

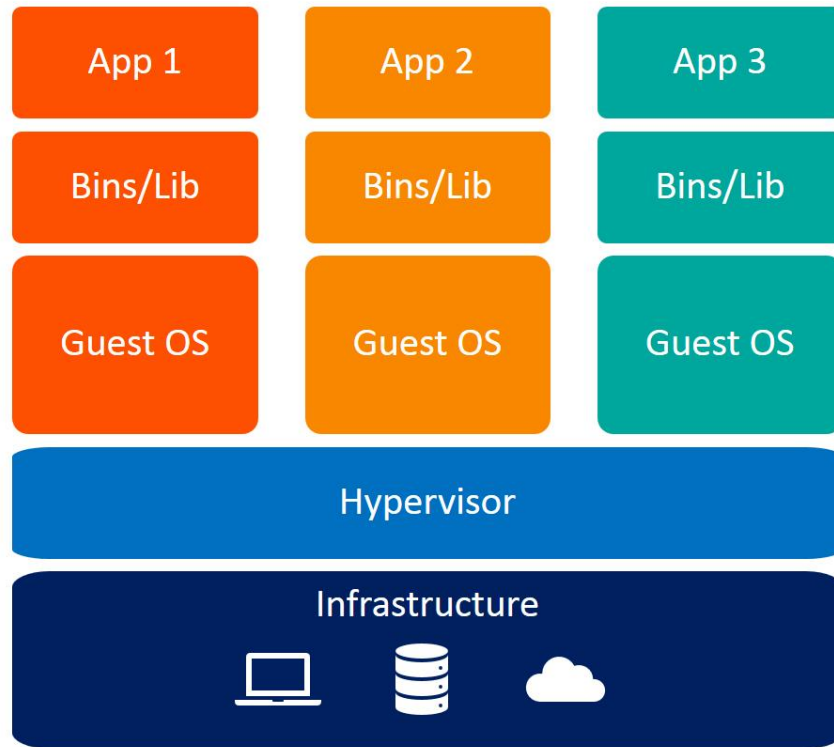
ML-centric cloud platforms

Jan 20 2026

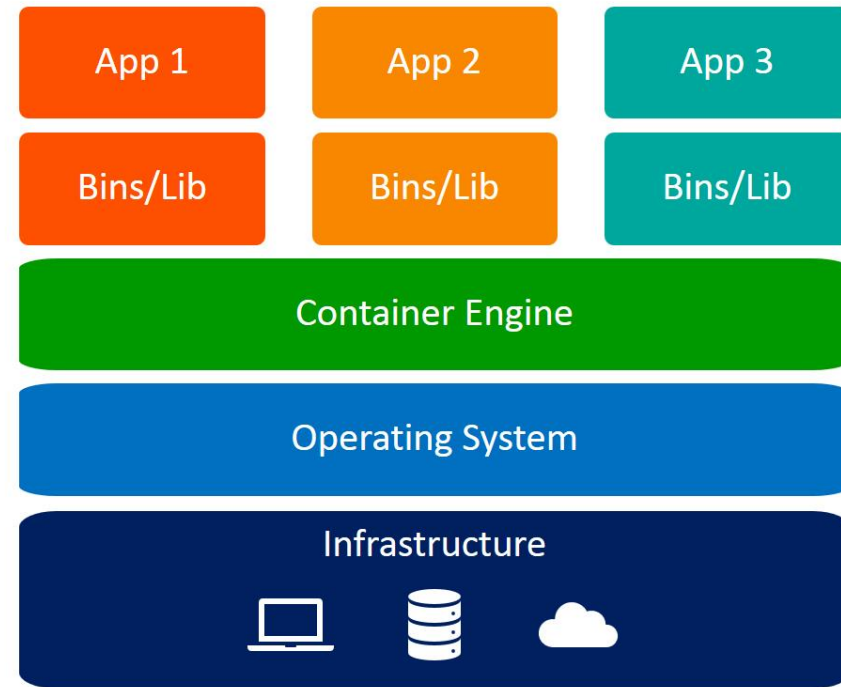
Administrivia

- Submit paper reviews on HotCRP before 9am
- Project dates and requirements will be posted this week

Cloud computing



Virtual Machines



Containers

What's in the cloud?

- Massive datacenters, ~1 million servers
- Heterogenous with clusters of homogeneity
- Two types of resource management:
 - Node agents – cpu scheduling, network allocation, power throttling
 - Cluster managers – VM placement, migration
- First-party vs third-party workloads

Cloud challenges

- Scale – need to manage millions of resources
- Difficult to track metrics + messy interactions
- Goals not well defined (e.g., what is an SLO?)
- Over-provisioning ➡ very low utilization

Cloud characteristics

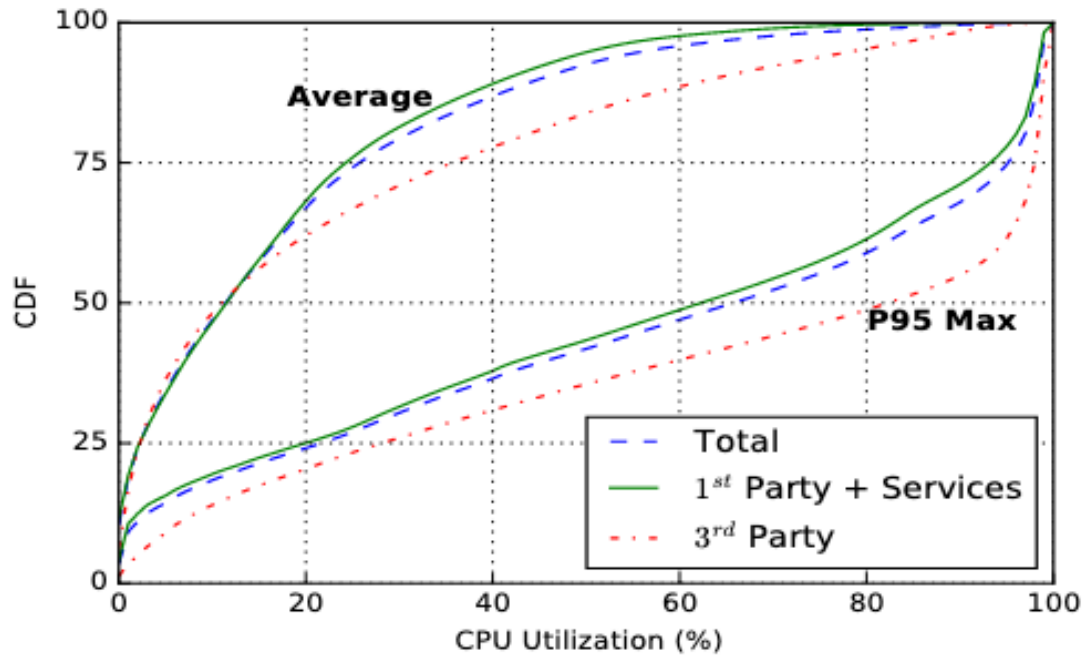


Figure 1: Average and P95 of max CPU utilizations.

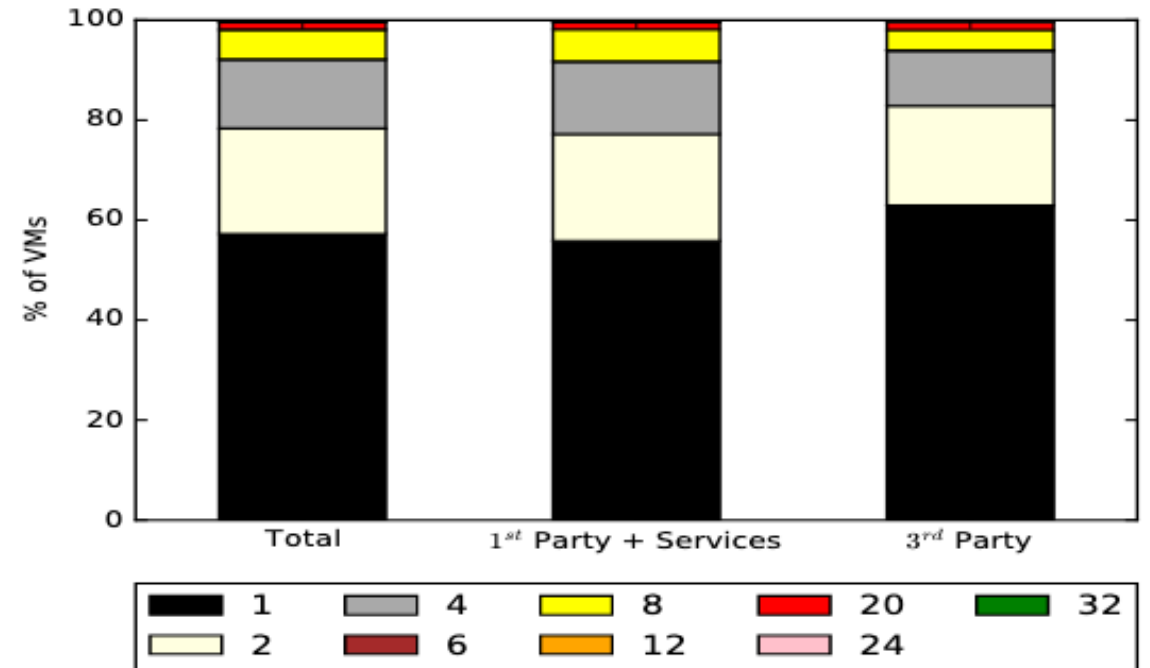


Figure 2: Number of virtual CPU cores per VM.

[1] Cortez et.al., "Resource Central: Understanding and Predicting Workloads for Improved Resource Management in Large Cloud Platforms" at SOSP'17.

Cloud characteristics

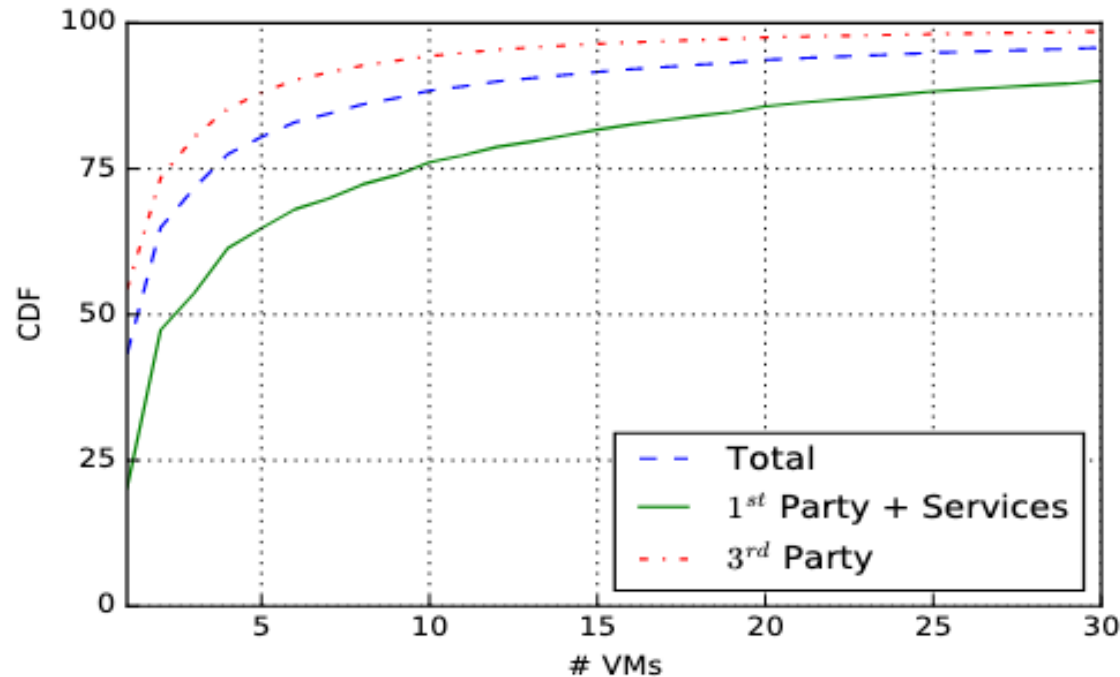


Figure 4: Max number of VMs in each deployment

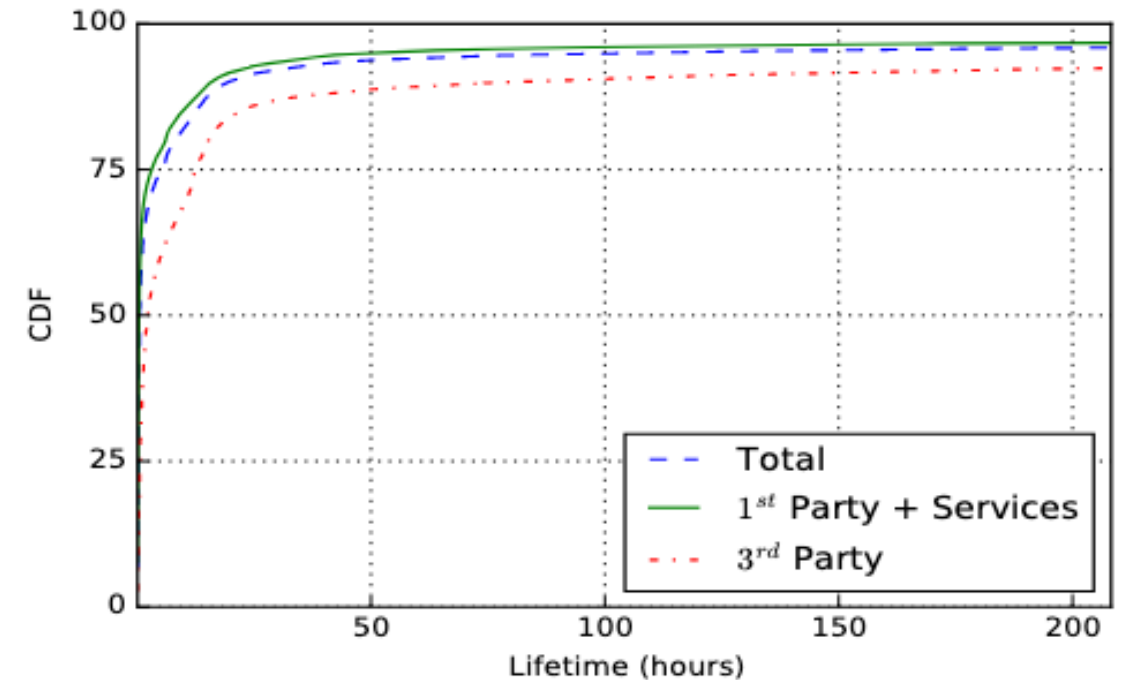


Figure 5: VM lifetime.

[1] Cortez et.al., "Resource Central: Understanding and Predicting Workloads for Improved Resource Management in Large Cloud Platforms" at SOSP'17.

ML for the cloud: recent works

- Workload + resource utilization prediction
- CPU: harvesting, power capping, interference prediction
- Memory: harvesting, cold data prediction, tiering
- Storage: failure prediction, admission, prefetching
- Microservices: Auto-scaling

Discussion

- Resource central uses a centralized “prediction service” that informs decisions across the fleet. What kind of policies are not supported by this architecture? When would “on-node” learning be beneficial?
- Ensuring robustness is critical in the cloud. What safeguards can cloud providers take to ensure ML is not leading to adverse behaviors?

Title: [Lecture 3]: Discussion