

# JoySim: Simulating Kubernetes Clusters at Scale

Yuan Chen, Qichao Lu, Jiangang Fan, Guang Zhou, Haifeng Liu

JD.com



# About JD.com

**China's largest online and overall retailer and biggest Internet company by revenue**

- 300 million+ active users
- 2018 revenue: \$67.2 billion

**China's largest e-commerce logistics infrastructure and fulfillment network**

- 550+ warehouses
- Covering 99% of population
- Standard same-and next day delivery

**First Chinese internet company to make the Fortune Global 500**

**Strategic partnerships** *Tencent* Walmart  Google



# JD Retail Technical Infrastructure Group

**Provide and manage hyperscale containerized infrastructures and platforms for all JD services**

- One of the earliest adopters of Kubernetes
- Run large scale Kubernetes clusters in production
- CNCF Platinum Member
- 2018 CNCF End User Award

**Support 2019 Singles Days Sales**

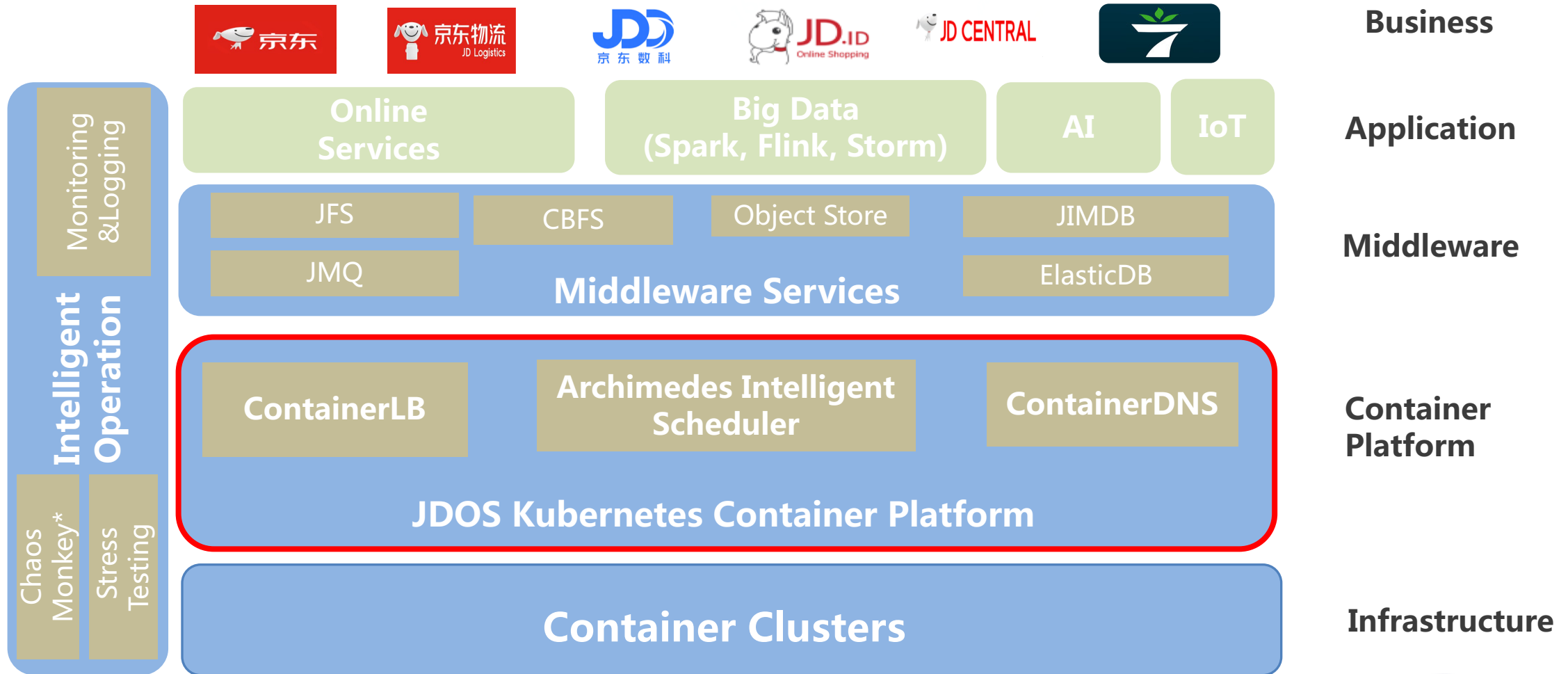
- ¥ 204.4 billion (about \$29 billion) !

<http://tig.jd.com/en>

*"We are thrilled to have JD... By sharing their Kubernetes experiences and investing directly in the project, JD.com is helping to spread cloud native computing throughout China", -Dan Kohn, Executive Director of the CNCF*



# JD Container Platform



# Why Do We Need A Simulator?

- Evaluate new configurations and features before deploying in a production Kubernetes cluster.
  - **Scheduling optimization performance and quality**
  - Kubernetes Master (e.g., APIServer, ETCD) performance and scalability
- **Time consuming and costly to perform such evaluations at scale.**
- Simulation is a useful technique and tool!



# Cluster Simulation Tools

	Complex Scenarios	Performance	Scheduling Quality	Visualization
Kubernetes	●	●	✘	✘
Mesos	●	✘	✘	●
Yarn	✘	✘	●	●

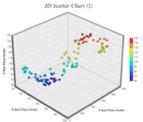
# Key Requirements



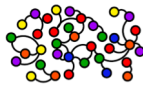
Simulate large-scale K8s clusters with a small amount of resources



Mimic real behaviors with a reasonable amount of confidence



Enable to evaluate scalability, performance and quality



Support realistic and complex scenarios



Comprehensive monitoring and visualization

# Outline

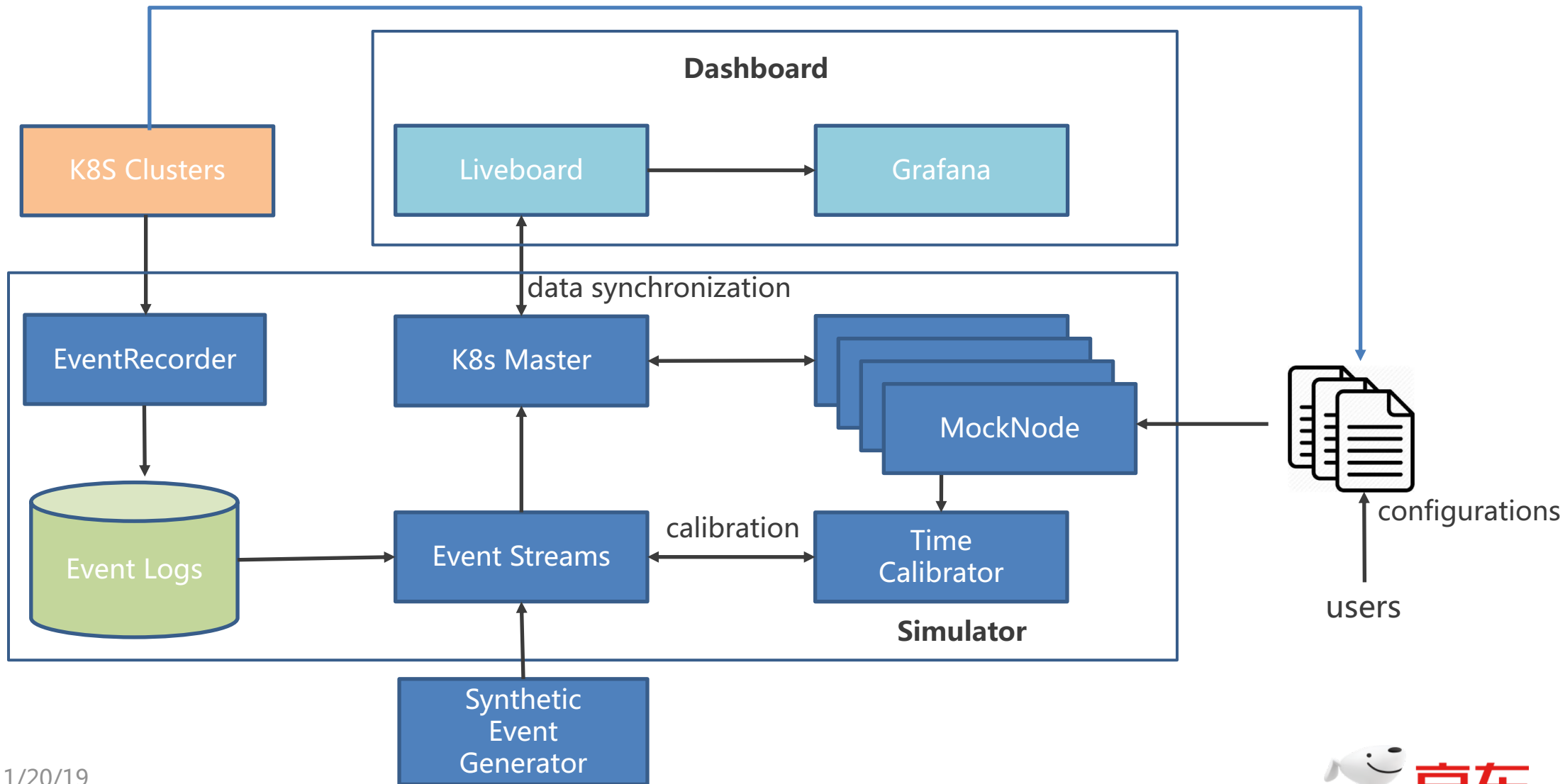
- Background and Introduction
- JoySim Design and Implementation
- Use Cases
- Conclusion
- Demo



# JoySim: A Kubernetes Cluster Simulator

- Run real K8s masters: APIServer, Scheduler, ETCD, ControllerManager.
- Simulate nodes (Kublet) using lightweight MockNodes.
- Configuration information from real Kubernetes clusters or configuration files.
- Replay events from real K8S clusters or customized scenarios.
- Comprehensive monitoring: resource utilization, scheduling traces, performance.
- Easy to deploy, configure, manage and scale using K8s.

# Architecture



# Events

- Collection from real Kubernetes production environments
  - Pods creation, deletion, preemption, rescheduling ...
  - Collection agent: **Telegraf**
  - Storage: **InfluxDB**
- Customized scenarios generated by a synthesizer
  - Cluster and node template: size and resource capacities
  - Pod template: request type and size
  - Event template: request arrival rate, inter-arrival time (e.g., deterministic, random, exponential distribution)

```
1 {
2   "type": "ADDED",
3   "object": {
4     "kind": "Configmap",
5     "metadata": {
6       "name": "newjdoslog-ht03",
7       "namespace": "group-manager",
8       "selfLink": "/api/v1/namespaces/group-manager/configmaps/newjdoslog-ht03",
9       "uid": "41ec6f68-54fc-11e8-9dec-f898efe7ce8c",
10      "resourceVersion": "445984",
11      "creationTimestamp": "2018-05-11T09:18:24Z",
12      "labels": {
13        "app": "newjdoslog",
14        "group": "ht03",
15        "system": "logsystem"
16      }
17    },
18    "timestamp": "2018-05-11T09:18:24Z"
19  }
20 }
21
22 {
23   "type": "ADDED",
24   "object": {
25     "kind": "Pod",
26     "metadata": {
27       "name": "6e7d5fbb-2e8d-4340-8437-23c72f7bcb8d",
28       "namespace": "capjdos",
29       "selfLink": "/api/v1/namespaces/capjdos/pods/6e7d5fbb-2e8d-4340-8437-23c72f7bcb8d",
30       "uid": "44ecaf13-571f-11e8-9dec-f898efe7ce8c",
31       "resourceVersion": "2905653097",
32       "creationTimestamp": "2018-05-14T02:34:04Z",
33       "labels": {
34         "app": "3056",
35         "app_id": "3056",
36         "env": "HT_CM0",
37         "group": "ht03",
38         "novaname": "3056-1526265234352-2003843108-f09ac252",
39         "novauuid": "6e7d5fbb-2e8d-4340-8437-23c72f7bcb8d",
40         "region": "ht03",
41         "res.lvm": "true",
42         "system": "capjdos"
43       },
44       "annotations": {
45         "containerEnvHash": "e047303e",
46         "containerResourceHash": "47ada25c",
47         "containerVolumeMountsHash": "7853f5cb",
48         "kubernetes.io/config.seen": "2018-08-22T20:25:10.883601119+08:00",
49         "kubernetes.io/config.source": "api",
50         "novaflavor": "c24caed4-5b06-41a8-9fdb-88f0a08ccaae",
51         "novaimage": "19954b7f-dc09-44bf-a643-82ffb5d14b0c",
52         "novazone": "HT_CM0"
53       }
54     },
55     "spec": {
56       "containers": [
57         {
58           "name": "container",
59           "image": "is.ht1.n.jd.local/cap2/base:search_bd_work_jdk6_20160519.171214",
60           "resources": {
61             "limits": {
62               "cpu": "8",
63               "memory": "16Gi"
64             },
65             "requests": {
66               "cpu": "2",
67               "memory": "10752Mi"
68             }
69           },
70           "terminationMessagePath": "/dev/termination-log",
71           "terminationMessagePolicy": "File",
72           "imagePullPolicy": "IfNotPresent",
73           "stdin": true,
74           "tty": true
75         }
76       ]
77     }
78   }
79 }
80 }
```

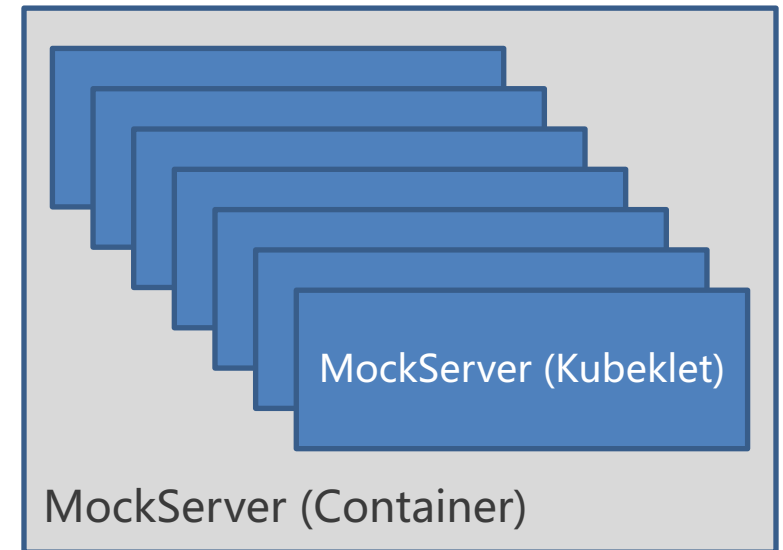
# Node Simulation

## MockNode

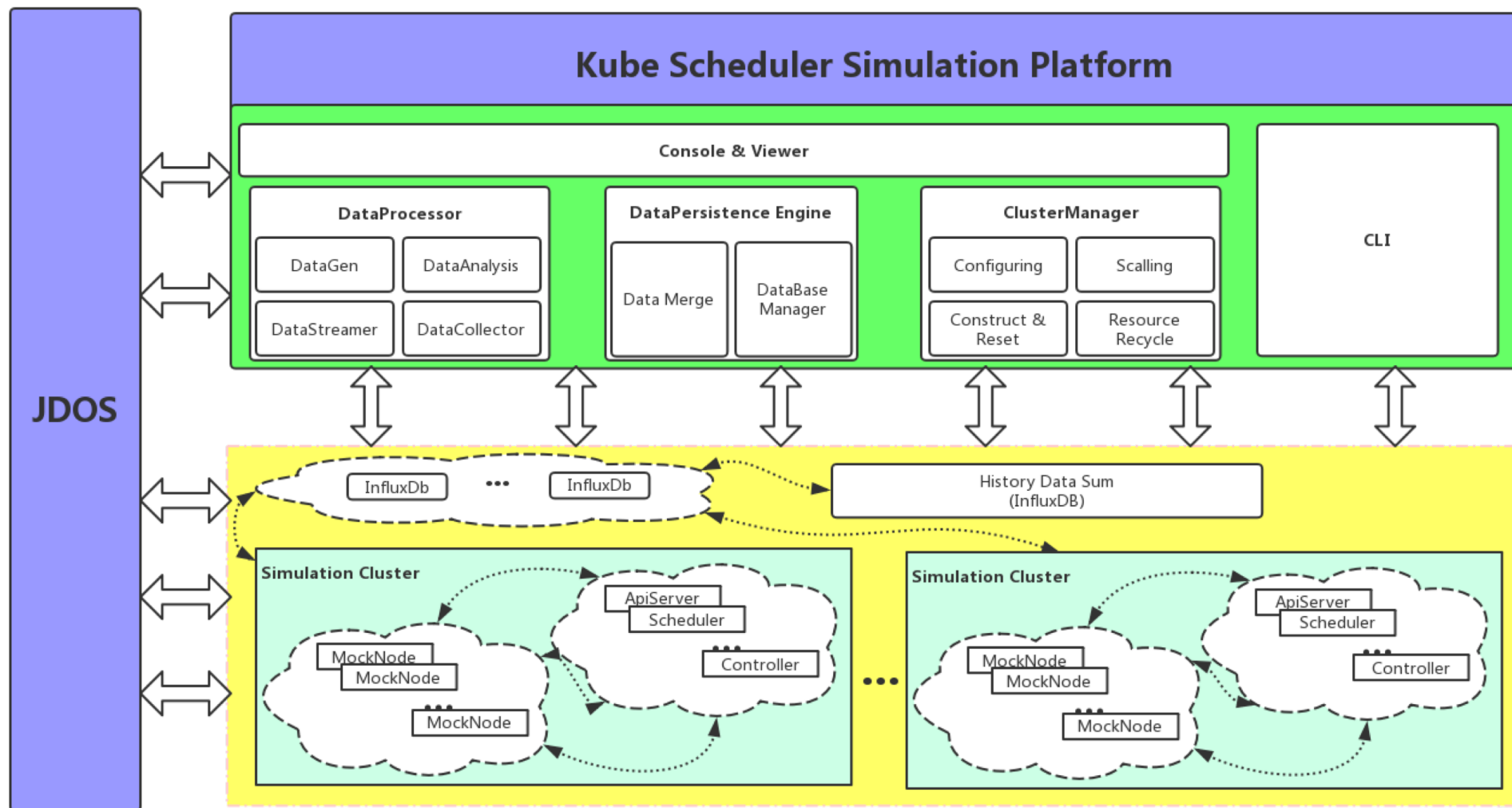
- Simulate Kubelet in K8s
- Resource template: customized memory, CPU disk...
- Watch pod scheduling events and update statuses
- Report resource usage

## MockServer

- Run in containers
- Deployed and managed by K8s
- **Simulate 100+ MockNodes per MockServer (8 cores and 16GB memory)**



# Implementation and Deployment



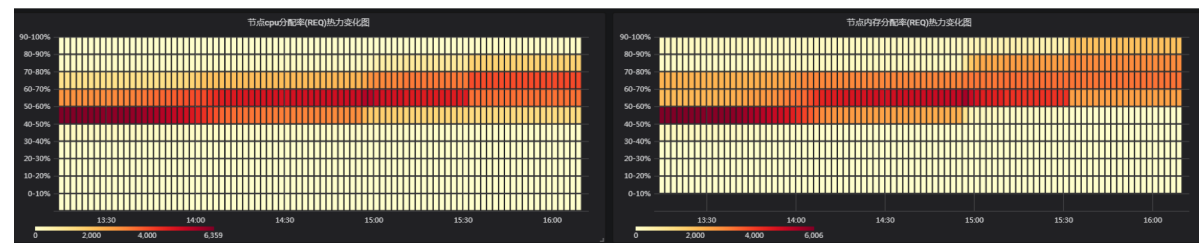
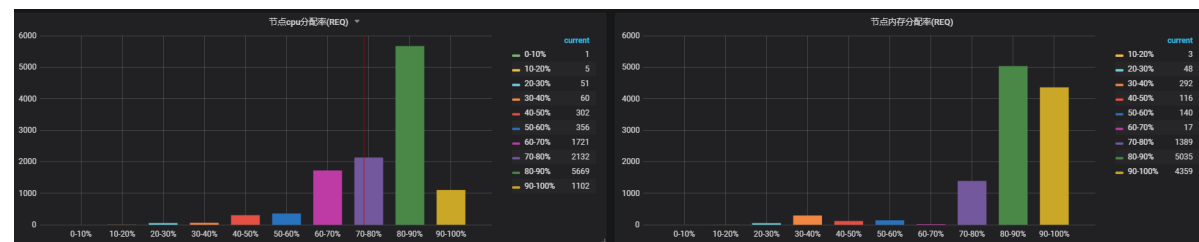
# Monitoring and Visualization

## Metrics

- APIServer performance
- Scheduling performance and quality: node CPU, memory, disk allocation rate

## Reporting and visualization

- Heat map, TS charts, statistics at multiple granularities
- Import and export



Last 30 minutes Refresh every 1m

**Custom range**

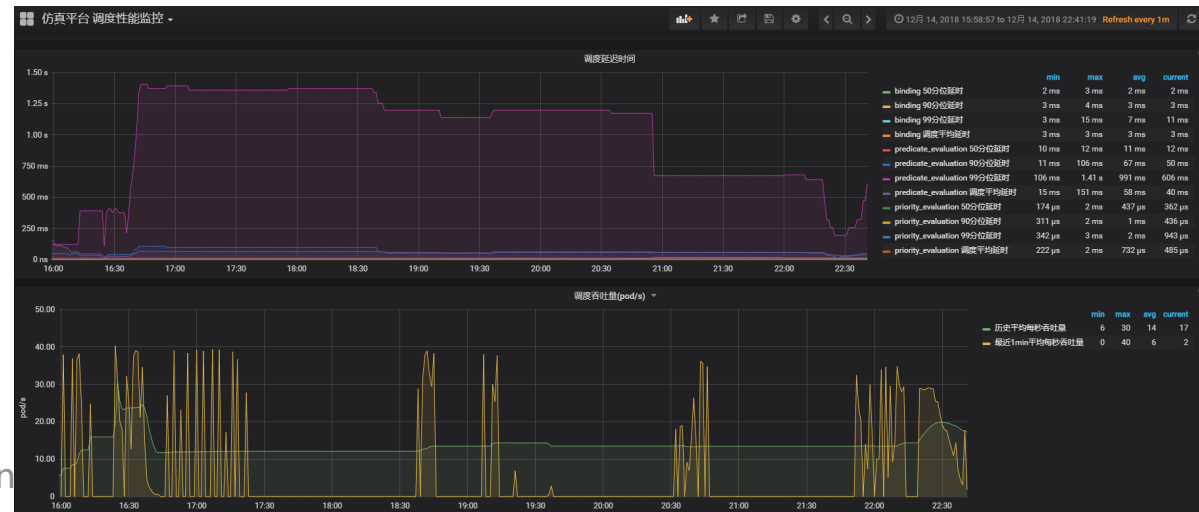
From:

To:

Refreshing every:

**Quick ranges**

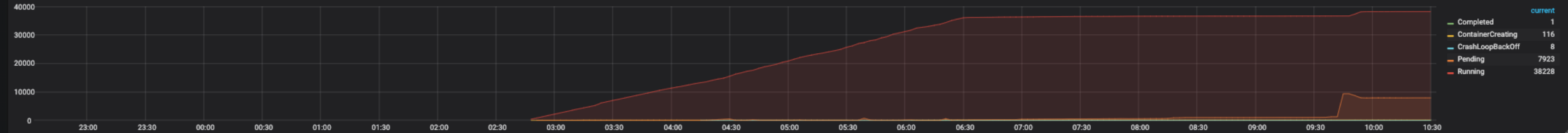
Last 2 days	Yesterday	Today	Last 5 minutes
Last 7 days	Day before yesterday	Today so far	Last 15 minutes
Last 30 days	This day last week	This week	Last 30 minutes
Last 90 days	Previous week	This week so far	Last 1 hour
Last 6 months	Previous month	This month	Last 3 hours
Last 1 year	Previous year	This month so far	Last 6 hours
Last 2 years		This year	Last 12 hours
Last 5 years		This year so far	Last 24 hours



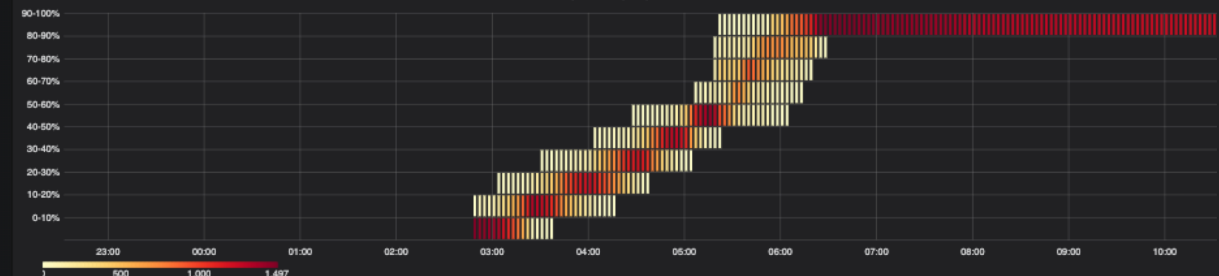
集群概况

region	物理机器数目	容器数目	总CPU(核)	总内存	已分配CPU(核)	已分配内存	CPU分配率(Limit)	内存分配率(Limit)	CPU分配率(REQ)	内存分配率(REQ)
simulation0	1500	38353	90000	366 TiB	89436	324 TiB	263.61%	126.99%	99.37%	88.35%

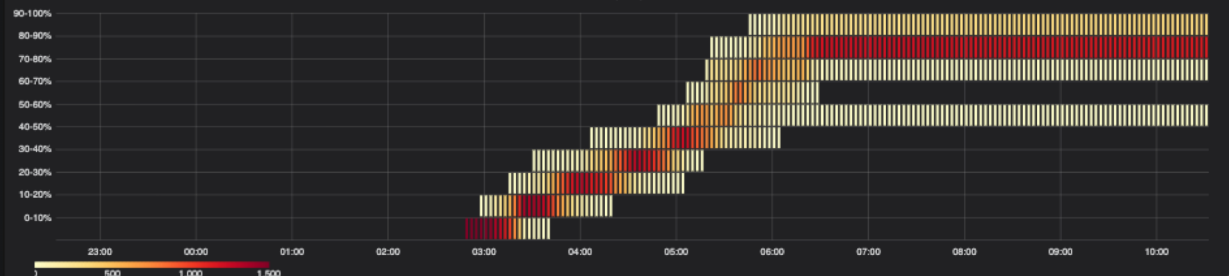
容器状态统计图



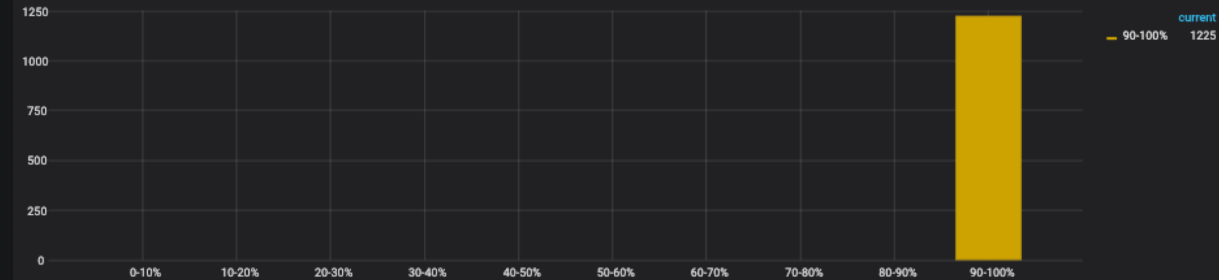
节点cpu分配率(REQ)热力变化图



节点内存分配率(REQ)热力变化图



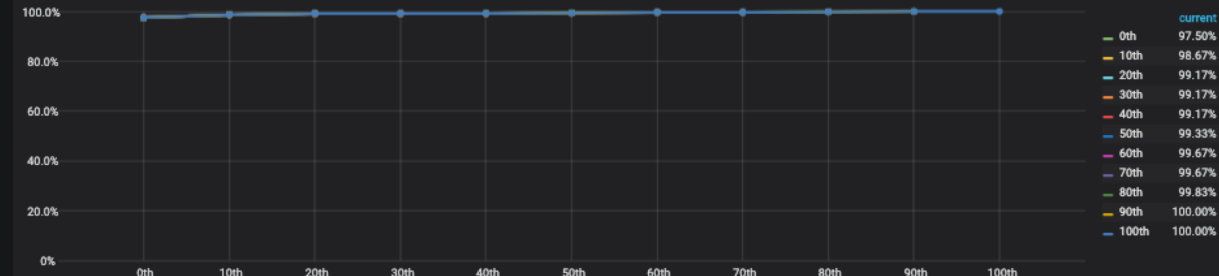
节点cpu分配率(REQ)分位图



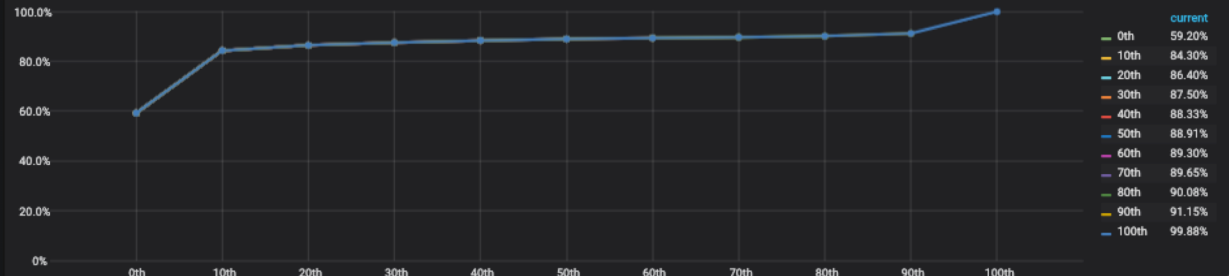
节点内存分配率(REQ)分位图

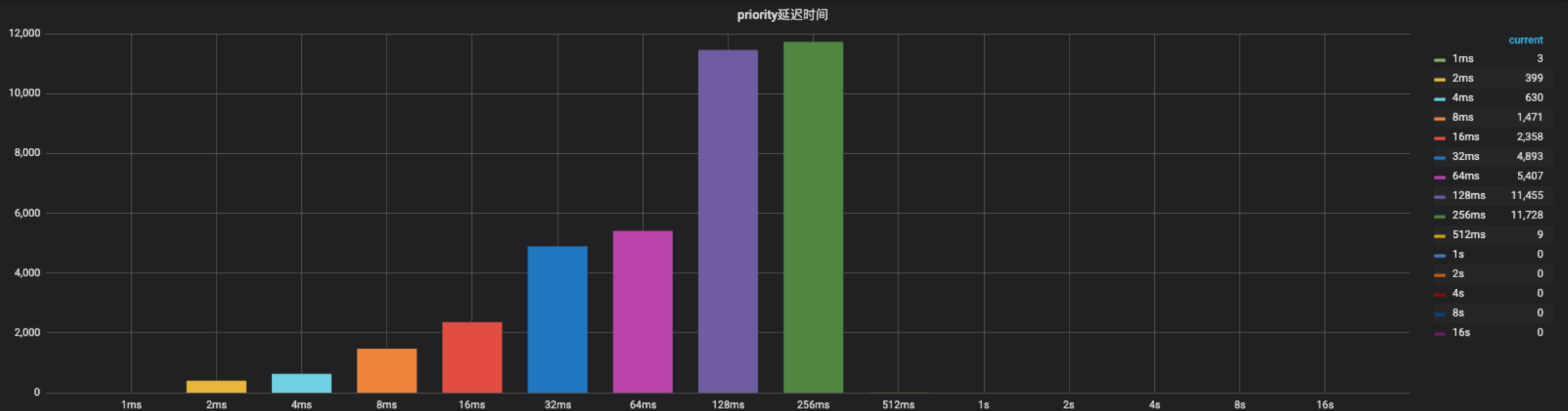
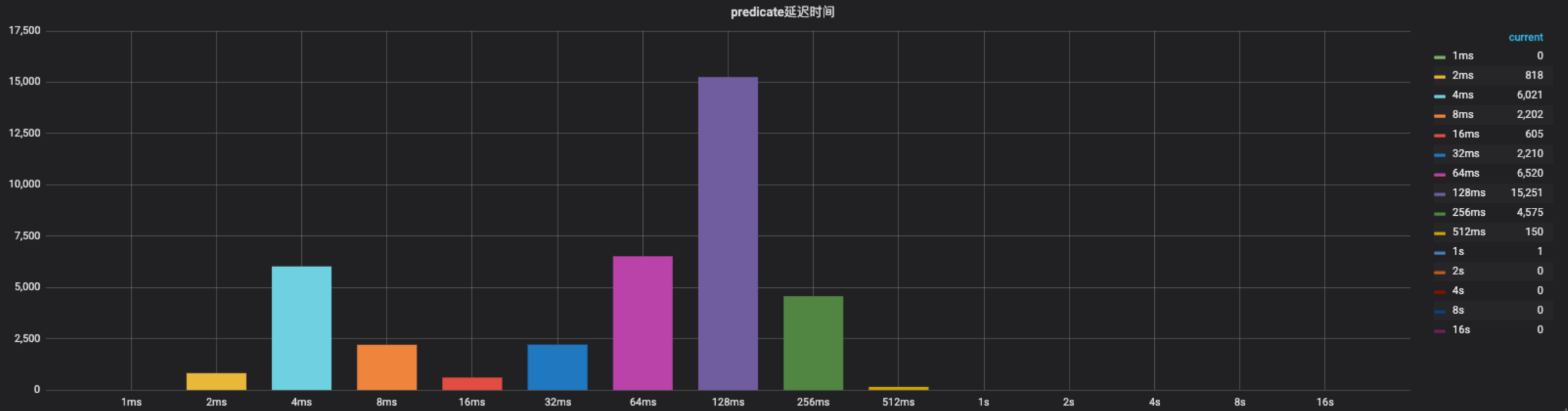


节点cpu分配率分位图



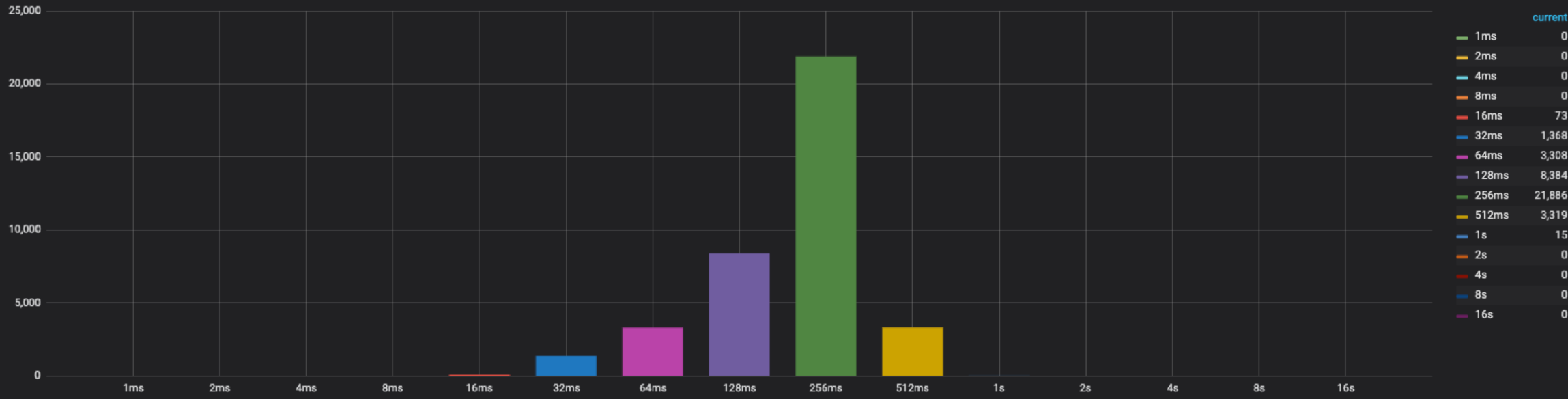
节点内存分配率分位图







调度算法延迟时间



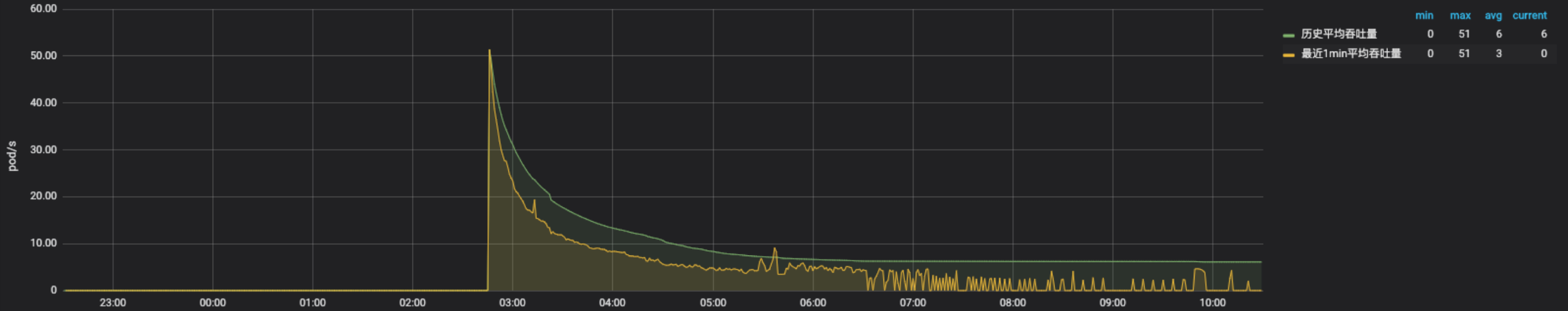
绑定binding延迟时间



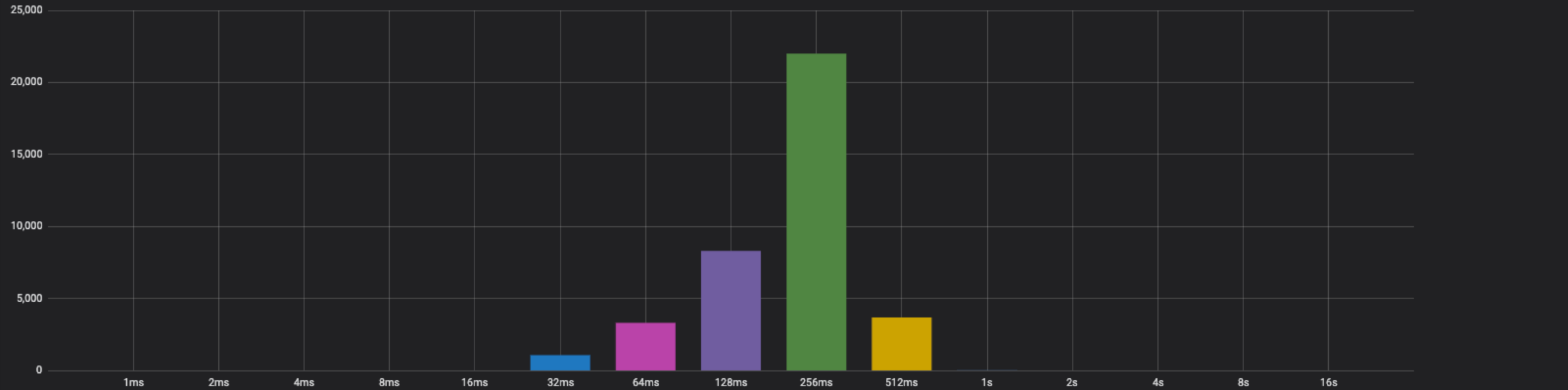


23:00 00:00 01:00 02:00 03:00 04:00 05:00 06:00 07:00 08:00 09:00 10:00

调度吞吐量(pod/s)

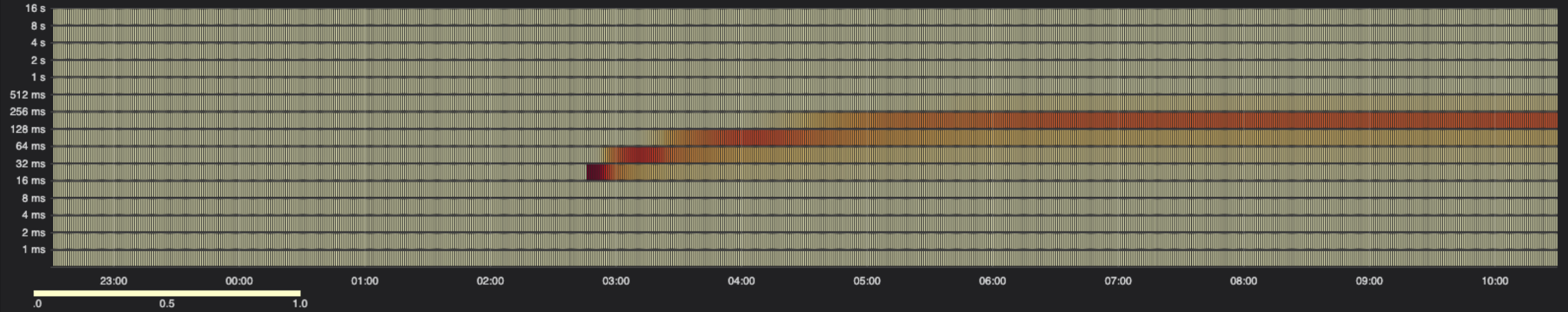


调度算法+绑定延迟时间

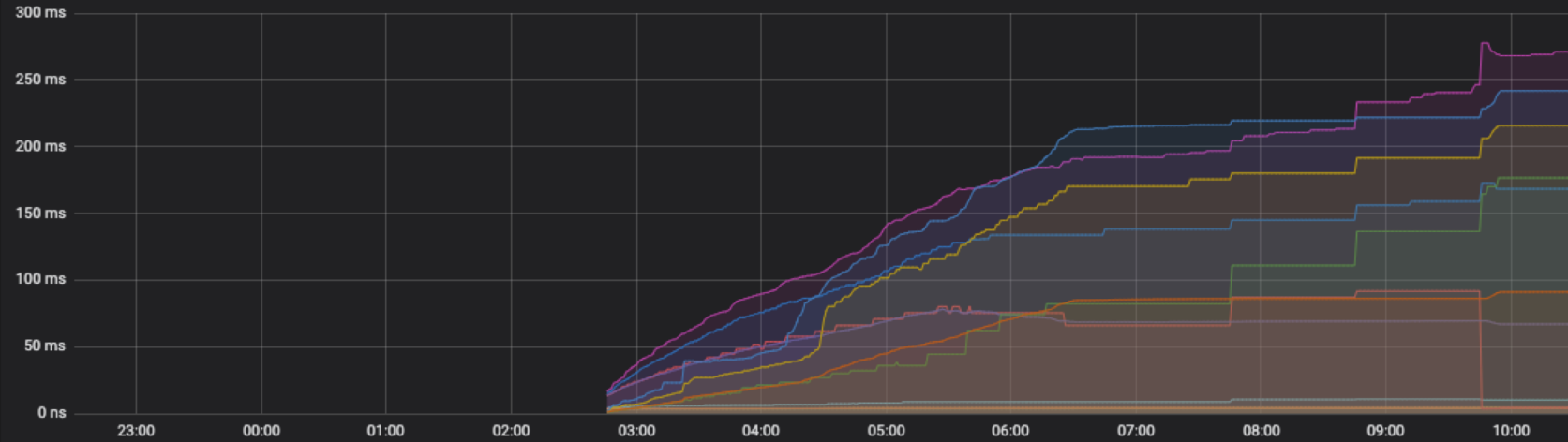




调度算法+绑定延迟时间热力图



调度延迟时间



	min	max	avg	current
binding 50分位延时	3 ms	4 ms	4 ms	4 ms
binding 90分位延时	4 ms	4 ms	4 ms	4 ms
binding 99分位延时	4 ms	11 ms	9 ms	10 ms
binding 调度平均延时	3 ms	4 ms	4 ms	4 ms
predicate_evaluation 50分位延时	3 ms	92 ms	62 ms	4 ms
predicate_evaluation 90分位延时	15 ms	173 ms	122 ms	168 ms
predicate_evaluation 99分位延时	17 ms	278 ms	171 ms	271 ms
predicate_evaluation 调度平均延时	14 ms	78 ms	63 ms	67 ms
priority_evaluation 50分位延时	1 ms	177 ms	78 ms	177 ms
priority_evaluation 90分位延时	1 ms	216 ms	134 ms	216 ms
priority_evaluation 99分位延时	1 ms	242 ms	163 ms	242 ms
priority_evaluation 调度平均延时	1 ms	91 ms	63 ms	91 ms

调度吞吐量(pod/s)

# Outline

- Background and Introduction
- JoySim Design and Implementation
- Use Cases
- Conclusion
- Demo

# Case Study: APIServer Evaluation and Optimization

## • Simulation testbed

- Event traces from multiple K8s clusters
- 1800 simulation containers, 6 Mocknodes (60 CPU cores, 300GB) per container
- **Simulate 11,400 nodes, 200,000+ container**

## • APIServer scalability evaluation

- 6 APIServers to support a 10,000-node cluster

## • ETCD optimization

- Two etcd clusters: one for Pod resources and one for other resource objects
- No read/write/storage bottlenecks

## • Scheduler optimization

- QPS: 2 → 40
- Quality: utilization 2X

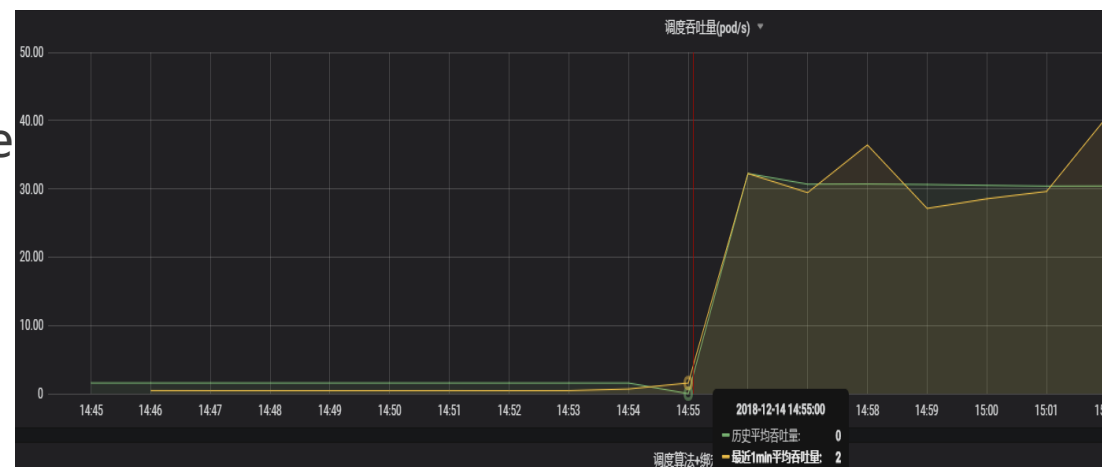


仿真平台 调度质量监控

集群概况

region	物理机器数目	容器数目	总CPU(核)	总内存	已分配CPU(核)	已分配内存	CPU分配率(Limit)	内存分配率(Limit)
simulation	11400	201162	684000	3 PiB	533875	2 PiB	192.02%	107.00%

容器状态统计图



# Case Study: Scheduler Performance Optimization

## Simulation testbed

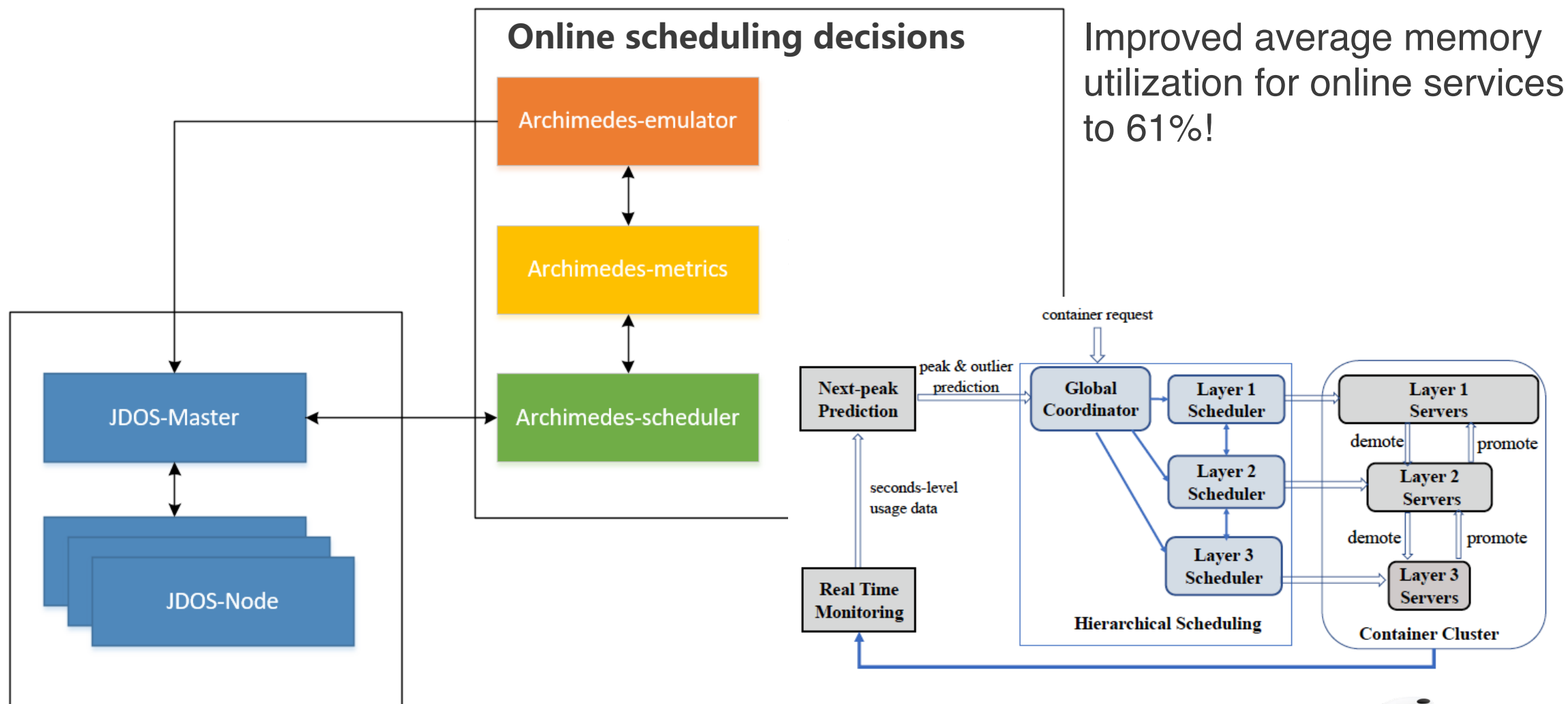
- 800 nodes , 8124 running pods, schedule a job of 500 pods

## Batch job scheduling

Scheduler	Execution Time	Scheduling Mechanisms
K8s Scheduler	16 sec	Per pod scheduling and binding
Volcano Batch Scheduler	4.72 sec	Per pod scheduling and group binding
JDOS Batch Scheduler v1	0.27 sec	Group scheduling and binding
JDOS Batch Scheduler v2	0.15 sec	Adaptive group scheduling and binding



# Case Study: Online Scheduling



# Summary

- **JoySim**: a simulator for large scale Kubernetes cluster simulation.
- Applications for scalability and performance evaluation and optimization of large scale Kubernetes cluster at JD.com.
- Planning to work with CNCF SIG-Scheduling for an open source release.



**Thank you!**



# Contact

Qichao Lu [luqichao@jd.com](mailto:luqichao@jd.com)