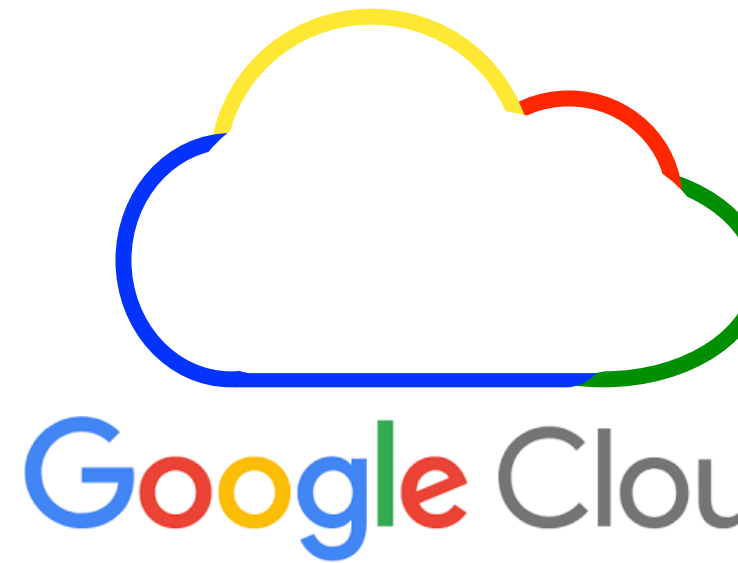# Measuring and Optimizing Tail Latency

**athryn S McKinley, Google**

Yang, Stephen M Blackburn,

d Haque, Sameh Elnikety, Yuxiong He, Ricardo Bianchini
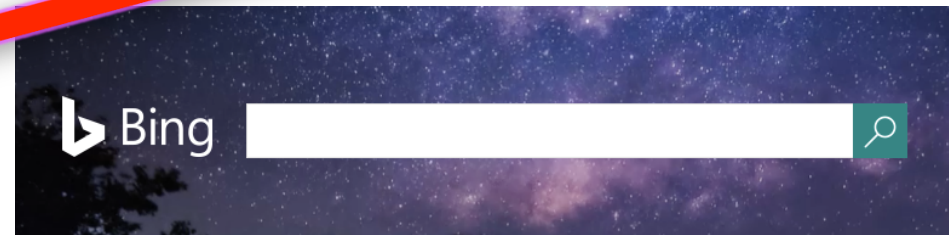
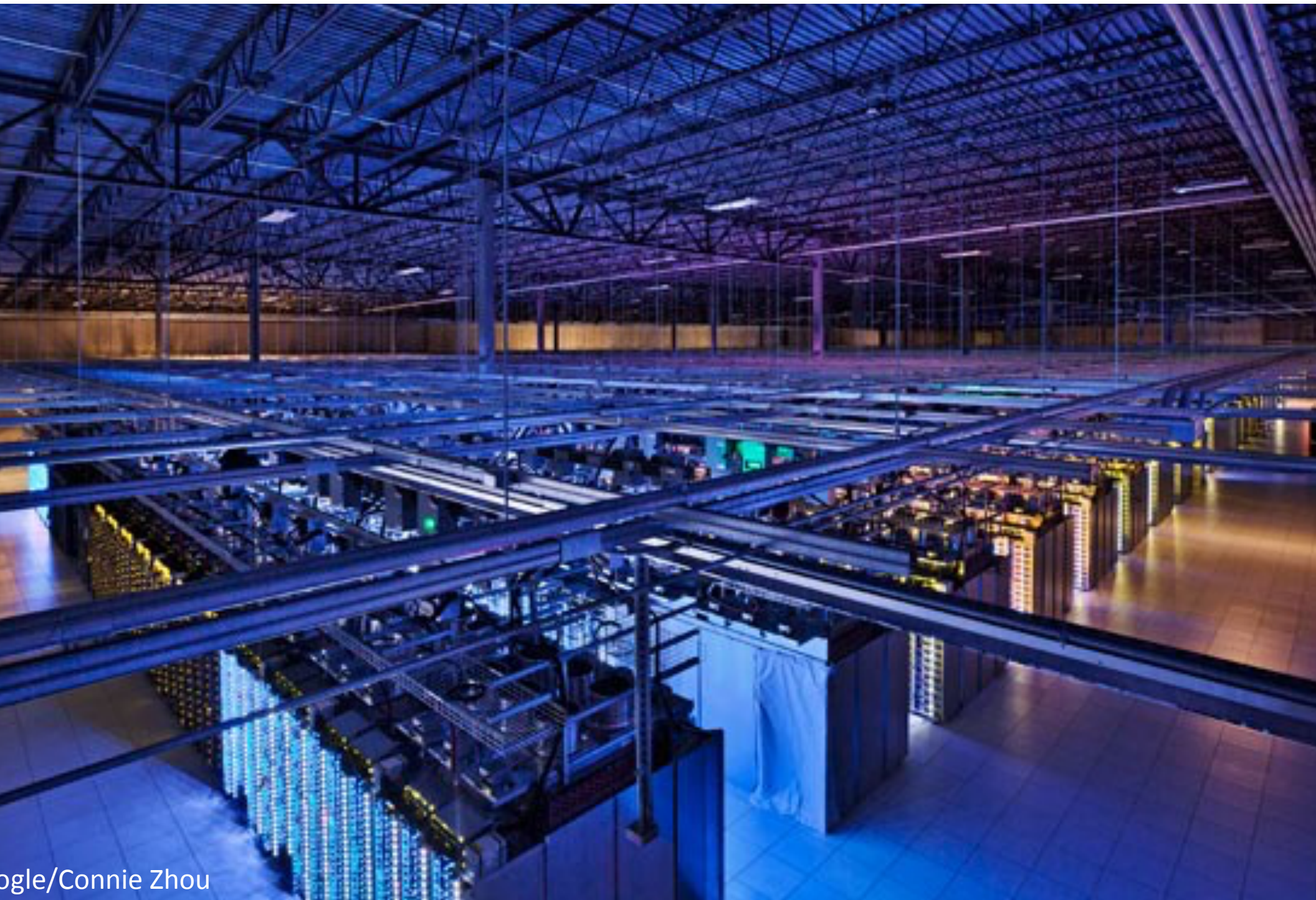Google Clou

# Tail Latency Matters

**TOP PRIORITY**

Google Australia

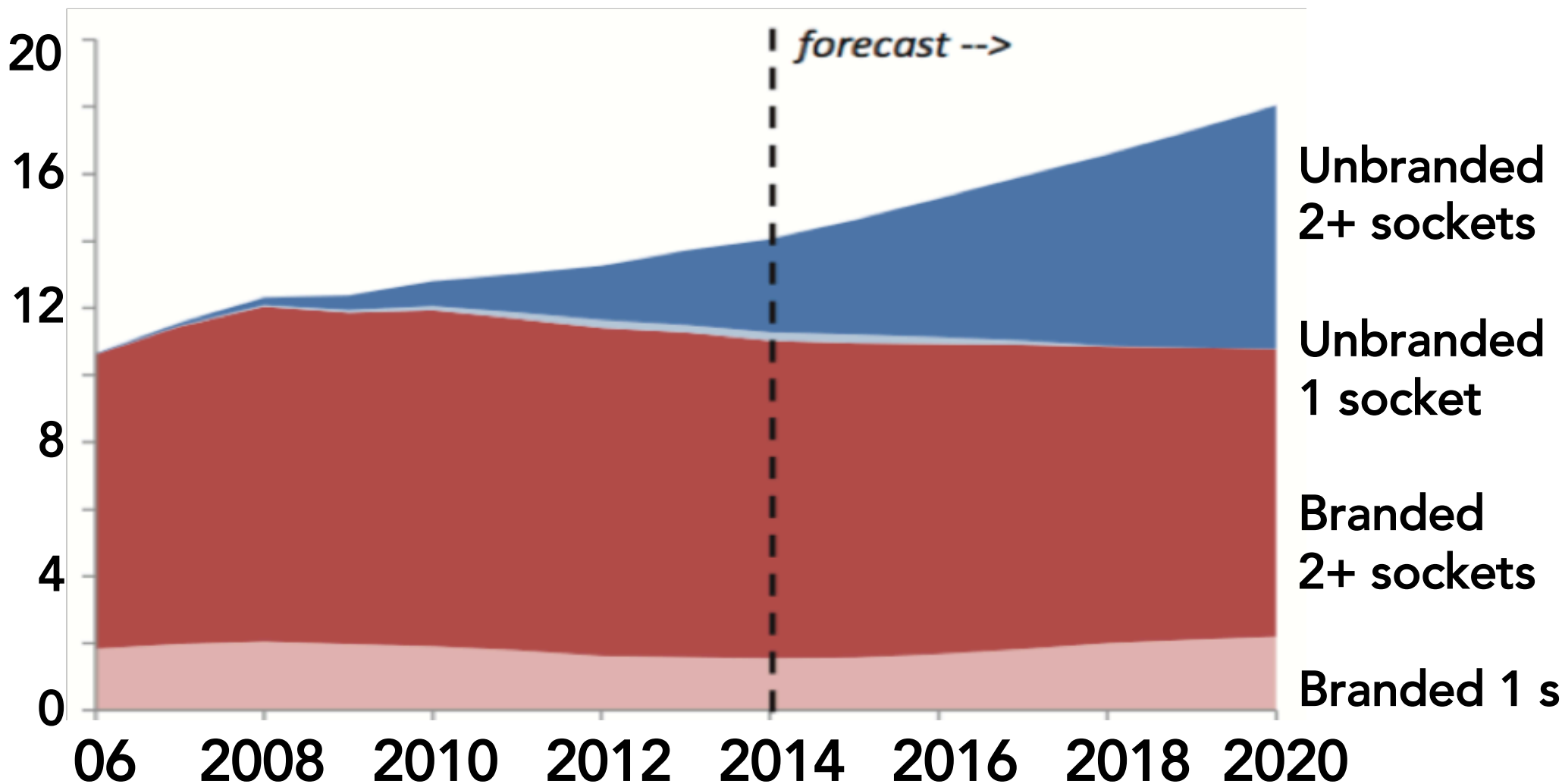00 millisecond delay decreased arches/user by 0.59%. [Jack Brutlag, Google]

Bing

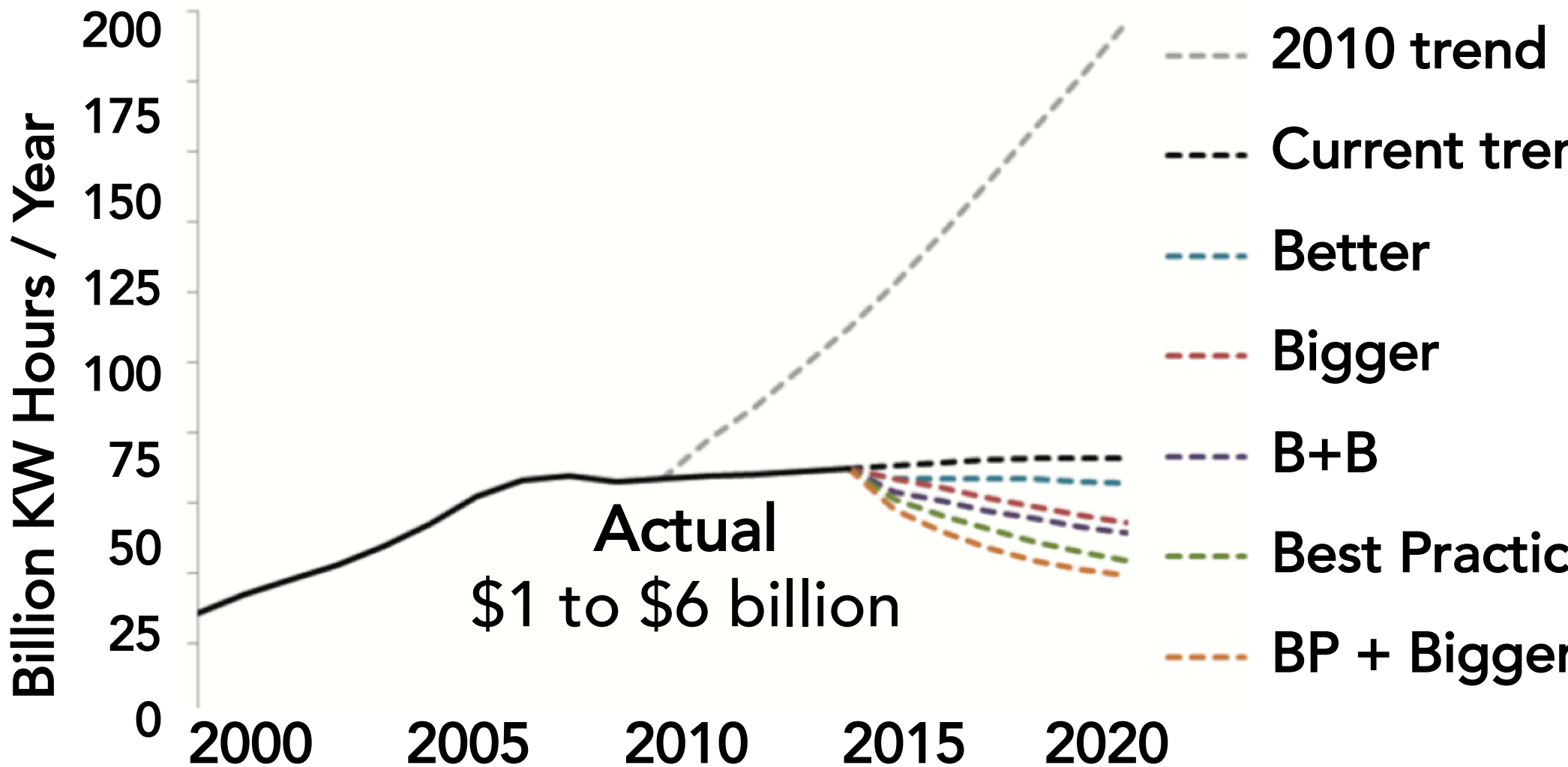Two second slowdown reduced revenue/user by 4.3%. [Eric Schurman, B...

# Servers in US datacenters



*Shehabi et al., United States Data Cen[...]
Usage Report, Lawrence Berkeley, 201[...]

# Electricity in US datacenters

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

# Datacenter economics quick facts*

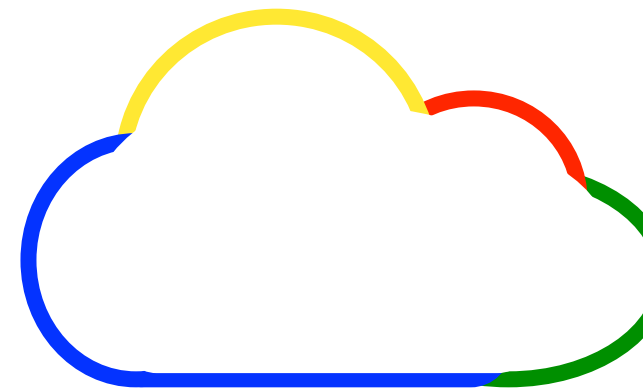| | |
|---|---|
| ~ $500,000 | Cost of small datacenter |
| ~3,000,000 | US datacenters in 2016 |
| ~ $1.5 trillion | US Capital investment to date |
| ~ $3,000,000,000 | KW dollars / year |
| ~ $30,000,000 | Savings from 1% less work |
| Lots more | by not building a datacenter |

*Shehabi et al., United States Data Cen
Usage Report, Lawrence Berkeley, 2016
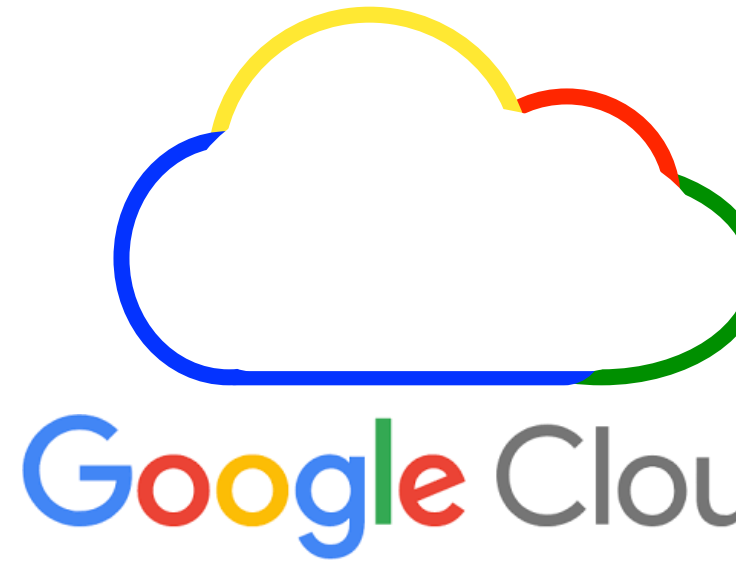
**TOP PRIORITY**

**Efficiency**
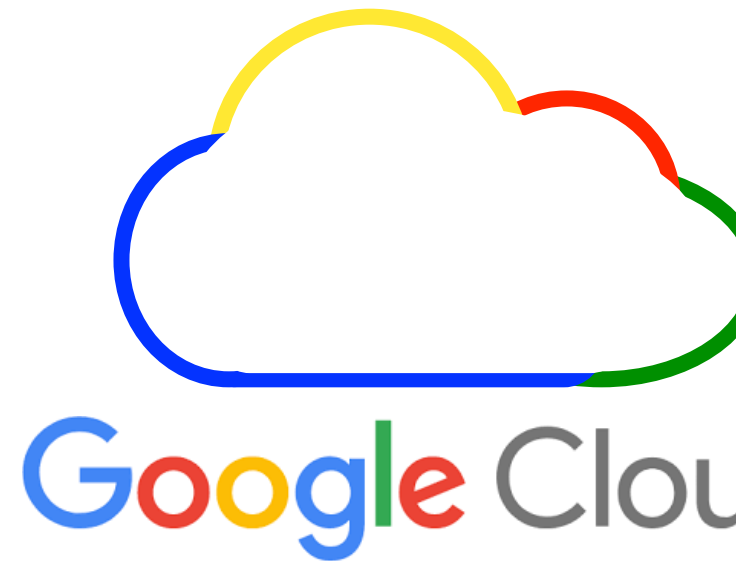
Google Cloud

il Latency

**TOP PRIORITY**

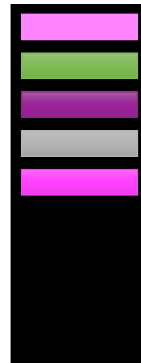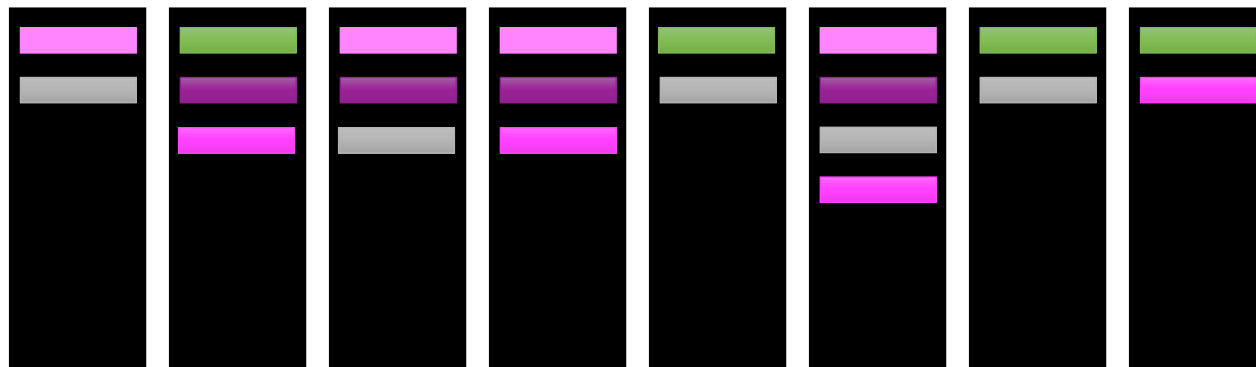**Efficiency**

Google Cloud

il Latency

**BOTH ?!**

**Efficiency**

Google Clou

# Server architecture

client
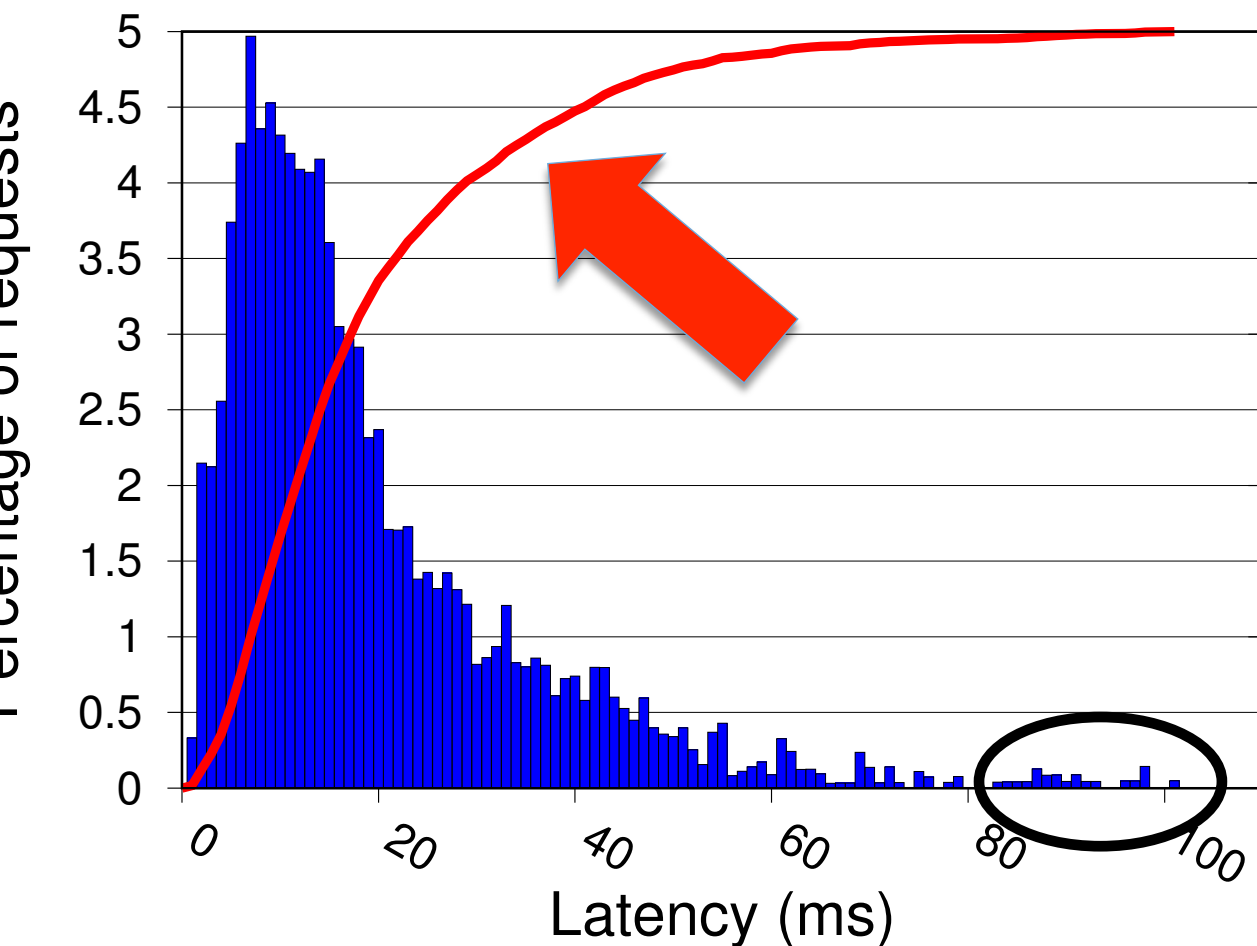
aggregator

workers

# Characteristics of interactive services



LC *Lucene*

Bursty, diurnal

**CDF** changes slowly

Slowest server dictates

Orders of magnitude di

average & tail - 99th %

# What is in the tail?

# Roadmap

Diagnosing the tail with continuous profiling

| | |
|---|---|
| **Noise** | systems are not perfect |
| **Queuing** | too much load is bad, but so is over provisionir |
| **Work** | many requests are long |

**Insights**  Use the CDF off line

Long requests reveal themselves, treat them specially

# Simplified life of a request

request    ↑ response    **client application / OS**

**network**

**aggregator SW / OS**

**network**

worker OS / VM    worker OS / VM    worker OS / VM

**Java VM**    **Java VM**    **·  ·  ·**    **Java VM**

**application**    **application**    **application**

# Prior state of the art

Dick Site, Google   https://www.youtube.com/watch?v=QBu2Ae8-8LM

# @ Google

- Hand instrument system
- % on-line budget
  - sample – but tails are rare…
- Off-line schematics
- Have insight
- Improve the system



LOCKs 0..22

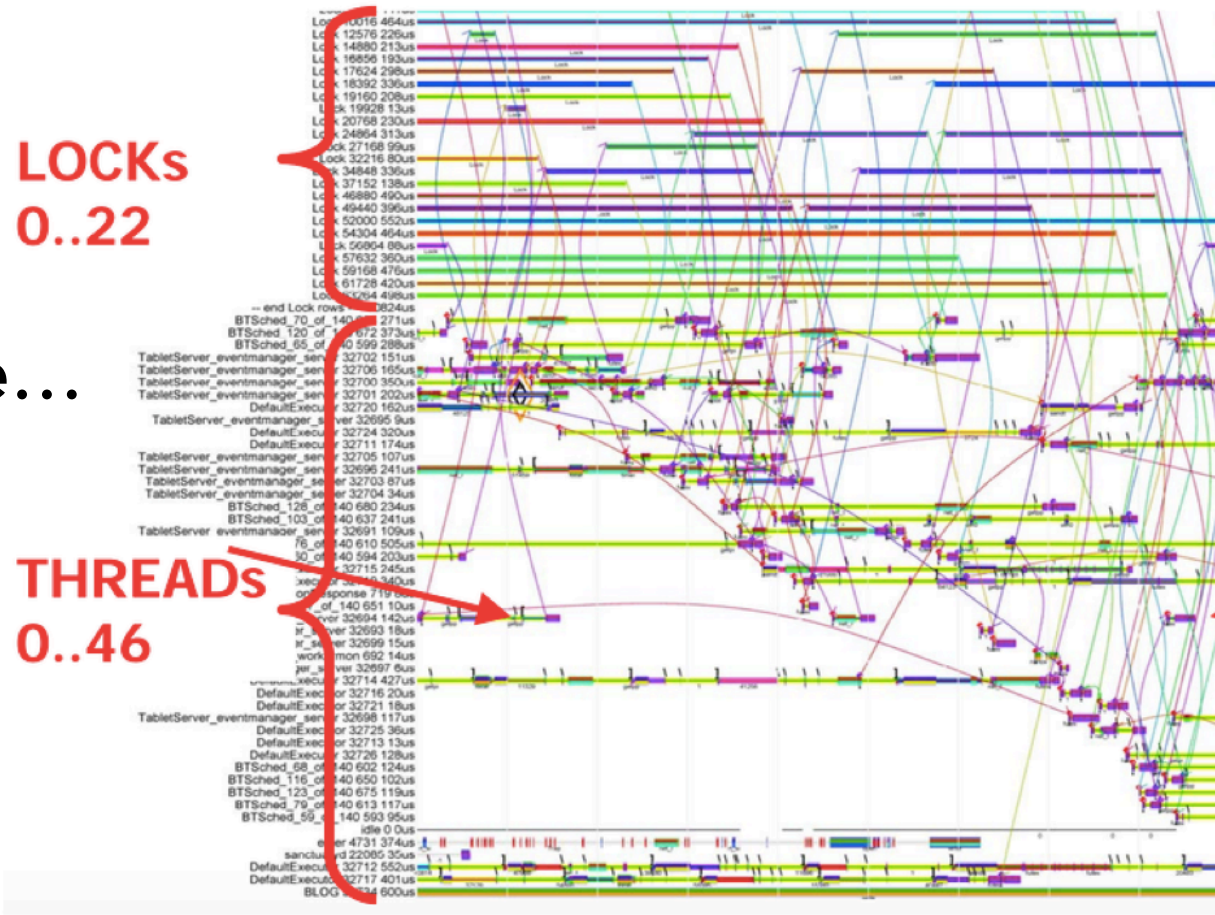THREADs 0..46

# Request profiling

Hand instrument system

% on-line budget

  sample – but tails are rare…

Off-line schematics

Have insight

Improve the system

# Request profiling

Hand instrument system     Automated instrumentation

% on-line budget     1% on-line budget

sample – but tails are rare…     continuous on-line profiling

Off-line schematics     Off-line schematics

Have insight     Have insight

Improve the system     Improve the system

    + On-line optimization

# Automated cycle-level on-line profiling

**Insight**    Hardware & software generate signals

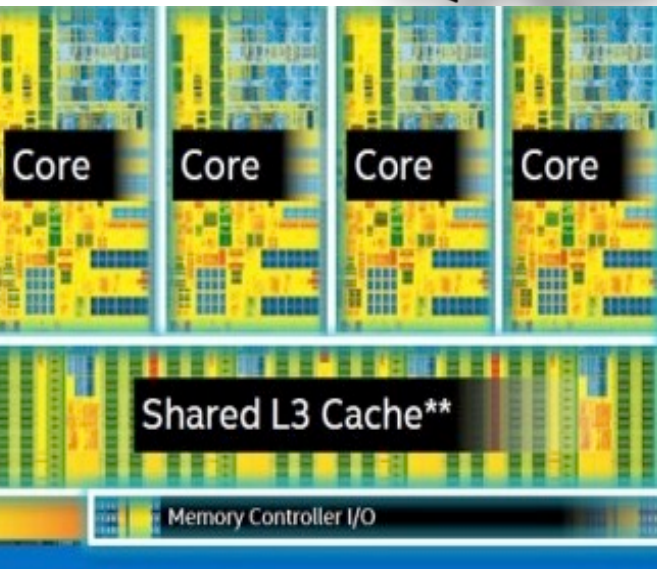|  | hardware signals<br>performance counters | software signals<br>memory locations |
|---|:---:|:---:|
| counters | ✓ | ✓ |
| tags | ✓ | ✓ |

# SHIM Design

SCA'15 (Top Picks HM), ATC'16

# Observe global state from other core

```
while (true):
    for counter in LLC misses, cycles:
        buf[i++] = readCounter(counter)
```
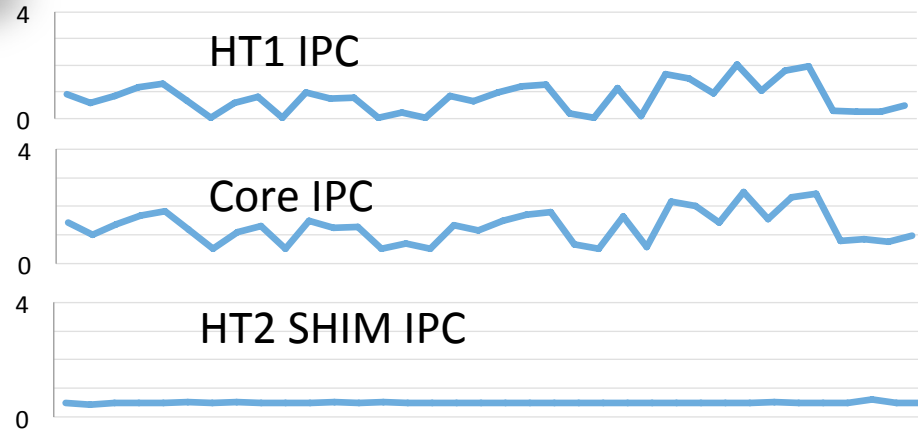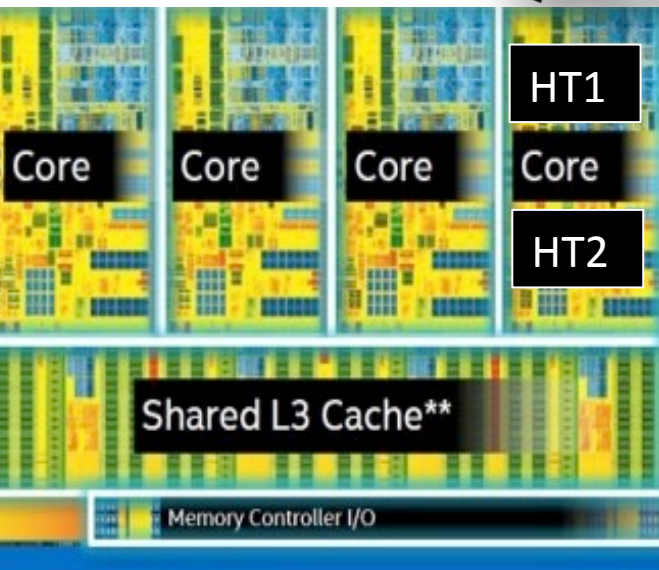


LLC misses per cycle
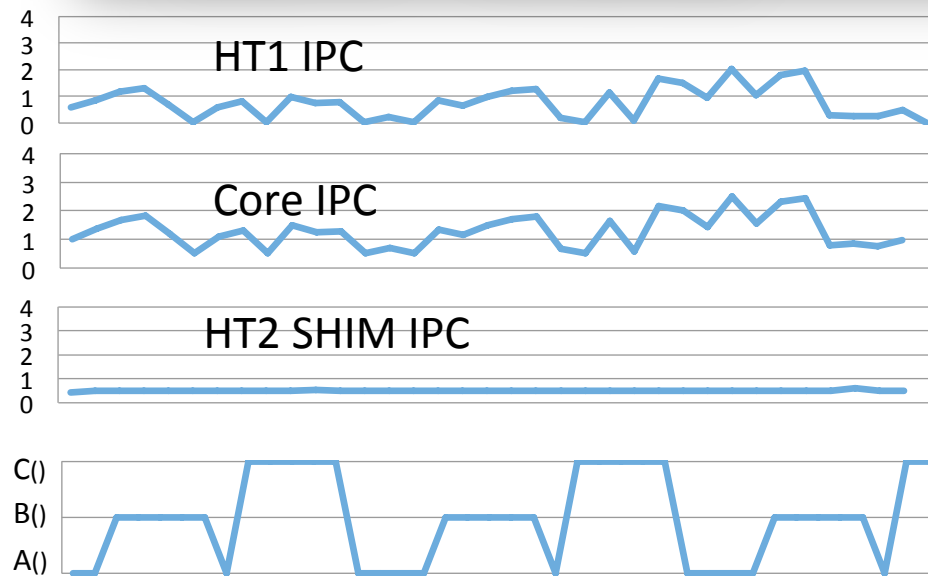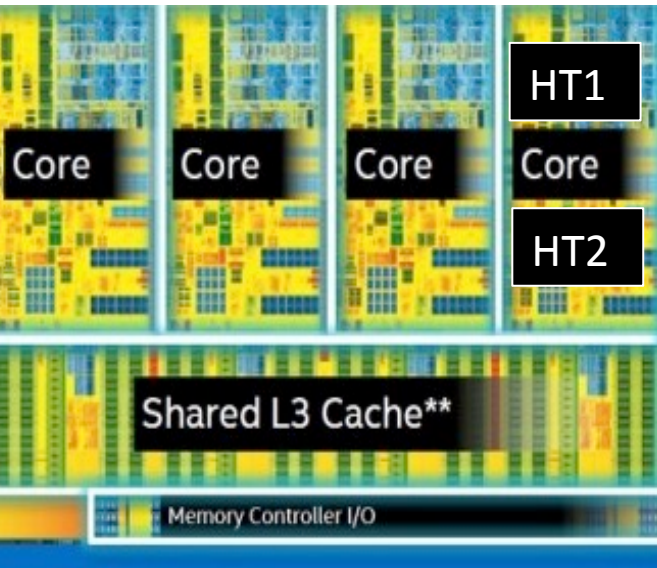
# Observe local state with SMT hardware

```
while (true):
    for counter in HT2 SHIM, Core, Cycles:
        buf[i++] = readCounter(counter);
```

HT1 IPC

Core IPC

HT2 SHIM IPC

HT1 IPC = Core IPC − HT2 SHIM IPC

# Correlate hardware & software events

```
while (true):
    for counter in HT2 SHIM, Core, cycles:
        buf[i++] = readCounter(counter);
    tid = thread on HT1
    buf[i++] = tid.method;
```
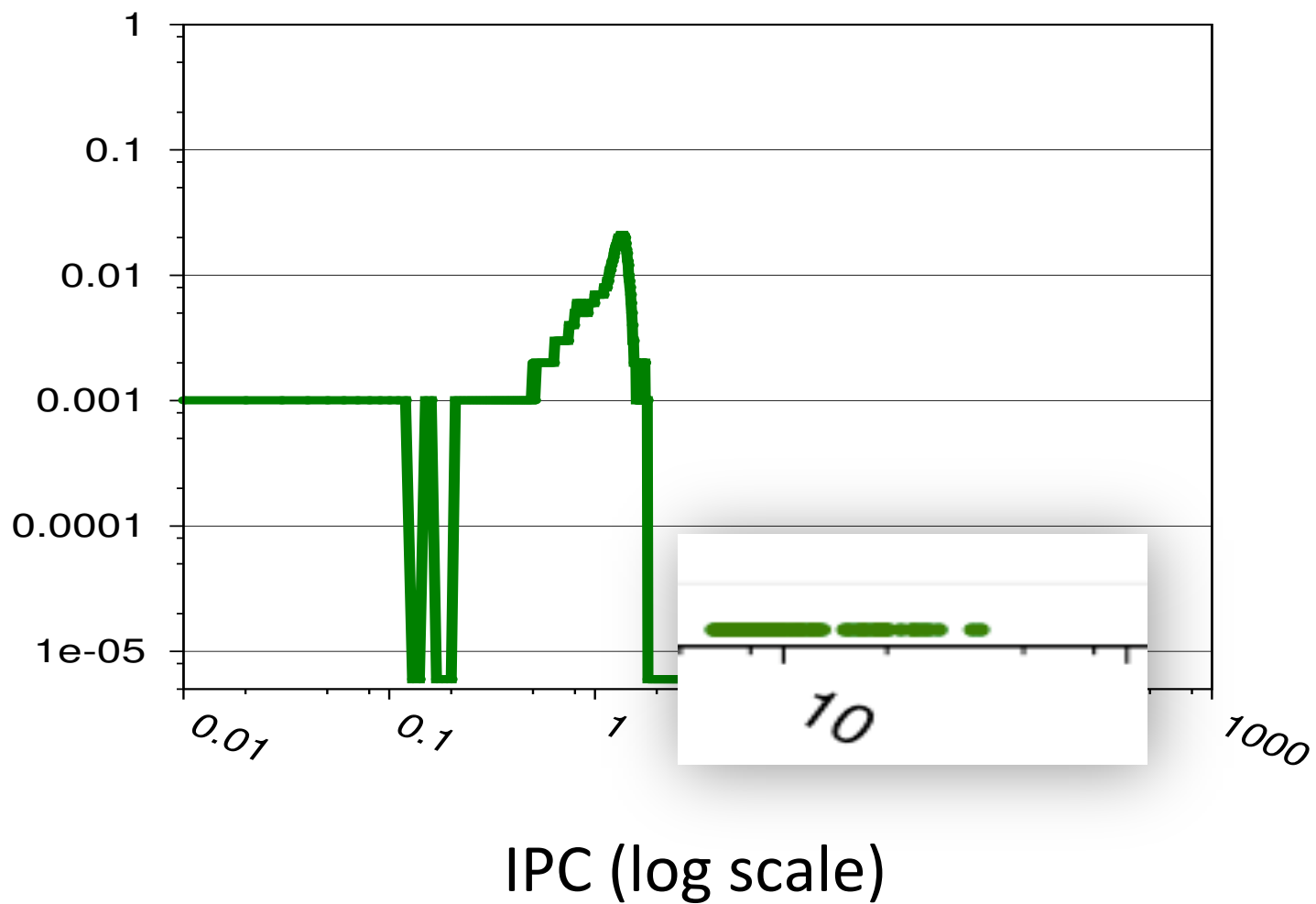
Fidelity
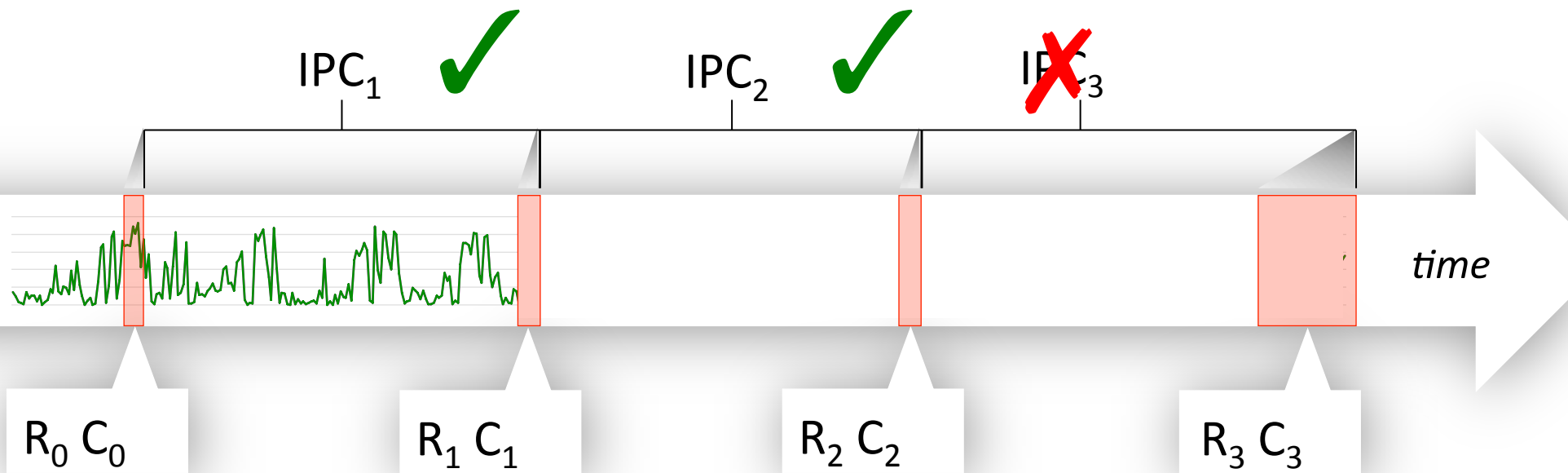
# Raw samples



samples
scale)

IPC (log scale)

# Problem: samples are not atomic

**Counters**   C: cycles    R: retired instructions
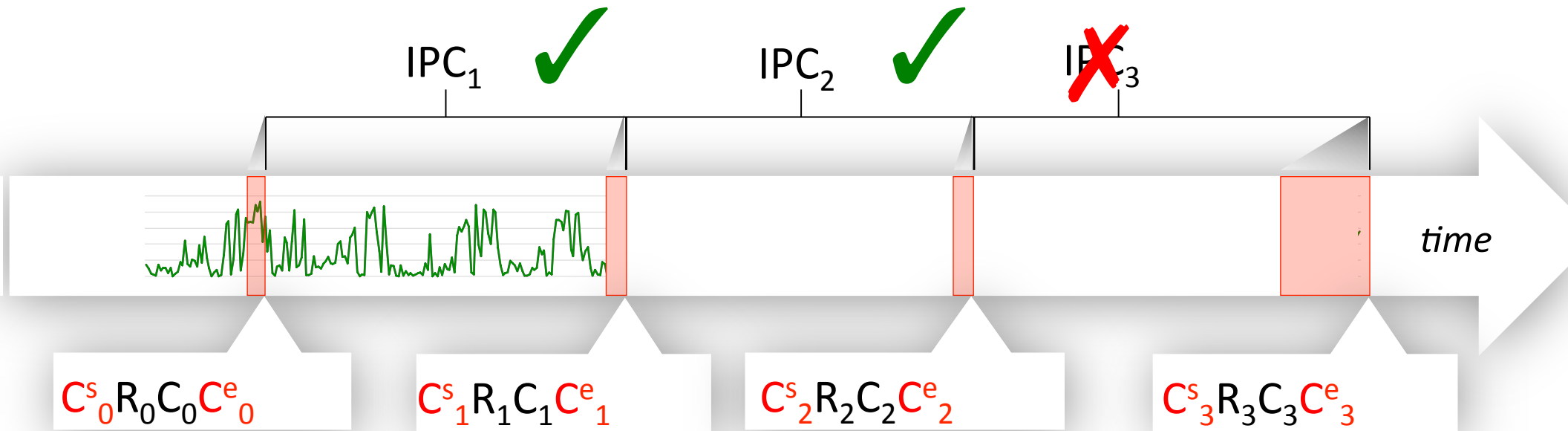
$$IPC = (R_t - R_{t-1}) / (C_t - C_{t-1})$$

# Solution: use clock as ground truth

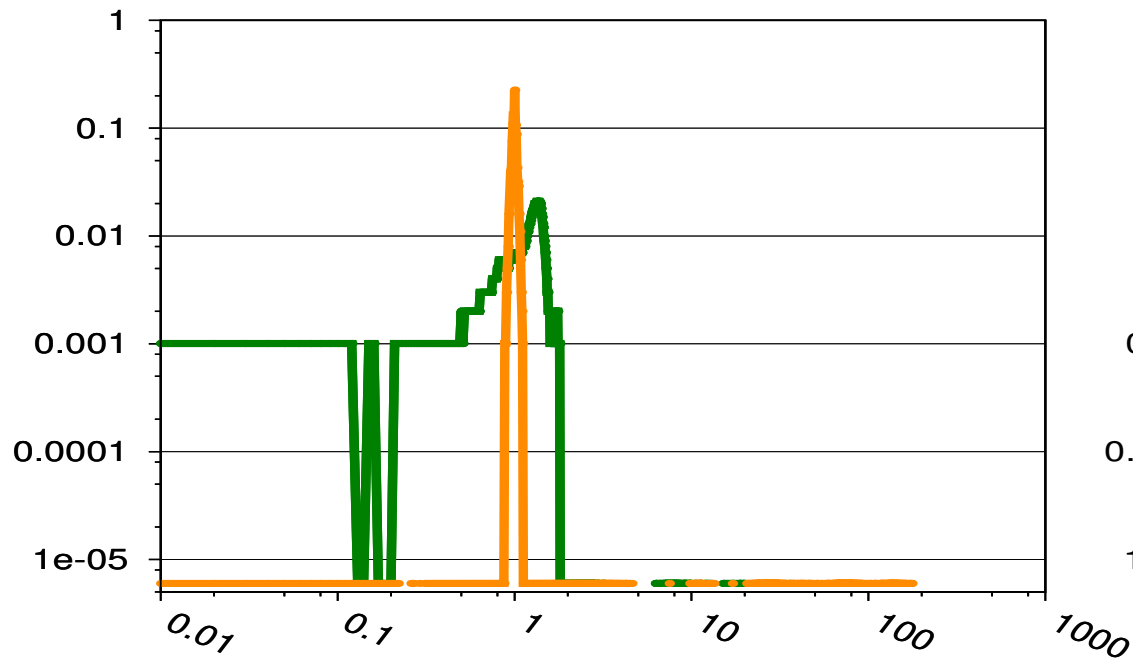$$CPC = (C^e_t - C^e_{t-1}) / (C^s_t - C^s_{t-1}) \quad \textit{this should be 1!}$$

CPC$_1$ = 1.0 +/- 1%        CPC$_2$ = 1.0 +/- 1%        CPC$_3$ != 1.0 +/- 1%

IPC$_1$    ✓        IPC$_2$    ✓        IPC$_3$    ✗

time

$C^s_0 R_0 C_0 C^e_0$        $C^s_1 R_1 C_1 C^e_1$        $C^s_2 R_2 C_2 C^e_2$        $C^s_3 R_3 C_3 C^e_3$
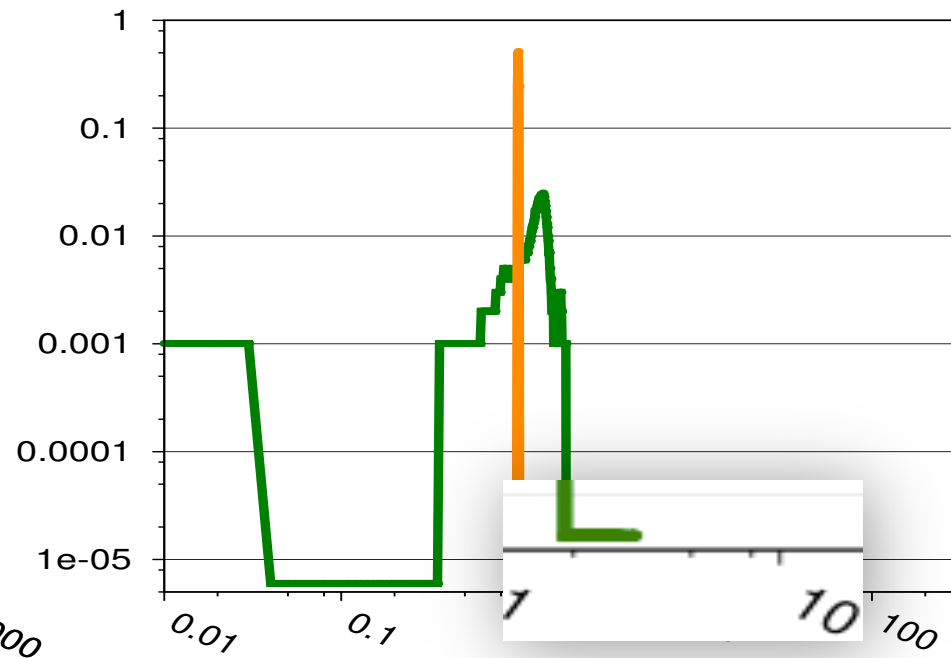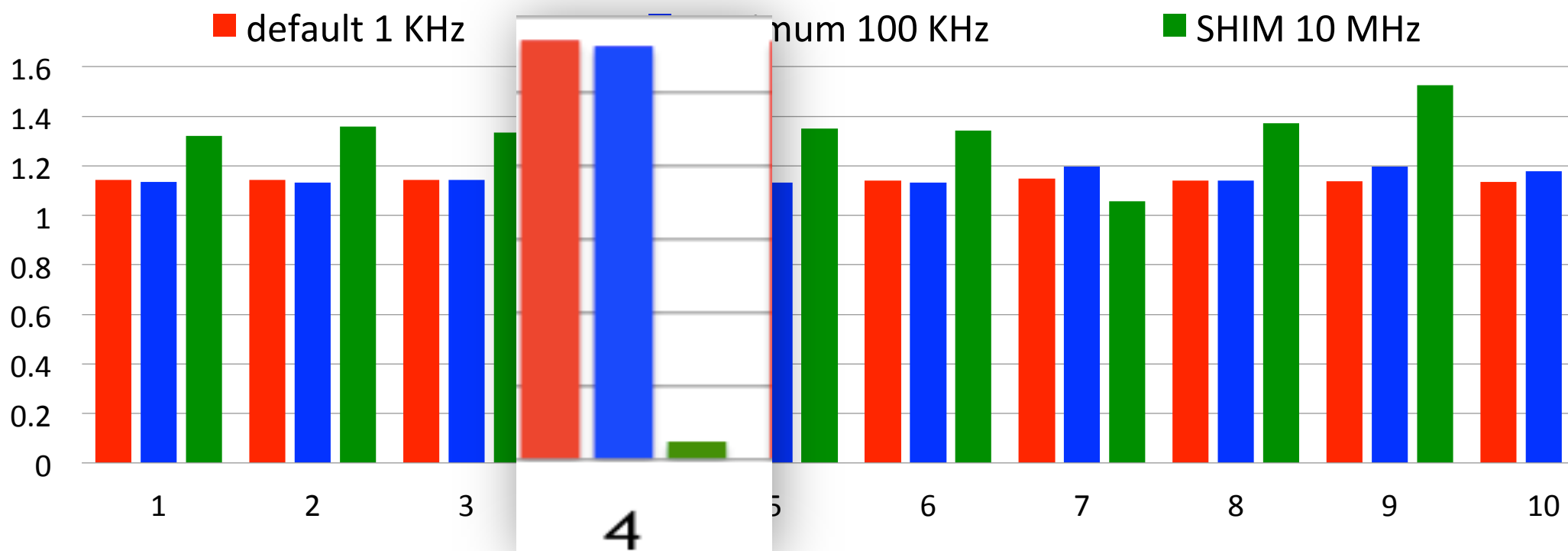
# Filtering Lusearch IPC samples



---- raw IPC

---- raw CPC

---- filtered IPC

---- filtered CPC in [0.99,1.01]

# PC of individual methods in Lucene
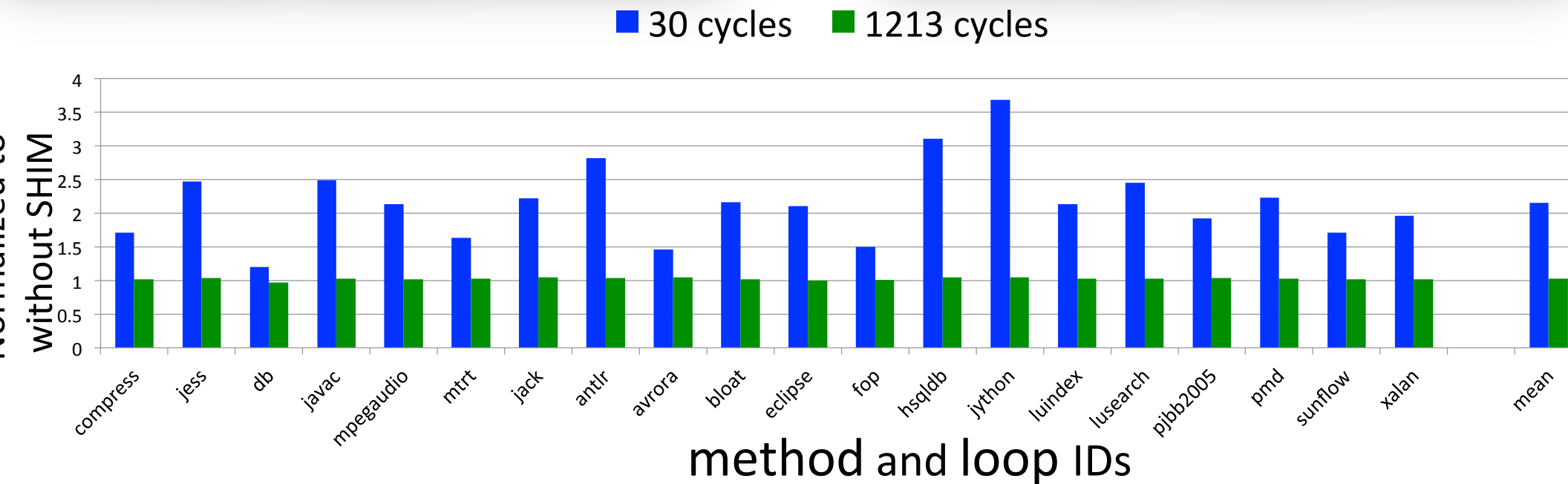


Legend: default 1 KHz | num 100 KHz | SHIM 10 MHz

top 10 methods (74% total execution time)

# Overheads from other core

113MHz: 3+ orders of magnitude over interrupt 'maximum'

3MHz: 1+ order of magnitude over interrupt 'maximum'

■ 30 cycles  ■ 1213 cycles



Overheads from write invalidations

# Understanding Tail Latency

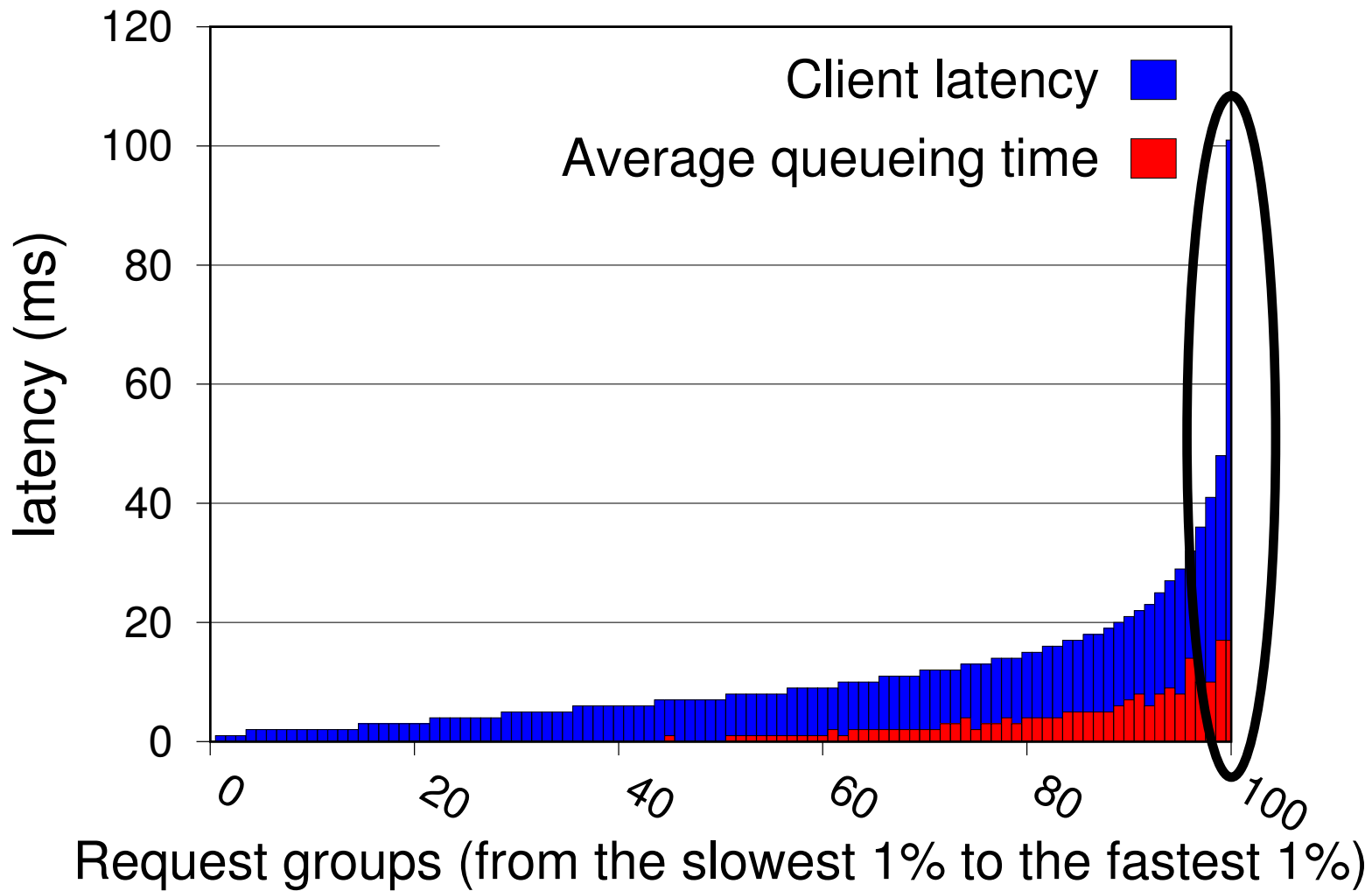# SHIM signals

Requests
- thread ids
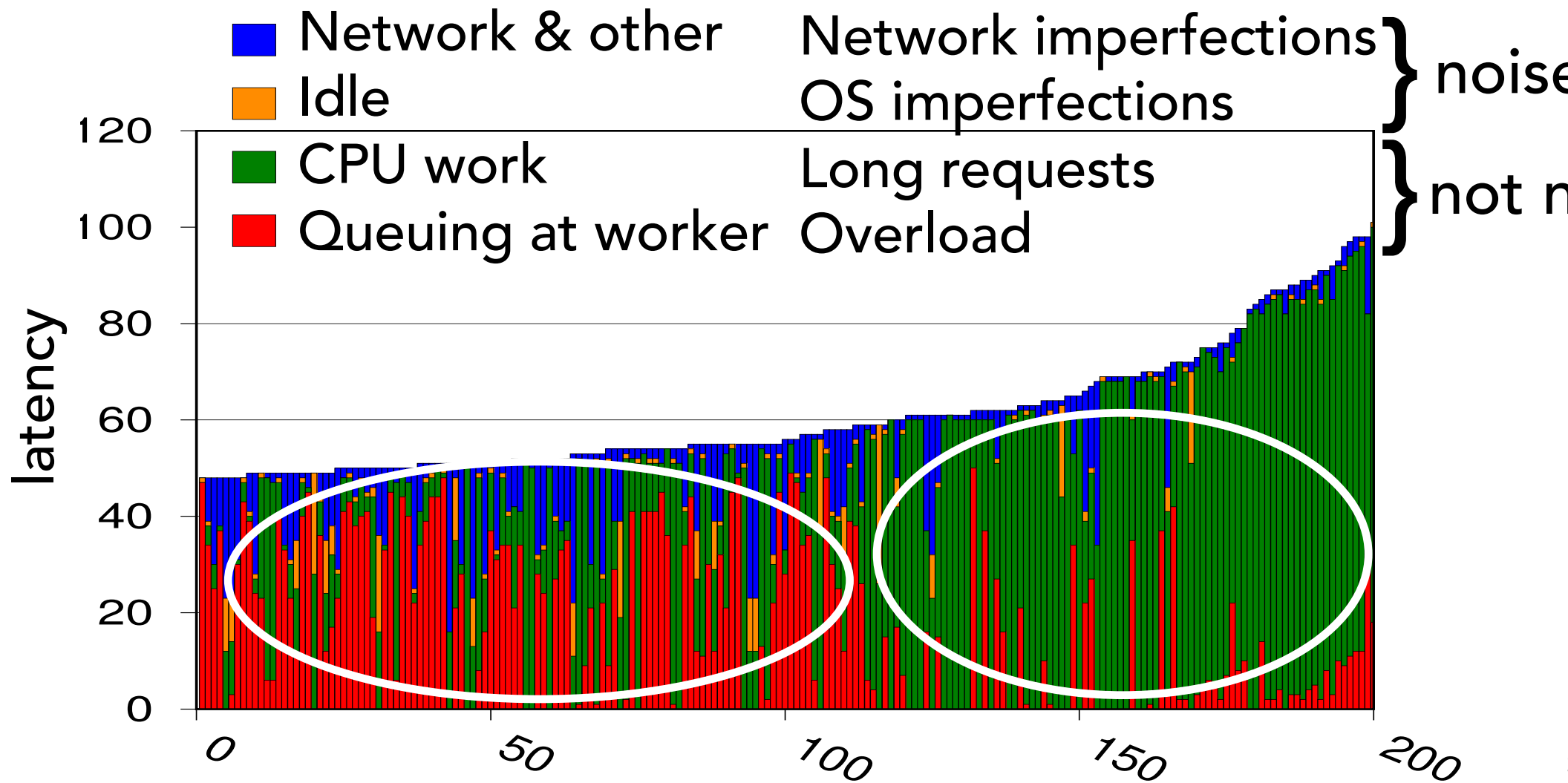- request id – configure
- time stamps, PC

System threads
- thread ids
- time stamp, PC

# All requests



latency (ms)

120
100
80
60
40
20
0

Client latency
Average queueing time

0          20          40          60          80          100

Request groups (from the slowest 1% to the fastest 1%)

# The Tail  Longest 200 requests

# Insight
## Long requests reveal themselves
Regardless of the cause

# Noise  Replicate & reissue

All requests?

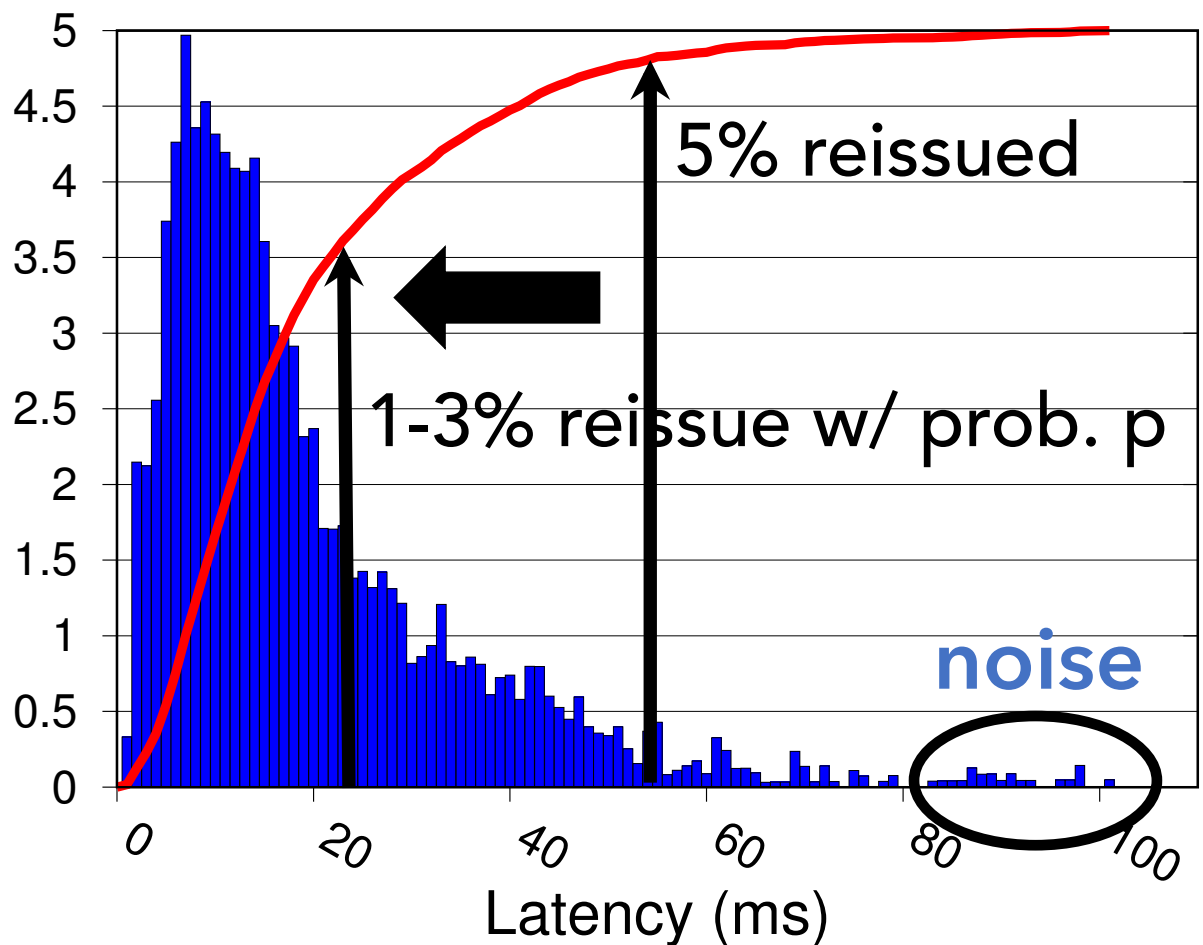5% reissued

10 % reissued

CFD for cost & pote

Fixed issue time

noise

# Probabilistic reissue

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



Adding randomness to reissue makes *one* earlier reissue time d (vs n) opt
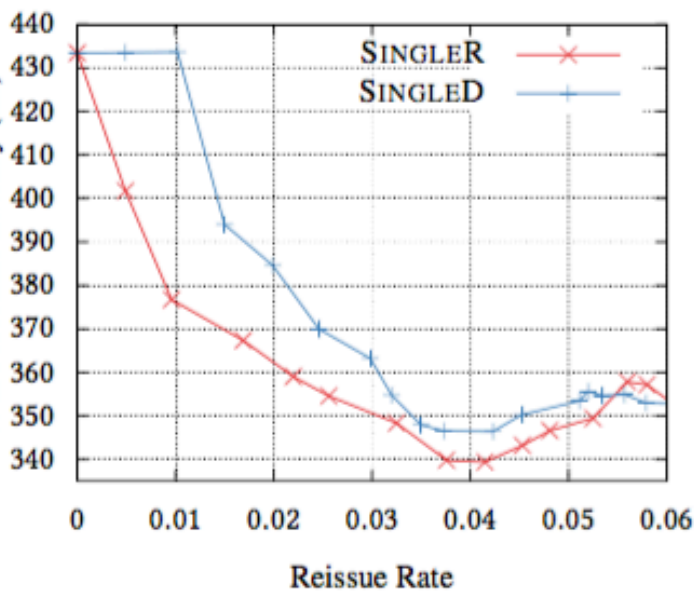
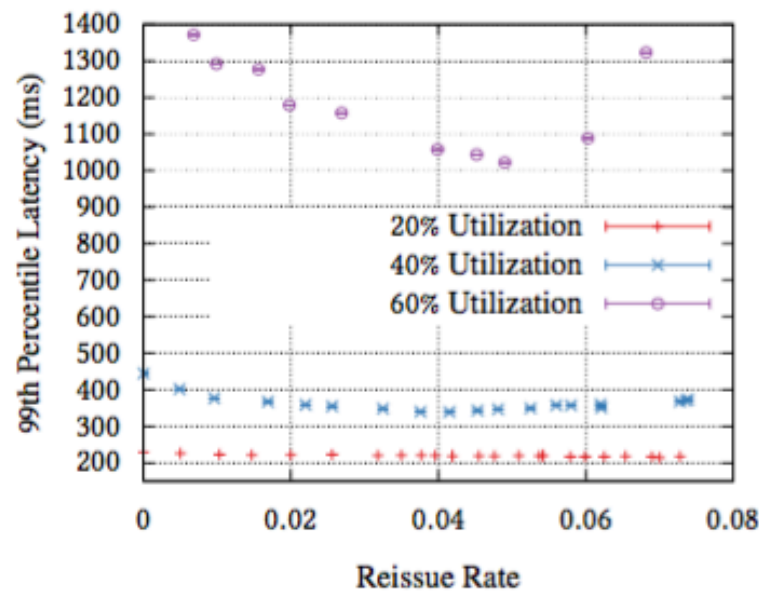Probability is proportion reissue budget & noise
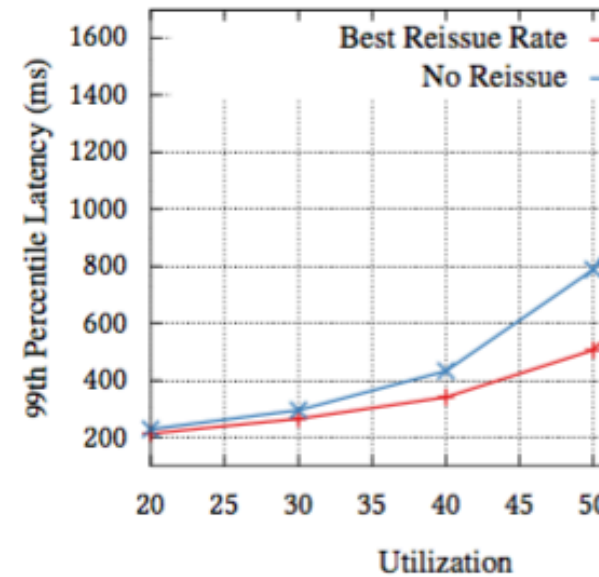
# Single R  Probabilistic reissue

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



(a) SINGLER vs SINGLED
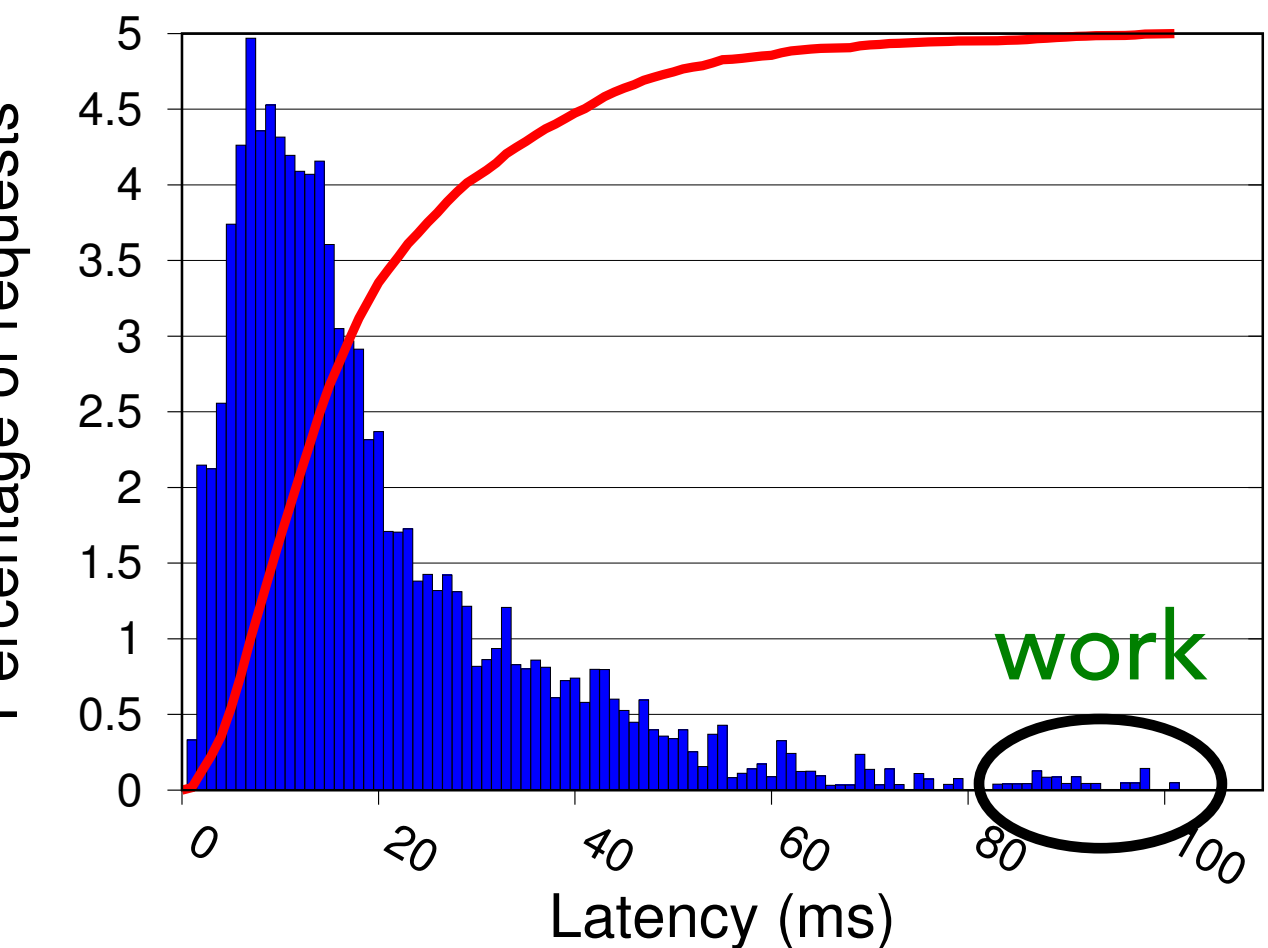
(b) Latency vs Reissue Rate

(c) Best Latency vs Utilization

# Work  Speed up the tail *efficiently*


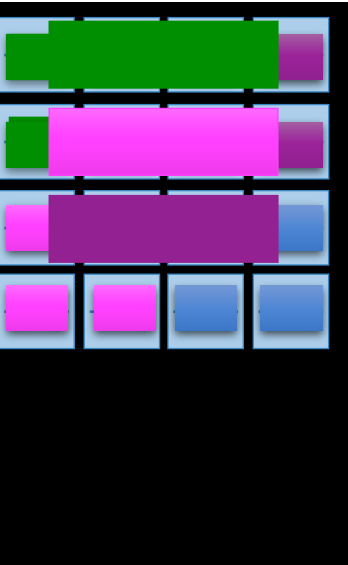
100 **Judicious parallelism**
[ASPLOS'15]

80 **DVFS faster on the t**
[DISC'14, MICRO'17]

60 **Asymmetric multicor**
[DISC'14, MICRO'17]

# Work Parallelism

Parallelism historically for **throughput**

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

# Queuing theory

Optimizing average latency maximizes throughpu

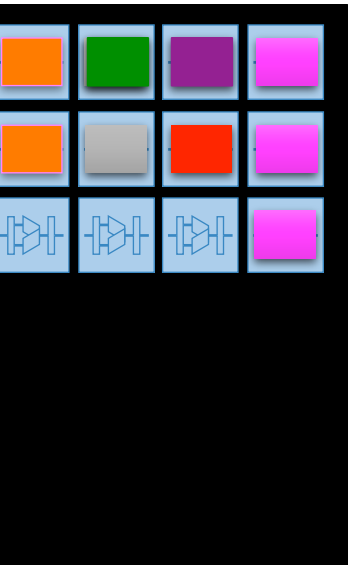But not the tail!

Shortening the tail reduces queuing latency

# Tail-to-Many Dynamic Parallelism [ASPLOS'15]

Parallelism historically for **throughput**

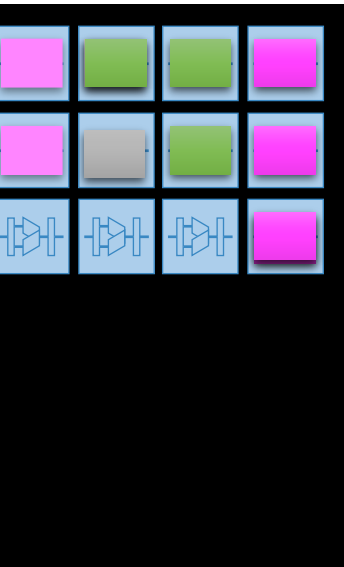**Idea** Parallelism for **tail latency**

**Insight** Long requests reveal themselves

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

# Few to Many at fixed delay d

## Add thread every d *ms*



Long delay good at high load

Short delay good at low load

**Sequential**
**4 way**
Fixed interval 20 ms
Fixed interval 100 ms
Fixed interval 500 ms

Tail latency *ms*

1500

1200

900

600

300

30  32  34  36  38  40  42  44
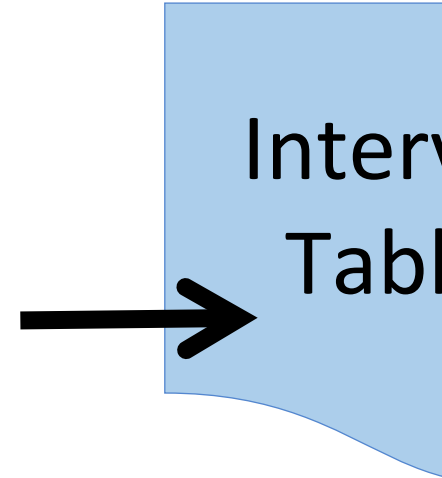
Lucene  RPS

be...
all lo...

# Offline

**Profiles**
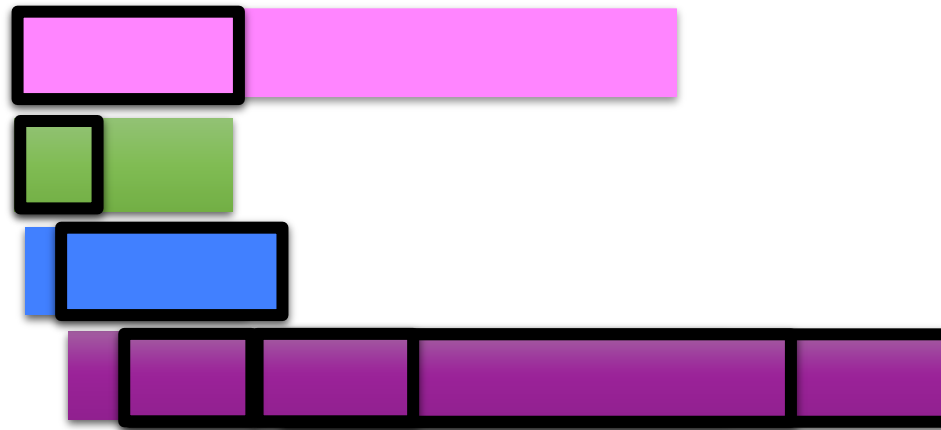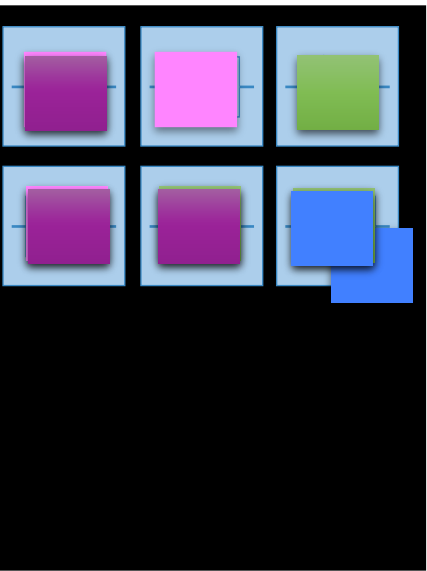Sequential & parallel demand distribution
Efficiency of parallelism

**Choose maximum target parallelism**
Utilize available hardware resources

**Exhaustively** explore parallelism given set of time
intervals $t$ & load find best tail latency & parallelism

Interv
Tabl

# Online self scheduling



| |requests| | Interval$_0$ = 0 | | Interval$_{1,2}$ = 50, 100 |
|---|---|---|---|
| ≤ 2 | @ 0 | parallelism = 3 | |
| 3 | @ 0 | parallelism = 1 | @ 50, parallelism = |
| 4 - 6 | @ 50 | parallelism = 1 | @ 100, parallelism = |
| ≥ 7 | @ exit | parallelism = 1 | @ 100, parallelism = |

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



**Sequential**  **Few to Many**

**21% fewer servers**

**or reduce tail by 28**

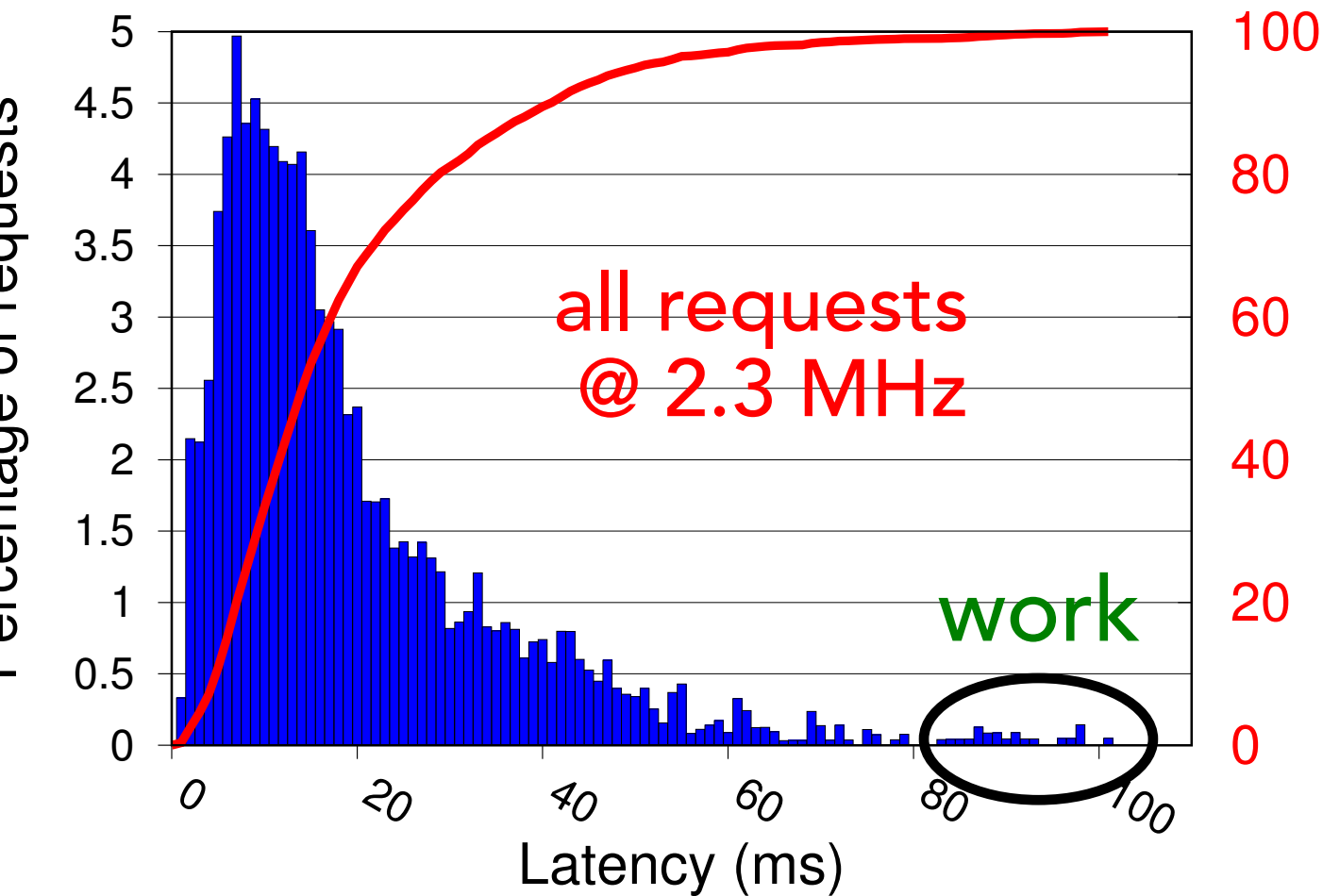**Requests per Second**

# Work speed up the tail *efficiently*
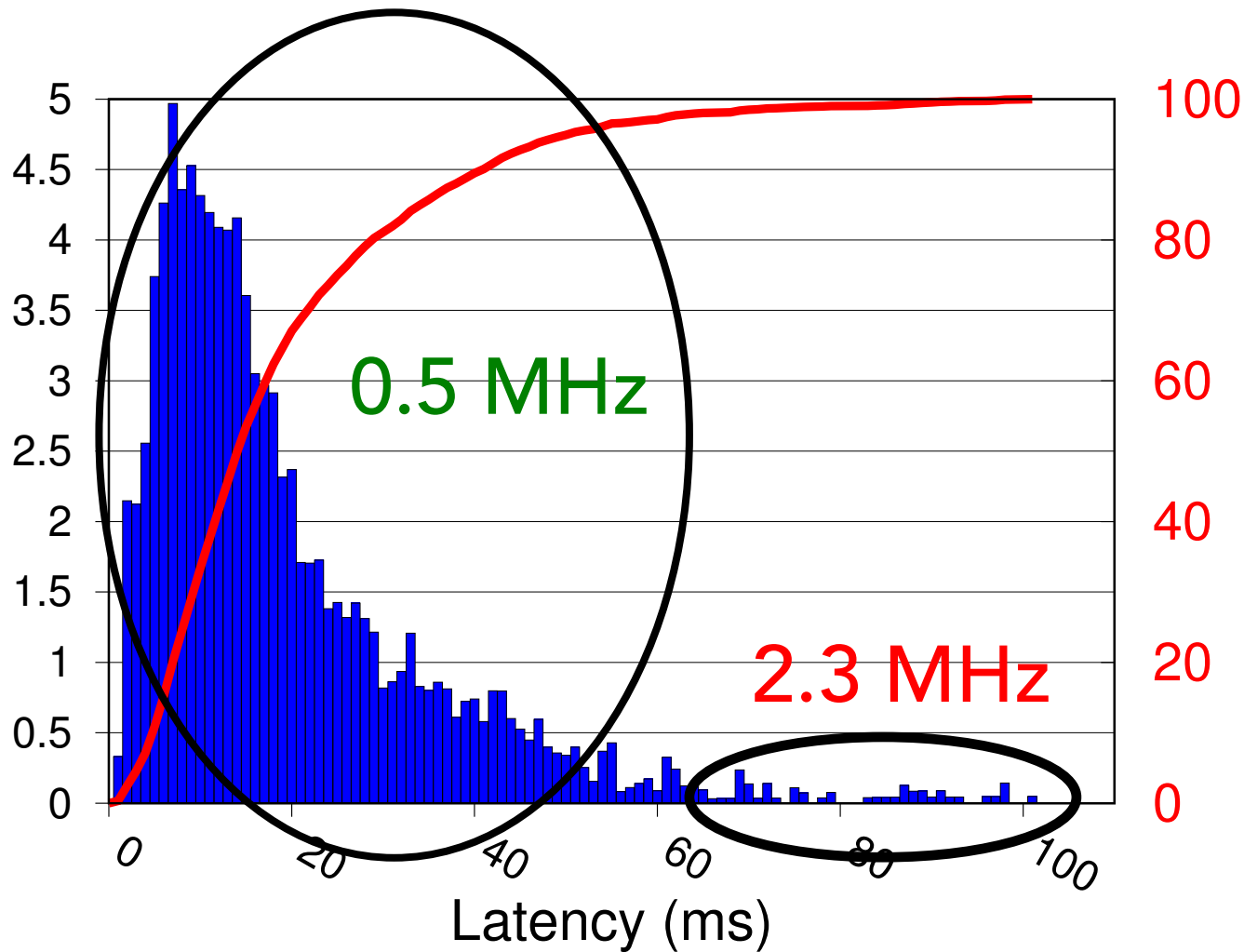


Judicious parallelis
[ASPLOS'15] ✔

work

# *Work* speed up the tail *efficiently*



DVFS faster on the
[DISC'14, MICRO'17]

Asymmetric multico
(AMP) [DISC'14, MICRO'17]

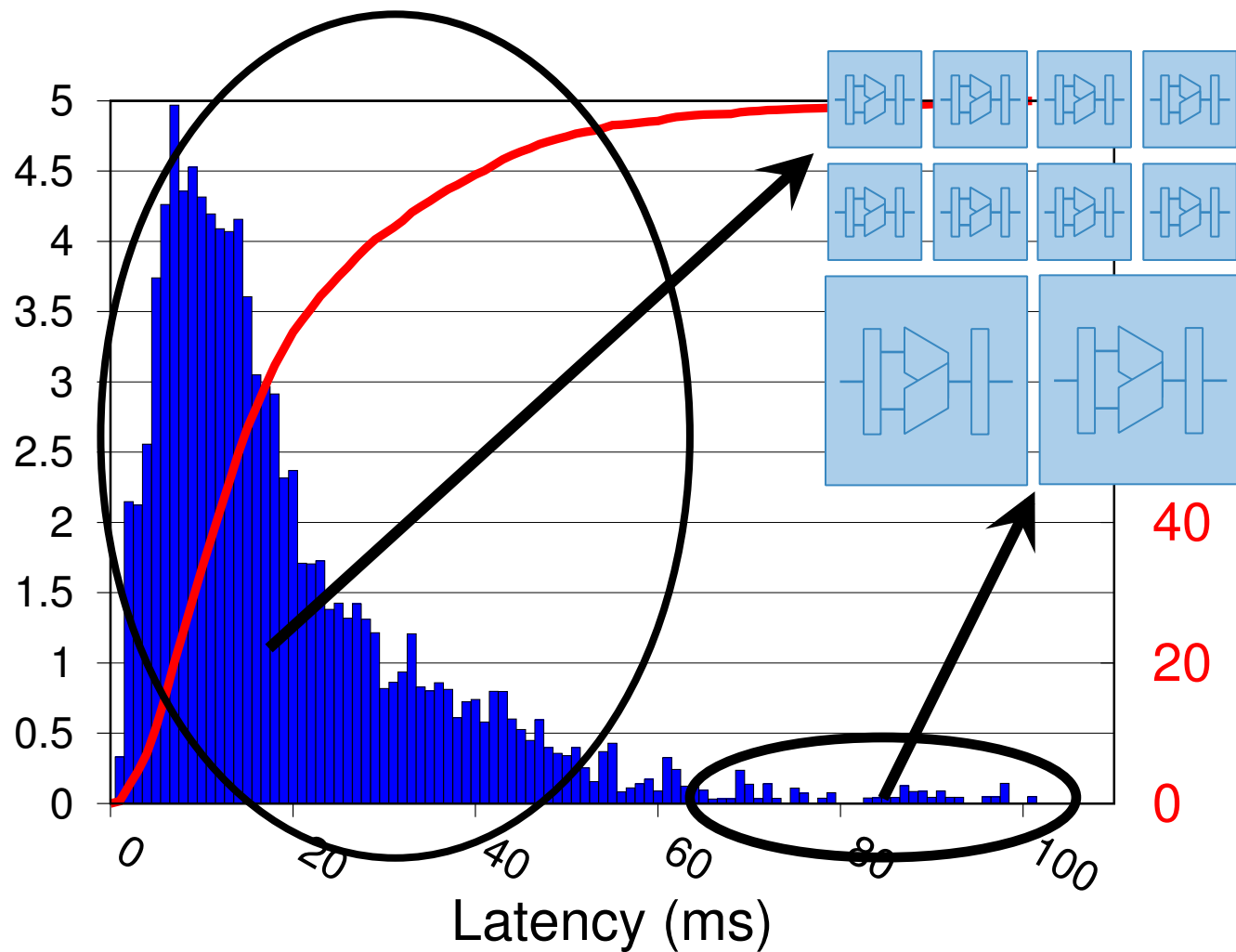# Speed up the tail *efficiently*



DVFS faster on the
[DISC'14, MICRO'17]

+ available in servers t

# Speed up the tail *efficiently*



DVFS faster on the
[MICRO'17]

+ available in servers t

Asymmetric multico
(AMP) [DISC'14, MICRO'17]

+ much more energy e

+ hyper-threading is a

- core competition

# Adaptive Slow to Fast Framework

Slow to fast migration is optimal [ICAC'13]
## Goal
Minimize energy consumption and satisfy a tail latency target
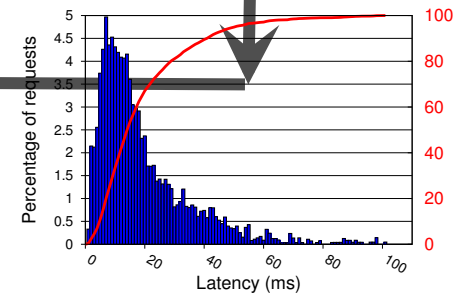
## Challenges
When to migrate?
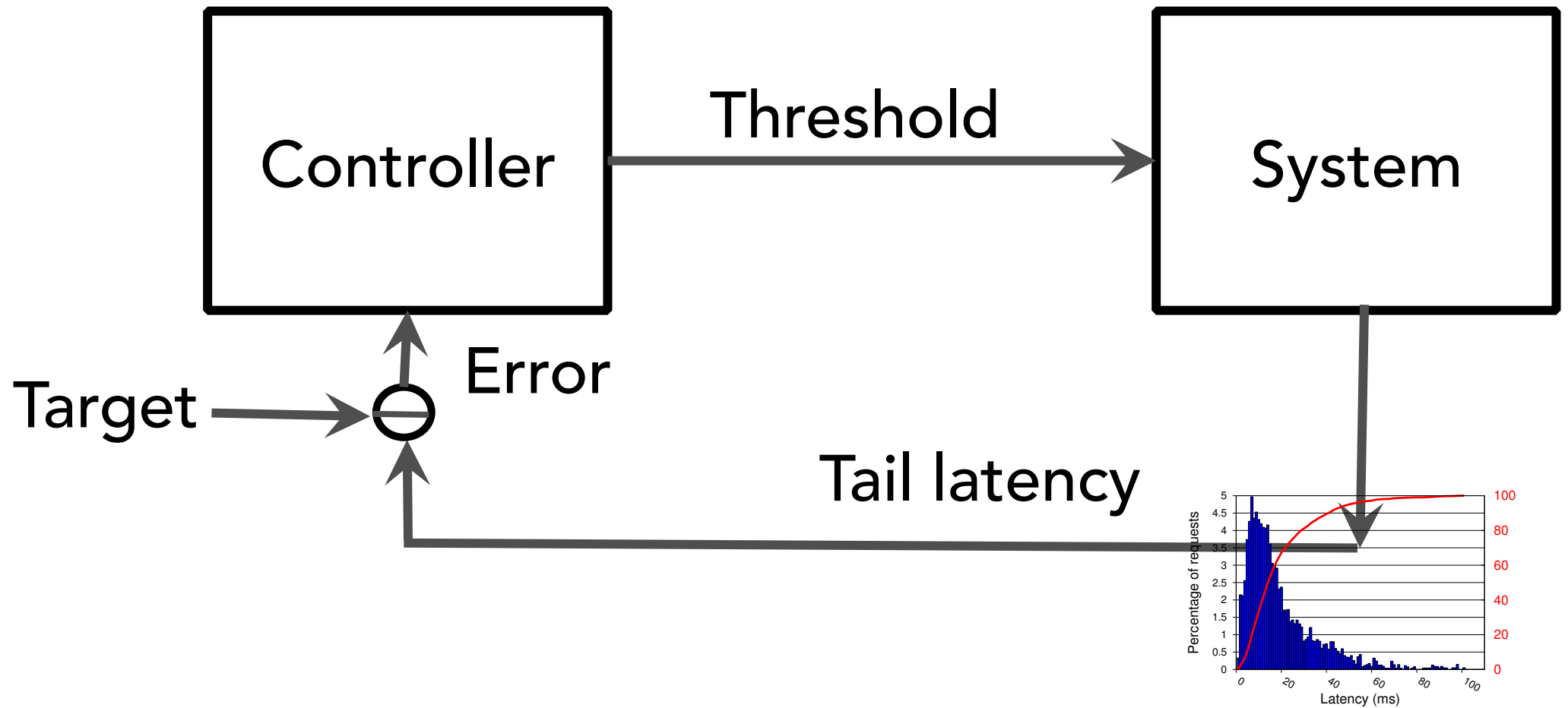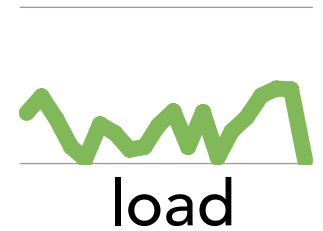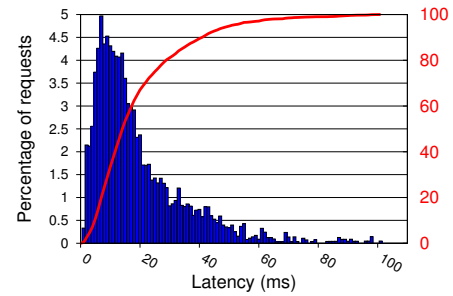What if the core speed is not available?
## Insight
Use **big core** just enough  $th+(l_{99}-th)/sp \leq target$
Migrate oldest first and migrate early under load!

# Controller design



load

Controller — Threshold → System

System — Tail latency → 

Target ⊖ → Controller

Error

# Policies

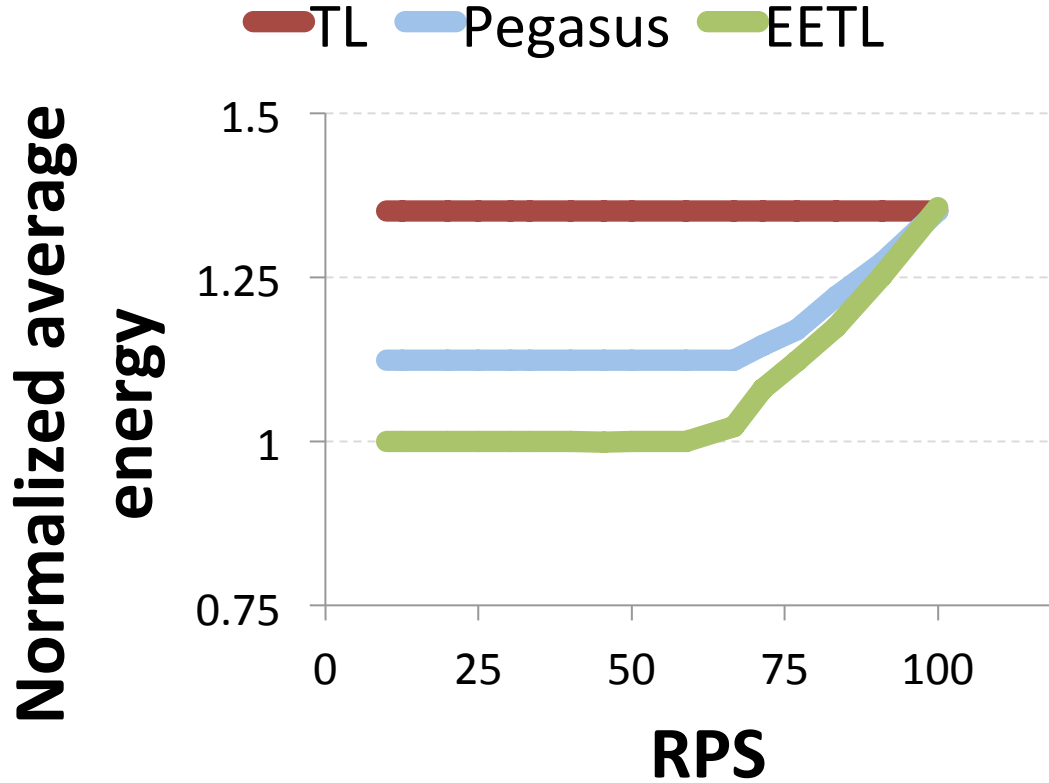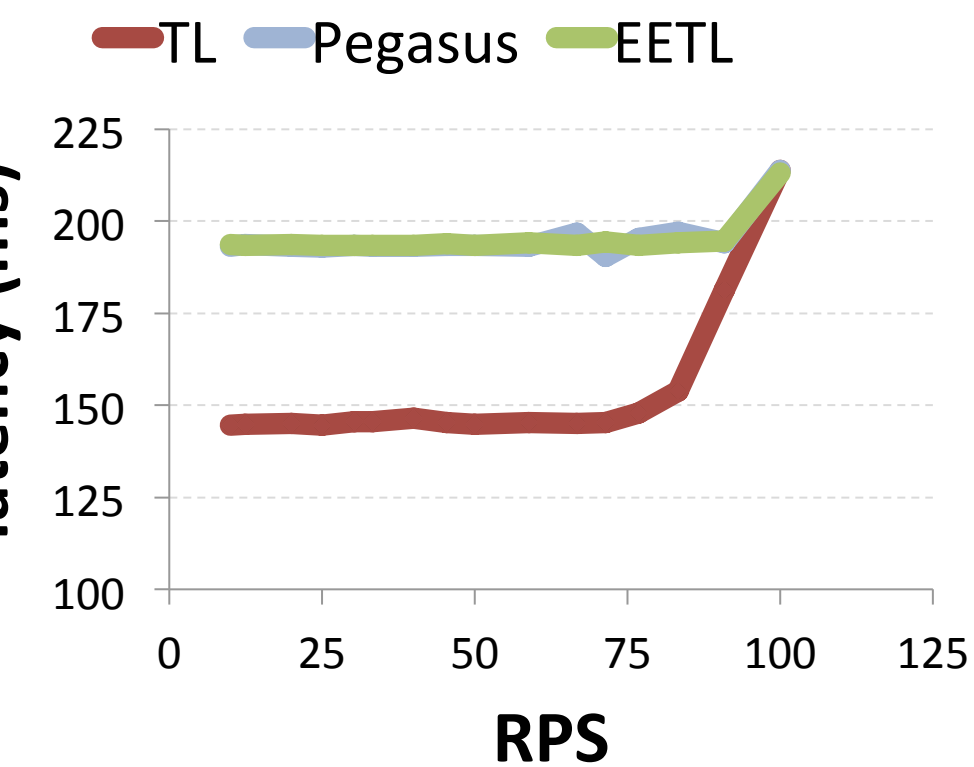## All-cores

**Pegasus** adjust all-core frequencies for load

[Towards Energy Proportional…, Google & Stanford, ISCA'14]
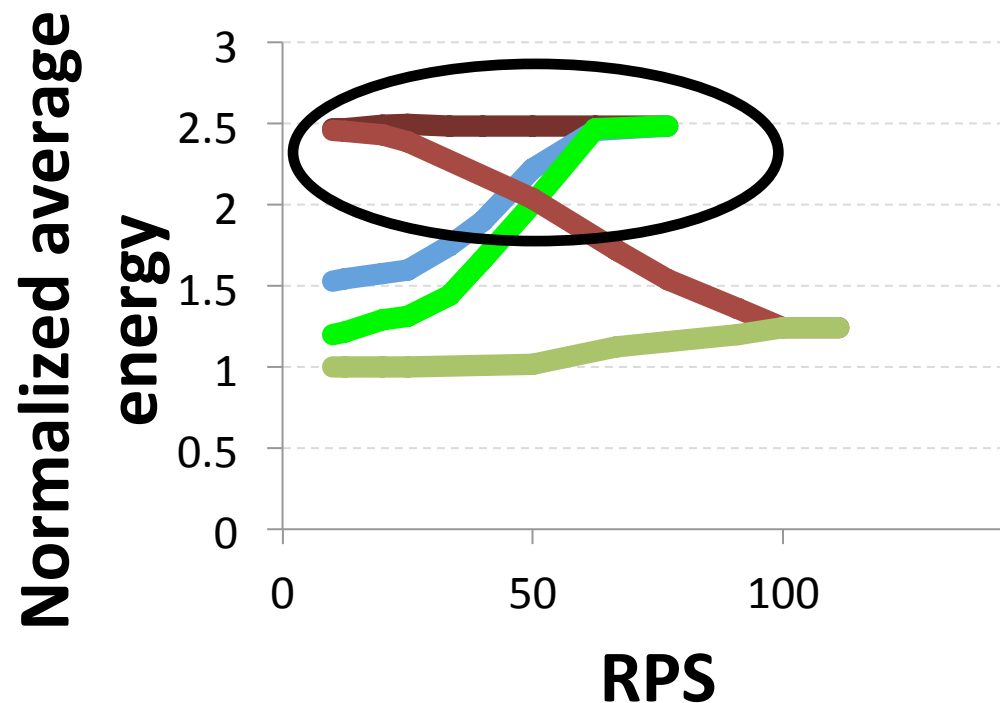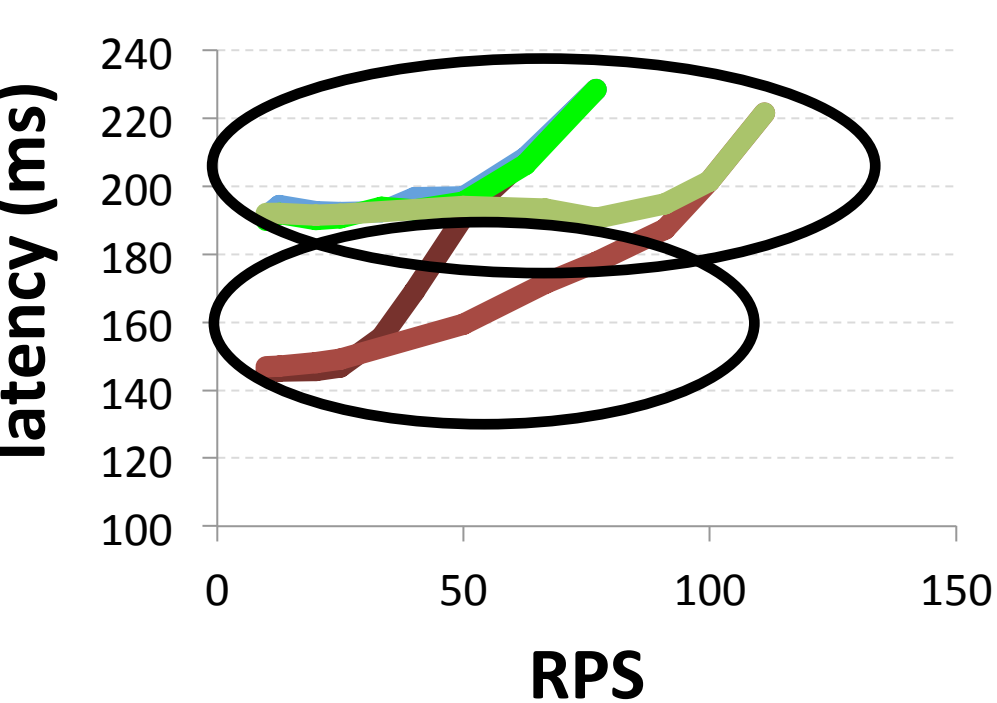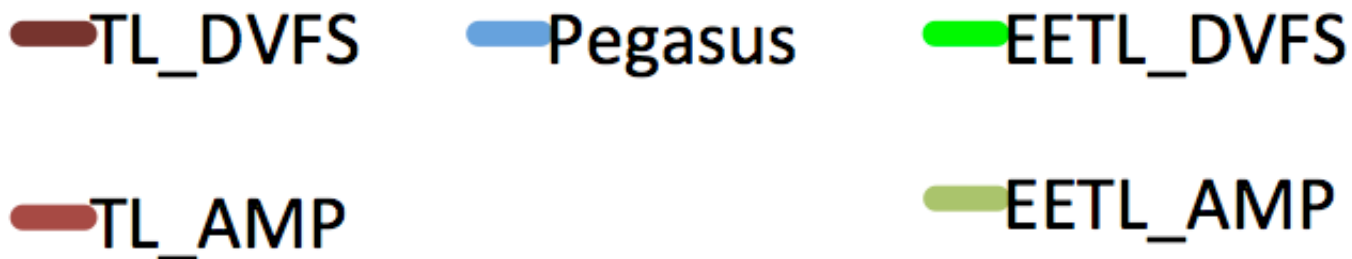
## Per-core approaches

**TL** minimize tail latency using 0 threshold

**EETL** energy efficient with target tail latency

# Lucene with DVFS on Broadwell

# Lucene on emulated AMP and DVFS

il Latency

BOTH !

Efficiency

Google Clou

# Efficiency at scale for interactive workloads

Diagnosing the tail with continuous profiling

      <span style="color:blue">Noise</span>      replicate, systems are not perfect

      <span style="color:red">Queuing</span>    not today!

      <span style="color:green">Work</span>      judicious use of resources on long requests

Request latency CDF is a powerful tool

Tail efficiency ≠ average or throughput

Hardware heterogeneity

<span style="color:blue">Th</span><span style="color:green">a</span><span style="color:red">n</span><span style="color:blue">k</span> <span style="color:green">y</span><span style="color:red">o</span><span style="color:blue">u</span>

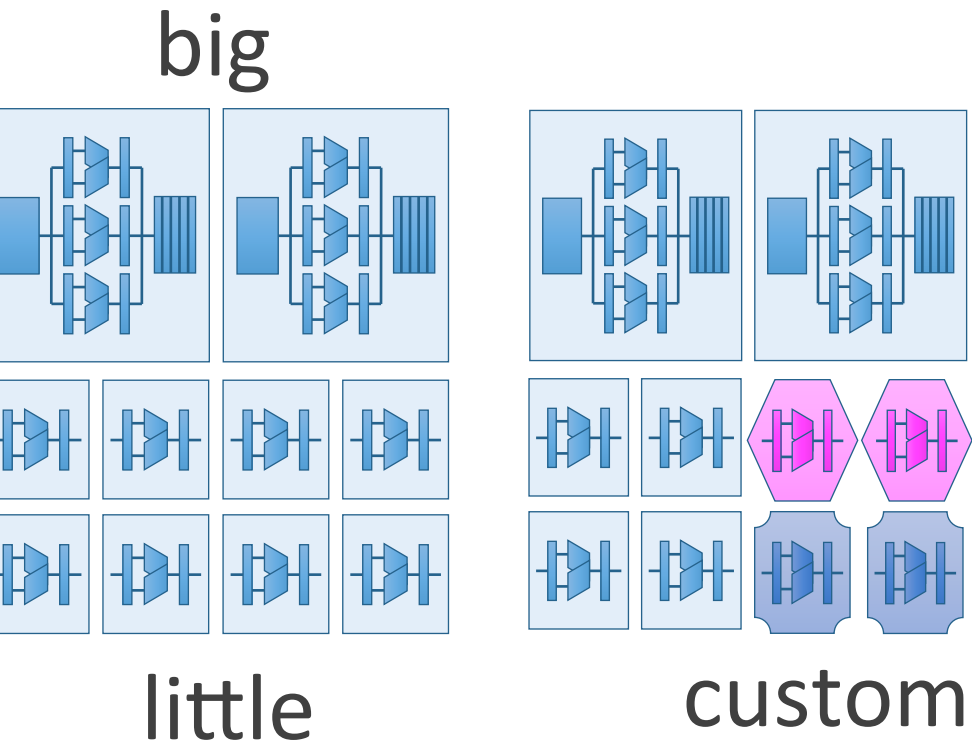# Requirements pull for heterogeneity!

[DISC'14, ICAC'13, submission]

**Heterogeneous hardware** dominates homogeneous hardware for throughput, performance, and energy with a fixed power budget & variable request demand

**Slow-to-Fast** sacrifice average a bit to reduce energy & tail latency
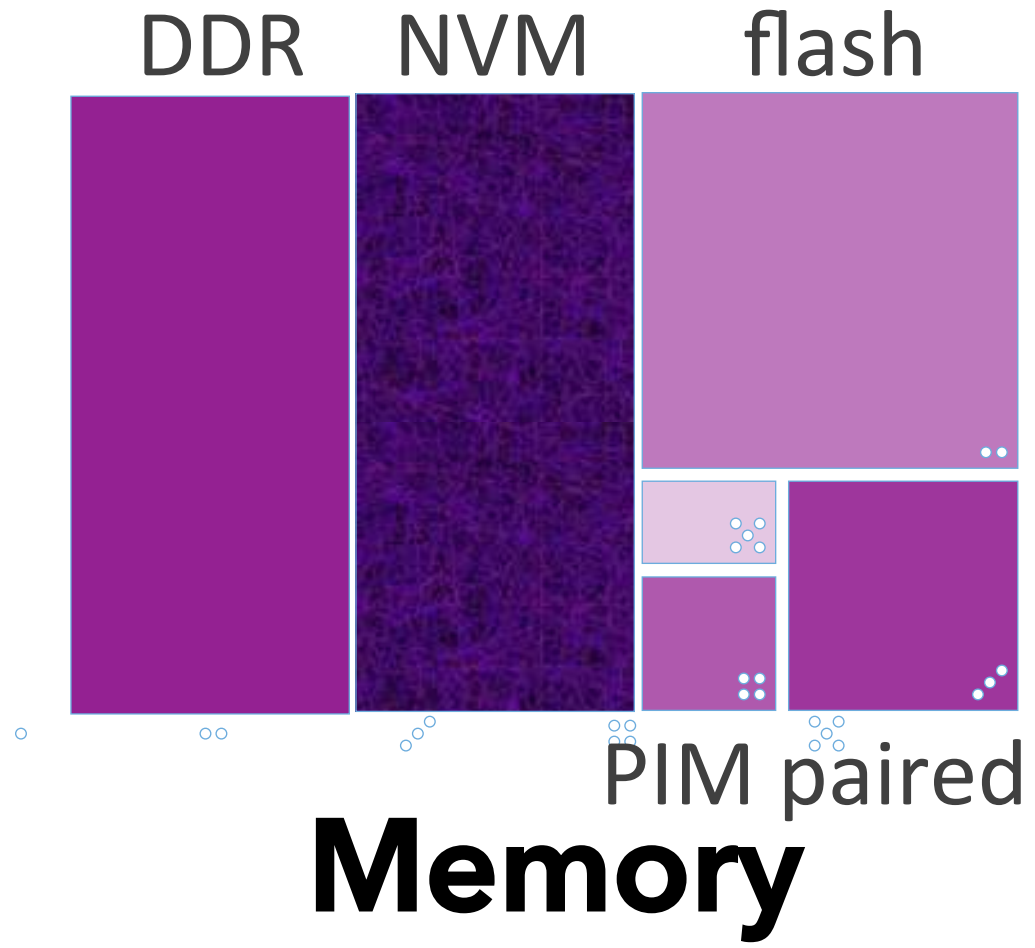
# Hardware heterogeneity – opportunity & challen...
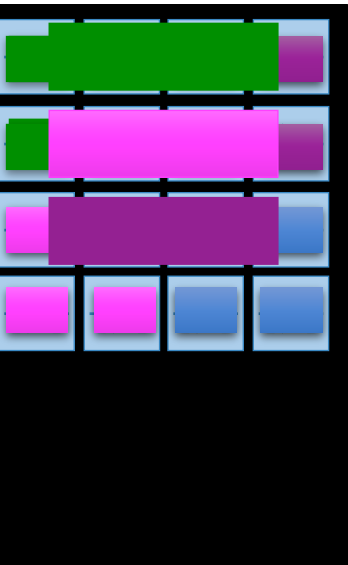


big

little

custom

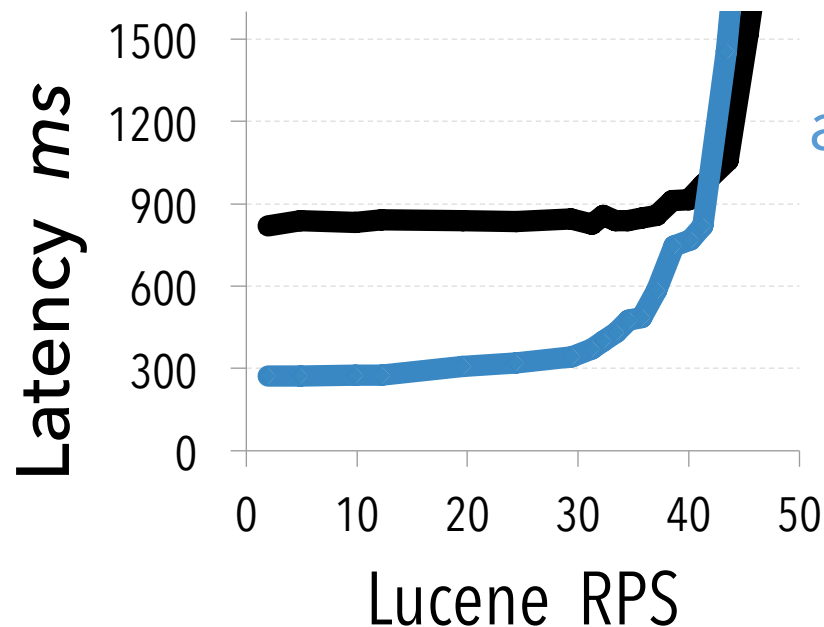**Processors**

DDR   NVM   flash

PIM paired

**Memory**

# Parallelism

Parallelism historically for **throughput**

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

# Software & hardware

**Lucene** open source enterprise search Wikipedia English
- 10 GB index of 33 million pages
- 10k queries from Lucene nightly tests

**Bing** web search with one Index Serving Node (ISN)
- 160 GB web index in SSD, 17 GB cache
- 30k Bing user queries

**Hardware** 2x8 64 bit 2.3 GHz Xeon, 64 GB Windows
- 15 request servers, 1 core issues requests
- Target parallelism = 24 threads

# Policies

**Sequential**

**N way**  single degree of parallelism for each request

**Adaptive**  Select parallelism degree when request starts using system load     [EUROSYS'13]

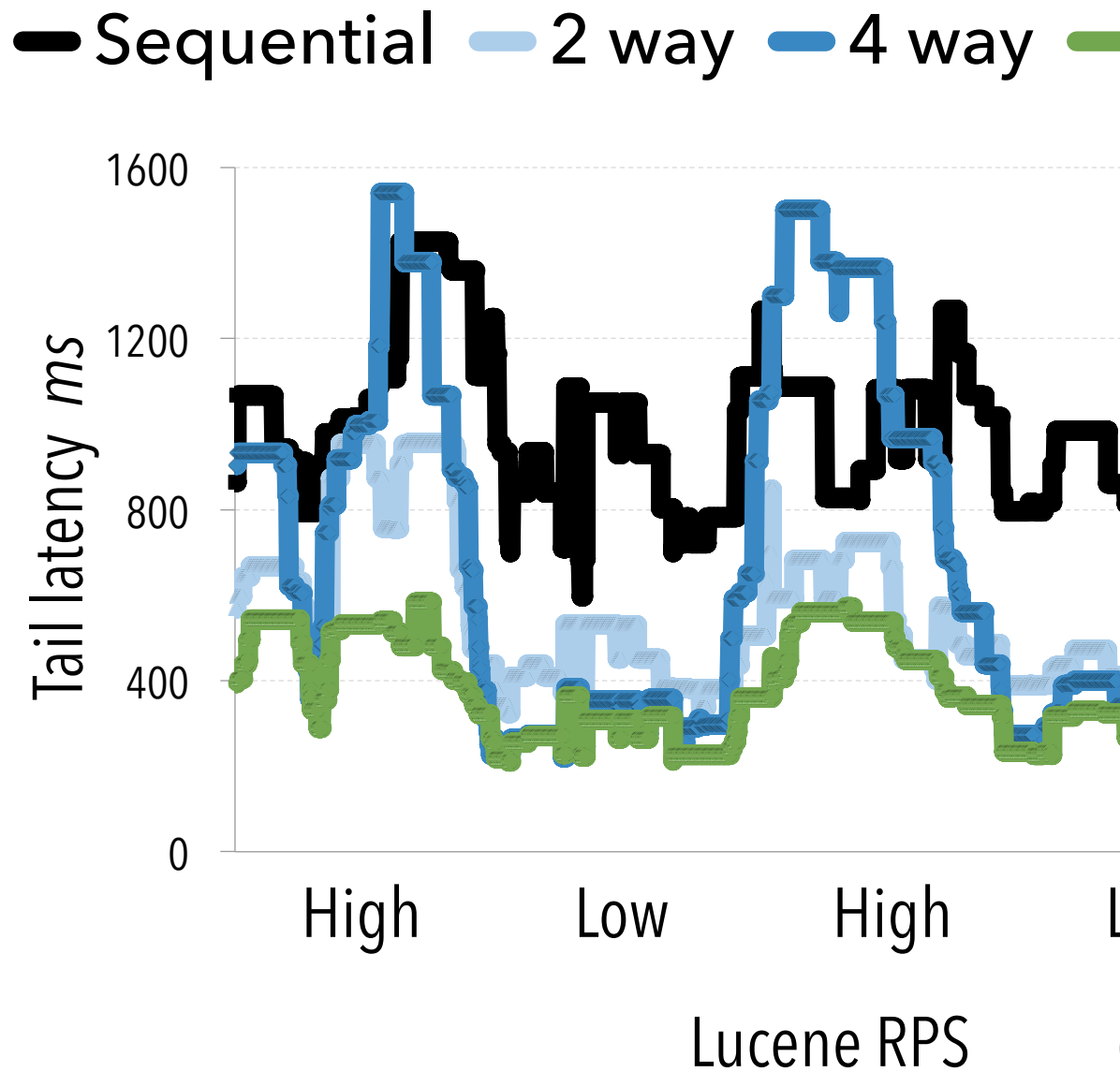**Request Clairvoyant**  parallelizes long requests by perfect prediction of tail

**FM**  Few to Many incrementally add parallelism

# Load variation

Alternate between high & low load

FM adapts to bursts with low variance

**Sequential** ▬ **2 way** ▬ **4 way** ▬



Tail latency *ms*

1600
1200
800
400
0

High    Low    High    L

Lucene RPS

# Fewer servers: Total Cost of ownership

Sequential ── FM ── Adaptive ── FM

**21%**

**9%**

Lucene RPS

Lucene RPS