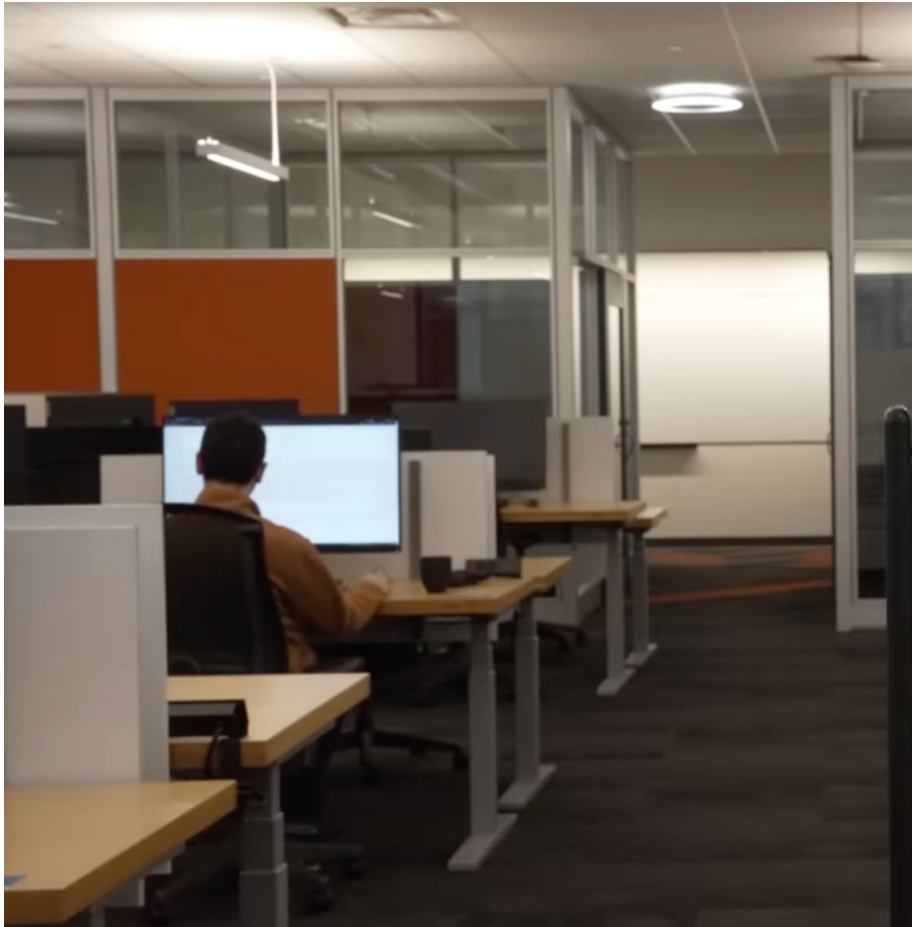


# Better latency with fewer servers

avoiding RAM noisy neighbors



**SCALE**



Site Reliability Team 😊

if (time == now()) get(Coffee());

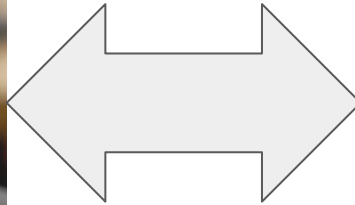
— I ♥ JAVA

DAYS SINCE  
LAST OUTAGE:

~~07~~ ~~42~~ ~~17~~ ~~6~~ 3

Q1 HEADCOUNT  
GOAL:

-4



## Why SRE?

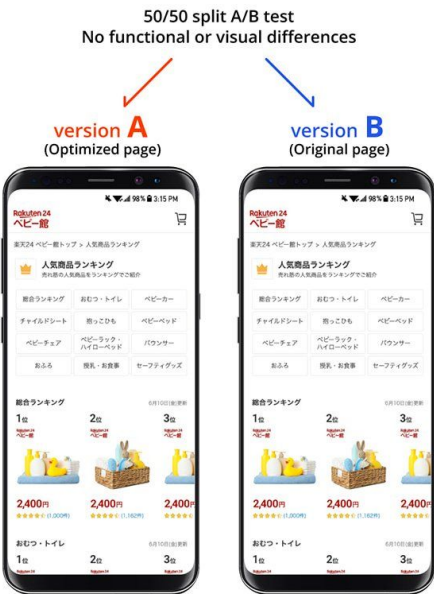


A scenic landscape featuring a calm lake in the foreground, reflecting the golden light of a sunset. The background is dominated by jagged, snow-capped mountains under a sky with soft, wispy clouds. The scene is framed by evergreen trees on the left and right sides, and some rocks are visible in the water and along the shore.

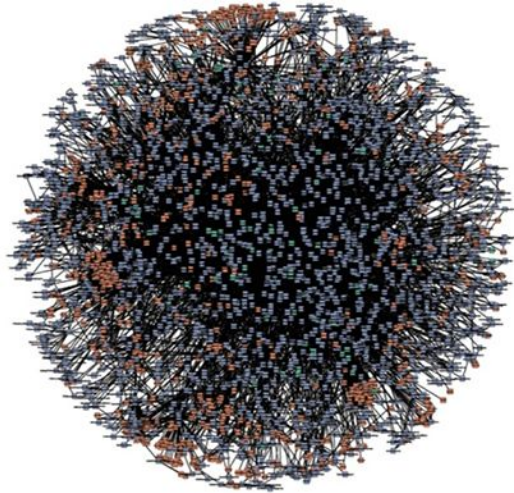
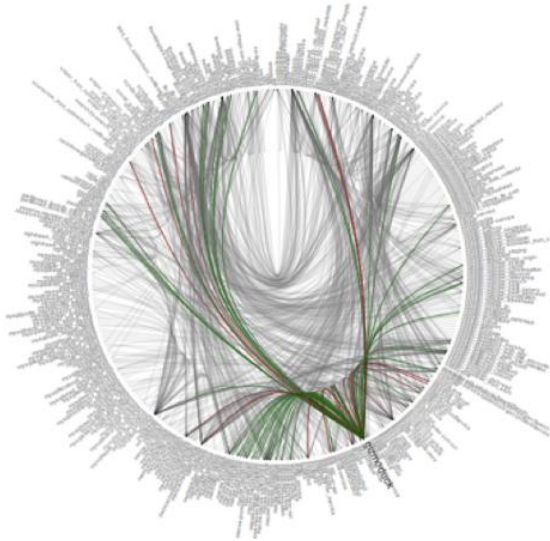
Your deployment is highly available and cost efficient  
Users delight in the product and its performance



# Performance matters

Company	Year	Impact	Source
Amazon	2019	<p>50/50 split A/B test No functional or visual differences</p> 	<a href="#">source</a>
Google	2019		<a href="#">source</a>
Akamai	2019		<a href="#">source</a>
Ebay	2019		<a href="#">source</a>
Vodafone	2019		<a href="#">source</a>
Yelp	2019		<a href="#">source</a>
Adobe Experience Cloud	2019		<a href="#">source</a>
Pfizer	2019		<a href="#">source</a>
Rakuten 24	2022		-0.4s latency → +53% revenue per visitor, +33% conversion, +15% average order value

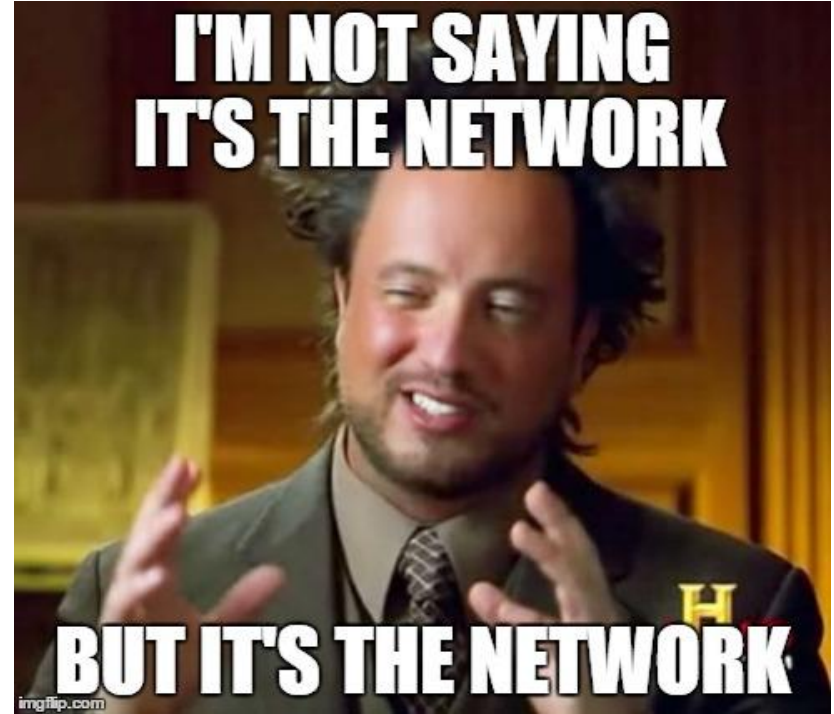
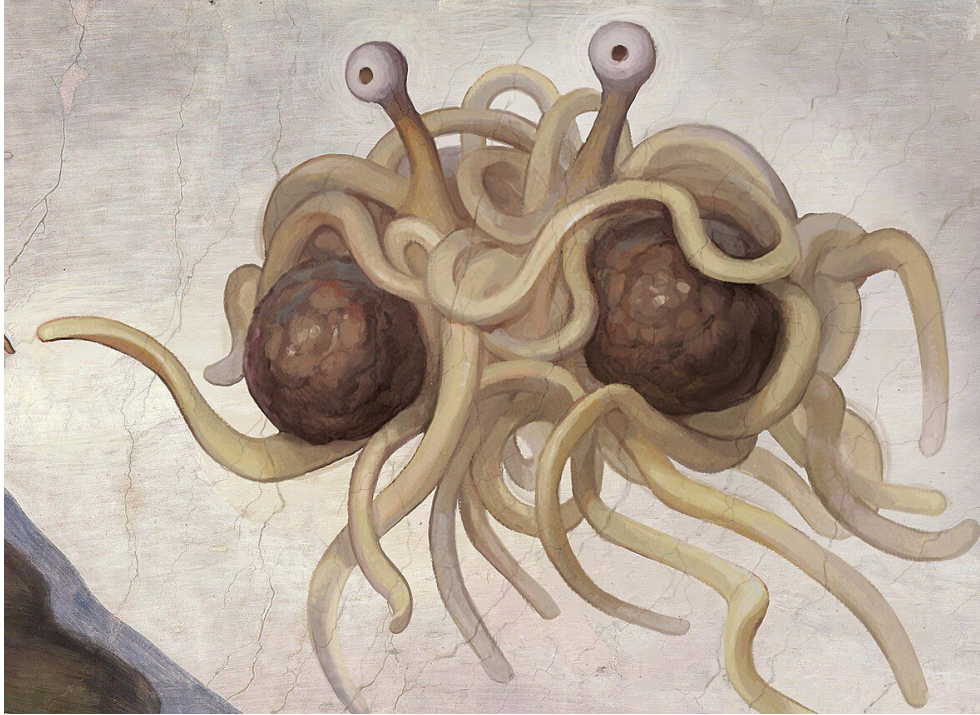
# Users want more functionality



amazon.com



Easy to blame “complexity”





Blame resource congestion!





Run 20%-50%  
more transactions

Reduce tail latency  
by 20-80%



2011

2013

2015

2016

2018

2019

2020

2021

2023



BubbleUp



Heracles



PerfIso



Caladan



Themis



iBench



Dirigent



PARTIES

Libra



Can be made widely useful

CPI<sup>2</sup>



CAT at Scale



Alita





# Talk Outline

The  
Problem

Why  
Solve

Available  
Tooling

Mitigation  
Systems

Current  
Efforts



Hi, I'm Jonathan Perry

[yonch@yonch.com](mailto:yonch@yonch.com)

@yonchco

- Ph.D.: *Network* Noisy Neighbor mitigation
- CEO @ Flowmill: Network observability
- CEO @ Unvariance



# The Problem

- What is memory noisy neighbor
- How does this affect my Pods?
- Do I have it in my cluster?

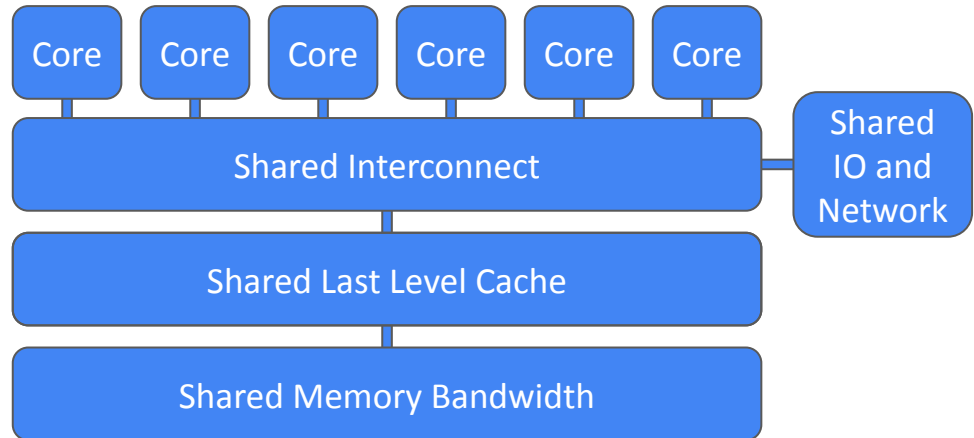


# Noisy neighbor: tragedy of the commons

Apps access physical resources



Shared resources are constrained



# Noisy neighbor: tragedy of the commons

Apps access physical resources

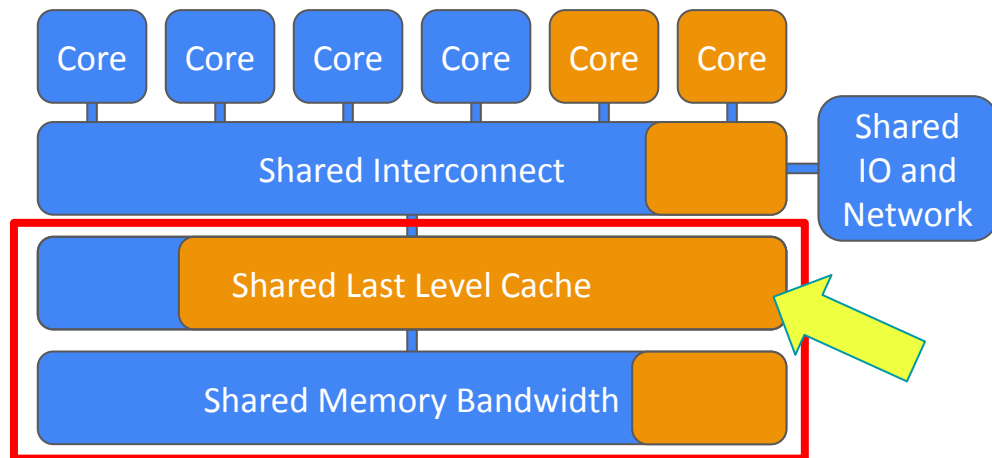


Shared resources are constrained

One App can use more than its fair share, degrading others

This talk:

- Cache
- Memory bandwidth

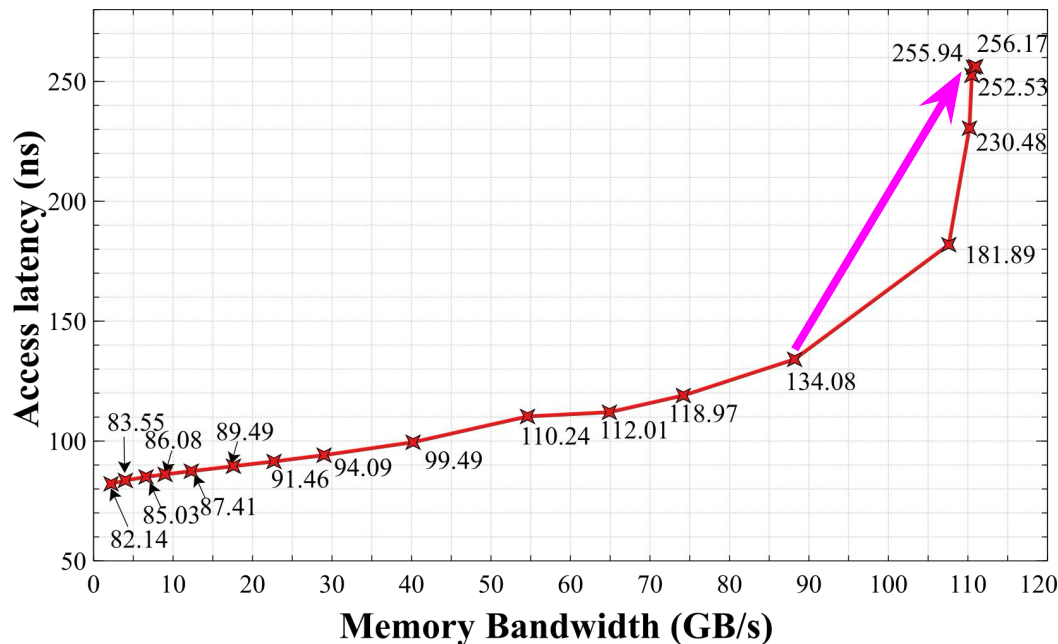


# Latency increases with memory bandwidth

Change memory bandwidth,  
Measure latency

Knee-point around 80%

80% → 100% bandwidth  
latency doubles!

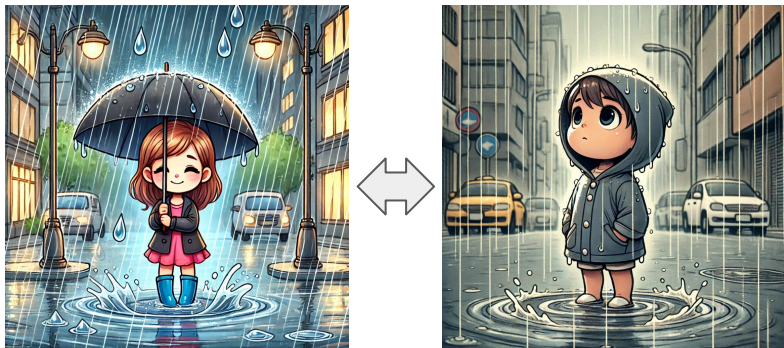


# Does memory latency cause slowdown?

CPUs try to hide memory latency:

- Prefetchers
- Reorder buffer (ROB)
- Caches

**Are they effective?**



We need to **measure slowdown**

Popular metric:

**Cycles Per Instruction – CPI**

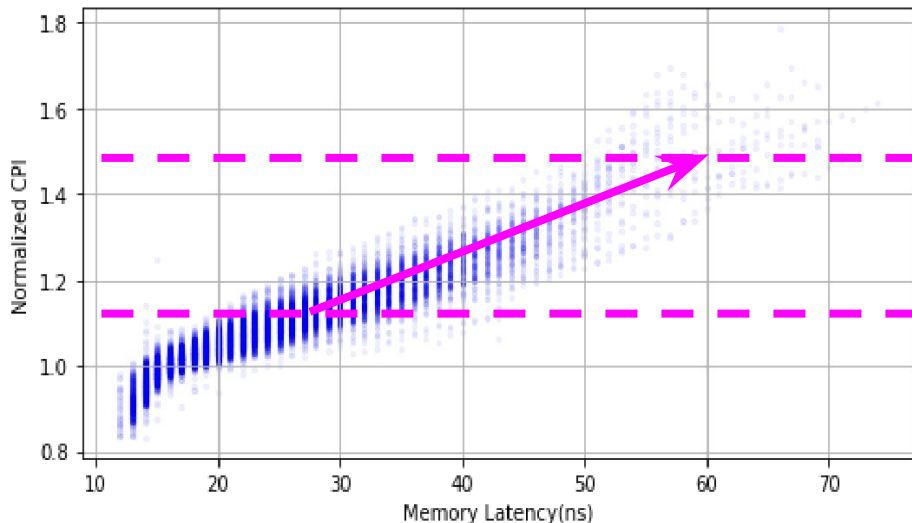
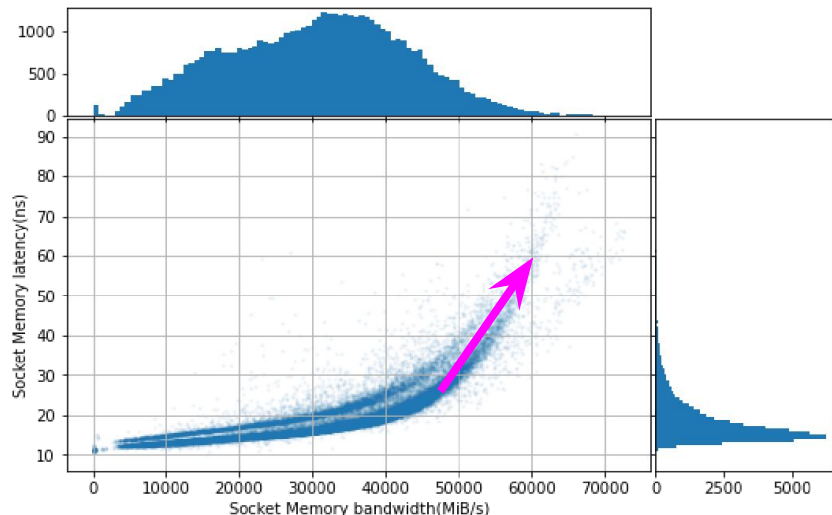
CPU waits for memory → Stall cycles

Many stalls → High CPI

Few stalls → Low CPI

# 80% bandwidth cap → 25% more compute-efficient

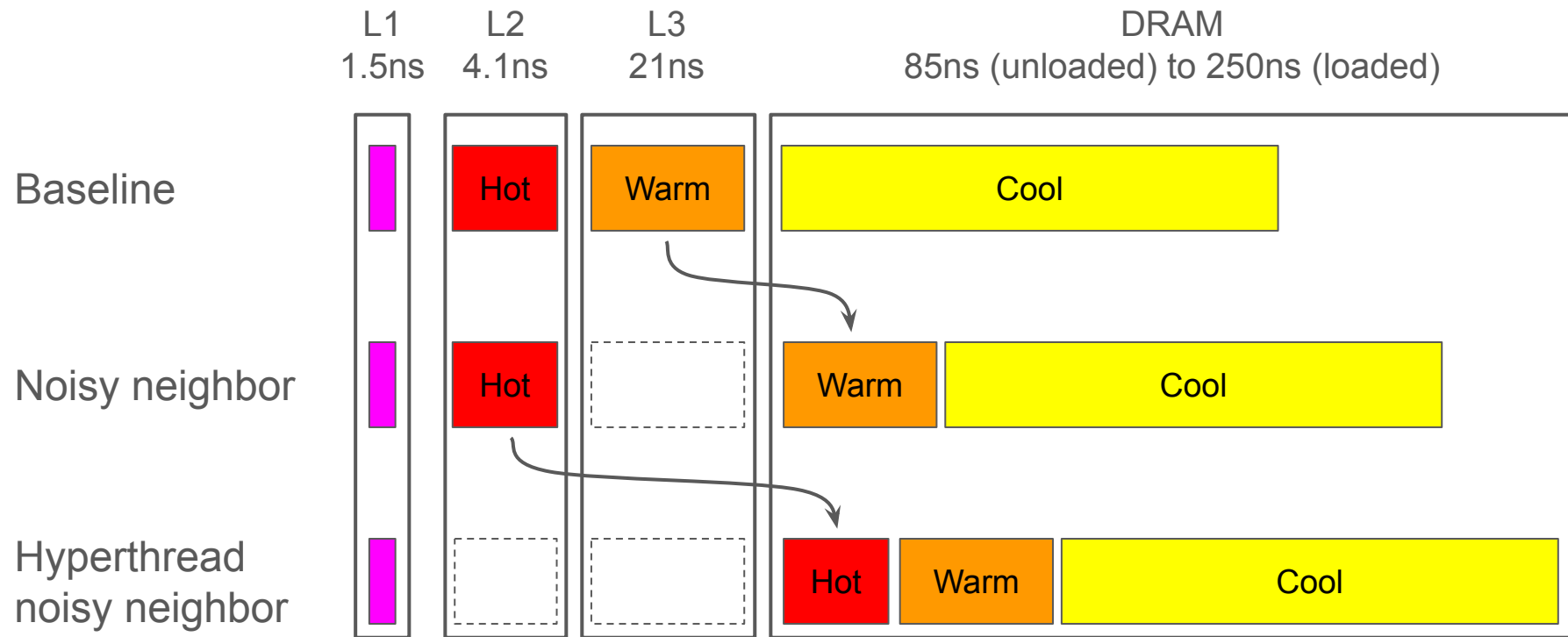
Alibaba Cloud production trace: 8k+ hosts, 1M+ containers, 24 hours



Memory bandwidth only – has cache mitigation



# Cache contention degrades access times



# Tail latency can explode with noisy neighbors!

websearch	10%	20%	30%	40%	50%	60%	70%	80%	90%
baseline(approx)	52%	57%	61%	60%	63%	62%	66%	77%	88%
LLC(small)	103%	96%	102%	96%	104%	100%	100%	103%	103%
LLC(med)	106%	99%	111%	103%	116%	110%	110%	125%	111%
LLC(big)	>300%	>300%	>300%	>300%	>300%	>300%	>300%	264%	123%
DRAM	>300%	>300%	>300%	>300%	>300%	>300%	>300%	270%	122%

> 4.75x

ml_cluster	10%	20%	30%	40%	50%	60%	70%	80%	90%
baseline(approx)	51%	58%	57%	59%	56%	58%	58%	60%	68%
LLC(small)	88%	84%	110%	93%	216%	106%	105%	206%	202%
LLC(med)	88%	91%	115%	104%	>300%	212%	220%	212%	205%
LLC(big)	>300%	>300%	>300%	>300%	>300%	>300%	>300%	250%	214%
DRAM	>300%	>300%	>300%	>300%	>300%	>300%	>300%	287%	223%

> 5.35x

memkeyval	10%	20%	30%	40%	50%	60%	70%	80%	90%
baseline(approx)	20%	20%	21%	21%	22%	29%	35%	42%	34%
LLC(small)	88%	91%	101%	91%	101%	138%	140%	150%	78%
LLC(med)	148%	107%	119%	108%	138%	230%	181%	162%	100%
LLC(big)	>300%	>300%	>300%	>300%	>300%	>300%	280%	222%	79%
DRAM	>300%	>300%	>300%	>300%	>300%	>300%	>300%	234%	103%

> 13.6x

Google, 2015

3 production services

get % of SLO target  
(99th / 95th percentile)

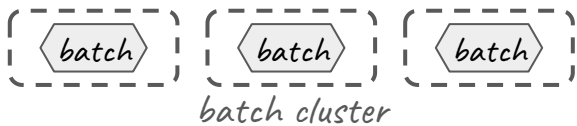
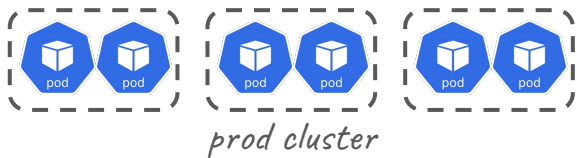
w/ synthetic noise  
generators

# Survey: Node size

Please raise hand if:

- Know what VMs or bare-metal are used in prod
- Never use fraction of physical CPU

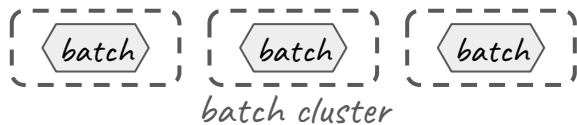
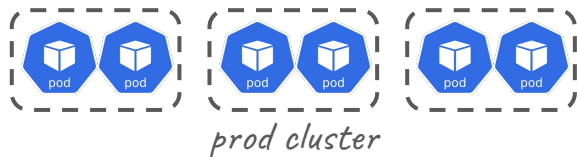
# Separating batch clusters



*What I think I'm running*

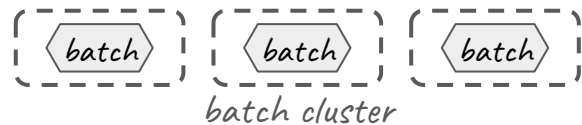
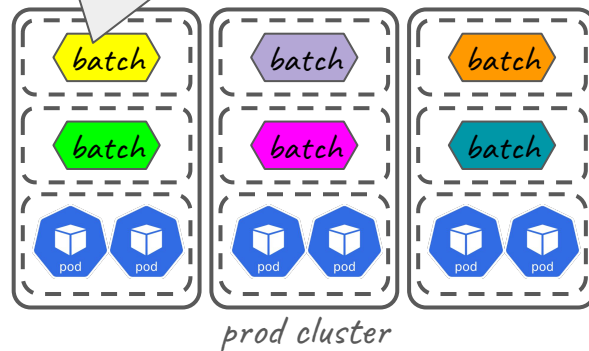


# Separating batch clusters



*What I think I'm running*

*jobs from some random tenants*

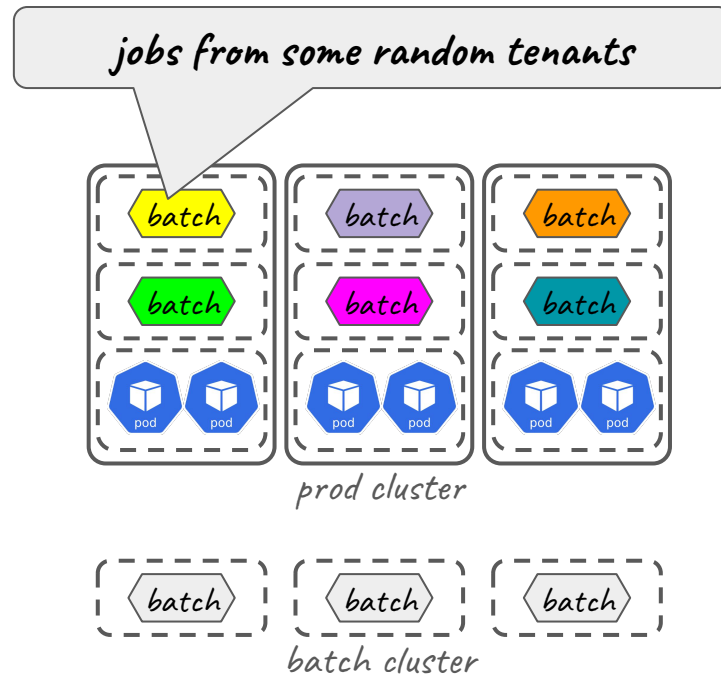


*What I'm actually running*

# Separating batch clusters



Does your provider protect your VMs from other VMs on the same machine?



*What I'm actually running*



Engineers spend years  
optimizing user experience







# Bare-metal: no cross-tenant noisy neighbor



A few 100-core machines



Lots of 8-core VMs

\* Need enough servers to handle server failures (“~10 servers”)

# Does my cluster have noisy neighbors?

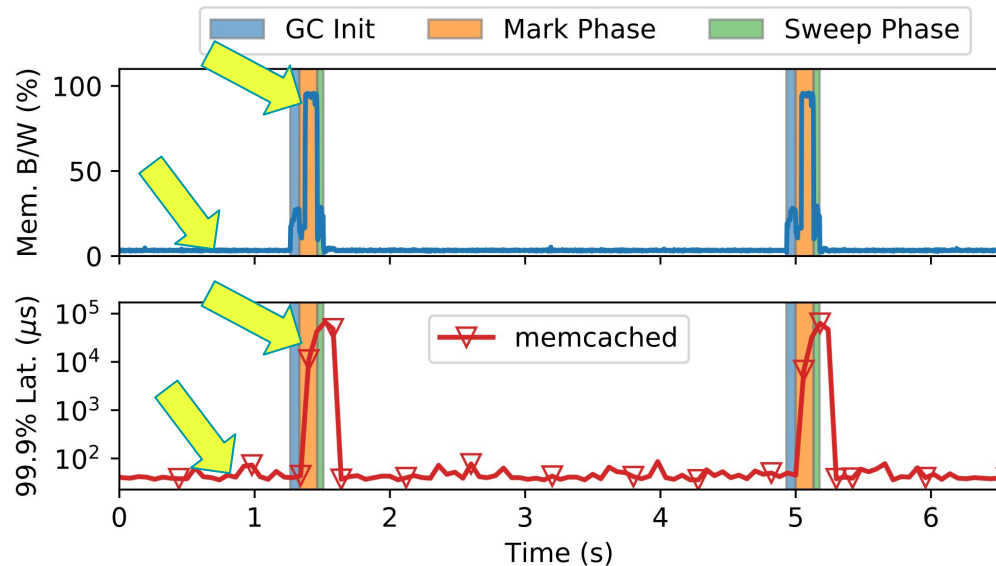
Run:

- Memcached
- garbage-collected workload

Mark Phase is memory-intensive,  
causes significant slowdown!

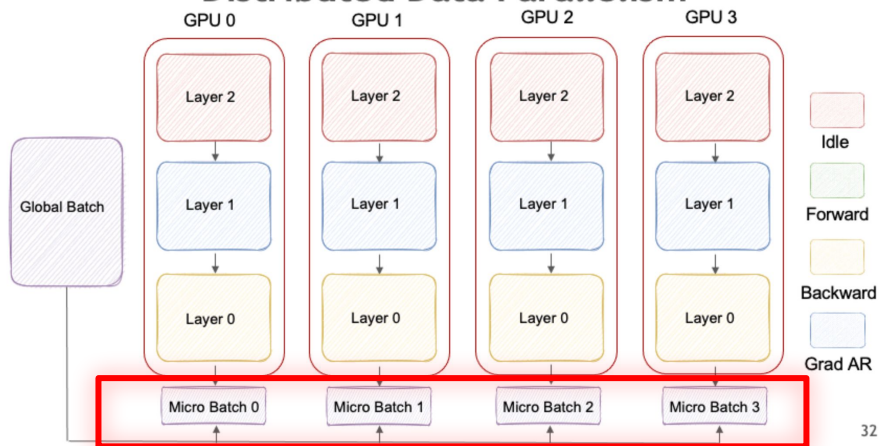
**Not only “big-data” is noisy**

Also: security scanning, video  
streaming, transcoding...



# Implications for AI Workloads

## Distributed Data Parallelism



32

### CPU at training:

- Read data
- Pre-process
- Transfer micro-batches → GPU

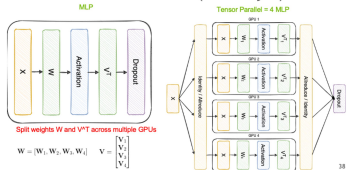
### CPU at inference:

- Input → GPU
- Output ← GPU
- Tokenization, batching
- User communication

These compete for memory bandwidth

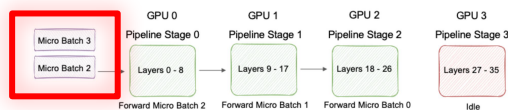
Looking for use-cases

### Tensor Parallelism (Intra-Layer)



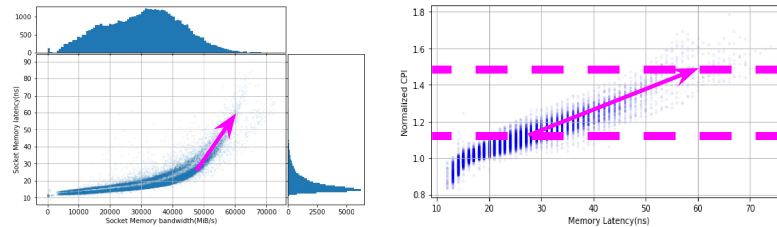
33

### Pipeline Parallelism (Inter-Layer)



# Recap: The Problem

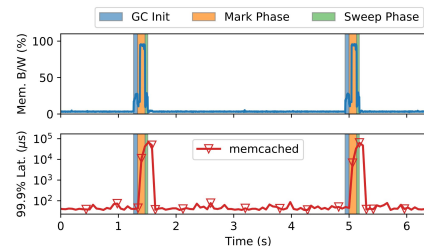
- High memory contention
- High memory latency
- High CPI (low efficiency)



Tail latency with noisy neighbor increases 4-13x



Many workloads can be noisy (e.g., garbage collection)

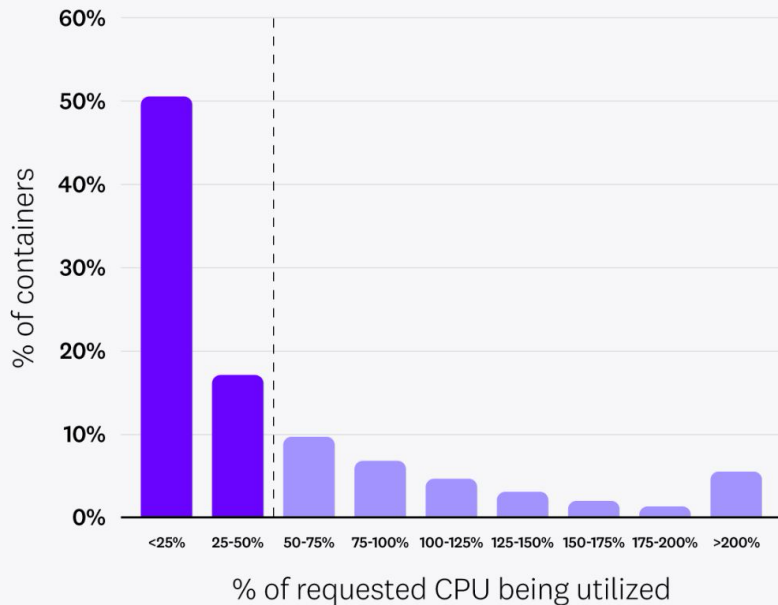


# Why Solve: Benefits

- State of cluster utilization
- How tail latency affects utilization
- Leveraging reduced tail latency



### Usage of requested CPU



**>65% OF  
CONTAINERS  
USE LESS  
THAN HALF OF  
REQUESTED CPU**

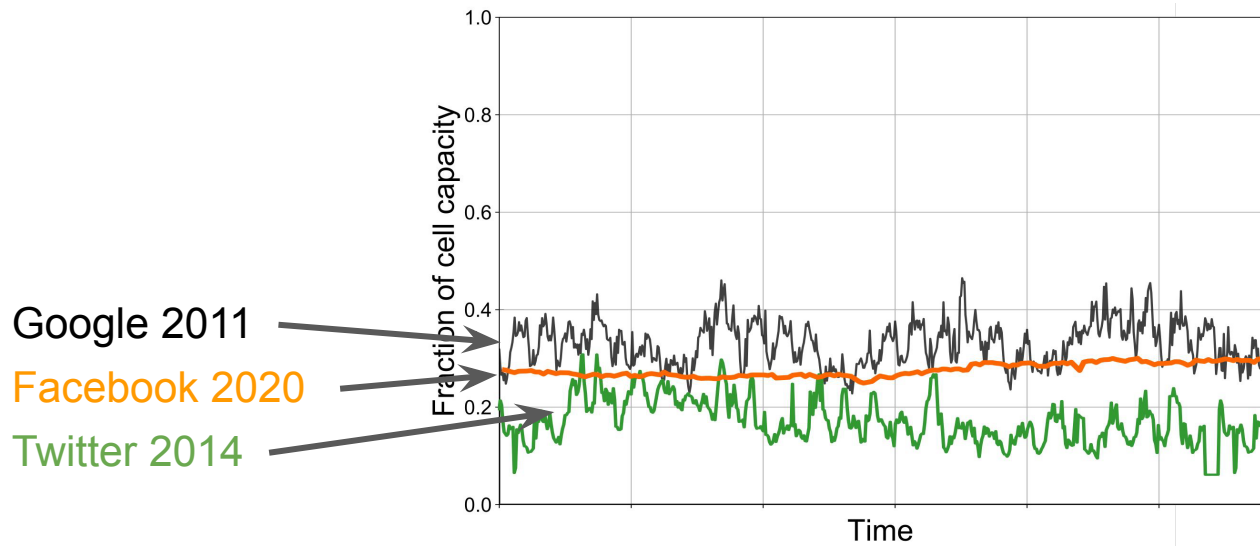
Source: Datadog

# Survey: cluster utilization

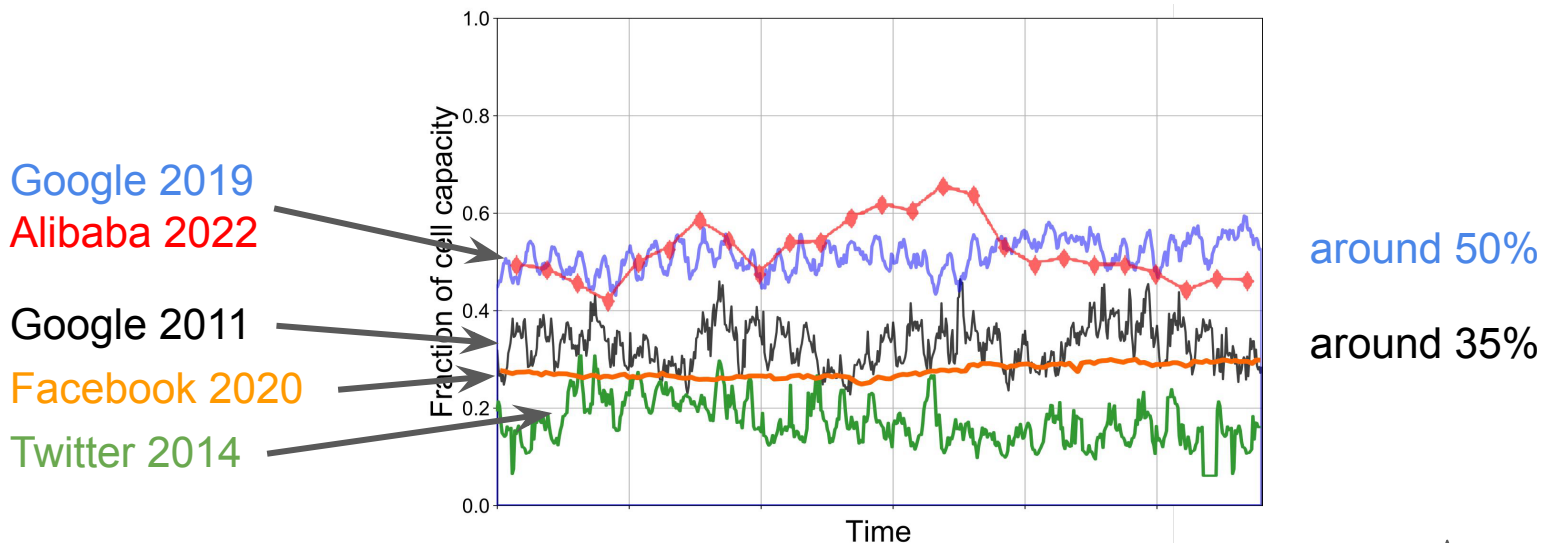
Raise hand if:

- Know prod cluster avg. CPU utilization
- Above 20%?
- Above 30%?
- Above 40%?
- Above 50%?

# Hyperscaler CPU was low...



# Then.. breakthrough?



## Vertical autoscaler

Google's Autopilot (EuroSys'20, not GKE)  
StormForge, PerfectScale, FairWinds, ...

Better packing → Less wasted resources

## Handle noisy neighbor

- Reduce cycles per request
- Improve tail latency

→ Improved SLA at higher utilization



# Tuning of Scaling Behavior

Adjusting scaling:

- Manual, (`spec.replicas``)
- HPA (Horizontal Pod Autoscaler)
  - CPU Utilization
  - Requests per second
  - Memory Utilization

How do we set thresholds?

So we meet SLOs

**Improve SLOs → Higher utilization OK**

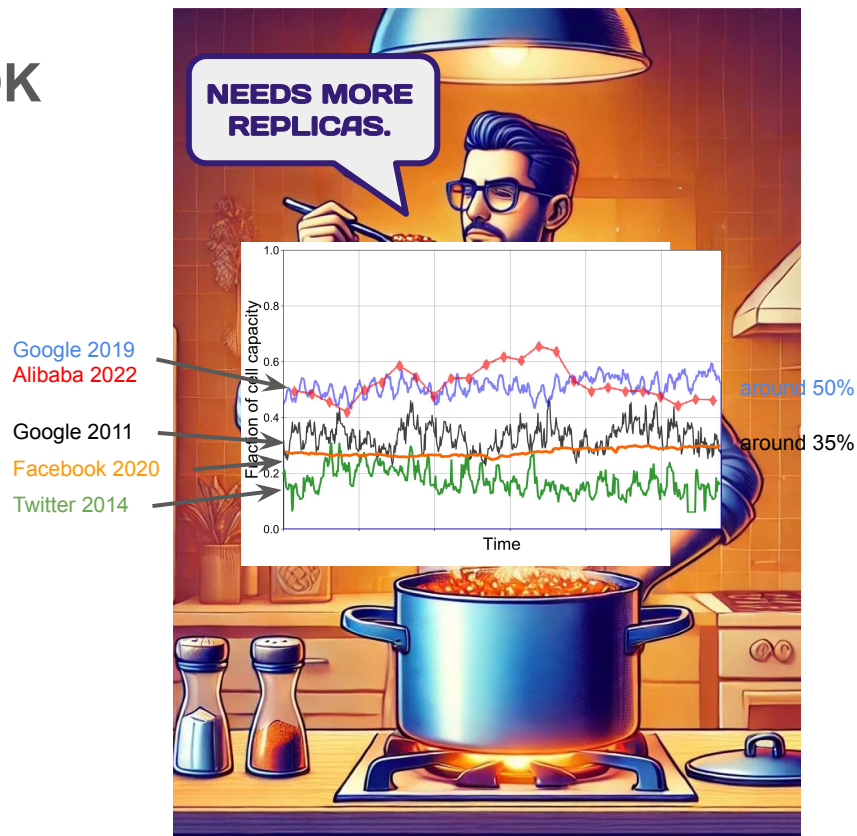




# Tuning of Scaling Behavior (cont'd)

## Improve SLOs → Higher utilization OK

- Reduce cost
  - same functionality
  - fewer replicas (at higher utilization)
- More features / business outcomes
  - added functionality
  - same replicas (at higher utilization)
- Improve SLO
  - Same functionality
  - Same replicas (same utilization)



A scenic landscape featuring a calm lake in the foreground, reflecting the golden light of a sunset. The background is dominated by jagged, snow-capped mountains under a sky with soft, wispy clouds. The scene is framed by evergreen trees on the left and right sides. The overall atmosphere is peaceful and majestic.

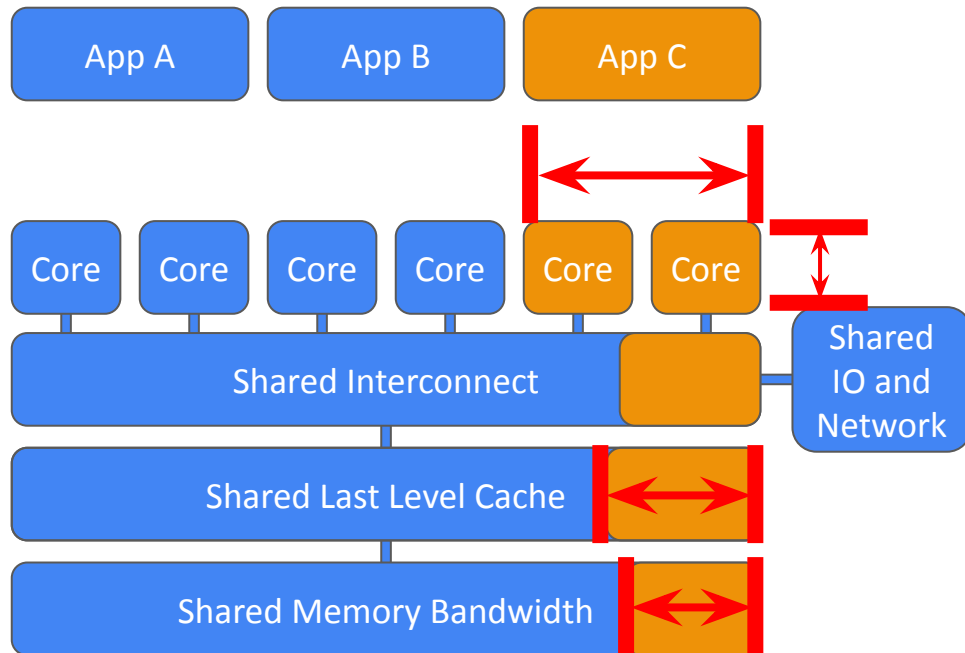
Your deployment is highly available and cost efficient

Users delight in the product and its performance

# Available Tooling

- Hardware and software mitigation
- Hardware interfaces, limitations
- Support in Linux

# Memory system mitigation knobs



Reduce demand (#cycles):

- Core pinning
- Frequency scaling

Direct control:

- Cache allocation
- Memory bandwidth limits

# Direct Control: Intel RDT's Interface

- CPU performs allocation using identifiers
  - For measurement, **RMIDs** (Resource Measurement ID)
  - For enforcement, **CLOSs** (Classes of Service)
- On every context switch, OS sets thread's (CLOS, RMID)
- Can read RMID's memory bandwidth, cache occupancy
- On EC2's m7i.metal-24xl:
  - 448 RMIDs
  - 15 CLOSs } Limited
- Supports GPU-DRAM bandwidth
  - **Intel RDT for Non-CPU Agents**



# Containers (not) to the rescue!

Different focus:

- cycles
- memory capacity (size)



<b>Cgroup v2</b>	<b>Responsibility</b>
CPU	Regulates distribution of CPU cycles
Memory	Controls distribution and accounting of memory usage
IO	Manages distribution of IO resources
PID	Limits the number of processes in a cgroup
Cpuset	Assigns specific CPU(s) and memory nodes
Device	Controls access to device files
RDMA	Regulates distribution and accounting of RDMA resources
HugeTLB	Limits the HugeTLB usage per control group
Perf_event	Allows perf events to be filtered by cgroup path
Misc	Provides resource limiting for scalar resources not covered by other controllers

# Linux supports allocation

egroup resctrl allocates:

- Memory bandwidth
- Cache space

Separate subsystem:

- Maintainers wanted self-tuning like cpuset
- Limited number of CLOS, RMIDs
- API through /sys/fs/resctrl

Demo:

- [Facebook Resource Control Demo](#)  
by Tejun Heo and collaborators

## Linux resctrl CPU support

	Kernel
<b>Intel RDT</b> Intel Resource Director Technology	<a href="#">v4.10</a> , <a href="#">v4.12</a> (2016)
<b>AMD QoS</b> AMD Platform Quality of Service	<a href="#">V5.0</a> (2018)
<b>ARM MPAM</b> Memory System Resource Partitioning and Monitoring	<a href="#">WIP</a> , <a href="#">x86→generic</a> Review needed



```
Welcome to fish, the friendly interactive shell
Type 'help' for instructions on how to use fish
htejun@slm ~-> █
```

I

Connecting

Activities Terminal Oct 2 6:02 PM demo@resctl-demo: -

### Facebook Resource Control Demo - 'q': quit

```

[ Running ] 2020-10-02 10:02:15 PM
[ config  ] satisfied: 17 missed: 0
[ oom     ] workload: +pressure -senpai system: +pressure -senpai
[ sideload ] jobs: 1/1 failed: 0 cfg_warn: 0 -overload -crit
[ sysload ] jobs: 0/0 failed: 0
[ workload ] load: 61.6% lat: 60ms cpu: 46.8% mem: 51.1g io: 7.3m

```

	cpu%	mem	swap	rbps	wbps	cpuP%	memP%	ioP%
workload	46.9	51.0g	43.5m	7.6m	-	4.4	0.0	0.0
sideload	26.6	4.4g	13.8m	558k	1.3m	88.7	0.0	0.0
hostcritical	0.8	604m	-	-	60.0k	0.4	-	-
system	0.0	140m	3.7m	-	-	0.3	-	-
user	-	-	-	-	-	-	-	-
-	78.7	56.6g	130m	8.2m	1.4m	89.7	0.0	0.0

### Workload RPS / Latency - 'g': more graphs, 't/T': change timescale

[intro.post-bench] Introduction to resource control demo - 'i': index, 'b': back

memory hog. The former will eat up as many CPU cycles as it can get its hands on along with some memory and IO bandwidth. The latter will keep gobbling up memory causing memory shortage and subsequent IOs once memory is filled up. The combination is a potent antagonist to our interactive rd-hashd.

[ Disable resource control and start the competitions ]

See the graph for the steep drop in RPS for hashd: That's the competitions taking away its resources: Not good.

Once workload's memory pressure (memP%) in the top right panel starts spiking, you might not have a lot of time before the whole system starts stalling severely. Let's stop them.

[ Stop the compile job and memory hog ]

Once RPS climbs back up and the memory usage of `workload` in the top right panel stops growing, start the same competitions but with resource control enabled and the compile job under the supervision of the sideloader:

### Management logs

```

[22:01:13 rd-agent] [ INFO] side: "compile-job" started
[22:01:14 rd-agent] [ INFO] svc: "rd-sysload-memory-hog.service" started (Running)
[22:01:32 rd-sideloader] OVERLOAD: end, resuming normal operation
[22:01:33 rd-sideloader] JOB: Starting rd-sideloader-compile-job.service
[22:01:33 rd-sideloader] Running as unit: rd-sideloader-compile-job.service
[22:01:54 rd-agent] [ INFO] svc: "rd-sysload-memory-hog.service" transitioned from Running to Other("deactivating:stop-sigterm")
[22:01:54 rd-agent] [ INFO] svc: "rd-sysload-memory-hog.service" stopped (NotFound)

```

[ Start the competitions under full resource control ]

Watch the stable RPS. rd-hashd is now fully protected against the competitions. The compile job and memory hog are throttled. The compile job doesn't seem to be making much progress. This is because sideloads (workloads under the sideloader supervision) are configured to have lower priority than sysloads (workloads under system). Don't worry about the distinction between sideloads and sysloads for now. We'll revisit them later.

Let's stop the memory hog and see what happens.

### Other logs

```

[22:02:16 rd-sideloader-compile-job] CC security/apparmor/match.o
[22:02:16 rd-sideloader-compile-job] CC crypto/scatterwalk.o
[22:02:16 rd-sideloader-compile-job] CC arch/x86/events/intel/bts.o
[22:02:16 rd-sideloader-compile-job] CC security/selinux/netlink.o
[22:02:16 rd-sideloader-compile-job] CC crypto/proc.o
[22:02:16 rd-sideloader-compile-job] AR arch/x86/kernel/fpu/built-in.a
[22:02:16 rd-sideloader-compile-job] CC arch/x86/kernel/irq_work.o
[22:02:16 rd-sideloader-compile-job] HDRTEST usr/include/drm/msm_drm.h

```

[ Stop the memory hog ]

fine and the compile job is now making reasonable forward progress. Sideloads are now sharing the machine safely and productively, as possible before.

Continue reading to learn more about the various components which make this possible.

[ Next: Cgroup and Resource Protection ]

Main workload is doing fine

# Hypervisors support allocation

- VMware in Telco Cloud Automation  
→ Telco + Finance?
- Static

## Intel RDT support in hypervisors

	Cache partitioning	Memory bandwidth
Xen	✓	✓
KVM	✓	
VMware	✓	Monitor (vSphere)
Hyper-V		Monitor PMU
ACRN	✓	✓

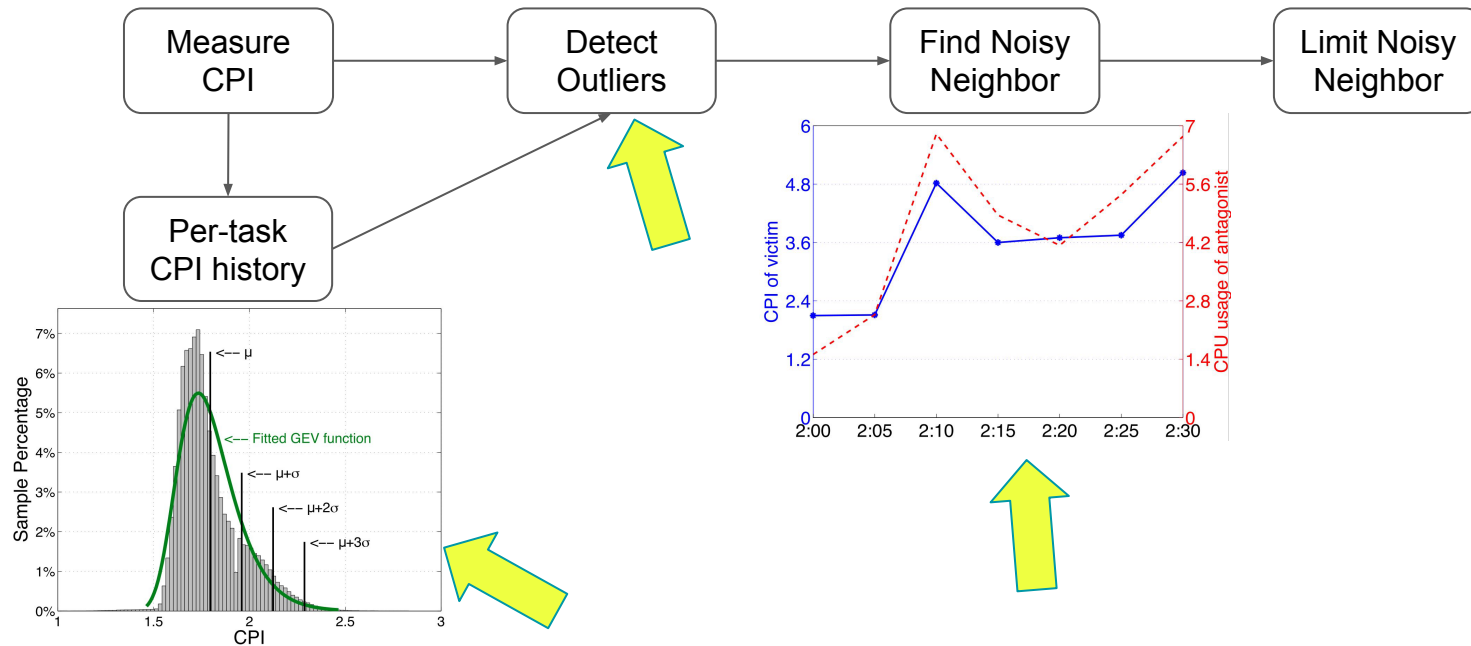


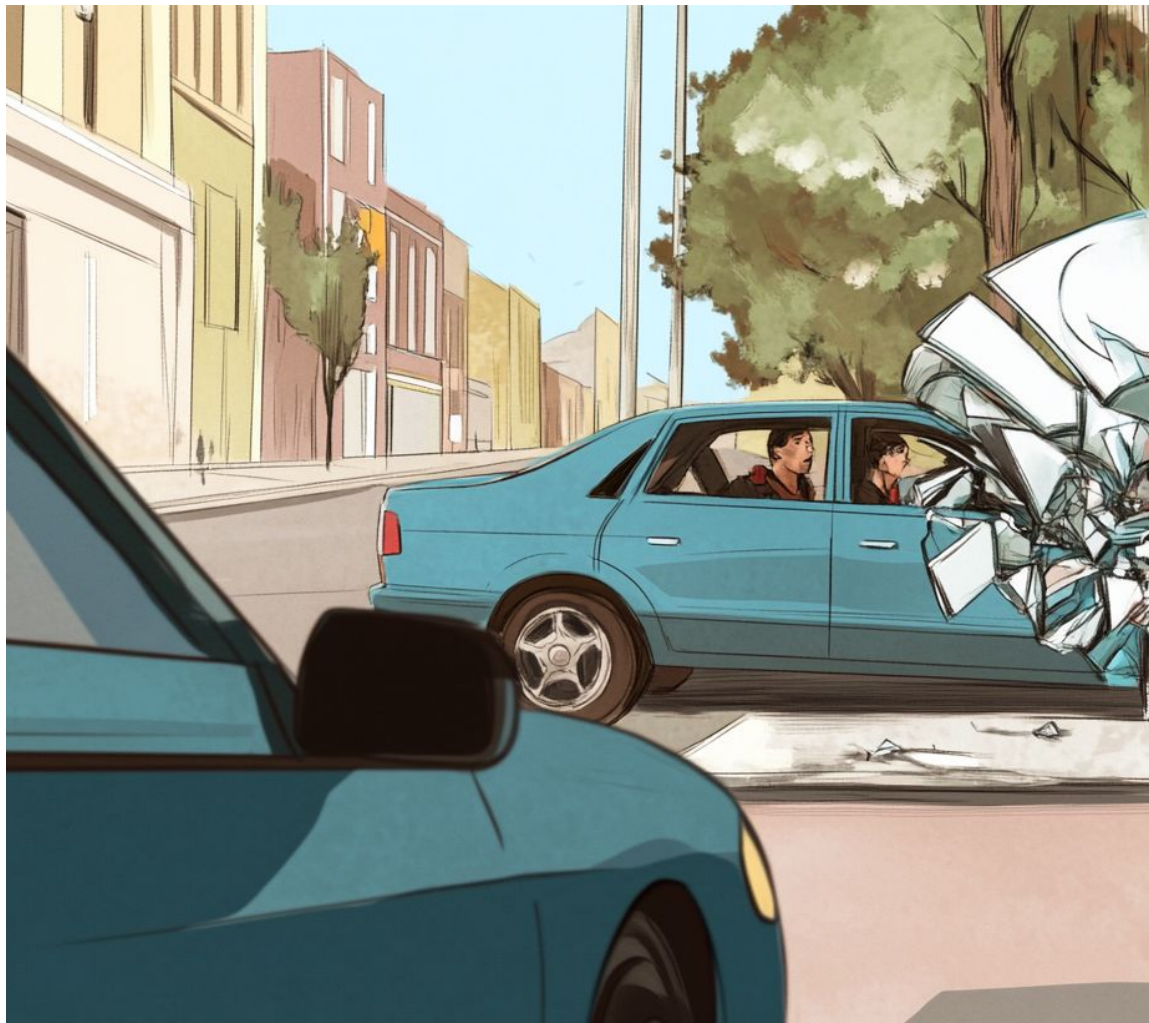
# Mitigation Systems

- Type 1: cycles per instruction (CPI)
- Type 2: latency control
- Type 3: usage control

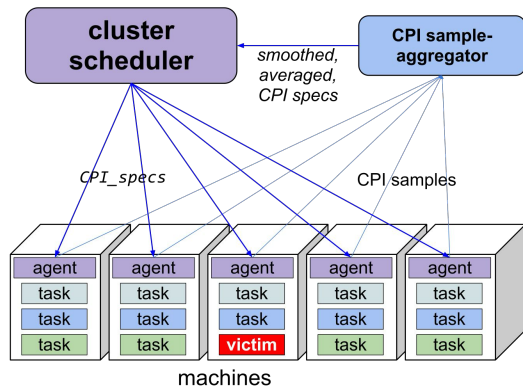
# Type 1: Cycles Per Instruction (CPI)

Uses: High interference → high CPI





# Type 1: Cycles Per Instruction



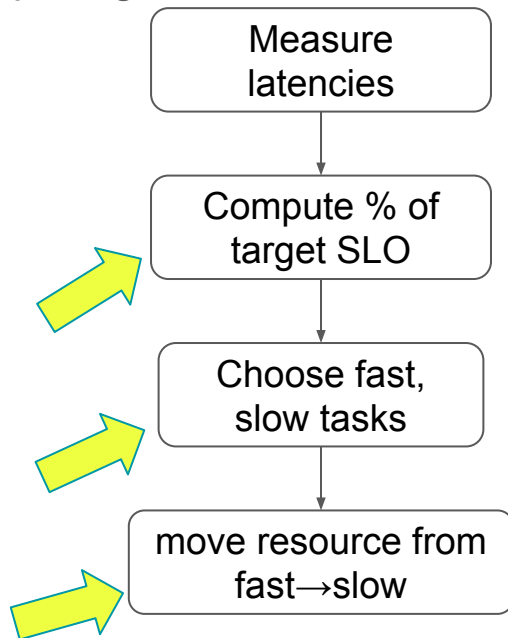
“We have rolled out CPI<sup>2</sup> to all of Google’s shared compute clusters.” – paper authors, 2013 @Google

	Type 1: CPI
Measurement	Cycles, Instructions
Averaging	High
Cluster components	Aggregator
Pros	<ul style="list-style-type: none"><li>• Simple to measure</li></ul>
Cons	<ul style="list-style-type: none"><li>• Complex deployment</li><li>• Averaging → Slow reaction</li></ul>

# Type 2: Latency Control

Use *application* latency

Example algorithm:



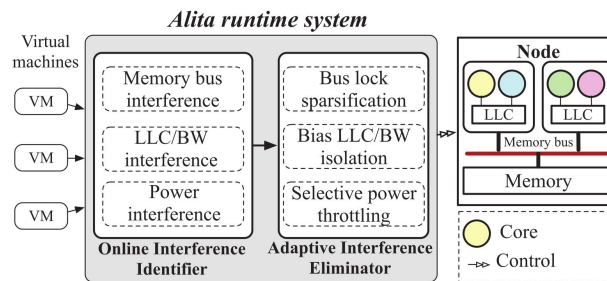
	Type 2: Latency
Measurement	App latency
Averaging	Medium
Cluster components	Node only
Pros	<ul style="list-style-type: none"><li>• No profiling → Node-local</li><li>• Control what you care about</li></ul>
Cons	<ul style="list-style-type: none"><li>• High developer effort</li><li>• Noisy signal → Averaging</li></ul>



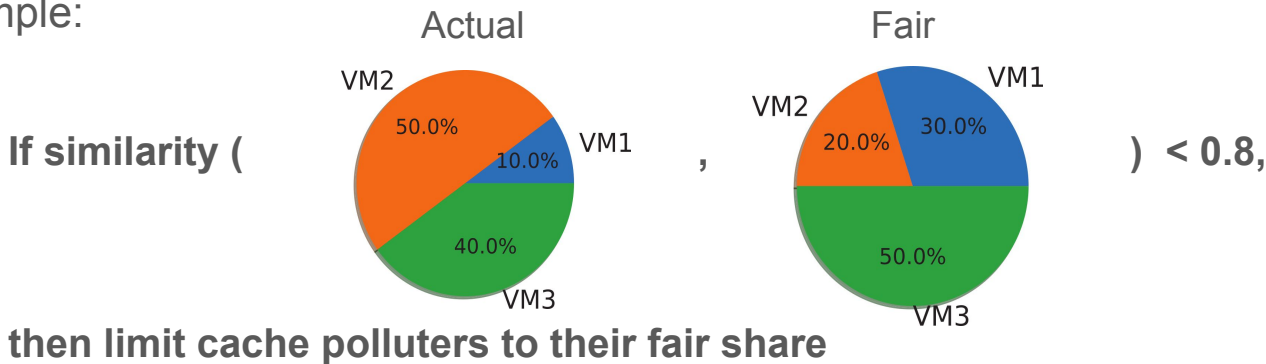
# Type 3: Usage Control

How:

- Measure per-app resource usage
- Find unfair allocation
- Limit offender



Cache example:



Do we really want a fair allocation?



# Type 3: Usage Control (cont'd)

	Type 3: Usage
Measurement	CPU counters
Averaging	Low
Cluster components	Node only
Pros	<ul style="list-style-type: none"><li>• Node-local</li><li>• Measure directly<ul style="list-style-type: none"><li>→ Fast reaction</li><li>→ Easy to reason</li></ul></li></ul>
Cons	<ul style="list-style-type: none"><li>• Do we really want fair allocation?</li></ul>

Deployed in production, > 2 years:

- 30k nodes (24 to 48 cores each)
- 250k VMs




- authors, 2020 @Alibaba Cloud

# Summary: Mitigation Systems

	Type 1: CPI	Type 2: Latency	Type 3: Usage
Measurement	Cycles, Instructions	App latency	CPU counters
Averaging	High	Medium	Low
Cluster components	Aggregator	Node only	Node only
Pros	<ul style="list-style-type: none"><li>• Simple to measure</li></ul>	<ul style="list-style-type: none"><li>• No profiling → Node-local</li><li>• Control what you care about</li></ul>	<ul style="list-style-type: none"><li>• Node-local</li><li>• Measure directly → Fast, Easy to reason</li></ul>
Cons	<ul style="list-style-type: none"><li>• Complex deployment</li><li>• Averaging → Slow reaction</li></ul>	<ul style="list-style-type: none"><li>• High developer effort</li><li>• Noisy signal → Averaging</li></ul>	<ul style="list-style-type: none"><li>• Do we really want fair allocation?</li></ul>

- The three “categories” of published systems
- Many good ideas → general purpose
- Usage control (#3) is promising

# Recap

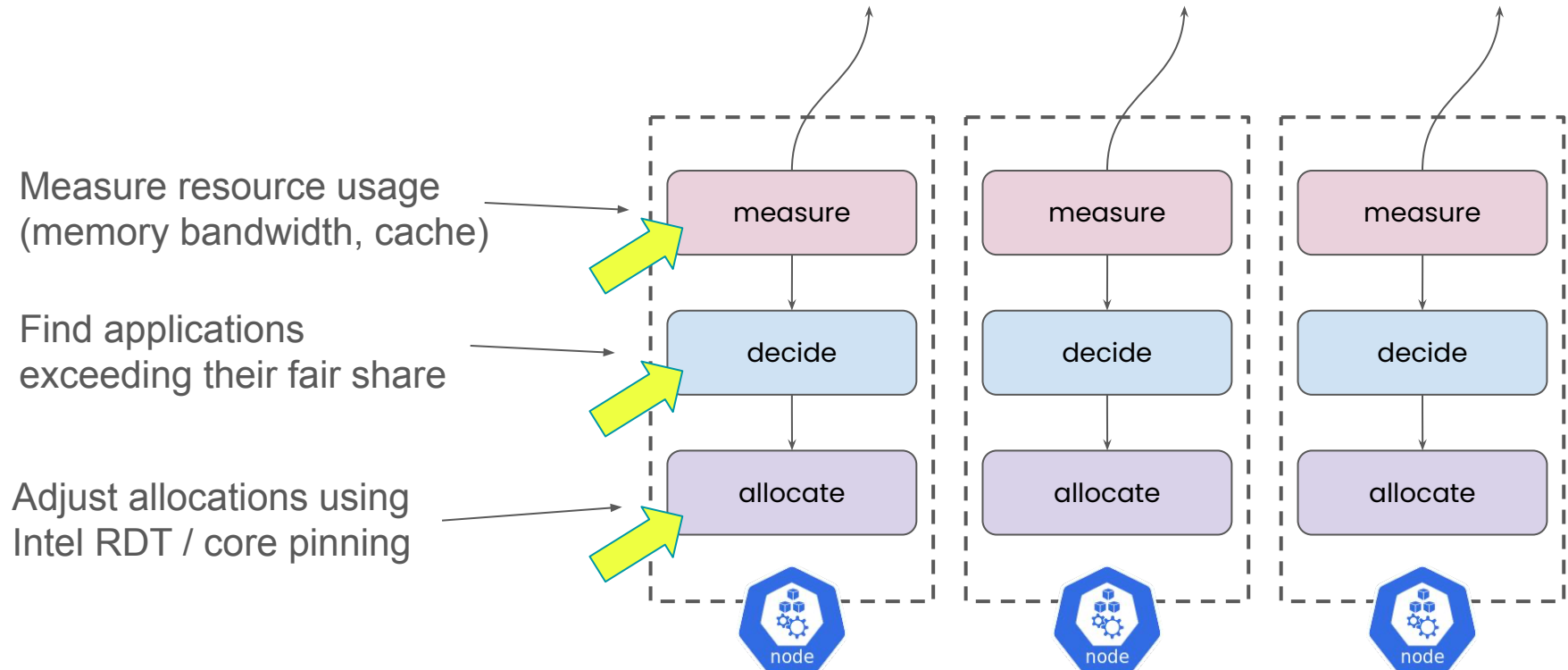
- Cache and bandwidth crunch can increase tail latency 4x - 13x
- Reducing tail latency: cost , functionality , response times 
- CPUs support monitoring and allocation, Linux has resctrl
- Systems that explicitly monitor and allocate: simple, fast reaction



# Current Efforts

- Kubernetes Deployment
- Major Challenge
- Collector architecture
- Goals, Open Questions
- Community Next Steps

# Kubernetes Deployment



# Traditional Observability: Too Slow!

Many collectors measure at 5 second intervals or more

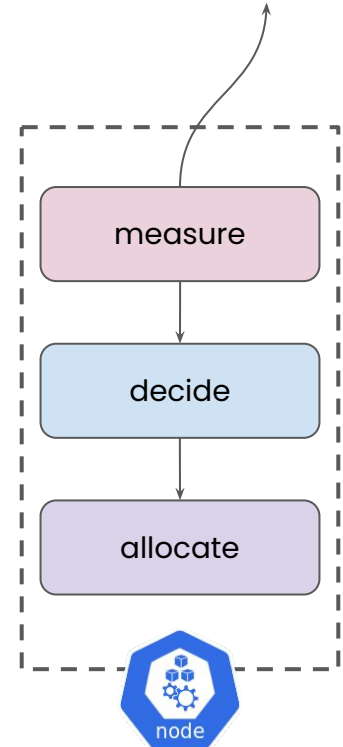
Garbage collection:

- Differs by language, application, heap size
- Minor GC: “1-10 milliseconds, every 0.1-10 seconds”
- Major GC: “10-100 milliseconds, every 10-100 seconds”

Heavy user transaction:

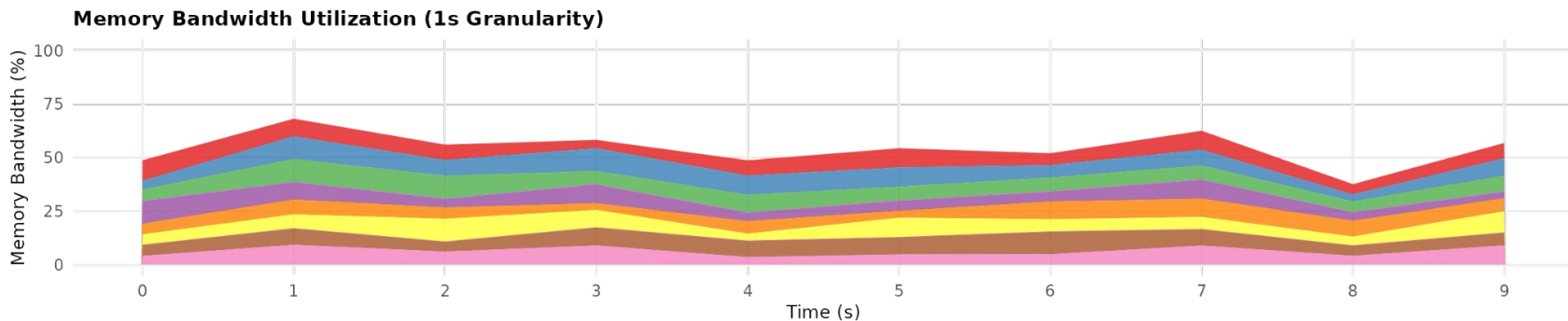
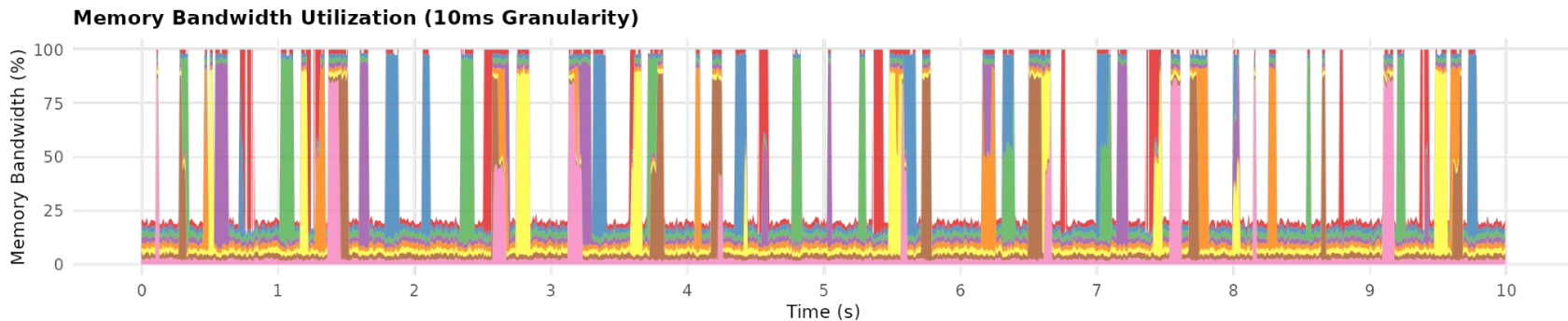
- “500 milliseconds every 10s of seconds”

5 seconds is too slow!



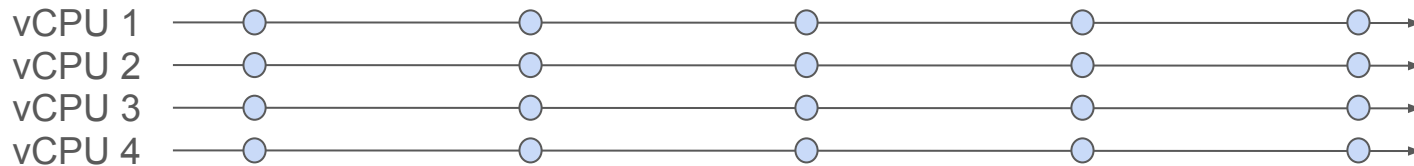
# Why we need frequent measurements

Simulation: 8 applications, each noisy for 10-100 ms every second

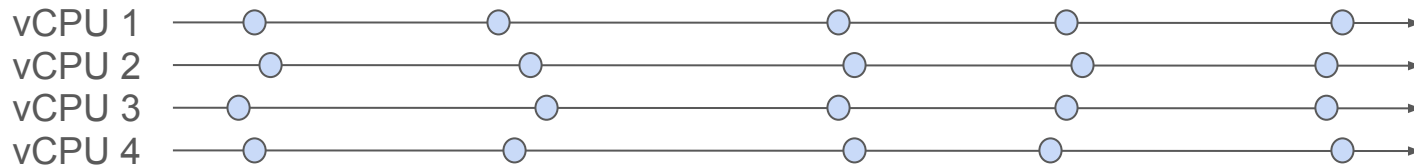


# Measuring at 1 millisecond intervals

Ideally:



In practice:



Q: How much jitter are we going to have?

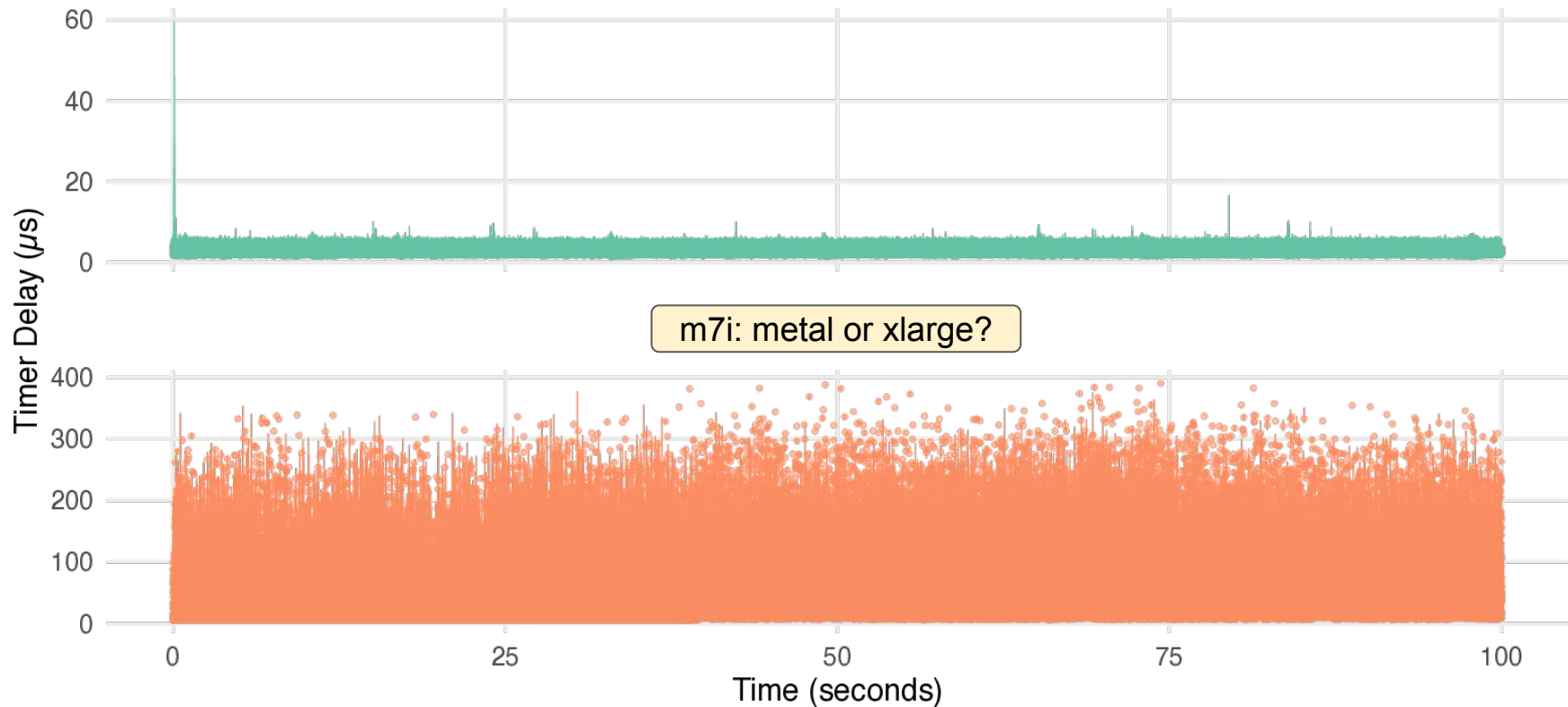


# Sync Timer Benchmark Results

Points show mean delay, vertical lines show min-max range

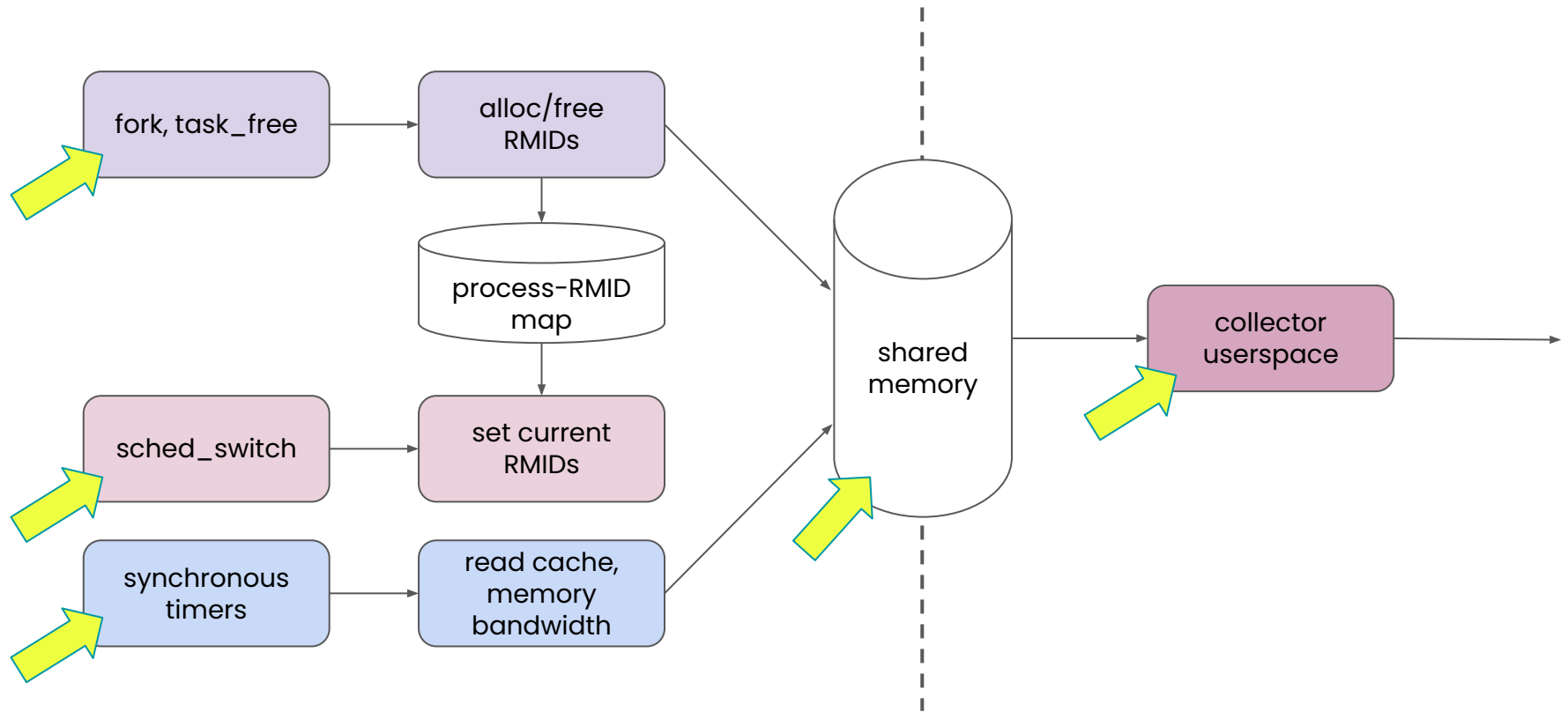
m7i: metal or xlarge?

m7i: metal or xlarge?

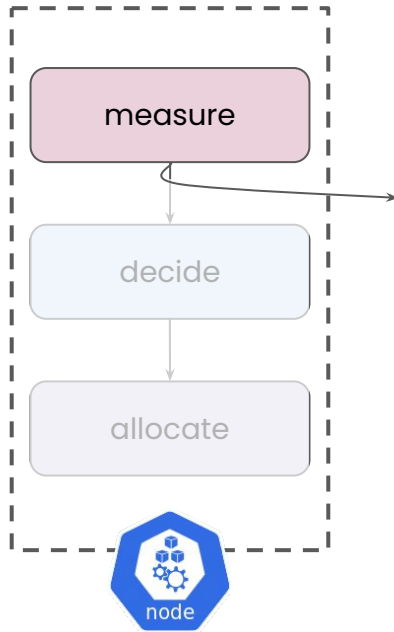


\* idle machines

# Collector architecture



# Data helps develop detection algorithms



Collect per-RMID, per-millisecond:

- Cache utilization
- Memory bandwidth
- Cycles per Instruction (CPI)
- Cache misses

**Goal: develop detection algorithms**

How:

- Calculate detector values on data
- Validate detection with CPI

Synthetic workloads ok; **looking for production data**

# Open Questions

- How to detect contention  
(bandwidth saturation)
- Is 1 millisecond frequent enough?  
(100 microseconds)
- Missing critical measurements?  
(CPU frequency)
- What if resctrl is unavailable?  
(use perf counters)

# Security and Privacy

- Schema contains only profiling data:
  - Process names, PIDs
  - Container names (soon)
  - Cycles, instructions, LLC misses
  - Memory bandwidth, cache utilization
- No PII captured
- Runs locally, produces parquet files
  - No external communication



# Overhead

- Designing for:
  - < 0.1% in-line overhead
  - < 1% userspace
  
- Even if more, might still make sense!

average	p95
40 ms	250 ms



average	p95
42 ms	75 ms

# Community Next Steps



## unvariance/collector



Ian Off  
Tarun Verghis  
Darshan Dedhia  
Nimrod Liberman

### Call for:

- Contributors
  - Collector
  - Kubernetes benchmarks
- Data contributions
  - Deploy in test/staging
  - Advance detector development

### Contact:

CNCF Slack: @Jonathan Perry  
[yonch@yonch.com](mailto:yonch@yonch.com)

Set up time: [yonch.com/collector](https://yonch.com/collector)

