



Kubernetes & YARN: a hybrid container cloud

Jian He & Bushuang Gao

About us

- Jian He
 - Staff Engineer @Alibaba cluster management team
 - Staff Engineer @Hortonworks
 - Hadoop Committer & Project Management Committee member
- Bushuang Gao
 - Senior Engineer @Alibaba

Agenda

- What/Why co-location
- Co-location @ Alibaba
- Kubernetes vs YARN
- Kubernetes & YARN: a hybrid architecture
- Resource Isolation
- Future

What/Why co-location

Data center Avg utilization

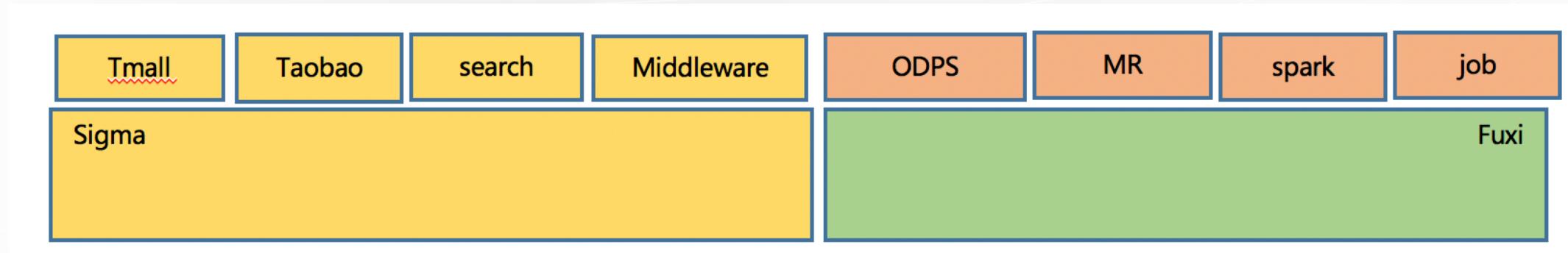


- Gartner has long talked about the "80% rule": that 80 percent of IT budgets get spent simply "keeping the lights on"
- The average data center cpu utilization is about 10%

Workloads in Alibaba data center

- Services
 - Commercial platform, taobao, Tmall, retail apps, search,
 - Payment platform: Alipay
 - Middleware: queue service, cache server
- Jobs
 - Big data batch jobs: MR, spark
 - Realtime jobs: flink

Before: separated Cluster

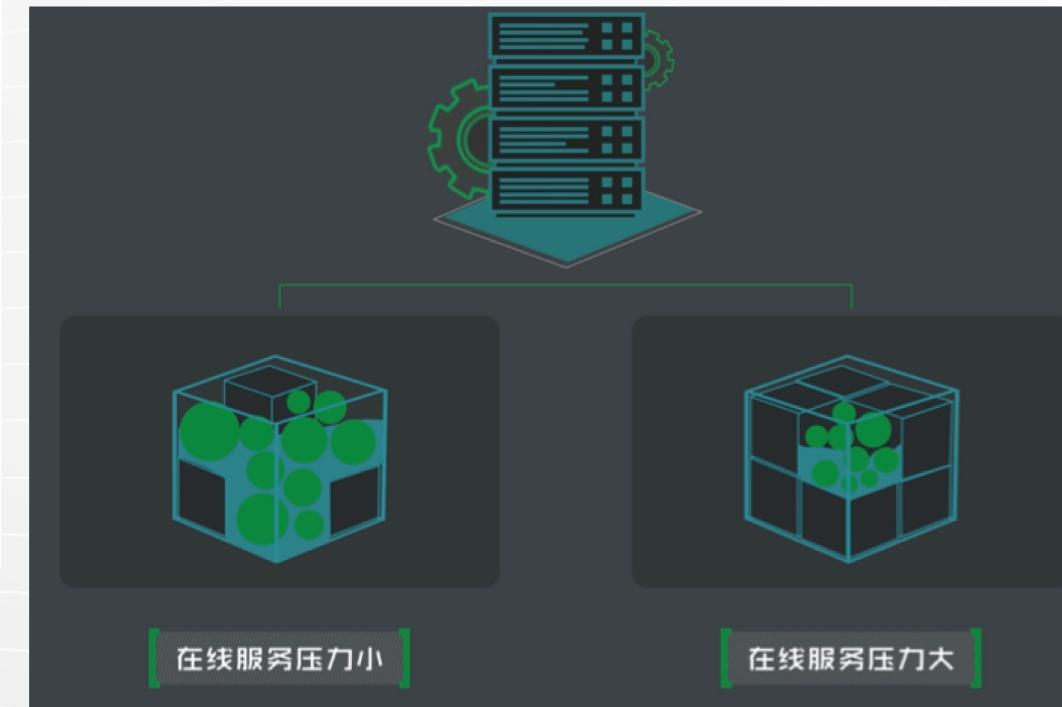


Services vs Jobs

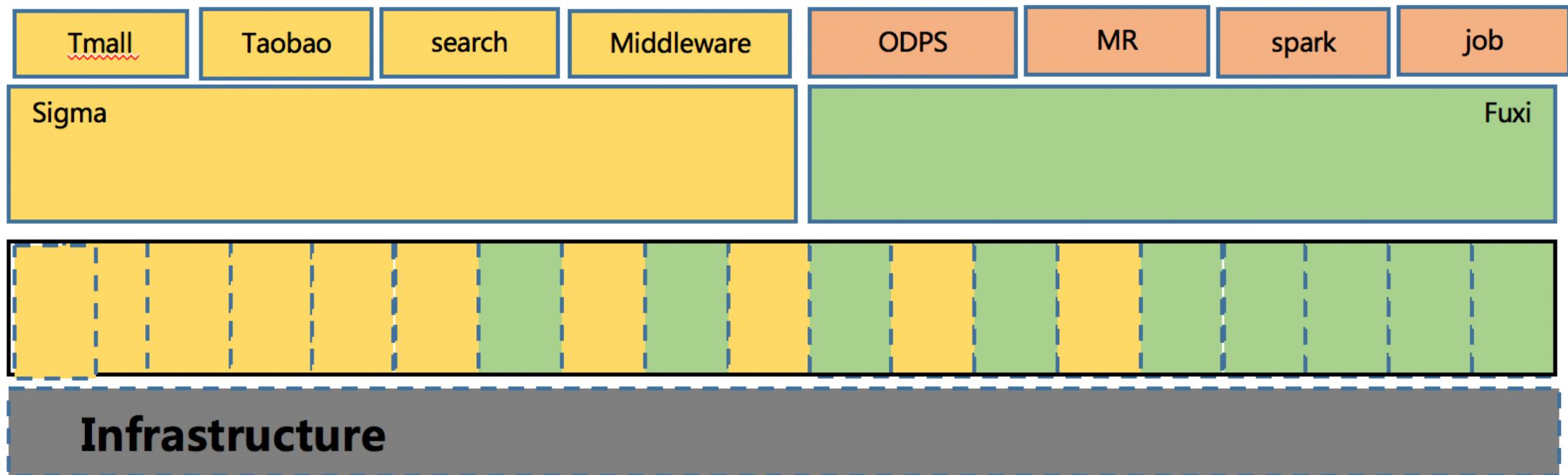
	Online service	Batch jobs
Category	Online shopping web apps, payment service	MR, spark, flink
Latency	Sensitive	Insensitive
Priority	high	low
Traffic pattern	Peak at day time	Peak at night time
Fault tolerance	should not fail	Fail and retry



Complementary !

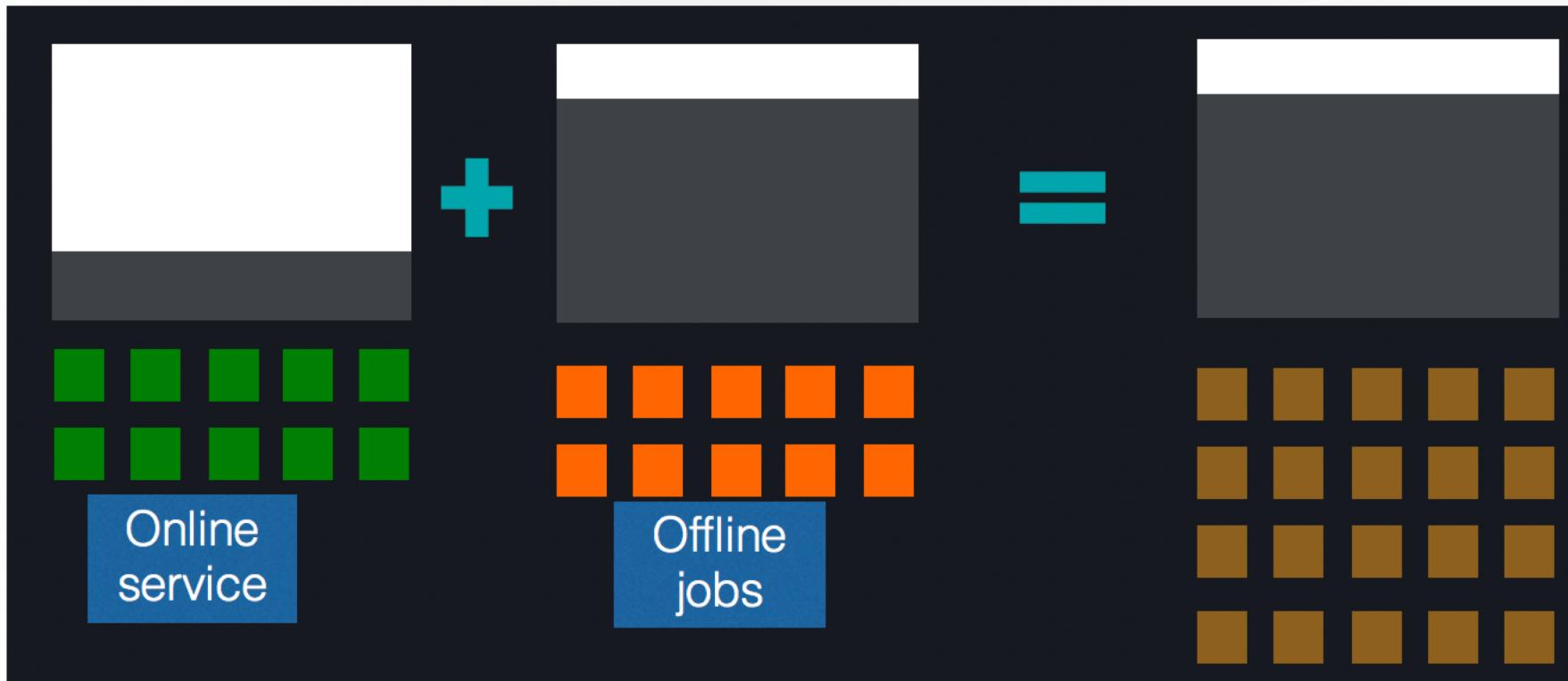


Now: co-locate clusters



What is co-location

- A hybrid deployment that supports both online services and offline jobs using the same hardware.



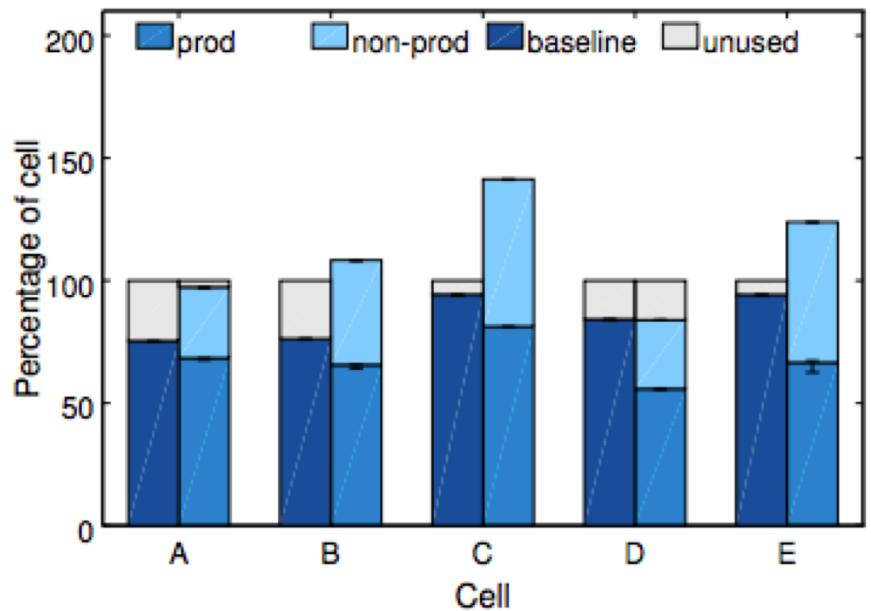
Key concepts in co-location

- Resource sharing: Make services and jobs share a single cluster and have them run on the same hardware
- Resource Isolation: Rely on scheduling efficiency and resource isolation to guarantee the SLA.

Why co-location

- Two separated cluster cause huge resource wastage. Increase resource utilization.
- 1 hour peak traffic on double 11 requires a lot of extra resources. However, those resources will be idle after the peak time.
- Lend those resources to computation jobs

Google Borg



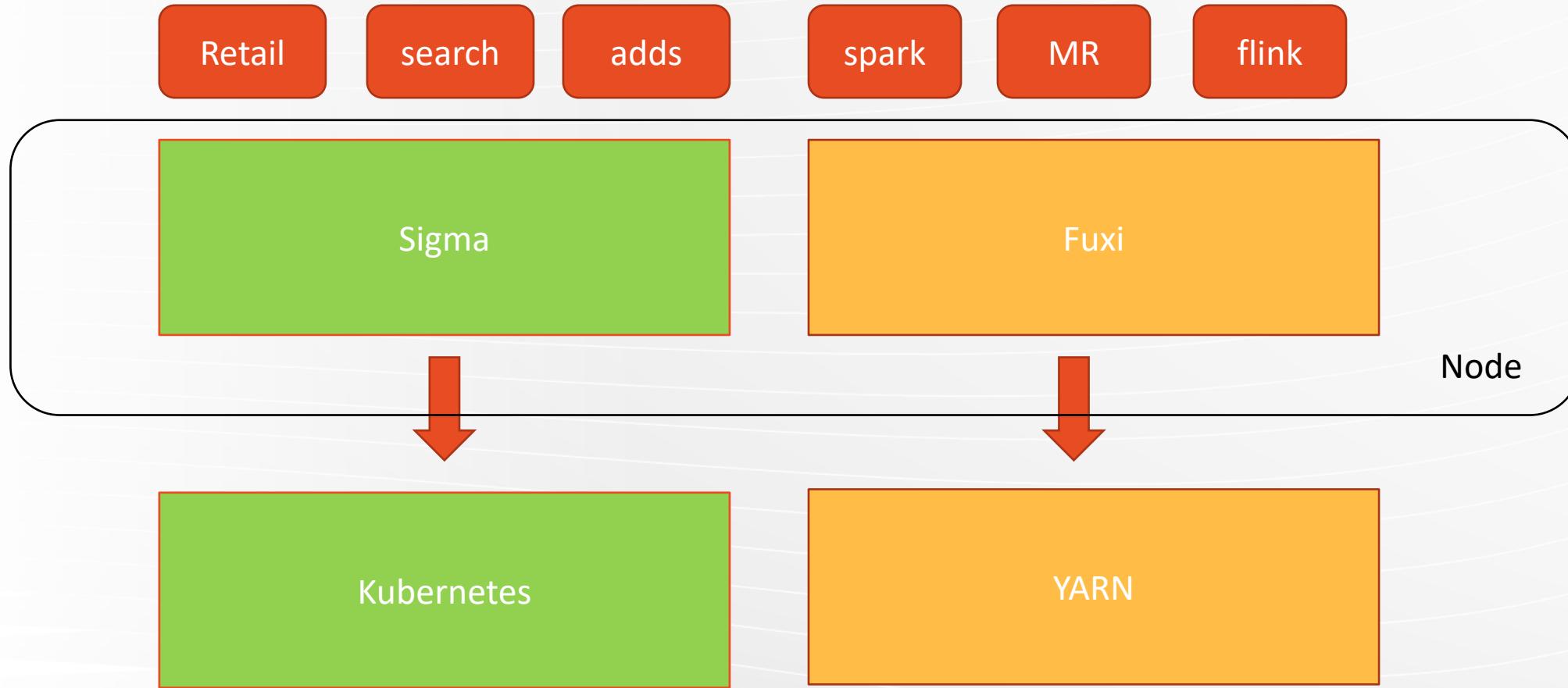
(a) The left column for each cell shows the original size and the combined workload; the right one shows the segregated case.

Borg paper mentions 20% - 30% more machines if
If segregating prod and non-prod workloads



Co-location @Alibaba

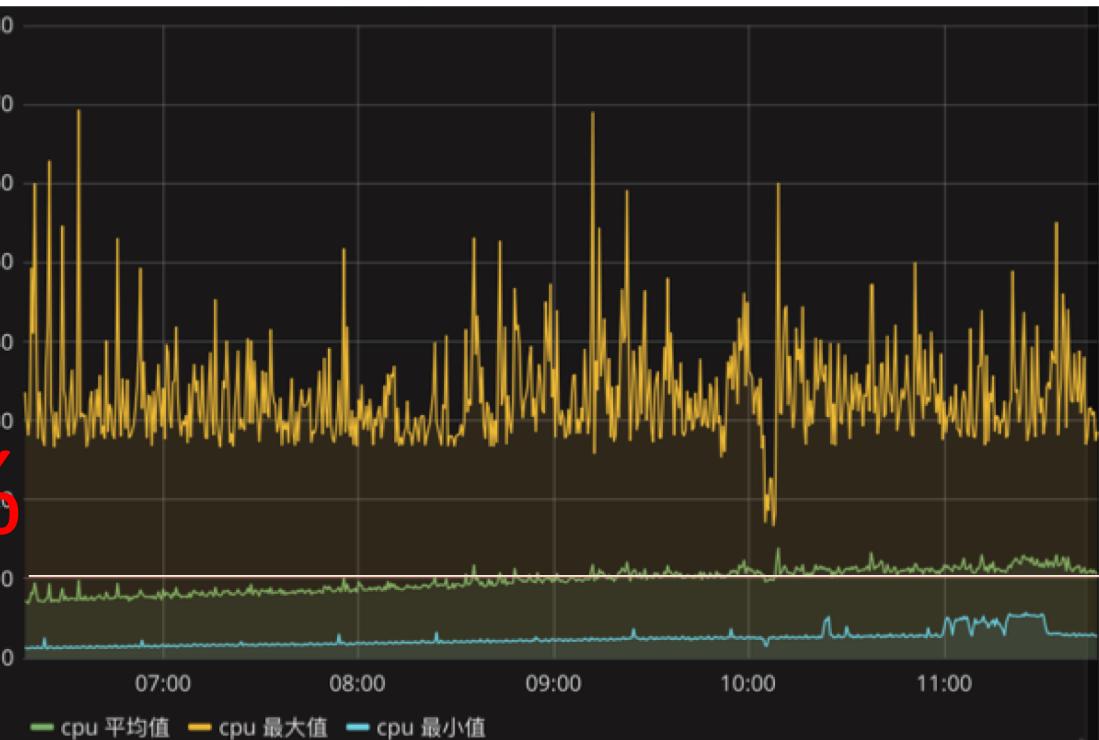
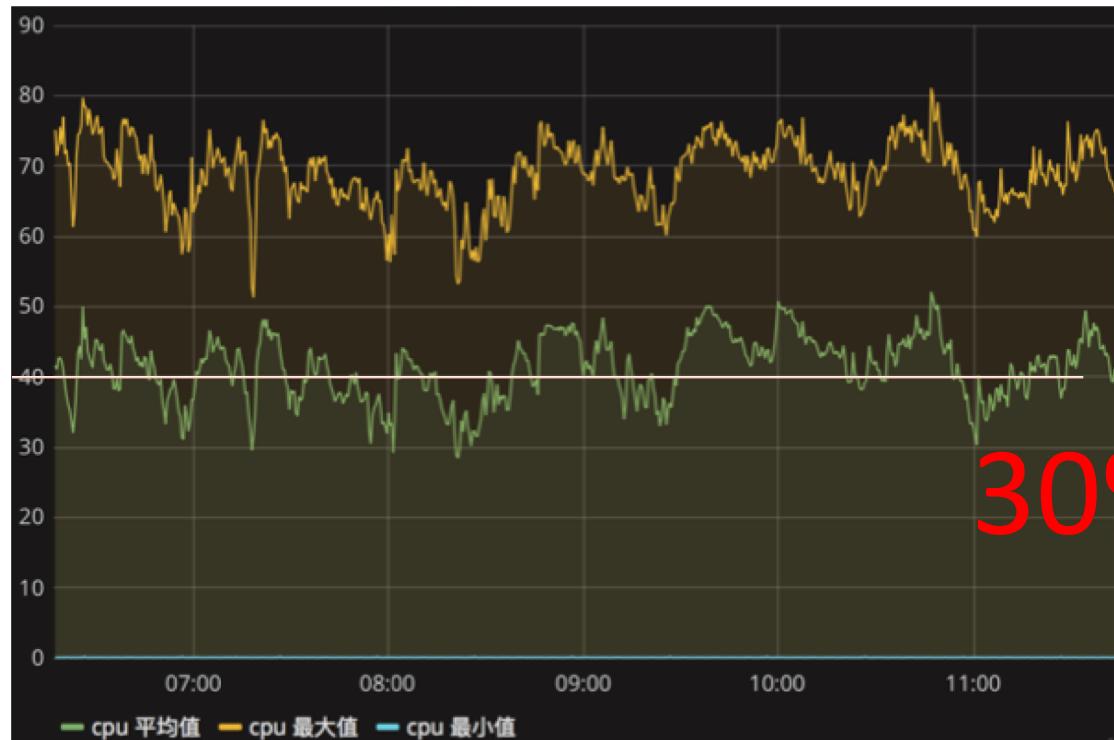
Co-location @Alibaba



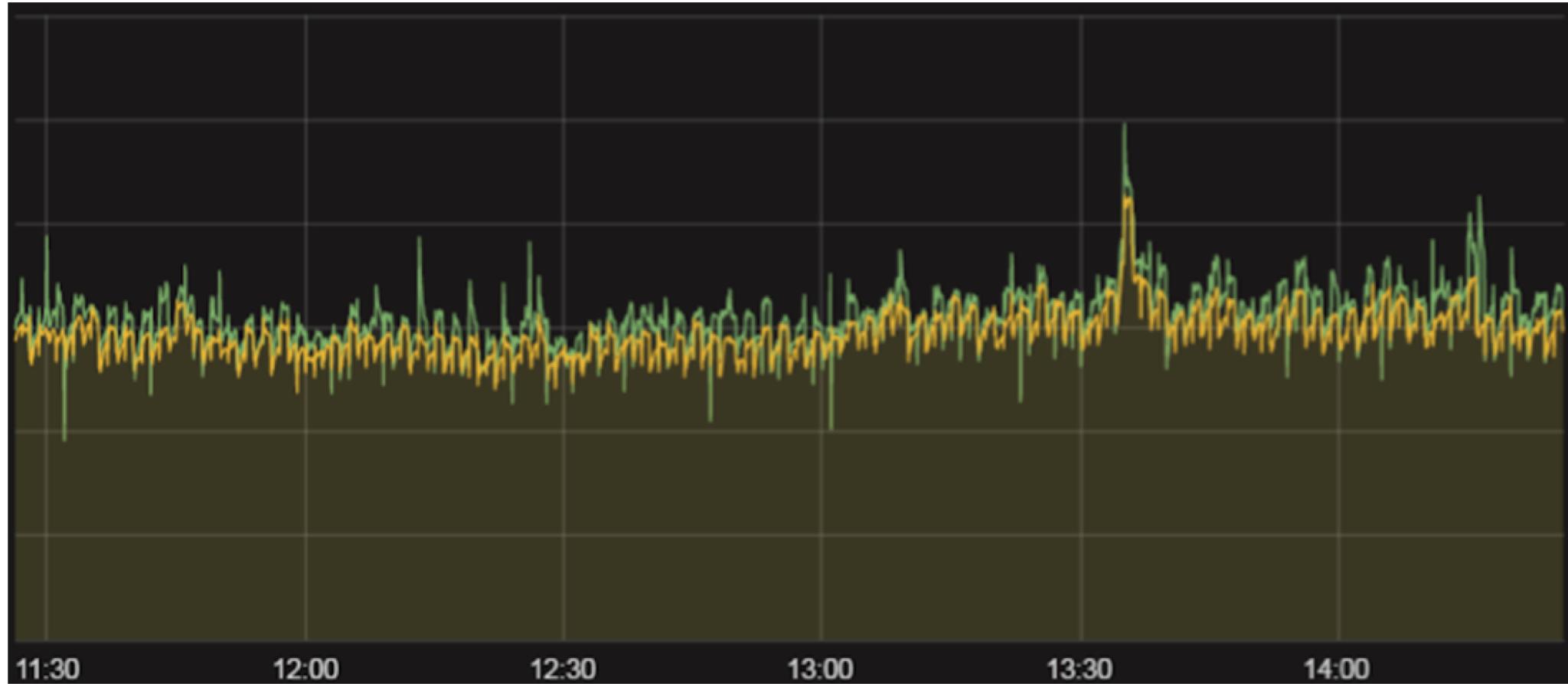
Co-located vs separated cluster CPU utilization

Co-located: **40%**

Seperated : **10%**



Impact on online service within 5%

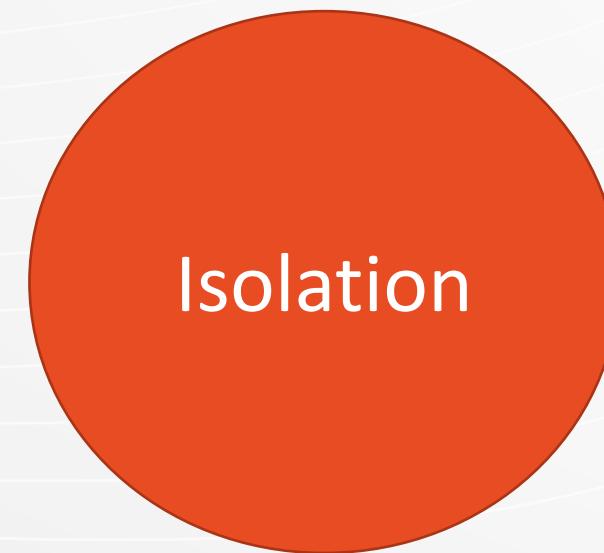


Two core concepts in co-location



Efficient placement of service container and tasks

Resource contention



When placed together, don't affect each other

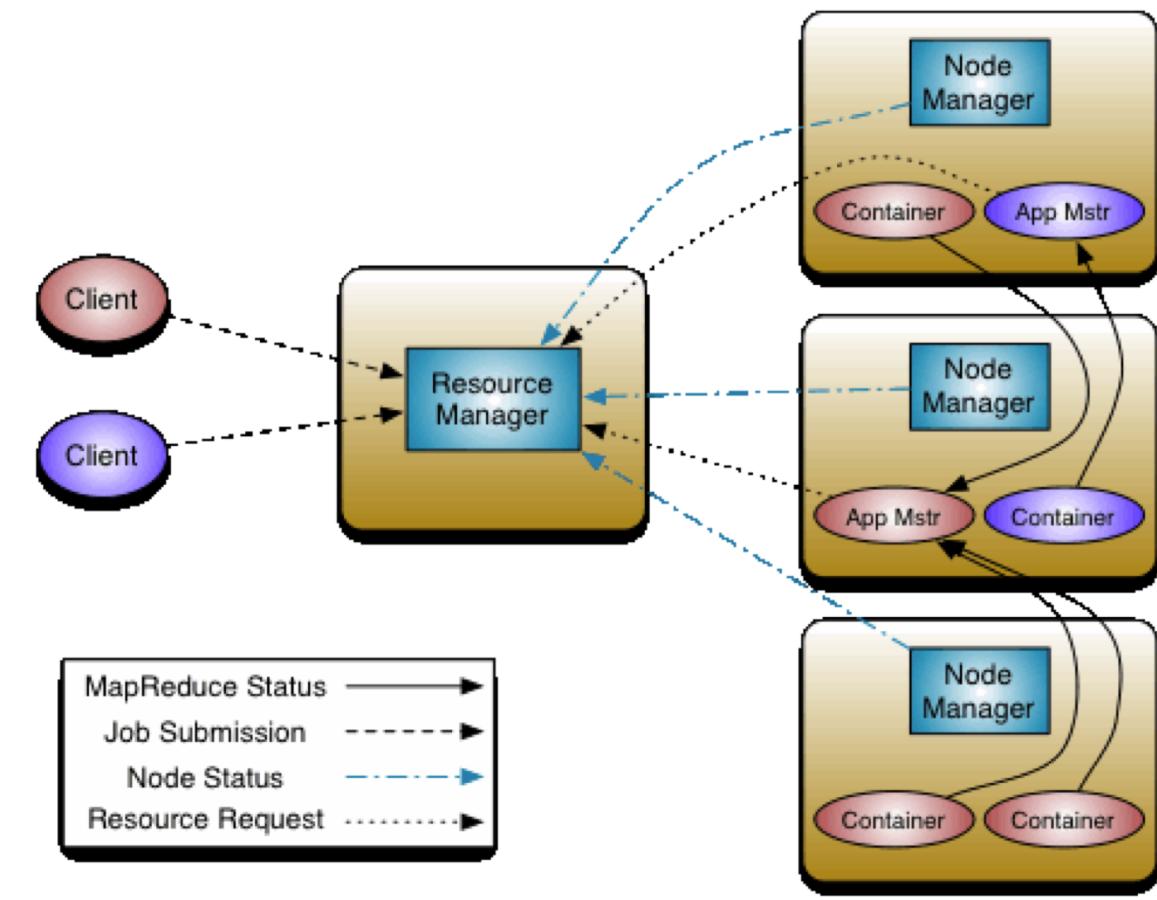
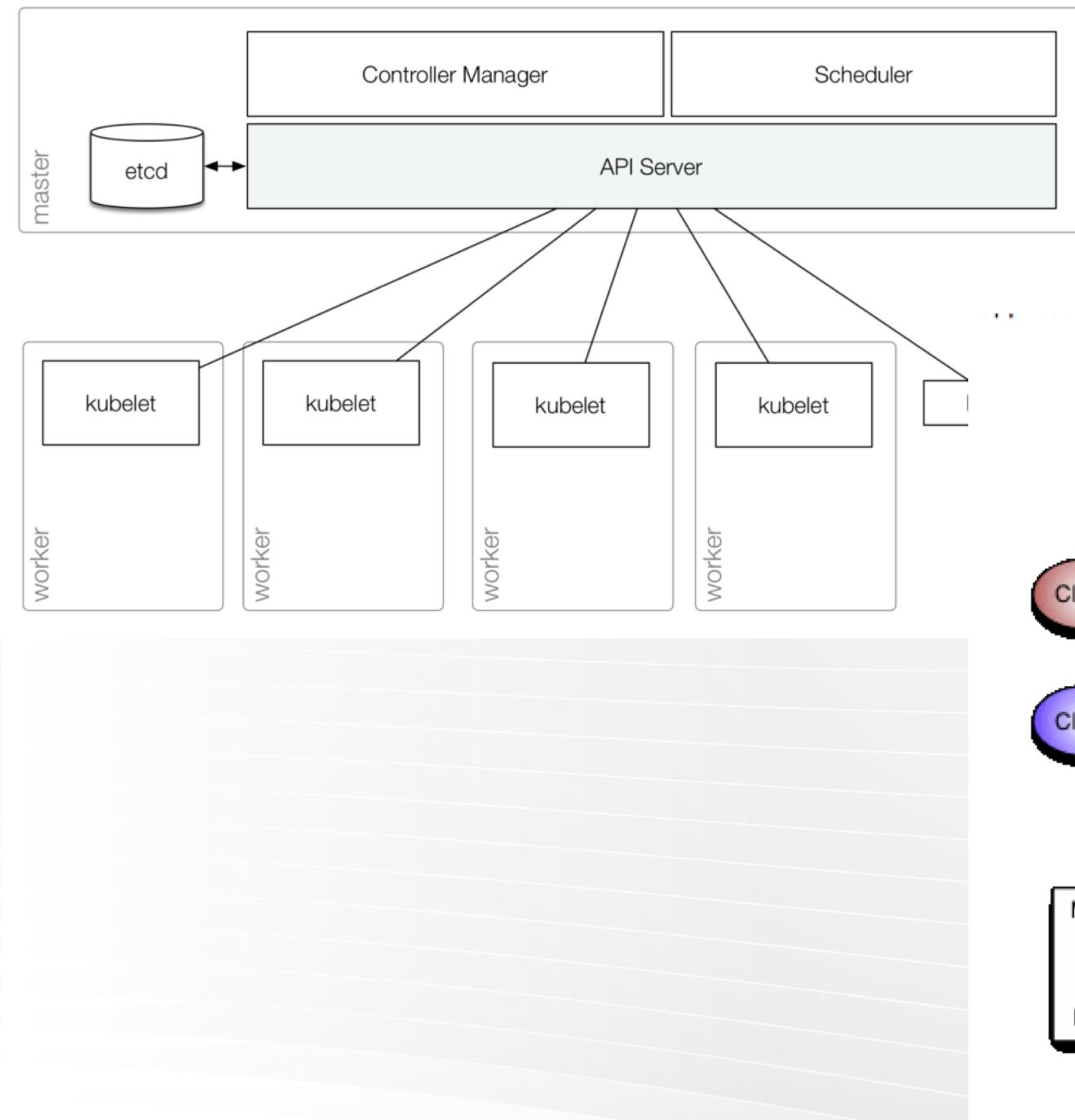
Time sliced scheduling



Normal day <-> Sales day



Kubernetes vs YARN



Design Principal

Kubernetes	Focus on long running service. Driving current state towards desired state with control loops
YARN	Focus on scheduling jobs

Scheduling Unit

Kubernetes	<p>Container centric – bottom up.</p> <p>Container is the primitive. Other primitives such as replicaset, deployment are built around containers.</p>
YARN	<p>Application centric: top down.</p> <p>Scheduling sequence: Queue -> user -> application -> container request</p>

Communication

kubernetes	Based on api-server watch mechanism Everything stored in etcd
YARN	Based on RPC Only application-level metadata persisted. Container data is not persisted. Recover from in-memory state from peers

Container Runtime

kubernetes

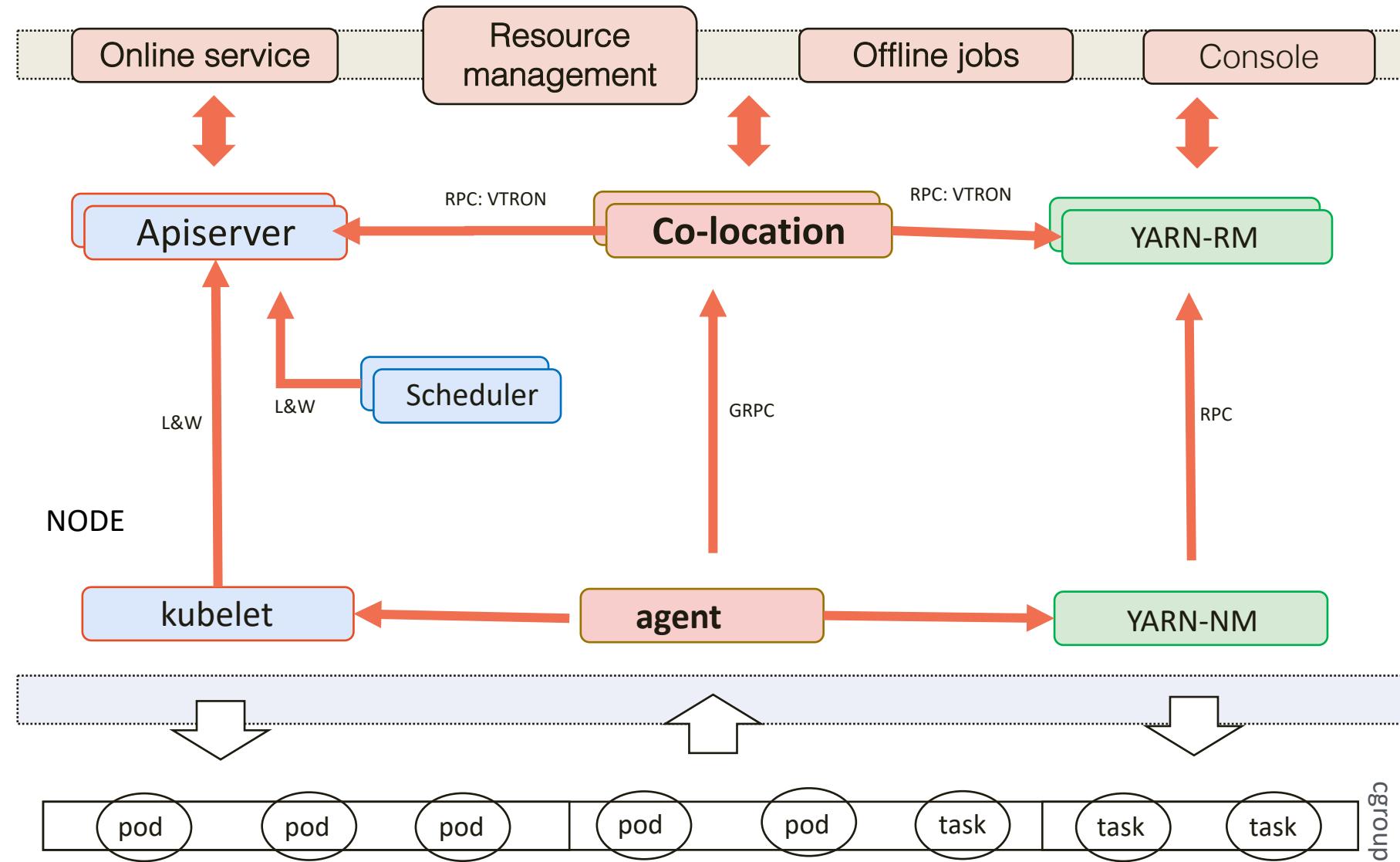
CRI compatible. Docker etc.

YARN

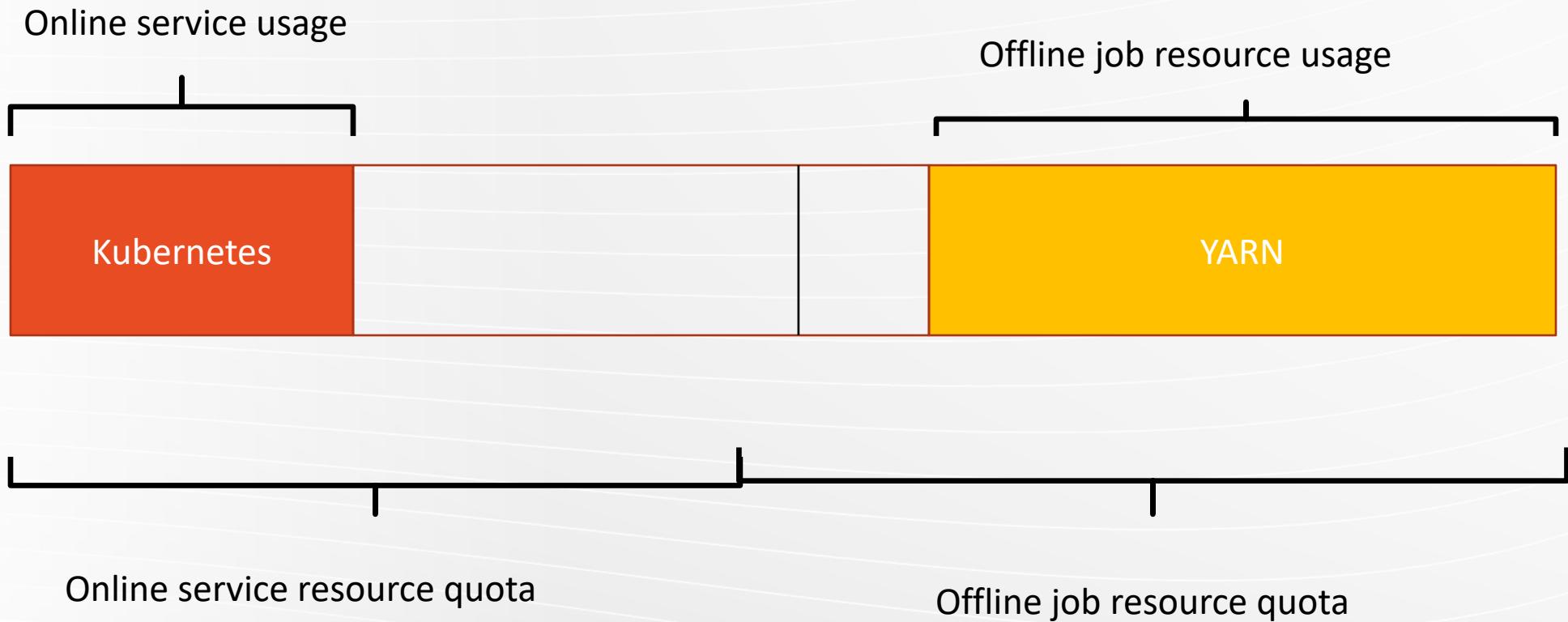
Docker + TAR ball

Kubernetes & YARN: a hybrid architecture

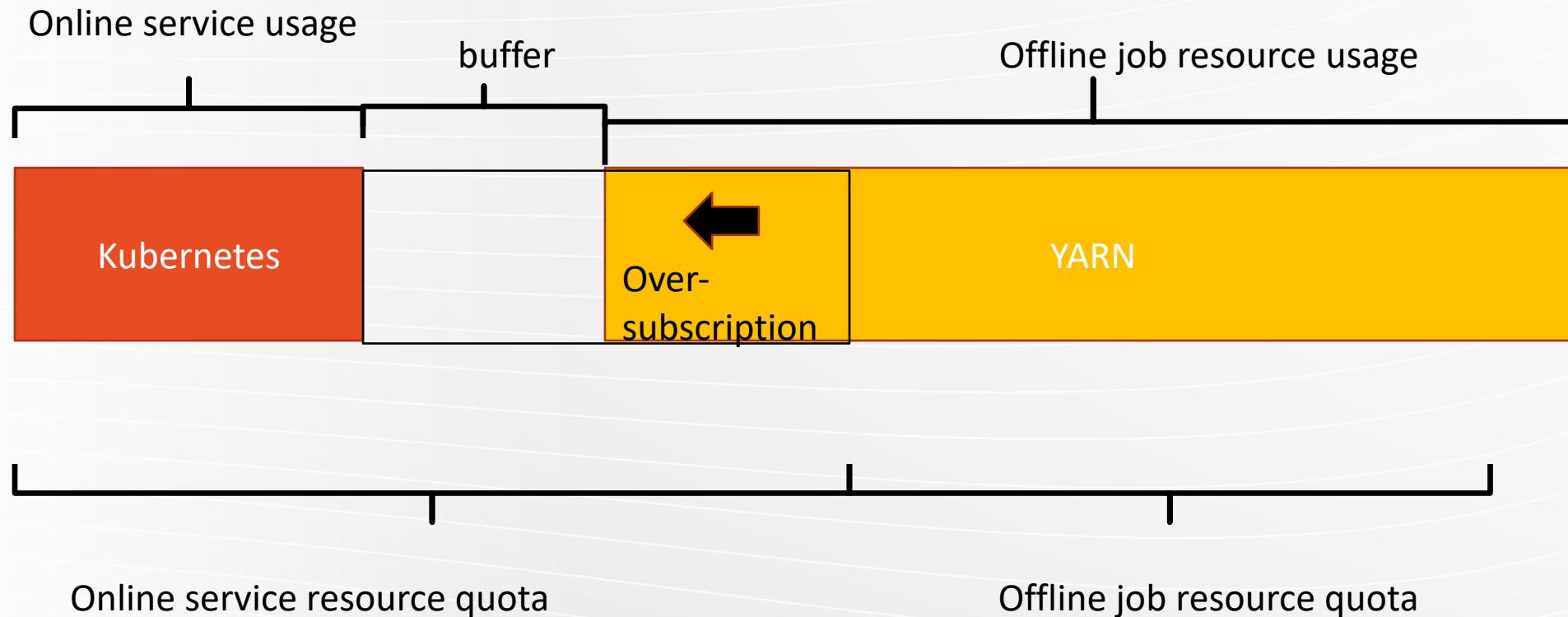
Architecture



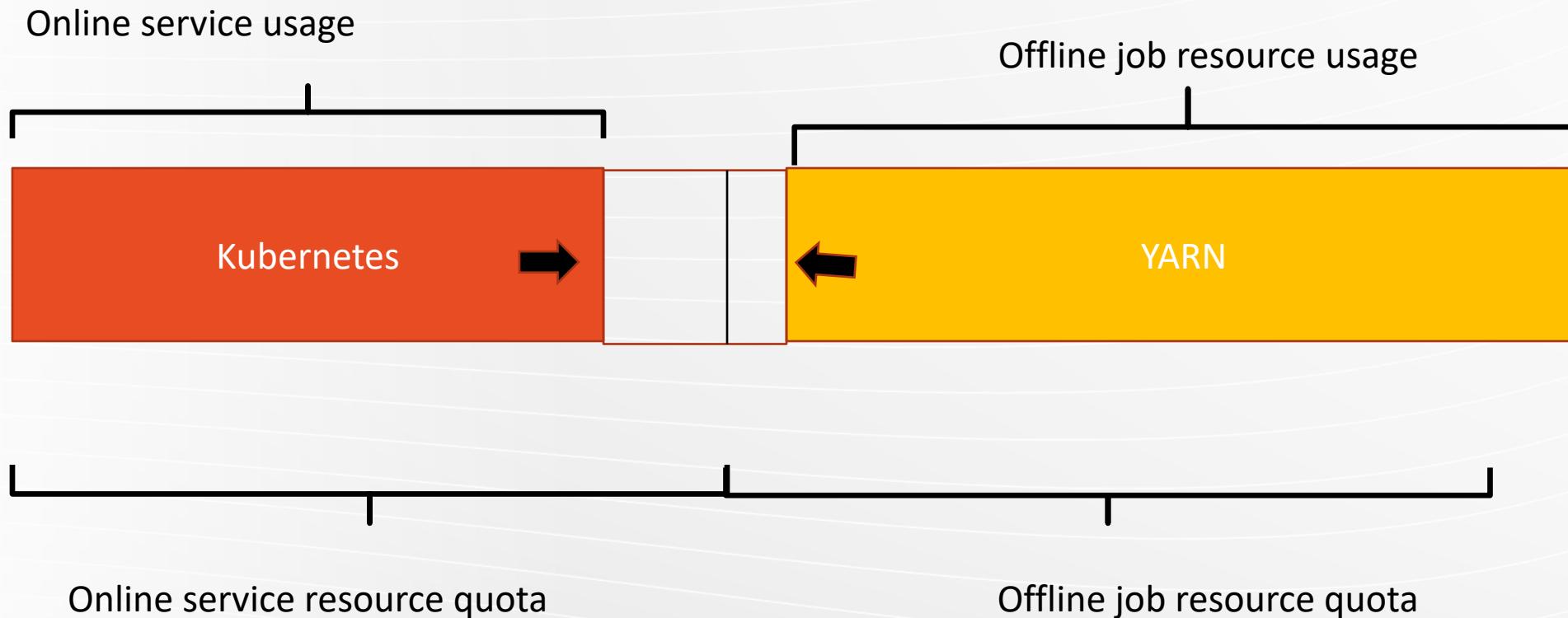
resource sharing



resource sharing

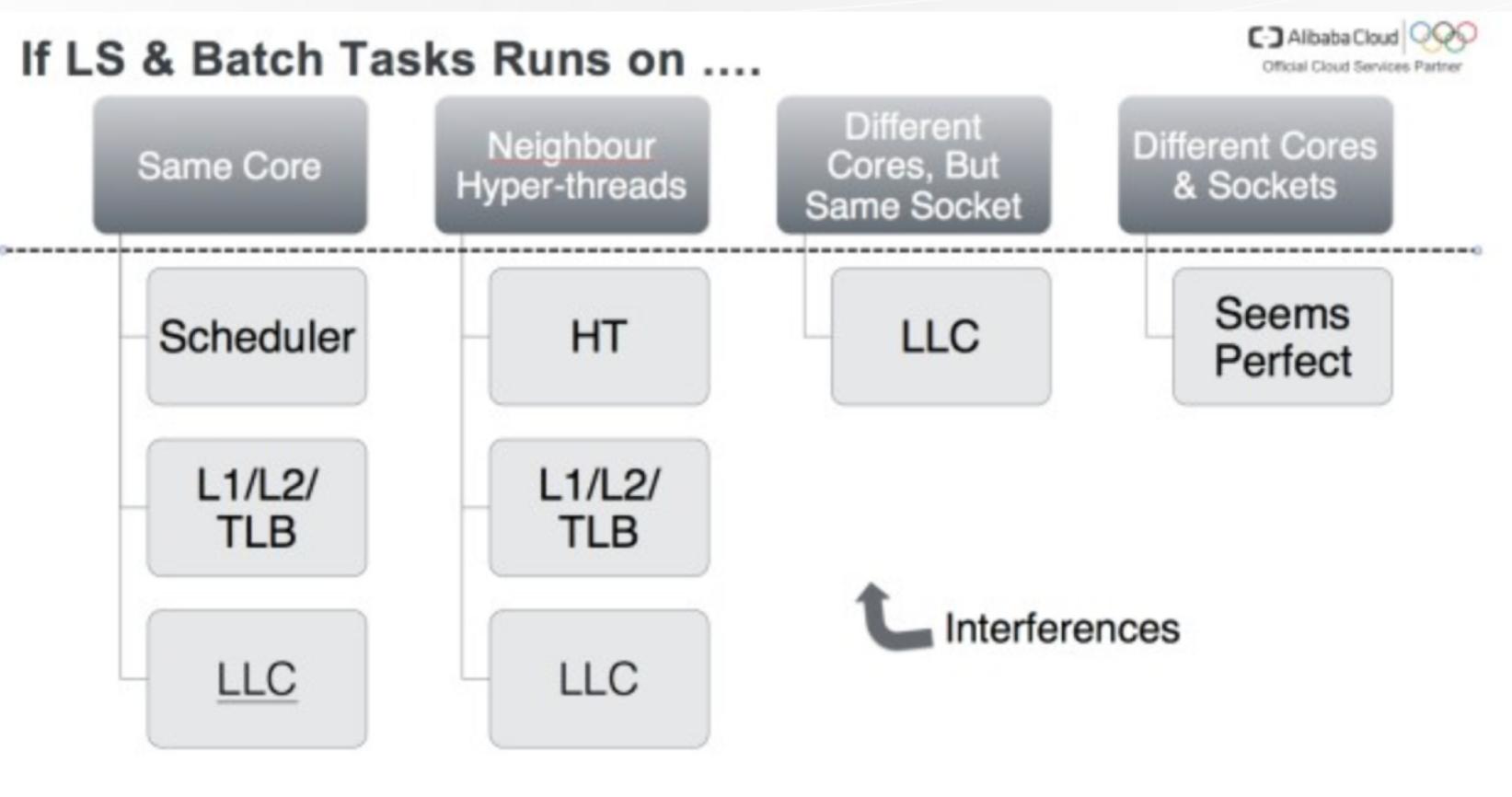


resource sharing



Isolation

Interferences



CPU Scheduler interference

- Affect two tasks on the same core
- Task preempt (Alikernal feature)
 - Online service can preempt offline job more easily
- cpu.shares
 - Online service has more shares than offline job

HT Interference

- 2 HT on two logic cores share L1-d/L1-i/L2/TLB Cache
- Eliminate the need to BIND the online service and offline jobs to exclusive cores

Core	Neighbor HT		
Idle	Idle	job run	
Idle	Service run	Job not run	
Job run	Service wakeup	Job leave	

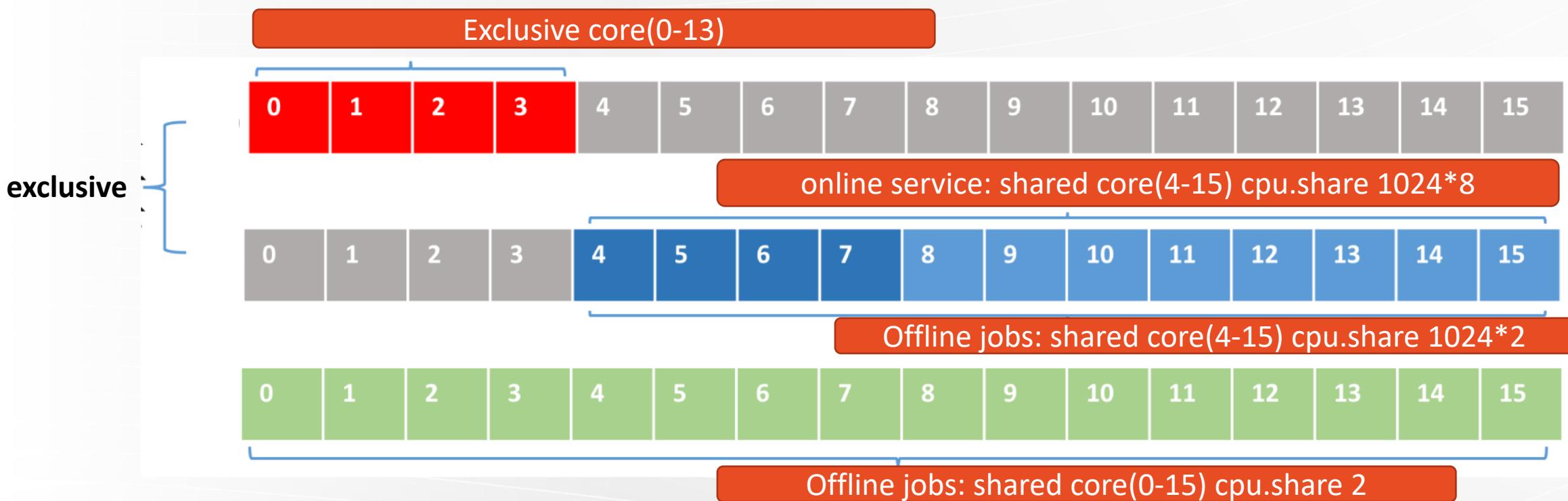
LLC(last-level cache) Interference

- Shared in the same CPU socket
- Intel CAT for LLC isolation
- 20 cache unit (16 for service, and 4 for job)

Memory Isolation

- Memory.priority
 - Online service memory reclaim will be less than offline job
- memory.use_priority_oom
 - Kill lower priority offline jobs when OOM
- memcg wmark_ratio
 - Similar to kswapd, when memory hits wmark_ratio, start memory reclamation

CPU

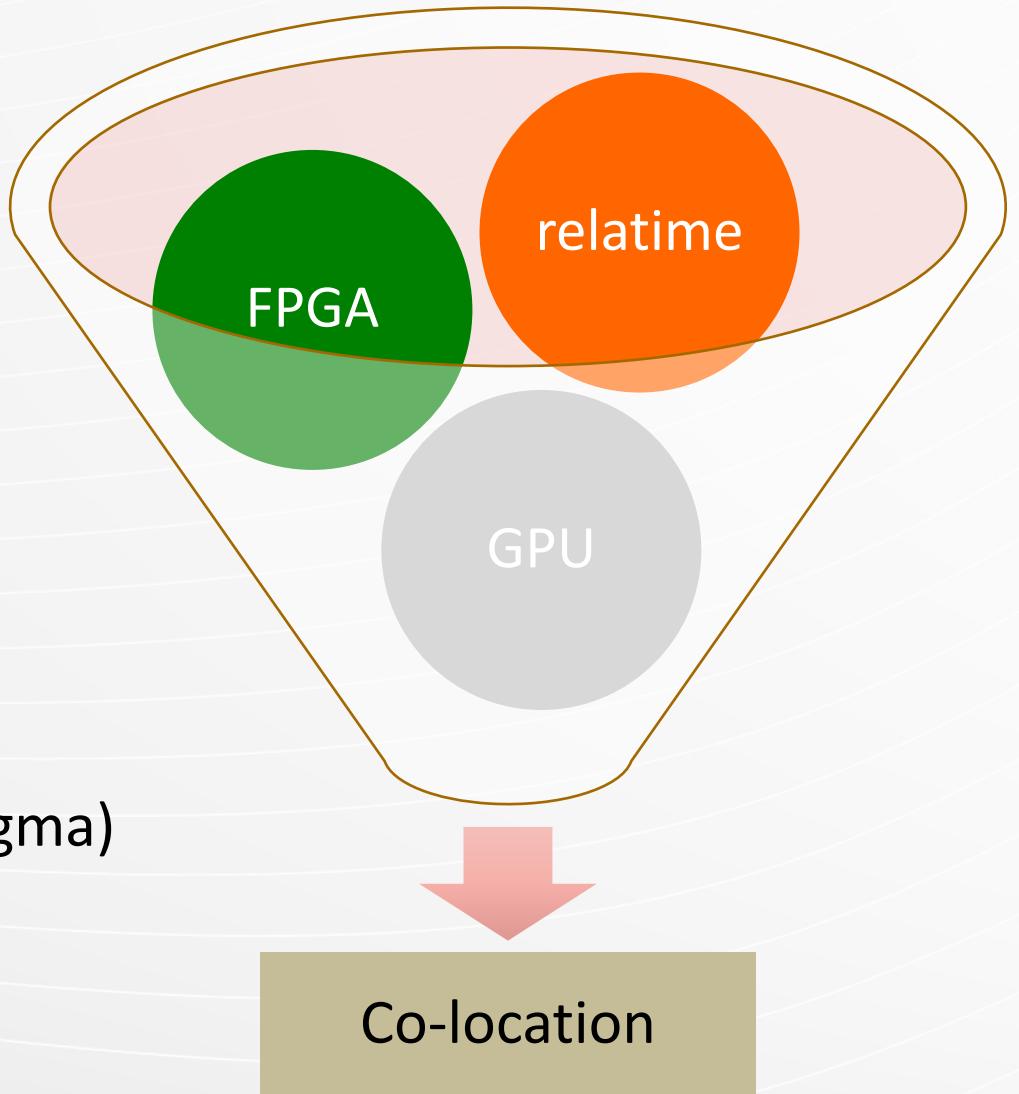




Future

Larger scale

- More resource dimension
- Expand Alibaba internal co-location scale (Fuxi & sigma)



Intelligent auto-rebalancing

- Learn online service and offline jobs and do prediction to auto-rebalance
- Adjust service and job cluster scale more efficiently on sales day events.

More

- Improve stabilities at large scale and high pressure.
- Consolidate the colocation technologies apply those in kubernetes & YARN co-location project