# AIStation
# Artificial intelligence development platform

Release AI computing power, accelerate intelligent evolution
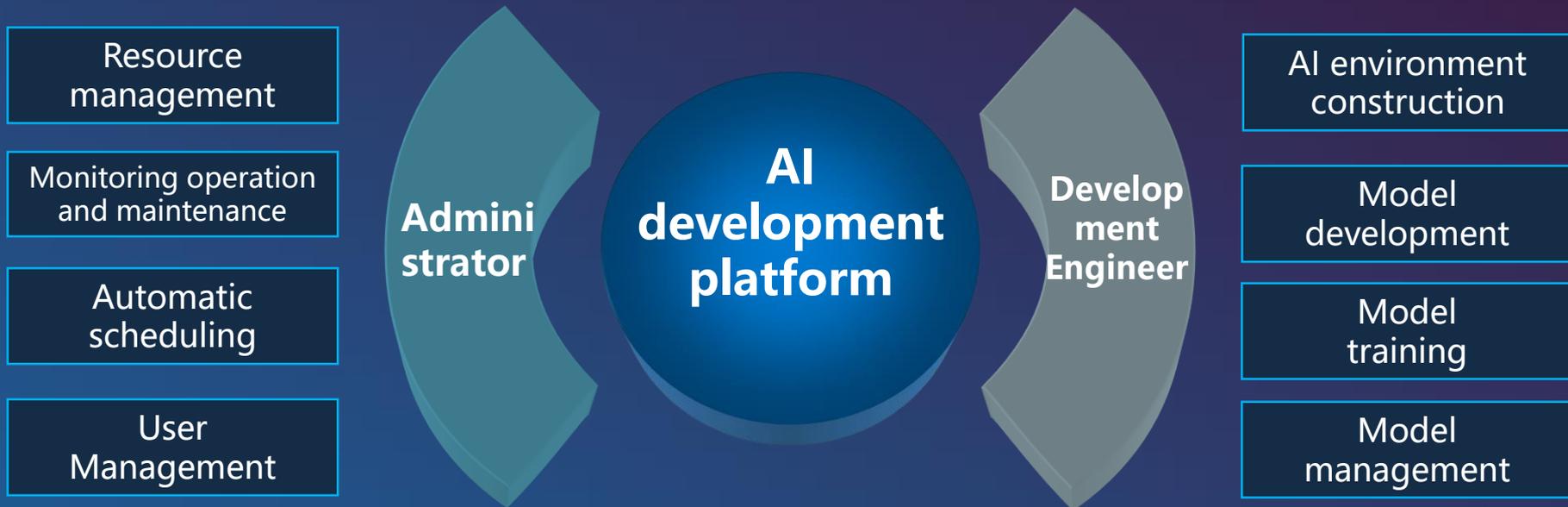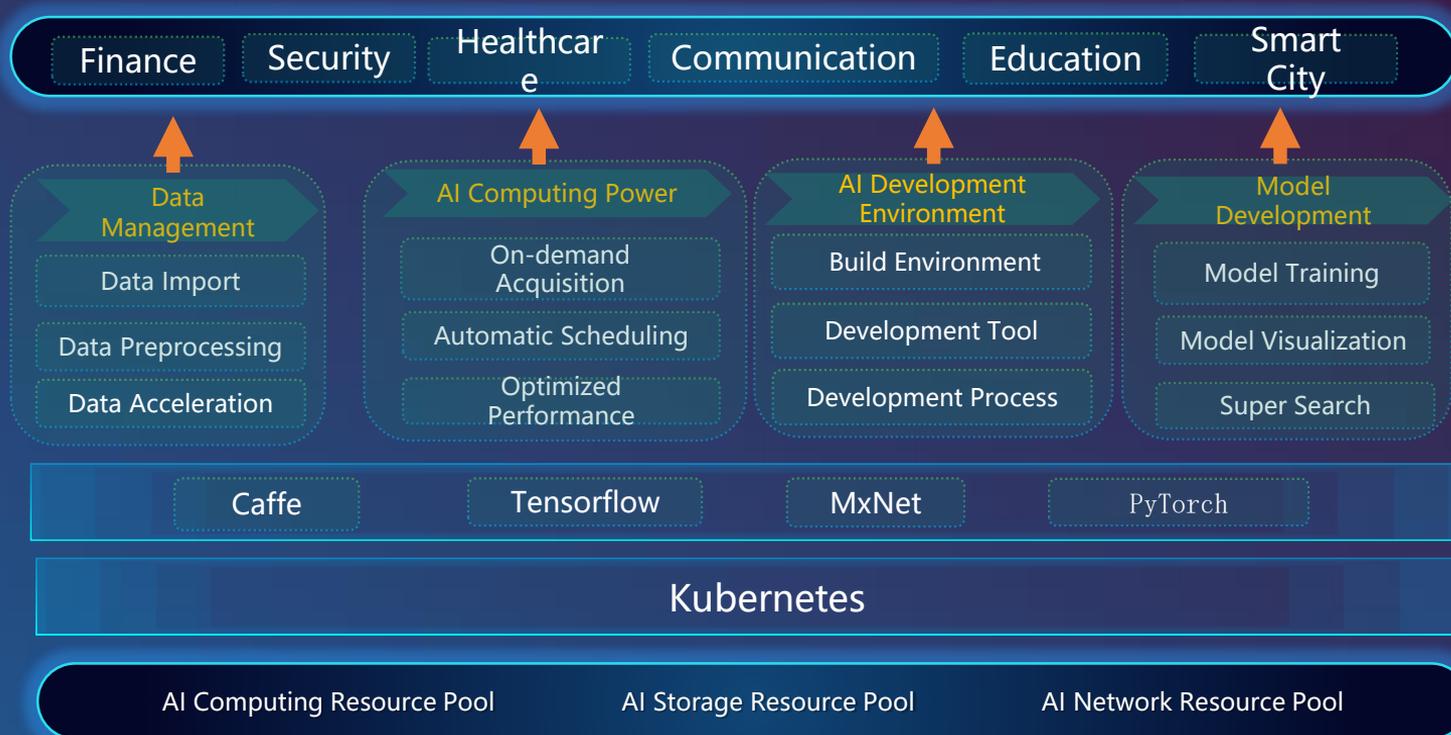
# contents

# 01 AIStation AI development platform

**Enterprise-level deep learning development scenarios**

**Unified management of scheduling computing resources, building AI development and training platform**

Resource management
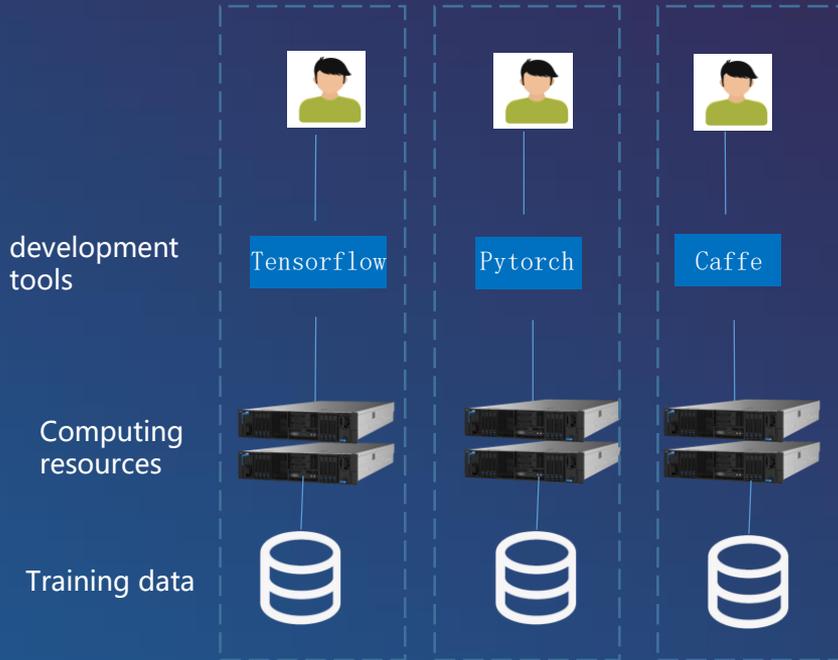
Monitoring operation and maintenance

Automatic scheduling

User Management

Administrator

**AI development platform**

Development Engineer

AI environment construction

Model development

Model training

Model management

# AIStation - 系统架构

# 01 AIStation AI development platform

## Distributed use of computing resources

development tools

Tensorflow  Pytorch  Caffe

Computing resources

Training data

Unable to complete statistics

Computing resources cannot run efficiently

System failure cannot be dealt with in time

**Admini strator**

# 01 AIStation AI development platform

Monitoring and operation and maintenance-unified monitoring and maintenance of AI resources and development services, and the AI platform operates continuously and efficiently

## Overall Monitoring

- Usage status of cluster resources such as GPU, CPU, and storage
- Computing node health and performance
- User task status and resource usage

## Resource Usage Statistics

- Cluster-level resource usage statistics
- Cluster-level task scale statistics
- User-level resource usage statistics
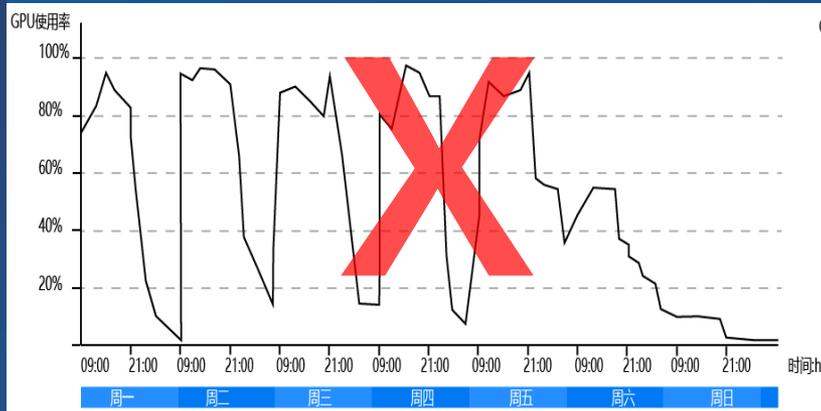- User-level task scale statistics

## System Alarm

- hardware malfunction
- System health status
- Computing resource utilization

# 01 AIStation AI development platform

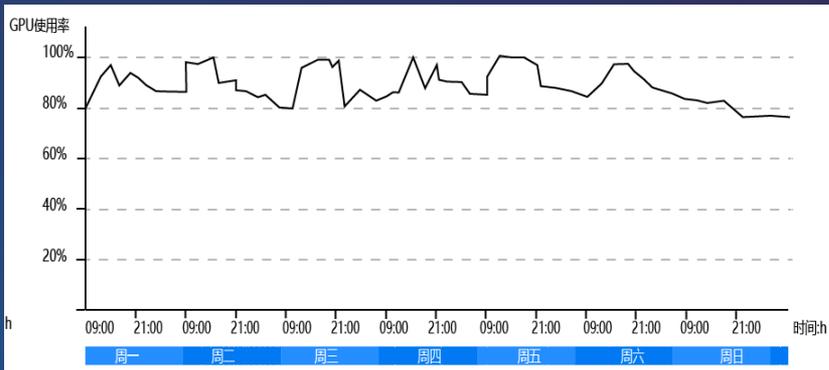◆ Low utilization of computing resources during holidays and rest hours



◆ Poor flexibility in computing resource allocation



Competing for computing resources

# 01 AIStation AI development platform

Resource management-automatic scheduling and allocation of AI computing resources, multi-user fair balance



Make full use of night and rest days for training tasks



**GPU sharing fine-grained distribution**

**Resource allocation management  Reduce over-occupancy**

- User GPU, CPU and storage resource quota limits
- Resources can be grouped according to different types of equipment

**Dynamic scheduling Improve resource reuse**

- Apply for computing resources on demand
- Auto release after calculation
- Training tasks are queued and hosted

**Intelligent resource scheduling to speed up the task**

- Multiple affinity scheduling
- GPU sharing fine-grained distribution

# 01 AIStation AI development platform

## Resource management-multiple resource scheduling strategies to improve resource utilization

**Overtime reminder of development environment to reduce long-term occupation**

The administrator can set the idle time and timeout policy (timeout reminder, whether to automatically stop), the user can manually restart the stopped timeout environment
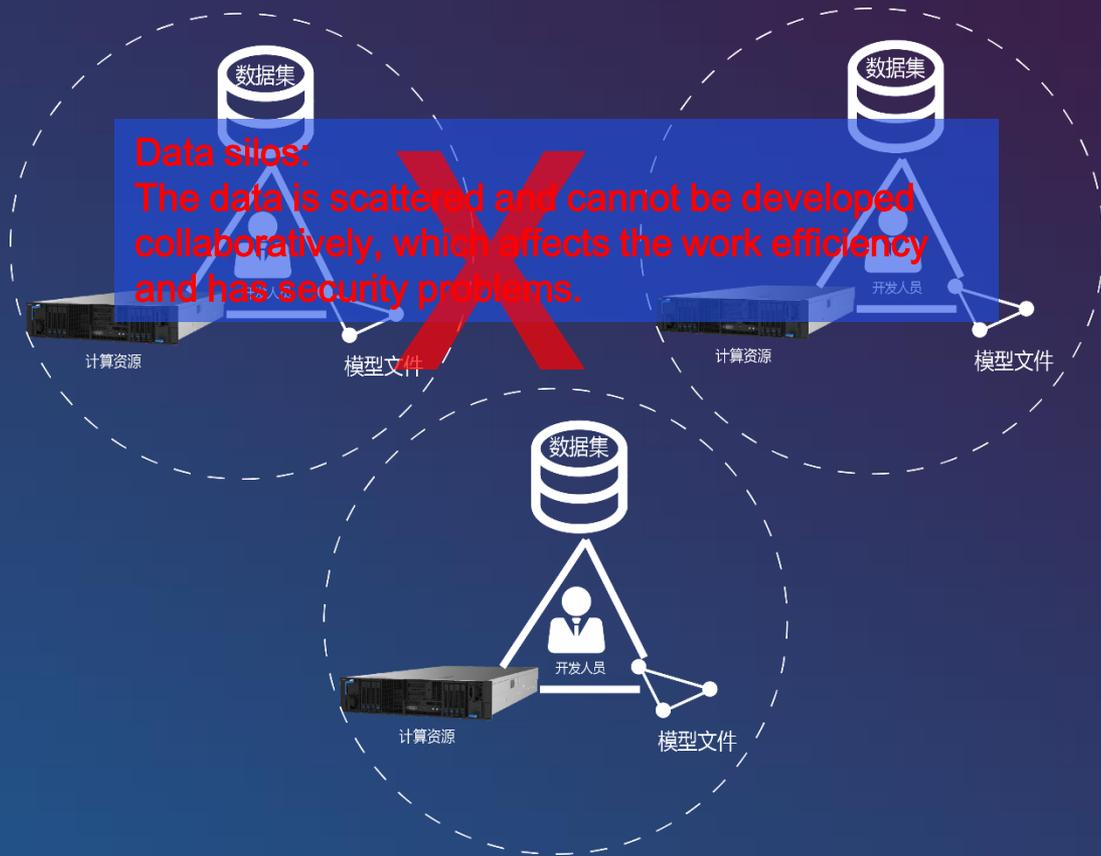
**Resource excess application reminder to improve task stability**

When the user application resource exceeds the physical specifications of the node or exceeds the user quota, it automatically reminds the user and automatically adjusts the number of applications

**Limit the number of development environments and resources to reduce over-occupancy**

Administrator settings can limit the number of platform users' development environment and the number of environment resource applications (GPU, CPU)

# 01 AIStation AI development platform



Data silos:
The data is scattered and cannot be developed collaboratively, which effects the work efficiency and has security problems.

# 01 AIStation AI development platform

## Data management centrally manages development data, with equal emphasis on reading speed and security

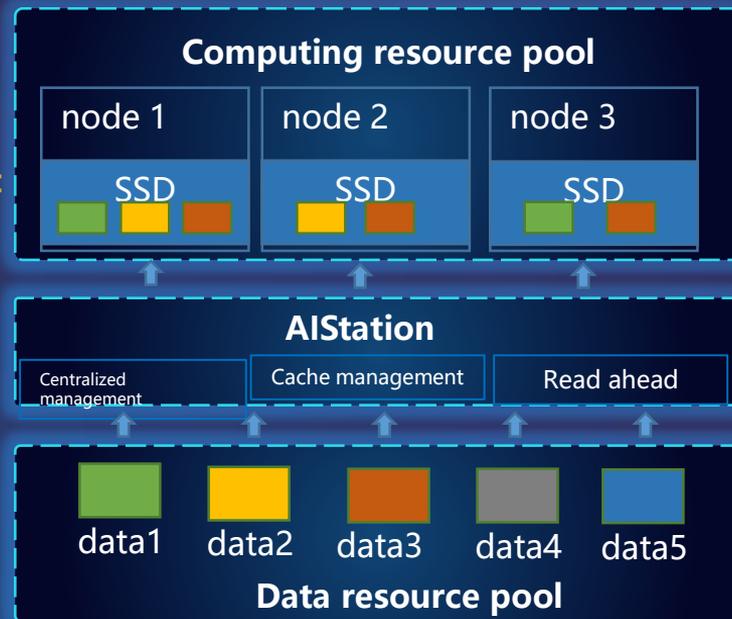**Accelerate data cache, effectively solve IO bottleneck**

- Automatic pre-reading of data sets
- Training tasks are first scheduled to nodes that cache data
- Cache data node automatic optimization and cleaning

**Unified data management to promote collaborative development**

- Personal data security isolation
- Collaborative development of data within the group
- Public data sharing application
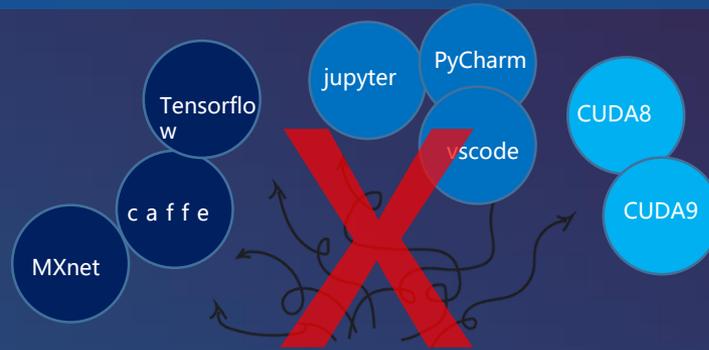- Unified management of data sets

**Data set security strategy**

- Flexible access control of data set access rights and download rights
- Multiple copies ensure safe data backup
- Support NFS, HDFS, BeeGFS, cloud storage system

### Computing resource pool

| node 1 | node 2 | node 3 |
|---|---|---|
| SSD | SSD | SSD |

### AIStation

| Centralized management | Cache management | Read ahead |
|---|---|---|

### Data resource pool

data1    data2    data3    data4    data5

# 01 AIStation AI development platform

Many development tools, slow task deployment
Development environment affects each other

Tensorflow

jupyter

PyCharm

CUDA8

vscode

caffe

MXnet

CUDA9

# 01 AIStation AI development platform

One-stop AI development environment to improve the efficiency of AI development engineers

AI Development Framework

AI Development components and tools

GPU Driver and development library

GPU computing resources

- Integrated mainstream AI development framework
- Import the image that comes with the installation package as needed
- Support tar package import and external mirroring of NGC and DockerHub platforms
- Open up data sets, computing resources, framework tools
- Provide a rich and complete AI development tool chain
- Connect to IDE tools such as pycharm and vscode
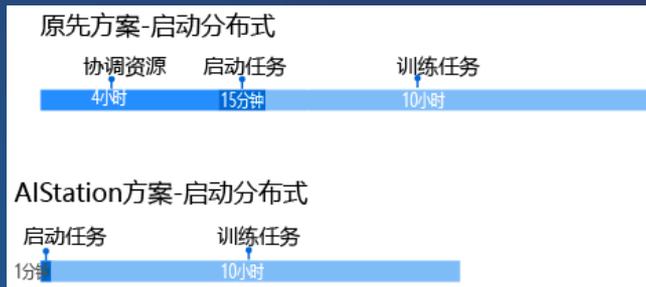
# 01 AIStation AI development platform

## Interactive model development Focus on model development

- Support one-click creation of historical development environment
- Support rapid online development
- Support one-click submission of training tasks
- Support batch submission of training tasks
- View task training progress at any time

# 01 AIStation AI development platform

**Automatically arrange and manage AI training tasks, accelerate AI development speed and shorten development cycle**

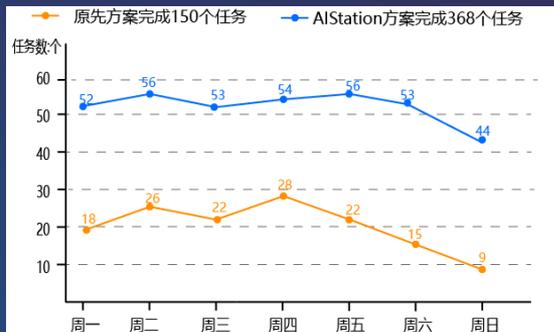Batch submit training tasks, model batch training

Training tasks are automatically queued, scheduled, and started

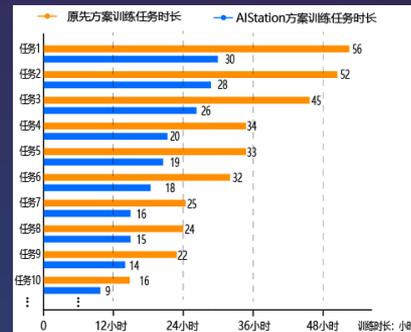Manage and monitor training process, progress and results

Distributed training automatic allocation resource orchestration process



Reduce time to deploy distributed tasks



Development efficiency increased by 2.3 times



Training speed is nearly double

# 01 AIStation AI development platform

Various training task management and scheduling, accelerate AI development speed and shorten development cycle

**Emergency task setting and priority scheduling strategy**

The administrator opens the emergency task permission for the user. After the user opens the permission, he can choose whether it is an emergency task when the task is submitted. Emergency tasks have the highest priority. Adjust the sequence of emergency tasks and designate an emergency task for priority scheduling

**User group task polling scheduling strategy fair sharing**

Support task scheduling according to user group polling, which can avoid centralized resource scheduling on a user group and further achieve the goal of fair sharing

# 01 AIStation AI development platform

## Fault tolerance mechanism to enhance the stability of the system

### Management node highly available

Supports health monitoring of the active and standby management nodes, HA status monitoring, and smooth switching of the active and standby machines, and the switching process does not affect the running business

### Abnormal warning and fault tolerance mechanism of computing server

Monitor the resource usage status and key service status of computing nodes to ensure the smooth operation of users' core business;

### Training task fault tolerance mechanism

System failure: Downtime, network disconnection, card dropout, the system can automatically start training tasks within 30 seconds, and can continue training from checkpoint

# XXSecurity case

## AI development platform

100+ training server, 800+ GPUs; 6 algorithm teams, 120+ developers

### Resource Management Strategy

Resource Use Grouping: Development: 32/Training: 700+; GPU Sharing Strategy: 4
Resource grouping strategy: P100, V100, 2080TI; resource quota policy: 24

| P100_shar | P100 | V100 | 2080Ti | 1080TI _dev |
|---|---|---|---|---|
| Quantity: 64 GPU<br>Sharing: 2<br>Uses: training<br>User: ALL<br>Quota: None<br>SSD cache | Quantity: 240 GPU<br>Sharing: none<br>Uses: training<br>User Behavior Analysis<br>Quota: 16<br>SSD cache | Quantity: 64 GPU<br>Sharing: none<br>Uses: training<br>User: ALL<br>Quota: 16<br>SSD cache | Quantity: 360 GPU<br>Sharing: none<br>Uses: training<br>User: Robot, face recognition<br>Quota: 16<br>SSD cache | Quantity: 32 GPU<br>Sharing: 4<br>Uses: development and debugging<br>User: ALL<br>Quota: None<br>SSD cache |

Sharing strategy: 32 GPU supports 120 people to develop and debug at the same time;

Task queuing: make full use of night and holiday time, increase utilization rate by 20%

Dynamic allocation: 4-5 tasks are trained at the same time, and the development cycle is shortened to 1/3

Resource utilization: increased from 70% to 90%

# One plus mobile phone case

## Business scene

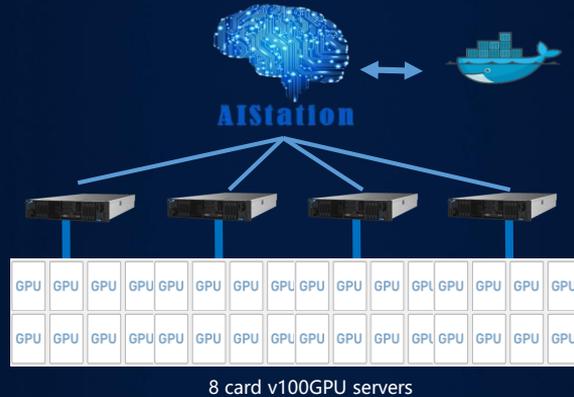| | |
|---|---|
| Object scene classification | Object detection |
| Target tracking | Hairline semantic segmentation |

## User issues

- Insufficient resources and low utilization
- Development environments conflict
- Inefficient creation of AI environments

## AIStation Solution

AIStation

8 card v100GPU servers

### Solve issues

- Centralized management and operation of computing resources, user quota restrictions
- GPU sharing strategy, creating multiple containers with a single GPU card
- Training tasks are queued for hosting, using night and holiday training tasks
- Quickly create development environments with Docker, and isolate each other without affecting each other
- Built-in various AI framework images, compatible with web open source images

## AIStation Value

**GPU resource usage 75% - 95%**

**The efficiency of cluster use has been greatly improved**

- Reduce user development resource usage; improve training task resource utilization
- The average daily GPU usage time is increased to 22-24 hours

**Improve developer productivity**

- Reduces a lot of repetitive operations
- Web UI operation makes the development process more convenient

# Case

## XX intelligent vision

**Business scene**

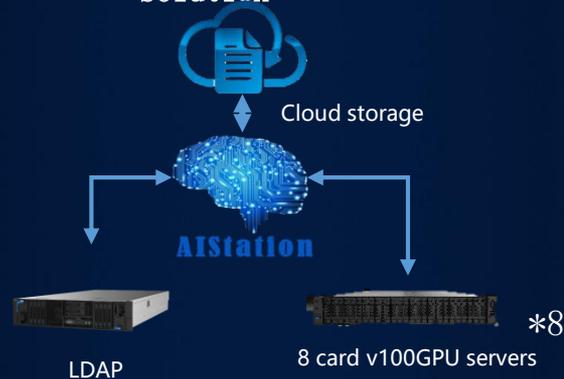| | |
|---|---|
| Computing vision | Cognitive science |
| Language dialogue | Machine learning |

### User issues

- GPU resources are tight, requiring a lot of manpower for maintenance and coordinated resource allocation;
- LDAP, cloud storage, and GPU servers do not have a unified connection. Each system is separately maintained and managed, which is inefficient;

## AIstation Solution

Cloud storage

AIStation

LDAP

8 card v100GPU servers   *8

### Solve issues

- GPU resource management and operation and maintenance, set 8 machines as 3 resources and set user quota;
- AIstation docks with cloud storage and downloads the data locally for sharing by multiple people;
- Set a data cache policy to regularly clean up the data;
- AIstation supports API interfaces such as LDAP and manages users in a unified manner;

## AIstation Value

**Improve manager productivity 60%**

**Reduce the cost of operation and maintenance personnel**

- Reduce the workload of two-thirds of operation and maintenance personnel
- Docking cloud storage, reducing operating costs and improving task training efficiency.
- Improve user management efficiency and reduce complexity after connecting to the LDAP system

# Case

## XX Bank

### Business scene

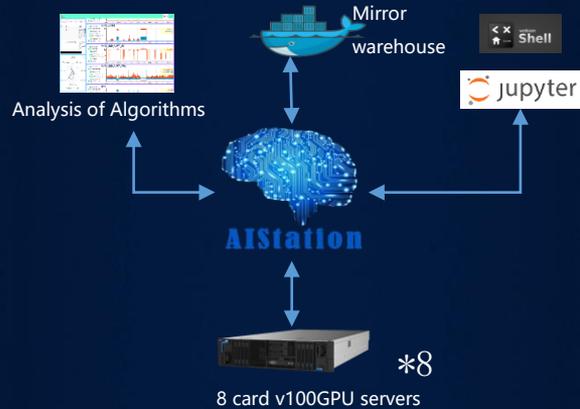| | |
|---|---|
| Computing vision | Target Detection |
| Language dialogue | Speech Recognition |

### User issues

- It is forbidden to access the external network. Various IDE tools cannot be installed directly.
- Algorithm error can not be quickly located

## Alstation solution



Mirror warehouse

Shell

Jupyter

Analysis of Algorithms

AIStation

*8

8 card v100GPU servers

### Solve issues

- Support for Web Shell tools and Jupyter tools
- Alstation private image repository built-in AI image of each framework
- Alstation supports training task process monitoring, occupancy of each resource and algorithm

## Alstation Value



**Developer productivity 40%**

### Improve developer productivity

- Help customers quickly deploy development work and quickly enter fast mode
- Quickly locate the bottleneck or bug of the algorithm and improve the training efficiency of the algorithm.

# Thanks