



**KubeCon**



**CloudNativeCon**

THE LINUX FOUNDATION



**AI\_dev**  
Open Source GenAI & ML Summit

---

**China 2024**

---



KubeCon



CloudNativeCon



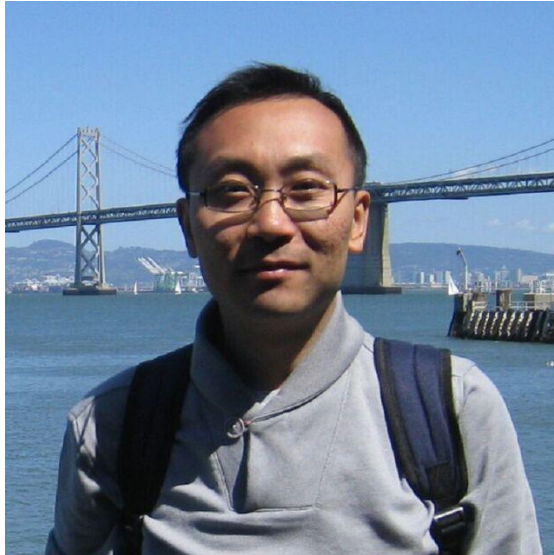
China 2024

# Boundaryless Computing: Optimizing LLM Performance, Cost, and Efficiency in Multi-cloud Architecture

# Who we are?



China 2024



[Kai Zhang \(wsxiaozhang@gmail.com\)](mailto:wsxiaozhang@gmail.com)  
Senior Staff Engineer,  
Alibaba Cloud Intelligence



[Jian Zhu \(jiazhu@redhat.com\)](mailto:jiazhu@redhat.com)  
Senior Software Engineer,  
RedHat



# Agenda



China 2024

- Challenges and solution of running LLM cross clouds/regions
- Accelerates LLM from the data perspective - Fluid
- Manages multiple clusters in the K8s way - OCM
- Demo - Deploy and scale LLM inference service crossing clouds quickly and easily
- Future works



# The Challenges & Solution

# Challenges to infrastructure brought by LLM



China 2024

- The emergence of AIGC/LLM has led to a significant increase in GPU resource consumption, especially during the pre-training phase of foundation models.
- Microsoft has hundreds of thousands of GPUs deployed in more than 60 data centers in Azure cloud for serving ChatGPT

Model	Parameters	GPU counts	Training days
Llama	7B	80 * A100	42
GPT3	175B	1K * A100	30
Llama 3.1	405B	16K * H100	54

- GPU resources in a single data center or cloud region cannot meet LLM workloads resource requirements
- Distributing, synchronizing, and managing model consistency and data security across multiple geographies is a challenge of efficiency and complexity



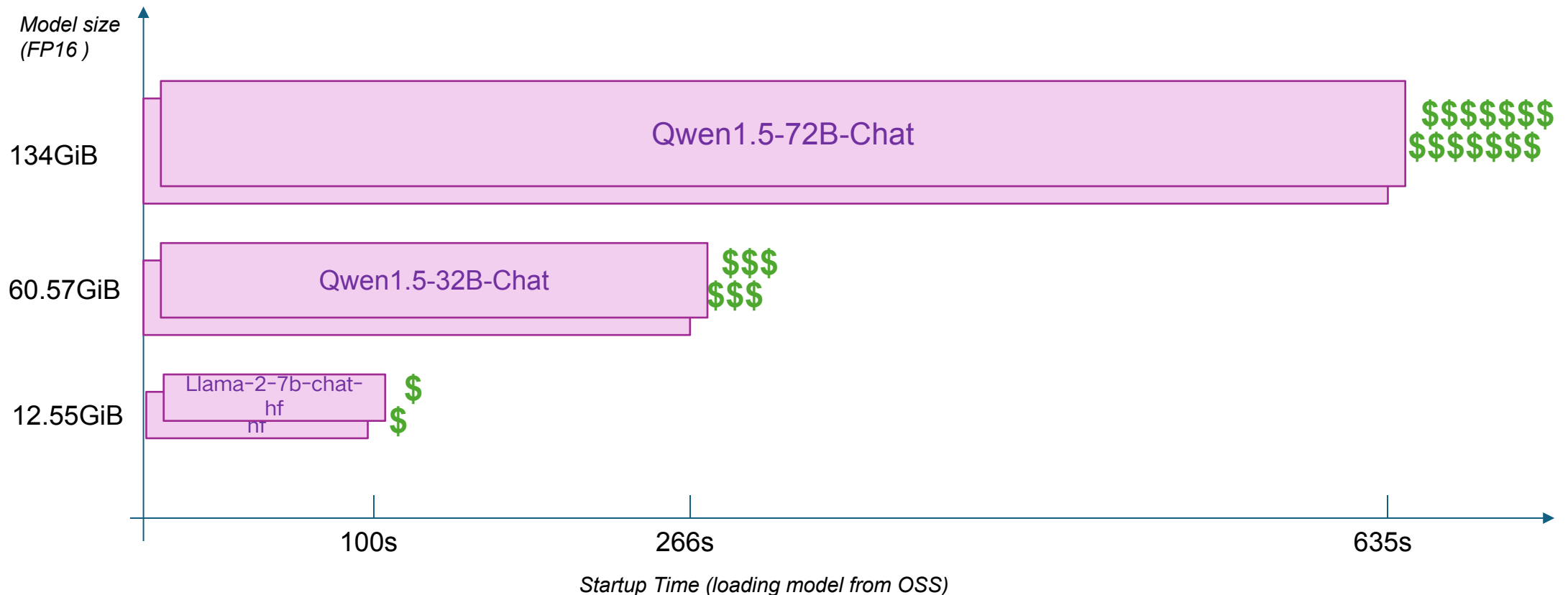
# Challenges to infrastructure brought by LLM



China 2024

- The large model causes the inference service to start very slowly, which seriously affects the elasticity and user experience
- Regional inference services, repeatedly pulling models from remote storage, rapidly driving up

bandwidth costs



# Optimization of LLM efficiency in multi-clouds and multi-regions



China 2024

**Optimize GPU  
resources  
scheduling**

1. Schedule GPU resources cross multiple Kubernetes clusters and dynamically adjust AI tasks distribution

**Optimize  
data/model  
access  
performance**

2. Automatic optimization of large model file loading process, accelerate LLM inference service startup and elastic scaling

**Optimize the  
ease of use of  
multi-  
geographic  
model services**

3. Cross-geographic models and data management and access acceleration, simplify the user experience

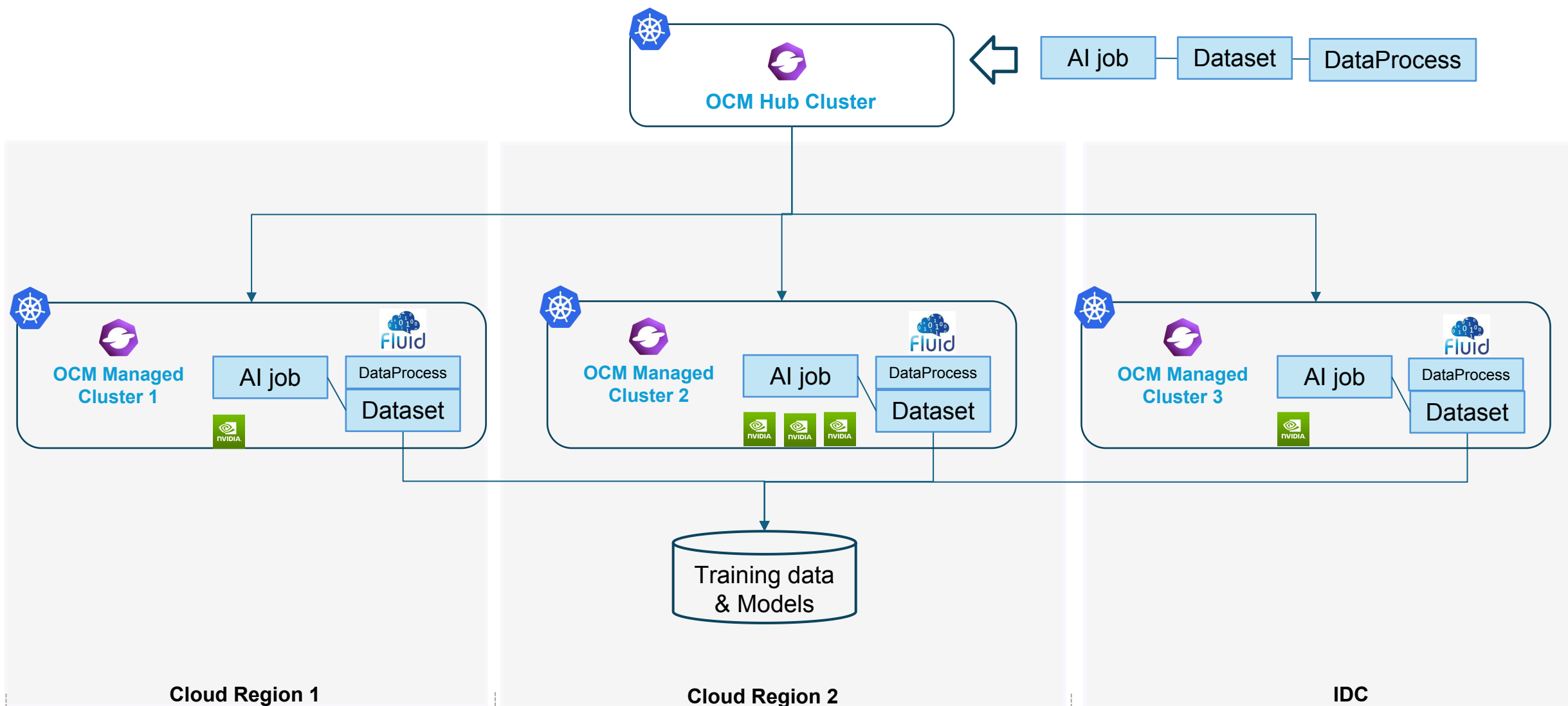





# Solution architecture: OCM + Fluid



China 2024





Fluid – Accelerates LLM from the data perspective

# What is Fluid



China 2024

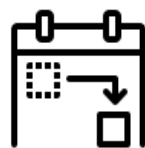
Defines the standard API for Kubernetes to access and manage data.



Dataset Abstraction



Dataset Acceleration

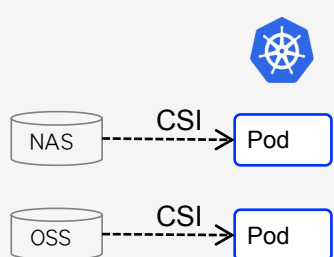


Dataset Process

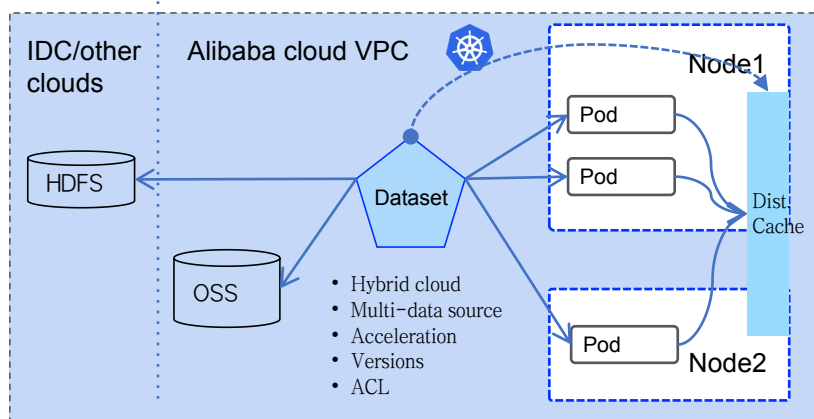


<https://fluid-cloudnative.github.io/>

Storage perspective of K8s



Data usage perspective of Fluid



• Unified arch.

serverful

serverless

Hybrid cloud

• Programmable

Fluid SDK

AI frameworks

• Standard API

Dataset

DataLoad

DataMigrate

• Extensible

CacheRuntime

ThinRuntime

DataProcess

Pipeline

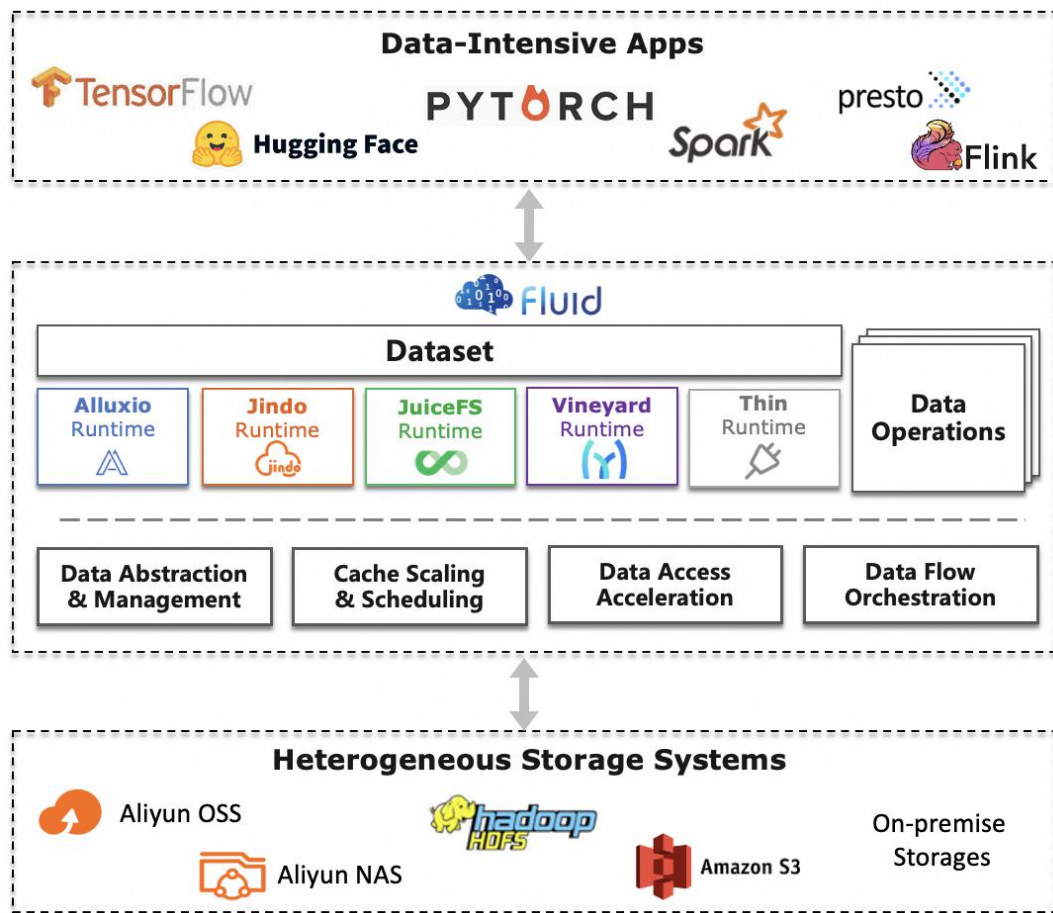
• Jindofs  
• Juicefs  
• Alluxio  
• EFC

• CubeFS  
• ChuboFS  
• Cefs  
• Proprietary storages

# Fluid: Data and Task Orchestrator in K8s



China 2024



- **Standardized:** K8s Native APIs for **data access** and **distributed cache management**.
- **Extensible:** Runtime plugins for different distributed cache and storage backends.
- **Elasticity:** Scale out and in the distributed cache on demand.
- **Performance:** Accelerate data access via elastic distributed cache
- **Automation:** Operation for Data like. *prefetching processing, migration and cache scaling*
- **Orchestration:** Data and task co-aware scheduling

Joint launched by Nanjing University, Alibaba Cloud and Alluxio

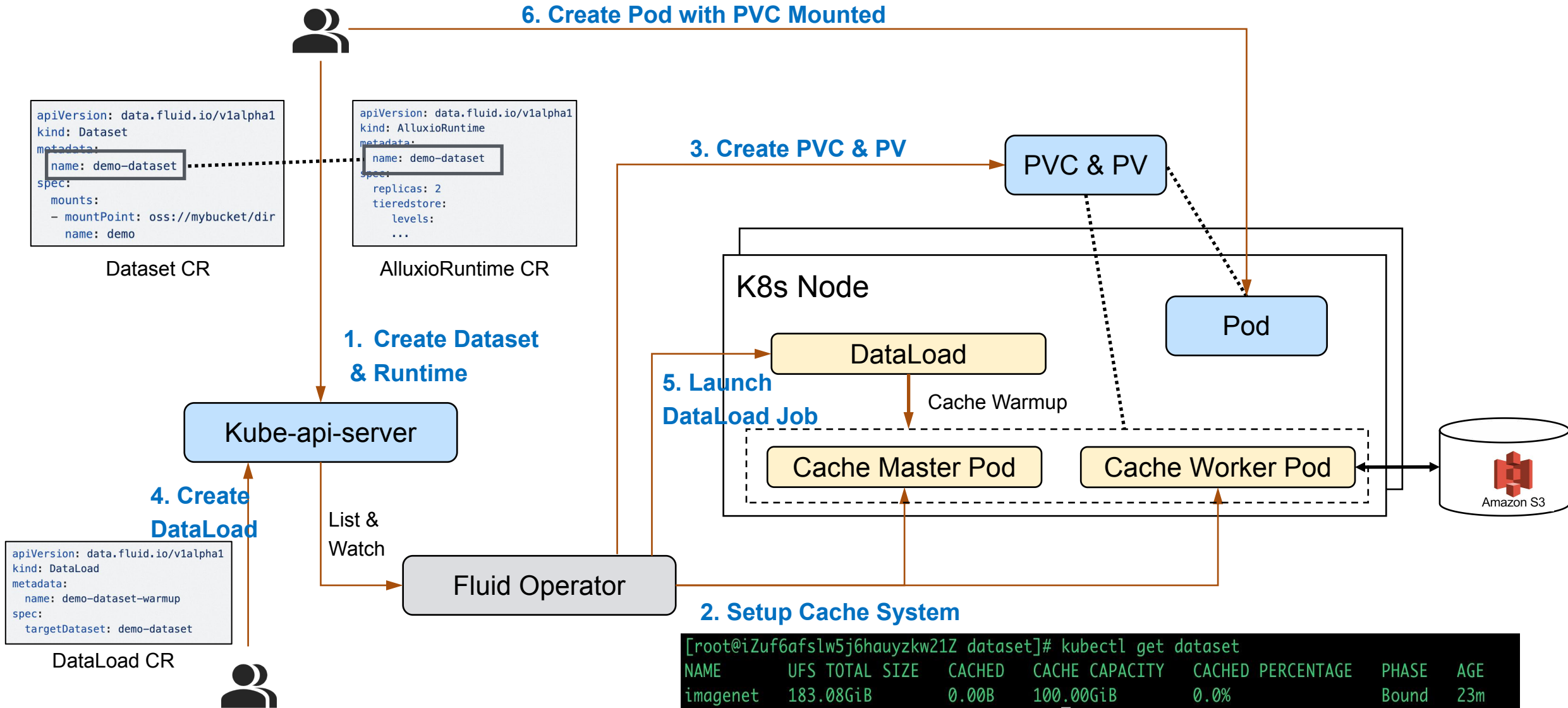
<https://github.com/fluid-cloudnative/fluid>



# Out-of-the-Box Distributed Cache



China 2024

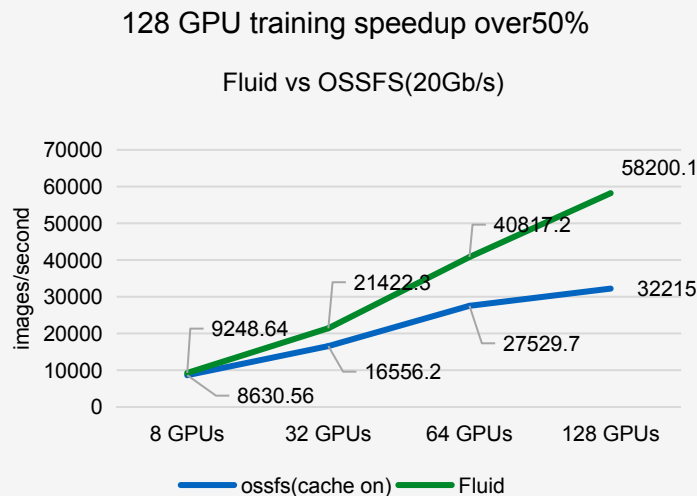


# Fluid user's scenarios

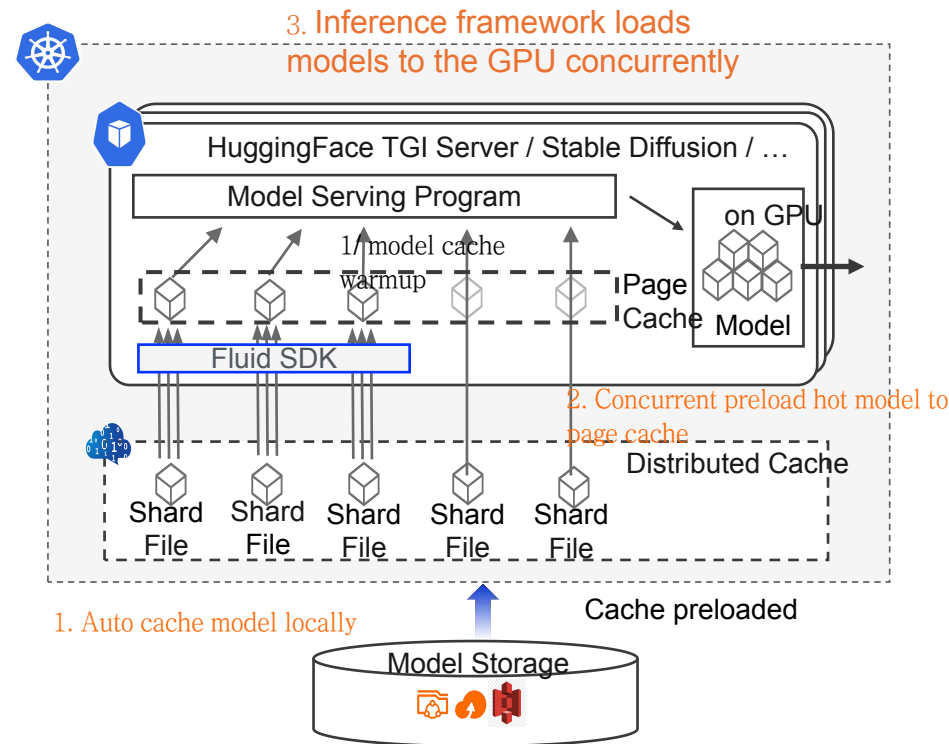


China 2024

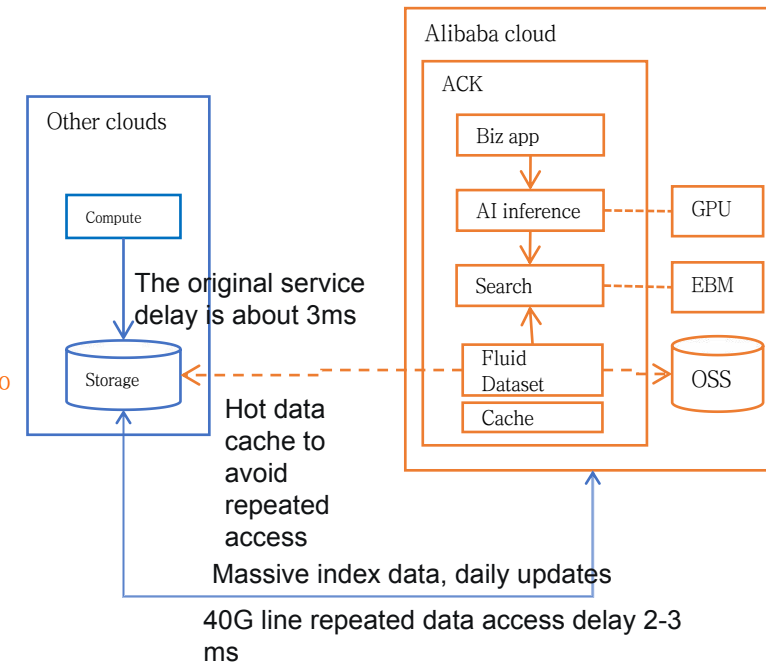
## 1. Accelerate AI distributed training and optimize scalability



## 2. Accelerate LLM inference service loading models



## 3. Hybrid cloud/multi-cloud data scheduling to accelerate access and optimize bandwidth costs



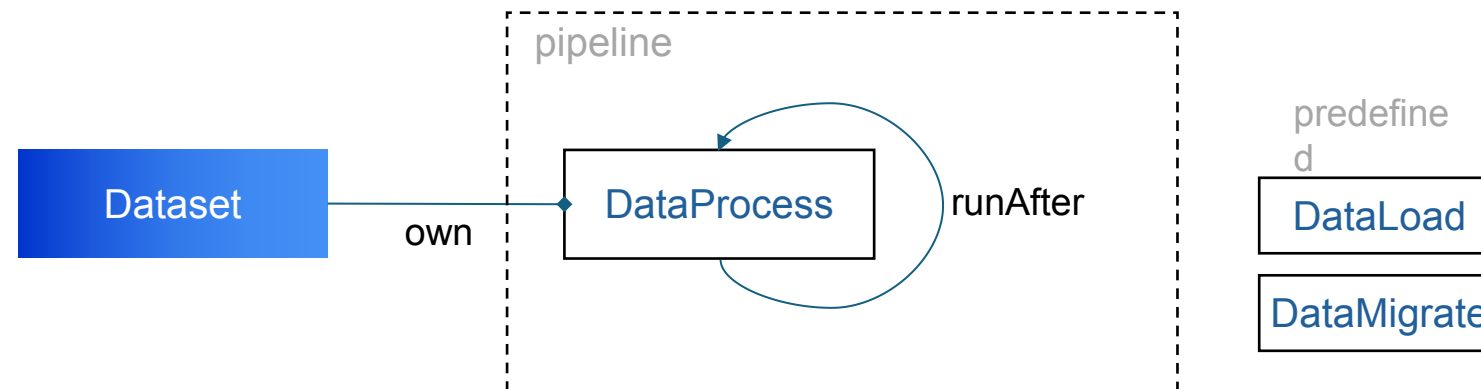
- Accelerate AI training by more than 30% and reduce the cold start delay of large model reasoning by 85%

# Fluid DataProcess

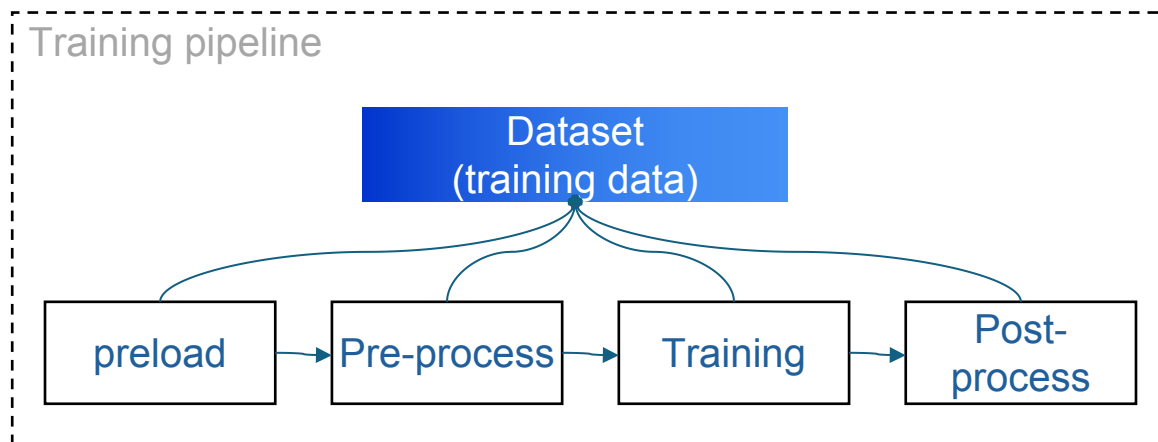


China 2024

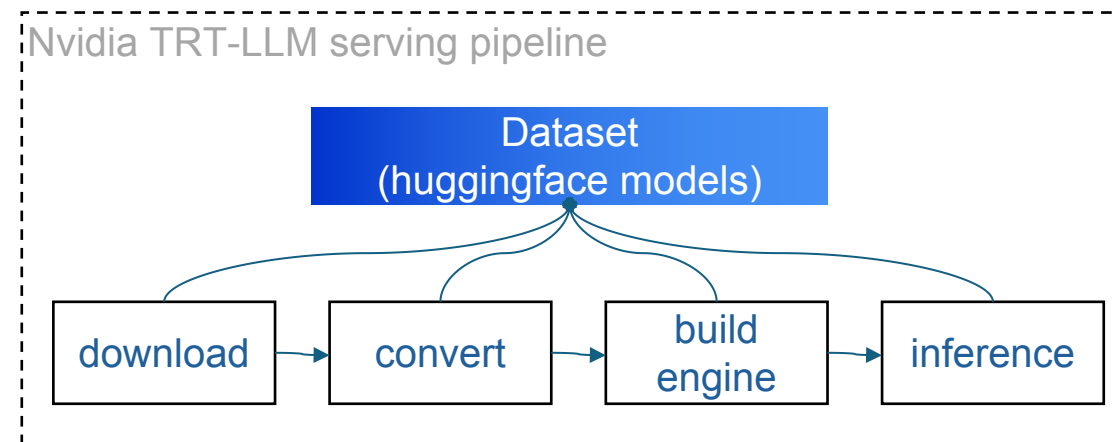
## Dataset centric processing pipeline



## Training pipeline



## Nvidia TRT-LLM serving pipeline





OCM – Manages multiple clusters in the K8s way



# What is OCM



China 2024



<https://open-cluster-management.io>

- An open-source CNCF Sandbox project
- Multi-cluster, multi-cloud Kubernetes orchestration, vendor neutral APIs
- Hub, spoke architecture with a centralized view of your entire fleet
- Modular and extensible
- Integration point for making Kubernetes capabilities multi-cluster aware

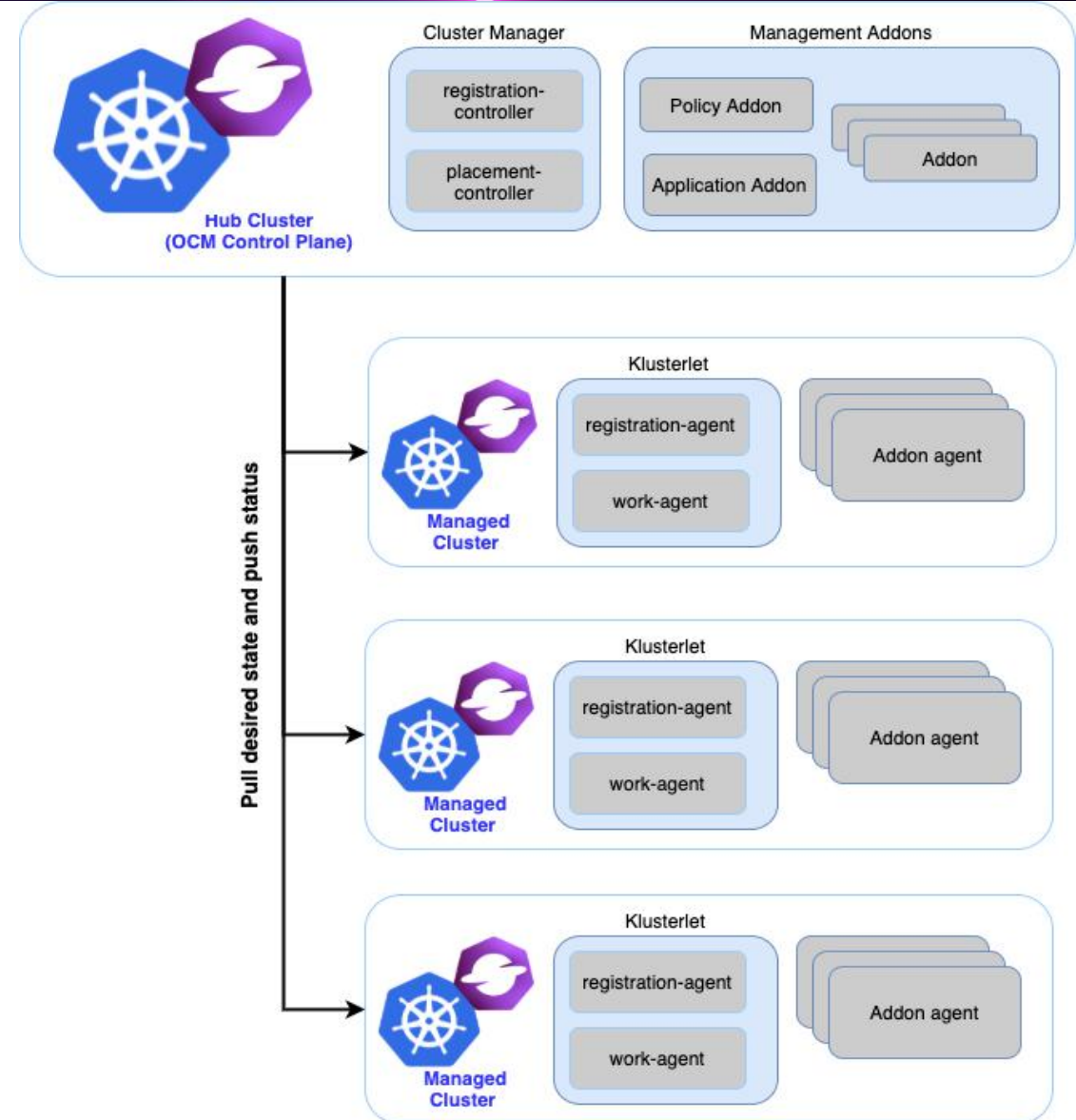


# OCM Architecture

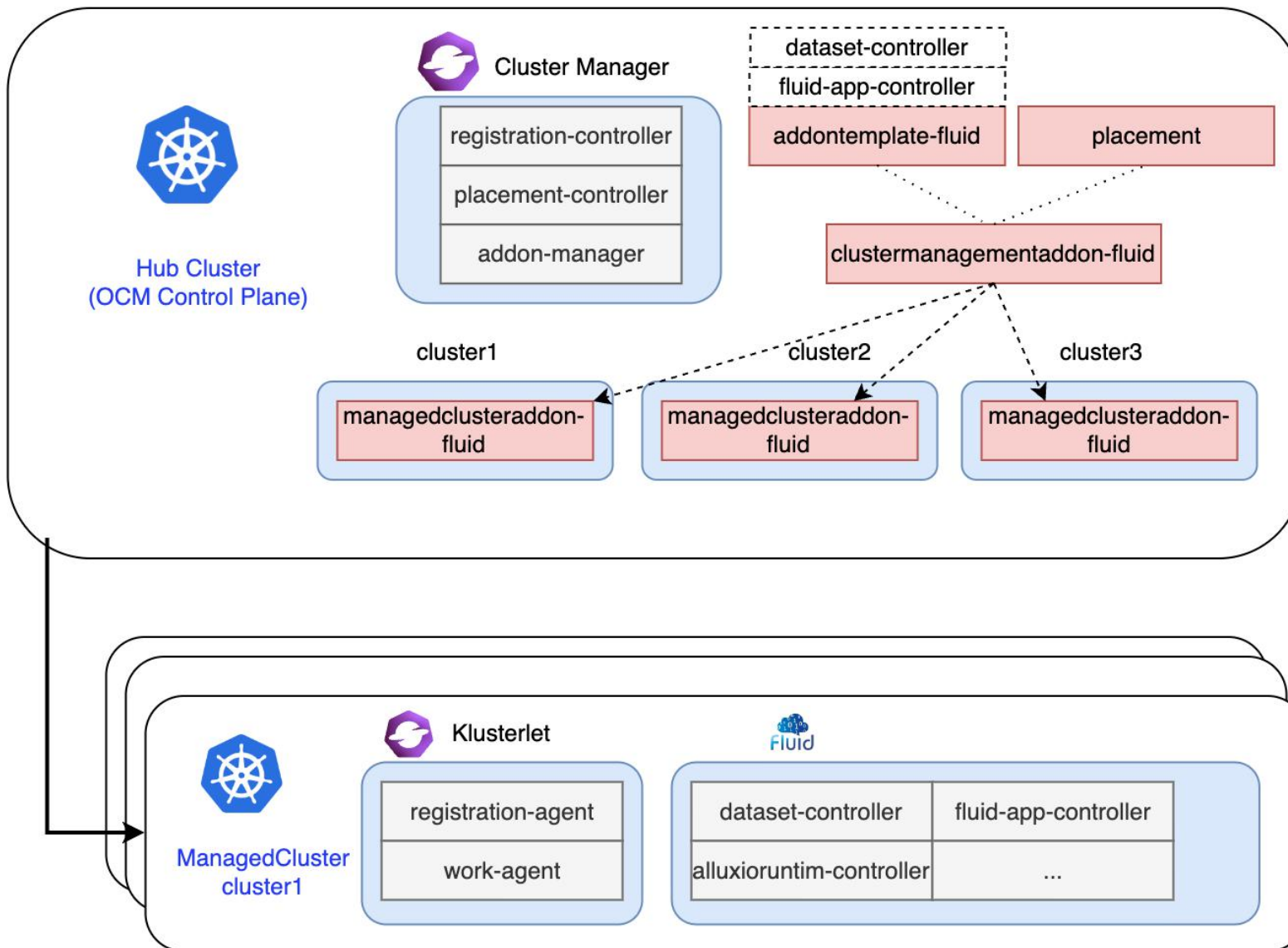


China 2024

- Registration
  - Managed cluster lifecycle manager
  - Initial double opt-in handshake
- ManifestWork
  - Deliver resources to one managed cluster
- Placement
  - Provide groupings of managed clusters based on cluster labels or claims
  - Extensible scheduling (CPU, GPU, etc) by the AddOnplacementScore
- ManifestWorkReplicaSet
  - Deliver resources to managed clusters selected by placements
- Add-ons
  - Policy Framework
  - Application lifecycle
  - Cluster Proxy
  - Managed serviceaccount
  - Fluid



# Integrate fluid into OCM



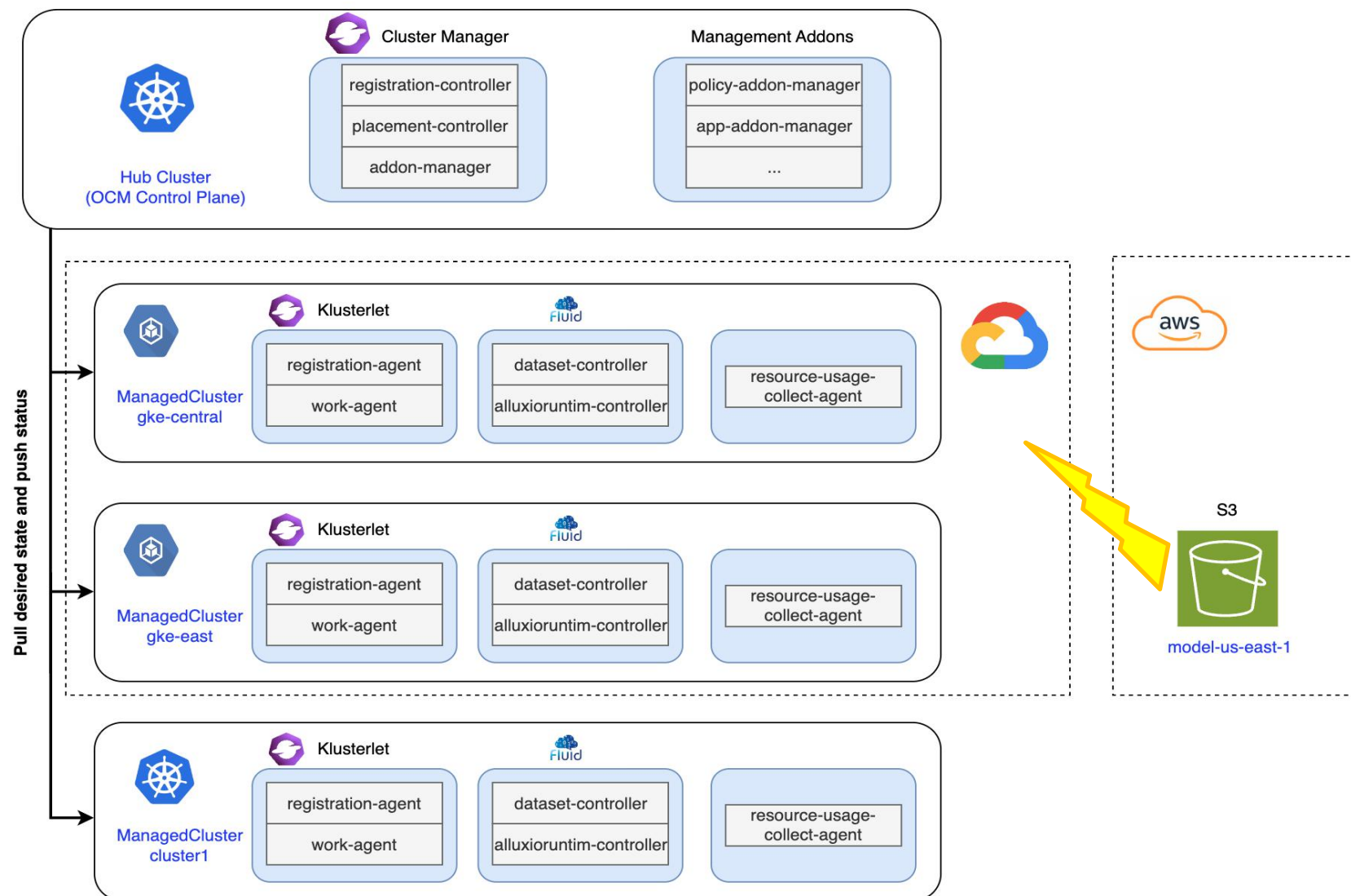


# Demo:

Deploy and scale LLM inference  
service crossing clouds quickly and easily



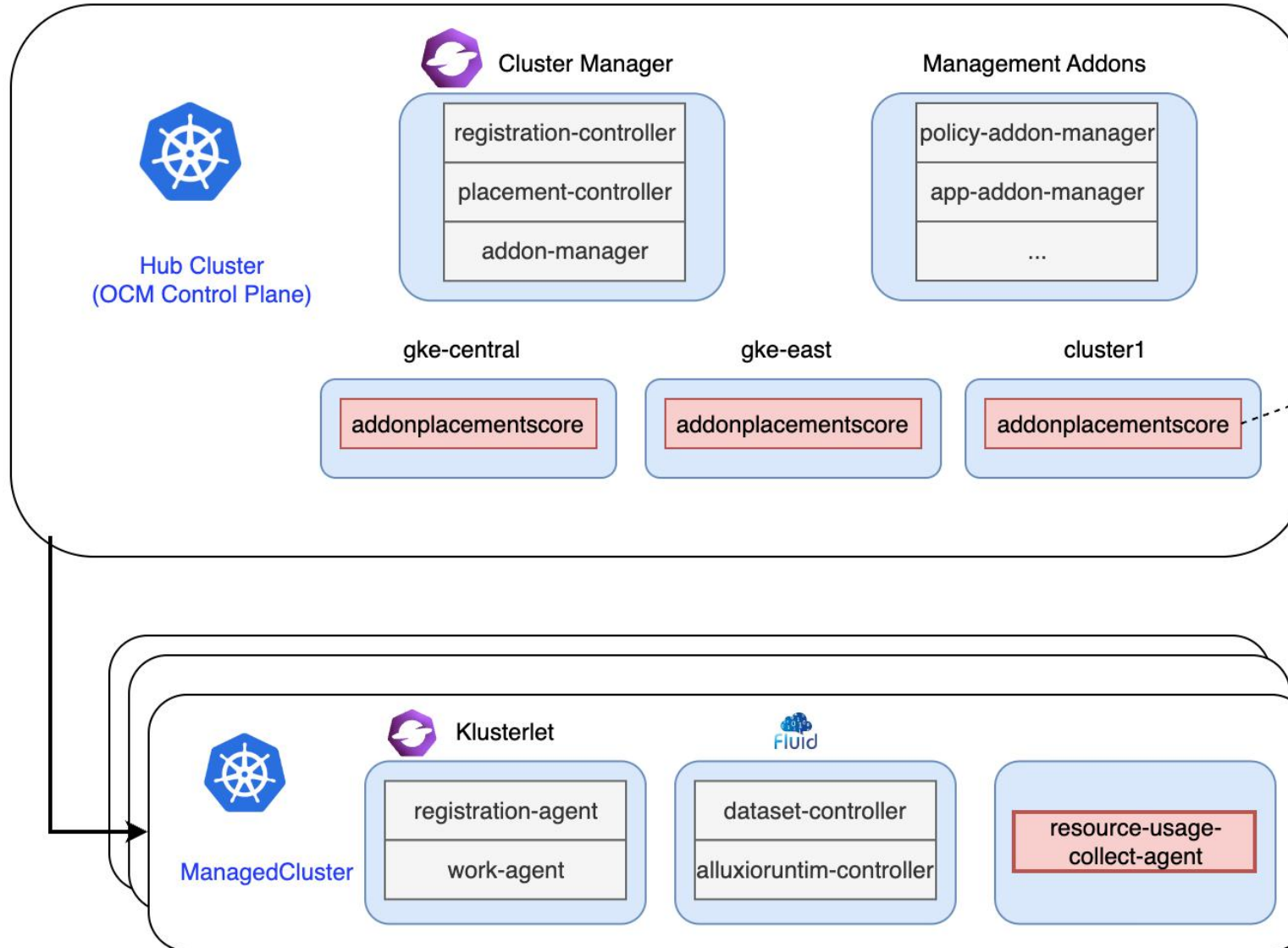
# Demo Scenario



# Demo Scenario



China 2024

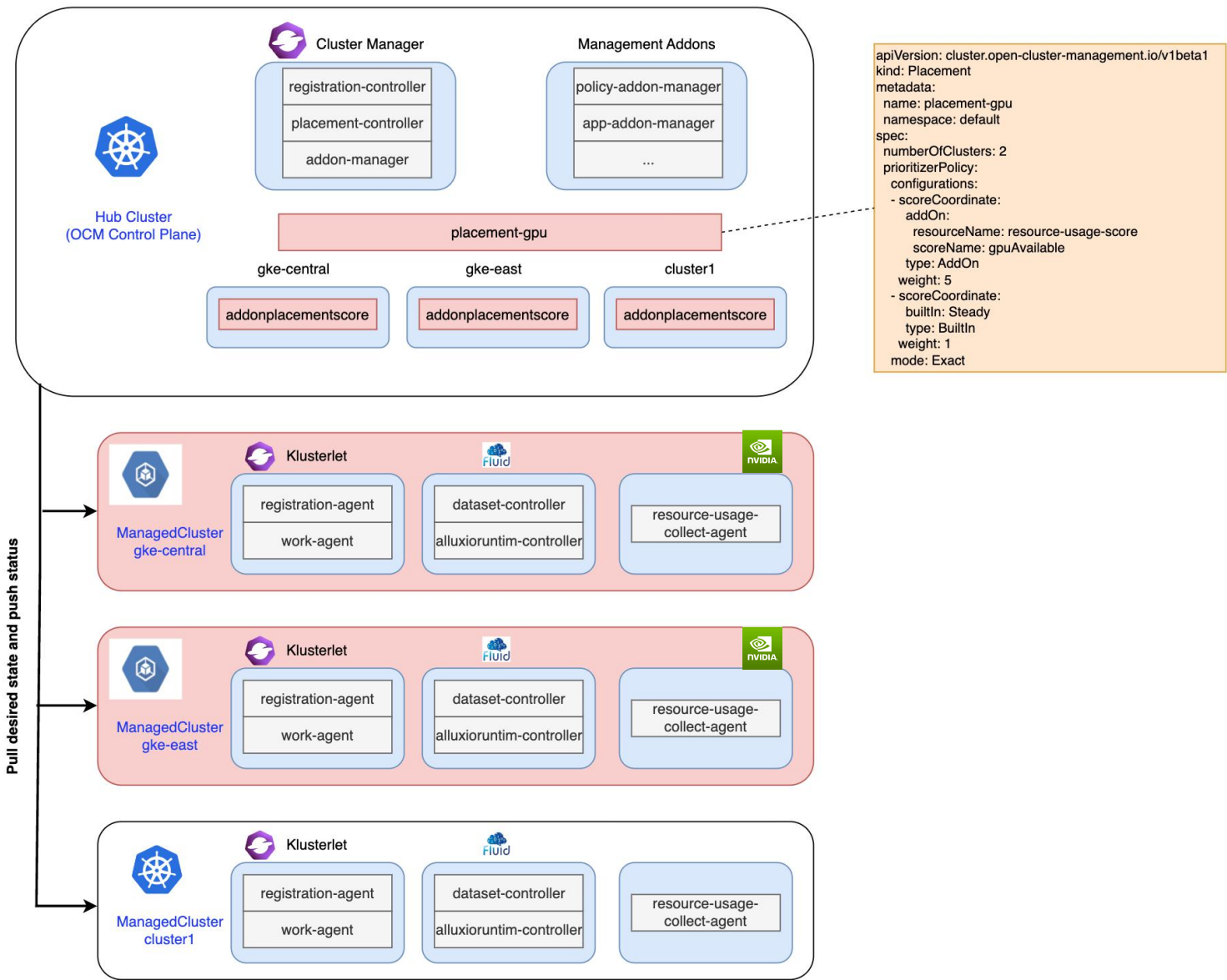


```
apiVersion: cluster.open-cluster-management.io/v1alpha1
kind: AddOnPlacementScore
metadata:
  name: resource-usage-score
  namespace: cluster1
status:
  scores:
    - name: cpuAvailable
      value: -70
    - name: memAvailable
      value: -96
    - name: gpuAvailable
      value: -100
    - name: tpuAvailable
      value: -100
```

# Demo Scenario



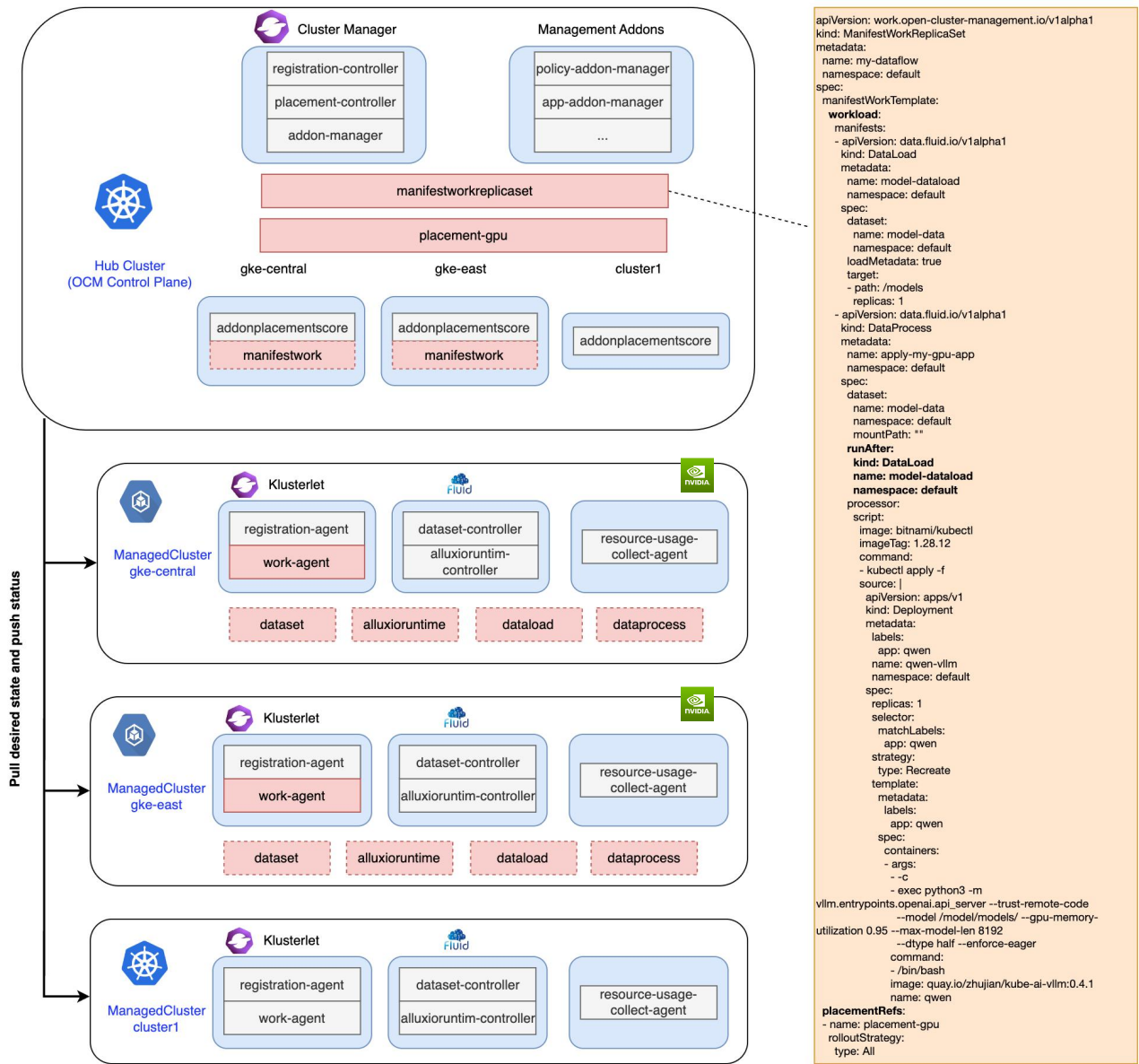
China 2024



# Demo Scenario

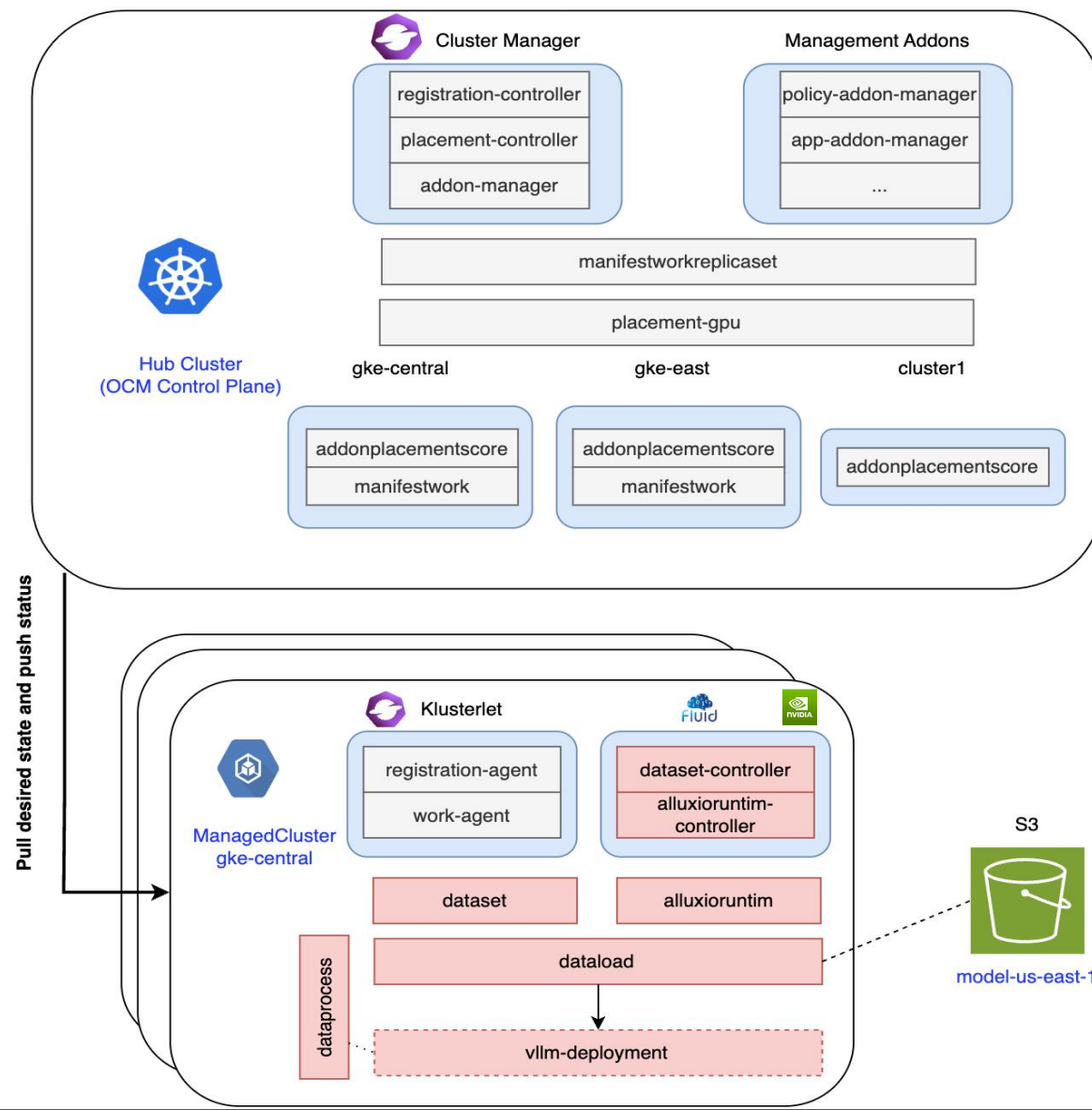


China 2024





# Demo Scenario



# Results Comparison



China 2024

1. When starting the vLLM application, using fluid to preload model from S3 takes only 50% time of loading with s3fs directly
2. When scaling the application, vLLM loads the fluid cached model locally, getting 5 times faster than the cold starting

	<b>dataload(s)</b>	<b>model loading(s)</b>	<b>server ready(s)</b> (including model loading time)	<b>sum(s)</b>
s3fs	n/a	132.6	147	147
fluid dataprocess	50	2.1	23.8	50+23.8=73.8
auto scale	n/a	2.1	14.5	14.5

# Future Works



China 2024

- GPU cost priority scheduling in multi-region and multi-cluster – [Skypilot](#)
- Multi-cluster training task scheduling based on priority queue– [Kueue](#), [Kube-queue](#)
- Unified traffic control and load balancing for LLM inference services across clusters – Service mesh, LLM Gateway
- Assign dynamic tuning/rescheduling of inference service instances across clusters



KubeCon



CloudNativeCon



China 2024

# Thank you!

&

# Question?