

Scalable and Efficient Scheduling of Large Scale Kubernetes Clusters

Yuan Chen



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 unrivalled nationwide fulfillment network

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

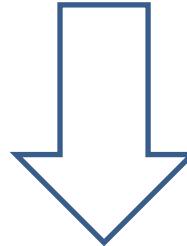
First Chinese internet company to make the Fortune Global 500

Strategic partnerships



Challenges

JD business is complex, growing and changing rapidly by the day.



The technical infrastructure must be Agile, Scalable, and Efficient!

JD Retail Technological Infrastructure Group

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

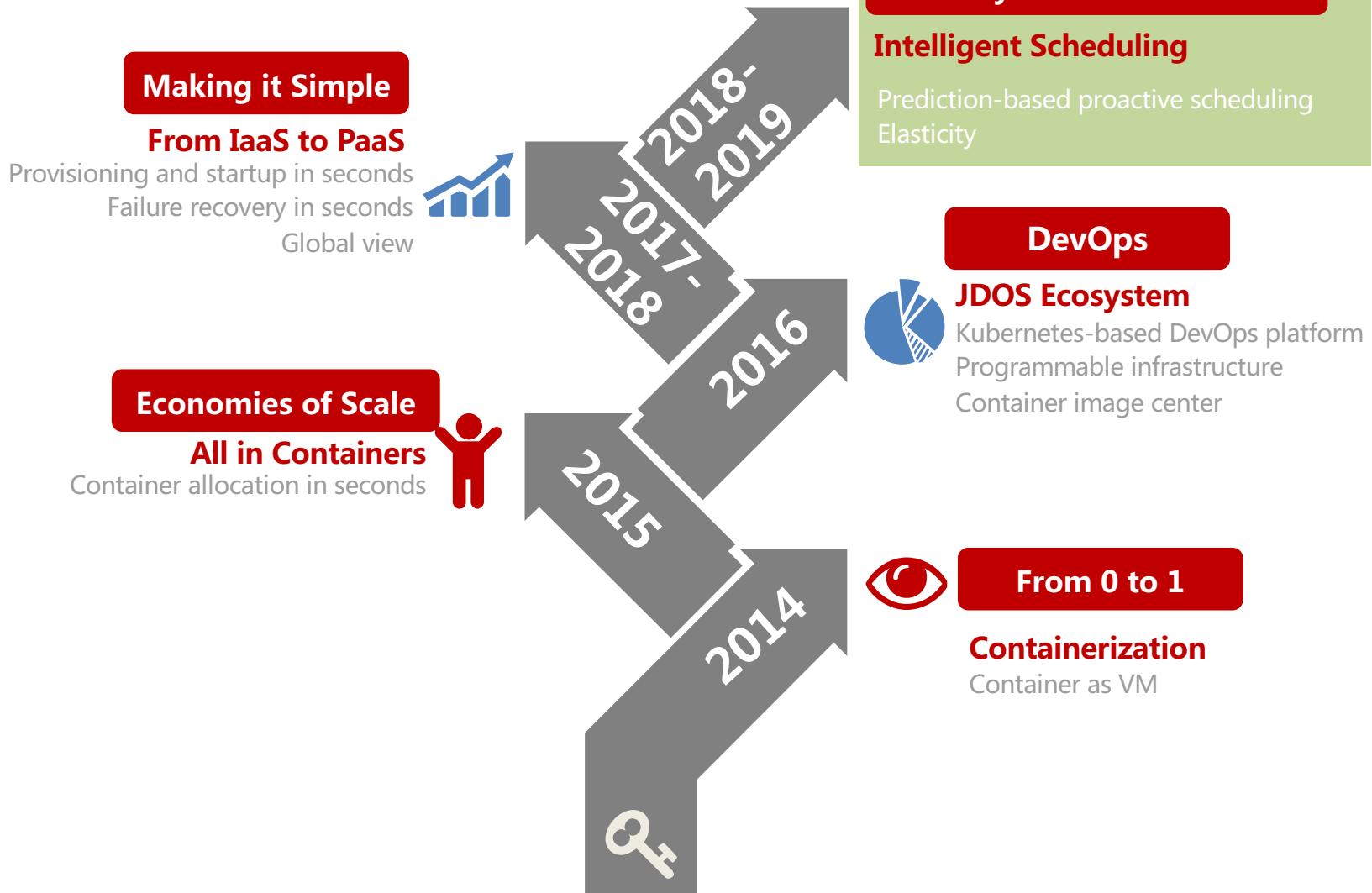
- One of the earliest adopters of Kubernetes
- Run one of the largest Kubernetes clusters in production
- CNCF Platinum Member
- 2018 CNCF End User Award

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

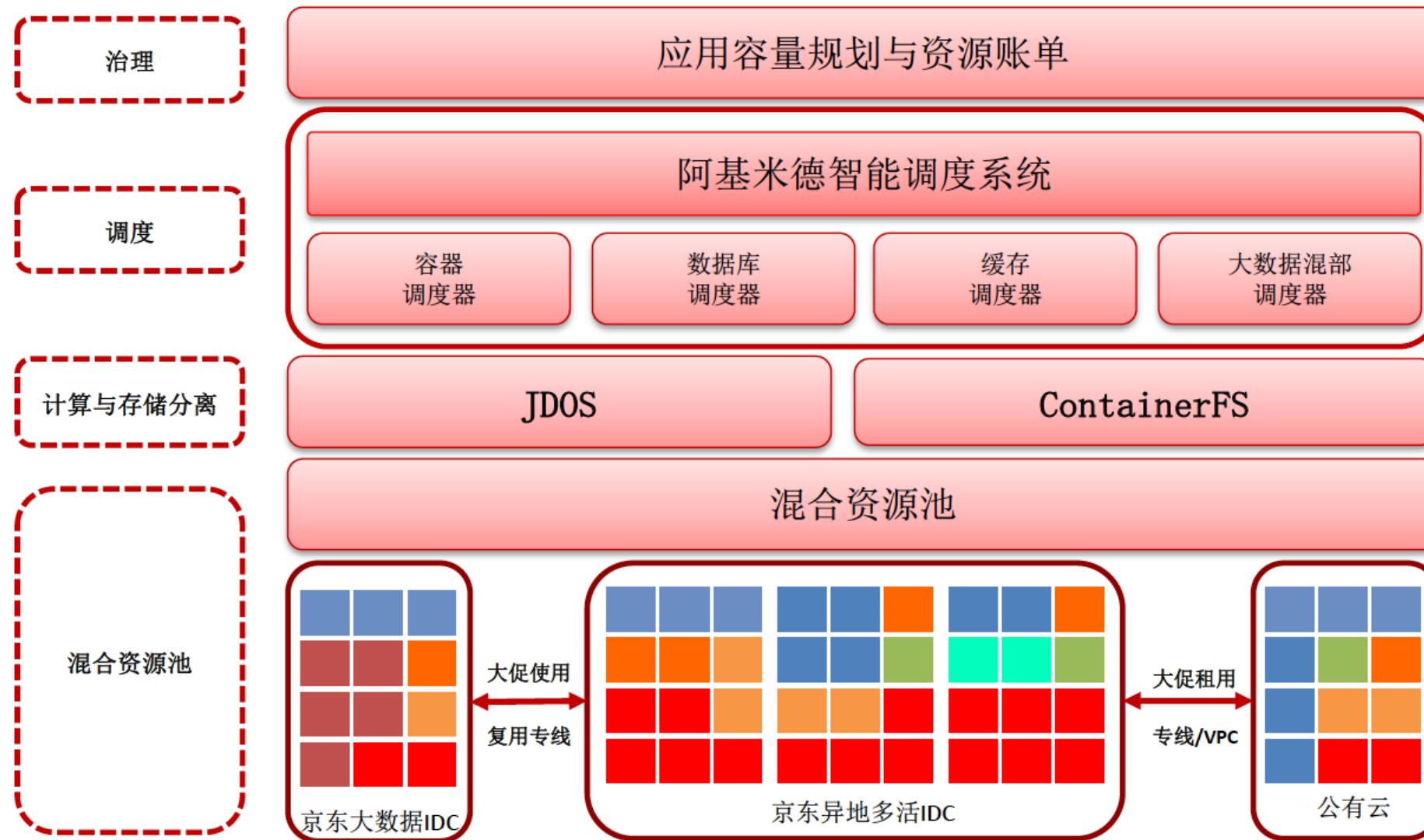
"We are thrilled to have JD... By sharing their Kubernetes experiences and investing directly in the project, JD.com is helping to spread cloud native computing throughout China", -Dan Kohn, Executive Director of the CNCF



Our Journey



JDOS (JD Data Center OS) Architecture



Resource Management: Challenges

Accuracy

- Users cannot accurately estimate their resource needs (often overprovisioning)
- Resource demand of online services varies a lot. **One size does not fit all!**
 - More -> waste
 - Less -> SLO violations

Complexity

- Stateful pods (e.g., databases)
- Multi-dimensional resources
- Heterogeneous workloads: online services, big data, AI, IoT

Scalability

- A single cluster with 10K+ nodes and 100K+ pods

Archimedes: Scalable and Efficient Scheduling

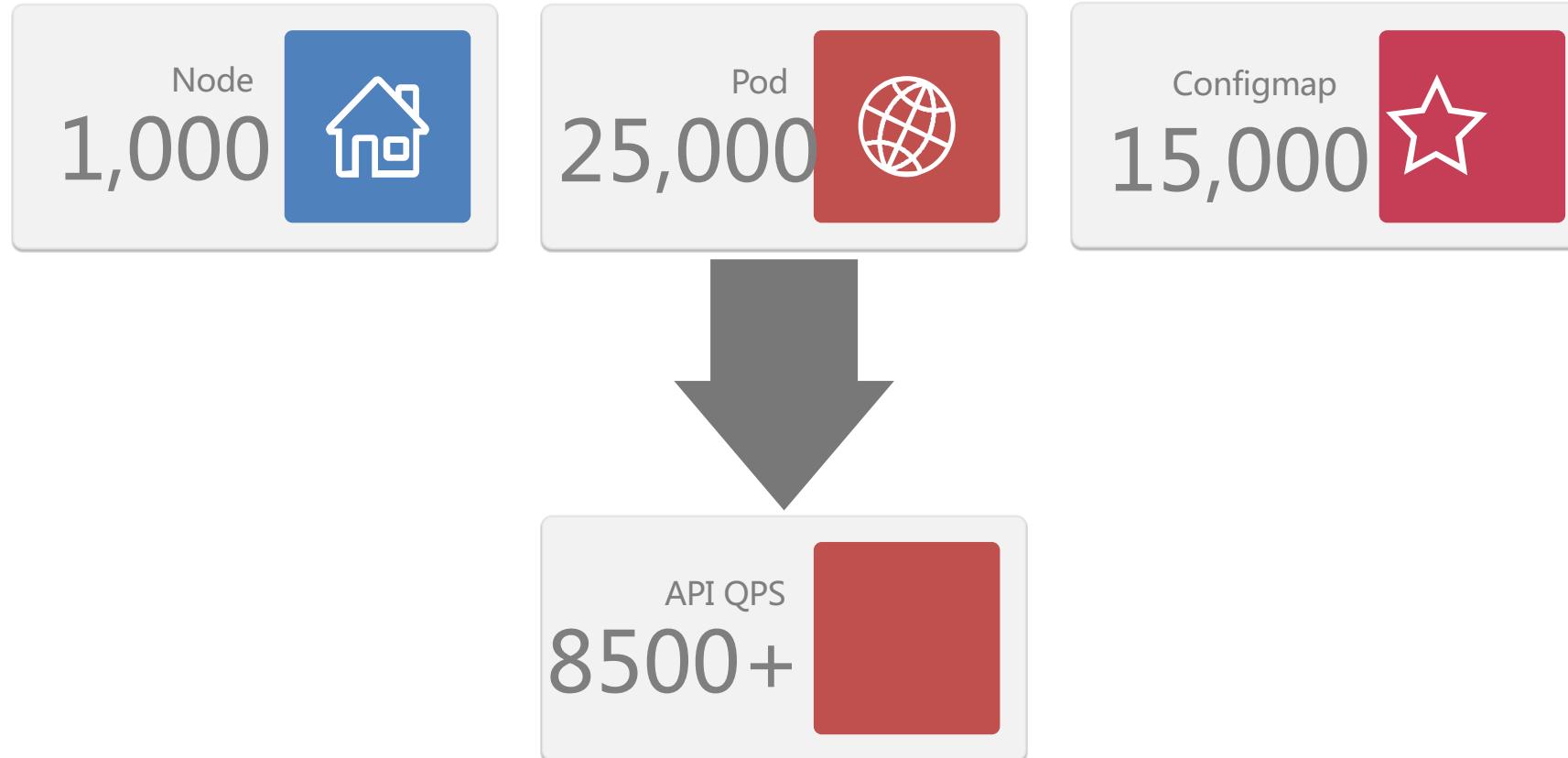
Performance and Scalability

- Customization: API servers, etcd, scheduler, image
- Advanced optimization
 - Preemption and admission control
 - Group scheduling and binding
 - Adaptive scheduling

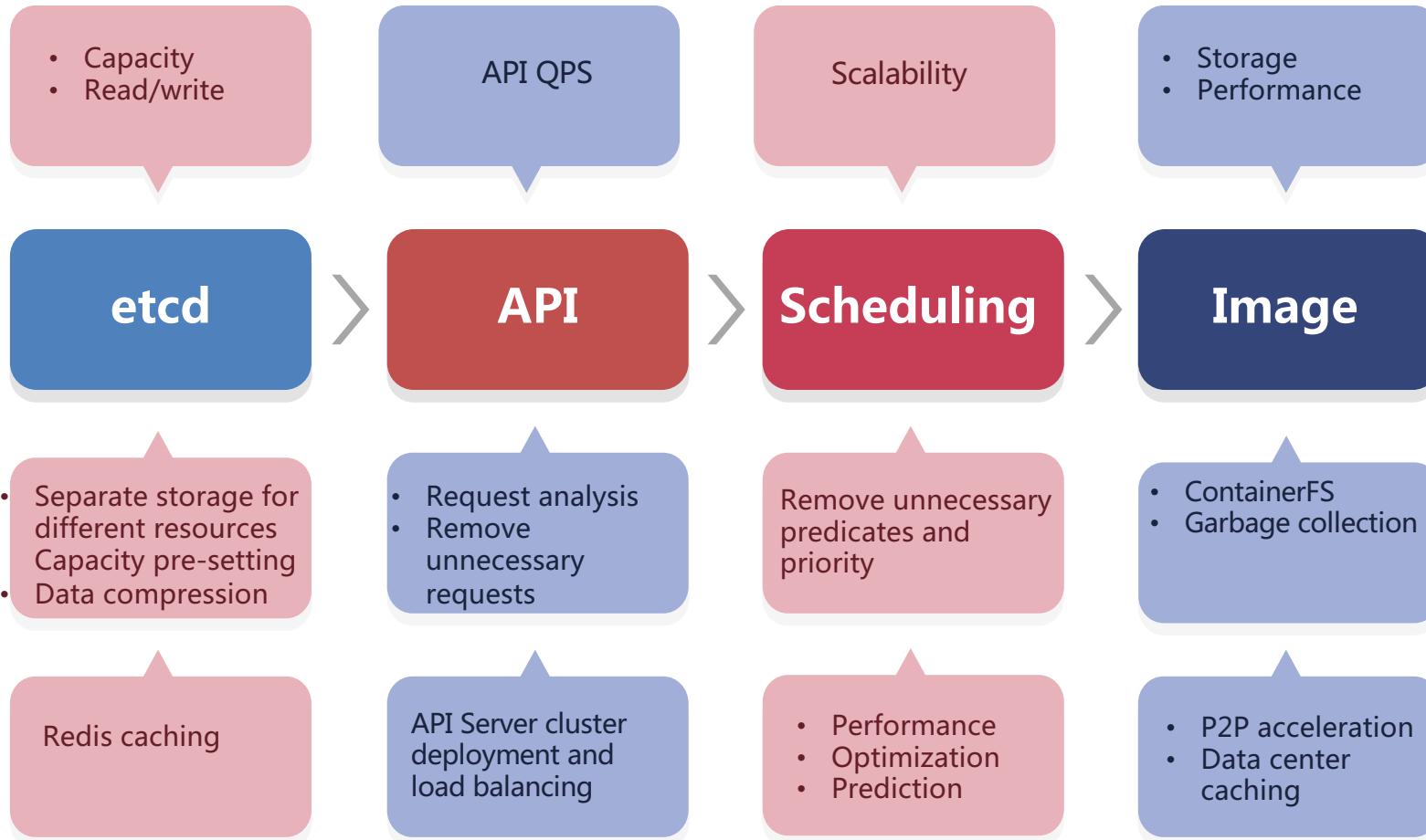
Quality and Resource Efficiency: data-driven intelligent scheduling

- Historical usage analysis and prediction-driven scheduling
- Mixed workload placement
- Auto scaling with hybrid resource pools

A Typical JDOS Cluster



Bottleneck Analysis and Optimization



Customization and Optimization

etcd capacity monitoring
compression



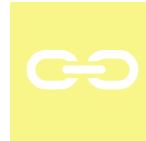
APIserver request analysis
configmap requests: 99%

different etcd for
different resource-types



GET requests
Redis cache for etcd

heartbeat optimization

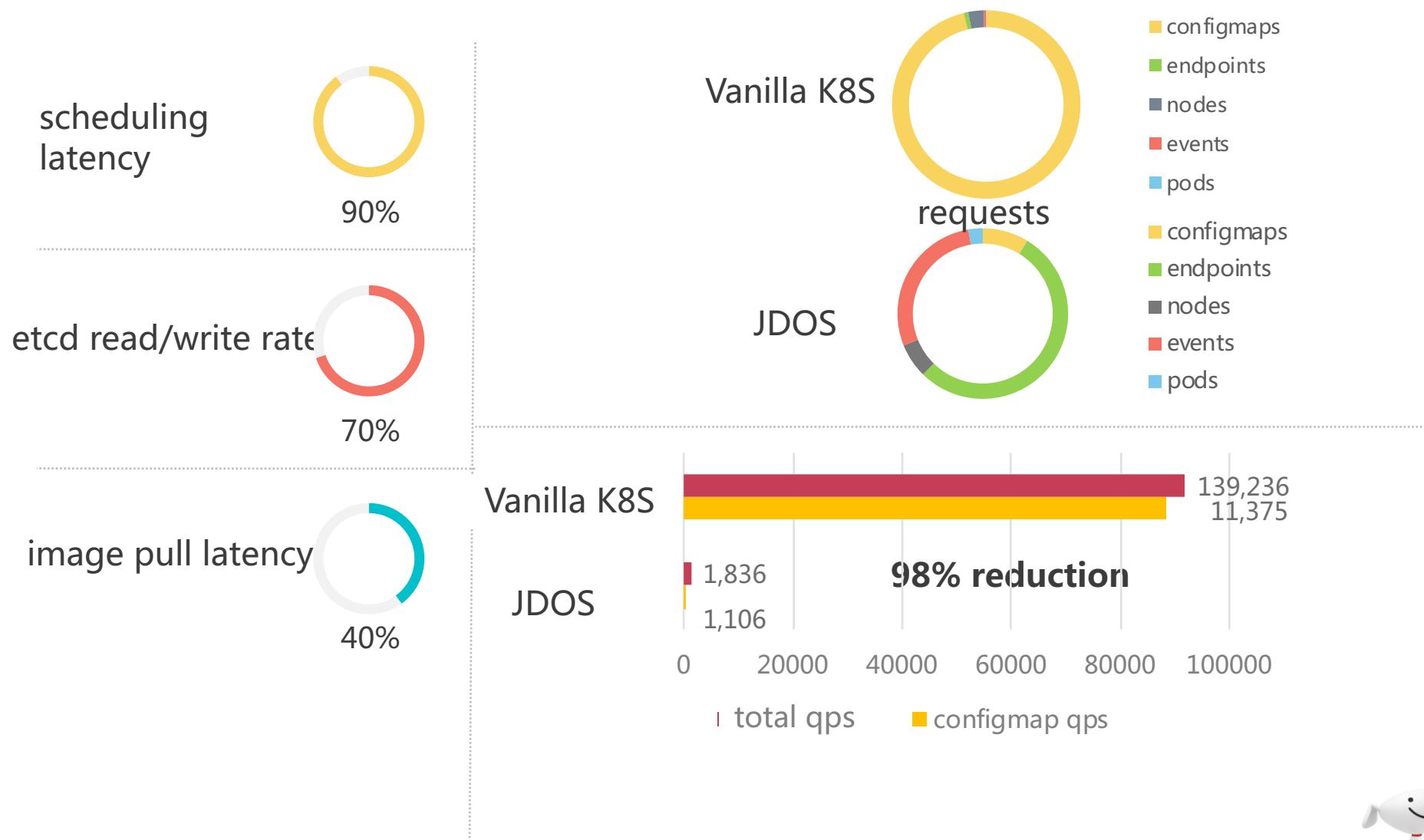


APIServer load balancing

- **Separate storage for different resource types**
- **Compression**
- **Redis cache**

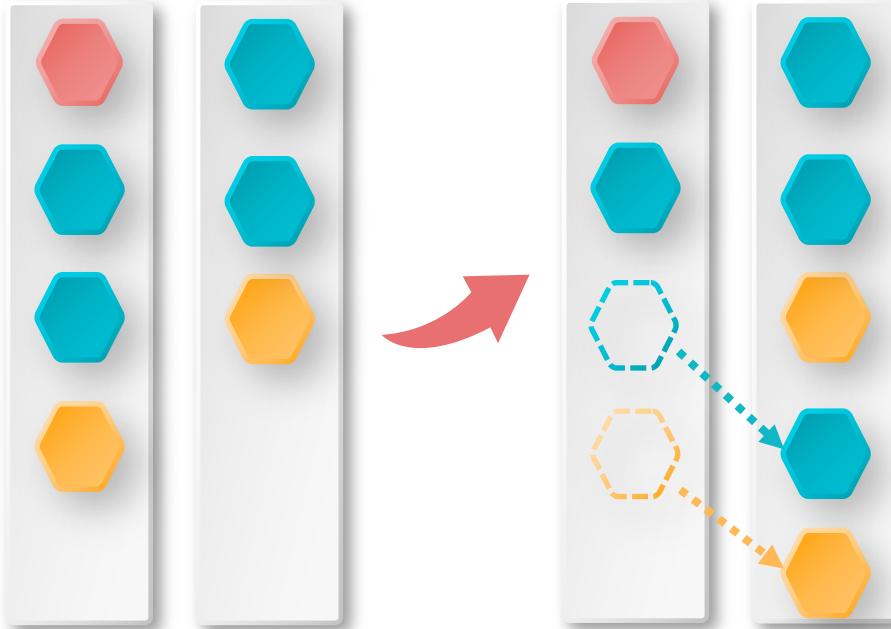
- **Remove unnecessary requests**
- **Load balancing**

Results



Preemption and Admission Control

Preemption Optimization



- Resource demand prediction
- Priority and SLOs-based preemption
- Proactive preemption
- Preempted apps in scheduling queues
- Overloaded host locked

Admission Control

- Reject low priority request.
- No requests after a certain number of retries.

Group Scheduling and Binding

Problem

- A large number of pod requests from an application (e.g., horizontal auto-scaling, big data batch jobs) cause a significant scheduling delay.

Solution : group scheduling

- Schedule an application group consisting of a number of pods with same/similar resource requirements at the same time.
- Use a representative pod resource requirement to evaluate nodes' priorities to obtain the top k ($k \leq m$) candidates.
- Assign and bind the m pods to the best k nodes simultaneously.

Other solution: Poseidon/Firmament

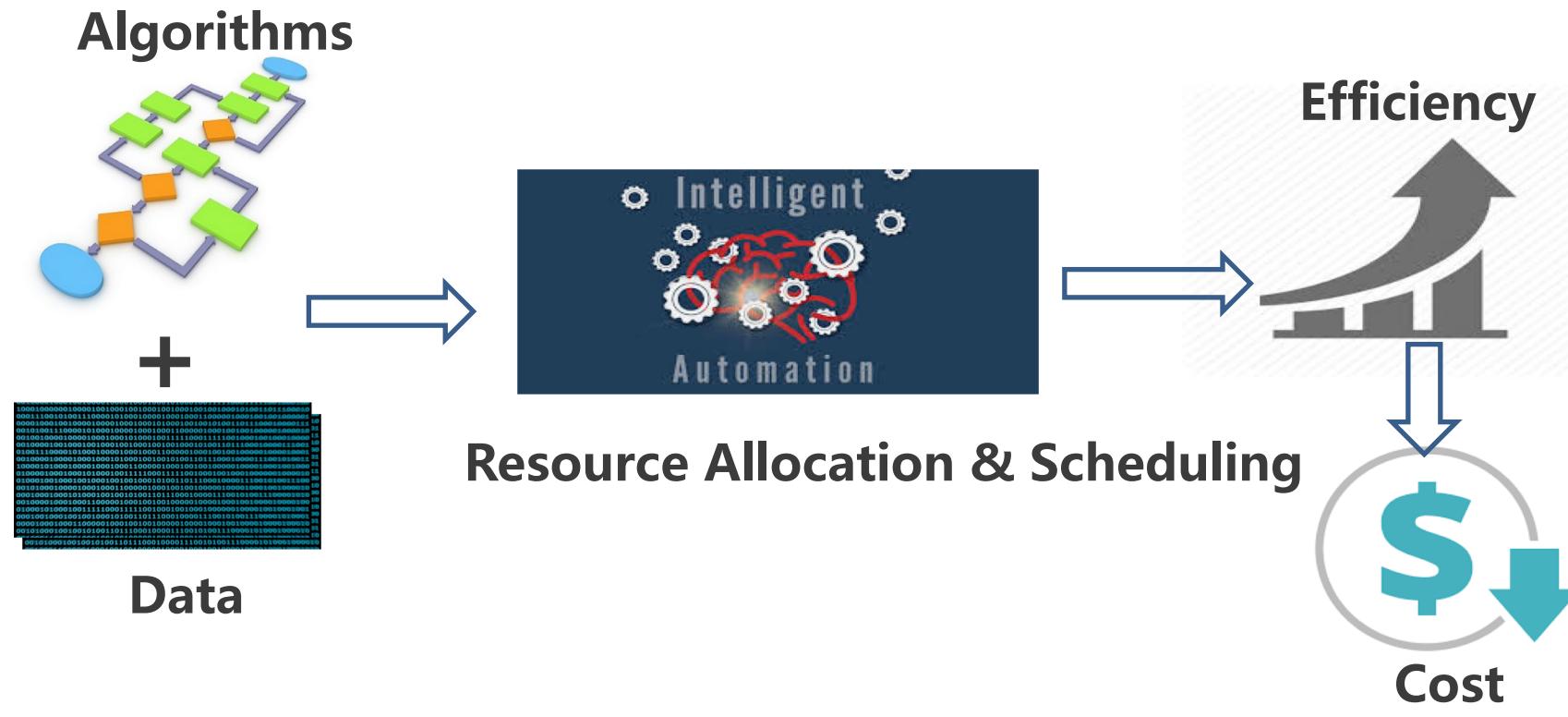
Adaptive Scheduling

Different algorithms and/or parameters for different workloads and scenarios.

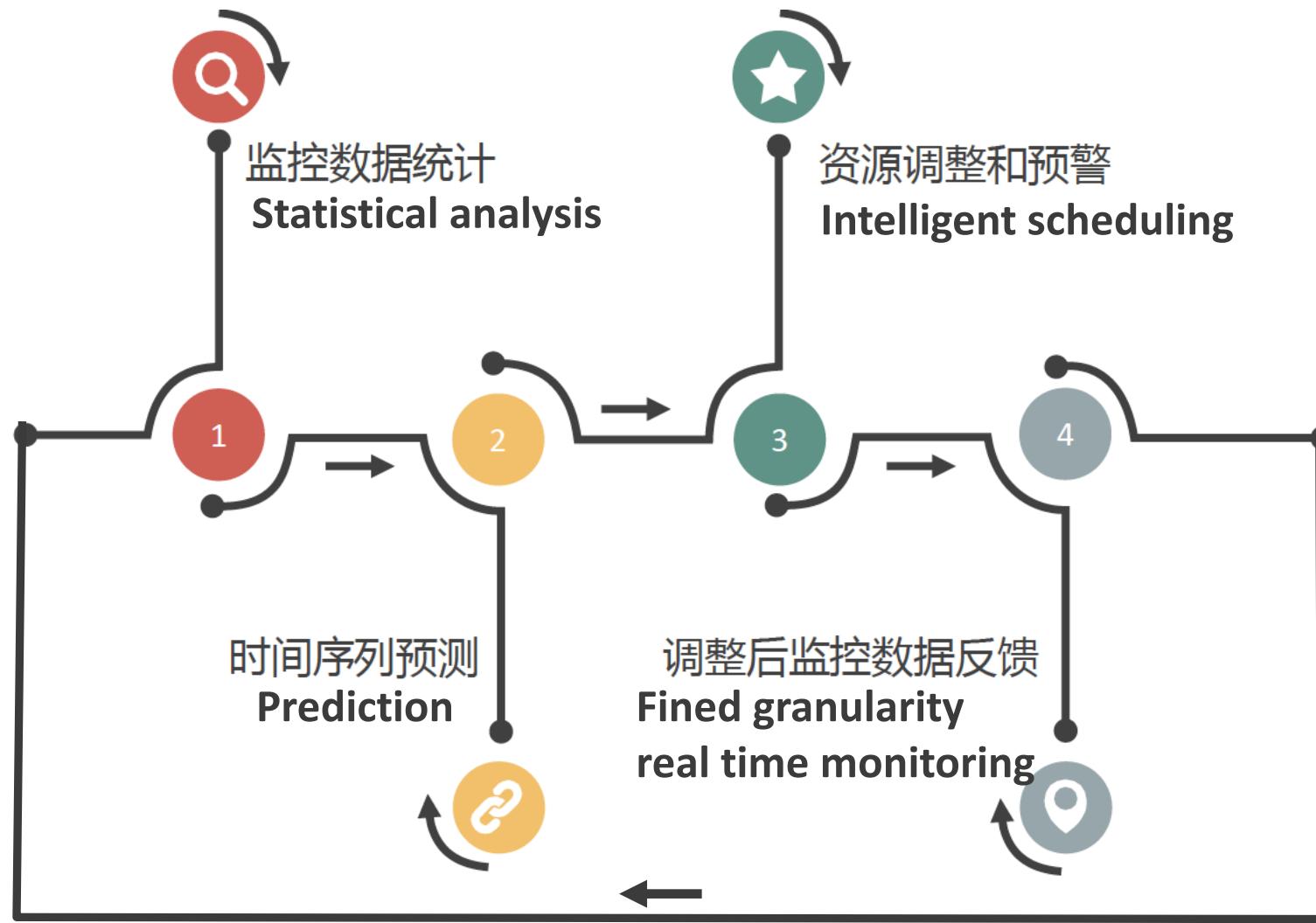
- Scheduling/rescheduling long running critical services
 - Advanced and accurate scheduling
 - A large number of candidate nodes
- Batch jobs
 - Simple and group scheduling
 - A smaller number of candidate nodes
- Cache scheduling decisions

Archimedes: Intelligent Resource Management

Utilize advanced model and machine learning algorithms to **automatically** and **accurately** allocate, schedule and adjust Pod resources to maximize the resource utilization while satisfying the SLOs.



A Closed Loop System



Advanced Scheduling

CPU and Memory sizing

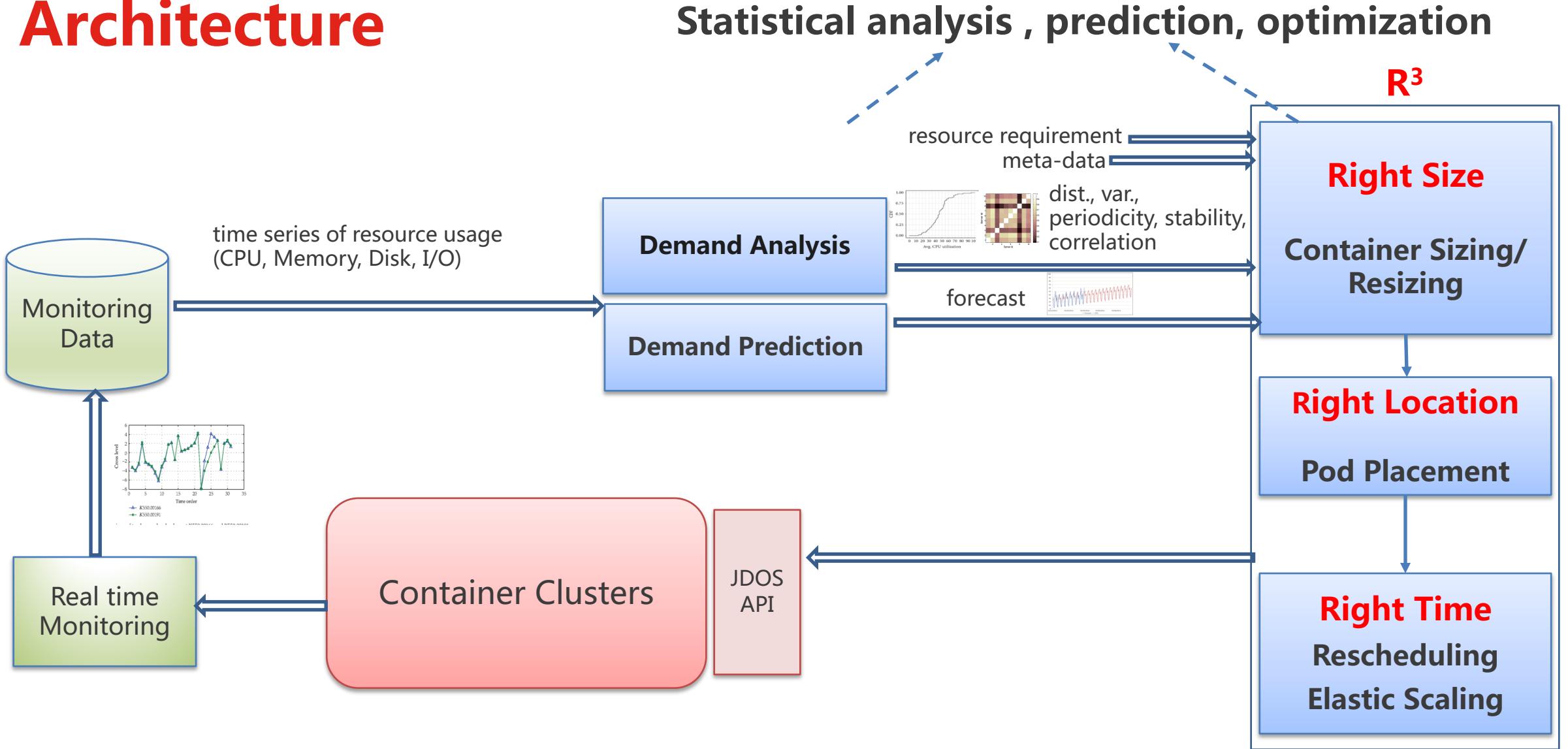
- Metadata data, historical resource usage statistics and future demand prediction
- Proactive adjustment

Host selection :resource balance, resource availability, pod affinity

Mixed placement of workloads

Resource scaling across multiple hybrid resource pools

Architecture



Right Sizing: Oversubscription with Performance Guarantees

- Estimate the workload demand.
 - New workload
 - Resource requirements
 - Meta-data based classification + group resource requirements
 - Existing workload
 - Historical usage analysis + prediction (ARIMA, LSTM, ...)
- Allocate or adjust the workload resource based on the estimates.
- Reschedule the workload if needed.

Right Sizing of CPU Resource

Adaptive tail-bound based resource request and limit

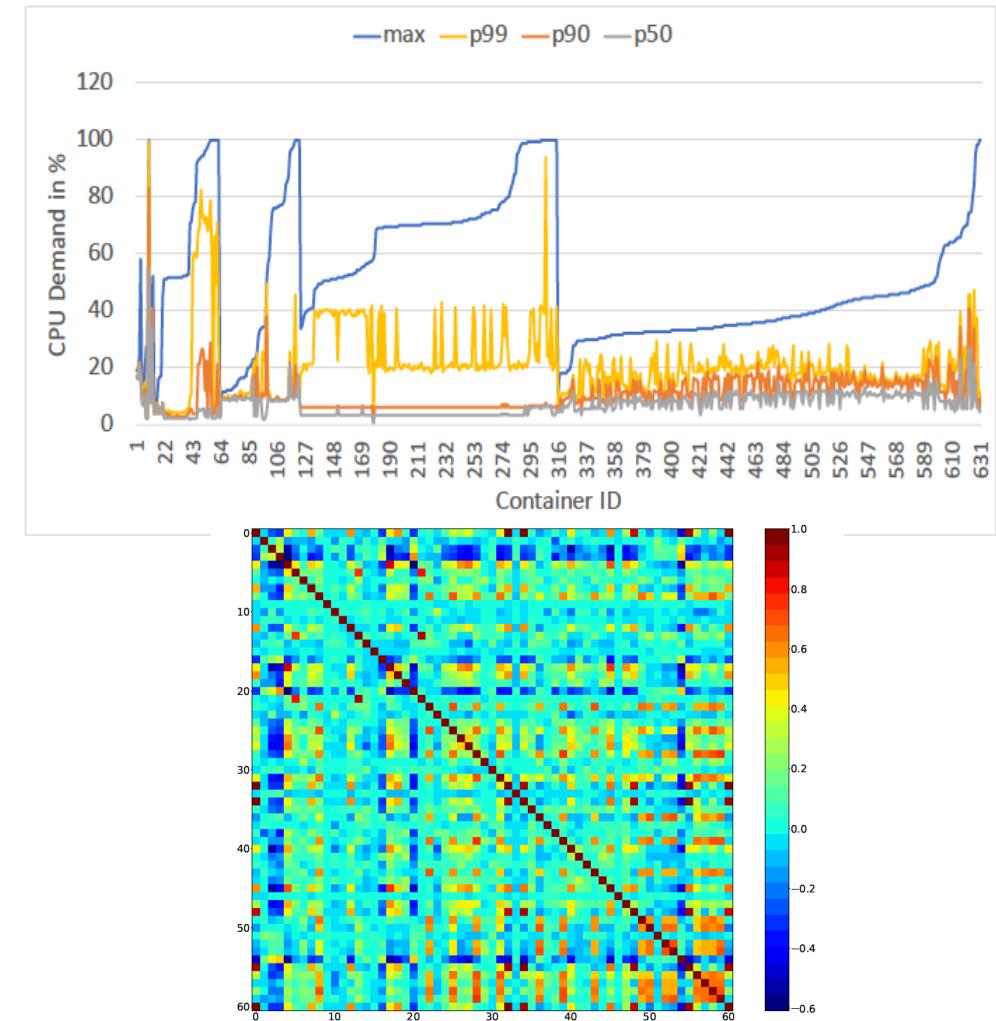
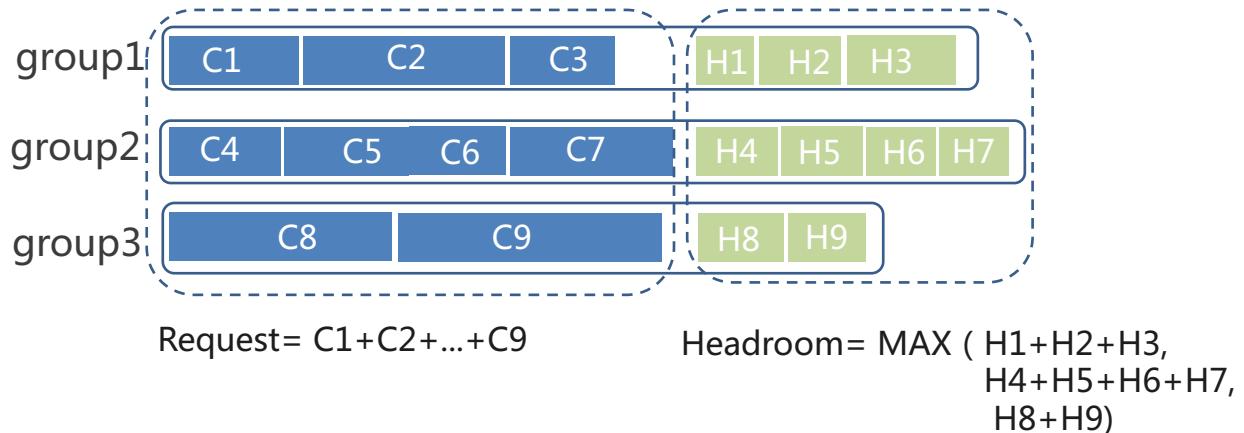
- Guaranteed resource (request): $request_i = 90\text{thile}$
- Maximum resource(limit): $headroom_i = maximum(99\text{thile})$
- Opportunistic headroom: $headroom_i = limit_i - request_i$

Shared headroom on host ($\max_of_sum < \sum_of_max$)

- Correlated containers: $headroom = headroom_i + headroom_j$
- Uncorrelated containers: $headroom = \max(headroom_i, headroom_j)$
- Shared headroom: $shared_headroom \leq \sum_{i=1}^N headroom_i$

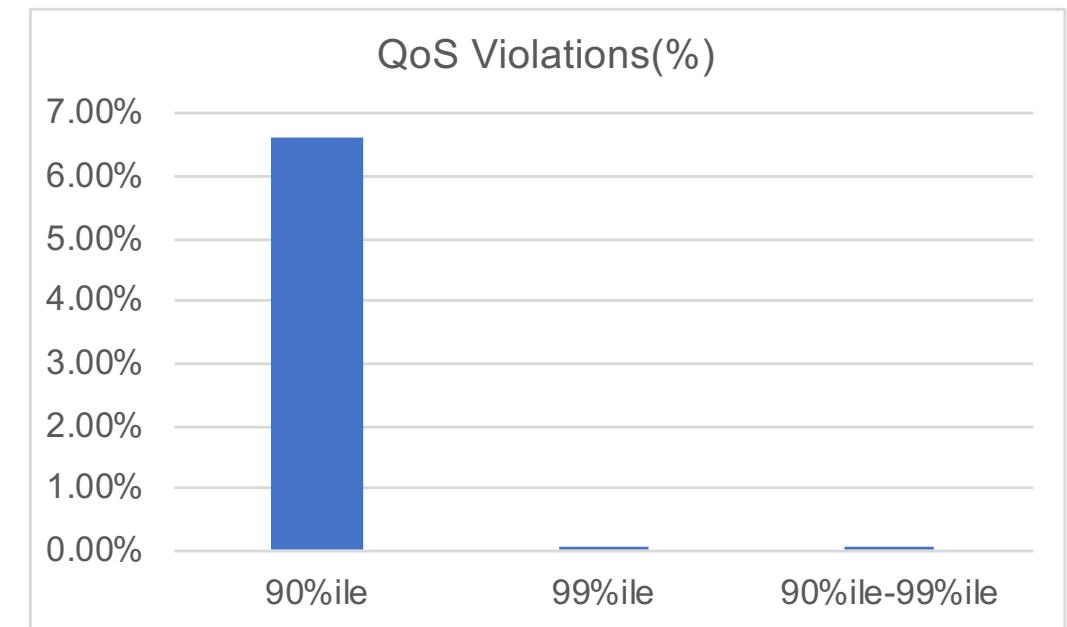
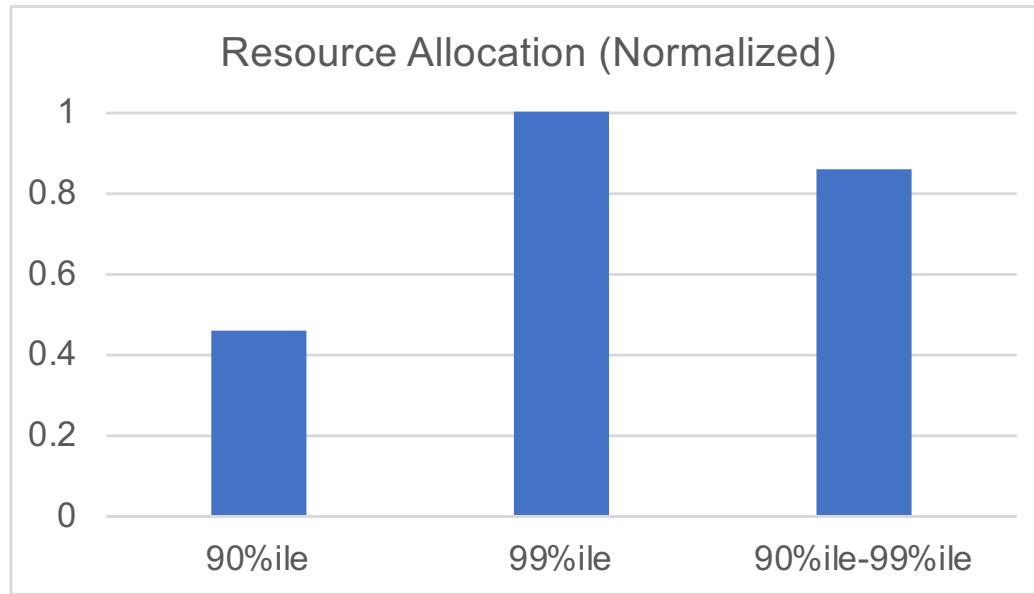
Scheduling constraint

- $\sum_{i=1}^M request_i + shared_headroom \leq capacity$



2-3X improvement without performance degradation!

Resource Usage and Performance Tradeoff



Right Sizing of Memory Resource

Maximize the memory utilization without OOMs

Estimate the next peak demand

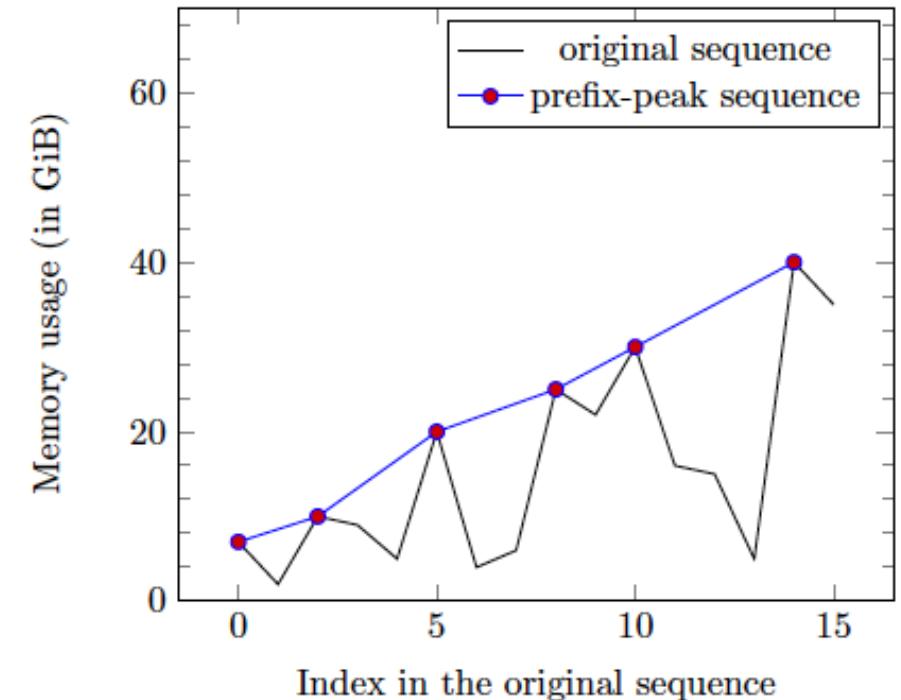
New workload

- Directly use memory limit specified in the submission.
- Classify the workload based on its meta-data and use the statistic measures of the group.

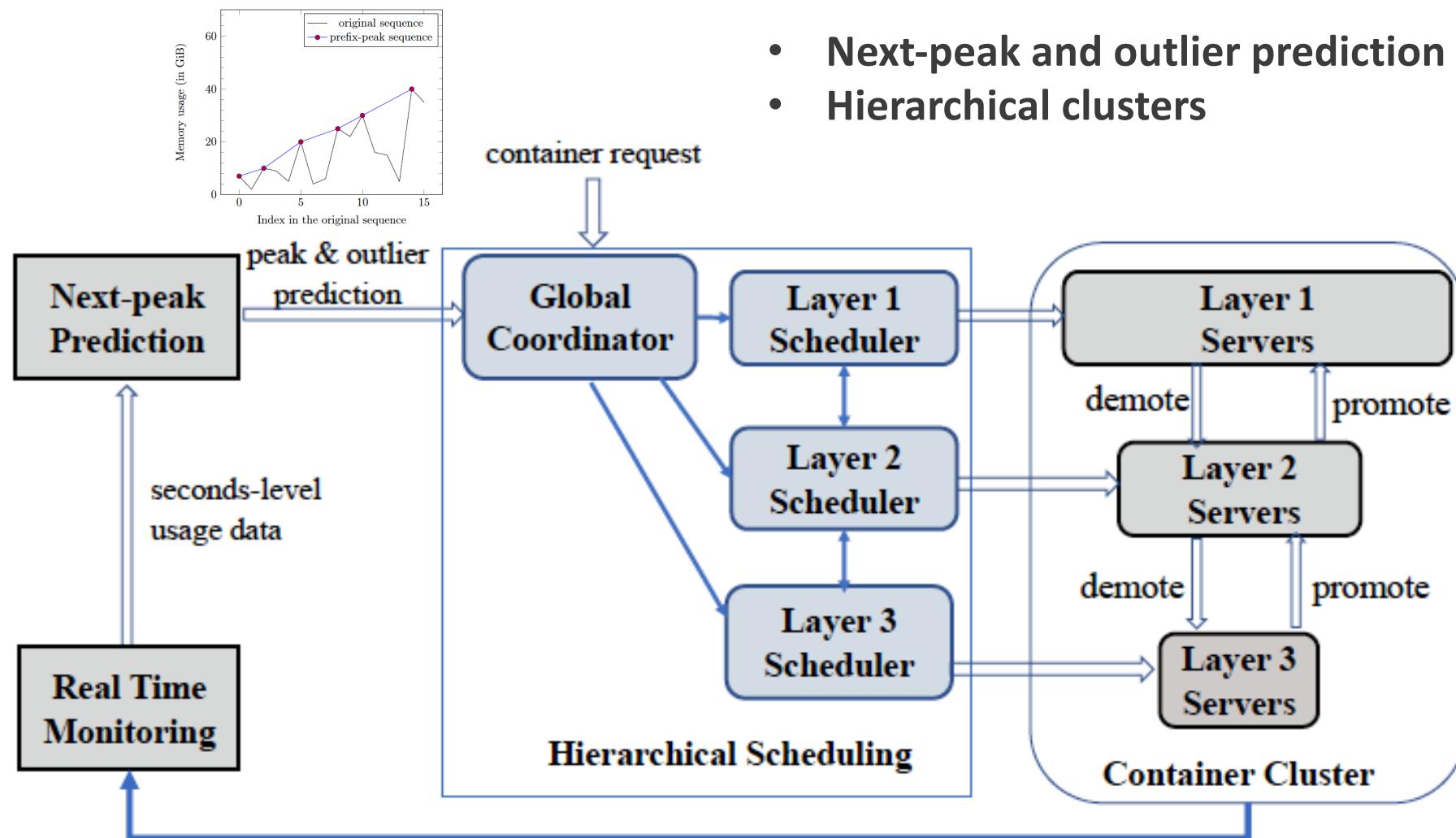
Existing workload

- Predict the next peak demand

Adjust or reschedule the workload if needed



Memory Scheduling for Online Services



Memory Utilization of Online Services

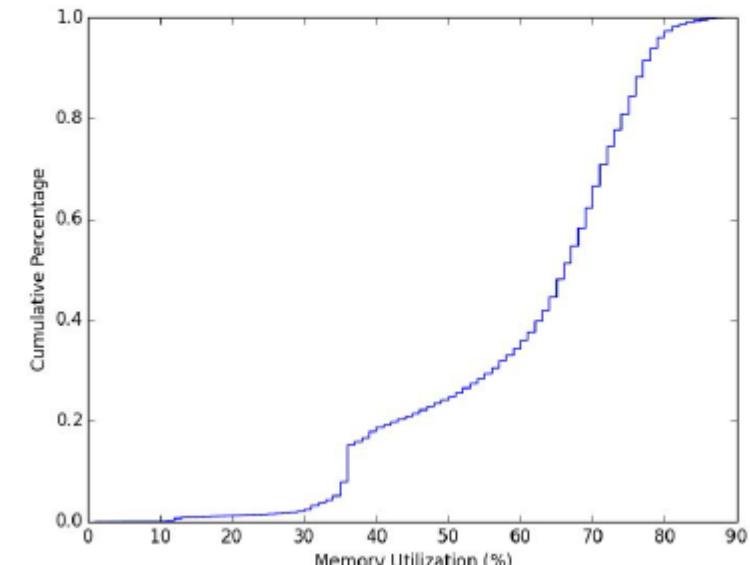
80% utilization without OOM

<i>memory utilization</i>	#OOM	$\frac{\#sched_fail}{\#workload}$	$\frac{\#migration}{\#promotion}$
80.45%	0	3.19%	2.26%
85.02%	3	5.32%	4.24%
89.48%	226	5.77%	4.76%

Real-trace Driven Simulation

61% utilization without OOM

#Servers	Memory Utilization				
	Avg.	Min	Median	99%ile	Max
8,126	61.0%	1.0%	66.0%	83.0%	90.0%



Production Systems

Host Selection Optimization

Host Priority	Vanilla Kubernetes	JDOS
Multi-resource balance	Variance between remaining CPU, memory and disk	Similarity between the host remaining resources and the request
Resource availability	Equal weights for CPU and Memory	Workload-aware variable weights Consider resource usage dynamics
Pod affinity	N/A	Pod demand correlation-aware

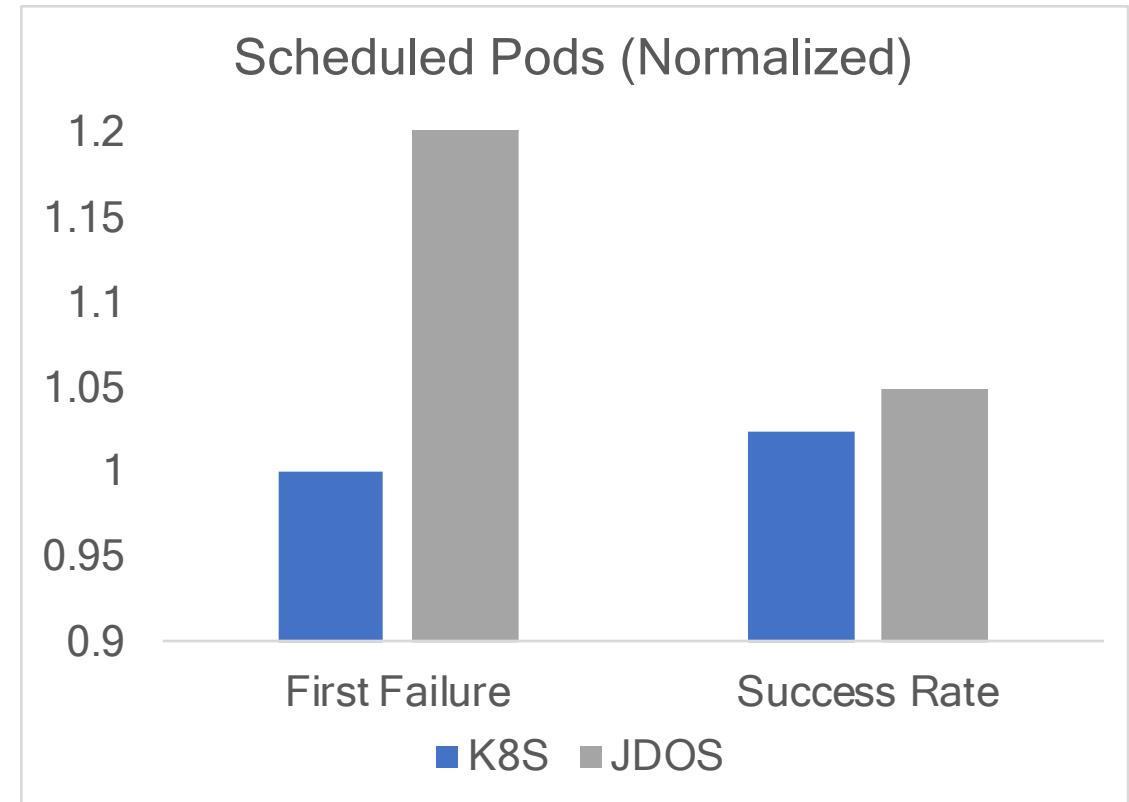
Comparison with K8S

Setup

- 4,857 servers (8 configurations)
- 25,000 pods (120 configurations)

Results

- Schedule 5-20% more pods than K8S.
- Comparable scheduling performance with K8S.



Mixed Workload Placement



Hybrid Resource Pools

Resource Scaling across hybrid resource pools



Co-location of online services and batch jobs is not a silver bullet.

- **Performance and fault isolation concerns.**
- **Memory safety margin:** batch jobs are not memory intensive.
- **Cost issue:** batch jobs can run in "cheap" and "remote" data centers, but internet services need to be located close to users with low latency.

June 18 Anniversary Sale (2019)

- \$29.2 billion in transaction volume
- Scheduled half million of containers (Pod) and millions of CPU cores
- Used 20% less servers

Conclusions

- JD is one of the earliest adopters of Kubernetes and container technology.
- We have customized Kubernetes and built a modern system on top of it.
- This entire ecosystem plus our own optimizations make our system scalable and efficient.
- Through advanced scheduling, we significantly improve resource utilization and reduce costs.

Acknowledgements

Thanks a lot for the following colleagues!

Jiangang Fan, Min Li, Yongcheng Bao, Xiaoping Liang, Huaxia Wang,
Zhengxuan Wu, Shugang Chen, Wanli Yang, Xinkun Xu, Feng Wang,
Huasong Shan, Liying Zhang, Shanshan Sun

Thanks for your time!

Contact:

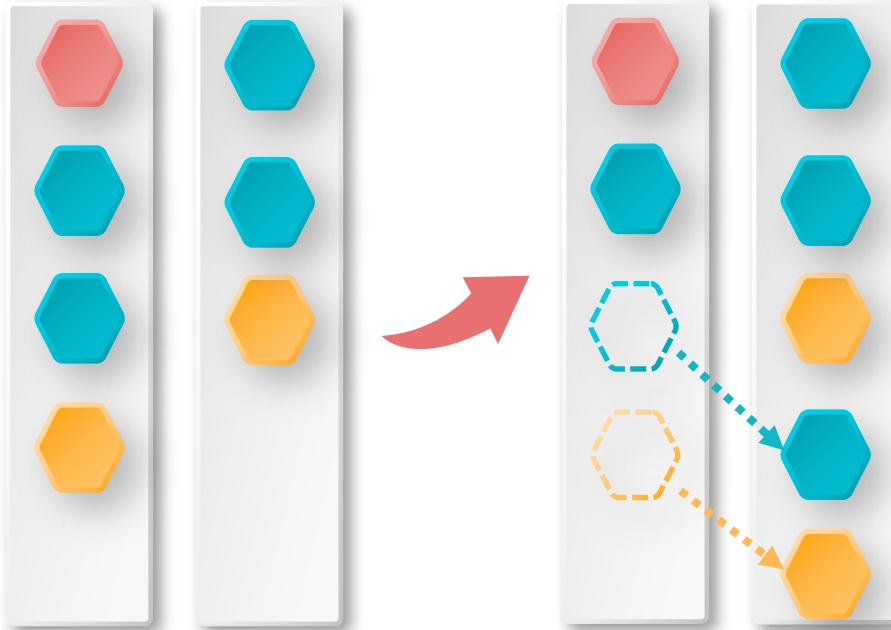
Yuan Chen (Email: yuan.chen@jd.com, WeChat: yuan_gt)

Haifeng Liu (Email: bjliuhafeng@jd.com)



Preemption and Admission Control

Preemption Optimization



- Resource demand prediction
- Priority and SLOs-based preemption
- Proactive preemption
- Preempted apps in scheduling queues
- Overloaded host locked

Admission Control

- Reject low priority request.
- No requests after a certain number of retries.

Right Sizing: Oversubscription with Performance Guarantees

K8S

- HPA alone does not solve the resource efficiency problem.
 - Right size is needed to save resources.
 - Not work for stateful pod.
- VPA implementation is too simple.
 - Request only
 - Recreate pod upon update
 - Thrashing
 - OOM

Archimedes

- Advanced models and algorithms for setting correct size for Pods
- Meta data based initial size
 - Statistical features of historical data + future demand prediction based (re)sizing
 - CPU: request and limit + headroom
 - Memory: peak and outlier prediction + hierarchical sub-clusters
 - Disk resource management for databases
 - **High utilization, low SLOs violations, no OOM, stable**

Right Size, Right Location and Right Time

R³: Right Size, Right Location, Right Time

- Pod sizing (**Right Size**)
- Host Selection (**Right Location**)
- Pod migration and rescheduling (**Right Time & Right Location**)
- Co-scheduling of mixed workloads across hybrid resource pool
(Right Time, Right Location, Right Size)

Host Selection: Multi-Resource Balancing

K8S: BalancedResourceAllocation = $10 - \text{abs}(\text{cpuRemainingFraction} - \text{memoryRemainingFraction}) * 10$

Optimization : similarity between the request resource and available resources

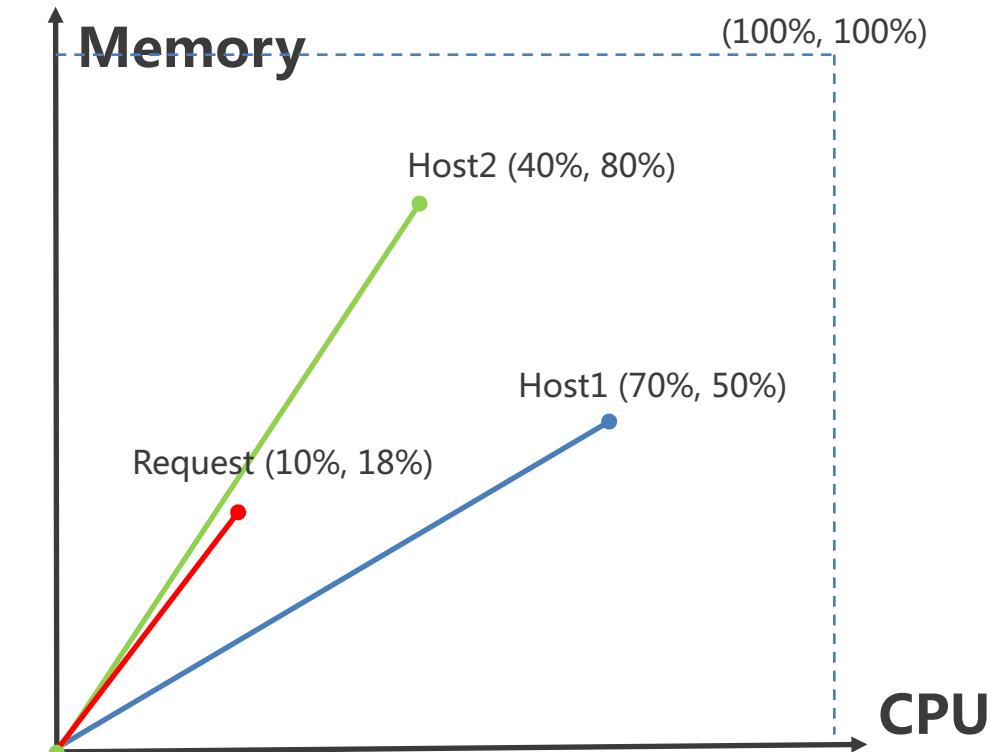
$$\text{similarity} = \cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}},$$

Host available resources : (U_{CPU}, U_{MEM})

Pod resource request : (R_{CPU}, R_{MEM})

Metric:

$$\frac{(R_{CPU}U_{CPU} + R_{MEM}U_{MEM})}{\sqrt{R_{CPU}^2 + R_{MEM}^2} \sqrt{U_{CPU}^2 + U_{MEM}^2}}$$



Host Selection: Resource Availability

K8S :

LeastRequestedPriority = $\text{cpu}((\text{capacity} - \text{sum(requested)}) / \text{capacity}) + \text{memory}((\text{capacity} - \text{sum(requested)}) / \text{capacity}) / 2$

Optimization: dynamic variable weights

$$w_{\text{cpu}} * \text{cpuAvailFraction} + w_{\text{mem}} * \text{memAvailFraction}$$

Algorithm 1:

Weighted sum of CPU : $\text{cpuFraction} = \sum_i f_i * \text{cpuFraction}_i$

Weight sum of memory : $\text{memFraction} = \sum_i f_i * \text{memFraction}_i$

$$w_{\text{cpu}} = \frac{\text{cpuFraction}}{\text{cpuFraction} + \text{memFraction}}$$

$$w_{\text{mem}} = \frac{\text{memFraction}}{\text{cpuFraction} + \text{memFraction}}$$

Algorithm 2: Hosts with stable resource usage are more usable.

MAD- Median Absolute Deviation

$$w_{\text{cpu}} = 1 - \text{median}(|\text{cpuUtil}_i - \text{median}(\text{cpuUtil}_i)|)$$

$$\cdot \quad w_{\text{mem}} = 1 - \text{median}(|\text{memUtil}_i - \text{median}(\text{memUtil}_i)|)$$

Host Selection: Correlation-awareness

X : existing pods resource usage

Y : new pod resource usage

Anti-affinity: If Y can be predicted by X, X and Y are highly correlated. We should avoid placing them on the same host.

$$\mathbf{X} = \begin{bmatrix} 1 & x_{1,1} & \cdots & x_{1,k} & \cdots & x_{1,K} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 1 & x_{n,1} & \cdots & x_{n,k} & \cdots & x_{n,K} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 1 & x_{N,1} & \cdots & x_{N,k} & \cdots & x_{N,K} \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \\ \vdots \\ y_N \end{bmatrix}$$

Multi-regression analysis $\mathbf{b} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad \hat{\mathbf{y}} = \mathbf{X}\mathbf{b}$

Correlation

$$\frac{\sum_{i=1}^n (y_i - m_Y)^2 (\hat{y}_i - m_{\hat{Y}})^2}{\sum_{i=1}^n (y_i - m_Y)^2 \sum_{i=1}^n (\hat{y}_i - m_{\hat{Y}})^2}$$