# Borg: the Next Generation

Muhammad Tirmazi,[1] Adam Barker,[2] Nan Deng, Md E. Haque, Zhijing Gene Qin,
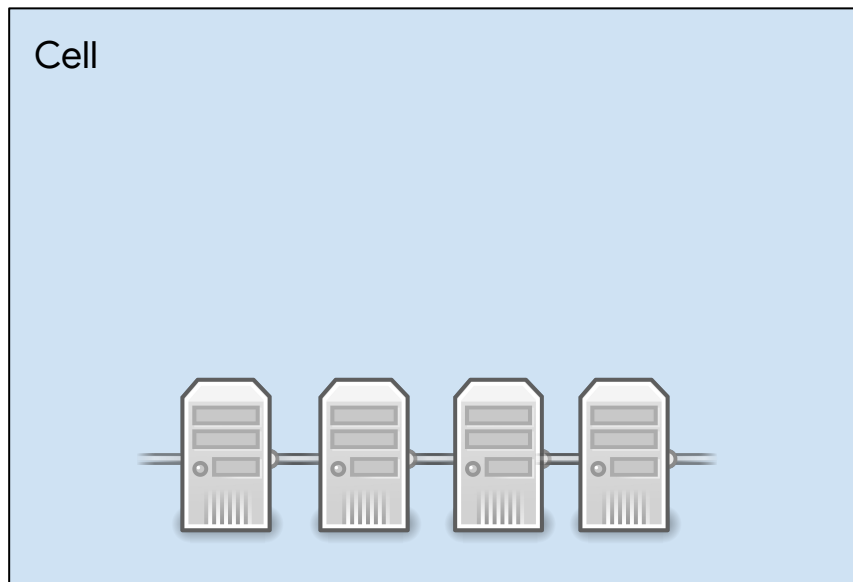Steven Hand, Mor Harchol-Balter,[3] John Wilkes

[1] Harvard University and intern at Google; [2] University of St Andrews and visiting researcher at Google; [3] CMU and visiting researcher at Google
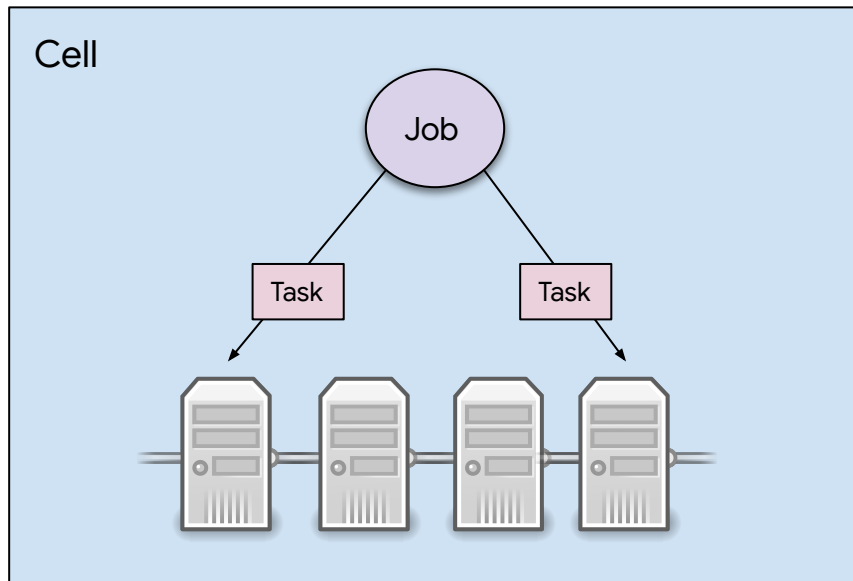
# Borg

Google's internal cluster manager.

**Cell**: a set of machines managed by Borg as one unit.



Cell

# Borg

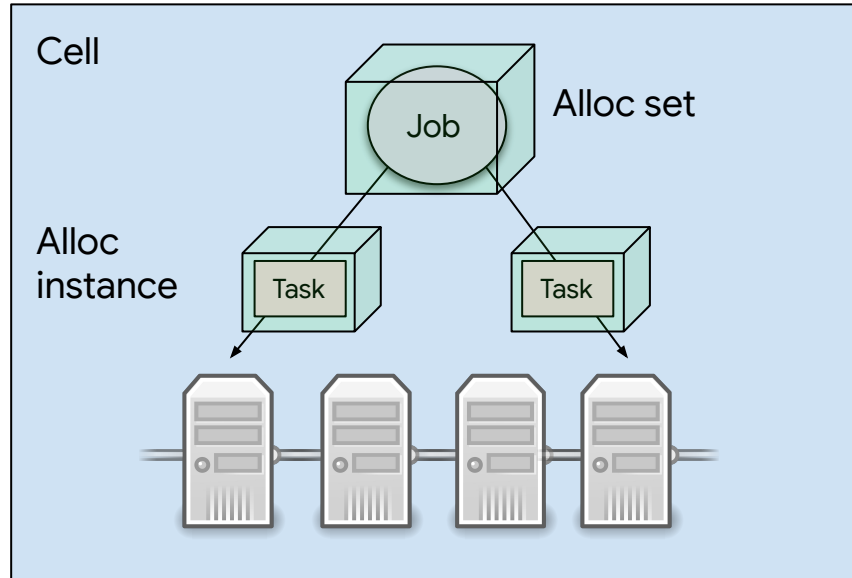Users submit work in the form of **jobs**

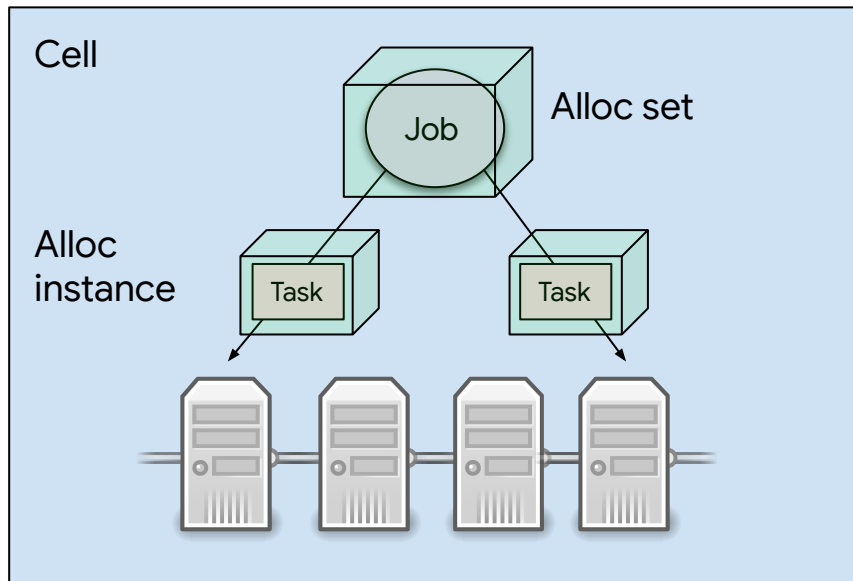each of which contains one or more **tasks**.

# Borg

A job may run in an **alloc set**

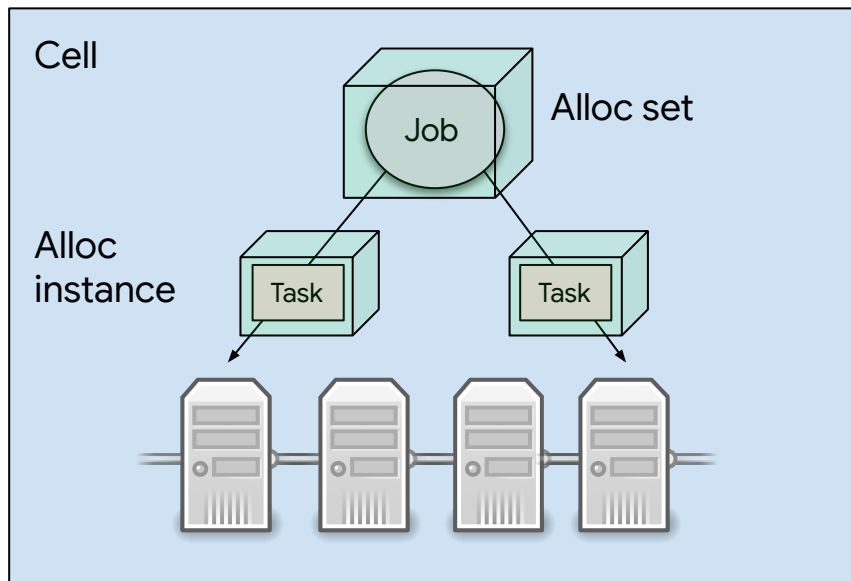     making each of its tasks run in an **alloc instance**

# Borg

Jobs have **tiers**: production, mid, best-effort batch, free.

# Borg

More info: "Large scale cluster management at Google with Borg" (EuroSys '15)
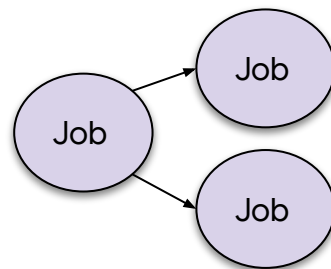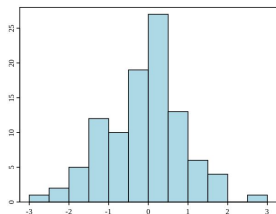
# Borg **traces**

A single Borg trace describes the **workload** in a Borg cell:

- {Jobs, tasks}, {alloc sets, alloc instances}
  - arrivals and departures: submit, update, finish
  - scheduling decisions: place, evict

- Resource allocations and usage

2011 trace: 1 cell from May, 2011

# Borg traces: what's **new**
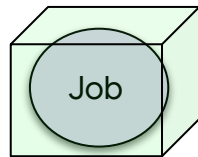
**2019 trace: 8 cells  for May 2019**

- ~96k machines in 3 continents
- CPU usage histograms
- Job-parent information
- Autopilot (see companion paper in session 5)

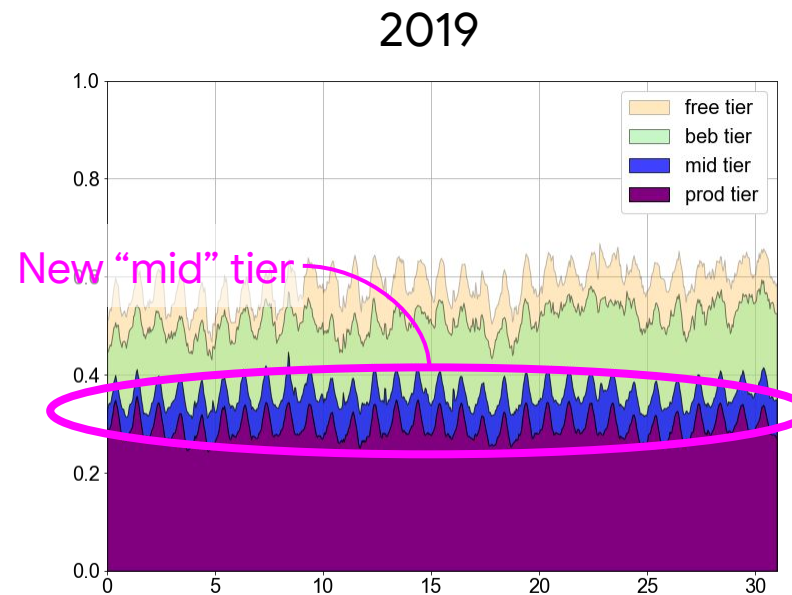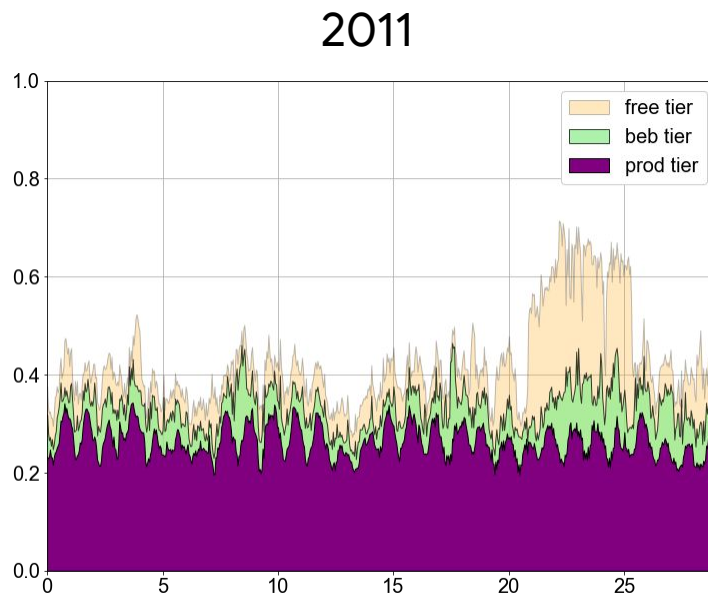**github.com/google/cluster-data**

# Resources used by jobs

**Two metrics:**



- **Resource used by job**

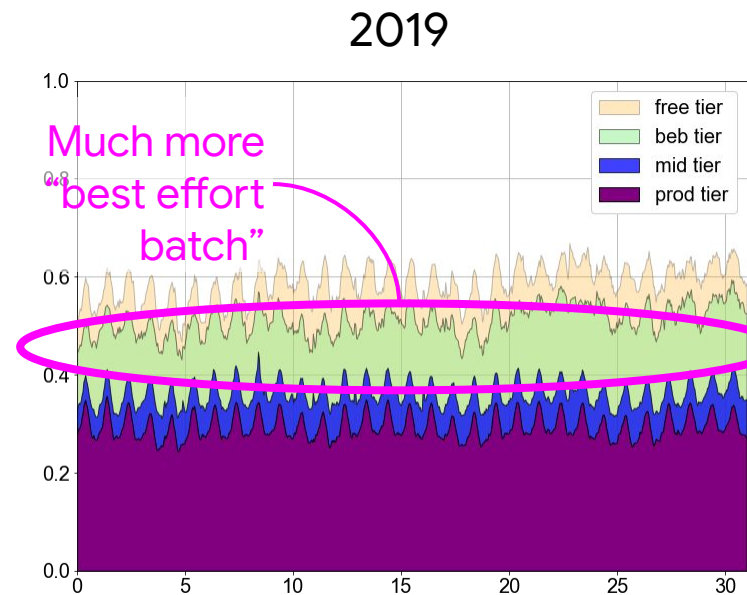- **Resource allocated to job**

# Compute **used** by jobs



Fraction of
cell capacity

2011

2019

New "mid" tier

Time (days)

# Compute **used** by jobs



Fraction of cell capacity

2011

2019

Much more "best effort batch"

Time (days)

# CPU + memory **used** by jobs

# Compute allocated to jobs

Fraction of
cell capacity

### 2011



| free tier |
| beb tier |
| prod tier |

### 2019



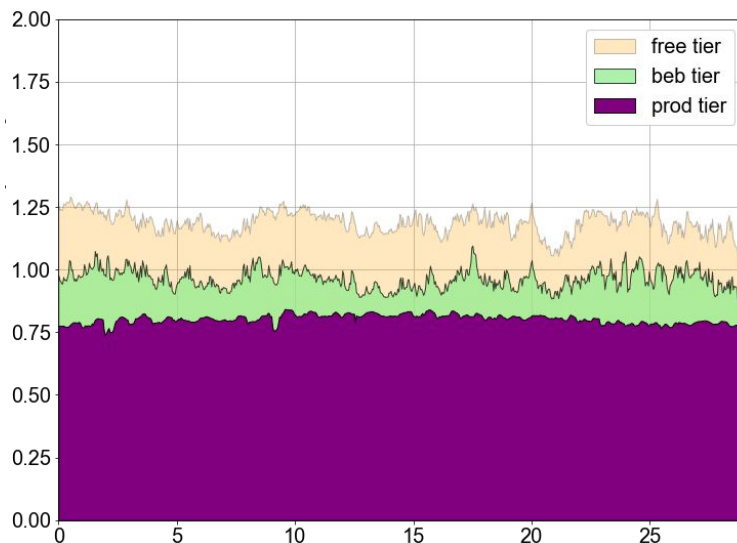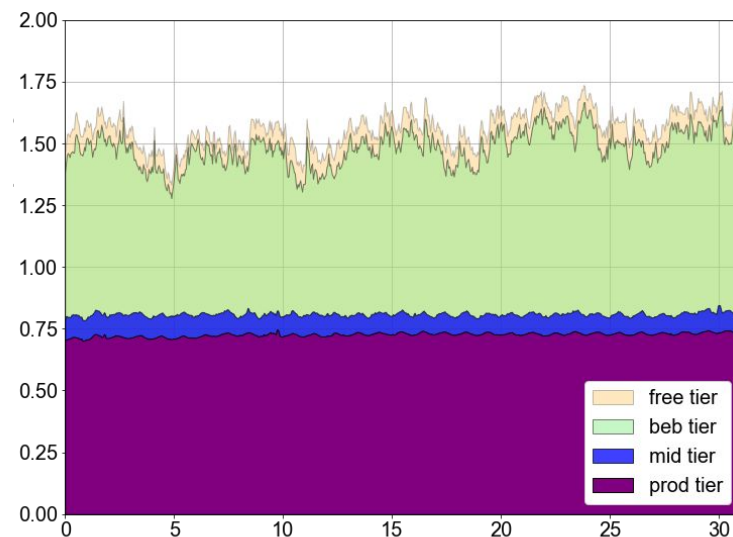| free tier |
| beb tier |
| mid tier |
| prod tier |

Time (days)

# **Memory** allocated to jobs
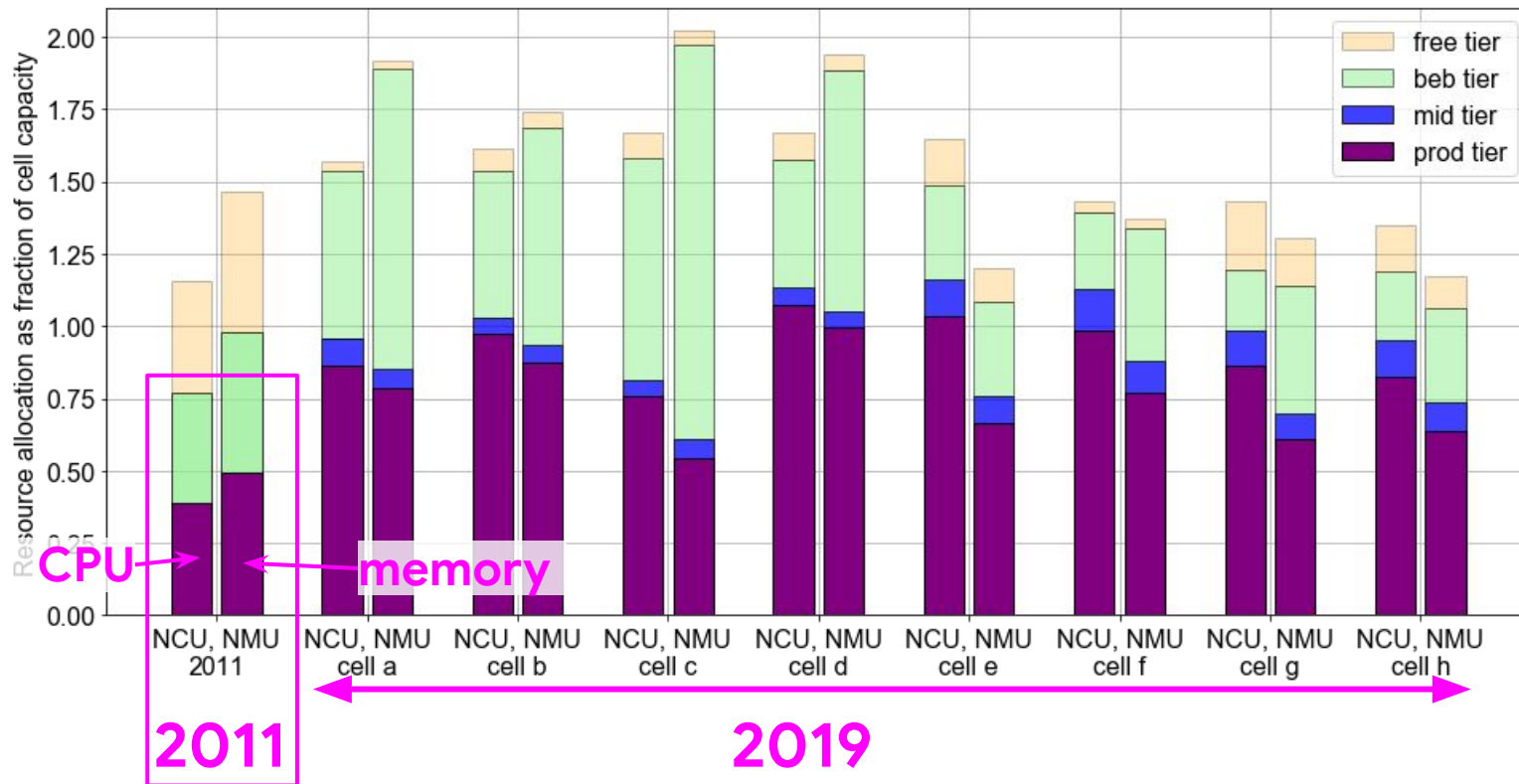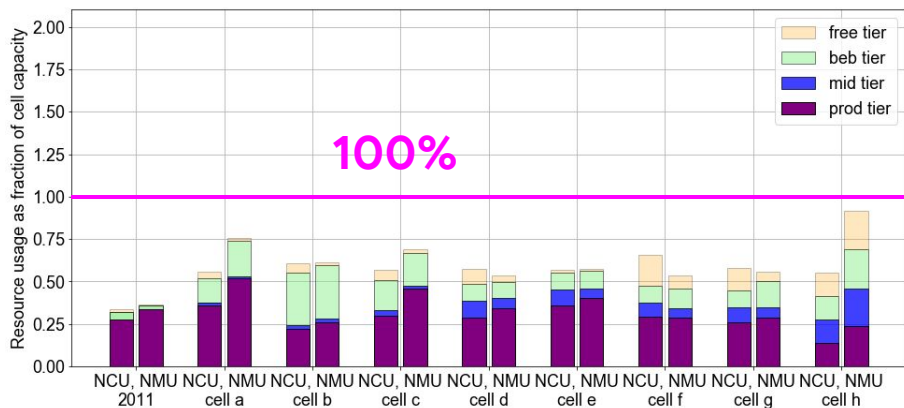


2011

2019

Fraction of
cell capacity

Time (days)

# CPU + memory **allocated** to jobs
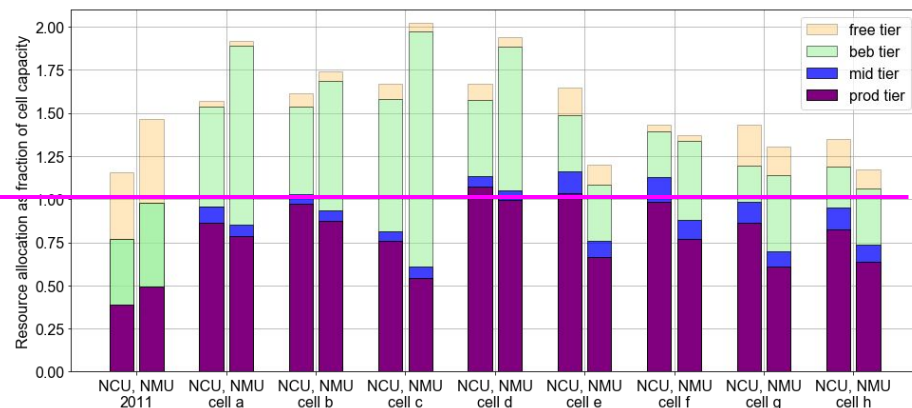
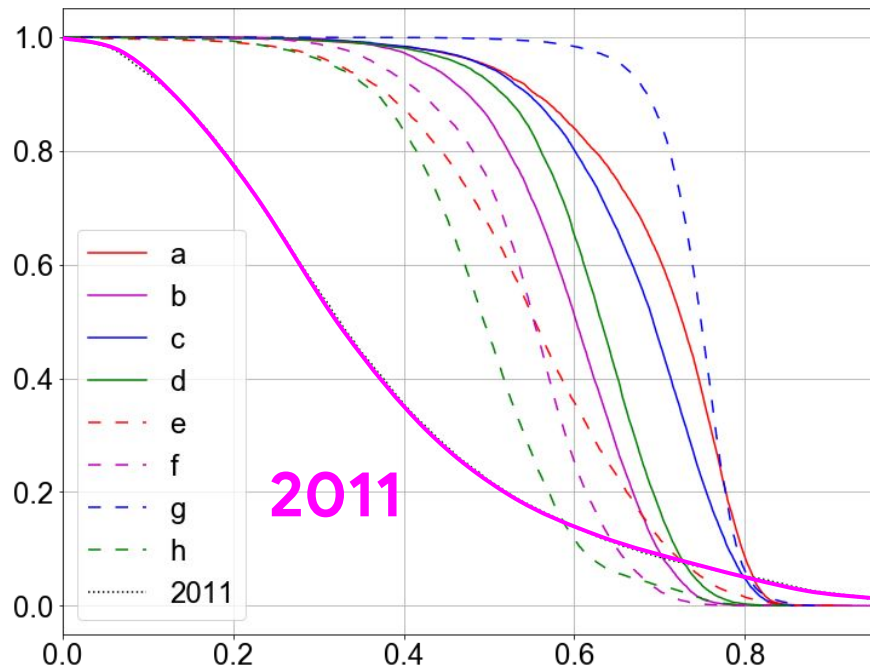# CPU + memory **used** vs **allocation**

## Resources **used** by jobs



## Resources **allocated** to jobs

# **Machines** used by jobs

P(utilization > x)



x - utilization

# Machines used by jobs

P(utilization > x)

Median machine **in 2011**:
~ 30% utilized

Median machine **in 2019**:

**50 - 77%** utilized



Median utilization is higher in 2019
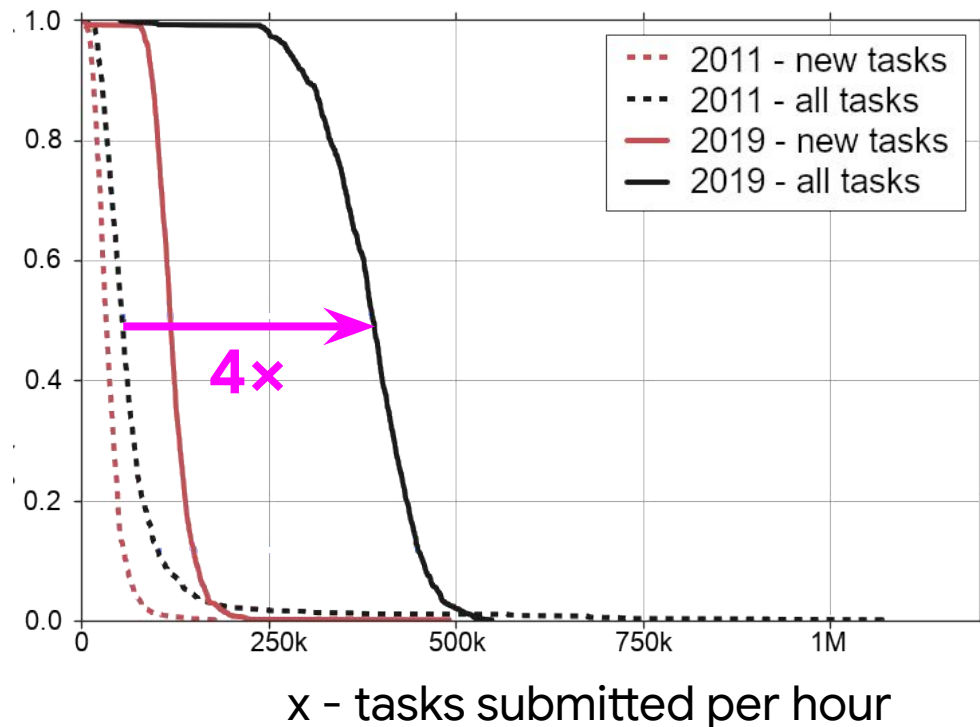
x - utilization

# Scheduler load is evolving

P(tasks submitted > x)

Scheduler load today:
~ **4 times** higher



**4×**

Legend:
- 2011 - new tasks
- 2011 - all tasks
- 2019 - new tasks
- 2019 - all tasks

x - tasks submitted per hour

# Job usage has **VERY** high variability

$C^2$ = variance / mean$^2$
for CPU-hours and memory-hours

- CPU-hours of UNIX jobs (1996): $C^2 \approx$ **50**
- CPU-hours of supercomputing jobs (2005): $C^2 \approx$ **250**
- CPU-hours of Google Borg jobs (2011): $C^2 \approx$ **8400**

# 2019 Google Borg trace: **23k**

# Hogs and mice

Largest 1% of jobs: **hogs**
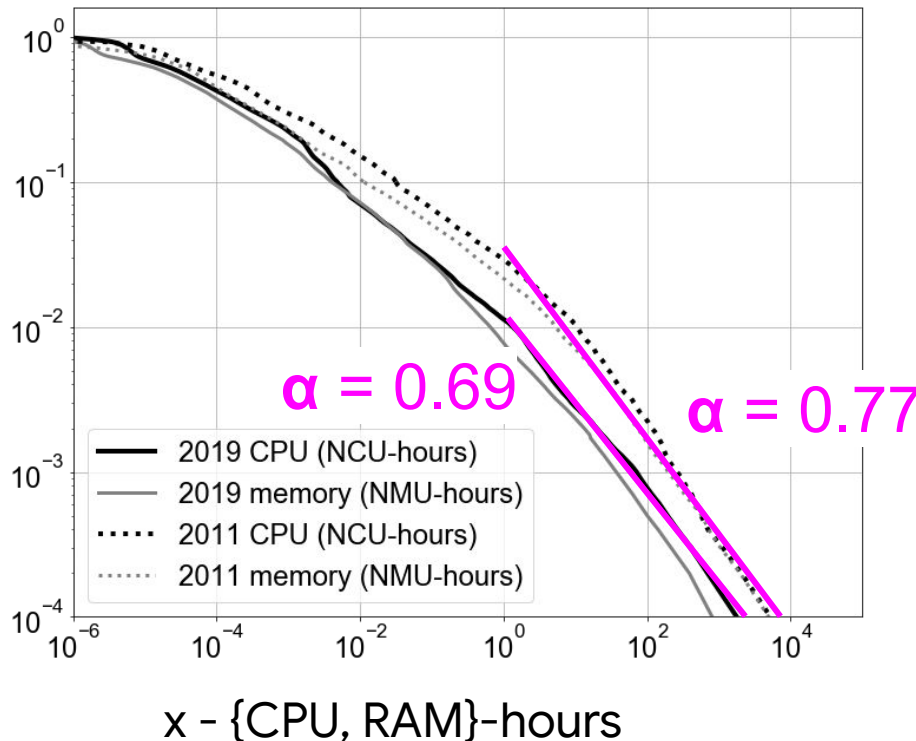
Remaining 99%: **mice**

Fraction of resources consumed by

- Prior work: 50%
- Google, 2011: 97.3%
- Google, 2019: **99.2%**

# Job usage is heavy tailed

Fraction of jobs where:
{CPU, RAM}-hours > x

Even more
heavy-tailed!



α = 0.69

α = 0.77

2019 CPU (NCU-hours)
2019 memory (NMU-hours)
2011 CPU (NCU-hours)
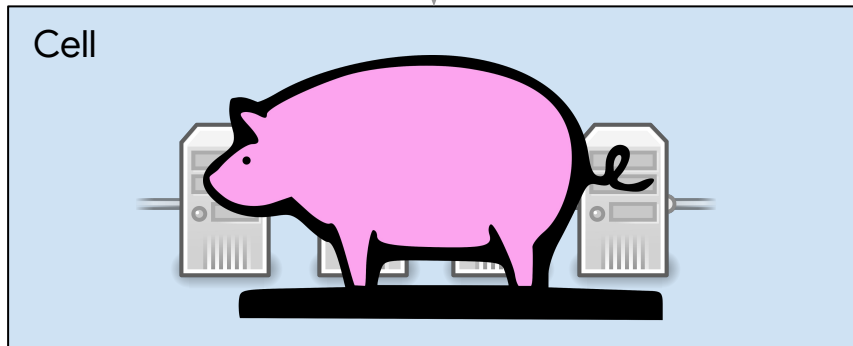2011 memory (NMU-hours)

x - {CPU, RAM}-hours

*Extremely*
heavy tailed

# Implications for **scheduling**

Since Google's workload has high $C^2$



Hogs can fill all the resources!

Cell

# Summary

- New Borg workload trace:
    - 8 cells for month of May 2019
    - 2.4TB data accessed via BigQuery
    - **github.com/google/cluster-data**

- Workload and machine utilization have increased

- Disparity between hogs and mice more extreme than any other reported trace
    - largest 1% of jobs consume >99% of resources