**Alibaba A.I. Labs**
阿里巴巴人工智能实验室

# AILabs & TVM

CONTENT

**Alibaba A.I. Labs**
阿里巴巴人工智能实验室

# ARM 32 CPU

# Result

CPU：MTK8167S（ARM32 A35 1.5GHz）
Model：MobileNetV2_1.0_224



- ■ TF Lite 8bit
- ■ NCNN 8bit
- ■ QNNPACK 8bit
- ■ MNN 8bit
- ■ TVM Overflow-aware
- ■ ACE Overflow-aware (Assembly)
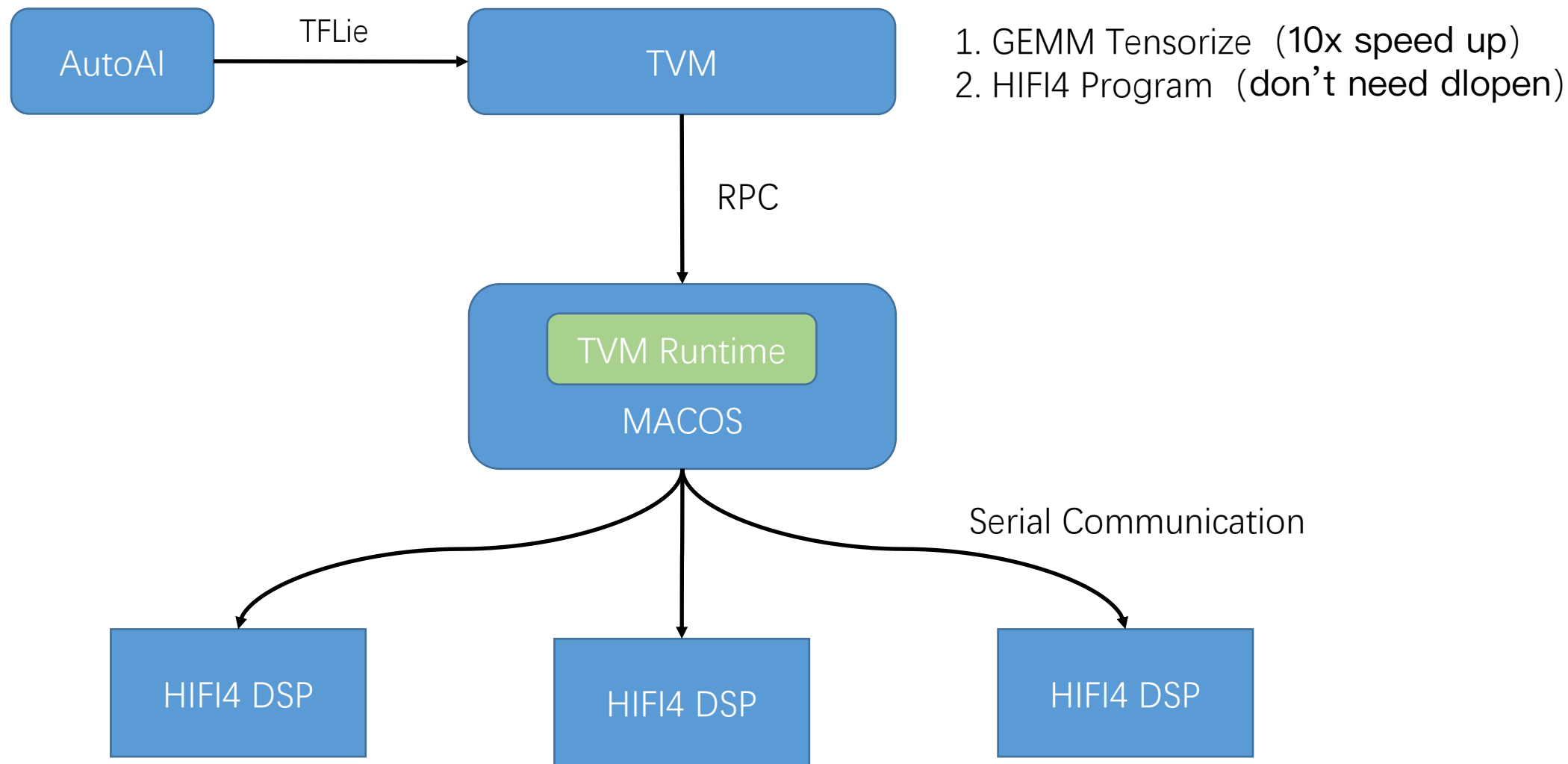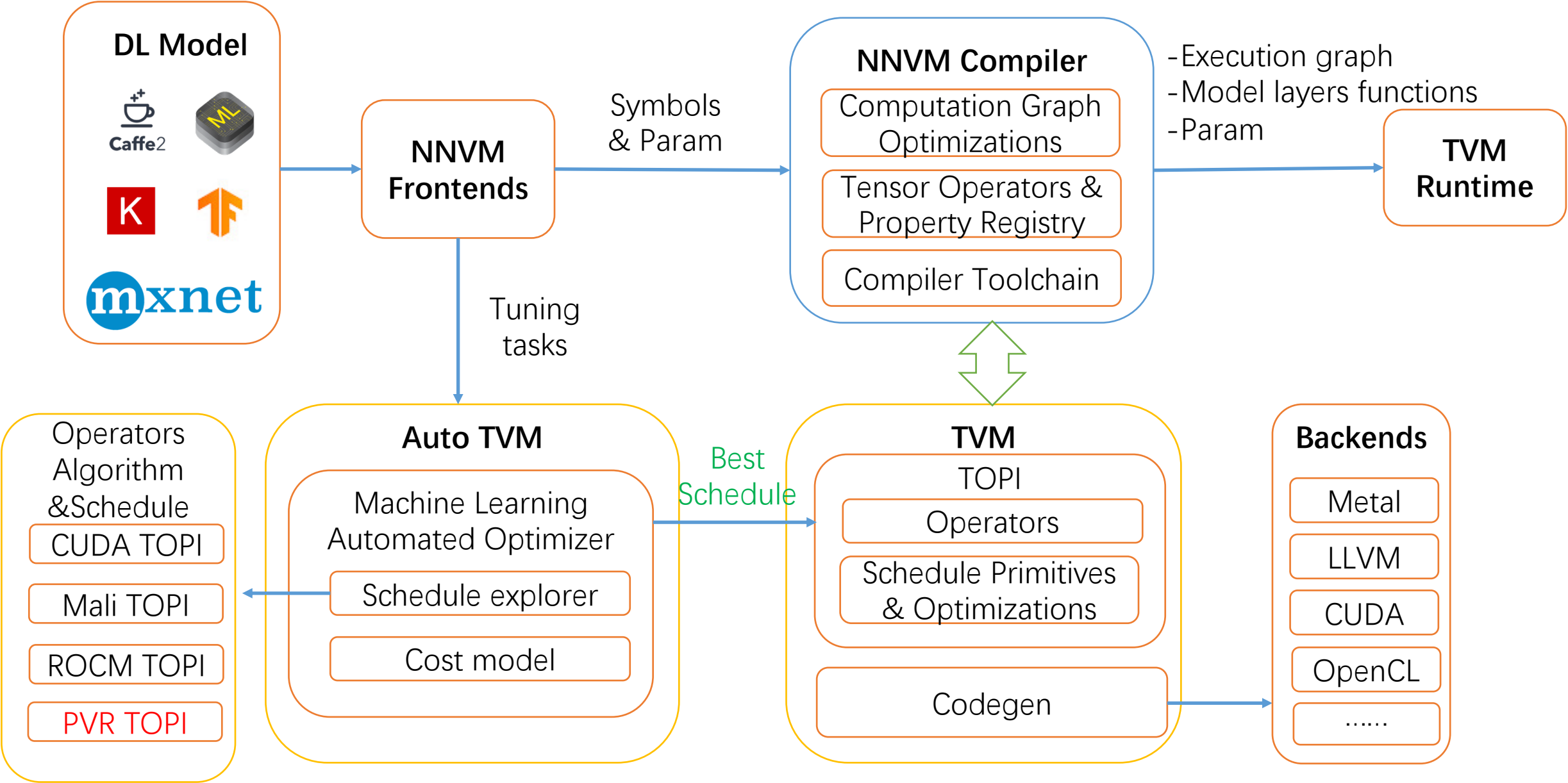
**Alibaba A.I. Labs**
阿里巴巴人工智能实验室

# HIFI 4

**Alibaba A.I. Labs**
阿里巴巴人工智能实验室

# PowerVR GPU

➢ TOPI for PVR, including what you want to compute and how to compute.

✓ What you want to compute
@autotvm.register_topi_compute(conv2d, 'pvr', ['direct'])
def conv2d_pvr(cfg, data, kernel, strides, padding, dilation, layout, out_dtype):
    #Describe algorithm with tensor expression language;
    #Return the out operation

✓ How to compute.
@autotvm.register_topi_schedule(schedule_conv2d_nchw, 'pvr', ['direct'])
def schedule_conv2d_nchw_pvr(cfg, outs):
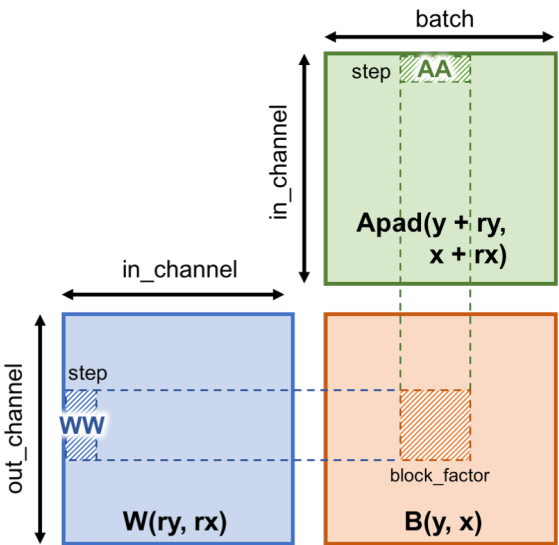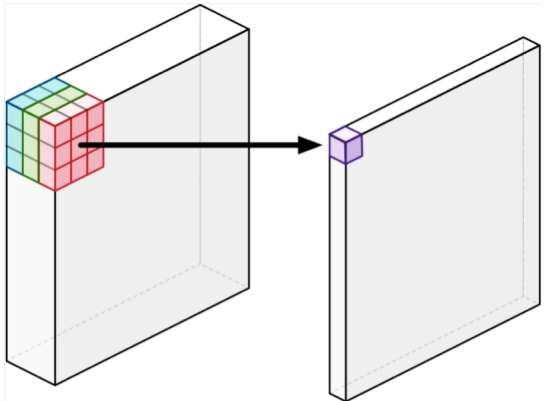    #Describe how to compute output by primitive

**Blocking**
splits the workload into thread
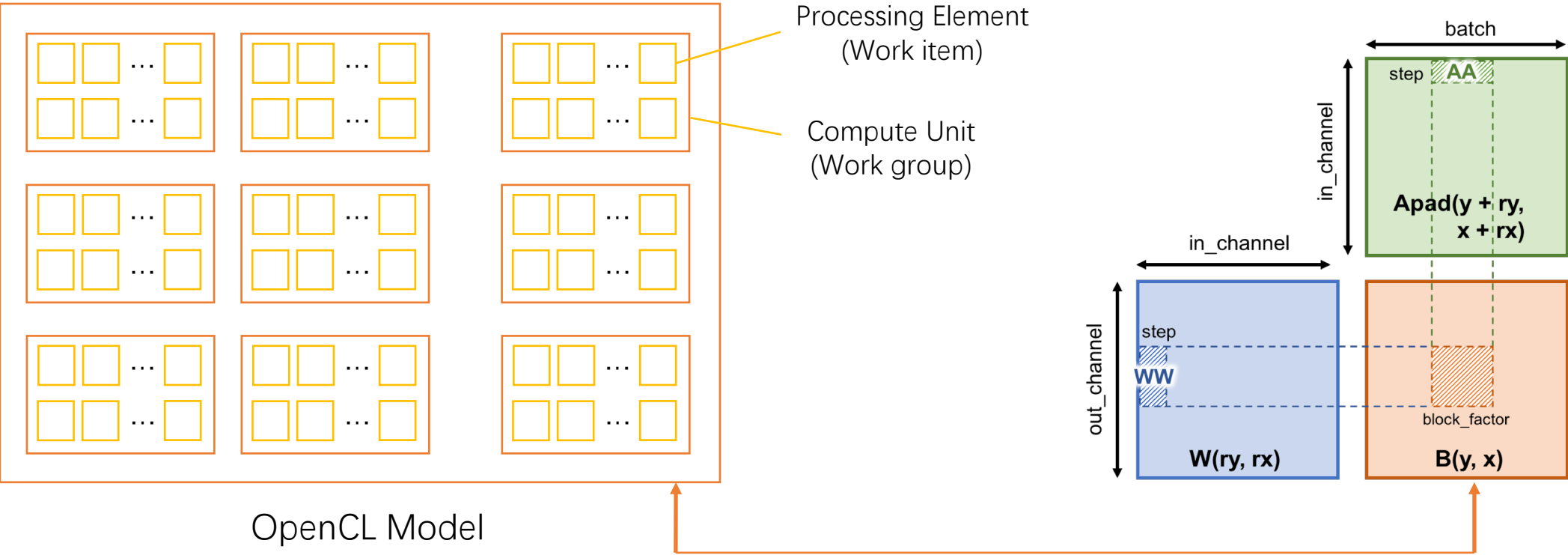blocks and individual threads

**Cooperative Fetching**
lets threads in the same thread block cooperatively
fetch dependent data



convolution



https://docs.tvm.ai/

## Blocking

Splits the workload into thread blocks (work groups) and individual threads (work items)



Processing Element
(Work item)

Compute Unit
(Work group)

OpenCL Model

## Cooperative Fetching

Lets threads (work item) in the same thread block (work group) cooperatively fetch dependent data

https://www.khronos.org/registry/OpenCL/specs/opencl-1.2.pdf

Alibaba A.I. Labs
阿里巴巴人工智能实验室

Thanks