

# Boost ML on Ray

**Enabling** Heterogeneous AI Accelerators

**Seamlessly**

Tiejun Chen  
VMware, OCTO  
6/18/2023

# Agenda

Empower AI everywhere

- Problem area on ML/AI
- ML on Ray
- How did we get to boost ML on Ray ?
- Demo
- What's likely next ?

# Problem area

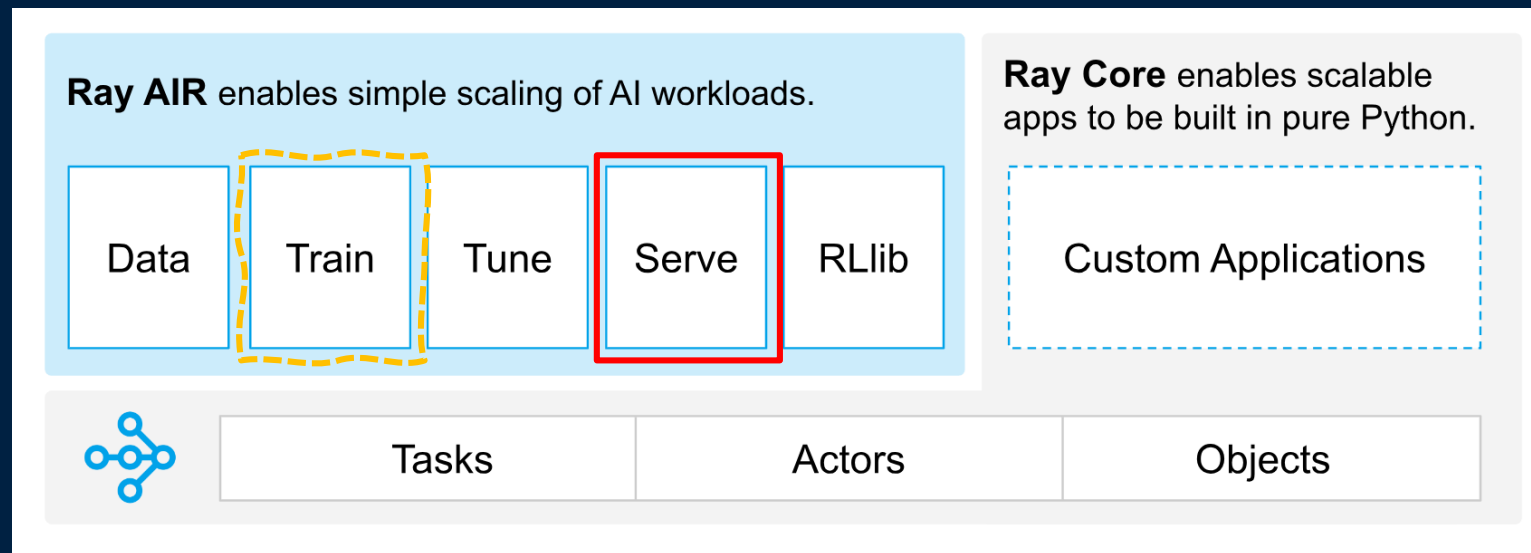
## Towards modern AI application centric platform

- Heterogeneous AI HW accelerators
- Various upstream ML frameworks
- Hard to exploit the best performance
- No such a modern AI platform with cloud native principle

# ML on Ray

## Brief

### ❖ *Flexible distributed Python for machine learning*



<https://github.com/ray-project/ray>

#### ➤ Ray Serve

- ❖ Be framework-agnostic
- ❖ Doesn't perform any model-specific optimizations

# How did we get to boost ML on Ray ?

## Project Yellowstone - Goal

- Build end-to-end ML service on Kubernetes from cloud to edge
  - ❑ Enable CRD based accelerators for ML serving
  - ❑ Boost ML by transparent backend acceleration

# How did we get to boost ML on Ray ?

Project Yellowstone - Enable CRD based local accelerators for ML serving

- Node Feature Discovery
- Device plugins
- Node selector
- Kubernetes scheduler

# How did we get to boost ML on Ray ?

## Graph compilers

- What ?

ML frameworks

Graph compiler

Lib/drivers

Computational Graph



Operations

- Graph compilers

- ☐ Apache TVM
- ☐ Nvidia TensorRT
- ☐ Intel OpenVINO
- ☐ AMD ROCM
- ☐ Xilinx vitis AI



# How did we get to boost ML on Ray ?

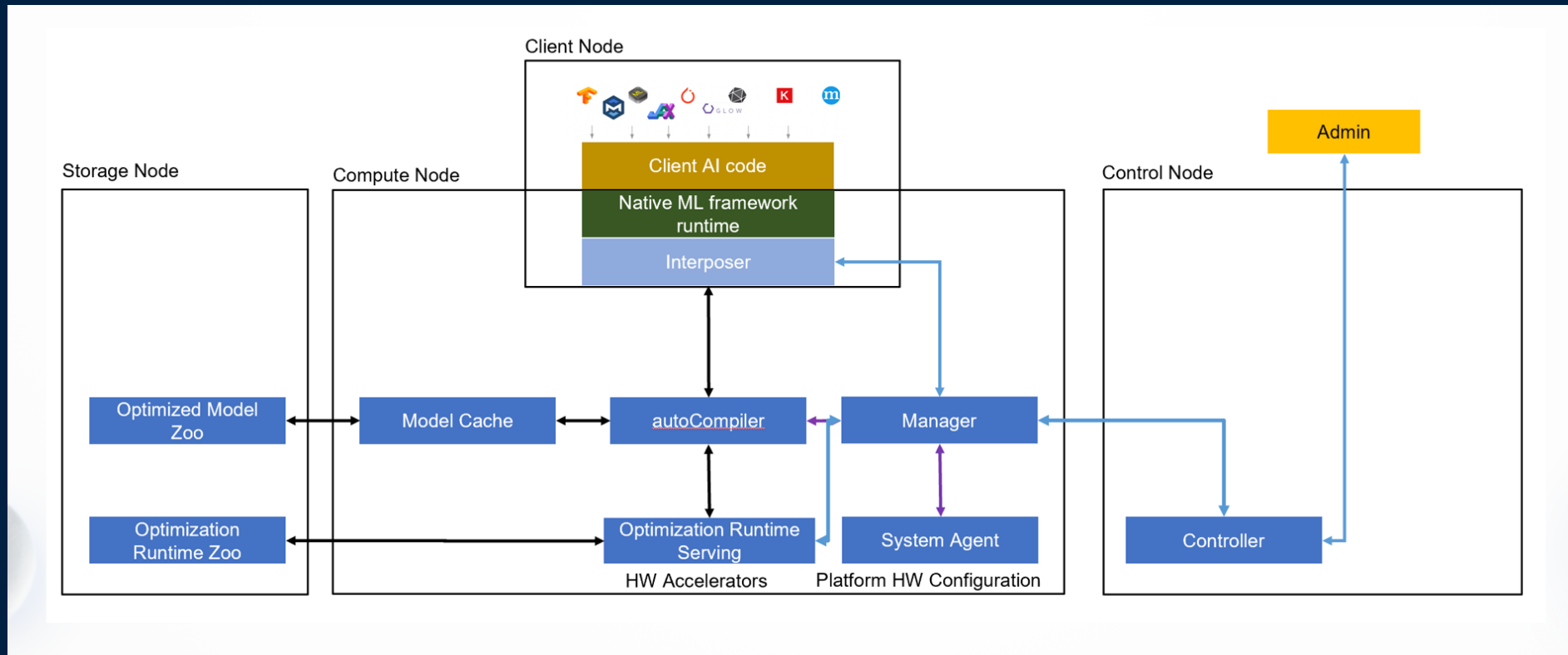
## Project Yellowstone - Overview

- Target
  - Boost ML/AI by enabling ML upstream frameworks seamlessly with graph compilers
- Design
  - Build ML Boost Serving System
    - ❑ Backend
    - ❑ Automated
    - ❑ Unified server architecture
- How
  - Interpose ML framework API
  - Built-in graph compilers processing – Auto {detecting, compiling, scheduling, inferencing, etc}



# How did we get to boost ML on Ray ?

## Project Yellowstone - Architect



# How did we get to boost ML on Ray ?

## Project Yellowstone – Seamless interposer

- Runtime interposer
  - ❑ Target to key APIs
  - ❑ Mapping between ML Frameworks APIs and Backend APIs
- Example
  - ❑ Tensorflow - Python
    - `tensorflow.keras.models.load_model = booster_load_model`
    - `tensorflow.keras.models.Model.predict = booster_predict`
  - ❑ Tensorflow Serving - C++
    - `predict --> session->Run()`
    - Hijack the process at runtime to call `booster_predict`

# How did we get to boost ML on Ray ?

## Project Yellowstone – Demo

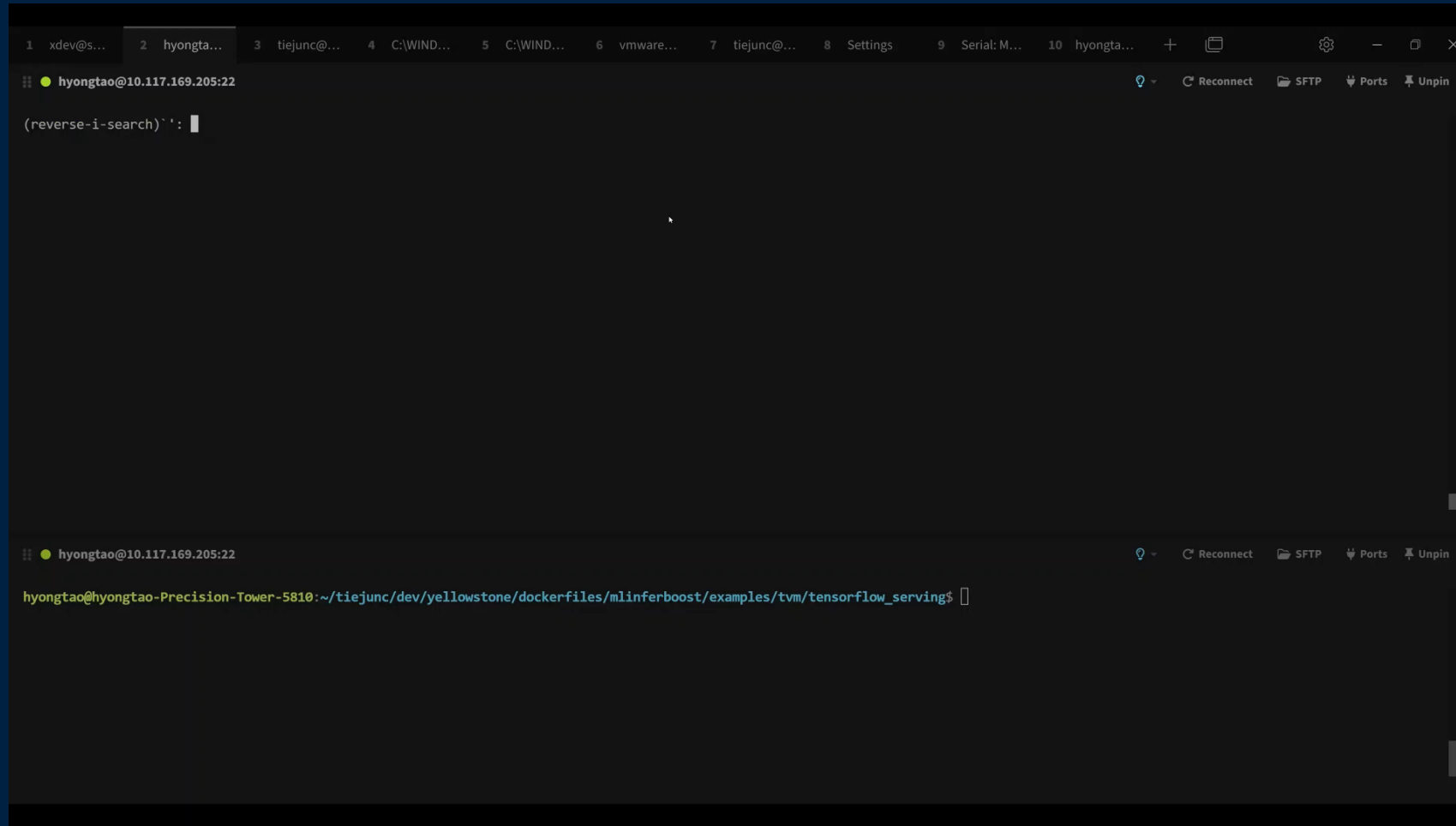
- TorchServe on GPU with backend TVM acceleration
- Tensorflow Serving on GPU with backend TVM acceleration
- Pytorch on GPU with backend TVM acceleration on Ray

# Demo I

## TorchServe on GPU with backend TVM acceleration

# Demo II

## Tensorflow Serving on GPU with backend TVM acceleration



The screenshot shows a terminal window with a dark background. The top bar displays a list of tabs: 1 xdev@s..., 2 hyongta..., 3 tiejunc@..., 4 C:\WIND..., 5 C:\WIND..., 6 vmware..., 7 tiejunc@..., 8 Settings, 9 Serial: M..., 10 hyongta..., followed by a plus sign, a document icon, a settings icon, and window control icons. The main area of the terminal shows a prompt `hyongtao@10.117.169.205:22` with a green dot icon to its left. Below the prompt, the text `(reverse-i-search)`':` is displayed, followed by a cursor. At the bottom of the terminal, another prompt `hyongtao@hyongtao-Precision-Tower-5810:~/tiejunc/dev/yellowstone/dockerfiles/mlinferboost/examples/tvm/tensorflow_serving$` is visible, also with a green dot icon to its left. The right side of the terminal window features a toolbar with icons for Reconnect, SFTP, Ports, and Unpin.

## Pytorch on GPU with backend TVM acceleration on Ray

[illegible]

# What's “here” now ?

## Project Yellowstone

- ML frameworks
  - ❑ Tensorflow/Pytorch/ONNX, Tensorflow Serving, TorchServe, KServe, etc
  - ❑ Ray
- Backend acceleration technologies
  - ❑ Apache TVM
  - ❑ Intel OpenVINO
  - ❑ Nvidia TensorRT
  - ❑ Xilinx vitis AI
- AI HW accelerators
  - ❑ {Nvidia, AMD, Intel} GPU, Xilinx FPGA
  - ❑ CPU

# What's likely next ?

## Project Yellowstone++

- Moving only ML Inference to ML {Training, Inference}
- Moving towards multi-AI cloud



# Empower AI everywhere

Thank you !



## Q & A

Tiejun Chen <[tiejunc@vmware.com](mailto:tiejunc@vmware.com)>