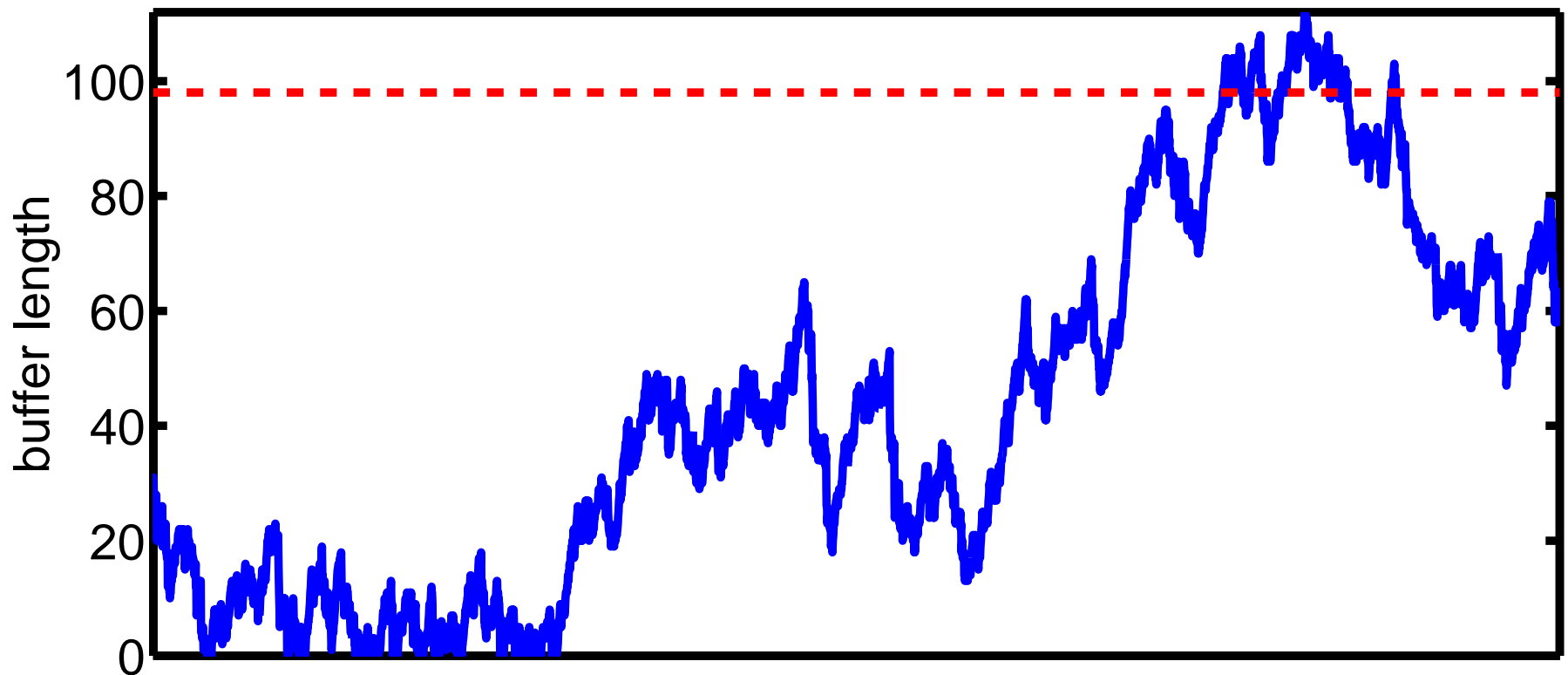
Impact of correlated measurements

EWMA: AR(1) process $Z_t = \alpha Z_{t-1} + (1 - \alpha)X_t$



How to apply here

- Perfect measurements (measurement error zero).
- variability comes from queueing delays
- are queueing delays correlated? **YES!**



M/M/1 queue

- Poisson packet arrivals (rate λ)
- Exponential service times (mean $1/\mu$)
- Average queue length

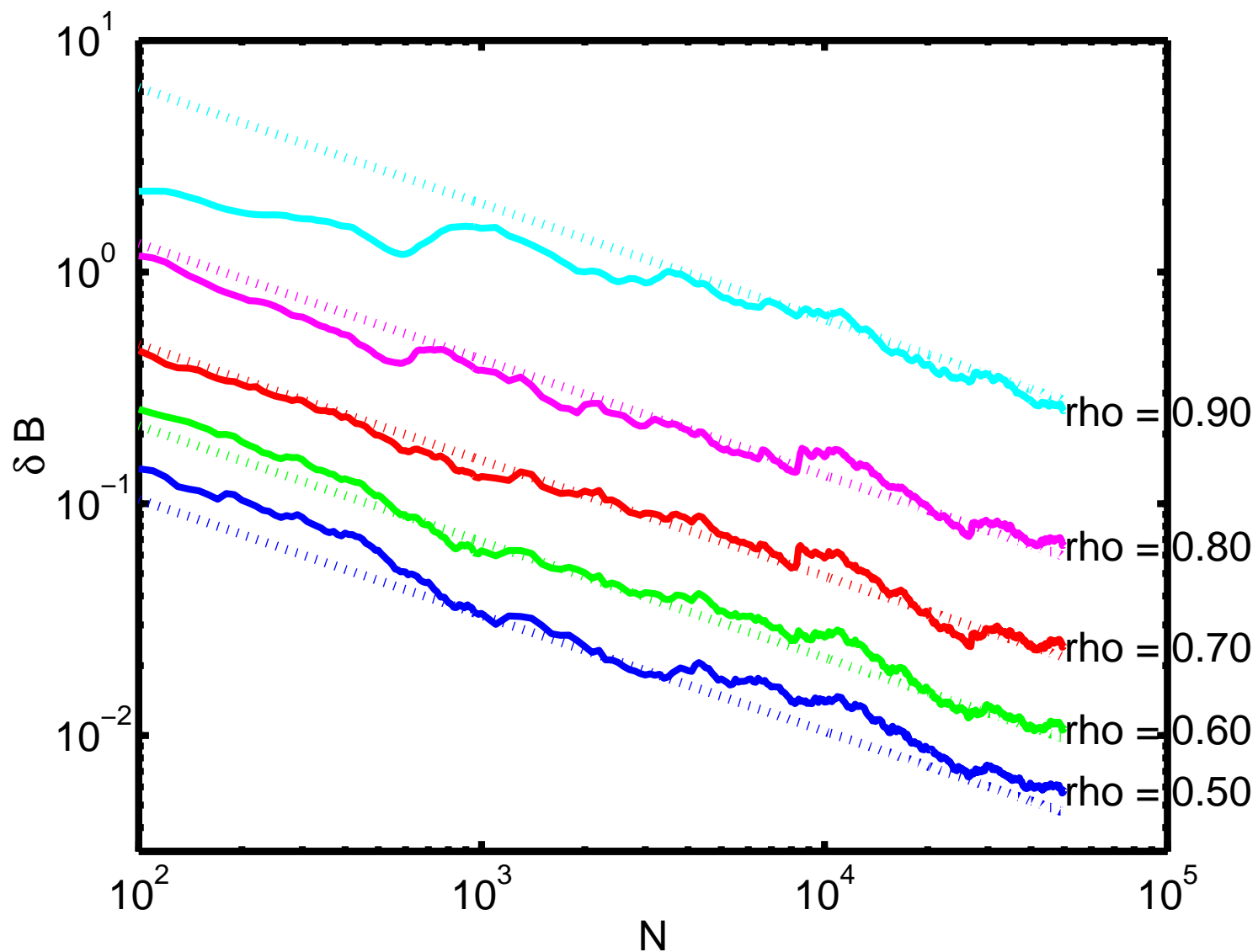
$$E[Q] = \frac{\rho^2}{1-\rho}$$

- asymptotic variance for M/M/1 (Whitt, 1989)

$$s^2 \simeq \frac{4\rho^2}{(1-\rho)^4}$$

- Correlations from excursions away from empty system
 - heavy-load \Rightarrow long busy periods
 - heavy-load \Rightarrow more correlation
 - s^2 is heavily load dependent

Results M/M/1



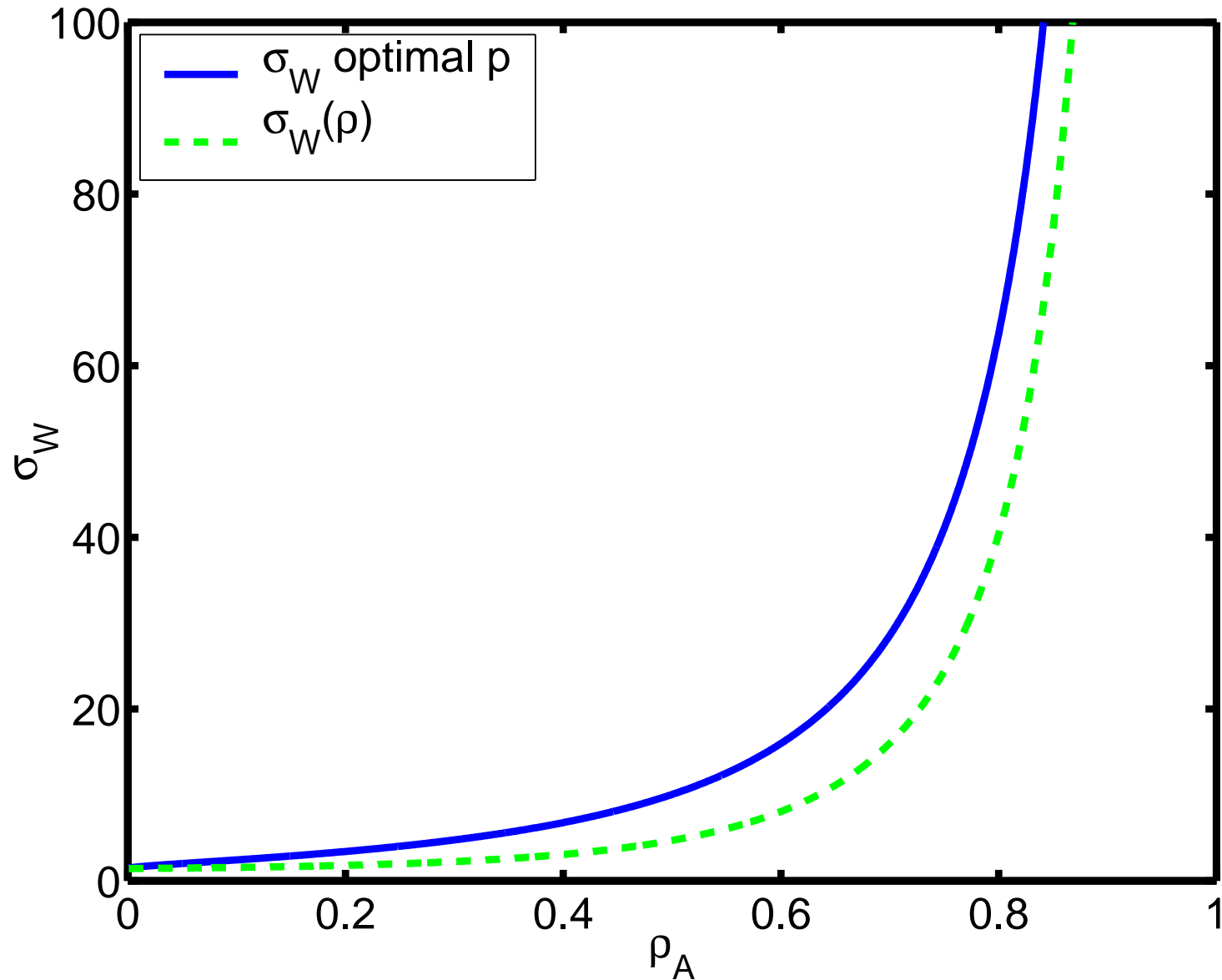
Implications

1. there is a **fundamental bound** on the accuracy with which we can estimate queueing delays,
 - it is dependent on the
 - length of the measurements interval
 - load on the queue

Active probing

- Everything until now has been passive
- Heisenberg effect
 - measurements impact the system
 - in turn this impacts the measurements.
- More rapid probing for more accuracy
 - increases queue load
 - increases correlations
 - reduces accuracy
 - can't be unravelled
- once again we can quantify
- we can compute optimal probe rate

Optimal Probing

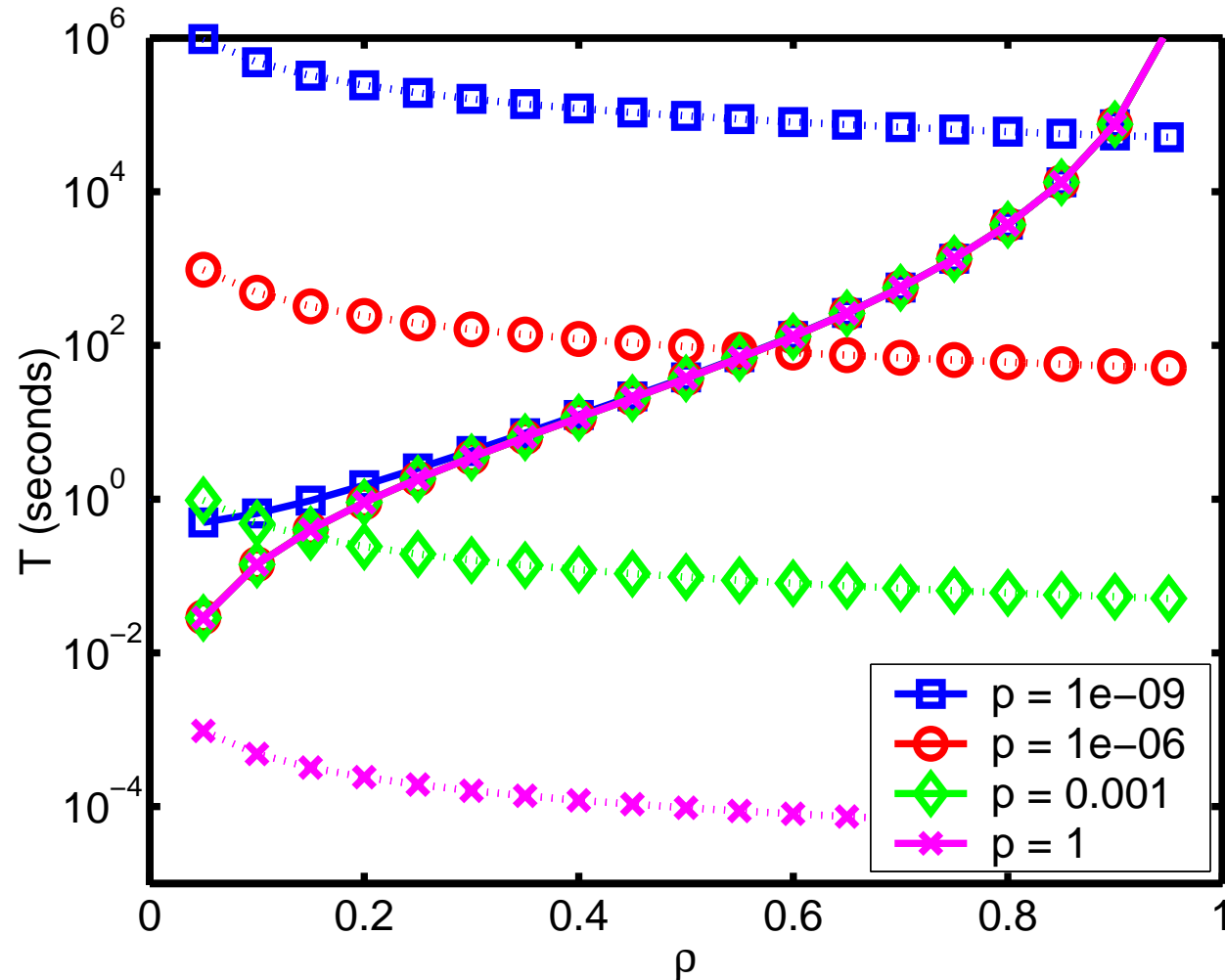


Implications

1. there is a **fundamental bound** on the accuracy with which we can estimate queueing delays,
 - it is dependent on the
 - length of the measurements interval
 - load on the queue
2. active probing increases the load
 - increases correlations
 - reduces the estimator accuracy.
3. **you can't do better by probing more quickly**
 - in fact you do worse
 - forms a bound like Heisenberg's uncertainty principle

The scale of the problem is big

- passive sampling
- M/D/1 queue
- OC48 (2.48 Gbps)
- 1500 byte packets
- p is proportion of arriving packets sampled
- ρ is normalized load
- desired accuracy $\pm 1\text{ms}$



Implications

- Faster measurements don't help much
 - Active probes should be fairly low rate
 - Passive delay measurement can sample
- TCP RTT measurements?
 - BSD only tried to get ± 500 ms
 - TCP Reno encourages large buffers
 - bad for Vegas & TCP Fast, in competition?
- load sensitivity is very bad
 - adaptive routing
 - will see oscillation for certain parameters
- problems for detecting network problems
 - can't do it quickly

Mitigation

- it's all OK for lightly loaded network
 - current networks
 - hence success for many experiments
 - maybe we should keep them lightly loaded
- ECN might be good
 - limit queue excursions
 - might just force correlations to edge
- Look at less correlated data
 - differences, not averages
 - e.g. look at queue growth
- Look at traffic, not queues
 - measure arrival rate, not queue

Conclusion

There are fundamental bounds that can't be broached

- need to understand for Internet measurement
- also need to understand for other Internet systems

Unanswered

- how important are local measurements vs global
- maybe congestion control only needs transient info?
- what do applications really need to know?
- what does this look like with real data?

[<matthew.roughan@adelaide.edu.au>](mailto:matthew.roughan@adelaide.edu.au)

Extra Slides

Discrete samples

- Correlations are not only a continuous time problem
- Discrete (uniform) samples (interval δt)

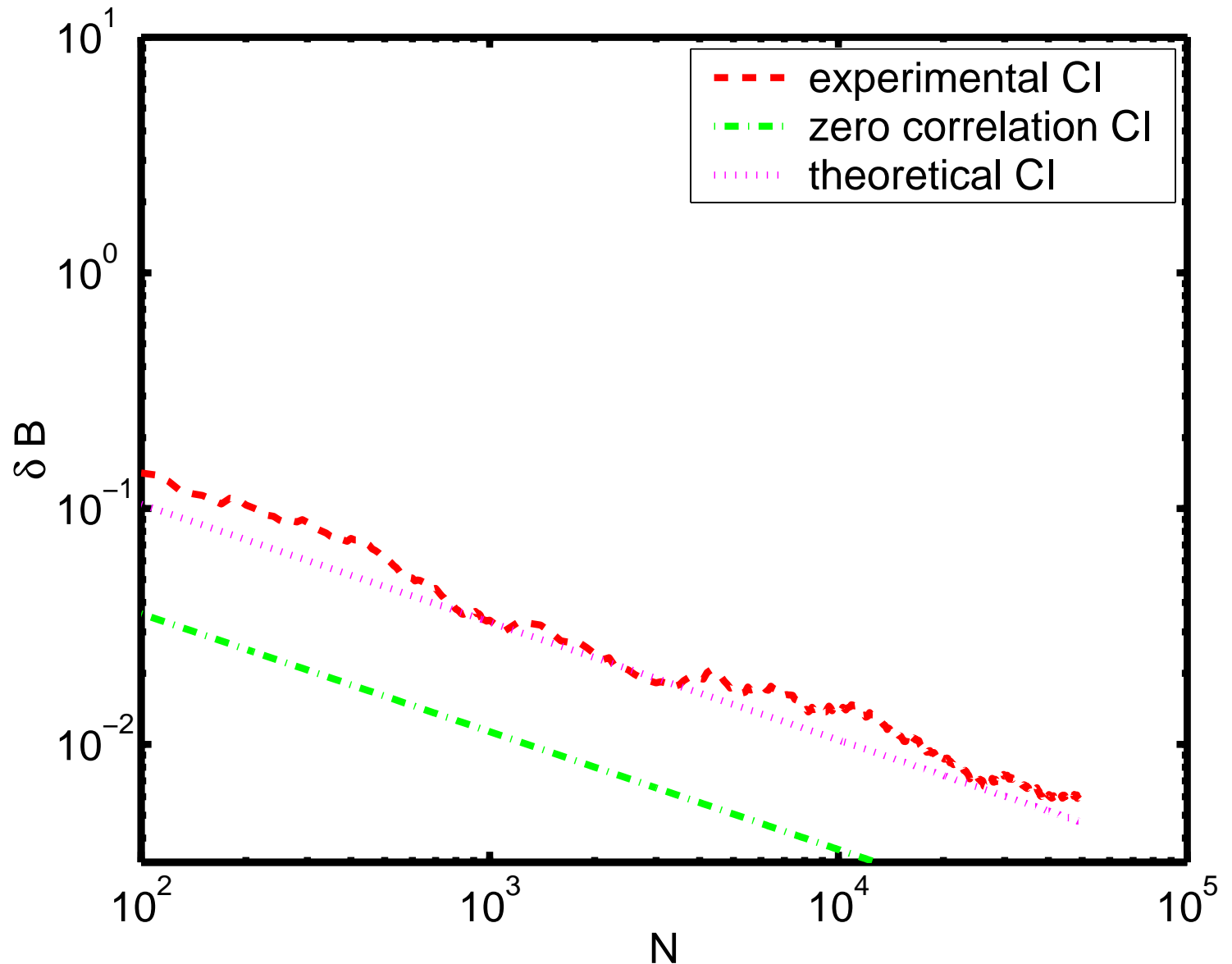
$$s^2 = \sigma^2 \left[1 + \sum_{i=1}^{\infty} r(k \delta t) \right]$$

- Poisson samples (rate λ)

$$s^2 = \sigma^2 \left[\frac{1}{\lambda} + 2 \int_0^{\infty} r(u) du \right]$$

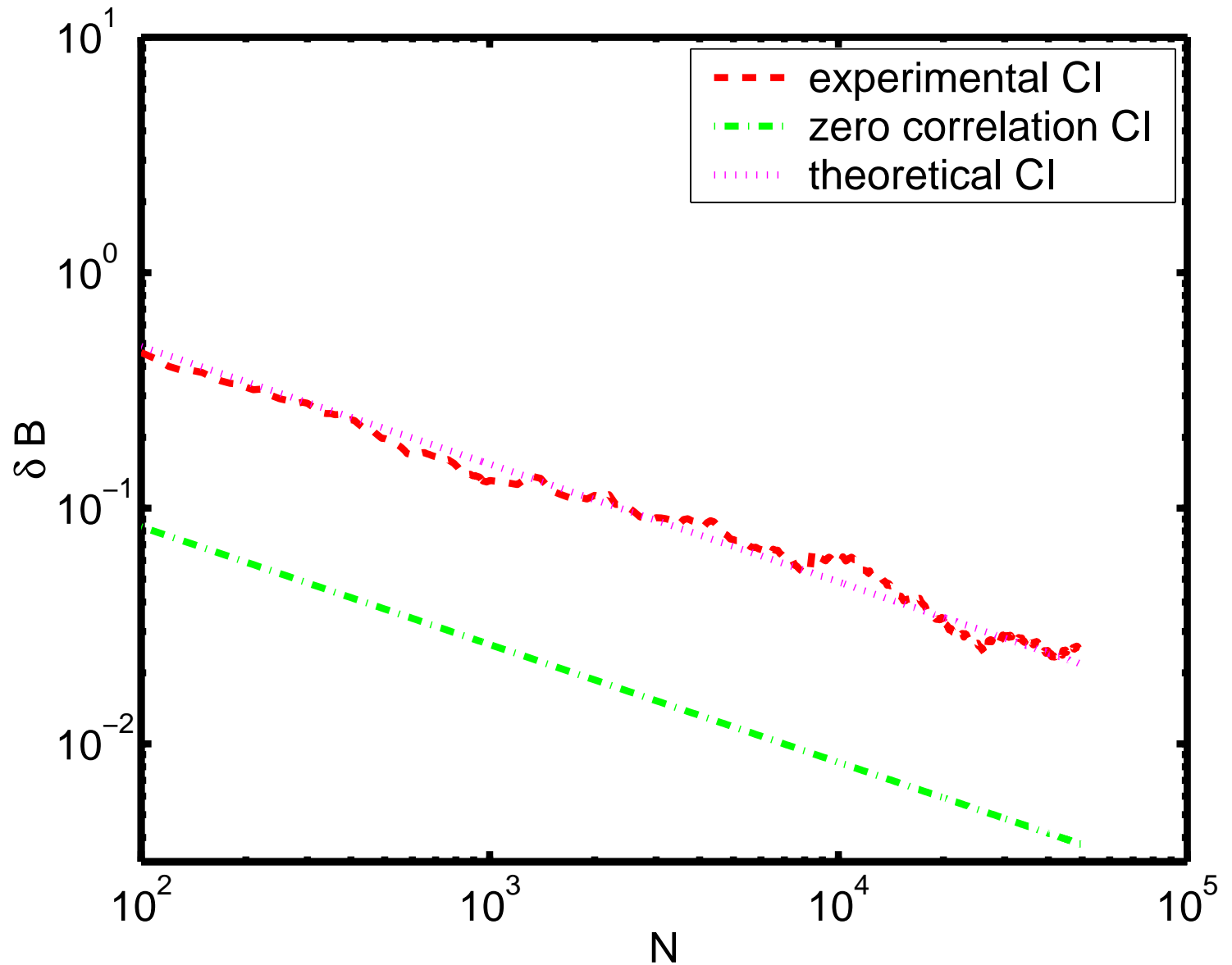
Results M/M/1

$\rho = 0.5$



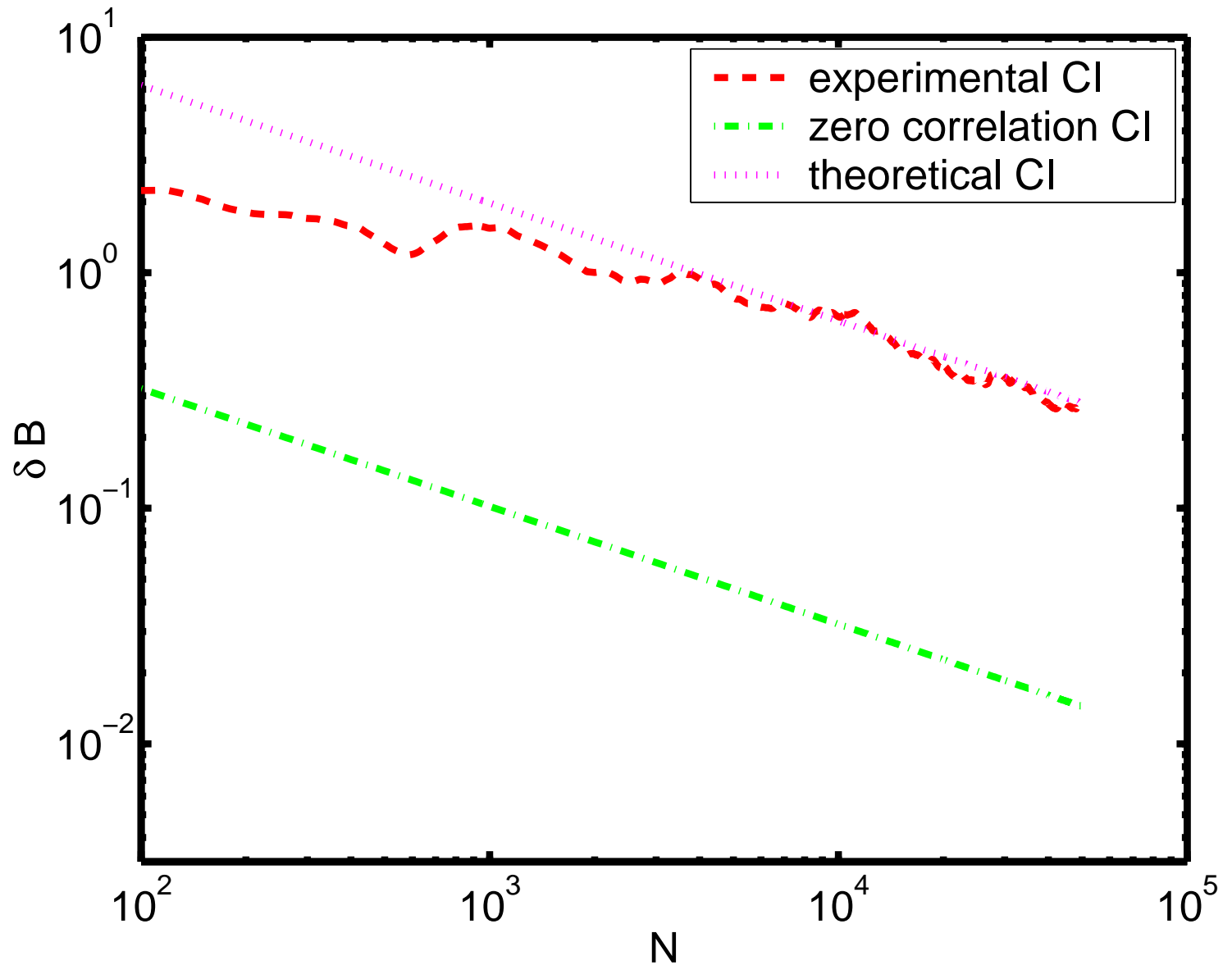
Results M/M/1

$\rho = 0.7$



Results M/M/1

$\rho = 0.9$



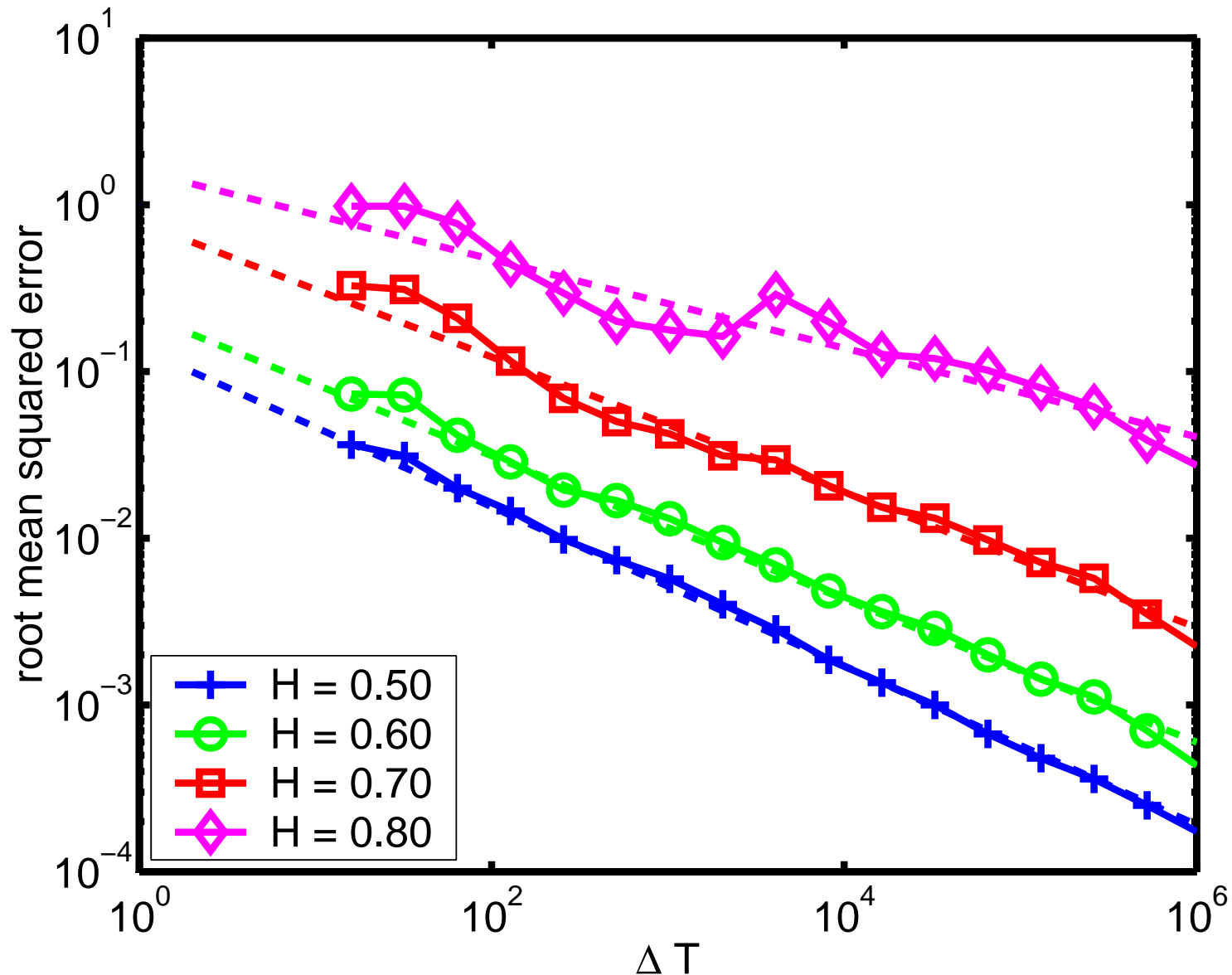
Generalizations

- M/G/1 queue (Whitt 1989)

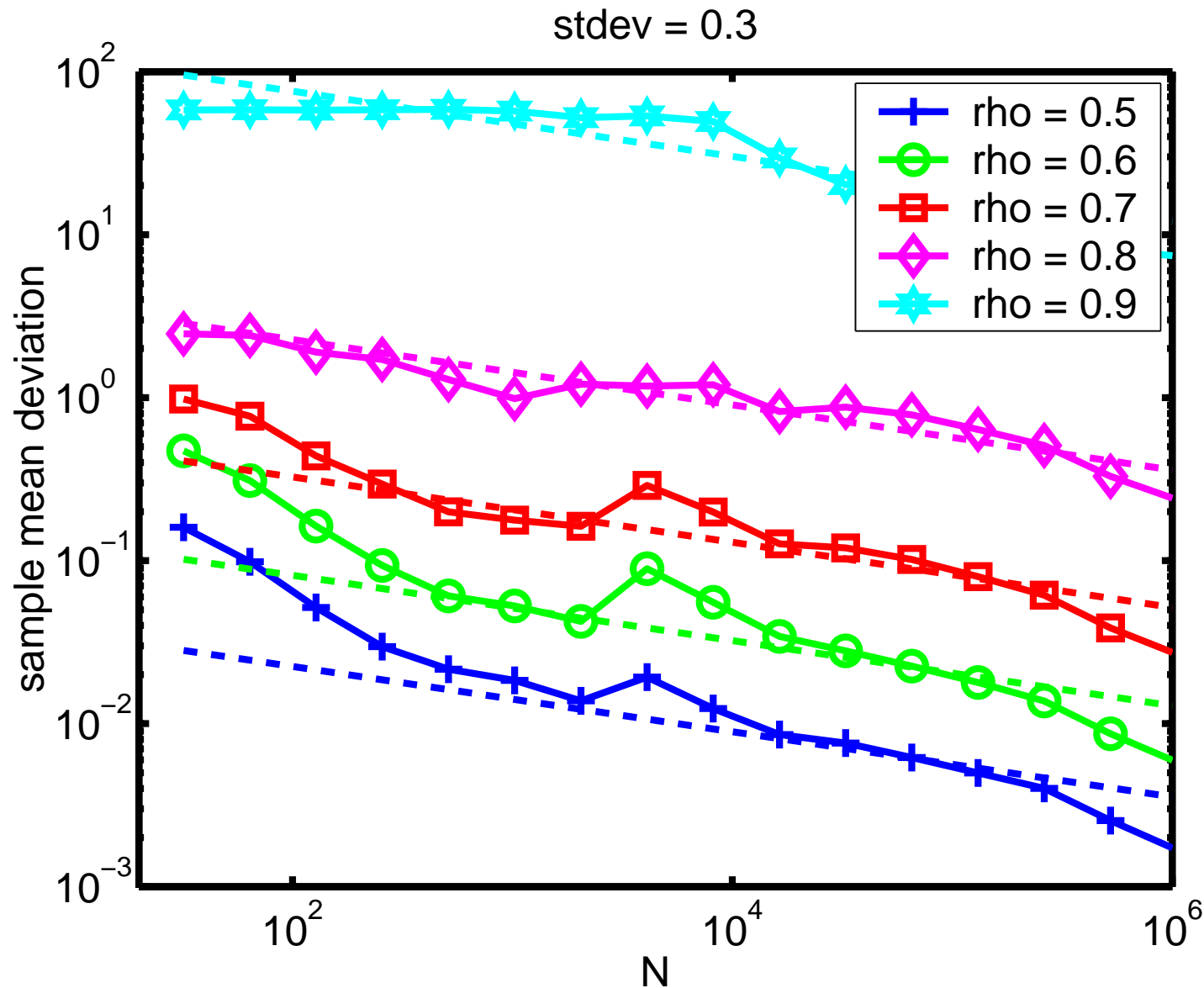
$$s^2 \simeq \frac{\rho[1 - (1 - \rho)c_s^2](1 + c_s^2)^3}{2(1 - \rho)^4}$$

- Networks: worst bottleneck
- RBM approximation (many queues)
- LRD traffic input to queues
 - generalized CLT
 - no known auto-correlations
(asymptotic results only)
 - let's use simulation

Simulation for LRD queue



Simulation for LRD queue



Optimal Probing

