



Large-scale Cluster Management at Google with Borg

By Abhishek Verma, Luis Pedrosa, Madhukar Korupolu, David Oppenheimer, Eric Tune and John Wilkes

Presented by Matthias Erdmann, Simon Buchacher, and Spencer Johann Thellmann

Overview

- I. Borg
- II. Evaluation
- III. Opinion

The background is a dark, textured surface. On the left, there are several overlapping geometric shapes: a light blue parallelogram, a dark grey parallelogram, and a medium blue parallelogram. The rest of the background is filled with a complex network of thin, curved lines in shades of blue and grey, connecting small circular nodes. Some nodes are highlighted with a white center. The overall aesthetic is technical and futuristic.

A Brief Introduction to Borg

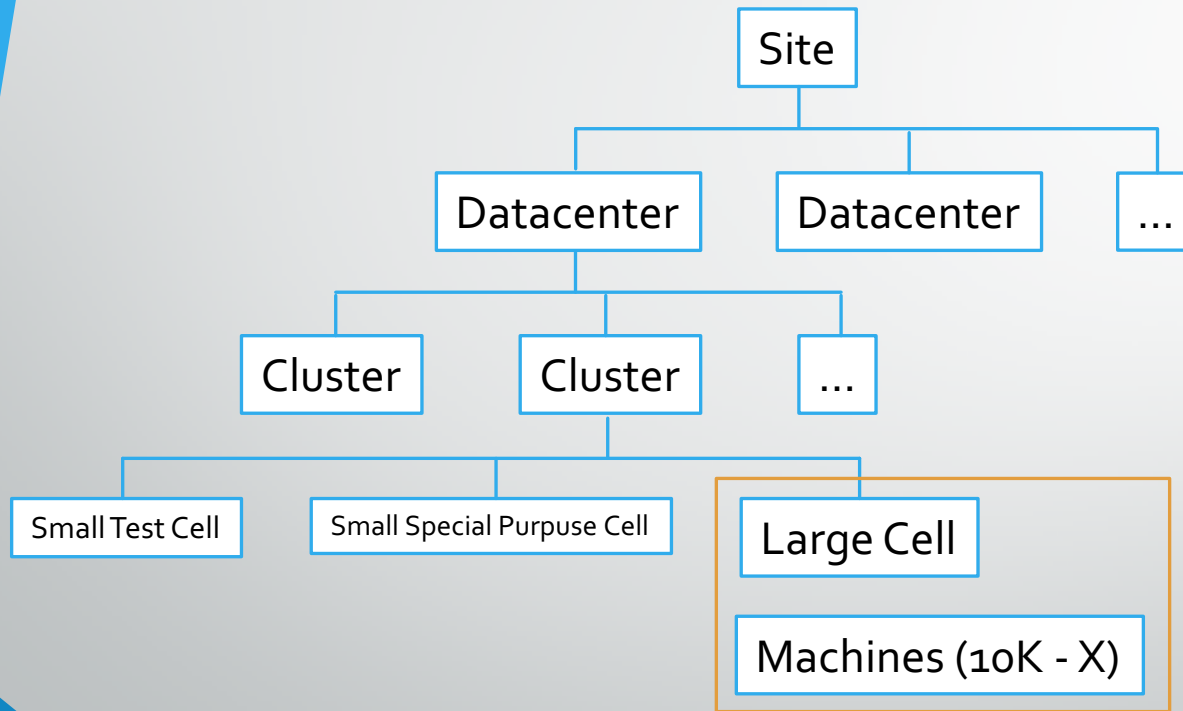
Analogy -

- Lectures
- Practicals
- Papers
- (essential & recommended) Readings
- Exams
- Free time

March 2021

	Mon 1	Tue 2	Wed 3	Thu 4	Fri 5
07:00	0000 Programming Practical	0000 Programming Practical	0000 Programming Practical	0000 Programming Practical	0000 Programming Practical
08:00	0000 C30044 - info Via Recordings	08:00 Breakfast	08:00 Breakfast	08:00 Breakfast	08:00 Breakfast
09:00	09:00 Sean Meeting	09:00 0000 Programming Practical	09:00 C30044 - Distributed Systems	09:00 0000 Practice presentation	09:00 C30044 - info - Sean...
10:00	10:00 0000 Programming Practical	10:00 C30044 - Distributed Systems	10:00 0000 project		10:00 0000 project
11:00	11:00 C30044 - SE Practice	11:00 C30044 - SE Practice		11:00 Sports	
12:00	12:00 Sports	12:00 Sports	12:00 Sports	12:00 C30044 - info - Sean...	12:00 Sports
13:00	13:00 Lunch Break	13:00 Lunch Break	13:00 Lunch Break	13:00 Lunch Break	13:00 Lunch Break
14:00	14:00 0000 Reading	14:00 0000 Reading	14:00 0000 Reading	14:00 0000 Reading	14:00 0000 project
15:00				15:00 Dinner Break	
16:00				16:00 0000 project	
17:00					
18:00					
19:00	19:00 Sports	19:00 Sports	19:00 Sports	19:00 Sports	19:00 Sports

Google infrastructure hierarchy



Cells: heterogenous workload



Long-running services

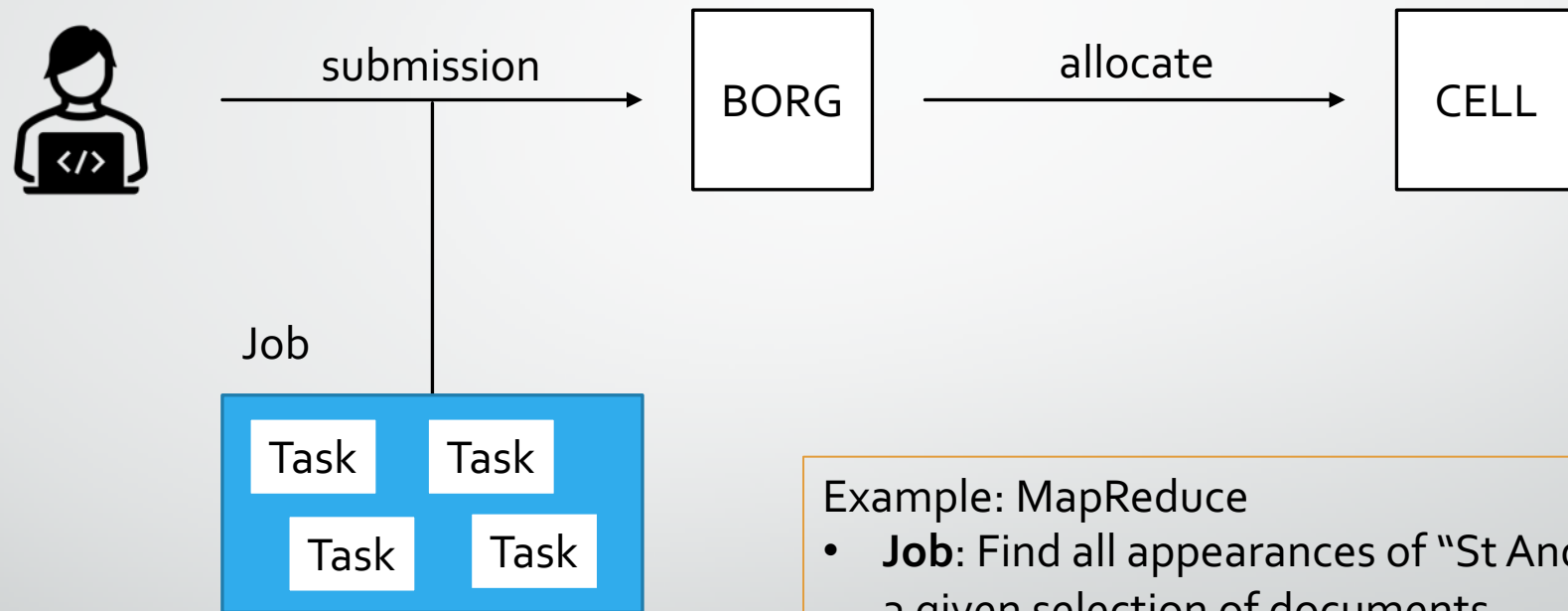
- (mostly) prod
- Constant availability

Batch jobs

- (mostly) non-prod
- Short to long jobs (seconds - days)

→ Variations in mix of prod and non-prod jobs

USP: Decoupled scheduling



Example: MapReduce

- **Job:** Find all appearances of “St Andrews” in a given selection of documents
- **Tasks:** Replicas of Mapping and Reducing Tasks + Master Task

Allocation of jobs

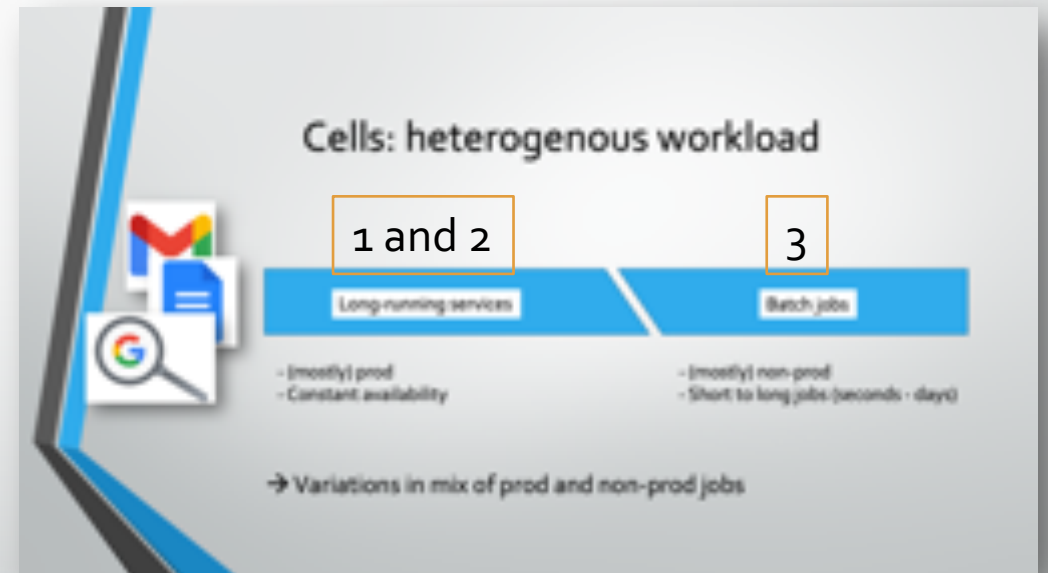
Alloc: resources on a machine that are obtained to run a job

- **Heterogeneity** in clusters: size, processor type, performance, capabilities
- **Importance of job** (→ priority)
- **Capacities** (→ quota)

Priority

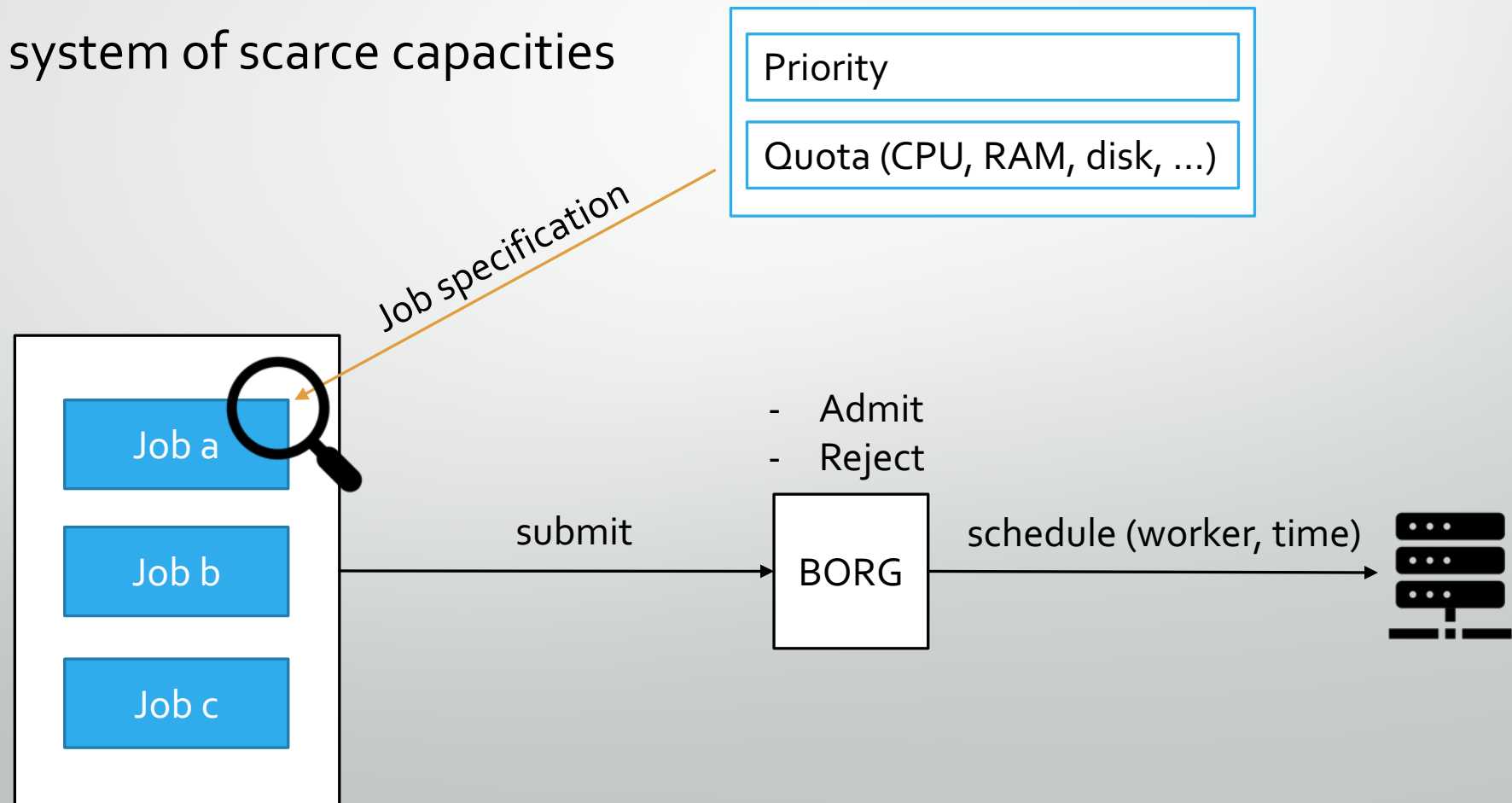
Non-overlapping priority bands:

1. Monitoring
2. Production
3. Batch
4. Best effort (testing)

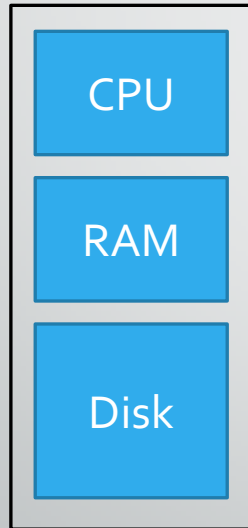


Admission control – an internal auction

→ In a system of scarce capacities



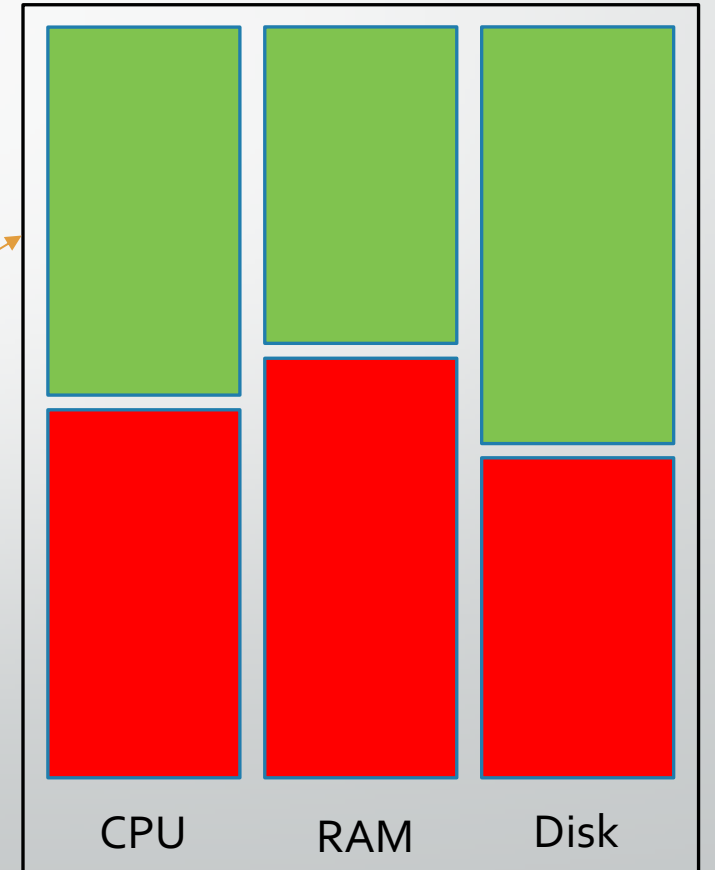
Scheduling jobs...



Job with sufficient **priority**

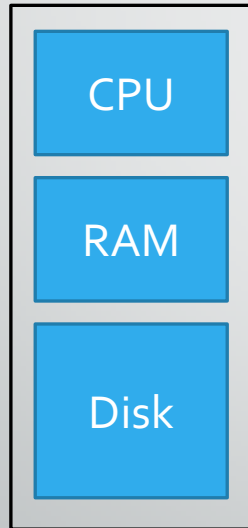


Capacity **elicitation**

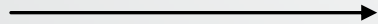


Worker: capacities for a certain **period of time**

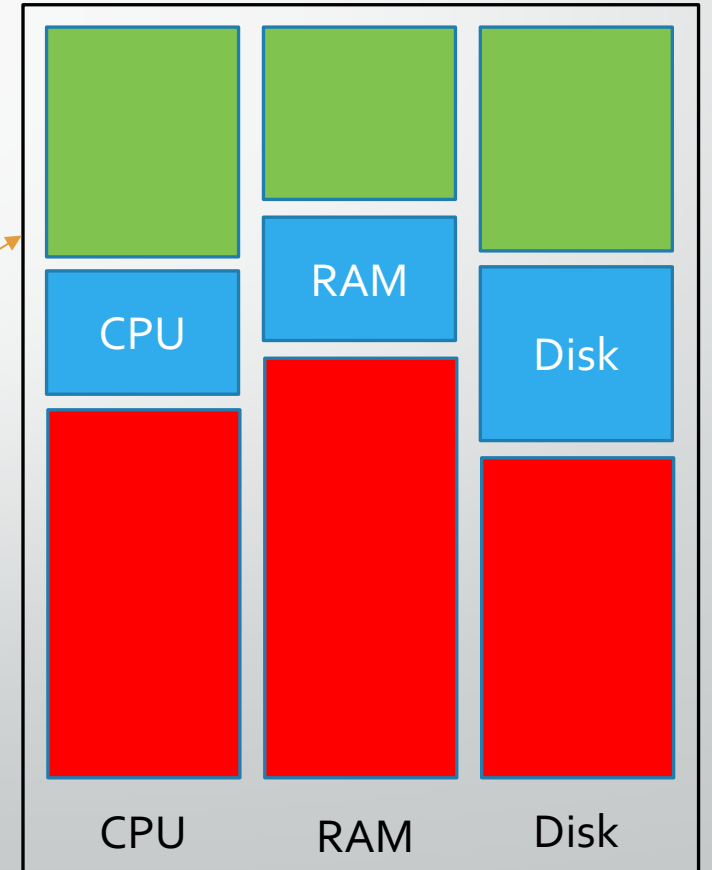
...Scheduling jobs



Job with sufficient **quota**



Capacity **allocation**



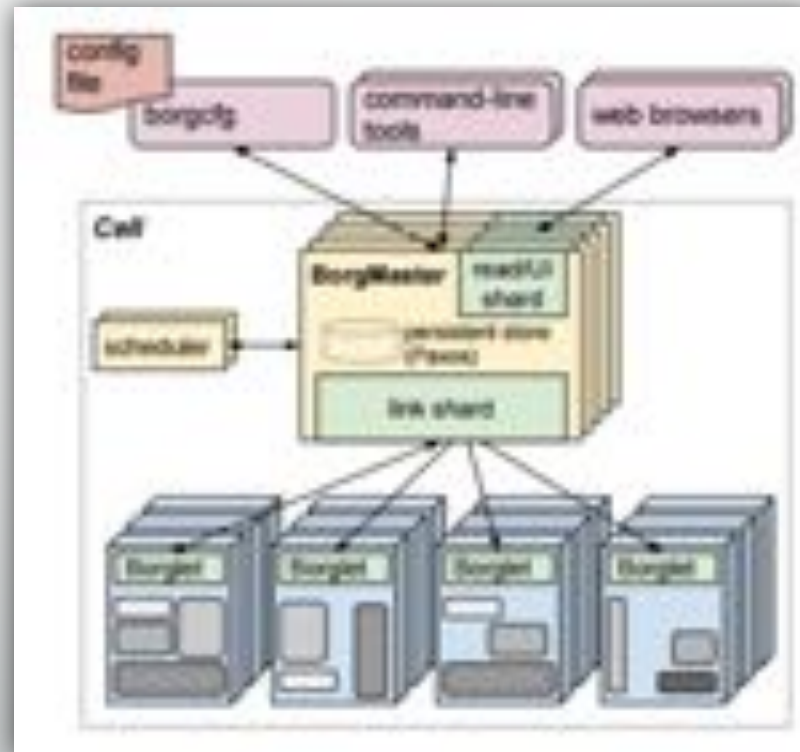
Worker: capacities for a certain **period of time**

User Tradeoffs:

- Specify priority: admission vs cost
- Specify quota: admission vs sufficient capacities

“Cost” – We assume that users receive a budget of real or some “virtual currency” they can spend on submitting jobs.

Borg Architecture



In a nutshell



BORG



The background features a dark blue field with a network graph of nodes and edges. On the left, there are three overlapping geometric shapes: a light blue parallelogram, a dark blue parallelogram, and a grey parallelogram, all slanted. The text "Evaluating Borg" is written in white on the right side.

Evaluating Borg

Questions



How much can you save by sharing machines between jobs?



What happens when you decrease cell size?



What happens when users choose memory and compute freely?



How much can you save by sharing machines between jobs?



How much can you save by sharing machines between jobs?

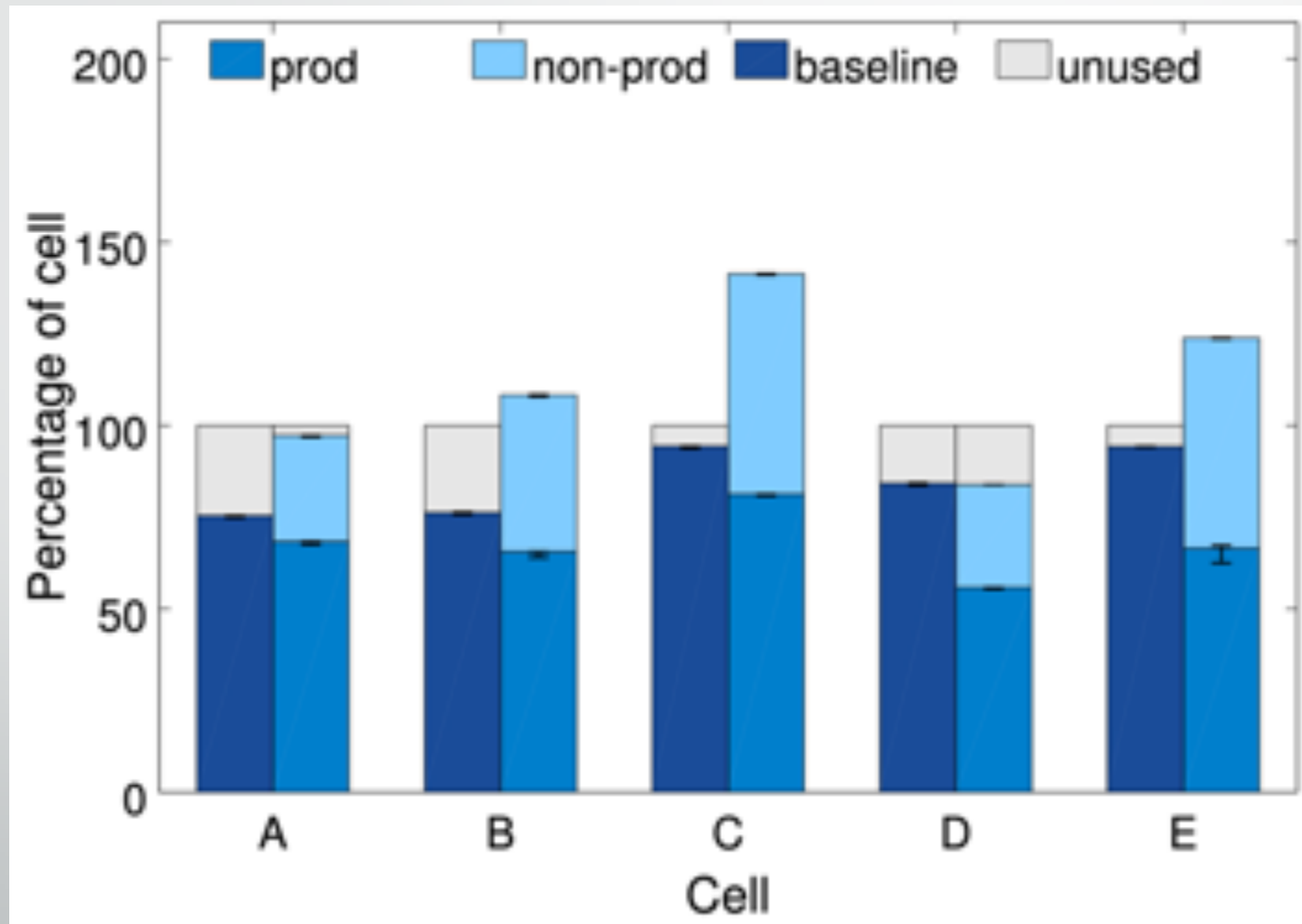
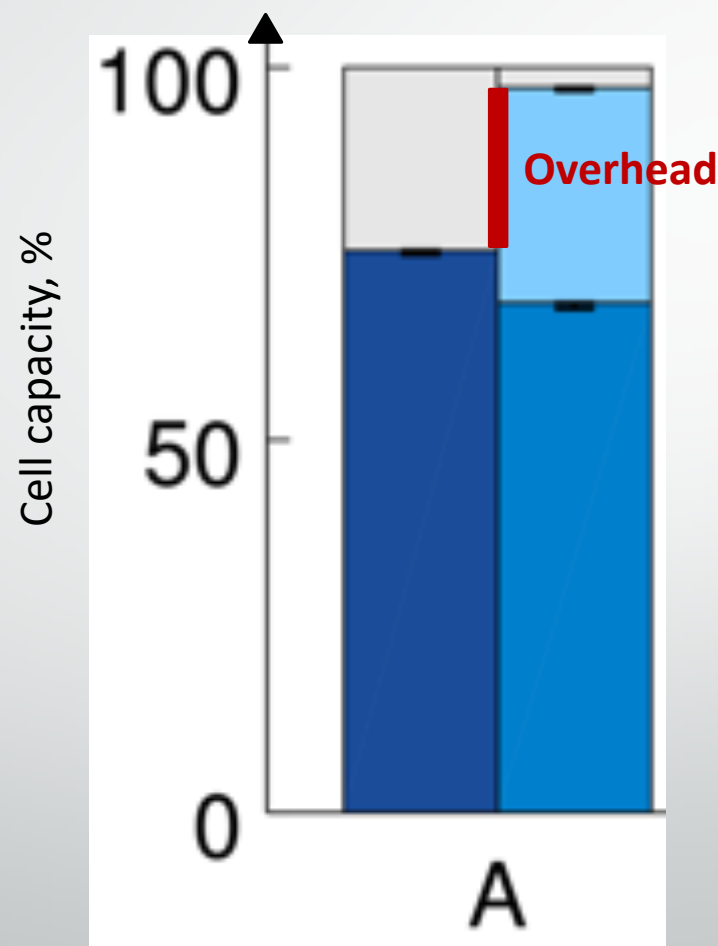


Figure 5(a)



How much can you save by sharing machines between jobs?





How much can you save by sharing machines between jobs?

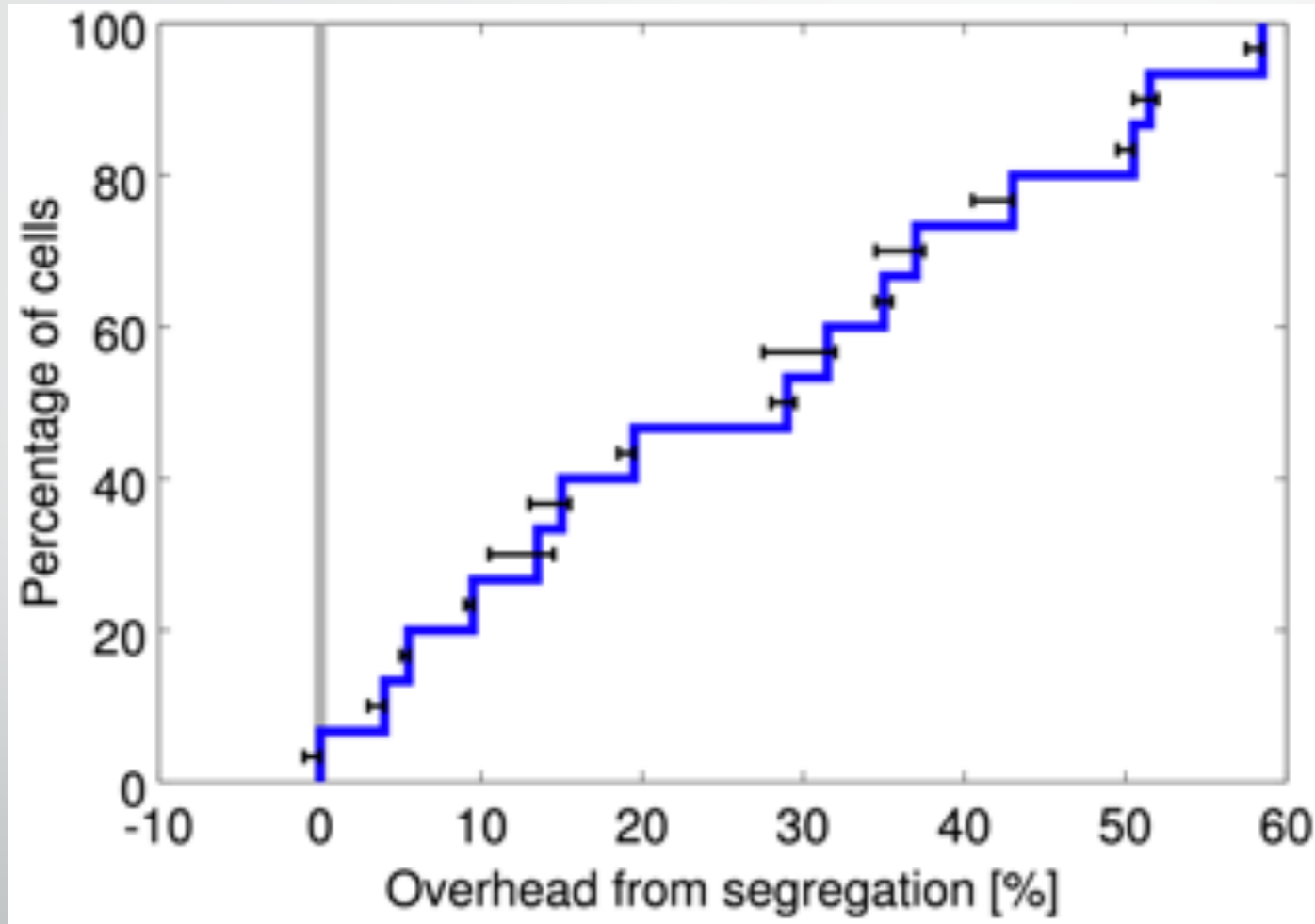


Figure 5(b)



How much can you save by sharing machines between jobs?

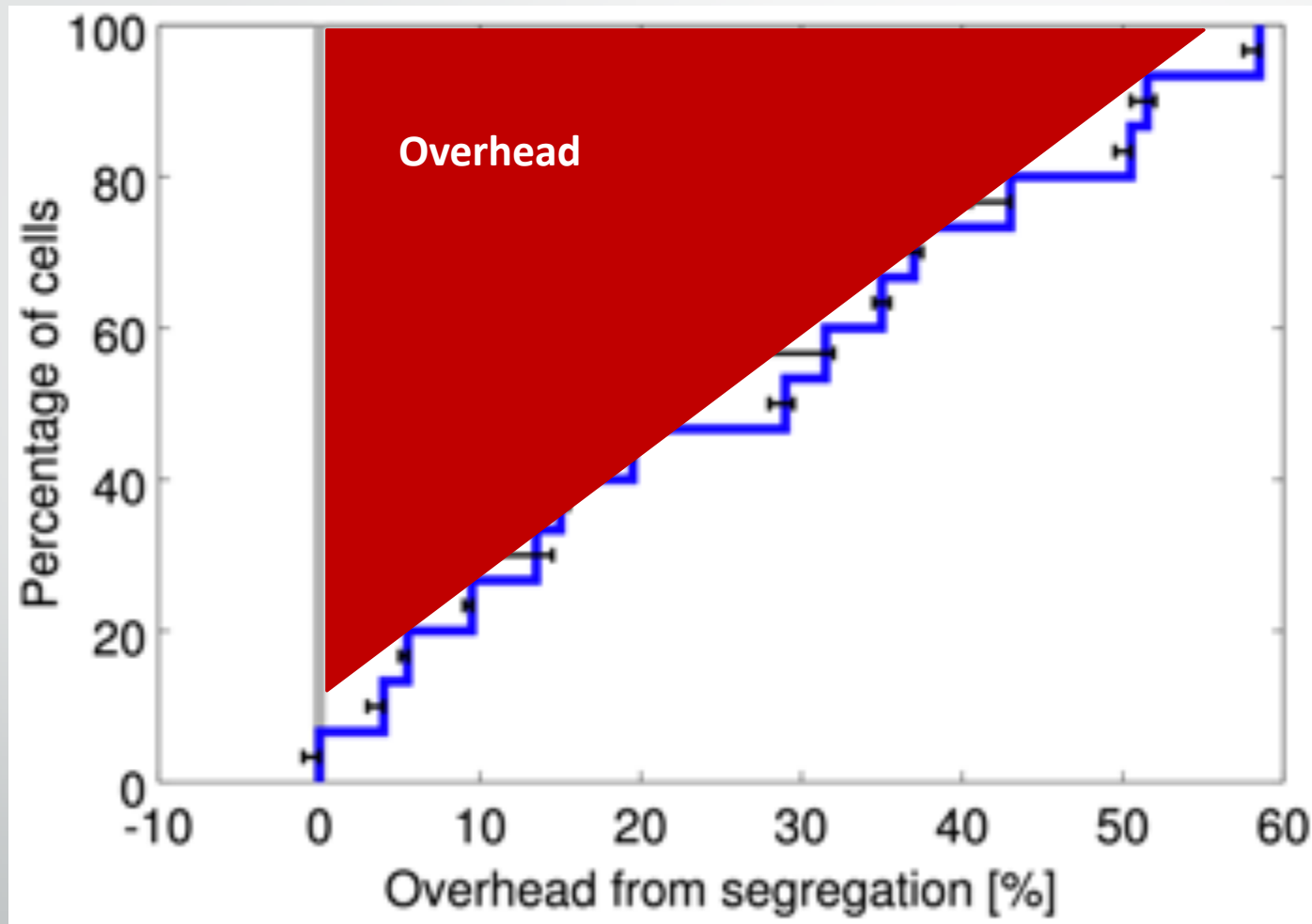


Figure 5(b)



What happens when you decrease cell size?



What happens when you decrease cell size?

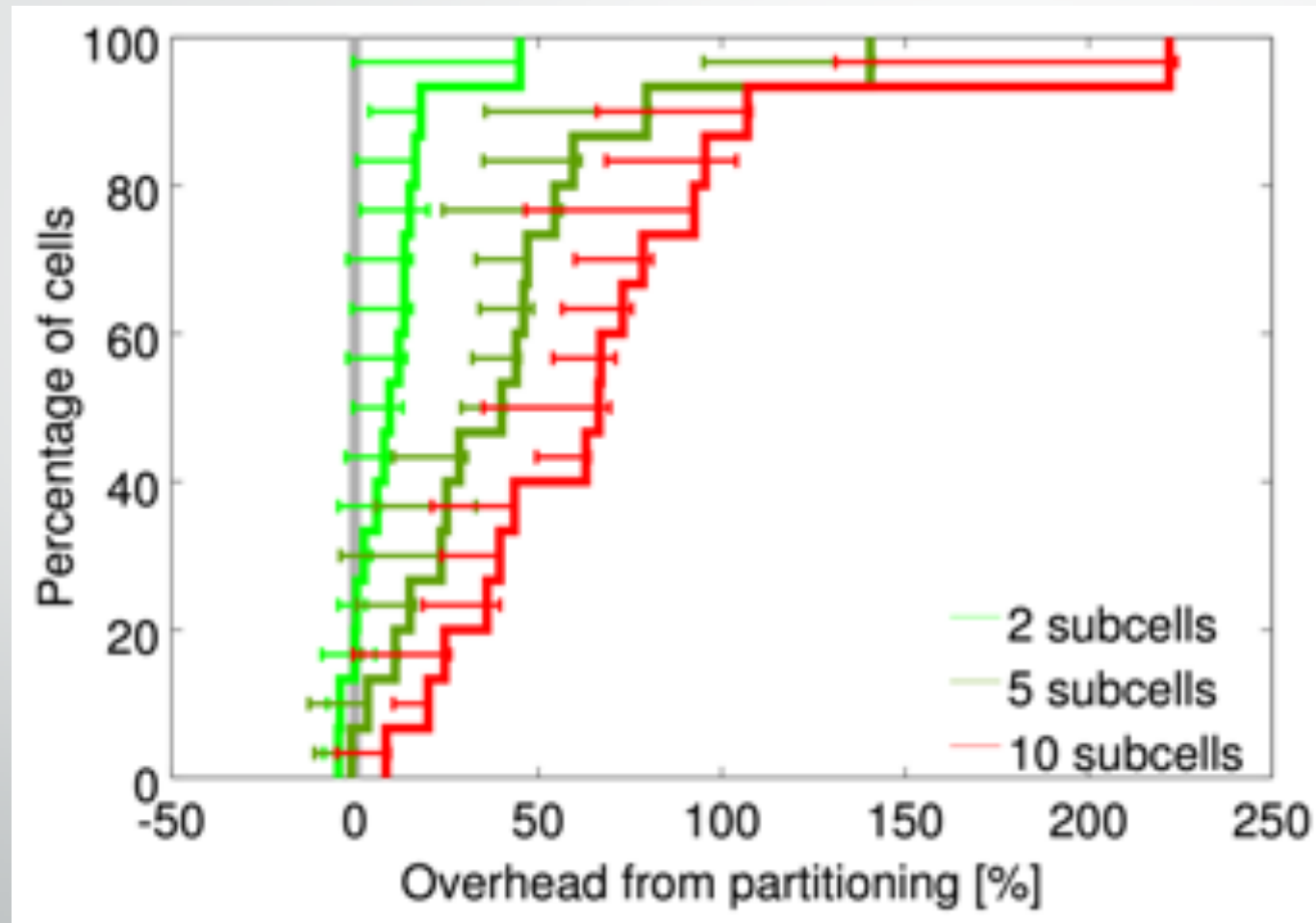


Figure 7(b)



What happens when you decrease cell size?

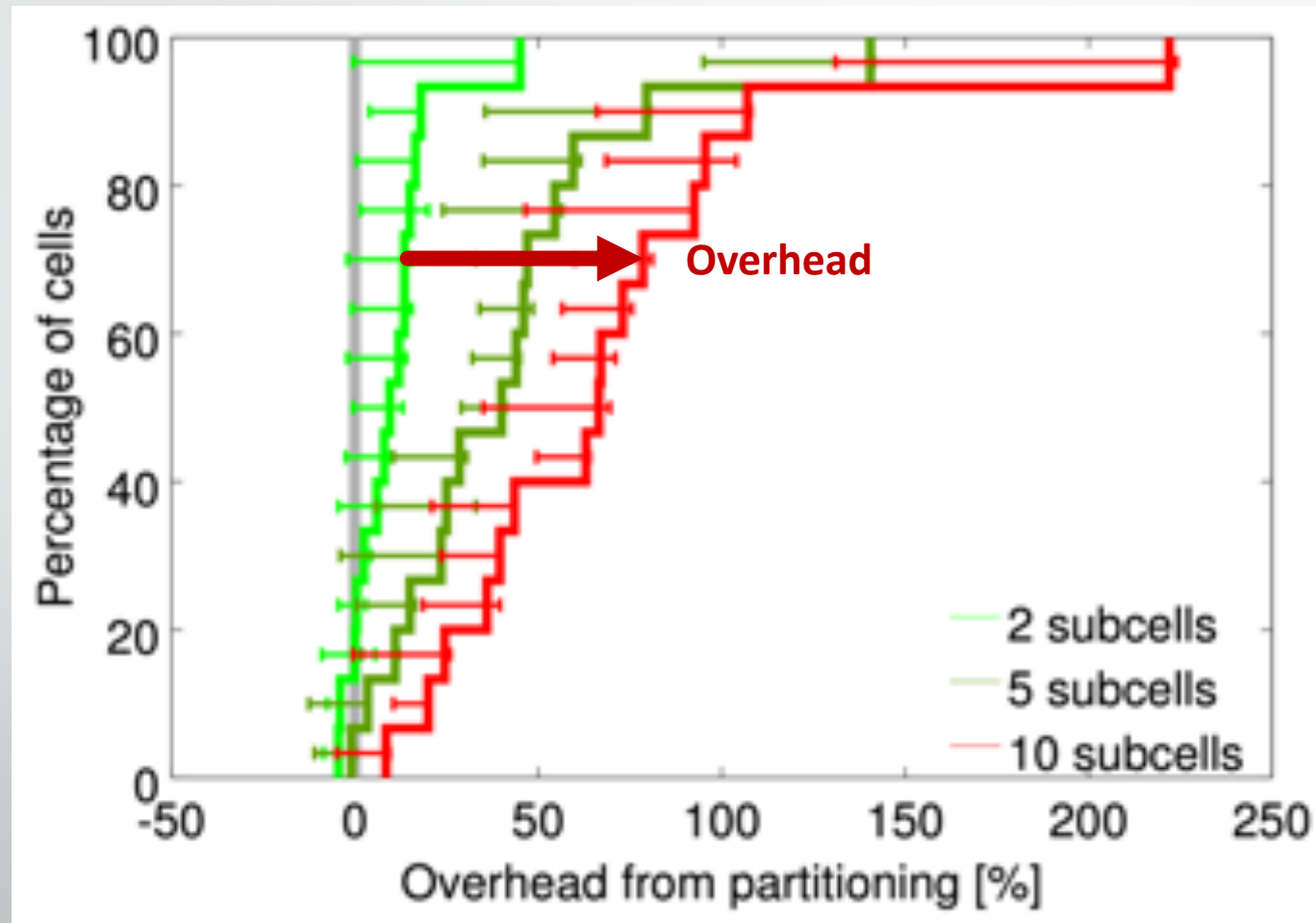


Figure 7(b)



What happens when users choose memory and compute freely?



What happens when users choose memory and compute freely?

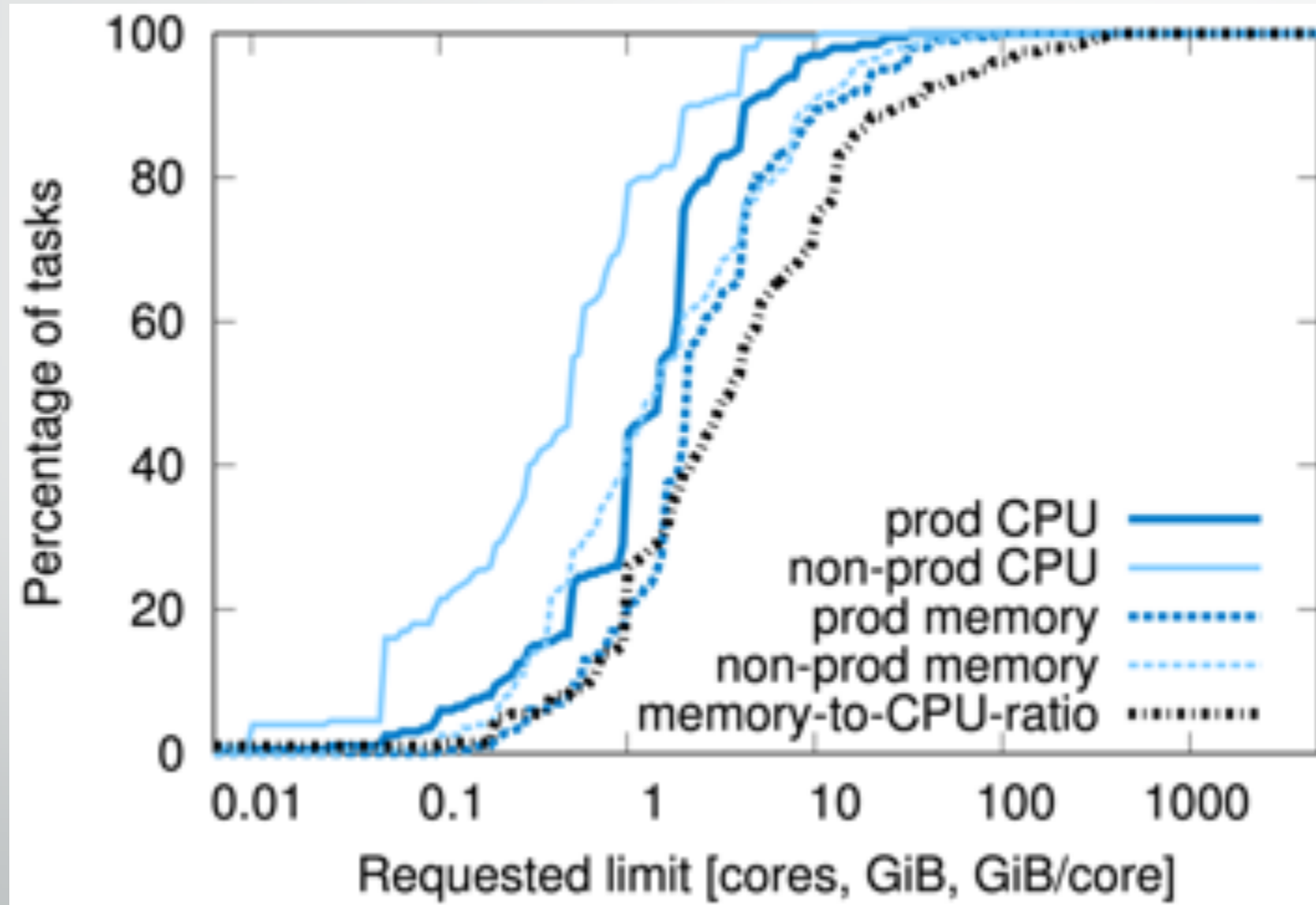


Figure 8



What happens when users choose memory and compute freely?

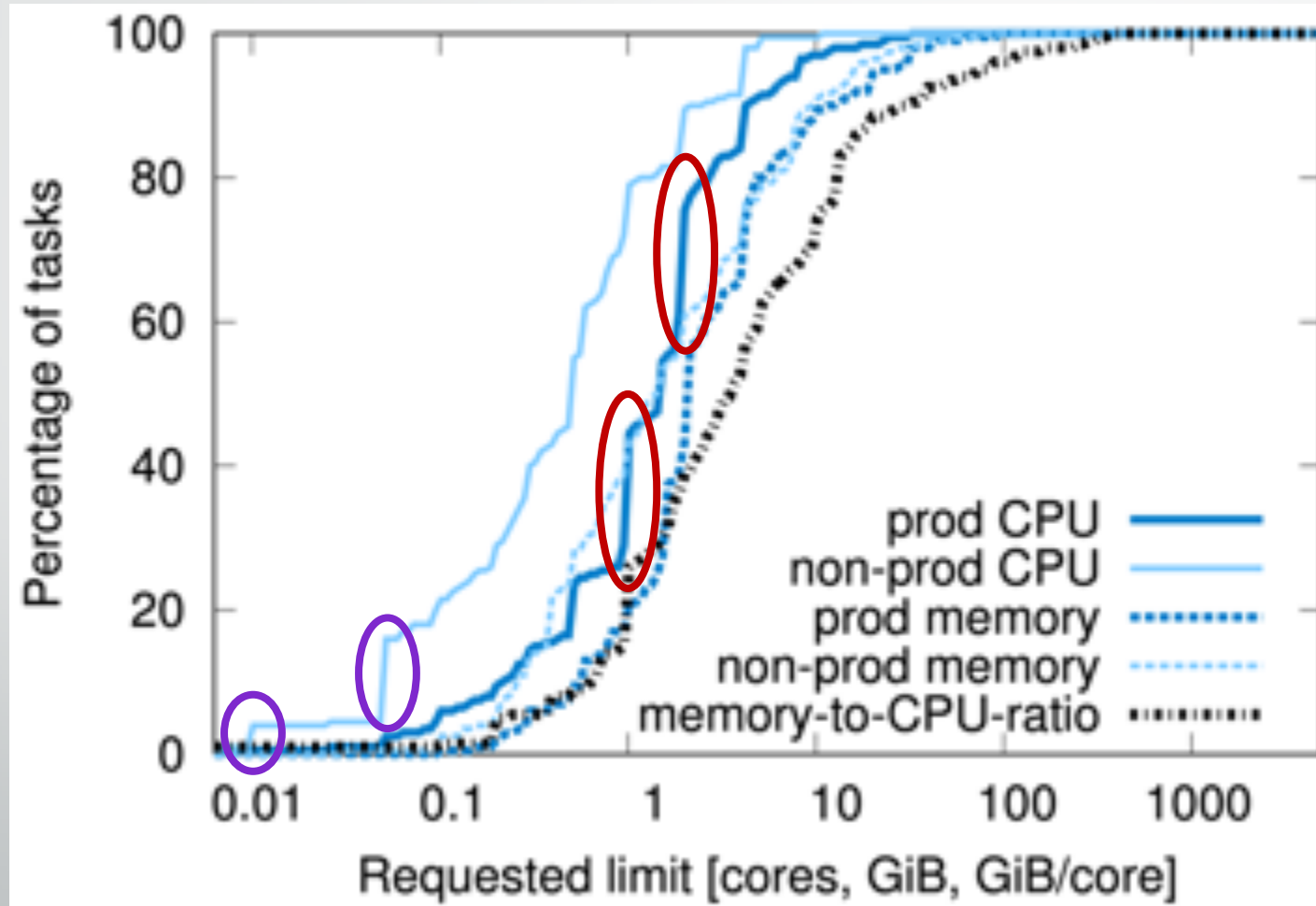


Figure 8



What happens when users choose memory and compute freely?

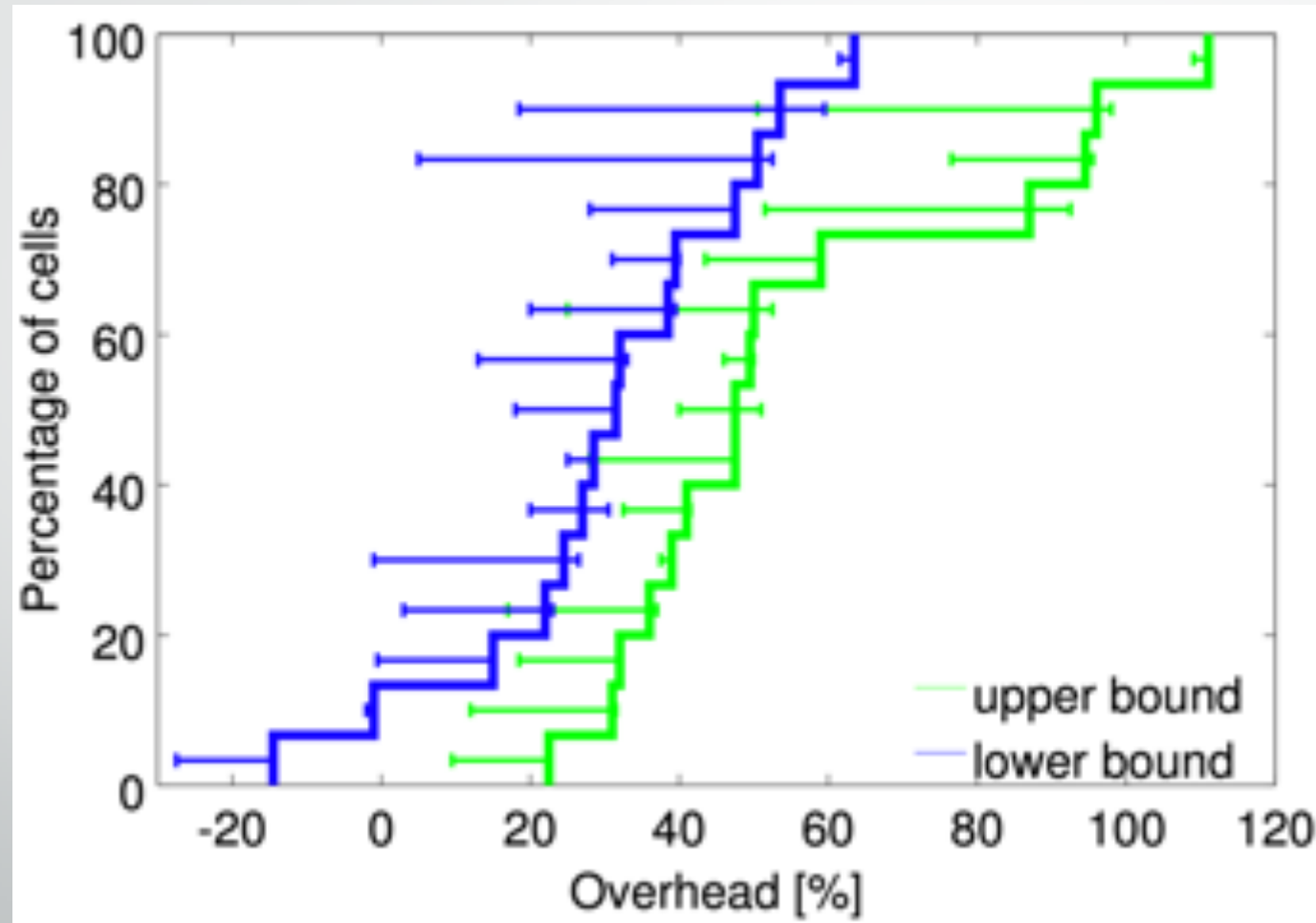


Figure 9



What happens when users choose memory and compute freely?



What happens when users choose memory and compute freely?

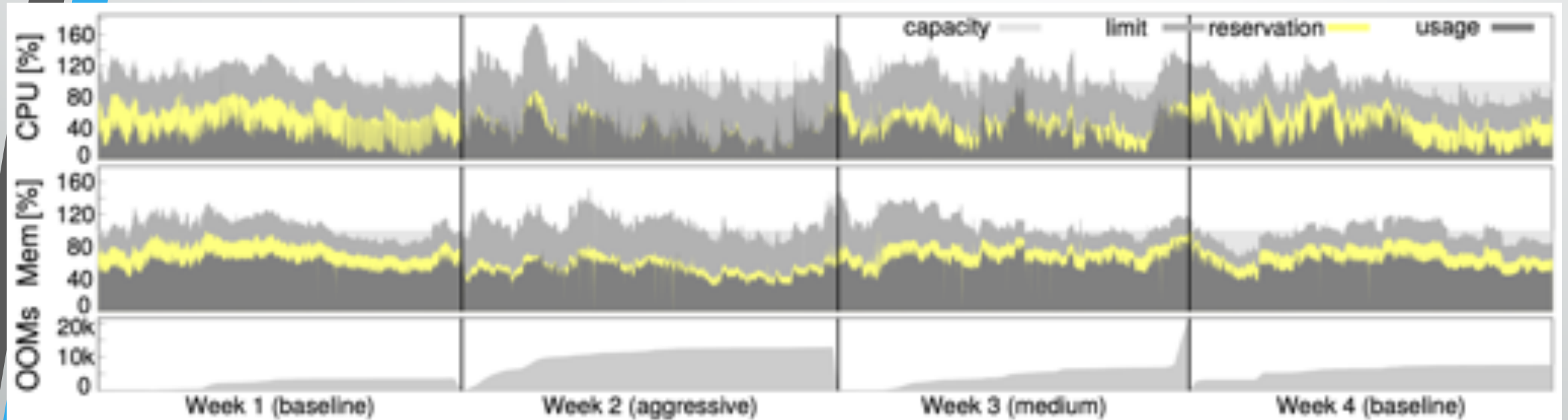
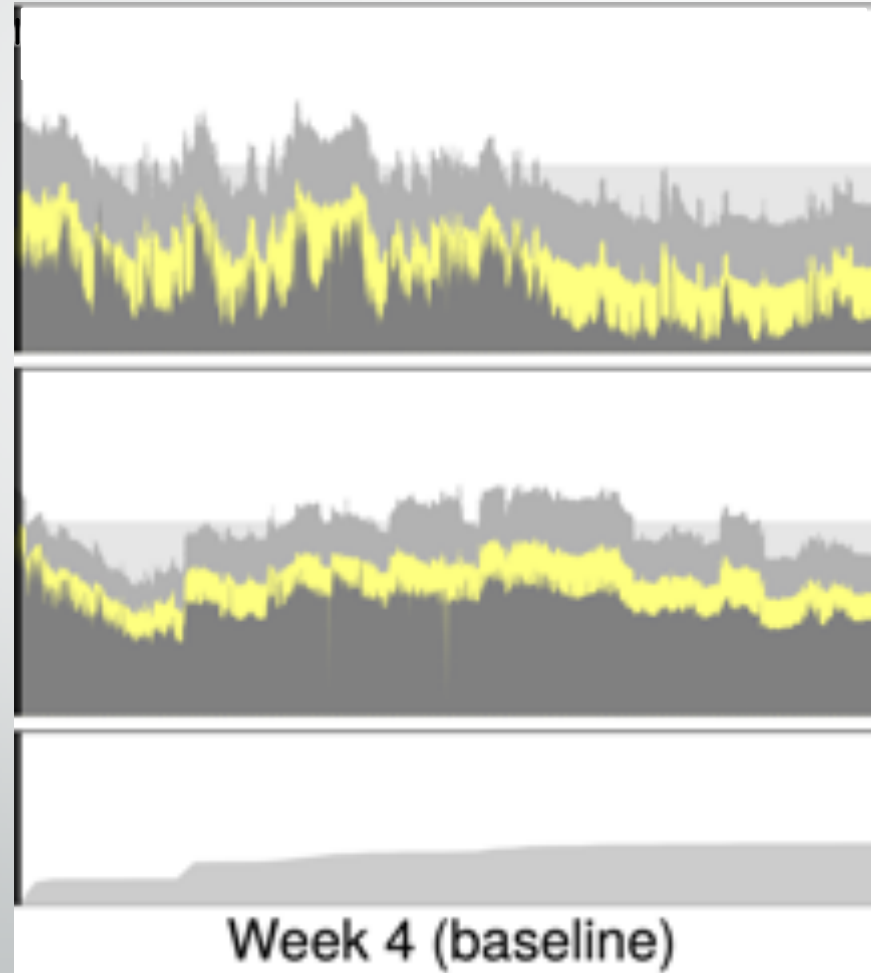


Figure 12

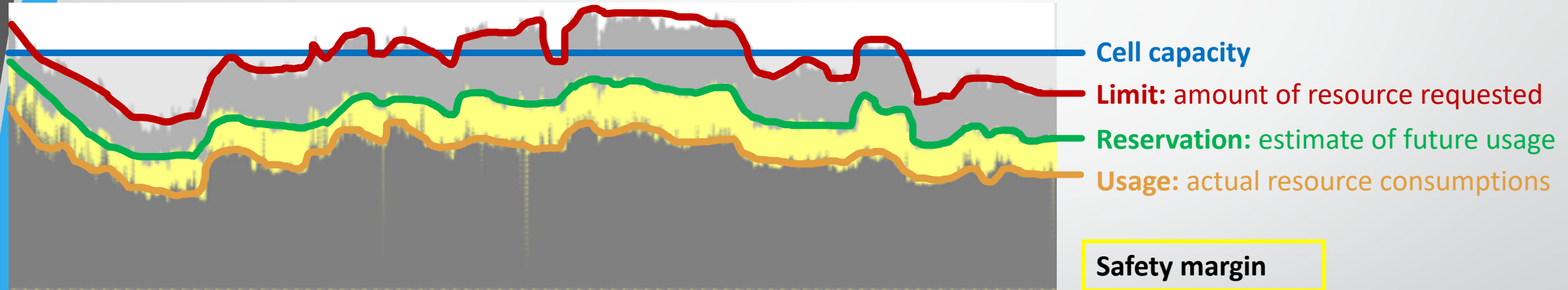


What happens when users choose memory and compute freely?





What happens when users choose memory and compute freely?





What happens when users choose memory and compute freely?

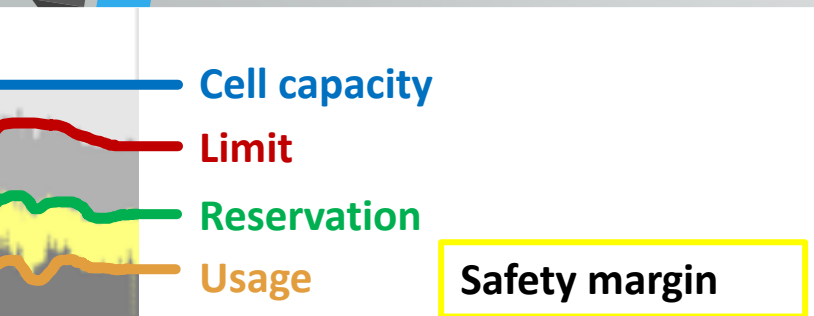
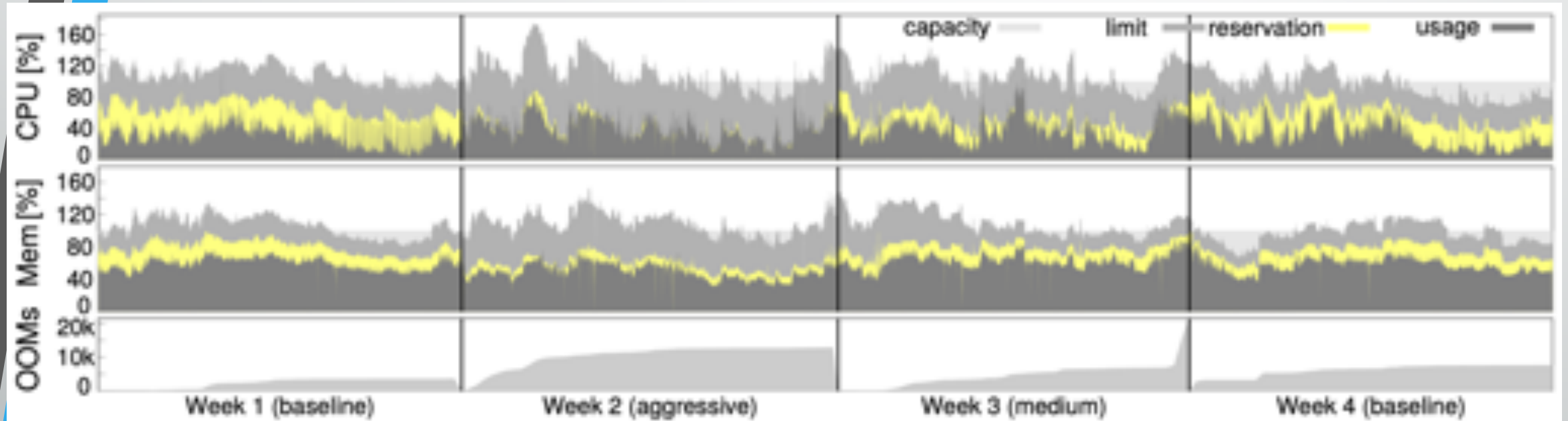


Figure 12



*Exception: region has 4 zones.

Source: <https://cloud.google.com/about/locations#regions>



The Future in Borg's Past

Kubernetes

- K8s is a set of tools which can be used to automate the deployment, scale, and orchestration of containers.
- K8s was released as OSS in 2014.
- Largely derived from the best practices and lessons learned from the Borg project.
- Google deploys Borg and Kubernetes internally.
- A managed Kubernetes service is available in Google Cloud Platform via Google Kubernetes Engine (GKE).



[This Photo](#) by Unknown author is licensed under [CC BY-NC-ND](#).

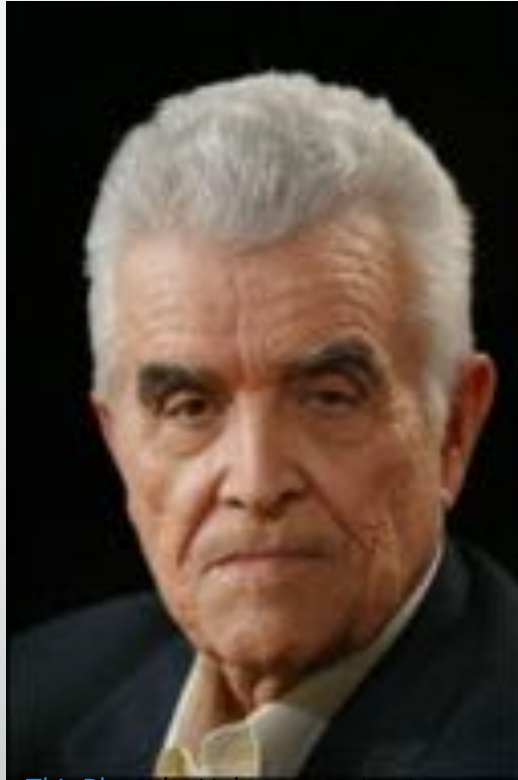
Kubernetes vs. Borg

- Kubernetes pods are largely derived from and equivalent to Borg allocs.
- Kubernetes transcends the Borg Job-Task paradigm through the labelling (key + value) of pods.
- Pod labels enable engineers to manage high-level services, instances of services, and pod subsets at a level of granularity which is not possible within the confines of the Borg Job-Task paradigm.
- Pod labelling enables significantly greater flexibility.
- Borg tasks carry the IP address of the machine they run on whereas Kubernetes pods can each allocated IP addresses.
- The scale attainable via Kubernetes deployments comes at a cost: massively complex configuration.
 - Solution: Google Kubernetes Engine (GKE) Autopilot.



[This Photo](#) by Unknown author is licensed under [CC BY-NC-ND](#).

René Girard and the Mimetic Theory of Desire



[This Photo](#) by Unknown author is licensed under [CC BY-NC-ND](#).

Site Reliability Engineering



[This Photo](#) by Unknown author is licensed under [CC BY-SA](#).

Borg and SRE Principles

- The development and maintenance of Borg largely gave rise to the principles of the Google Site Reliability Engineering methodology.
- Core SRE Principles:
 - Embrace risk
 - Define and adhere to service level objectives
 - Eliminate toil through automation
 - Monitor distributed systems
 - The release process is an engineering problem
 - Embrace simplicity

Strengths

- SRE != SWE: abstracting away infrastructure details.
- Built from first principles to meet Google's unique availability, scale, and workload requirements.
- Clearly defined evaluation metrics.
- Thorough comparison to existing large-scale server cluster systems.

Weaknesses

- Deeply entrenched proprietary systems increase switching costs.
- Borg tasks inherit the IP address of their host machine.
- Job-Task paradigm lacks flexibility.
- Human incentive to over-estimate resource consumption.
- Little insight into the operational dimensions of Borg within Google.

Further Reading on K8s, SRE, and Philosophy

- <https://kubernetes.io/>
- <https://static.googleusercontent.com/media/research.google.com/en//pubs/archive/44843.pdf>
- <https://sre.google/sre-book/part-ii-principles/>
- <https://sre.google/sre-book/production-environment/>
- <https://kubernetes.io/blog/2015/04/borg-predecessor-to-kubernetes/>
- <https://cloud.google.com/kubernetes-engine>
- https://www.theregister.com/2021/02/25/google_kubernetes_autopilot/
- <https://iep.utm.edu/girard/>

Sources

- Google Site / Datacenter: <https://www.datacenterknowledge.com/google-alphabet/google-spend-11-billion-new-data-centers-netherlands>
- Google Datacenter Cluster: <https://www.datacenterknowledge.com/archives/2012/10/17/google-reveals-its-data-centers>
- Borg Architecture: <https://www.nextplatform.com/2015/05/05/google-omega-to-become-part-of-borg-collective/>
- Cluster Management at Google with Borg, John Wilkes: <https://youtu.be/oW4gz8hVnok>

Q & A

