

# Measuring and Optimizing Tail Latency

*Kathryn S McKinley, Google*

*CRA-W Undergraduate Town Hall  
April 5<sup>th</sup>, 2018*



**CRA-W**

Computing Research Association  
Women

# Speaker & Moderator



*Kathryn S McKinley*

Dr. Kathryn S. McKinley is a Senior Research Scientist at Google and previously was a Researcher at Microsoft and an Endowed Professorship at The University of Texas at Austin. Her research spans programming languages, compilers, runtime systems, architecture, performance, and energy. She and her collaborators have produced several widely used tools: the DaCapo Java Benchmarks (30,000+ downloads), the TRIPS Compiler, Hoard memory manager, MMTk memory management toolkit, and the Immix garbage collector.

She served as program chair for ASPLOS, PACT, PLDI, ISMM, and CGO. She is currently a CRA and CRA-W Board member. Dr. McKinley was honored to testify to the House Science Committee (Feb. 14, 2013). She is an IEEE and ACM Fellow. She has graduated 22 PhD students.



*Lori Pollock*

Dr. Lori Pollock is a Professor in Computer and Information Sciences at University of Delaware. Her current research focuses on program analysis for building better software maintenance tools, software testing, energy-efficient software and computer science education. Dr. Pollock is an ACM Distinguished Scientist and was awarded the University of Delaware's Excellence in Teaching Award and the E.A. Trabant Award for Women's Equity.



**CRA-W**

Computing Research Association  
Women

# Measuring and Optimizing Tail Latency

**Kathryn S McKinley, Google**

Xi Yang, Stephen M Blackburn,

Md Haque, Sameh Elnikety, Yuxiong He, Ricardo Bianchini

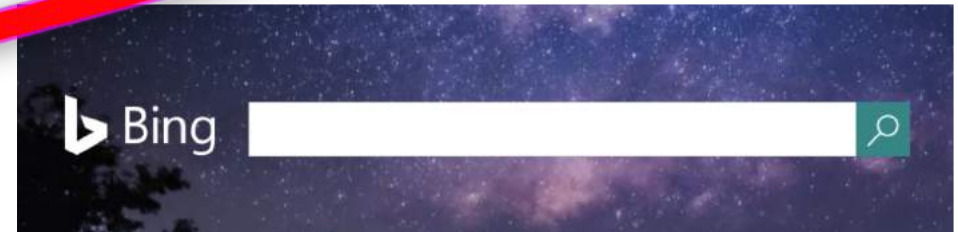




# Tail Latency Matters



400 millisecond delay decreased searches/user by 0.59%. [Jack Brutlag, Google]



Two second slowdown reduced revenue/user by 4.3%. [Eric Schurman, Bing]



# Datacenter economics quick facts\*

~ \$500,000 Cost of small datacenter

~3,000,000 US datacenters in 2016

~ \$1.5 trillion US Capital investment to date

~ \$3,000,000,000 KW dollars / year

~ \$30,000,000 Savings from 1% less work

Lots more by not building a datacenter

---

\*Shehabi et al., United States Data Center Energy Usage Report, Lawrence Berkeley, 2016.



**Tail Latency**

**TOP PRIORITY**

**Efficiency**



**Google Cloud**



**Tail Latency**

**BOTH ?!**

**Efficiency**

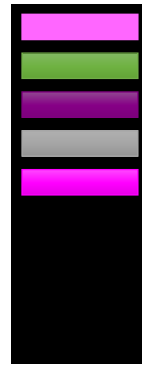


**Google Cloud**

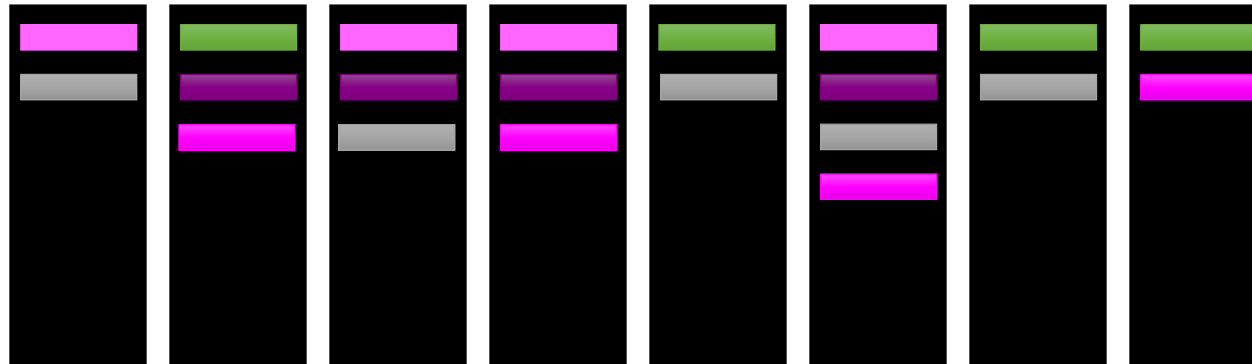


client

# Server architecture

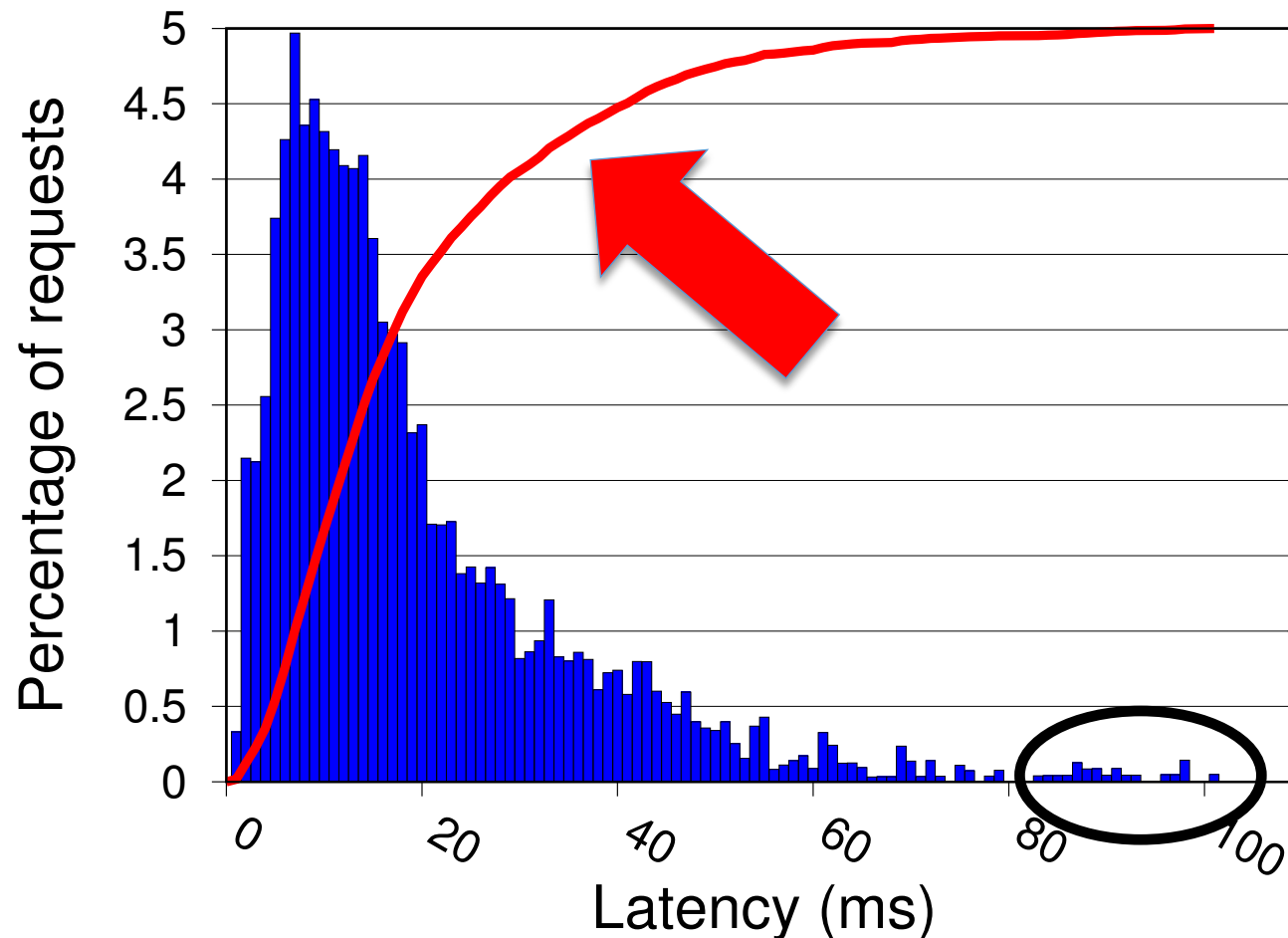


aggregator



workers

# Characteristics of interactive services



100



80

Bursty, diurnal

60

**CDF** changes slowly

40

Slowest server dictates tail

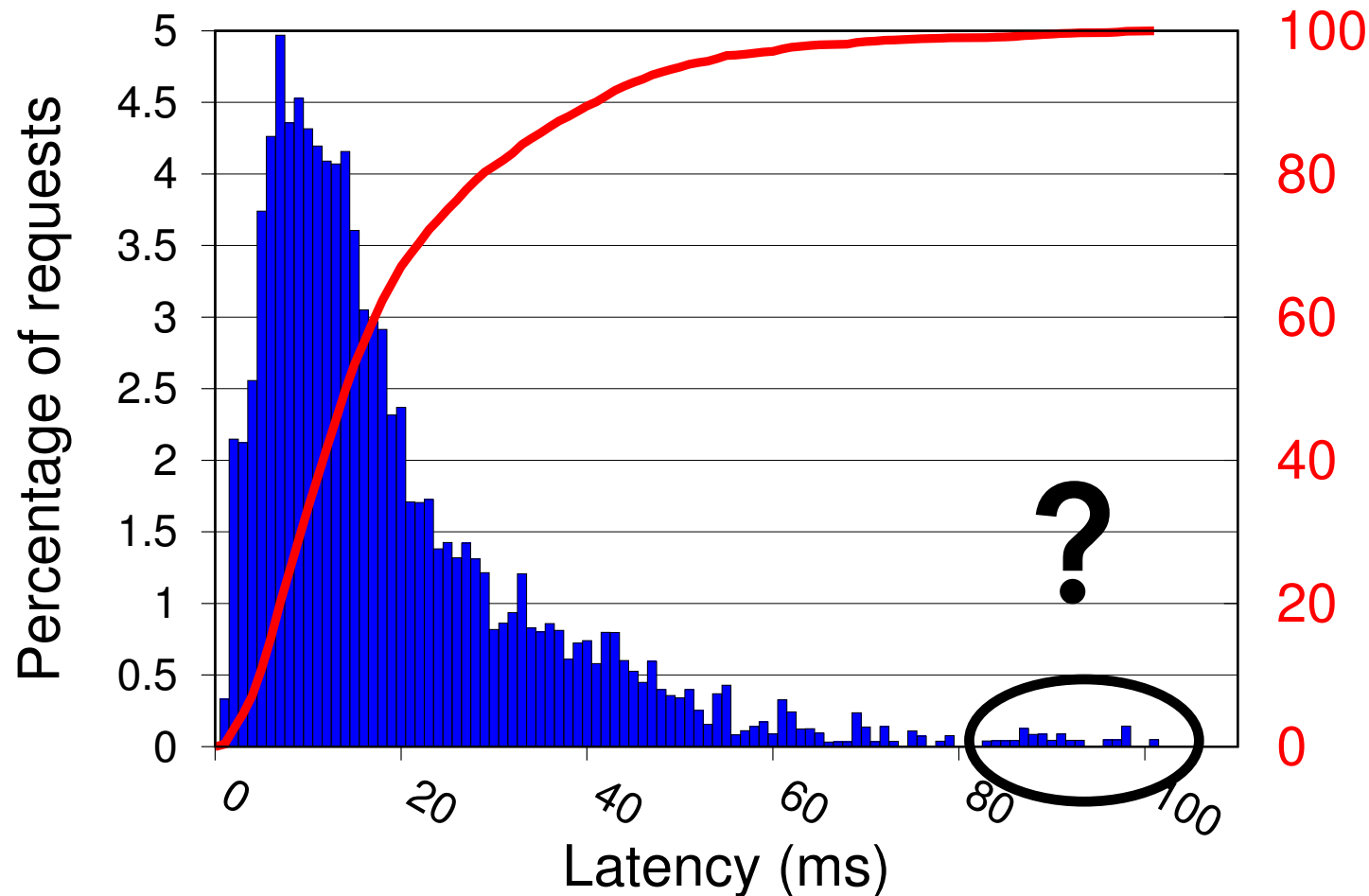
20

Orders of magnitude diff

0

average & tail - 99th %tile

# What is in the tail?



# Cycle-level on-line profiling tool

[ISCA'15 (Top Picks HM), ATC'16]



**Insight** Hardware & software generate signals without instrumentation



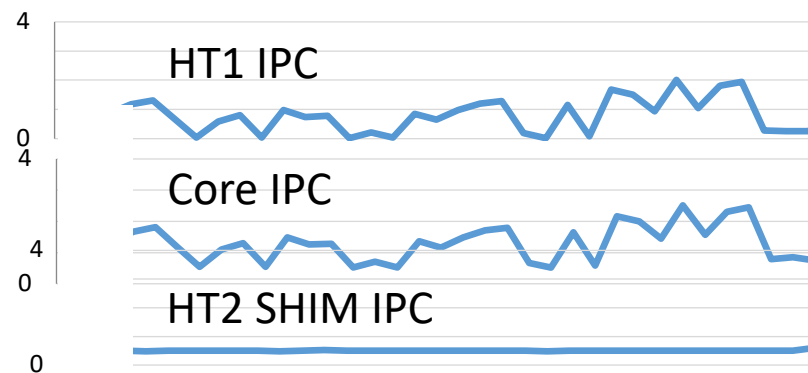
performance  
counters

memory  
locations

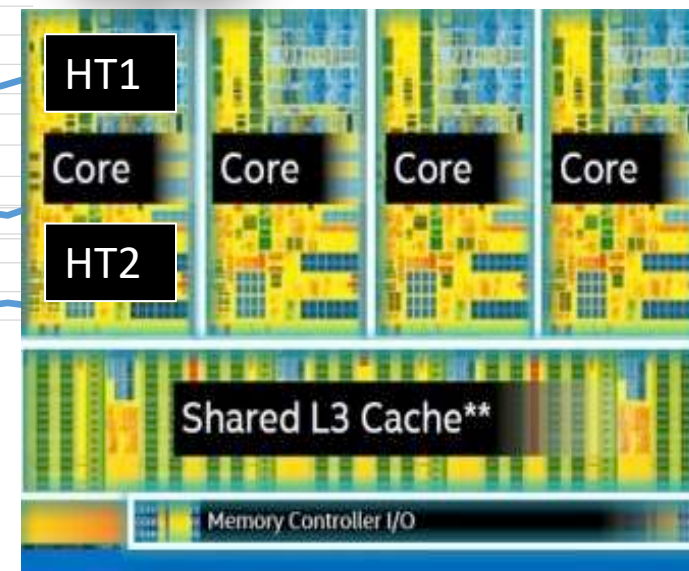
counters



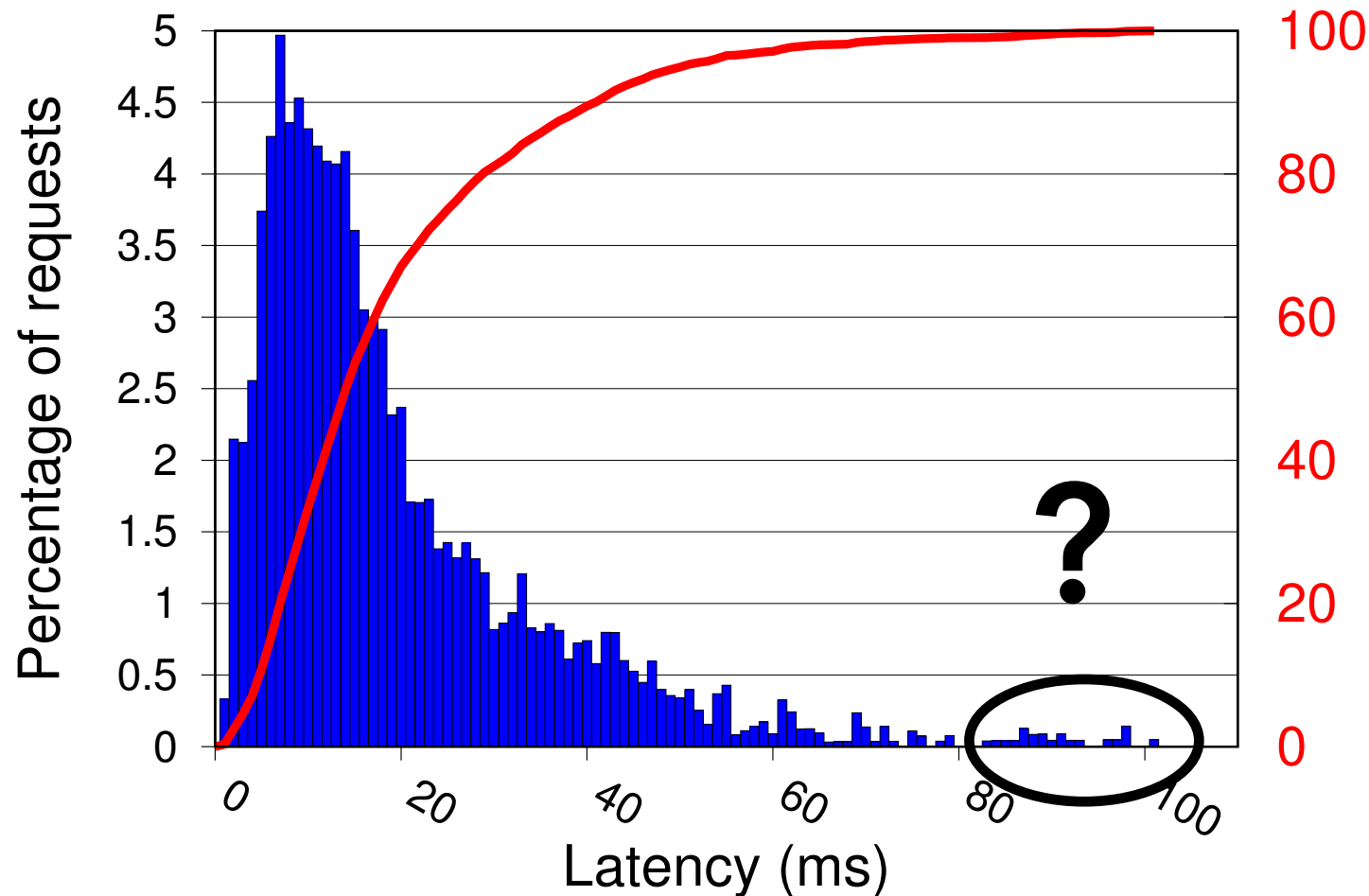
tags



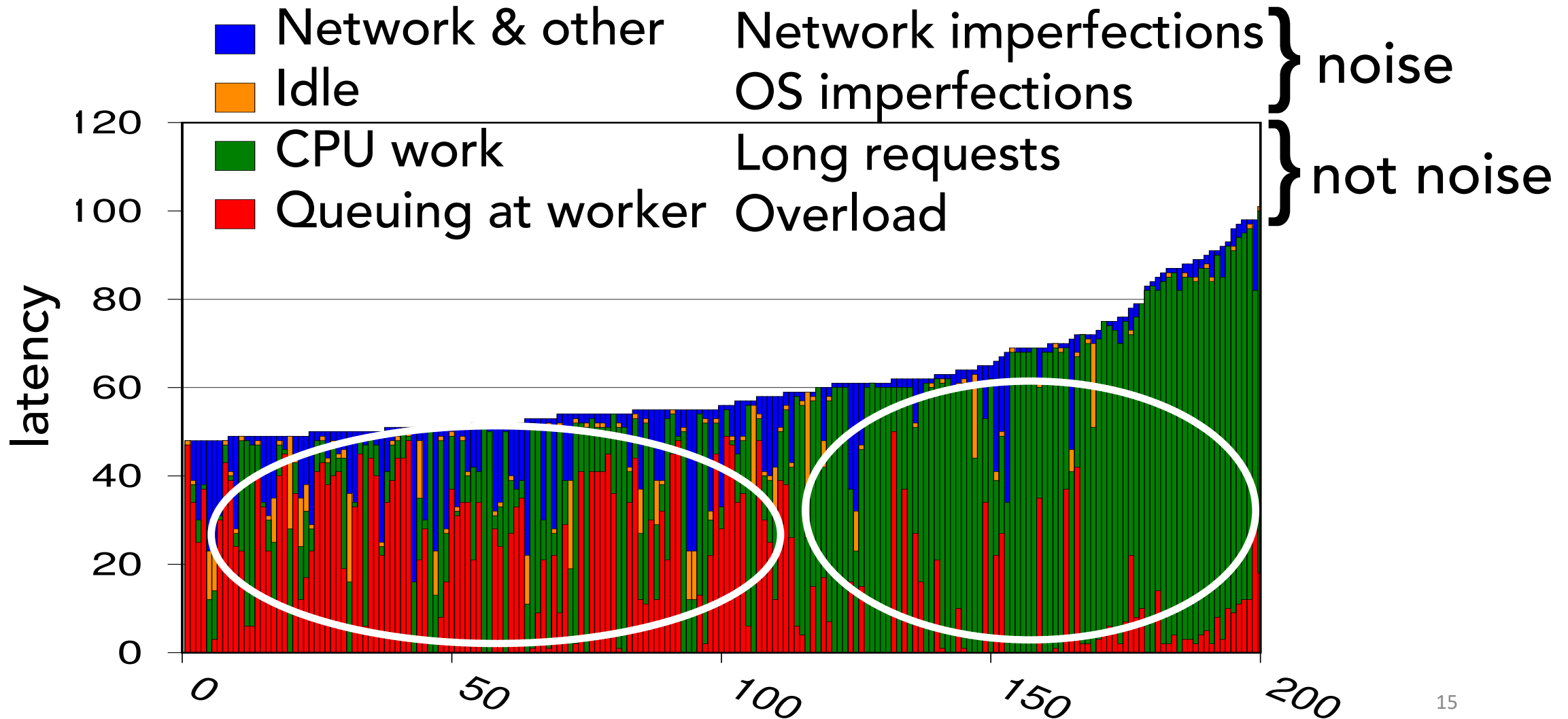
$$\text{HT1 IPC} = \text{Core IPC} - \text{HT2 SHIM IPC}$$



# What is in the tail?



# The Tail Longest 200 requests



# Optimizing the tail

Diagnosing the tail with continuous profiling

Noise

systems are not perfect

Queuing

too much load is bad, but so is over provisioning

Work

many requests are long

Insights Use the CDF off line

Long requests reveal themselves, treat them specially

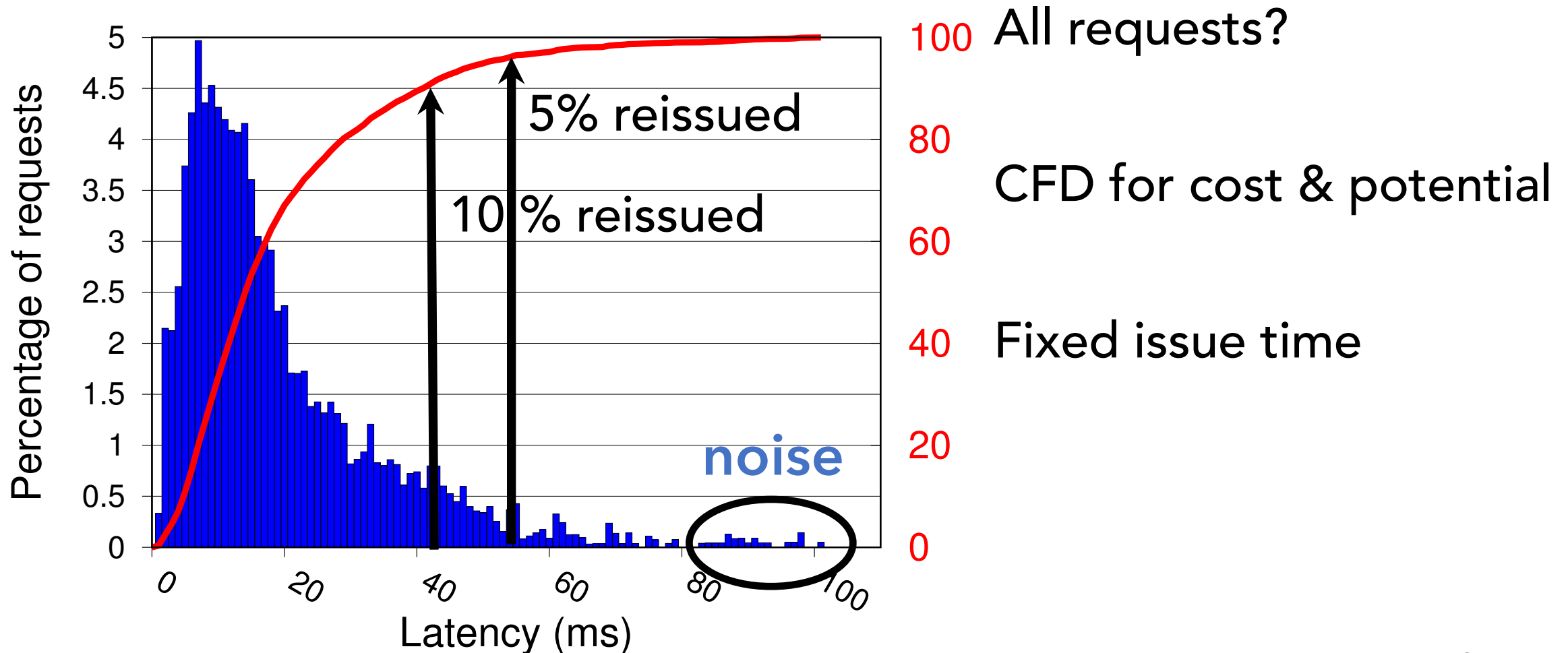
Insight

**Long requests reveal themselves**

Regardless of the cause

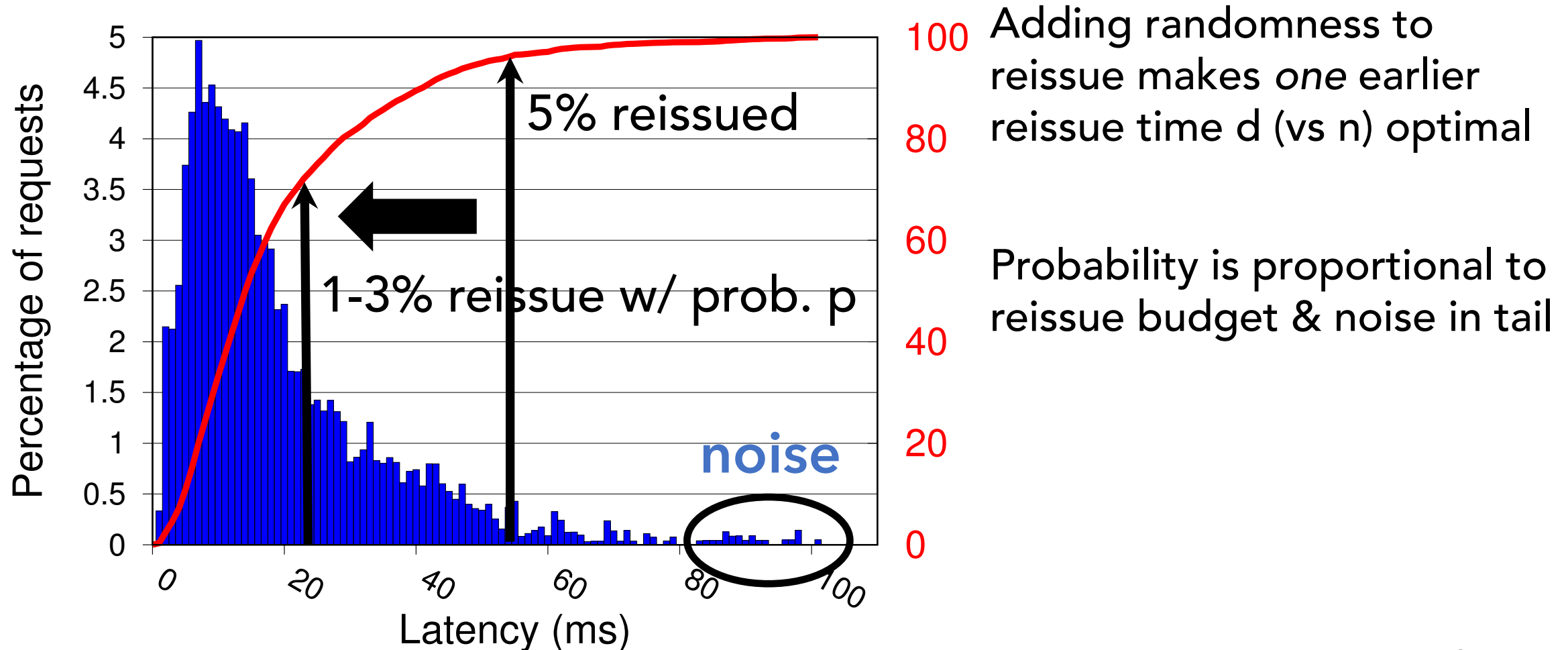
# Noise Replicate & reissue

The Tail at Scale, Dean & Barroso, CACM'13



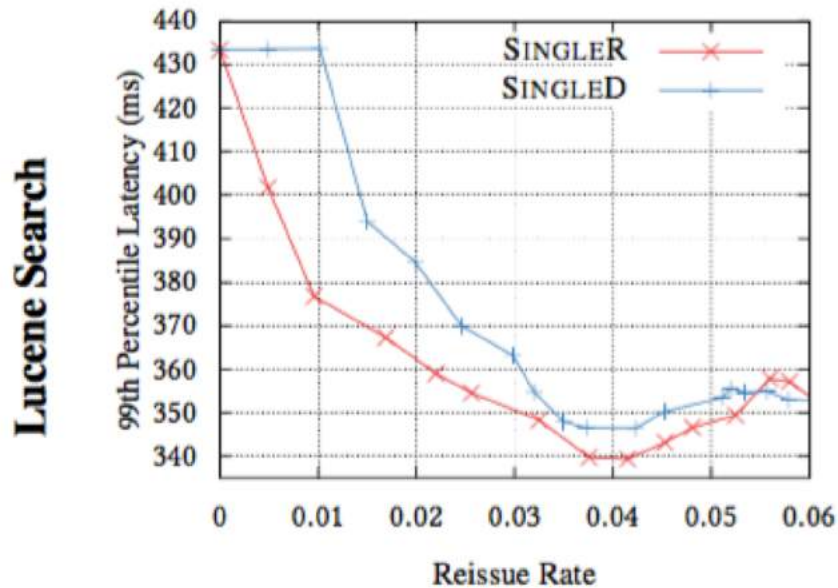
# Probabilistic reissue

Optimal Reissue Policies for Reducing Tail Latencies, Kaler, He, & Elnickety, SPAA'17

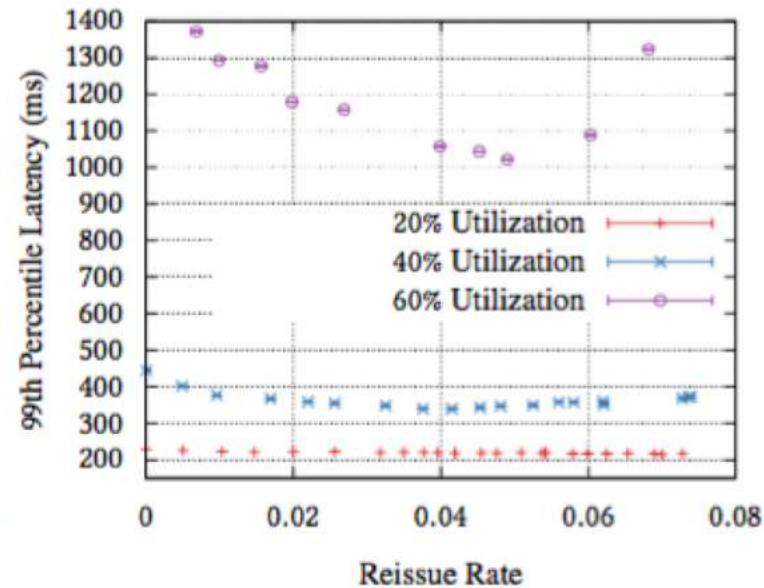


# Single R Probabilistic reissue

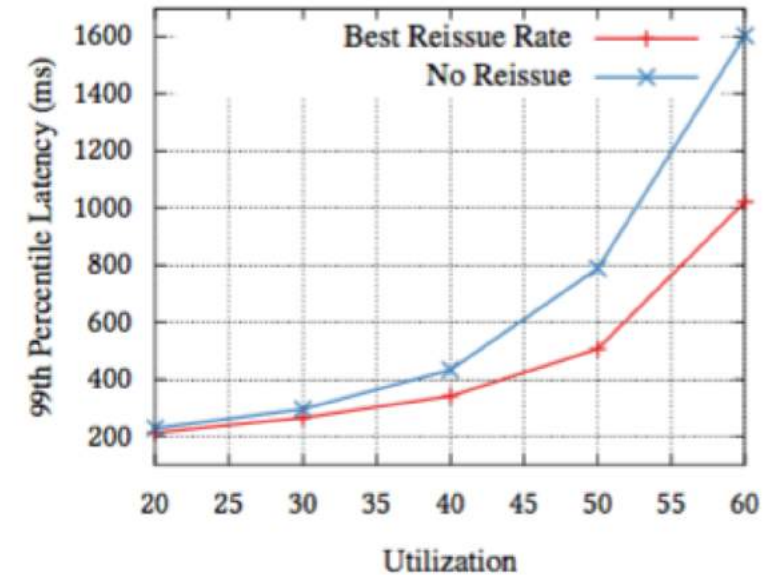
Optimal Reissue Policies for Reducing Tail Latencies, Kaler, He, & Elnickety , SPAA'17



(a) SINGLER vs SINGLED

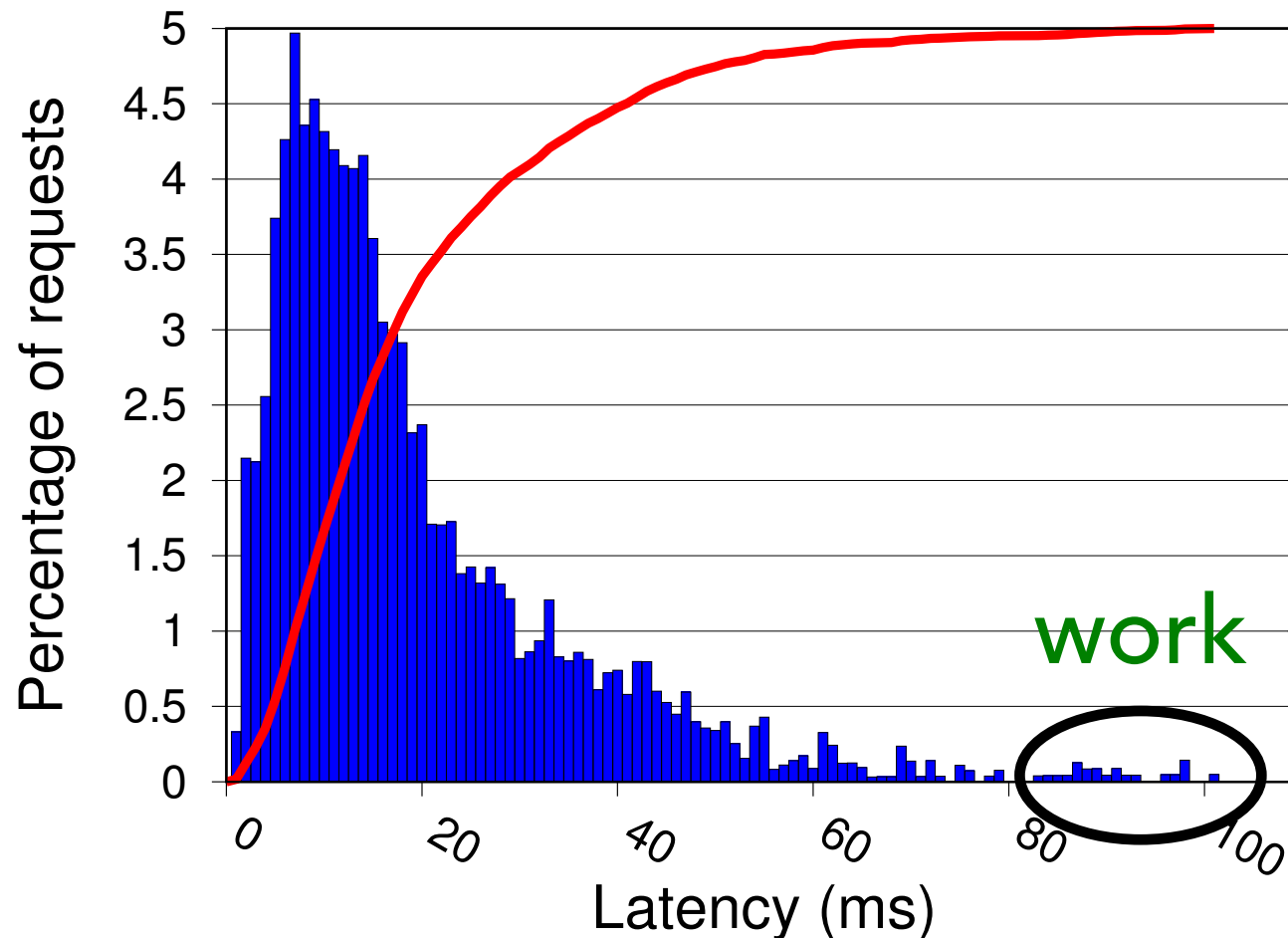


(b) Latency vs Reissue Rate



(c) Best Latency vs Utilization

# Work Speed up the tail *efficiently*

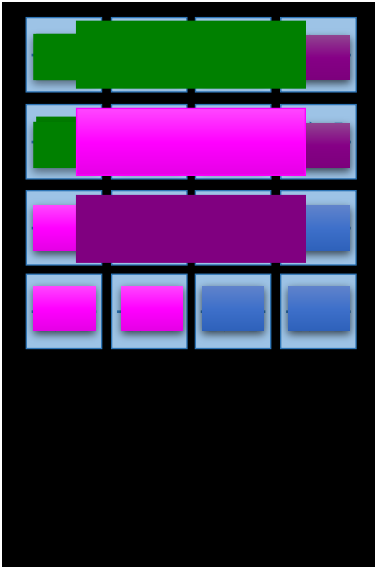


- 100 Judicious parallelism  
[ASPLOS'15]
- 80 DVFS faster on the tail  
[DISC'14, MICRO'17]
- 60 Asymmetric multicore  
[DISC'14, MICRO'17]

# Work Parallelism

Parallelism historically for **throughput**

Idea Parallelism for **tail latency**



# Queuing theory

Optimizing average latency maximizes throughput

But not the tail!

Shortening the tail reduces queuing latency

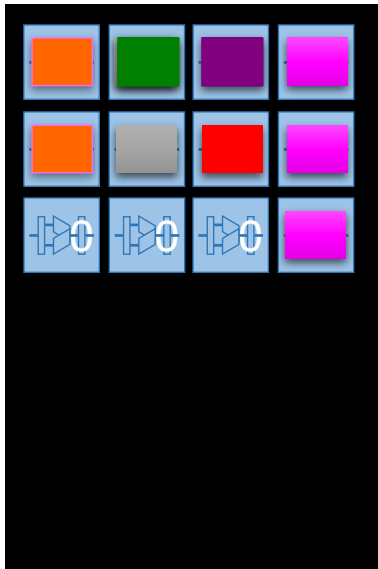
# Parallelism

Parallelism historically for **throughput**

**Idea** Parallelism for **tail latency**

**Insight** Long requests reveal themselves

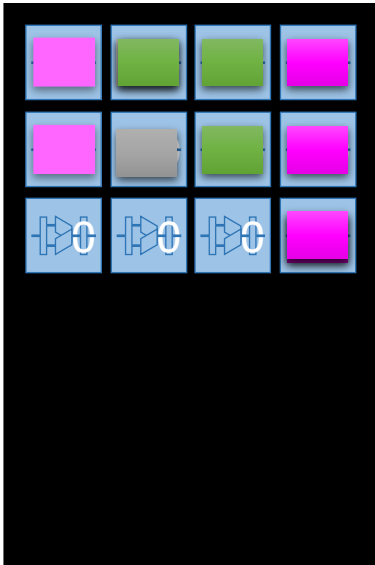
**Approach** Incrementally add parallelism to long requests – the tail – based on request progress & load



# Few to Many

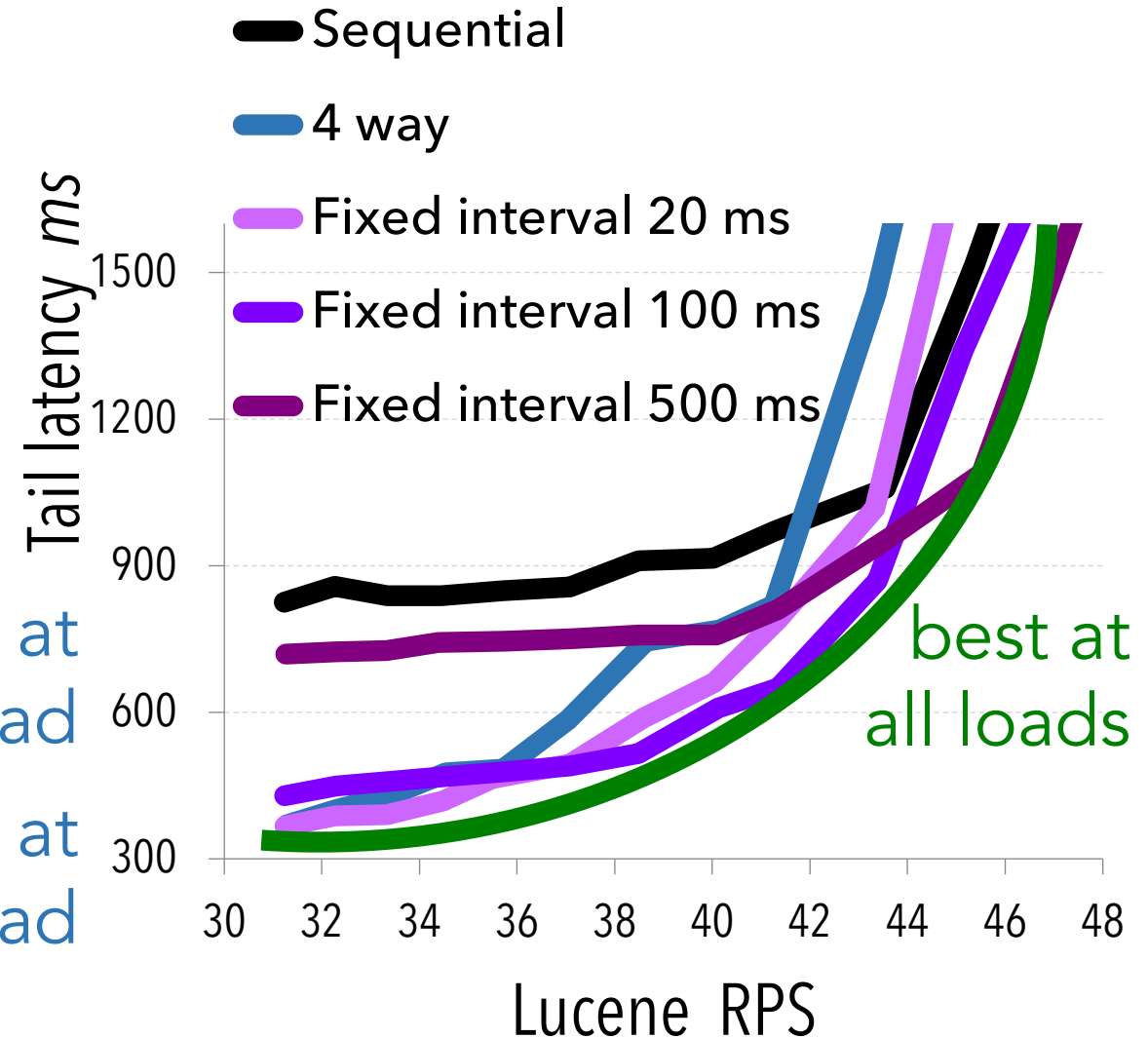
Fixed: add thread every  $d$  ms

Dynamic: use load

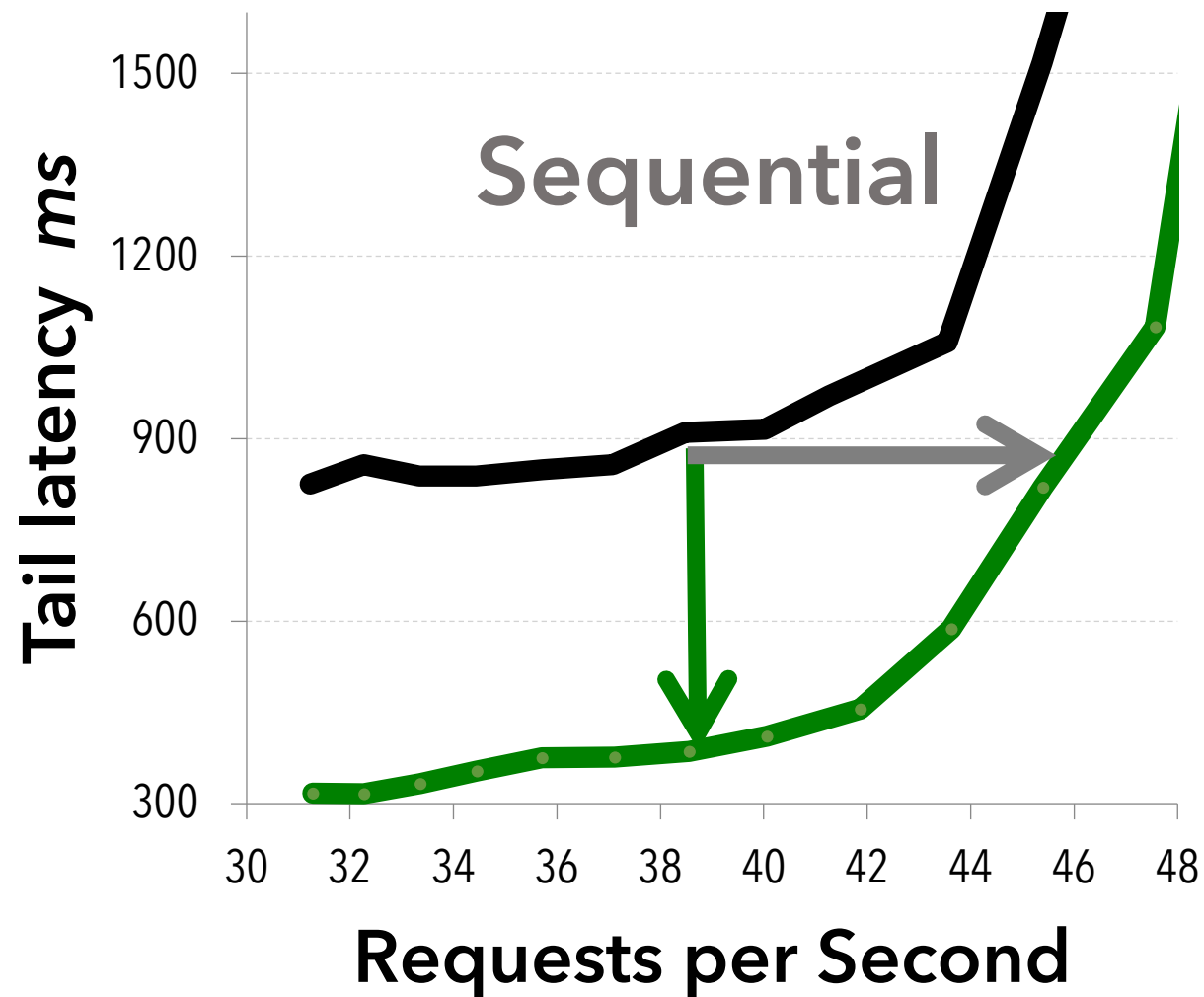


short delay good at  
low load

long delay good at  
high load



# Evaluation 2x8 64 bit 2.3 GHz Xeon, 64 GB

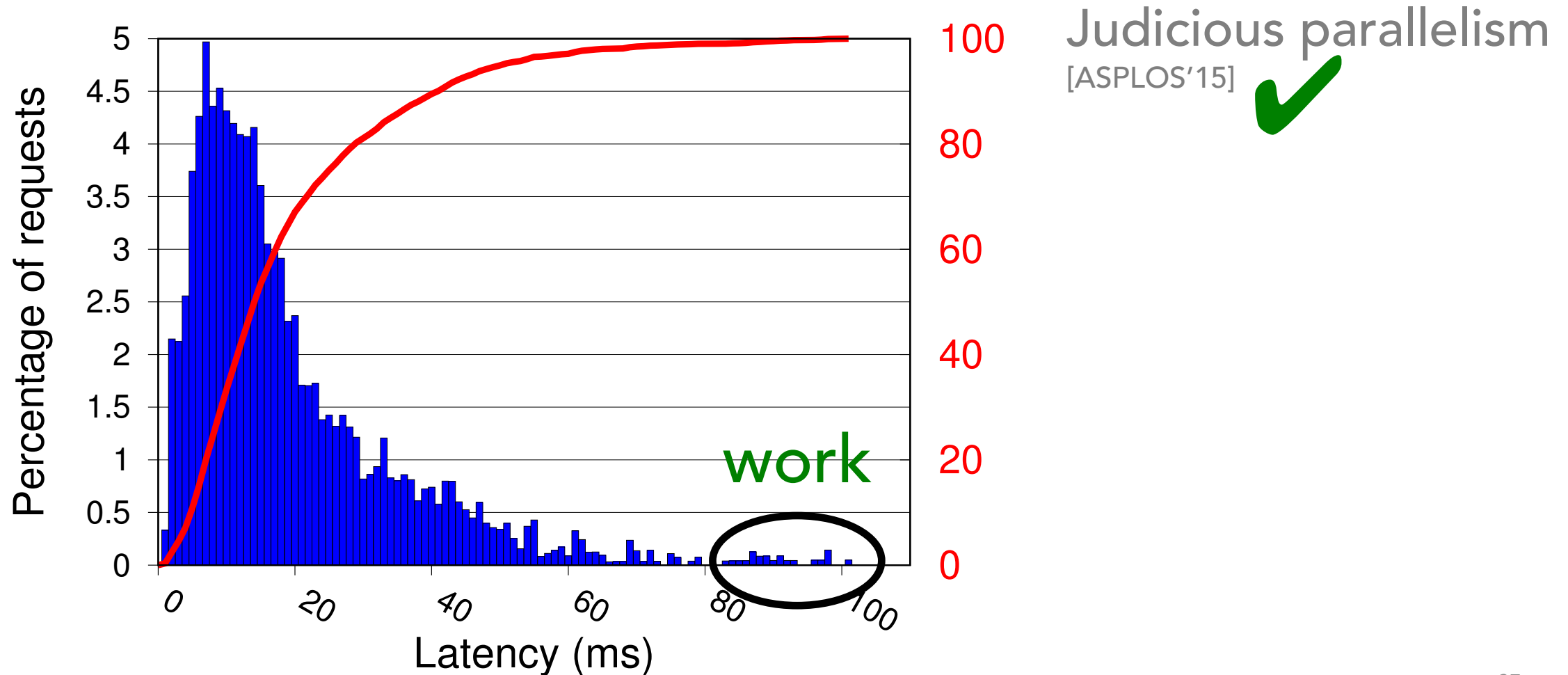


**Dynamic parallelism**  
**Few to Many**

21% fewer servers

or reduce tail by 28%

# Work speed up the tail *efficiently*





**Tail Latency**



**Efficiency**



**Google Cloud**

# Efficiency at scale for interactive workloads

Diagnosing the tail with continuous profiling

**Noise** replication, systems are not perfect

**Queuing** replication + judicious choice

**Work** judicious use of resources on long requests

Request latency CDF is a powerful tool

Tail efficiency  $\neq$  average or throughput

Hardware heterogeneity

## Questions?

# **Professional and Research Relationships**

# Your Academic Village

- Peer students
- Students senior & junior to you
- Teaching assistants
- PhD students
- Faculty



**CRA-W**

Computing Research Association  
Women

# My Professional Village

- Researchers in all career stages
  - Undergrads, PhD students, post docs
  - Faculty, industrial researchers, staff, administrators
- Industrial village
  - Software engineers in all career stages
  - Managers, directors, admins,
  - in/out my management chain



**CRA-W**

Computing Research Association  
Women

# Faculty Mentors

Don Johnson



My Professor

Ken Kennedy



PhD Advisor

Dave Stemple



Dept. Chair



**CRA-W**

Computing Research Association  
Women

# Building a Village



**CRA-W**

Computing Research Association  
Women

# Networking is....

## *Building and sustaining professional relationships*

- Participating in an academic / research community
- Finding people you like and you learn from, and building a relationship



**CRA-W**

Computing Research Association  
Women

# Networking is *not*....

- Using people
- A substitute for quality work



**CRA-W**

Computing Research Association  
Women

# But I am Horrible at Small Talk

- You have CS in common
- Networking is not genetic
- It is a research skill
  - **Practice**
  - **Meet people**
  - **Learn**
  - **Go places**
  - **Volunteer!**
  - **Sustain your relationships**



**CRA-W**

Computing Research Association  
Women

# With whom do you network?

- People you like
- People senior to you, who can show you the way
- People at different career stages, so you can anticipate
- Your peers



**CRA-W**

Computing Research Association  
Women

# Peer Mentors

Mary Hall



Doug Burger



Margaret Martonosi



**CRA-W**

Computing Research Association  
Women

# Your Village Will

- Write letters for grad school, jobs, etc.
- Help you solve problems
- Point you in good directions
- Encourage you
- Choose you for important roles
- You will do the same or more for them
- Make your life and work more fun and meaningful



**CRA-W**

Computing Research Association  
Women