



Cluster management at Google

LISA 2013

john wilkes (johnwilkes@google.com)
Software Engineer, Google, Inc.

We own and operate data centers around the world



✗ Not a data center



✓ Data center

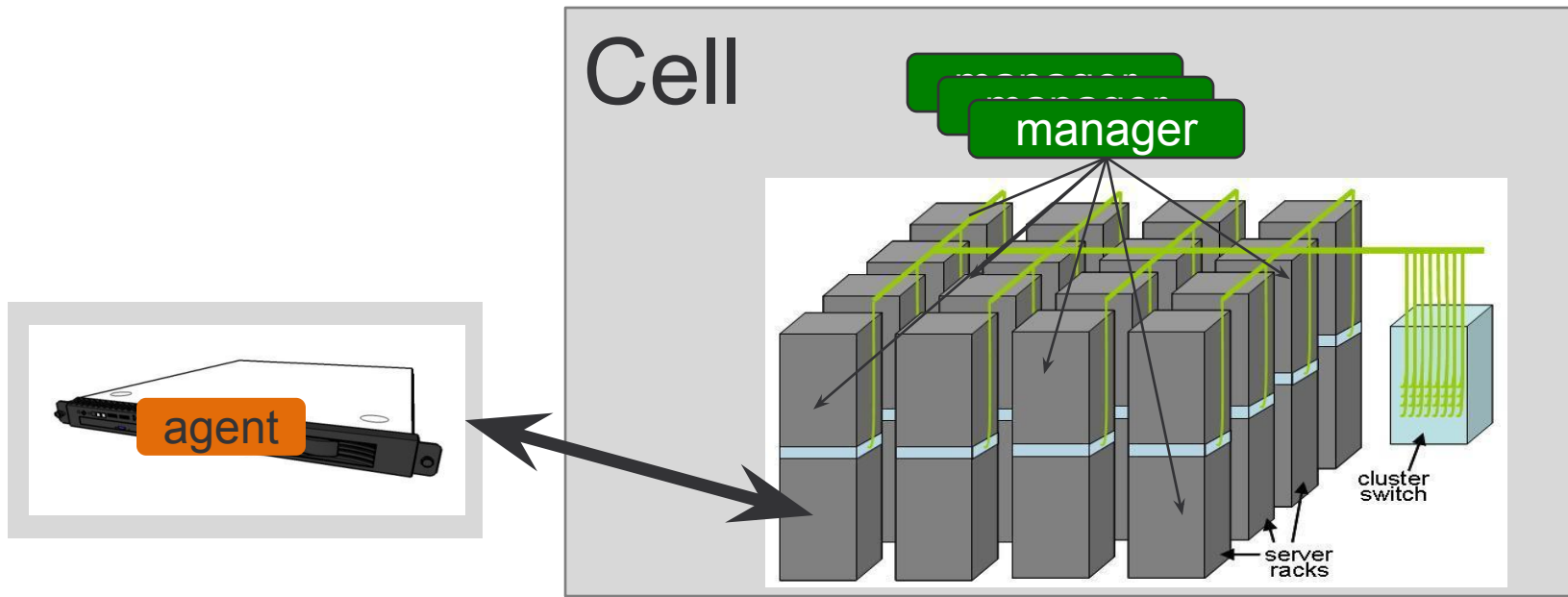




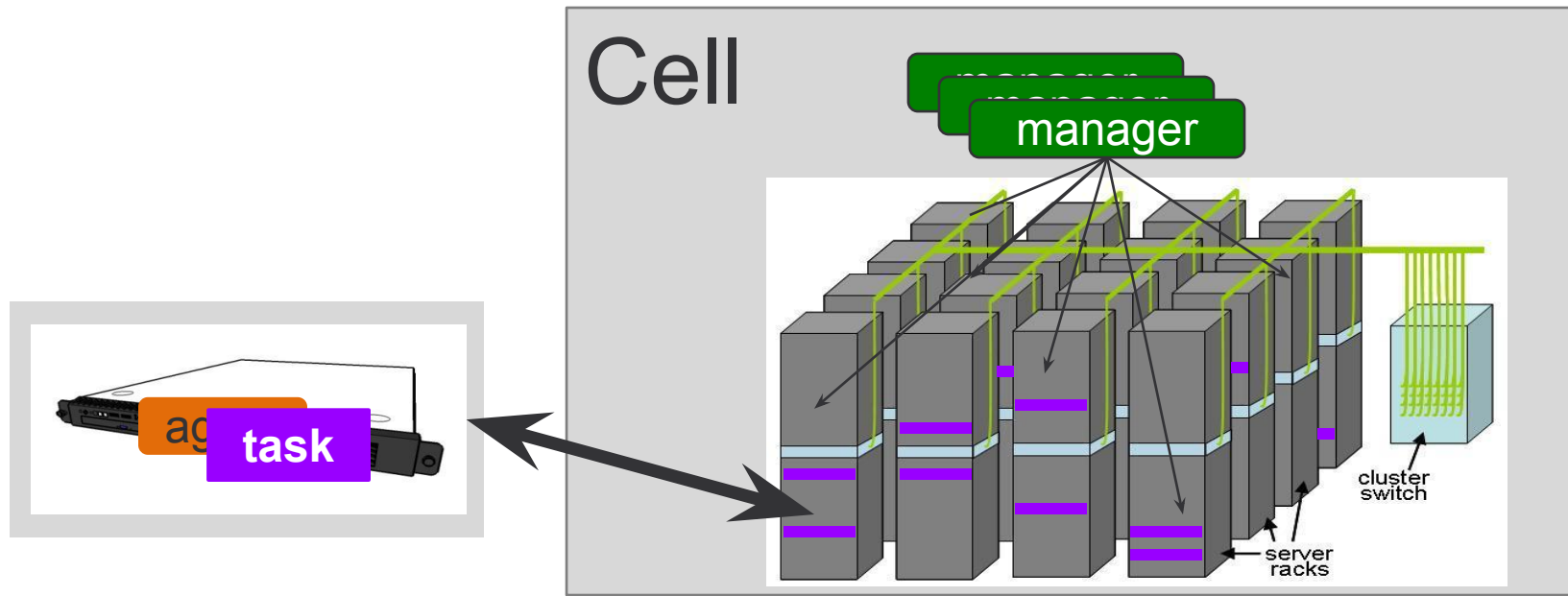




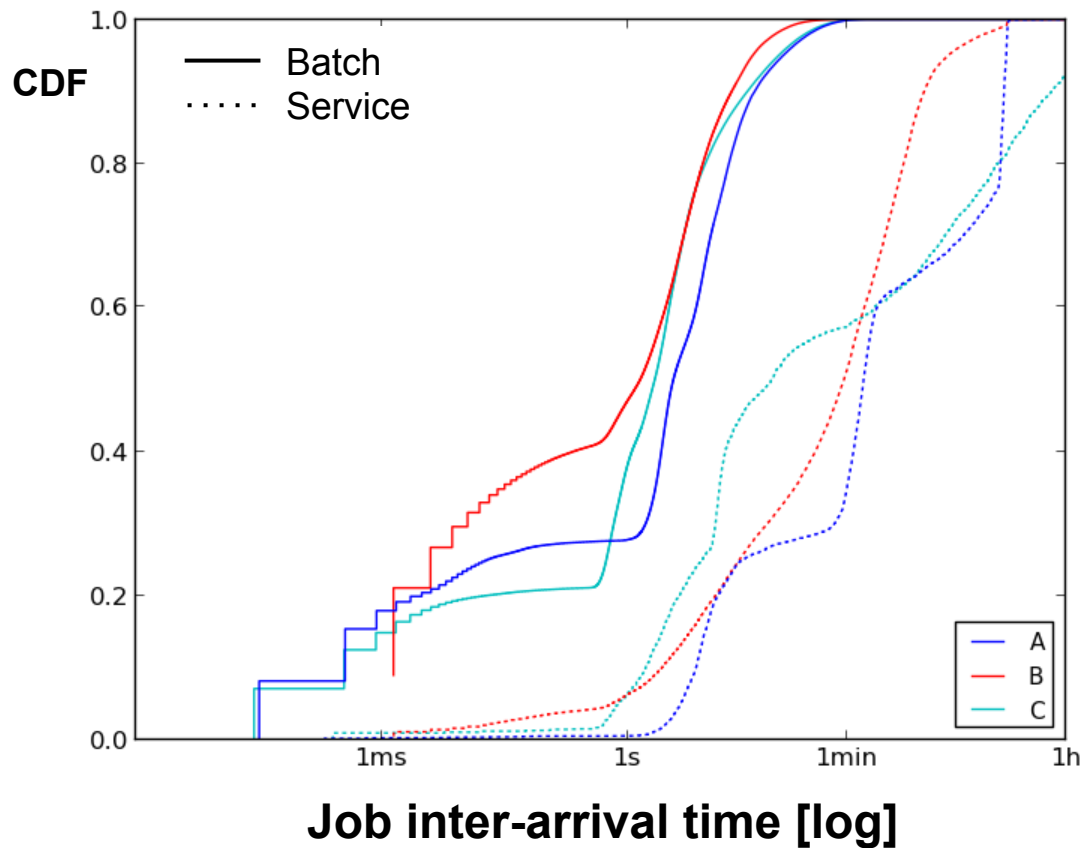
A cluster is managed as 1 or more **cells**



A cell runs **jobs**, made up of (1 to thousands) of **tasks**, for internal users



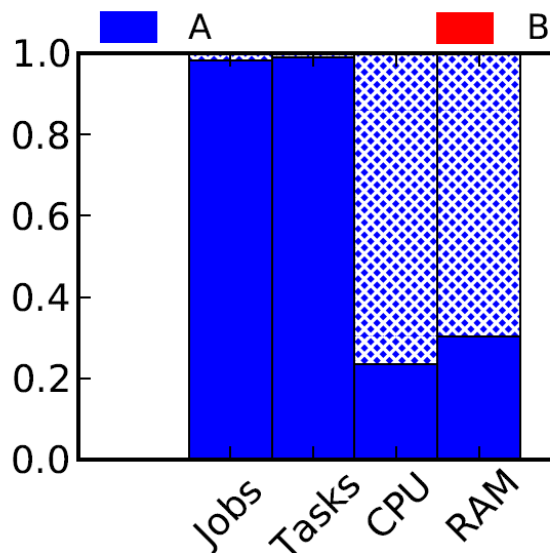




Cluster A

Medium size

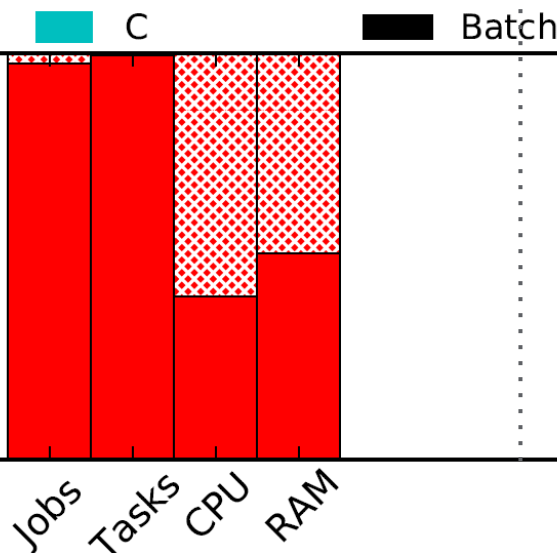
Medium utilization



Cluster B

Large size

Medium utilization

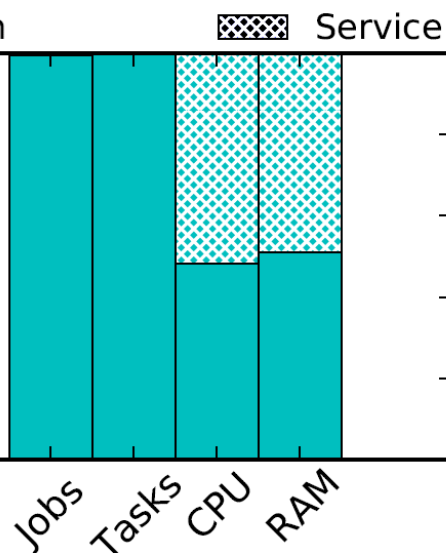


Cluster C

Medium (12k mach.)

High utilization

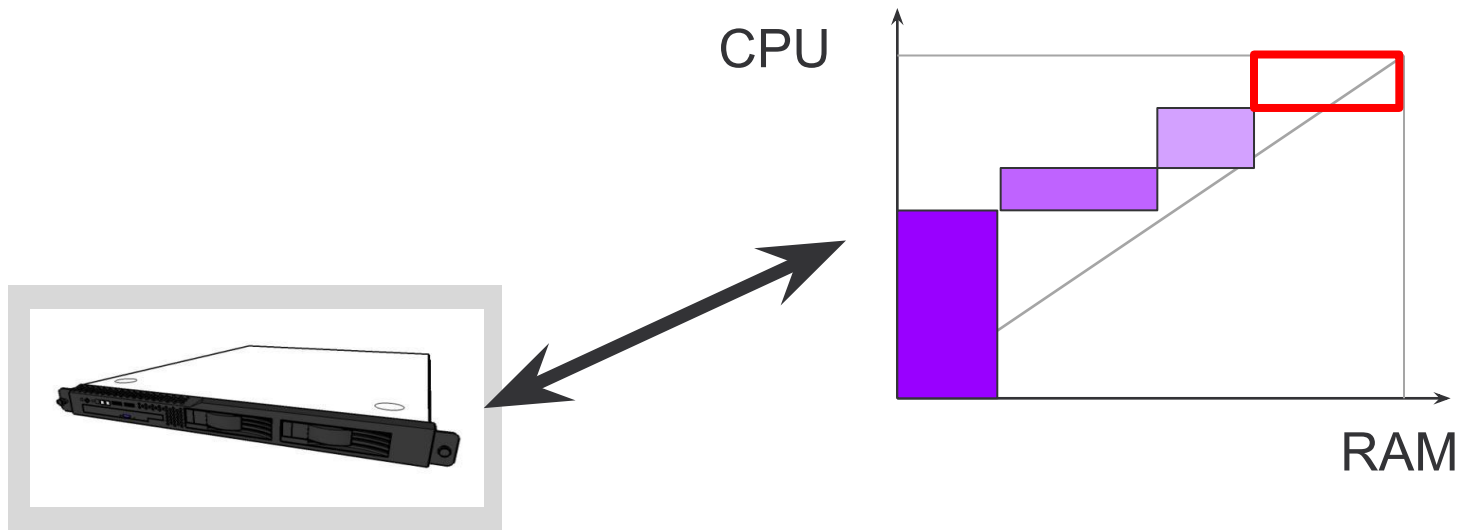
Public trace

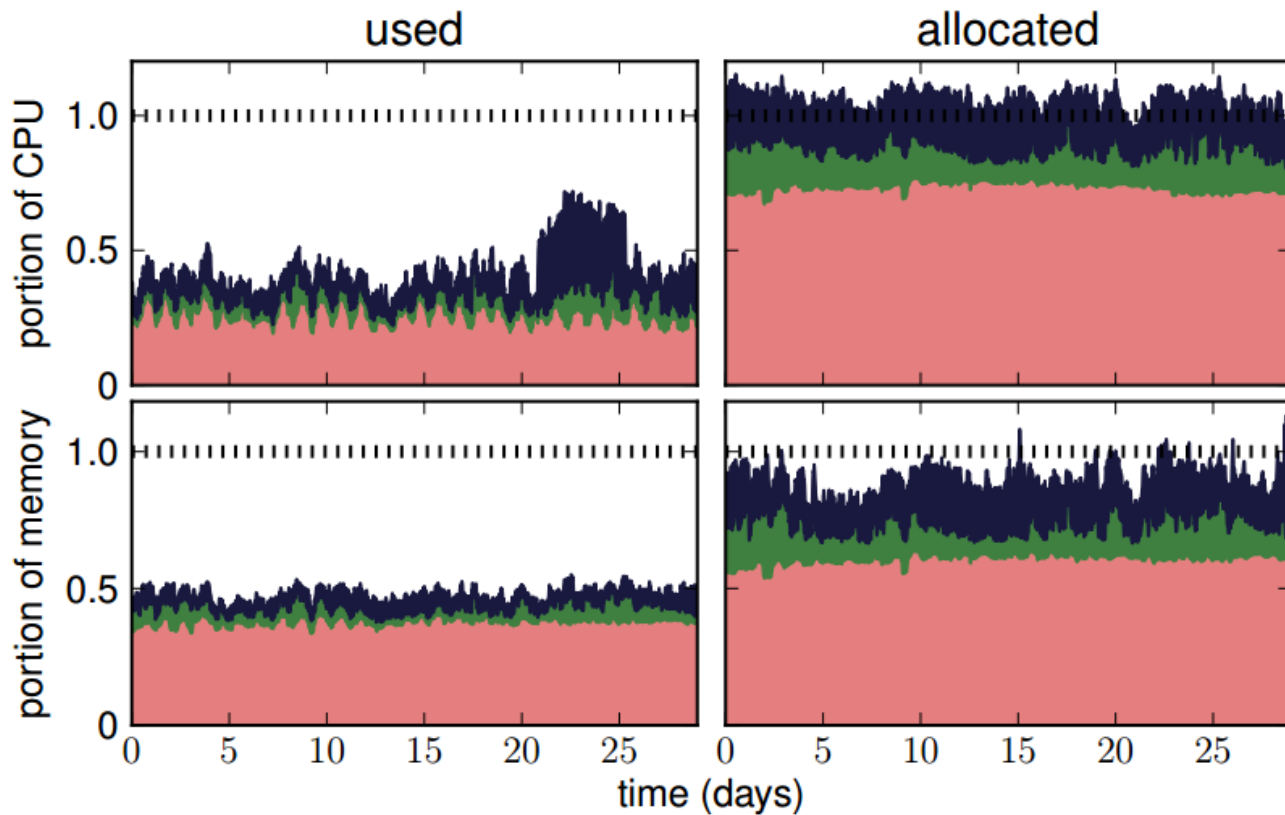


Jobs/tasks: counts

CPU/RAM: resource seconds [i.e. resource job runtime in sec.]

Placing tasks onto machines is just a *knapsack problem*, with constraints



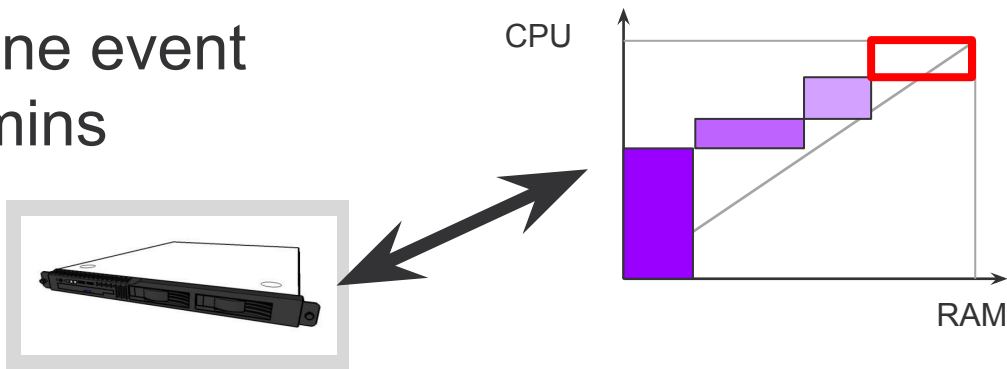


Want to try it yourself?

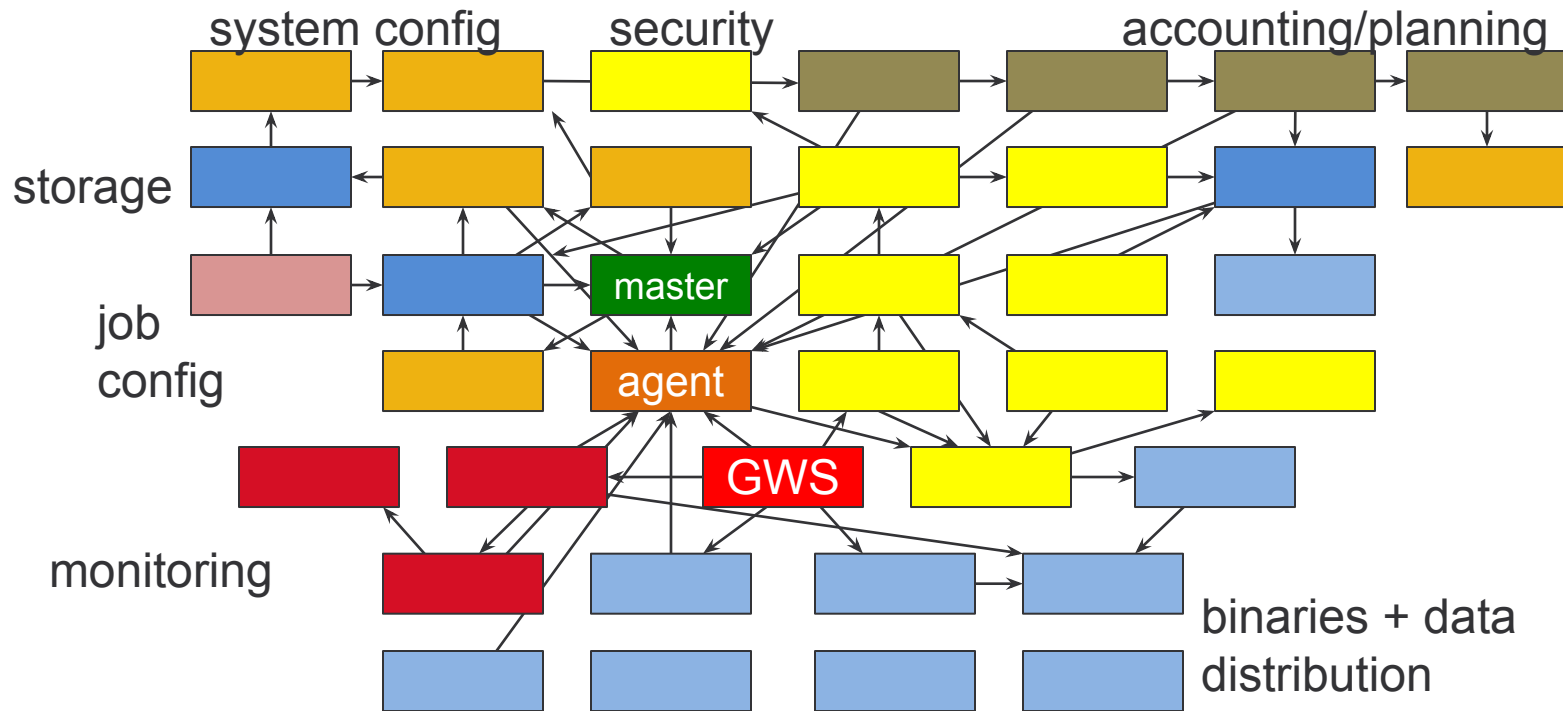
- search for [[john wilkes google cluster data](#)]

29 day cell workload trace from May 2011

- ~12.5k machines
- every job/task/machine event
- task usage every 5 mins
- (obfuscated data)





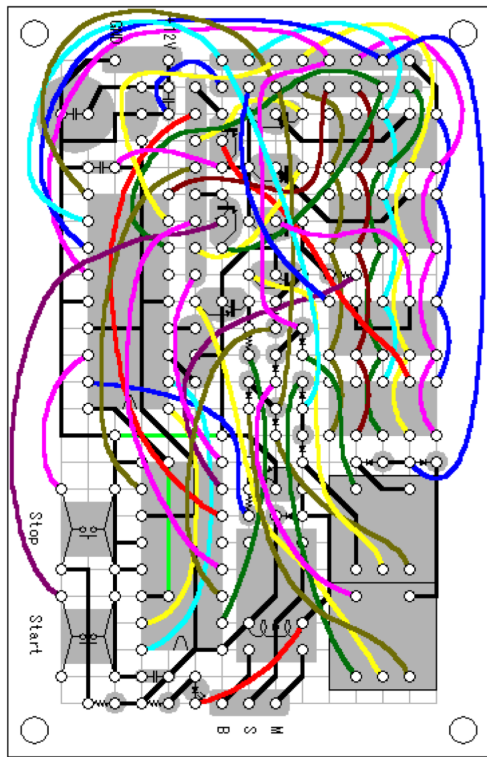


Make an app work right once: ***simple***

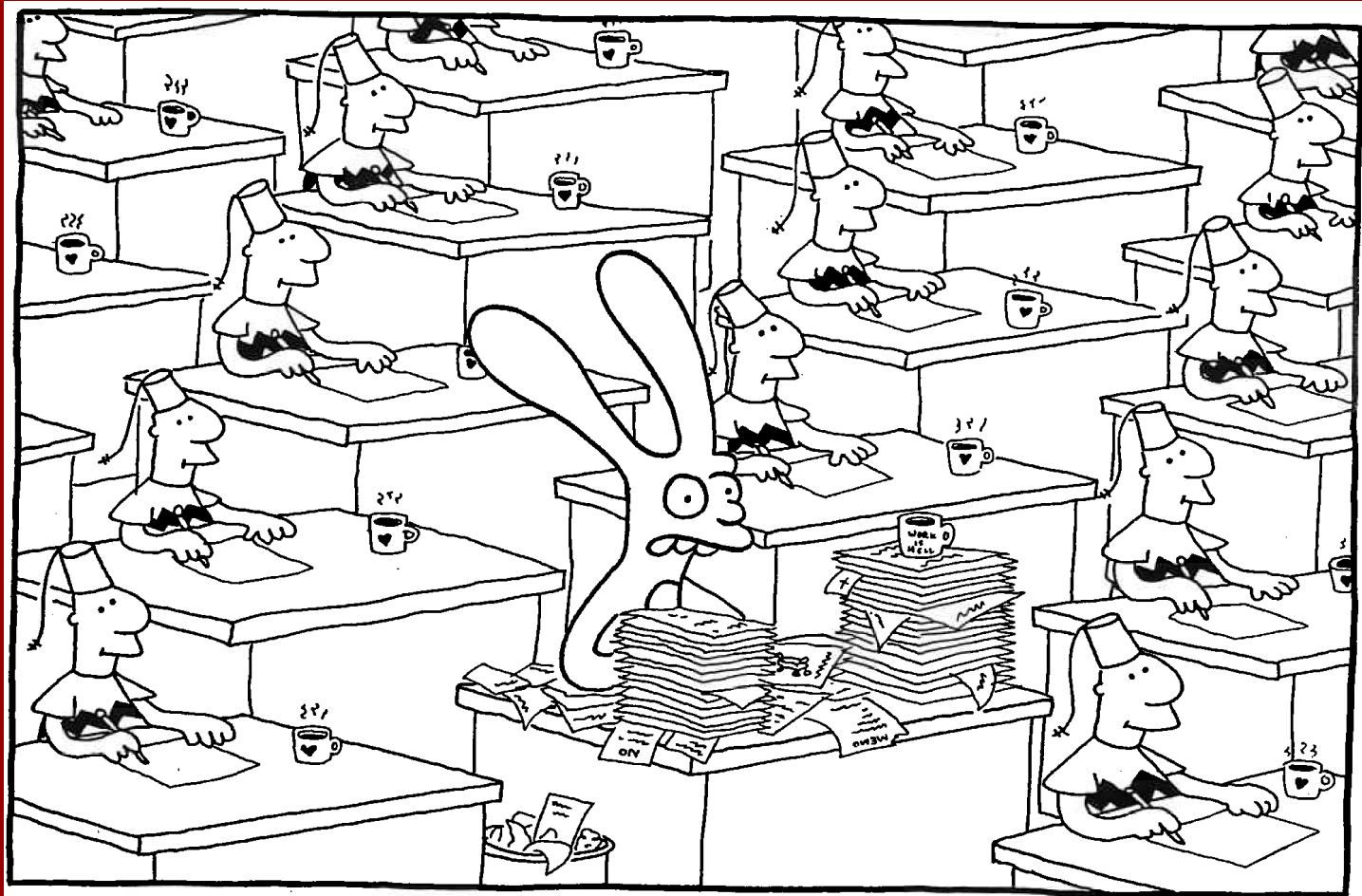
Google Docs uses ~50 systems and services

Make it run in production: ***priceless***

- run it in 6 cells
- fix it in an emergency
- move it to another cell



If you aren't measuring it, it is out of control ...



1. Availability

- % time master is up
- max outage duration

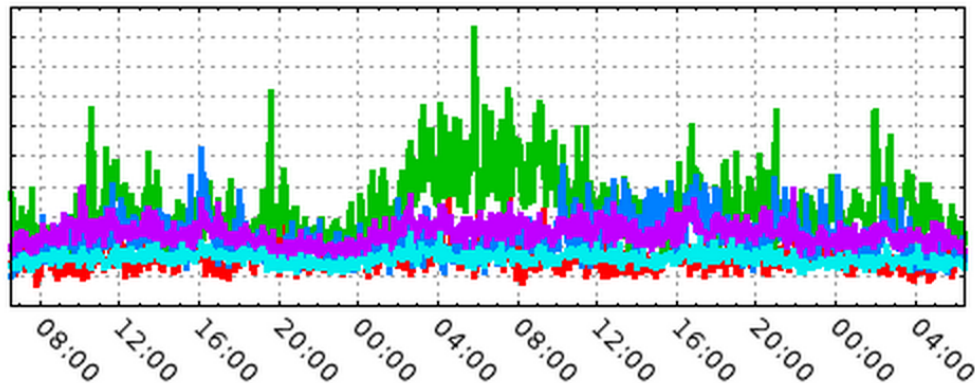
2. Performance

- RPC request latencies
- Time to schedule+start job

3. Evictions

- % jobs with too high an eviction rate
- % jobs with too many concurrent evictions
- % tasks evicted without proper notice

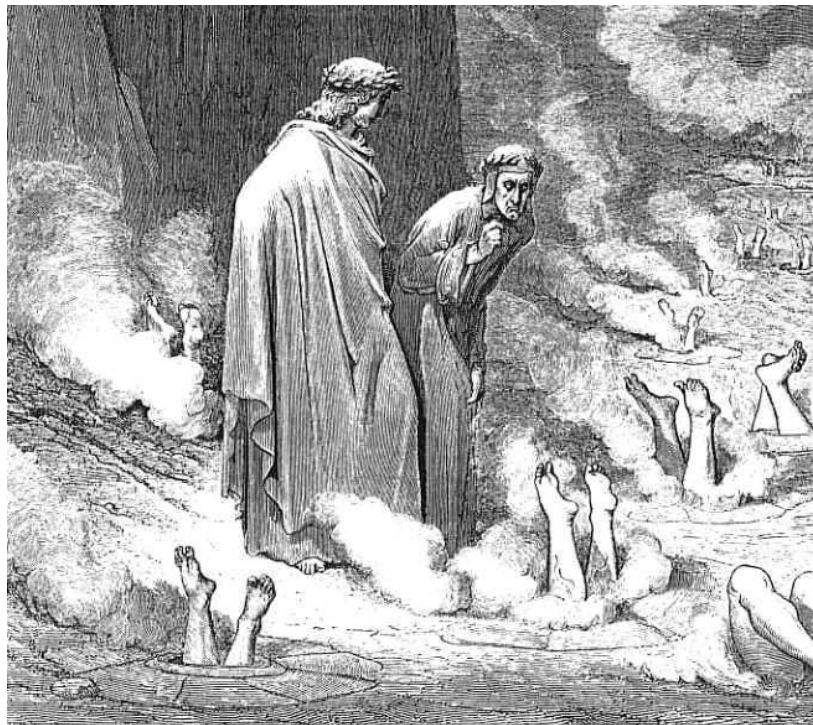
Mutating RPC latency (overall percentiles and top 3 cells)

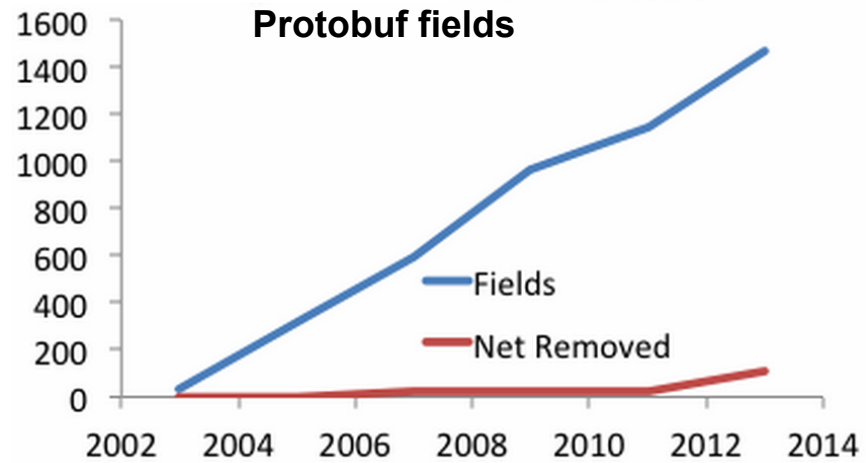
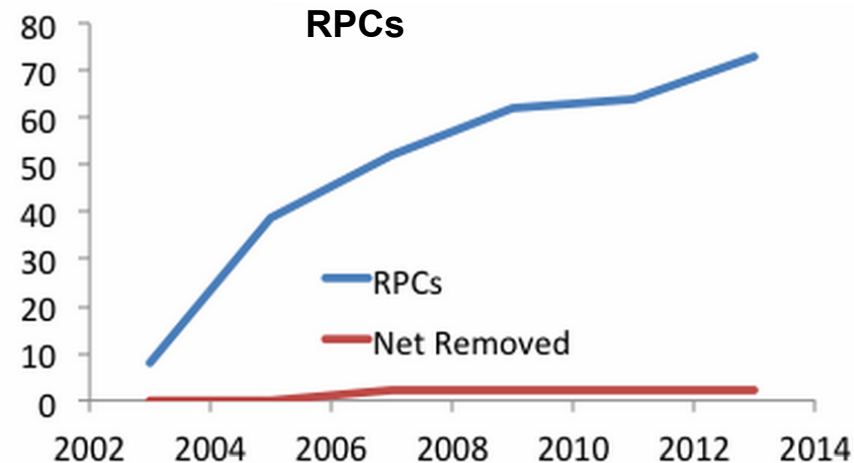


The current system was built 2003-4. It works pretty well 😊

But ...

- growing scale (#cores)
- inflexibility (adding features)
- internal complexity (adding people)

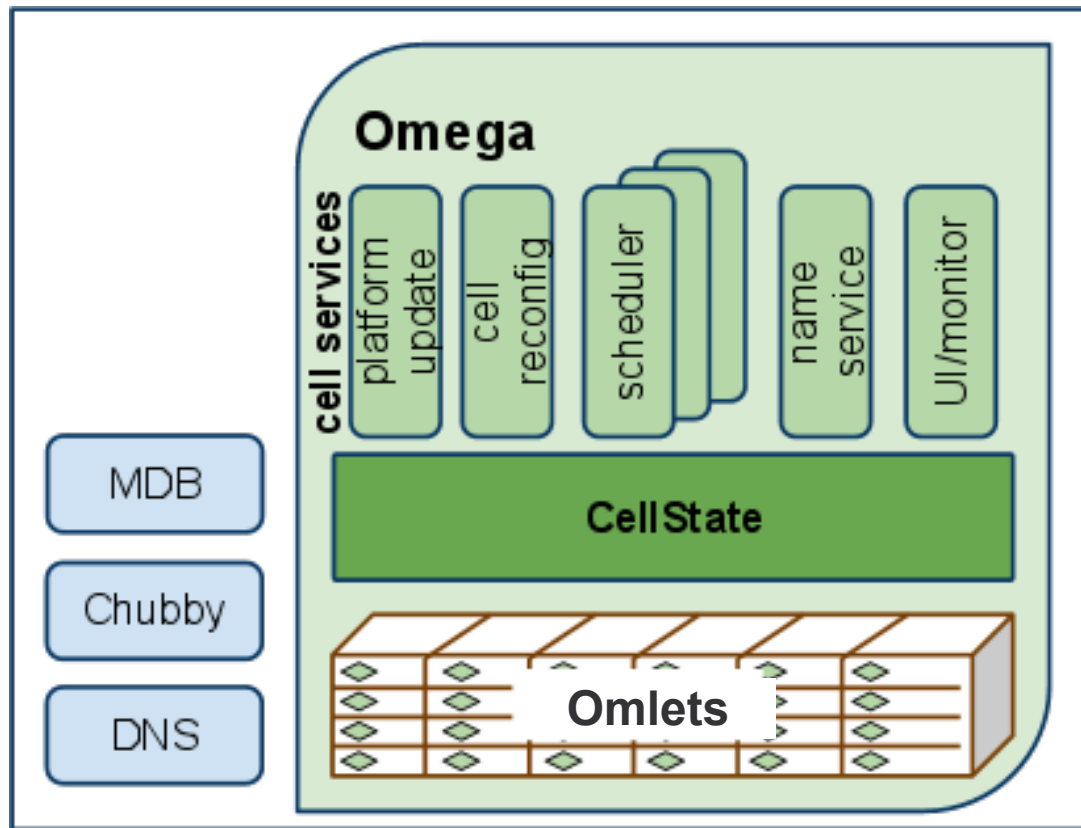




The “second system” ...

- Main *user* goals: predictability & ease of use
- Main *team* goal: flexibility

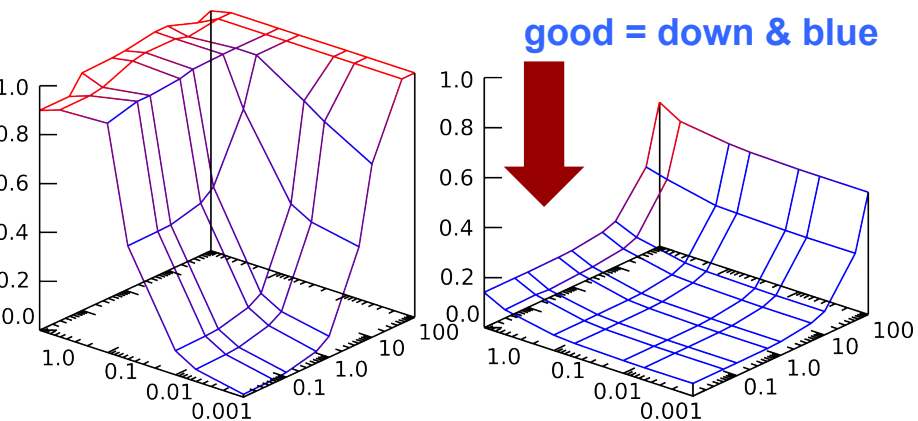
Omega is currently being deployed + prototyped



CellState:

- minimum necessary data
- no policies
- just enforces invariants

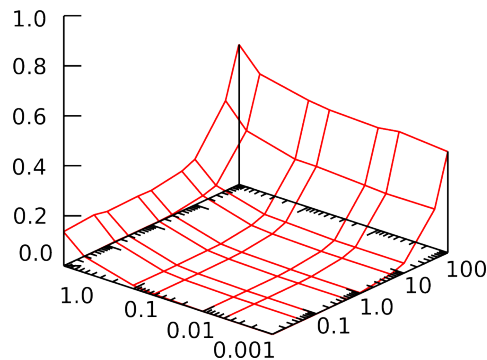
Calendar: everything is a scheduling event!



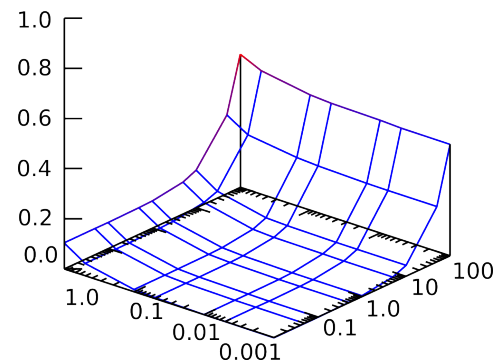
one path

fast batch path

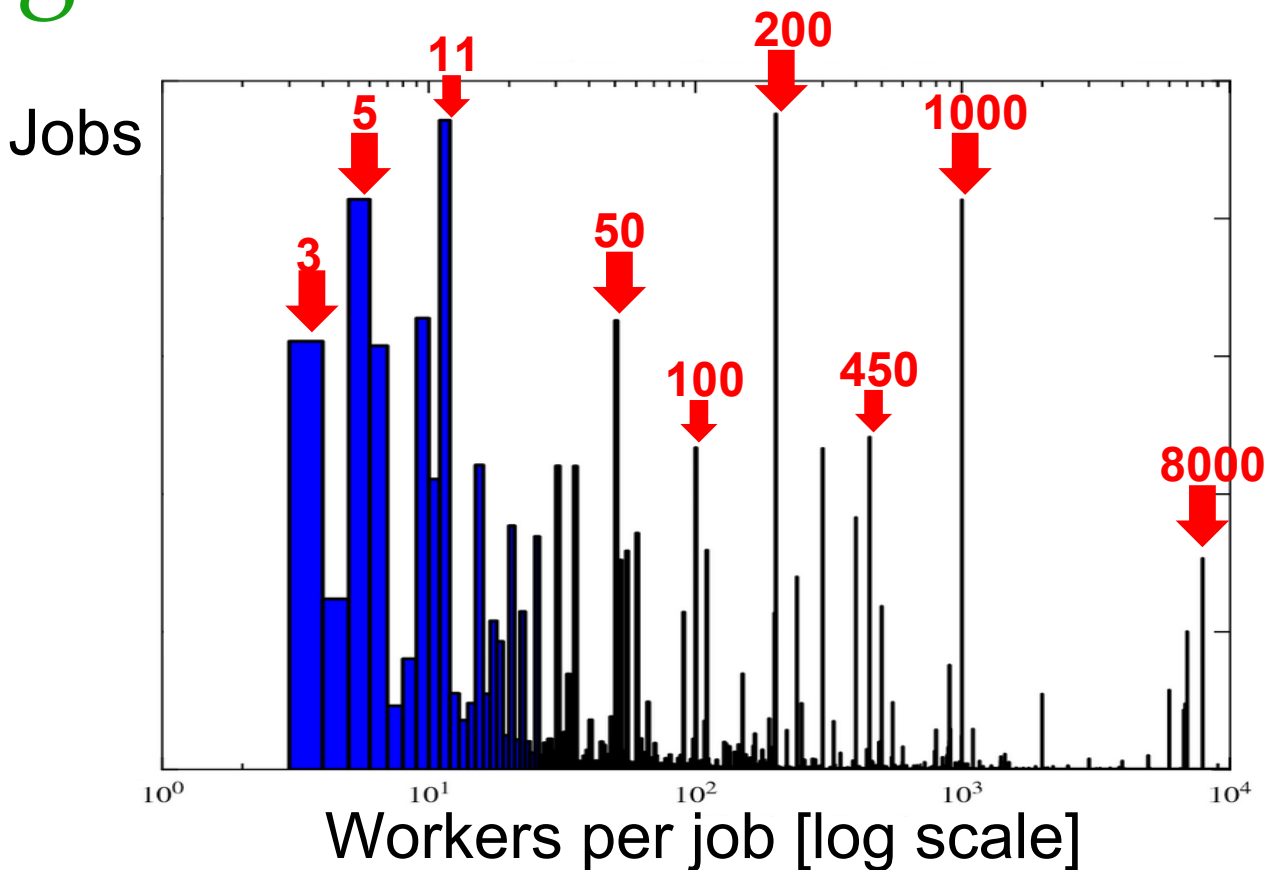
Monolithic scheduler

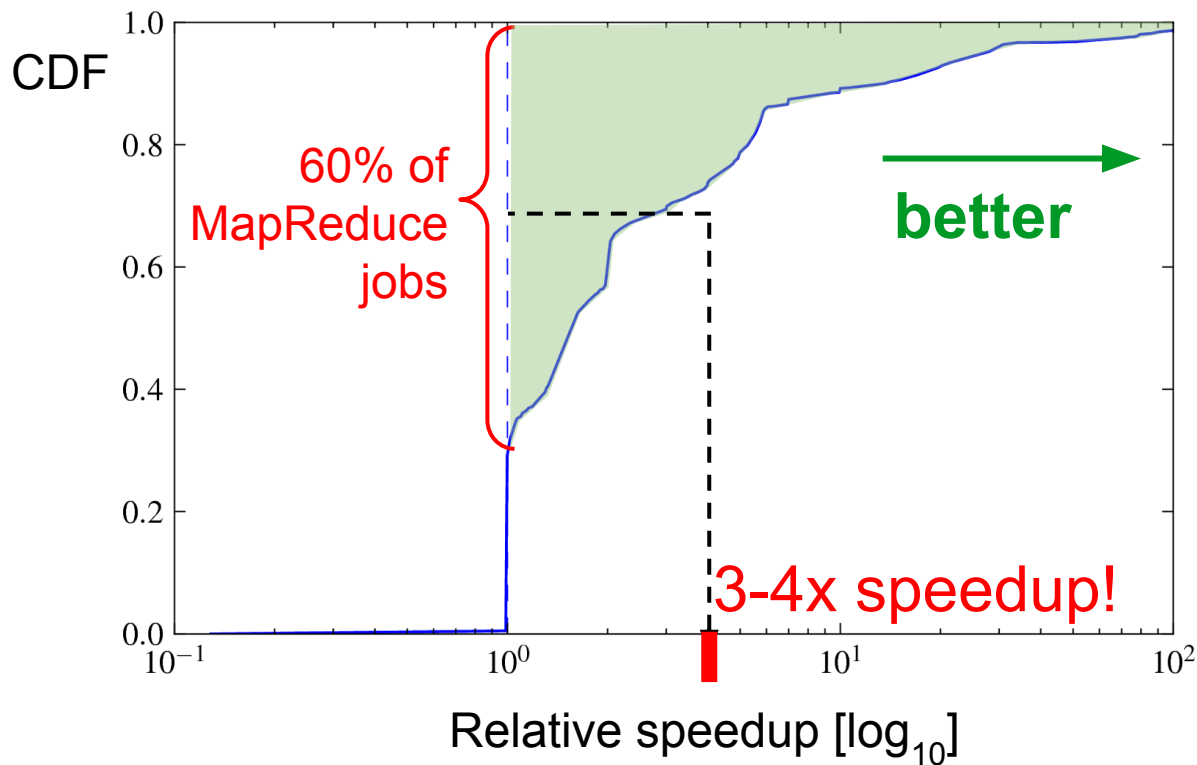


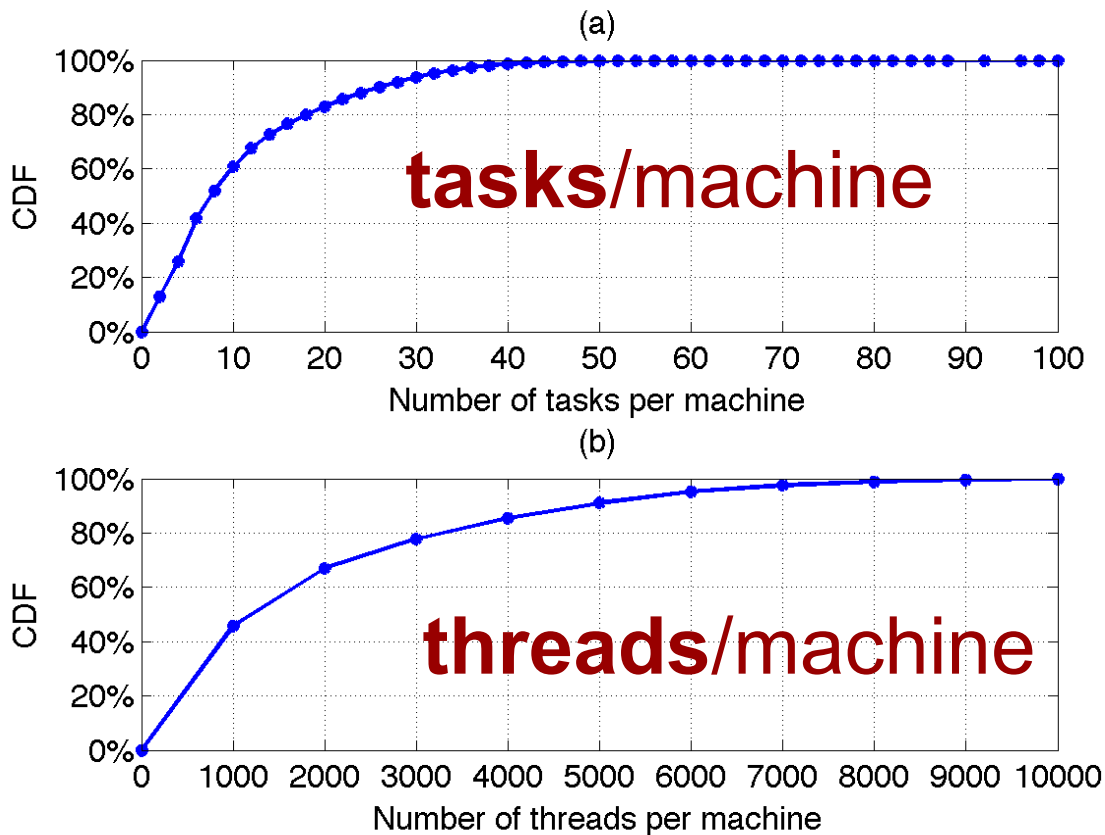
**UCB
Mesos**

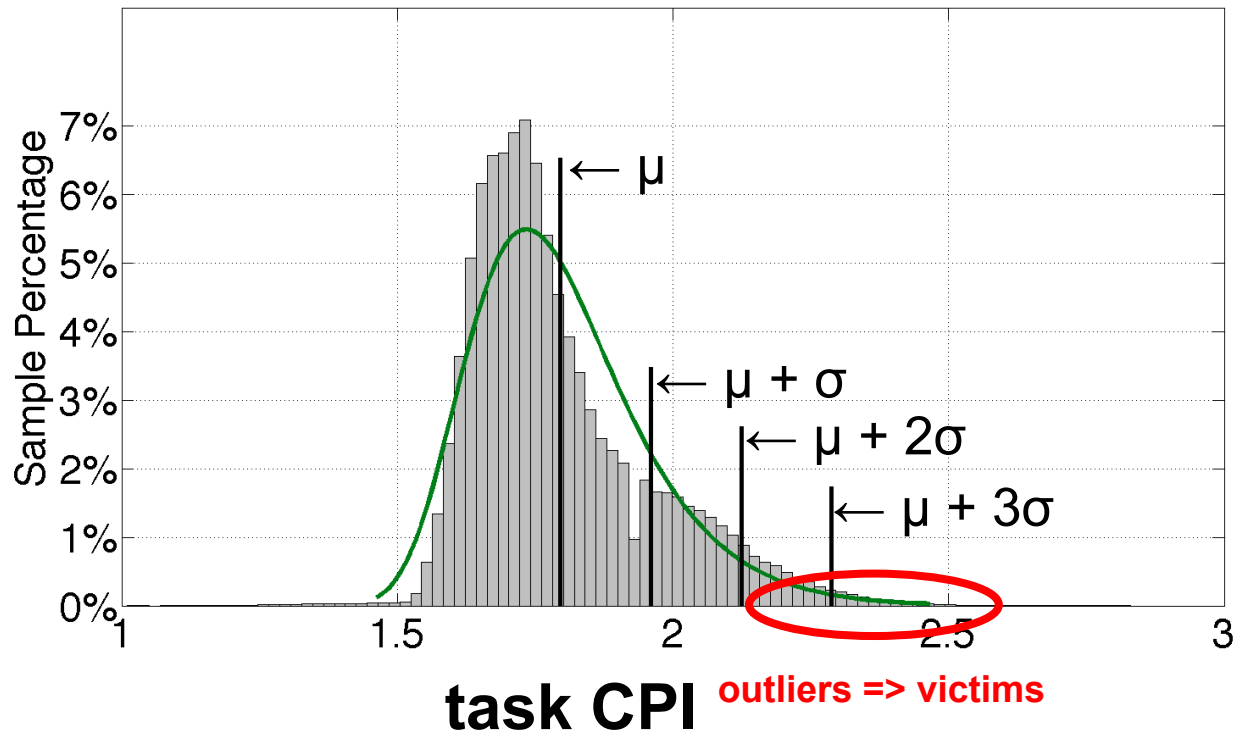


Omega



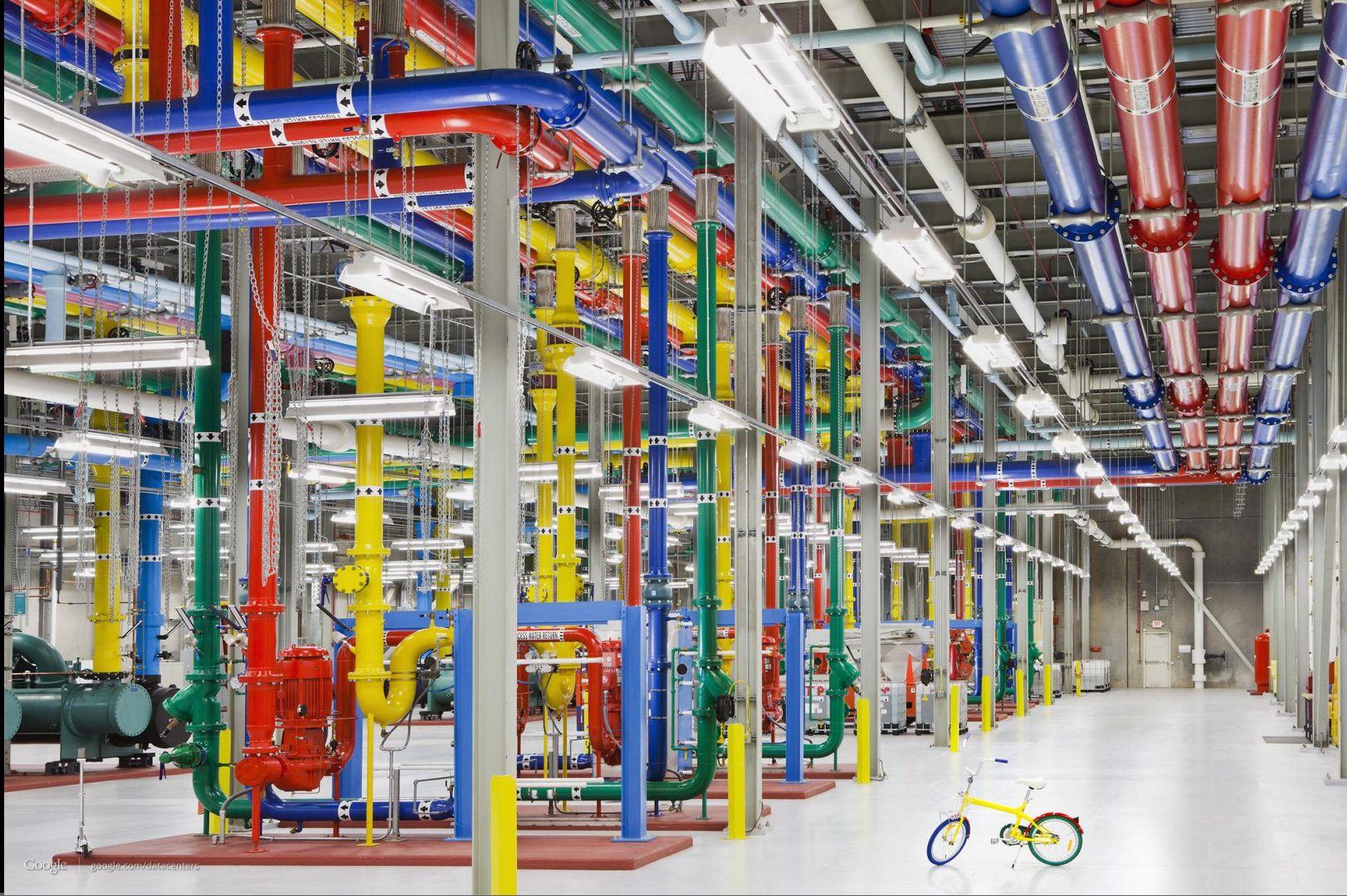






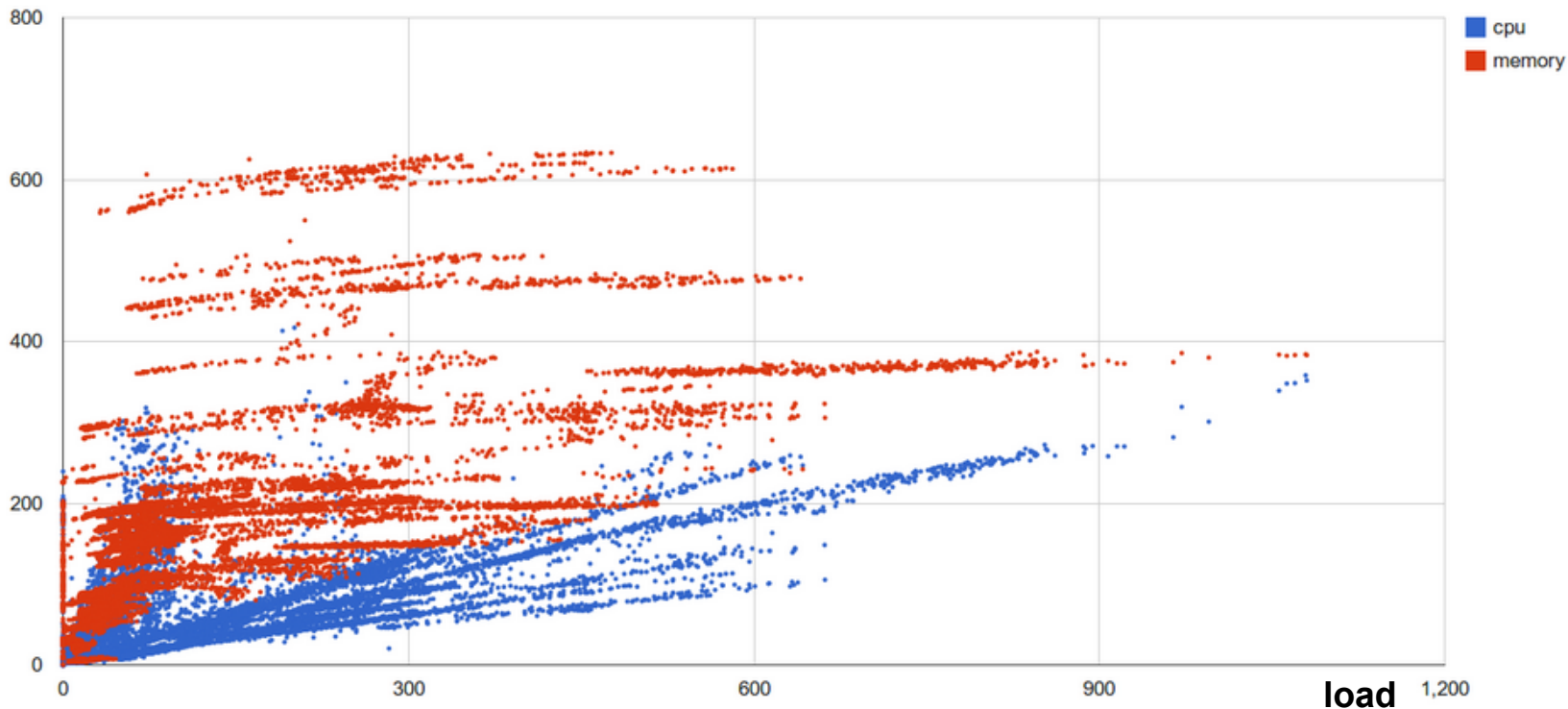
CPI data for all the tasks in a job

- per-cluster
- per-platform
- mean (μ)
- stddev (σ)



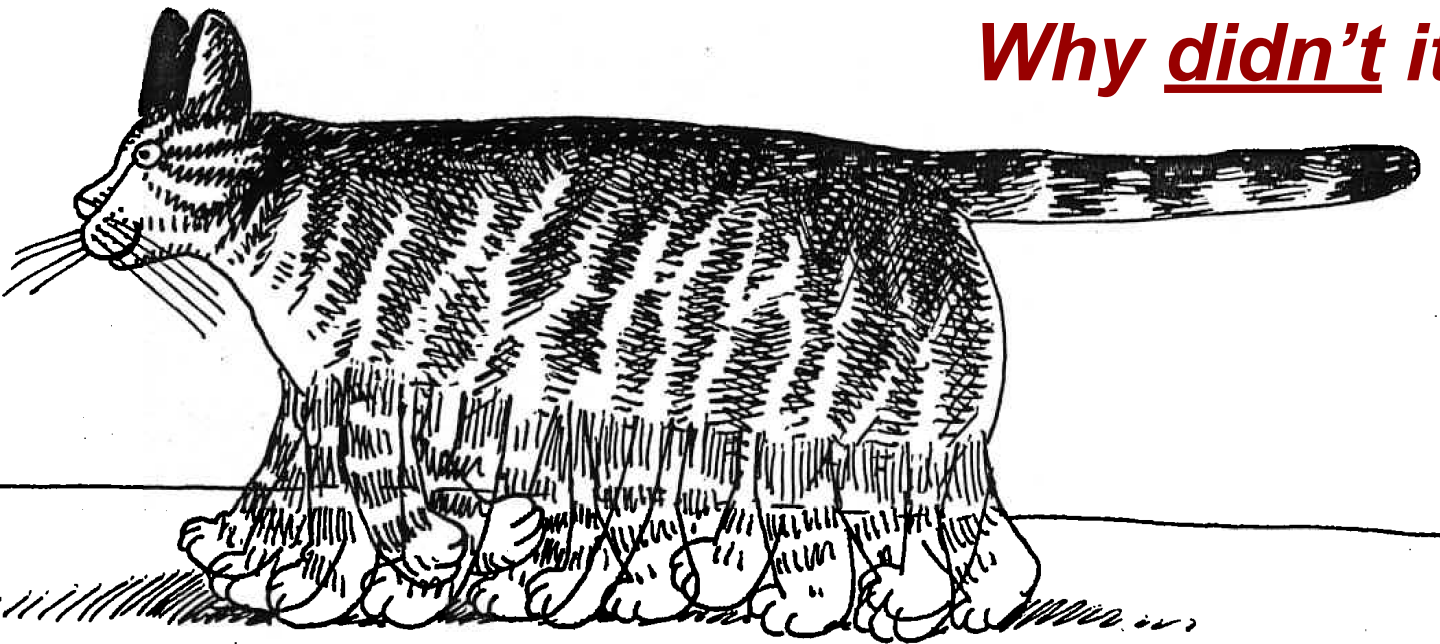
Workload variation

CPU and RAM usage as a function of offered load



“Here’s a binary ... run it”

Why didn't it run?



Open issues: *smart explanations*

It's 3am and your pager goes off

- are we in trouble?
- are we about to get into trouble?

→ what should you do about it?

Large-scale systems are exciting beasts

- scale → failures
- QoS → new designs
- configuration may be **the** next big challenge

DO COOL THINGS THAT  MATTER

We're hiring! www.google.com/jobs

johnwilkes@google.com