



KubeCon



CloudNativeCon



China 2024

JDCloud cross-cluster large-scale application management practice

JDCloud

2024



KubeCon



CloudNativeCon



China 2024

JDCloud
2024



XiaoFei Wang

JDCloud
CloudNativeEngineer
wangxiaofei67@jd.com

CONTENTS

01

Development Road of JD Cloud Containers

02

Federation Cluster

03

Federation Auto-scaling

04

Summary



KubeCon



CloudNativeCon



China 2024





KubeCon



CloudNativeCon



China 2024



JDCloud
2024

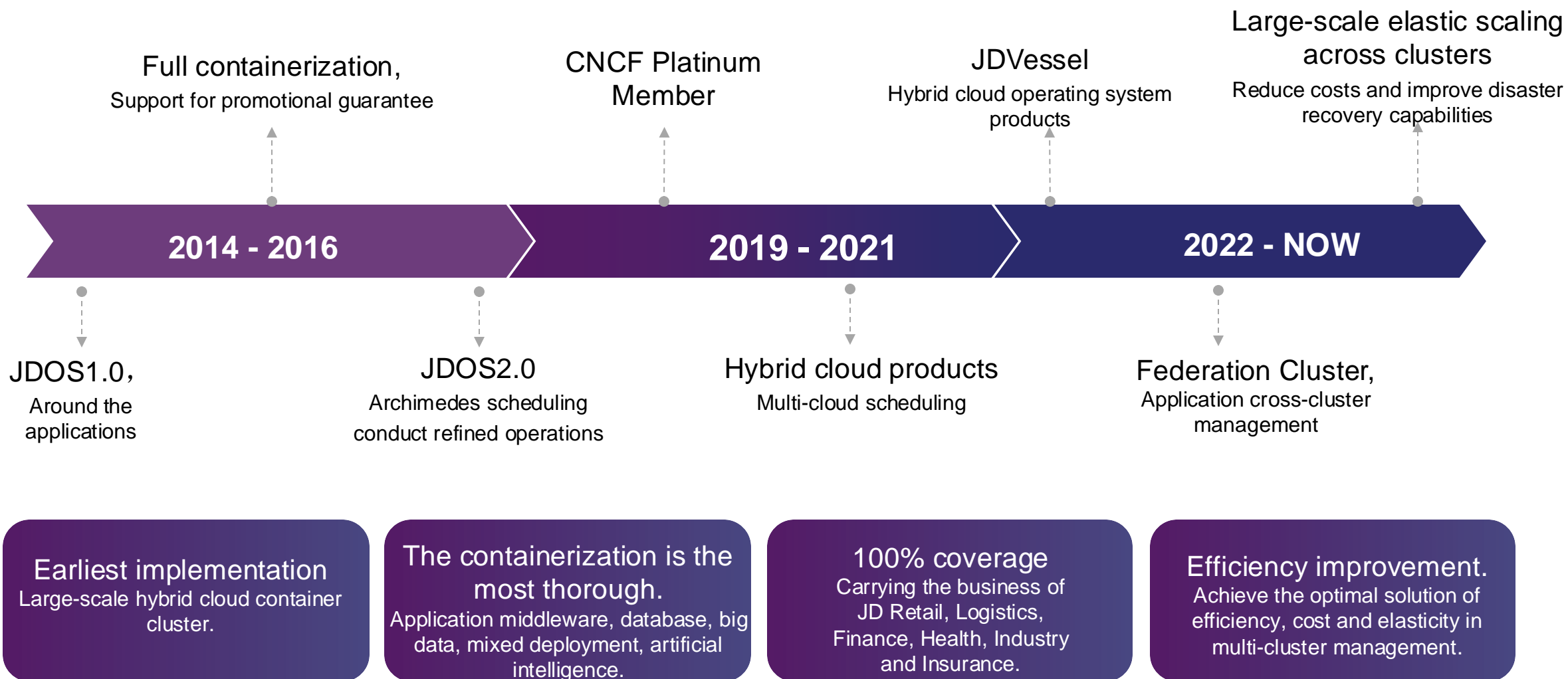
01

Development Road of JD Cloud Containers

JD Containerization Development History



China 2024



Everything around the application.



China 2024

Single-Cluster limitation

- Fault explosion radius is too large, affecting the application of SLA.
- The application of cross-cluster auto-scaling is difficult, manual operation and maintenance.

Multi-cluster Scheduling

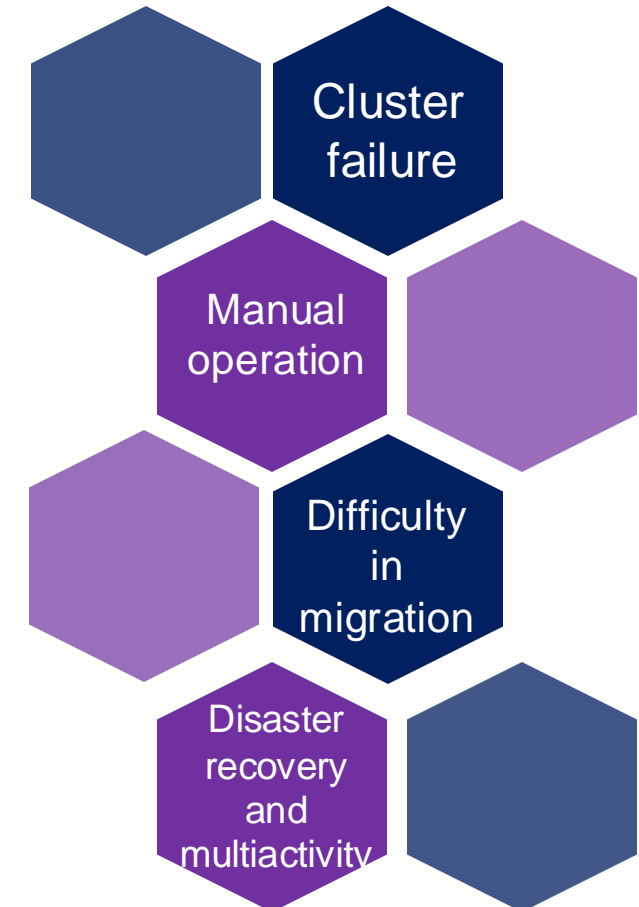
- The Unified resource scheduling
- Multi-group flexible scheduling strategy

Multi-cluster auto-scaling

- Cross-cluster elastic scaling
- One-click migration

Multiple and HA

- Multiple applications
- Cross-cluster network





KubeCon



CloudNativeCon



China 2024



JDCloud
2024

02

Federation Cluster

Product



China 2024



JDCloud-JDVessel

<https://www.jdcloud.com/cn/products/yunjian>



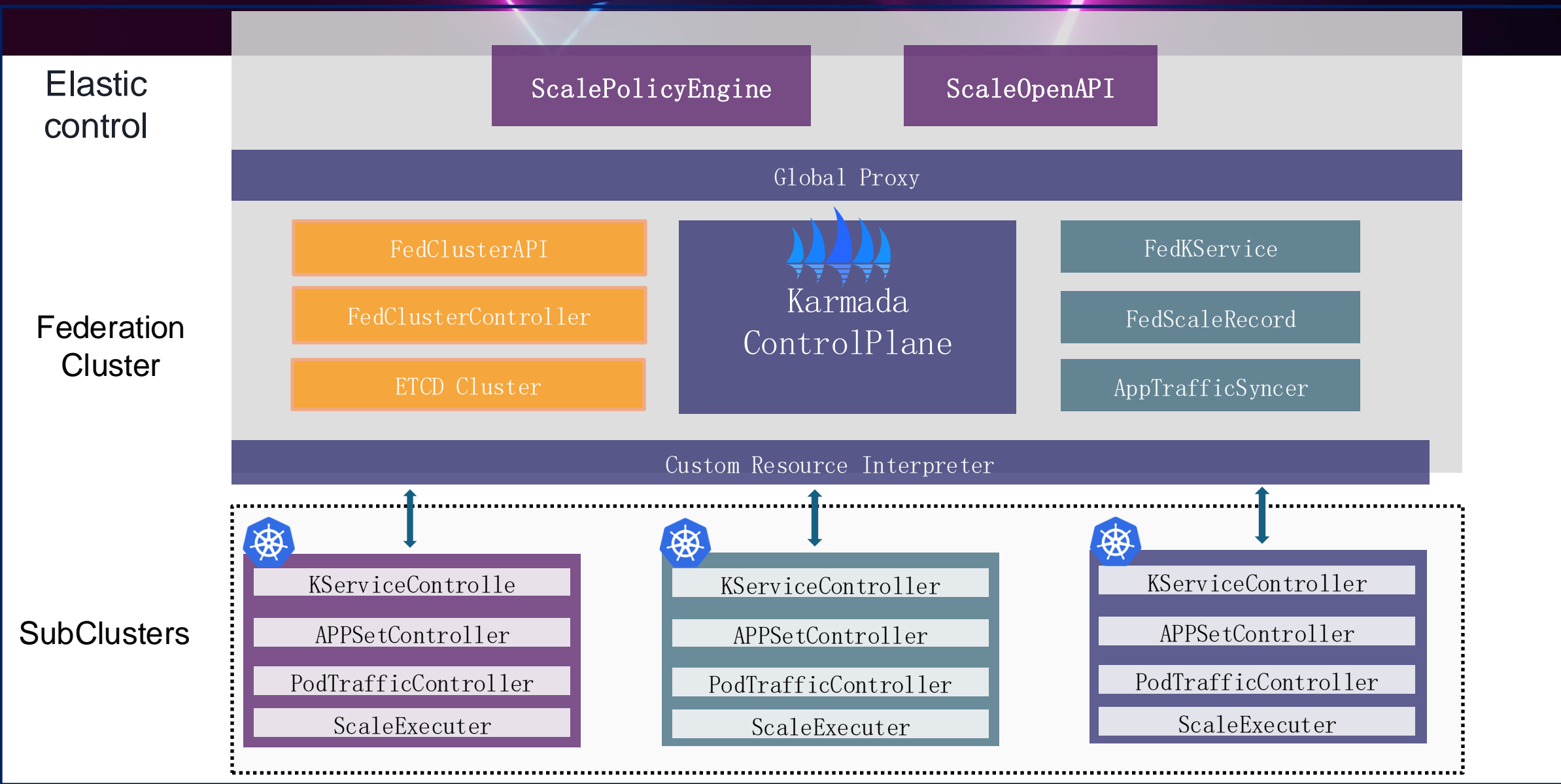
In-depth development and enhancement based on Karmada

<https://github.com/karmada-io/karmada>

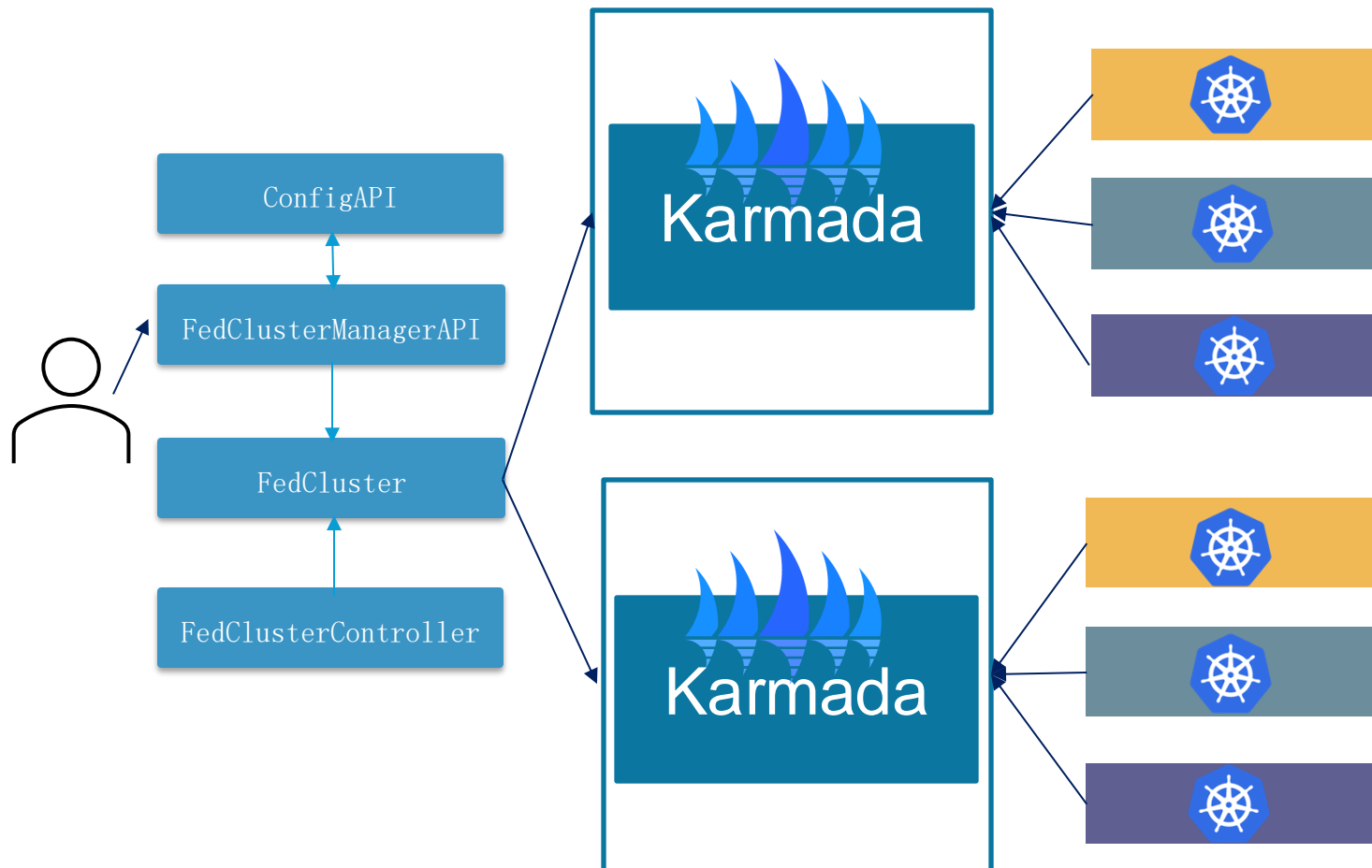
Architecture



China 2024



Federation cluster deployment



Create

The FedClusterManagerAPI handles user requests to create a Federation cluster and interacts with the ConfigAPI service to obtain the creation context and parameters. Create the FedCluster CR.

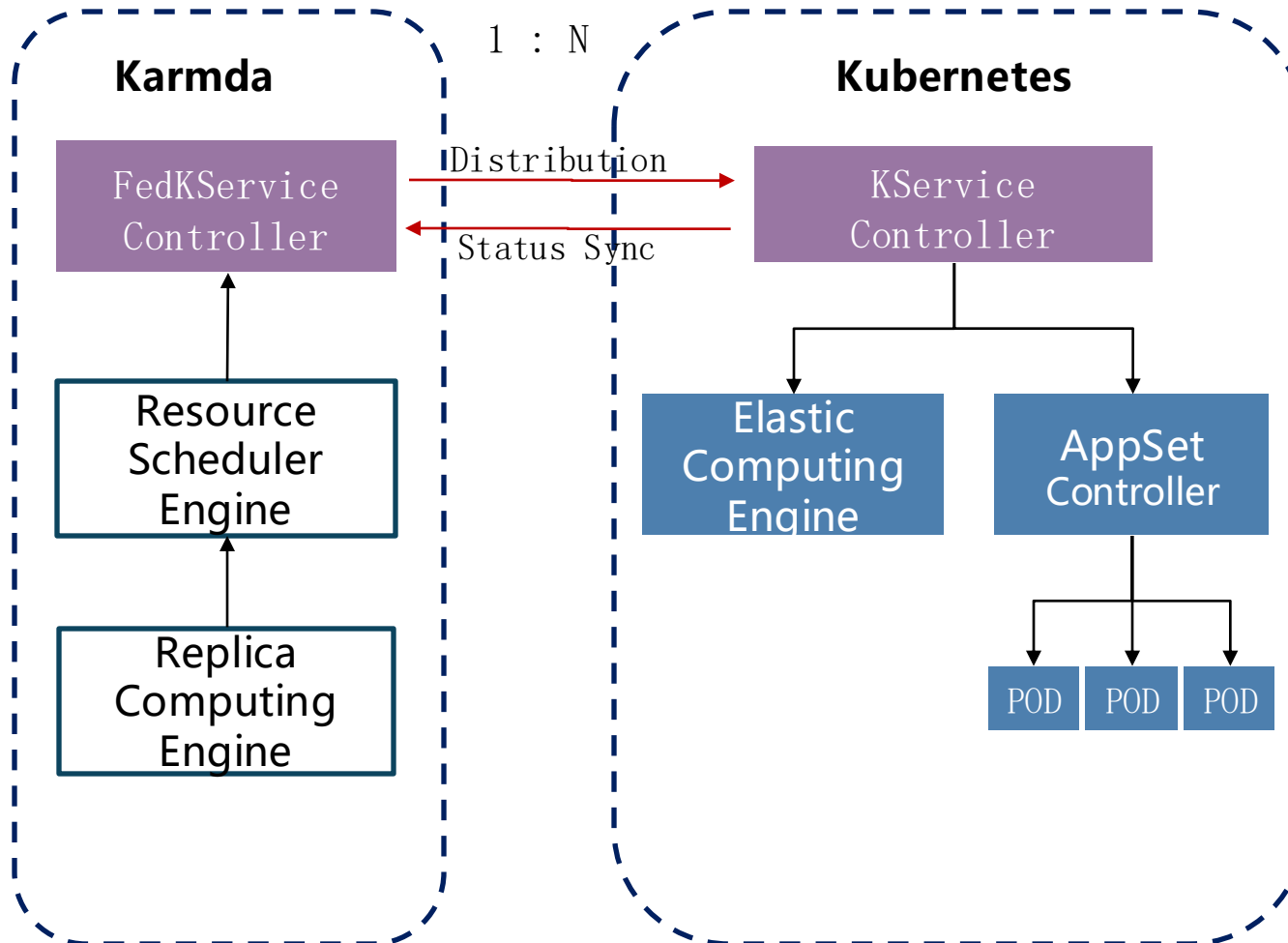
Watch

FedClusterController monitors FedCluster, according to the Spec statement, in the specified Kubernetes cluster creates a federated cluster ControlPlane.

Manage

FedClusterController monitors FedCluster and declares Spec to join the specified federation cluster. The federation cluster begins to manage multi-cloud Kubernetes clusters.

FedKService Controller



Resource Scheduler Engine

Based on the Hippo service, conduct precise calculation of multi-cluster resources to obtain resource portraits.

Replica Computing Engine

According to the resource portrait and scheduling strategy, perform replica multi-cluster splitting.

Elastic Computing Engine

Execute elastic strategy, execute scaling plan.

AppSet Controller

Directly control the workload of pods.

Federation Enhance



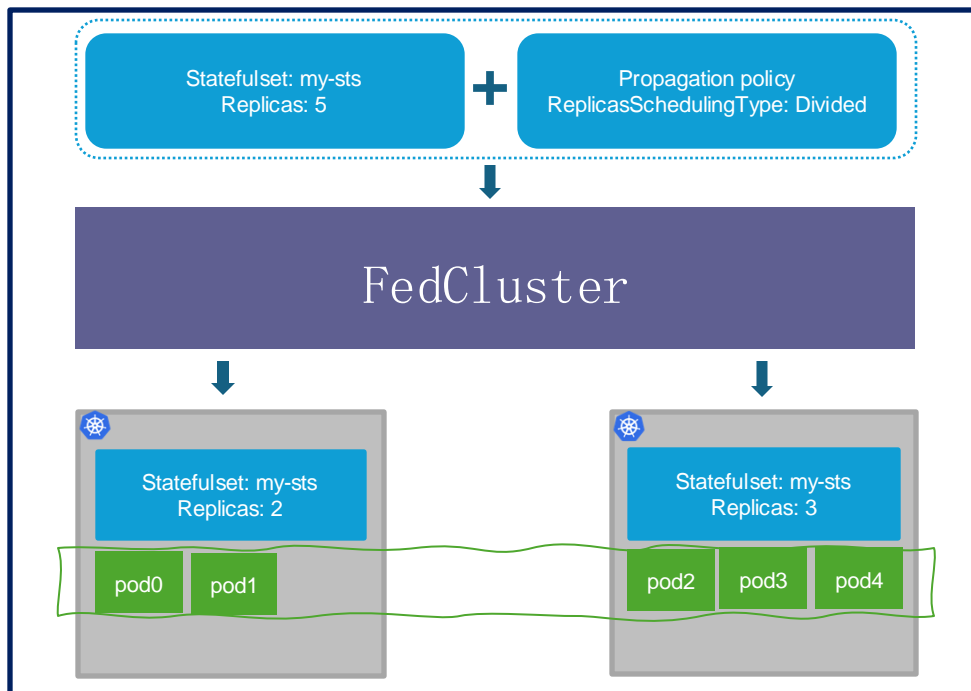
China 2024

KService status sync

Solve the problem of KService status resource synchronization and merging from the sub cluster to the federation cluster.

StatefulSet starting number in multi-cluster control

The behavior of using StatefulSets in the federation cluster and single Kubernetes clusters with sequence number control is the same.

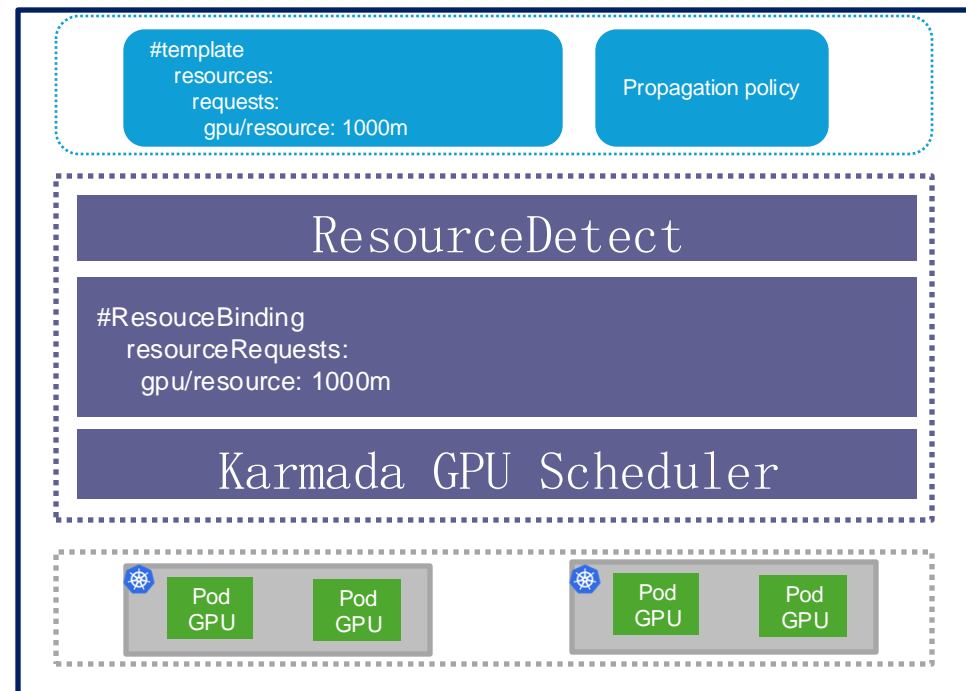


Global proxy performance improvement

1. Use the PB protocol for sub-cluster queries.
2. Modify from separate queries to concurrent queries.
3. Optimize stored procedures to achieve higher data access performance.
4. Performance improvement of 30%.

Karmada scheduler enhance

Support GPU cluster scheduling. Support heterogeneous cluster scheduling.





KubeCon



CloudNativeCon



China 2024

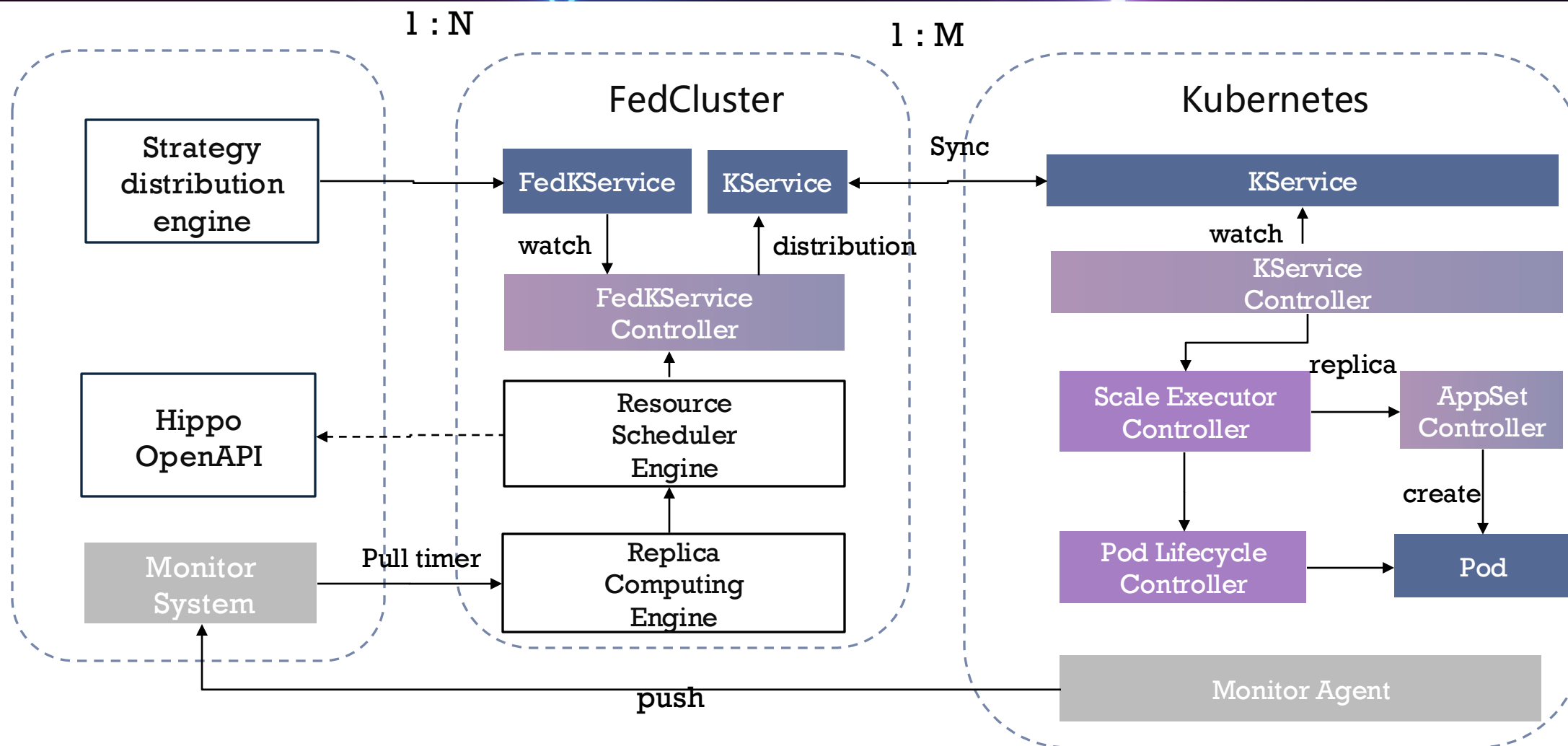


JDCloud
2024

03

Across-cluster Auto-scaling

Auto-scaling Architecture



Auto-scaling Process



China 2024

Strategy
Calculation

Federation

Real-time matching auto-scaling strategy, calculating the number of replicas required

Resource
distribution

Federation

Resource estimation and distributing replicas to distinct clusters

Resource
Sync &
Desync

Federation

Synchronize and desynchronize resource by Karmada

Scaling
execution

Cluster

Repeatedly scaling up and down according to the auto-scaling strategy until the num of mounted Pods reaches the expected value

Pod
Manage

Cluster

Configurable Pod lifecycle management
Abnormal Pod GC mechanism

扩容



缩容



Auto-scaling Scenario



China 2024

Automatic failover

Cluster node failover
Cluster offline

Disaster recovery migration across data centers

Indicator monitoring scaling

Support machine performance and method call
performance indicators

Personalized expansion and contraction rules



Big promotion of hot standby of pressure test container

When traffic is low, the hot standby container is compressed to a low specification

When traffic is high, low-specification hot standby containers are restored

Scheduled task scaling

Support Cron expressions

资源: CPU使用率 聚合方式: Avg

组内容器平均值连续: 次 大于等于: %时进行扩容

组内容器平均值连续: 次 小于等于: %时进行缩容

扩缩方式: 数量

在 将实例数 扩容 到 个

且在 将实例数 缩容 到 个

Practice: Switching between auto-scaling and manual deployment



China 2024



Auto-scaling scenarios

- Enabled by default, auto-scaling
- Only adjust the number of replicas
- Not responsible for updates of image versions, environment variables, configuration files, etc.



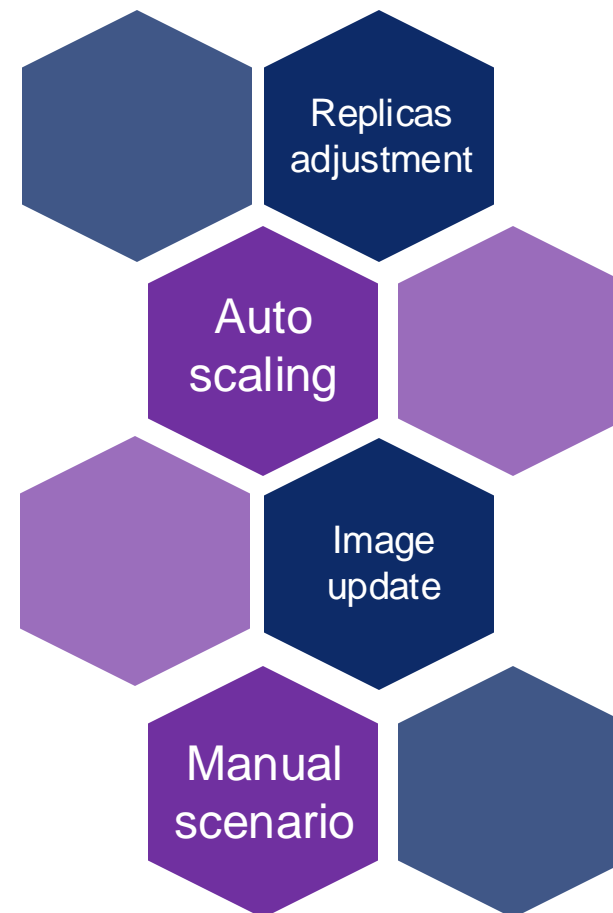
Manual scenario

- First online deployment
- Updates to images, configurations, environment variables, etc.
- Specify a single Pod update, online regression verification



Scene contradiction

- When deploying manually, you want the number of replicas to be stable and stop auto-scaling
- During auto-scaling, manual deployment will interfere with the execution effect of scaling



Three modes



China 2024

Manual

Manual deployment scenario

Allow manual deployment

Shield auto-scaling control

Manual

Auto

Auto

Auto-scaling scenarios

Allow auto-scaling

Shield manual deployment control

None

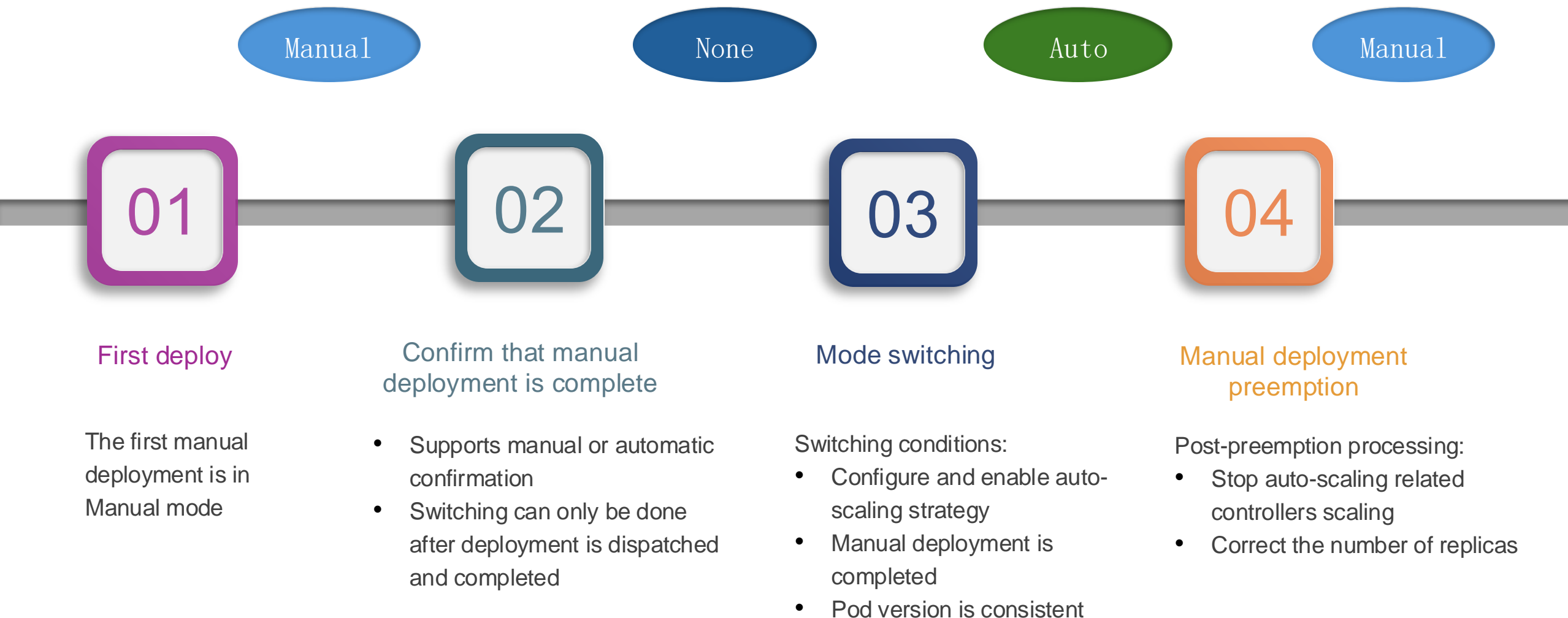
None

Modeless state, isolated scenarios

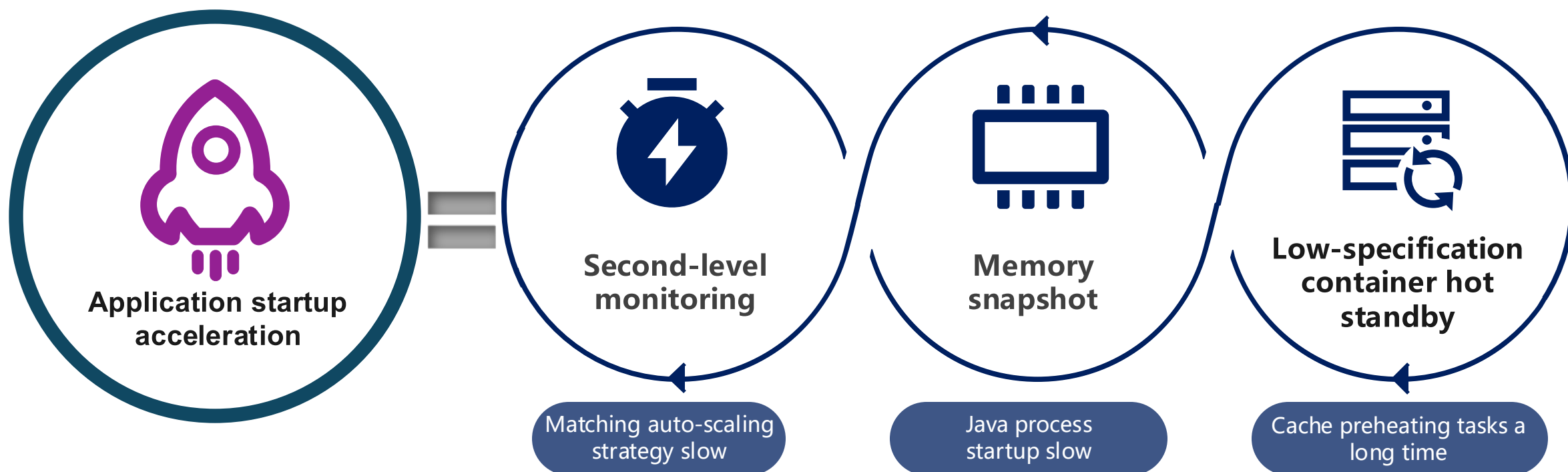
Shield auto-scaling control

Shield manual deployment control

The whole process



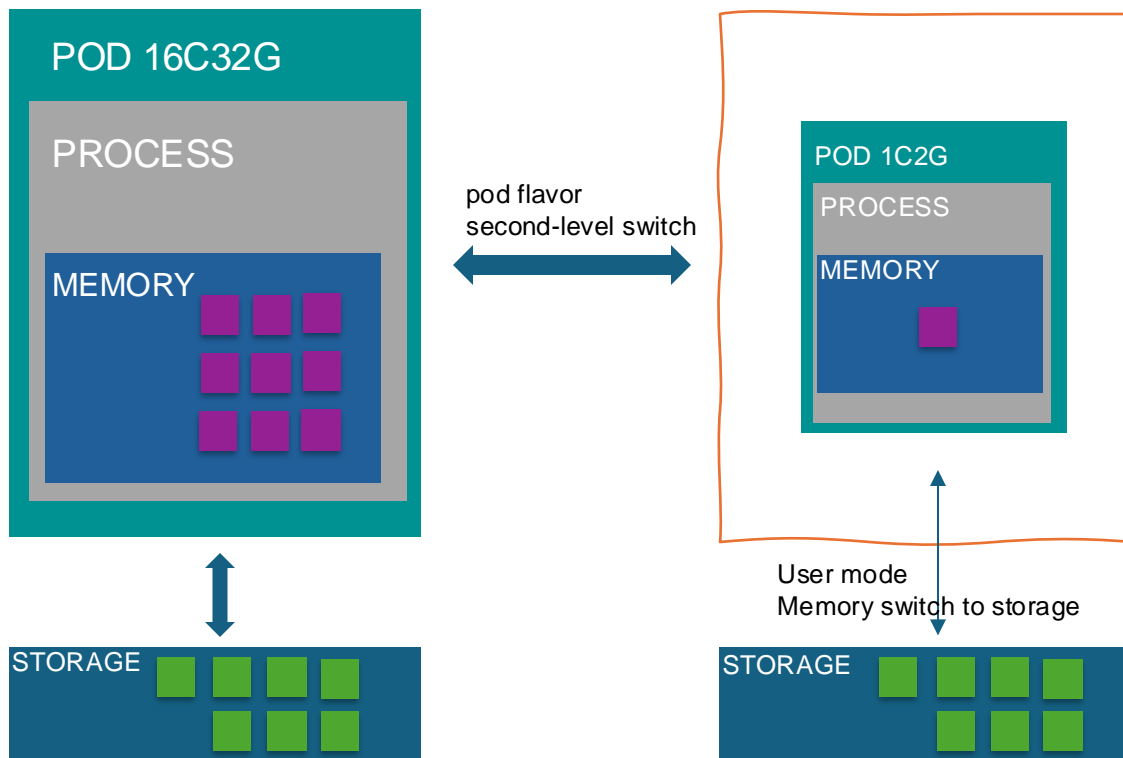
Practice: Application startup acceleration



Low-specification container hot standby

Solve scene problems

Before the promotion stress test, hot standby containers need to be cache-prepared in advance, which takes up a lot of resources.



Advantages

- Lossless, no need to modify applications
- Quick Pod startup, startup time reduced by more than 80%
- Java applications can handle traffic without preheating

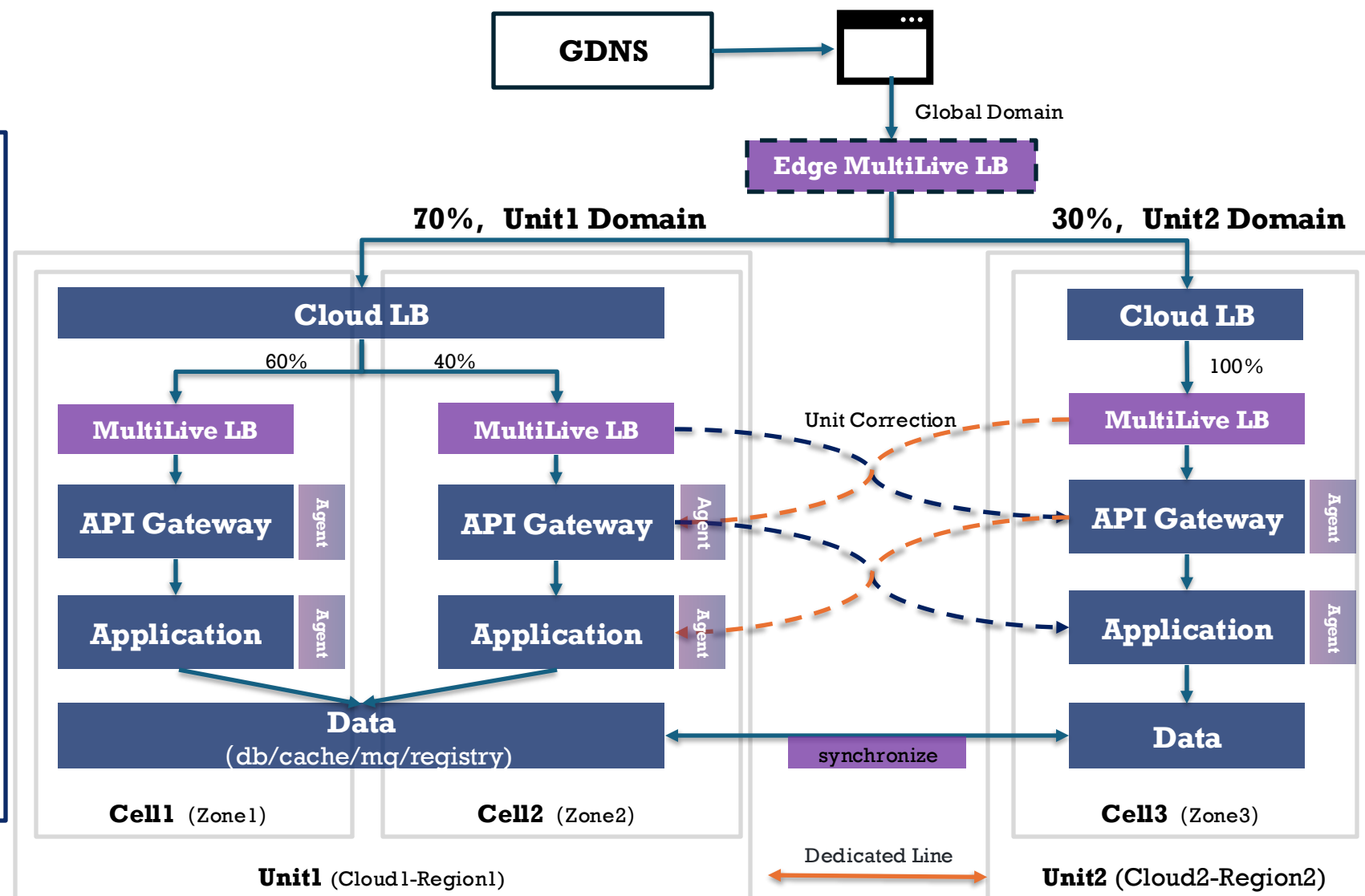
Disadvantages

Need to deploy low-specification instances in advance, occupying a small amount of resources

Scenarios

The workload is mainly triggered by traffic, and the CPU utilization is low after traffic is removed.

- **Service governance framework**
A microservice traffic governance framework for application multi-active and unit-based based on bytecode enhancement
- **High availability improvement**
Integrates Joylive to provide unit traffic management and disaster recovery recovery capabilities
- **Github**
<https://github.com/jd-opensource/joylive-agent>





KubeCon



CloudNativeCon



China 2024



JDCloud
2024

04

Summary

Production use



China 2024

JD Spring Festival Gala Red envelopes Project

Supported JD Spring Festival Gala Red envelopes
Support hundreds of millions of times to grab red envelopes

E-commerce promotion

Supported multiple JD618 and JD11.11 e-commerce promotions, and supported the scaling of trillion-level transaction volume business

The main way to deploy services in the future

Continuous promotion and it will be one of the main ways to services deployment and automatic migration in the future
Application 200,000 cpu+



KubeCon



CloudNativeCon



China 2024

JDCloud
2024

Thank you

Contact us

JDCloud Dev

Tech Group

