# Project MLInferBooster
*Agenda*

- Background – our focus

- Project MLInferBooster Introcution

- Next

# ML/AI upstream frameworks
## *Problem area*

- Limited AI accelerators enablement
  - Types - GPU/FPGA
  - Vendors - Nvidia, AMD and Xilinx
  - ❖ Other accelerators ?

- Limited AI performance
  - Focuses on training
  - ❖ Differentiated optimization technologies ?
    - ❑ Training vs Inference

- Limited to native environment
  - Host = Target
  - Running on real AI accelerator
  - ❖ Cross arches? No HW accelerators?

# SW Accelerator
## *Graph compiler*

- The levels of the ML pyramid
    - The low-level libraries
    - Deep learning frameworks
    - Compiler
- The Graph compiler
    - What
    - The goal
    - Projects
        - ❏ TensorRT
        - ❏ XLA
        - ❏ Glow
        - ❏ TVM
        - ❏ …

# Project MLInferBooster
## *The backend accelerator – TVM*

- TVM
  - "Apache TVM is a compiler stack for deep learning systems."

- Why
  - Open source
  - TVM supports most AI/ML frameworks
  - TVM targets various types of AI accelerators
    - Including CPU
  - Cross-compiling
    - Host /= Target
  - Good performance
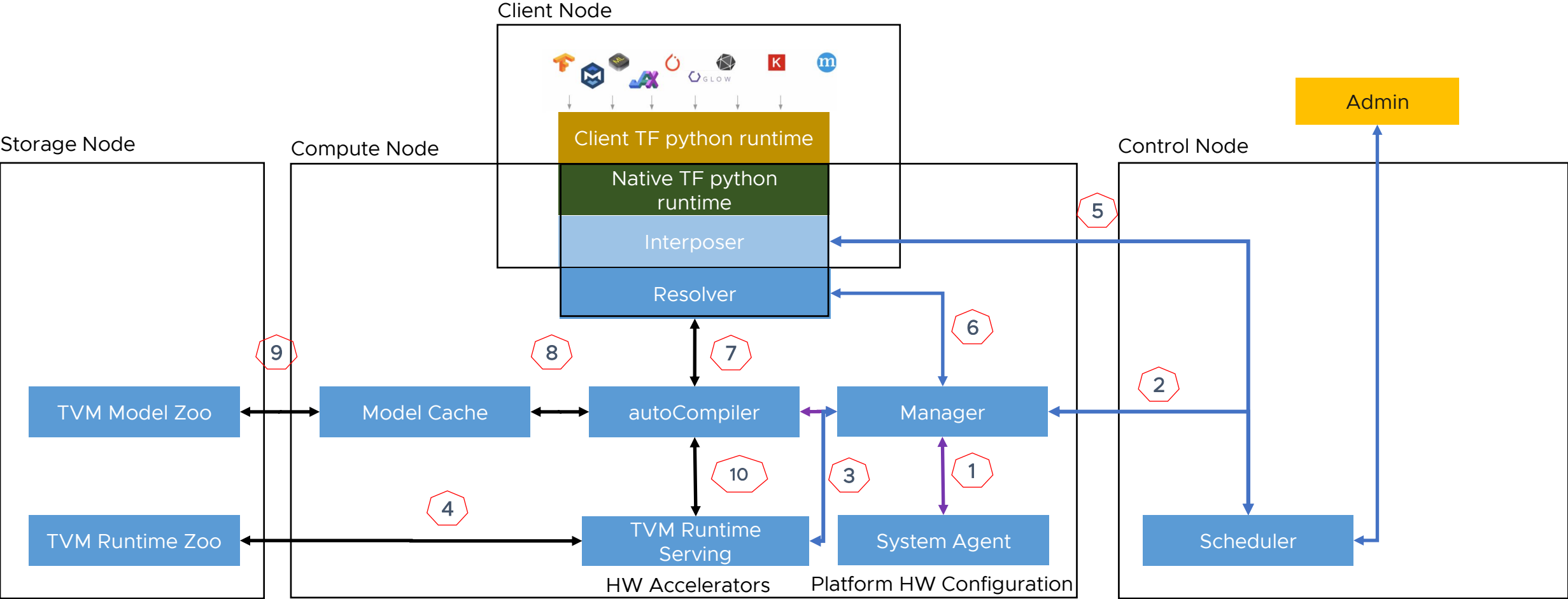    - *Only for inference now*

# Project MLInferBooster
*Solution*

- Focus
  - ML Inference
- Target
  - Ml Inference Acceleration by TVM
- Goal
  - Build a TVM Serving System
    - ❑ Backend
    - ❑ Automated
    - ❑ Unified server architecture
- How
  - Interpose ML framework python API
  - TVM progressing – Auto {detecting, compiling, scheduling, inferencing, etc}
    - ❑ HW accelerator type
    - ❑ Model & Mode info – {input, output} layer, shape, etc
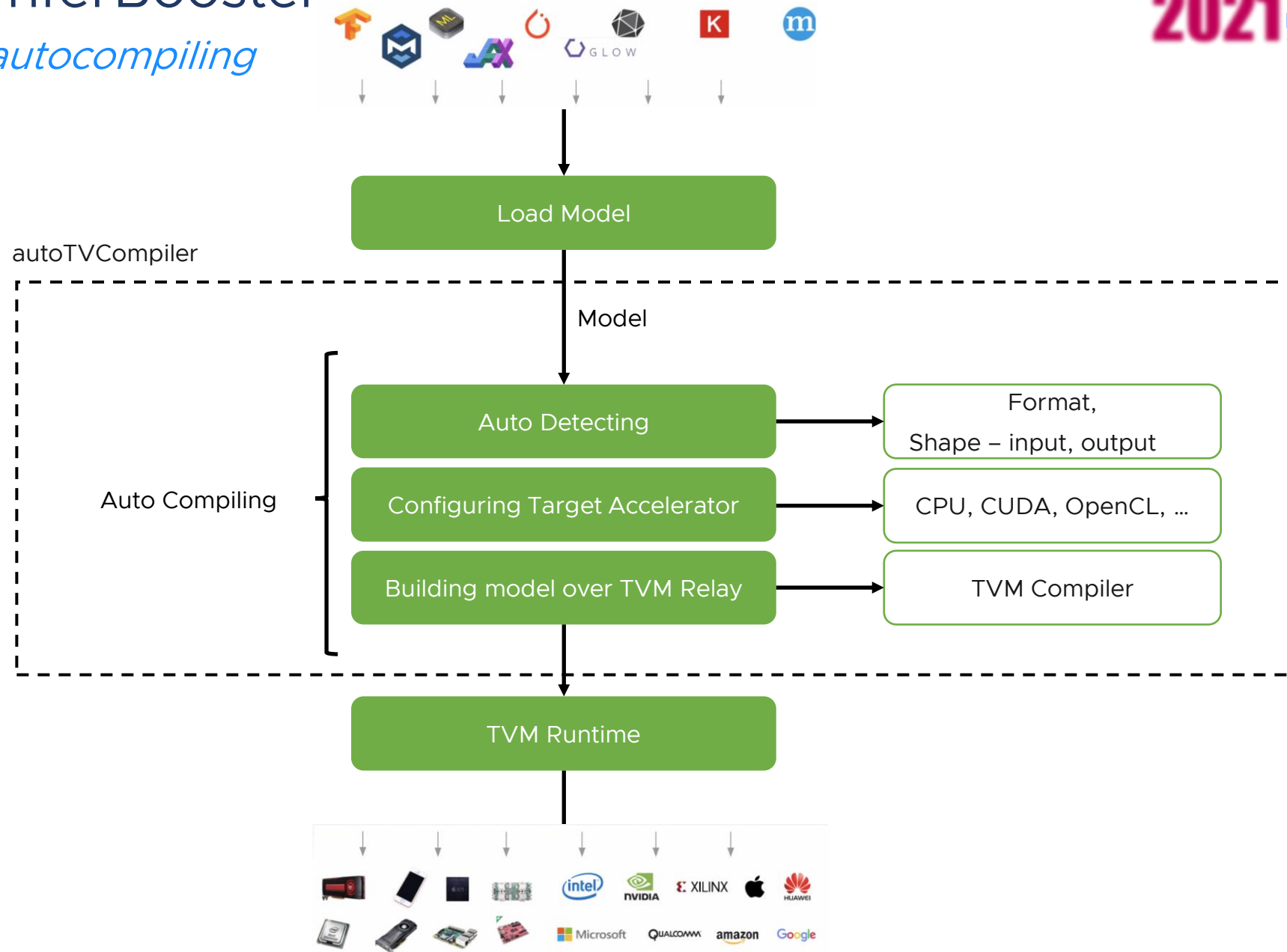    - ❑ TVM API : ML framework API

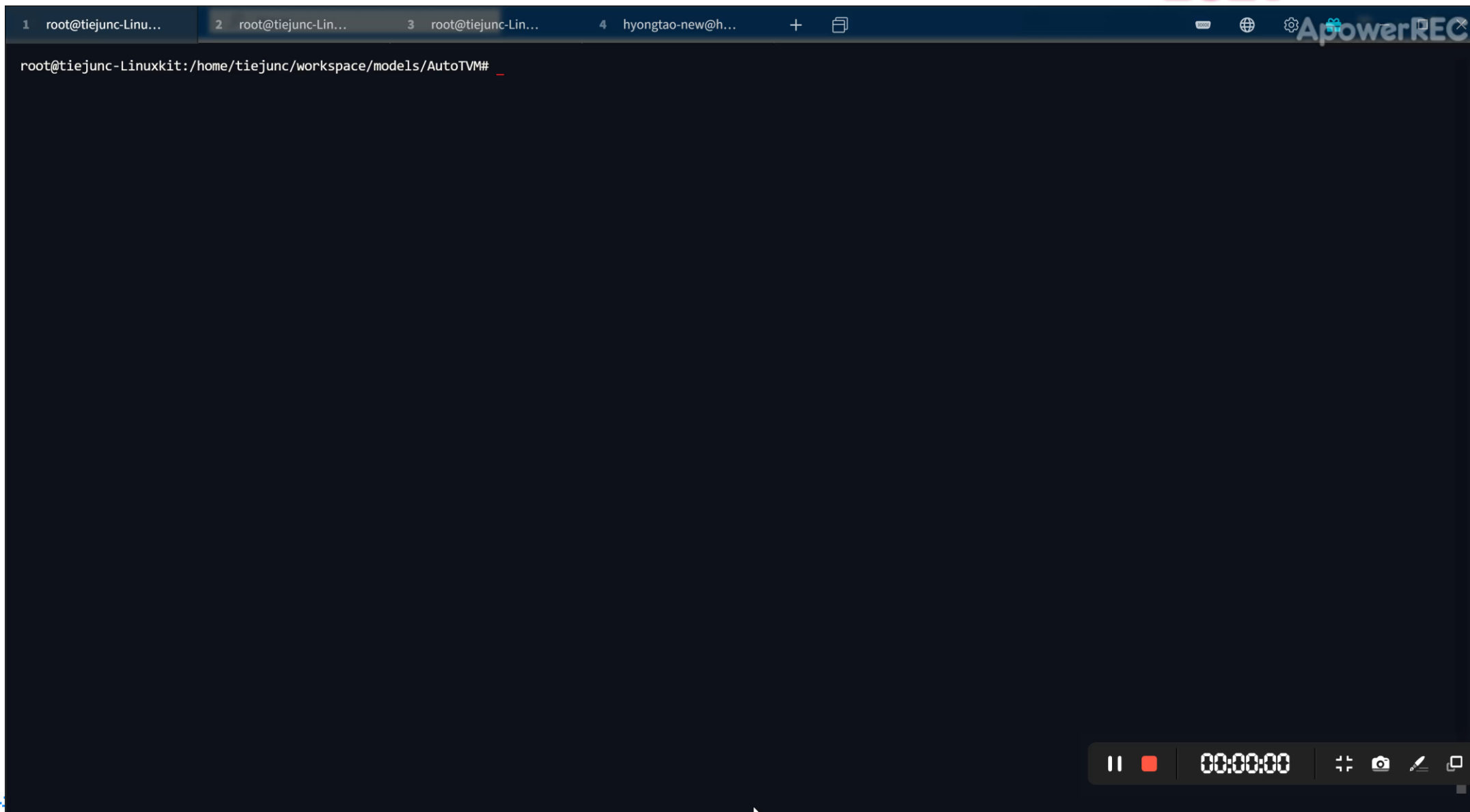# Project MLInferBooster
*Architecture Overview*

# Project MLInferBooster
## *Components - autocompiling*



autoTVCompiler

Load Model

Model

Auto Compiling

Auto Detecting → Format, Shape – input, output

Configuring Target Accelerator → CPU, CUDA, OpenCL, …

Building model over TVM Relay → TVM Compiler

TVM Runtime

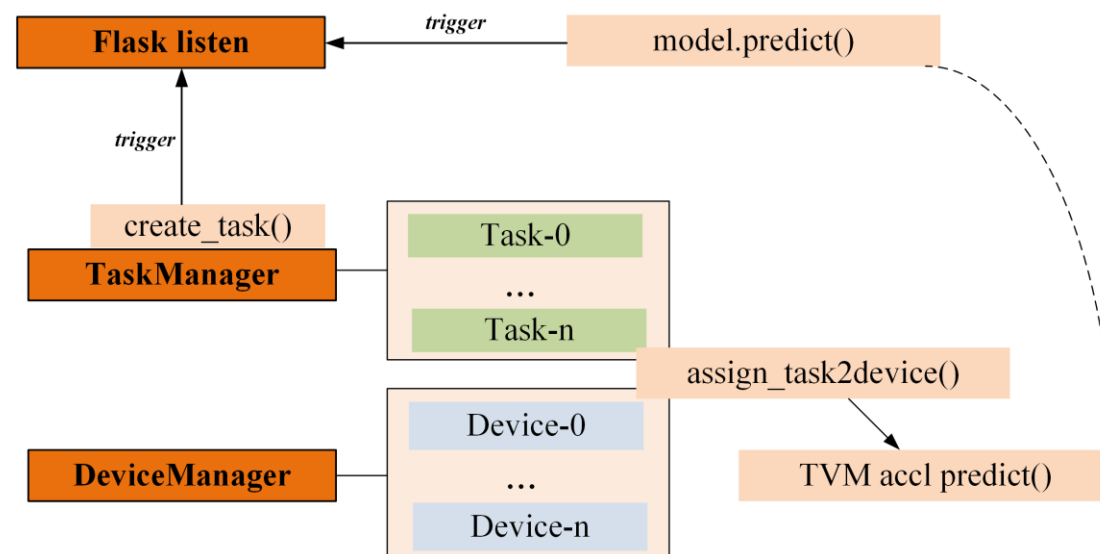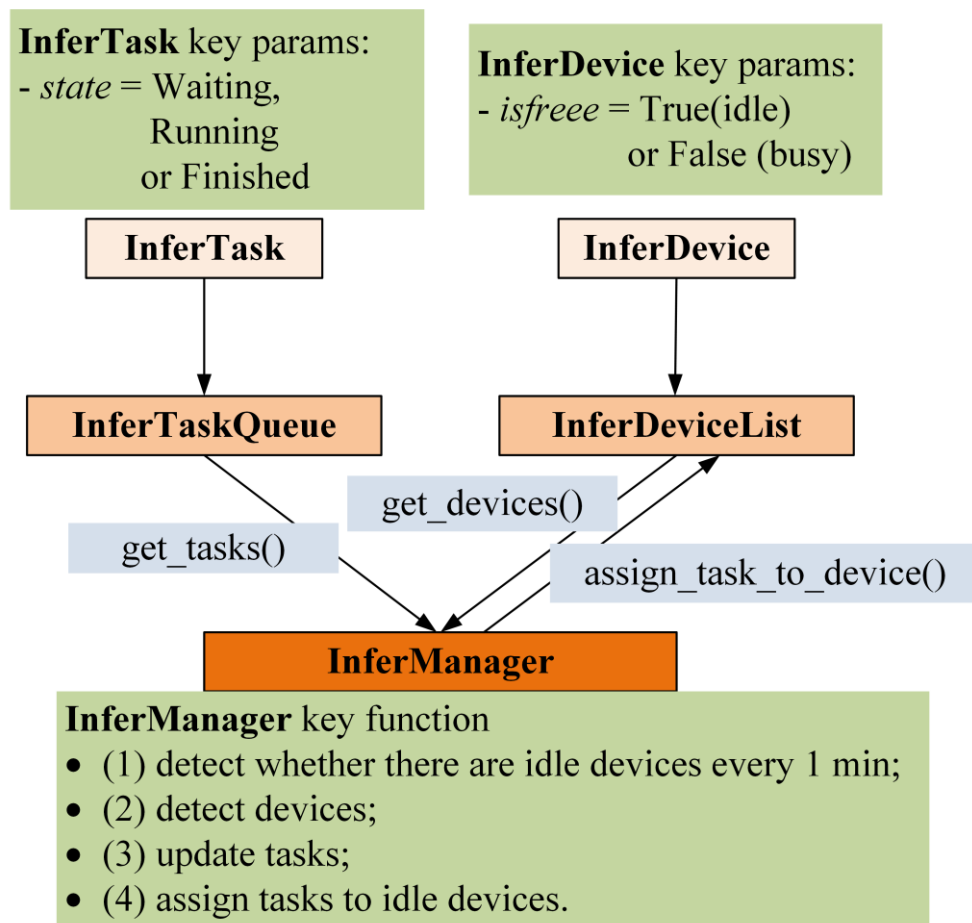# Project MLInferBooster

*Demo A1*

边缘计算分论坛

# Project MLInferBooster

边缘计算分论坛

# Project MLInferBooster
## Components - Scheduler

**InferTask** key params:
- *state* = Waiting,
        Running
        or Finished

**InferDevice** key params:
- *isfreee* = True(idle)
            or False (busy)

**InferTask**

**InferDevice**

**InferTaskQueue**

**InferDeviceList**

get_devices()

get_tasks()

assign_task_to_device()

**InferManager**

**InferManager** key function
- (1) detect whether there are idle devices every 1 min;
- (2) detect devices;
- (3) update tasks;
- (4) assign tasks to idle devices.

**Flask listen**

*trigger*

model.predict()

*trigger*

create_task()

**TaskManager**

Task-0
…
Task-n

**DeviceManager**

Device-0
…
Device-n

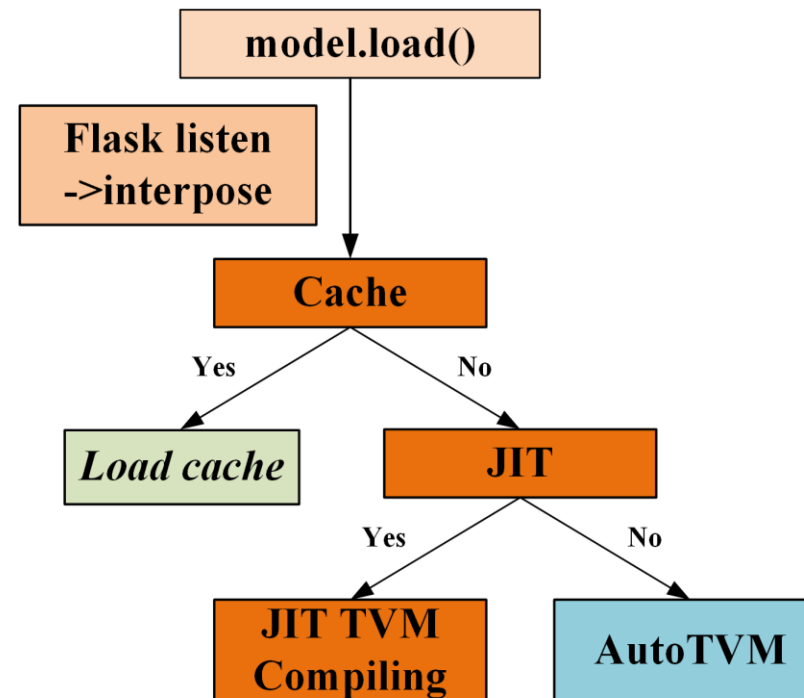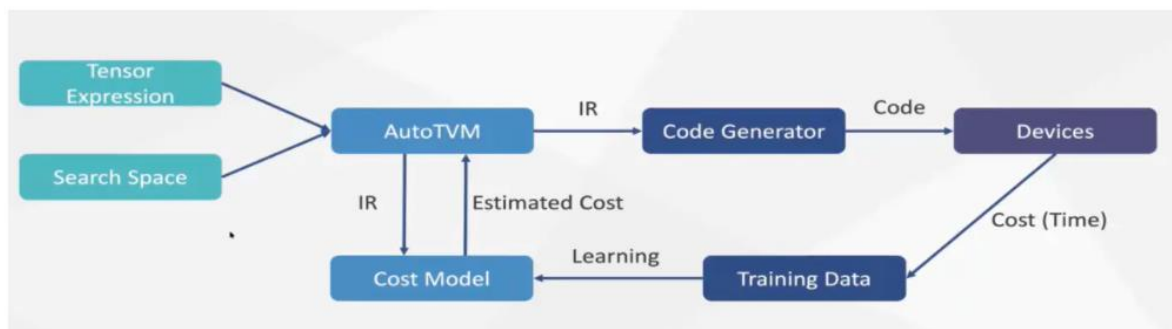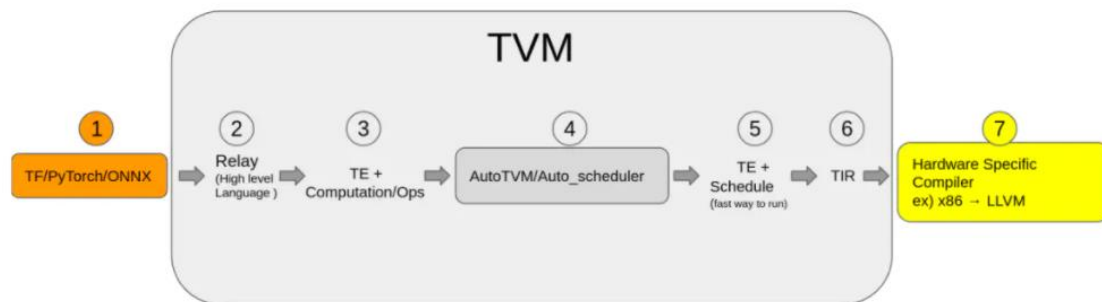assign_task2device()

TVM accl predict()

# Project MLInferBooster

*Demo Auto-Schedule*

# Project MLInferBooster
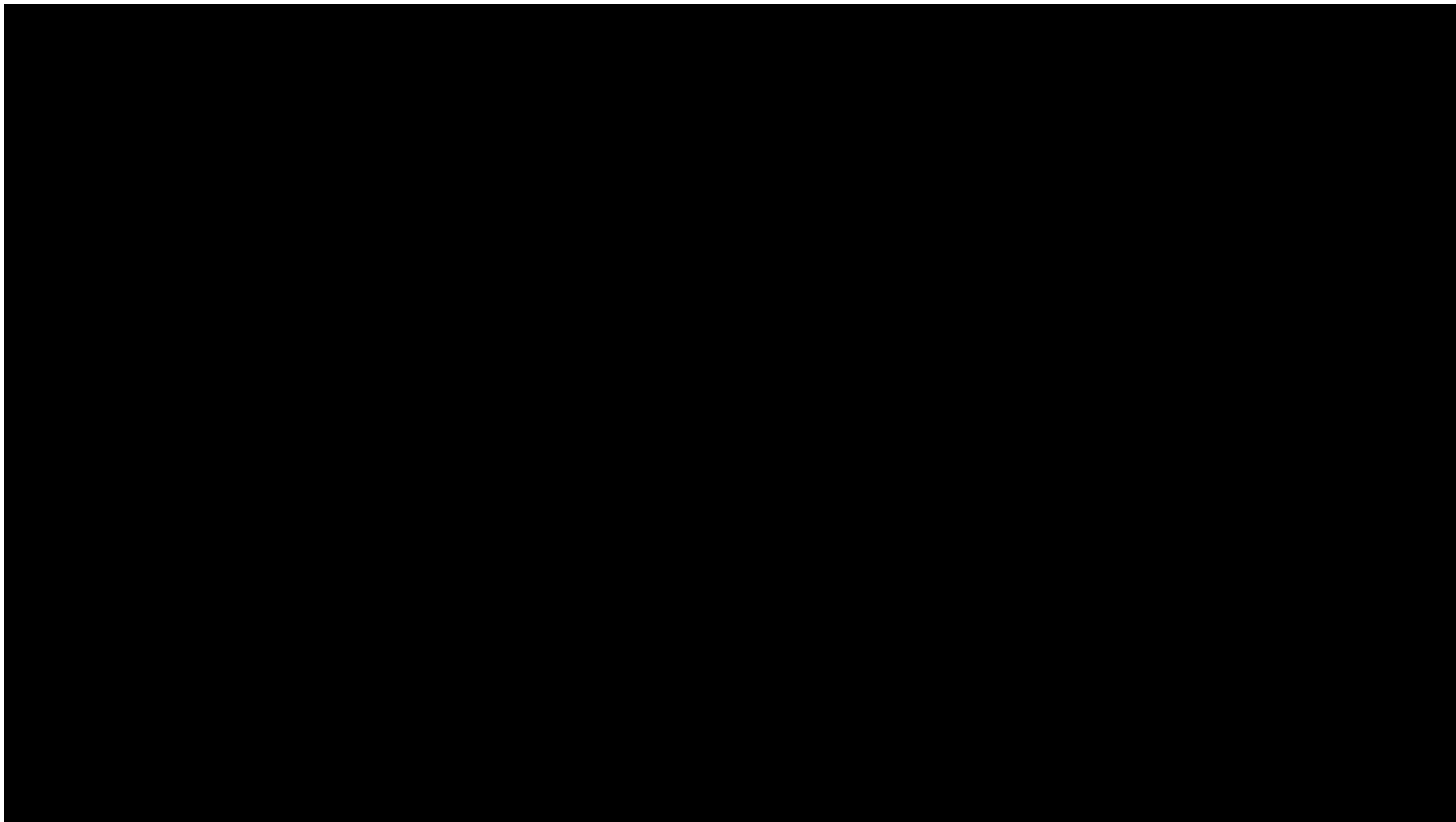## *Components - AutoTVM*

参考资料：
[1] 《Learning to Optimize Tensor Programs》NIPS-2018，陈天奇。
[2] https://zhuanlan.zhihu.com/p/37181530　陈天奇 知乎。

边缘计算分论坛

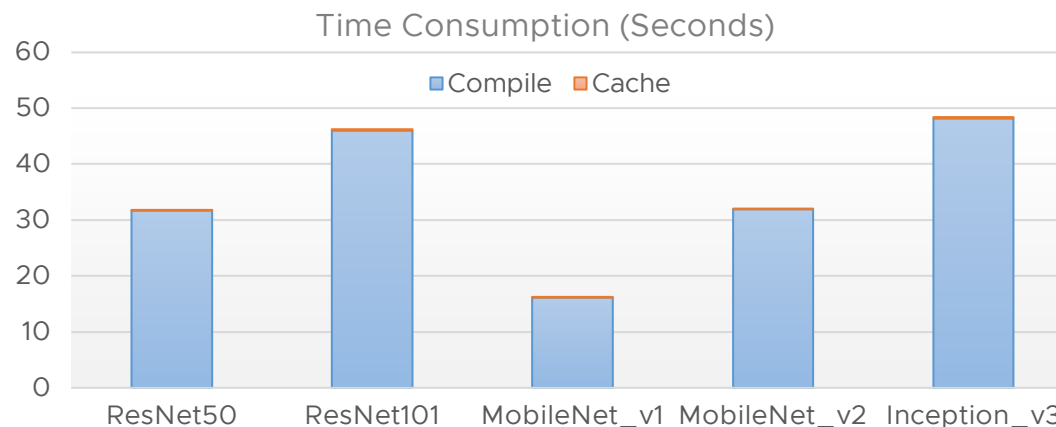# Project MLInferBooster

*Demo AutoTVM tuning*

边缘计算分论坛

# Project MLInferBooster

*Components - Model Cache*

- Objectives
  - Cache the compiled model information
  - Mapping mechanism
  - Least Frequently Used (LFU) cache replacement policy

- Benefits
  - Avoid recompile of the same model and save time
  - Apply efficient strategy to prevent cache overflow
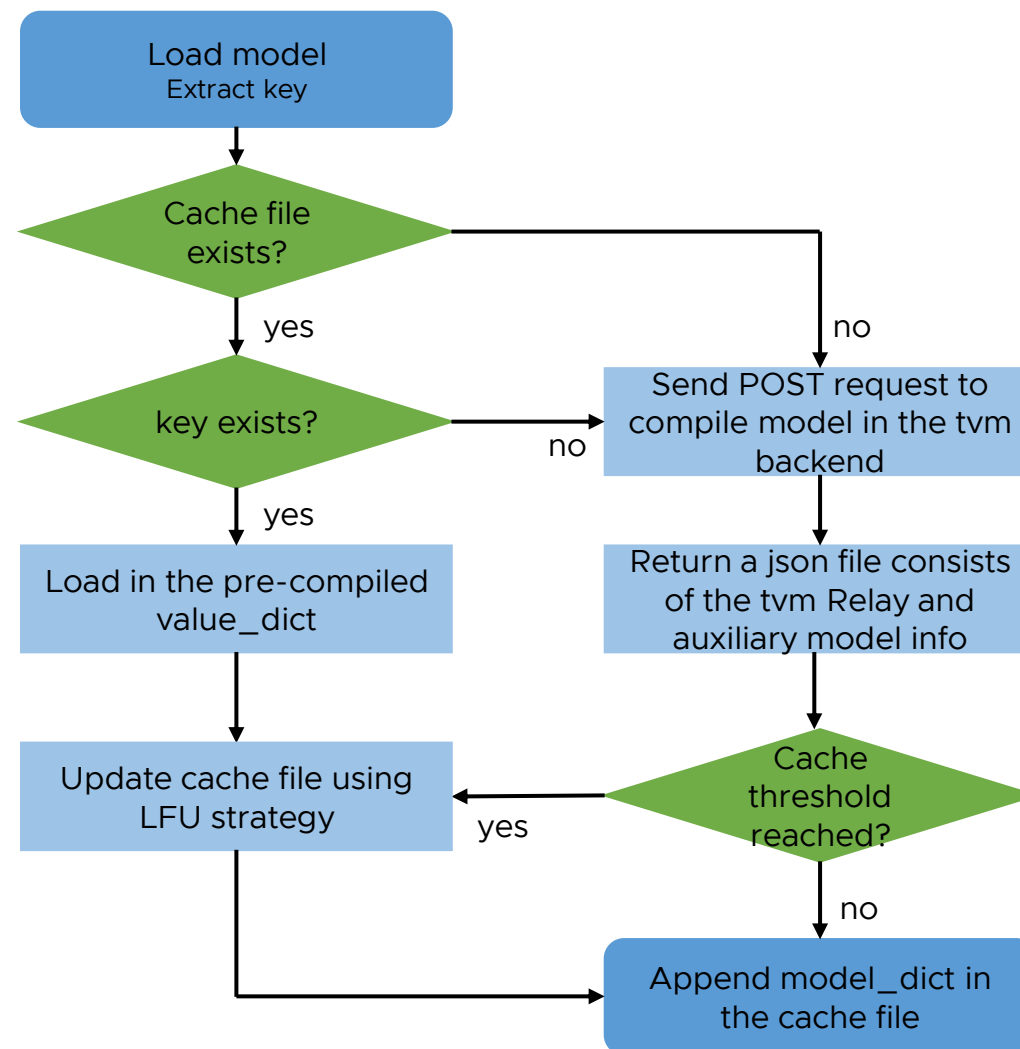


Time Consumption (Seconds)

边缘计算分论坛

# Project MLInferBooster

*Simple mapping mechanism & workflow*

Model_dict = {key : value_dict}

key : model_name#model_createtime

value_dict : { 'tvm_relay' : -----.so,
        'input_layer' : -----,
        'input_layer_dtype' : -----,
        'output_layer' : -----,
        'target' : -----,
        'freq' : 1}

*LFU cache replacement policy*

Model is already compiled
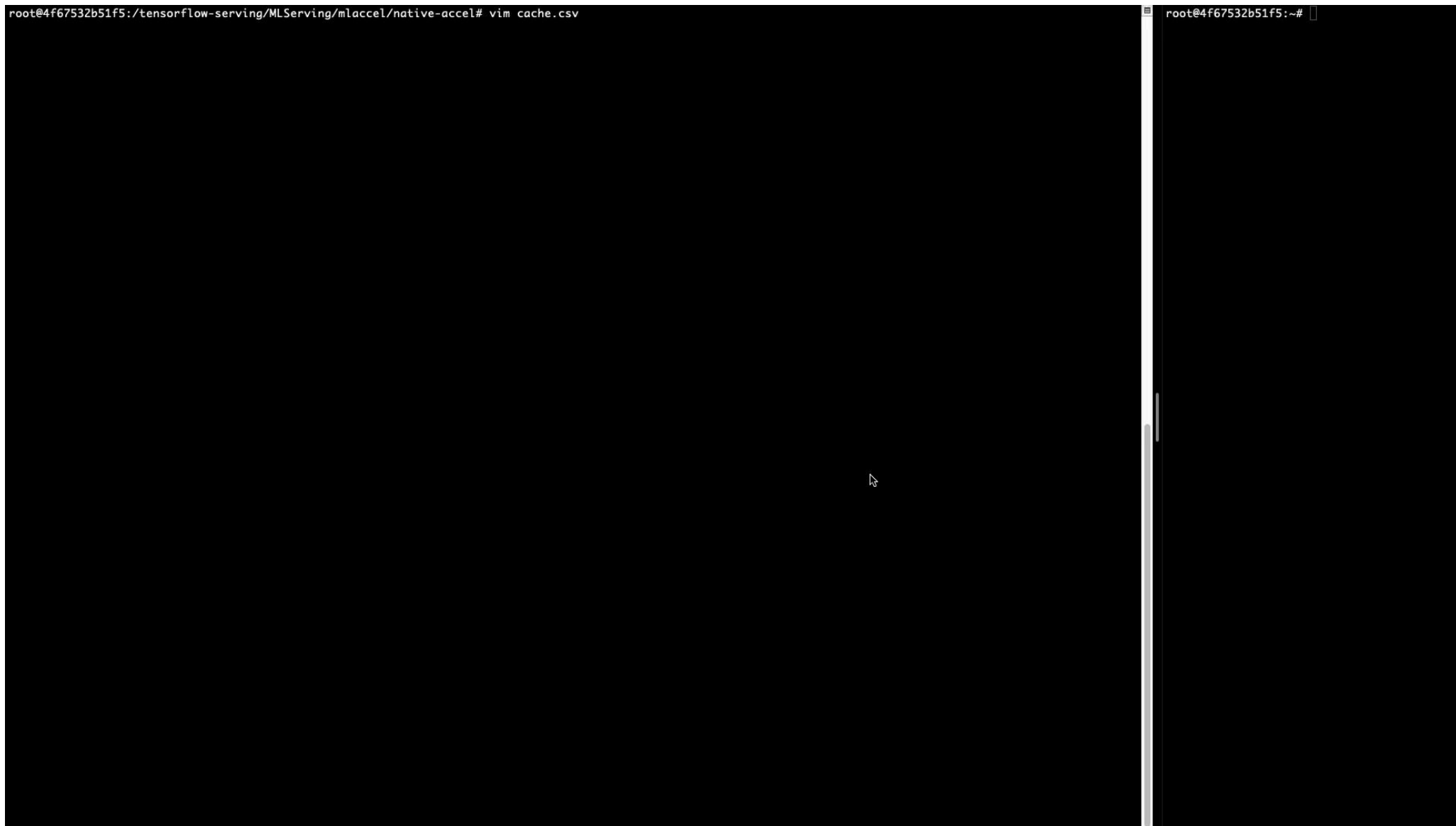
| key | value_dict | | | | | |
|-----|------|-------|--------|-------------|--------------|------------------|
| | freq | Relay | target | input_layer | output_layer | input_layer_dtype |
| model_1#date | 2 | -- | -- | -- | -- | -- |
| model_2#date | 1 | -- | -- | -- | -- | -- |
| model_3#date | 1 | -- | -- | -- | -- | -- |

**drop** ←

**append**    freq = freq +1

| key | value_dict | | | | | |
|-----|------|-------|--------|-------------|--------------|------------------|
| model_1#date | 3 | -- | -- | -- | -- | -- |

A new model comes in when cache threshold has been reached

| key | value_dict | | | | | |
|-----|------|-------|--------|-------------|--------------|------------------|
| | freq | Relay | target | input_layer | output_layer | input_layer_dtype |
| model_2#date | 1 | -- | -- | -- | -- | -- |
| model_3#date | 1 | -- | -- | -- | -- | -- |
| model_1#date | 3 | -- | -- | -- | -- | -- |

**popitem (FIFO)** ←

**append**

| key | value_dict | | | | | |
|-----|------|-------|--------|-------------|--------------|------------------|
| model_4#date | 1 | -- | -- | -- | -- | -- |

# Project MLInferBooster
*Demo Model Cache*



`root@4f67532b51f5:/tensorflow-serving/MLServing/mlaccel/native-accel# vim cache.csv`

`root@4f67532b51f5:~#`

# Project MLInferBooster

*Next*

- More ML upstream frameworks

- More ML serving system

- More ML models supported
  - ❏ Model conversion (TensorFlow/PyTorch → ONNX)

- K8s integration
  - KFServing