

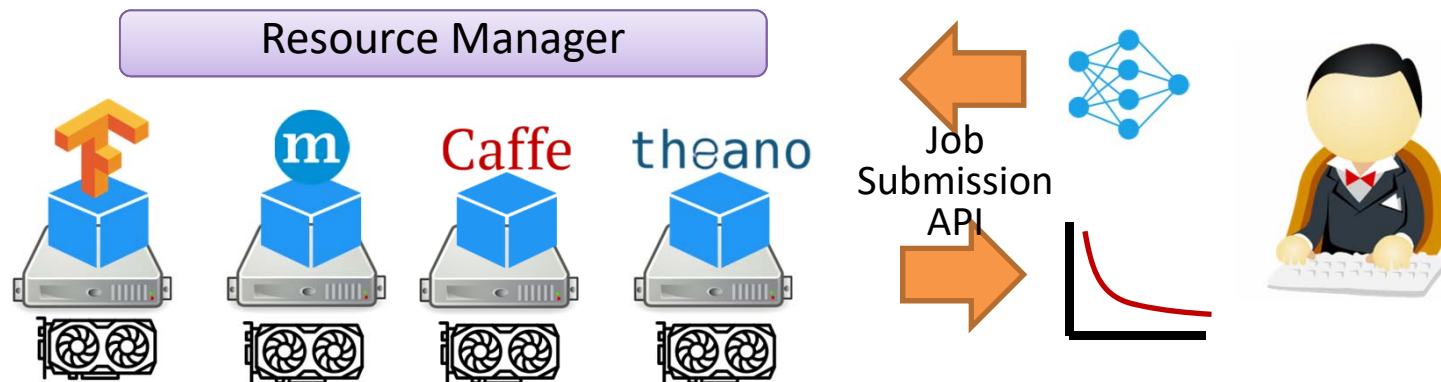
AI Cloud Service Platform

- Objective: Build an AI cloud service platform for deep learning computations
- Motivation:
 - **Expensive** computing resources
 - **Complex** resource management & job tuning
 - **Diverse** training model and hardware devices
- Approach:
 - An deep learning computing **job submission API**
 - Container and GPU based **shared resource pool**
 - Elastic, fault-tolerant and **self-managed computing service**
 - Resource usage and job execution **optimization**

Training hour 100



150,000 USD



System for AI

- Customize the system design for AI workload from *hardware resource allocation, framework optimization* to *computation parallelization*
 1. Deploy a **container-based GPU** cluster computing service in a cloud platform like **OpenStack**
 - **A proto-type system has been developed**
 2. Guarantee resource allocation **fairness and performance isolation on a shared GPU** device among containers
 - **Modifying the scheduler of CUDA driver**
 3. Develop an **elastic deep learning framework** to adjust the resource allocations of a training job at run time
 - **Modifying TensorFlow based the solution from Baidu**
 4. Improve the **parallel efficiency of distributed deep learning computations** across compute nodes or devices
 - **Modeling the time of operations and communications for model parallelism**

AI for System

- Apply deep learning to **solve fundamental system management problems**, such as
 - 1. Predict the execution time of a (DAG) job** when it is running with varied job configurations and hardware devices
 - **Achieved 90% accuracy for Hadoop and Hive jobs for job execution optimization**
 - 2. Auto-tune job parameters** and resource **configurations**
 - **Reduce the cost and time of running Hadoop jobs by 4x**
 - 3. Detect machine failures** or job **execution anomaly**
 - **Cowork with Ill to solve noise neighbor problem in cloud**
 4. Re-design **job scheduling and resource allocation algorithms** using reinforcement learning
 - **Developing the scheduler for our DL computing cluster**