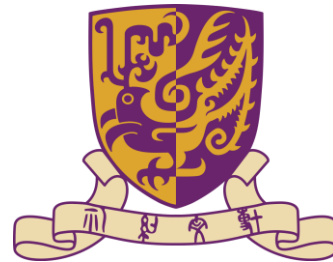


BlastNet: Exploiting Duo-Blocks for Cross-Processor Real-Time DNN Inference

Neiwen Ling¹, Xuan Huang¹, Zhihe Zhao¹, Guan Nan², Zhenyu Yan¹, Guoliang Xing¹

¹The Chinese University of Hong Kong,

²City University of Hong Kong



Real-time Deep Learning on the Edge

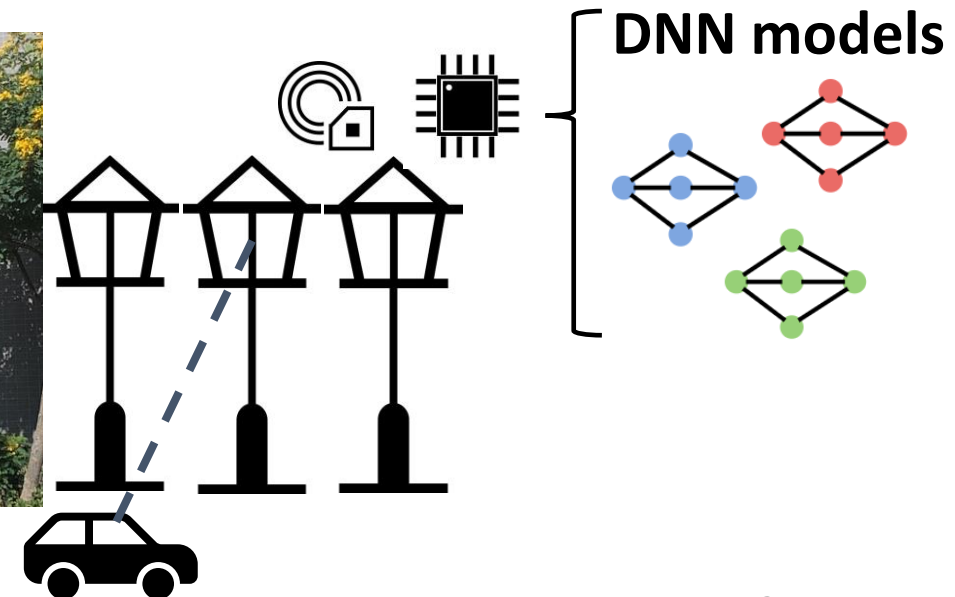
Autonomous driving
(e.g., Tesla)



Image source: tukuppt

Image source: Bernard Marr

Smart roadside infrastructure
(CUHK smart lampposts)



Heterogeneous Processors on Edge Platforms



Image source: EDN Taiwan

Edge Platform

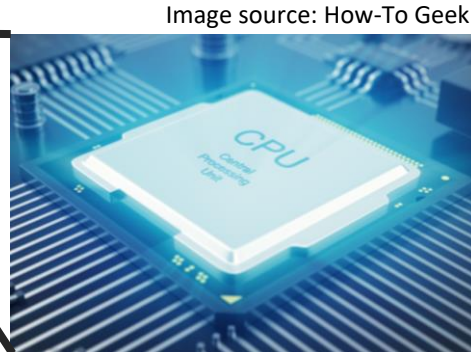


Image source: How-To Geek

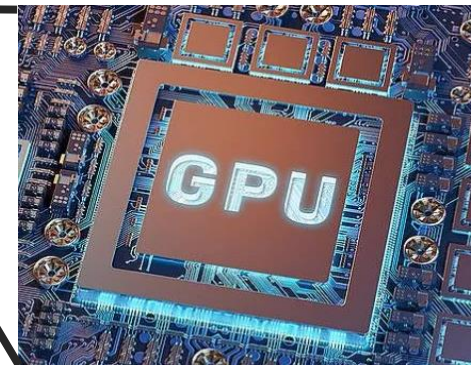
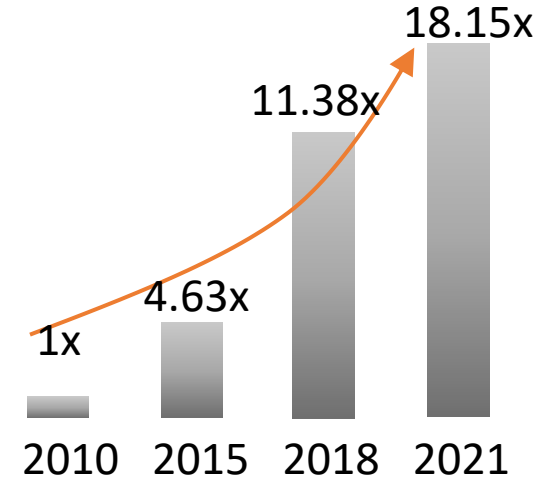
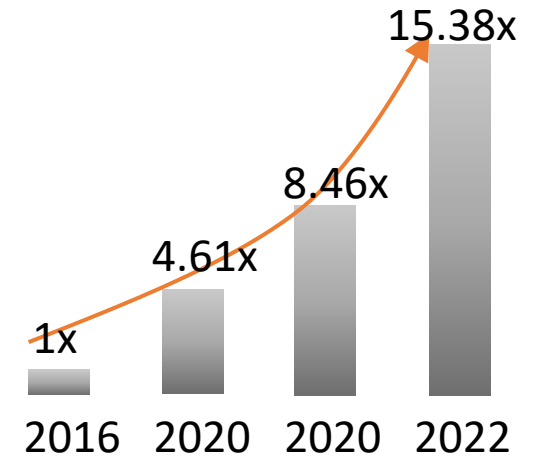


Image source: OFweek

CPU

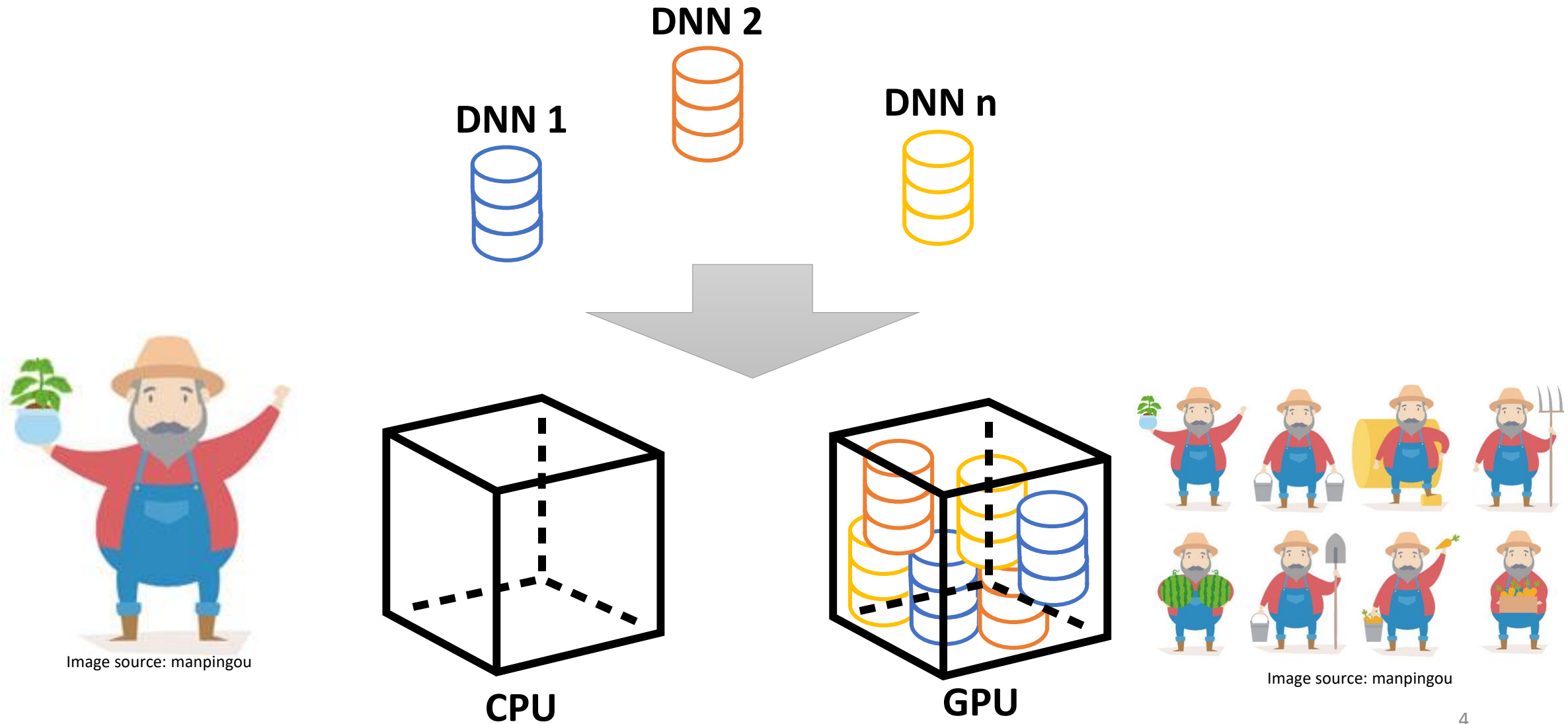


GPU



Compute Power

Untapped Compute Power



How to utilize both edge **CPU** and **GPU** for running **multiple** DNNs efficiently?

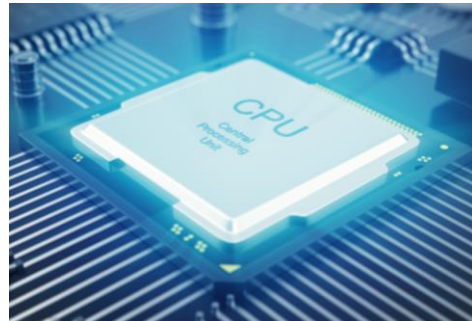


Image source: How-To Geek

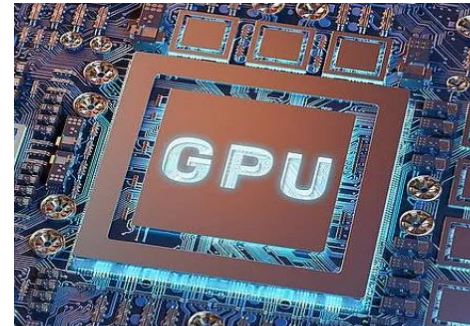
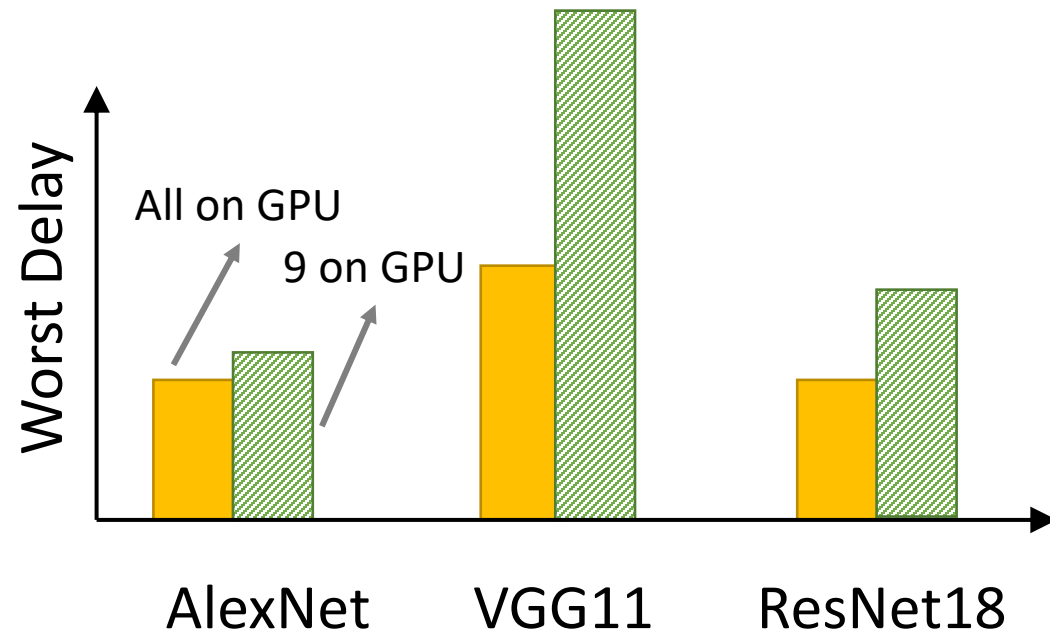
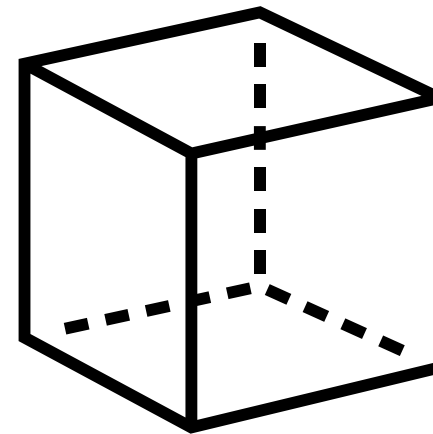


Image source: OFweek

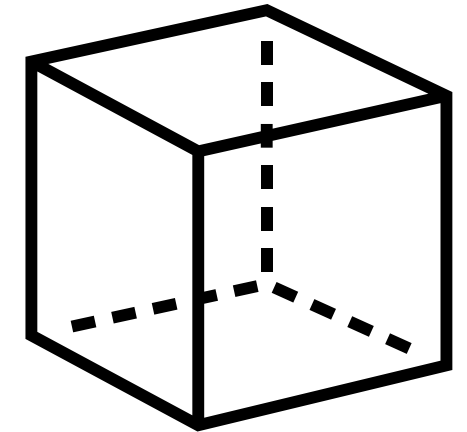
Offloading models?



DNN models



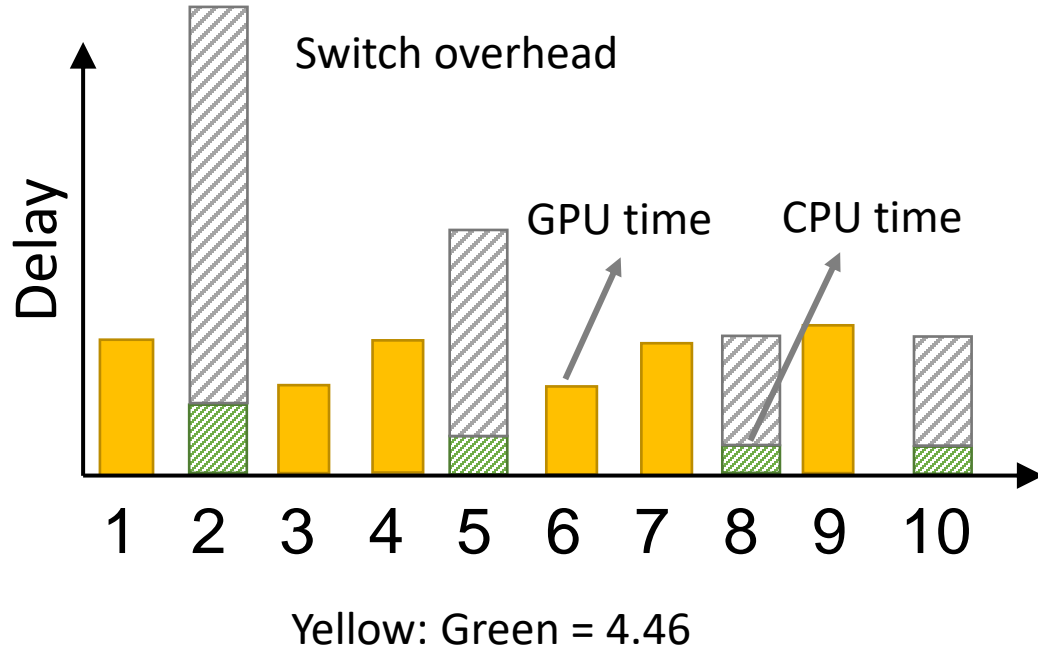
CPU



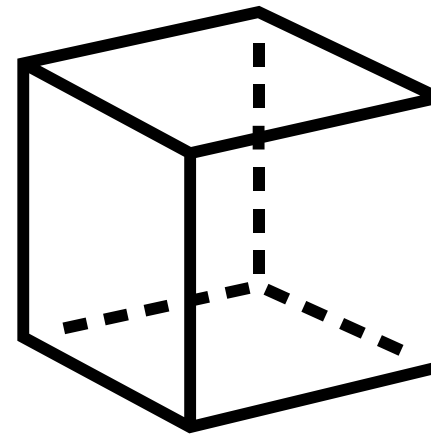
GPU

Too coarse for efficient resource utilization

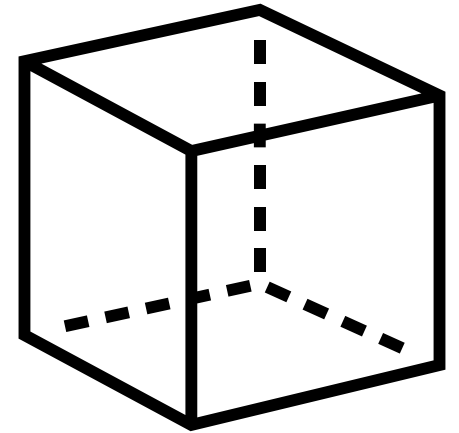
Offloading layers?



DNN layers



CPU



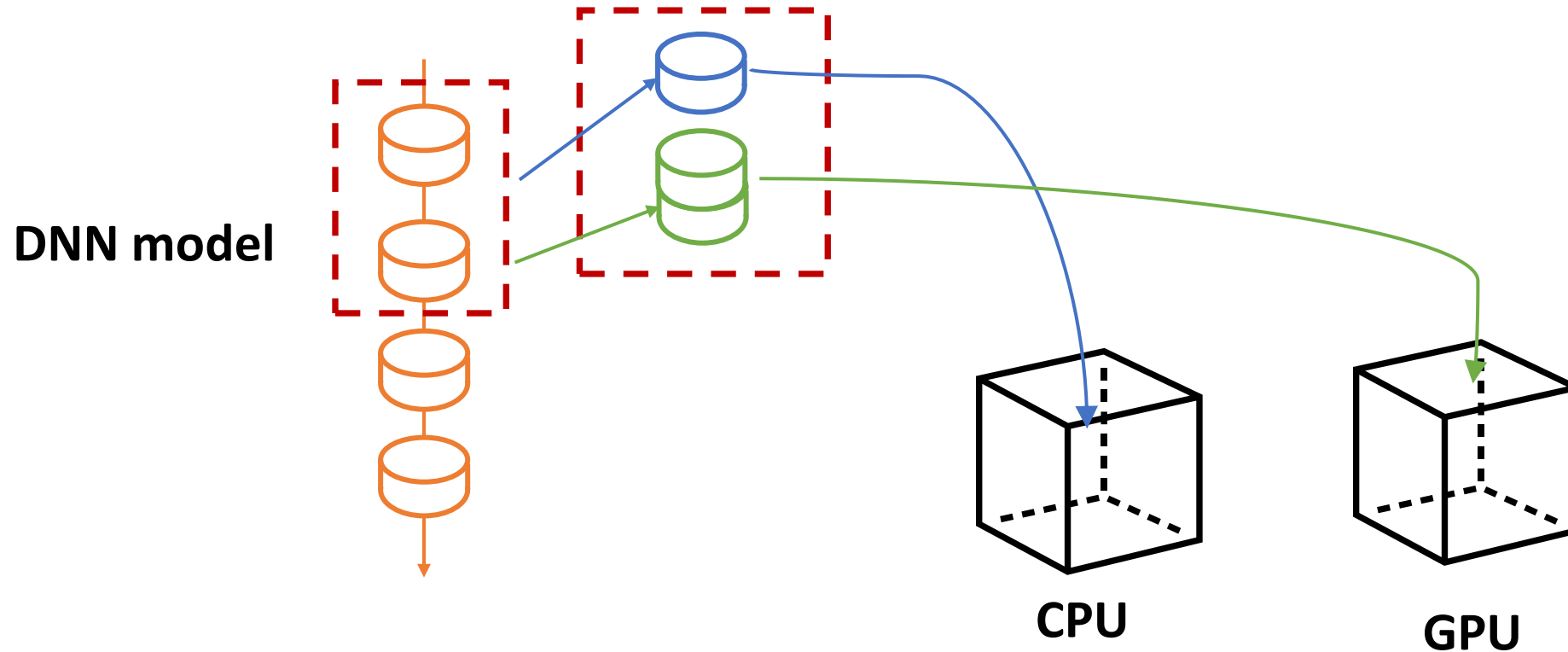
GPU

Significant idle time and switching delay

BlastNet: Block-Level model optimization and Scheduling system

Duo-Block

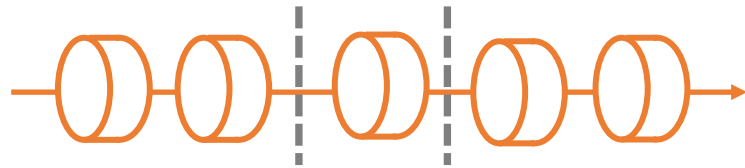
a new model
inference abstraction



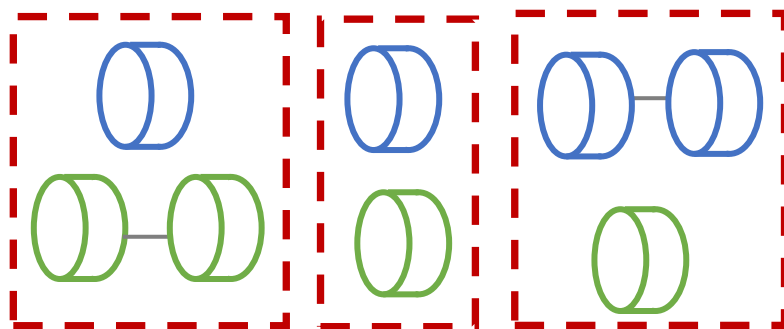
BlastNet: Block-Level model optimization and Scheduling system

- Duo-block generation

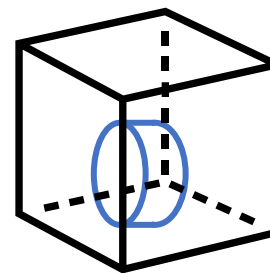
DNN model



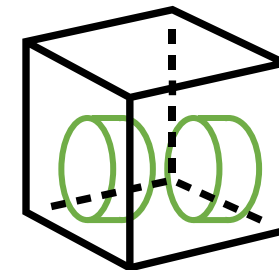
Duo-block1 2 3



- Dynamic cross-processor DNN scheduling



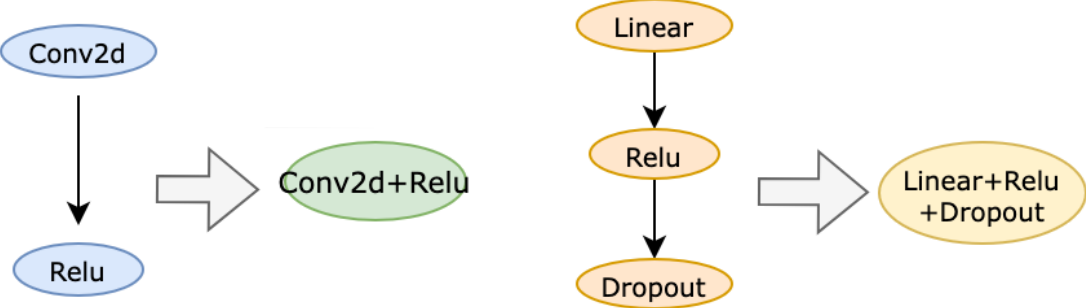
CPU



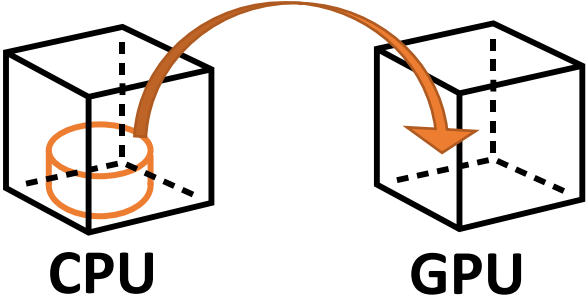
GPU

Duo-block Generation

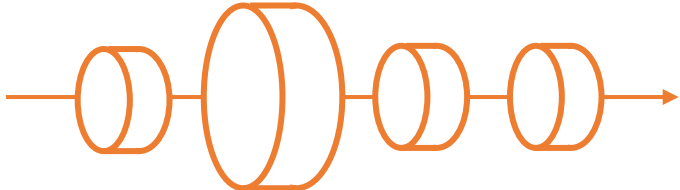
1. Preserve operator fusion



2. Switching < Execution?

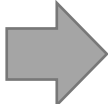
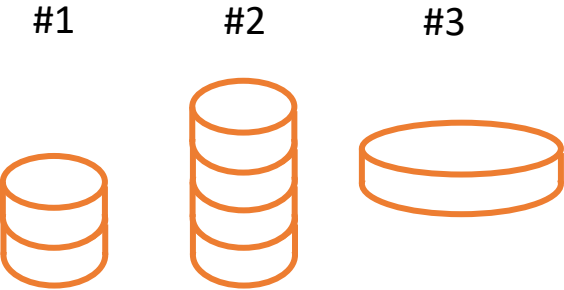


3. Optimizing bottleneck blocks



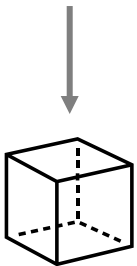
Duo-block Generation with Neural Architecture Search

1. Candidate blocks generation with processor-friendly operators



2. Optimize search space based on profiling

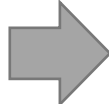
Candidate Block



Processor



Time

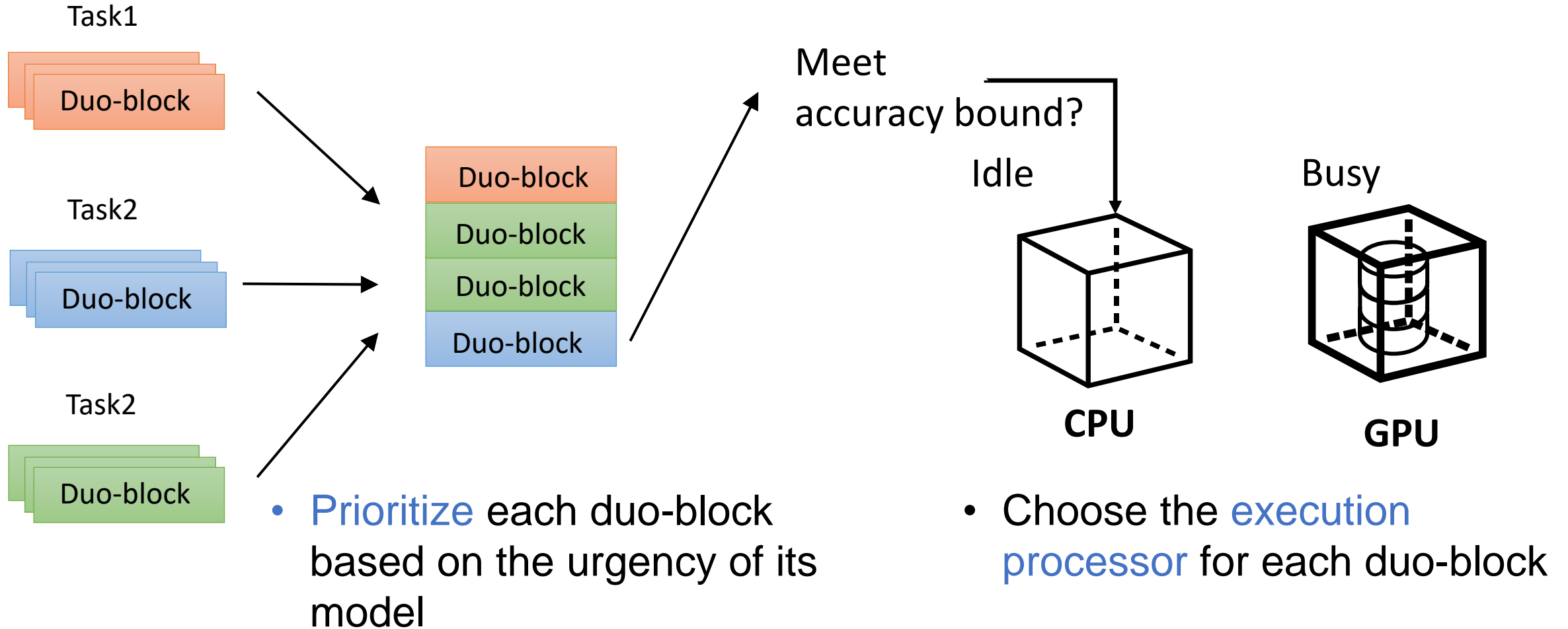


3. Neural Architecture Search



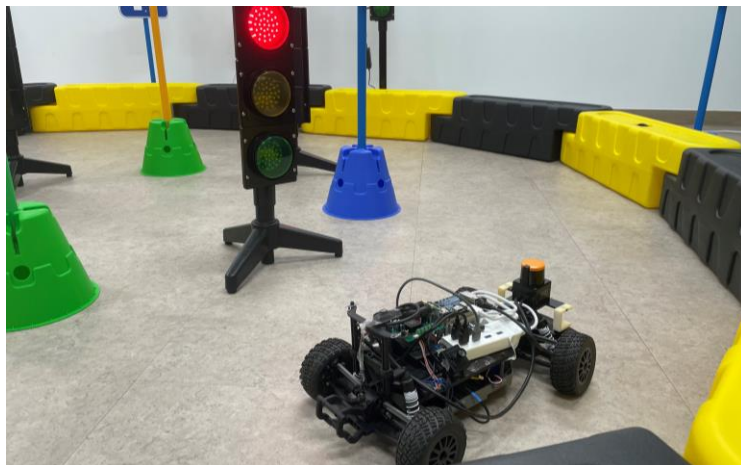
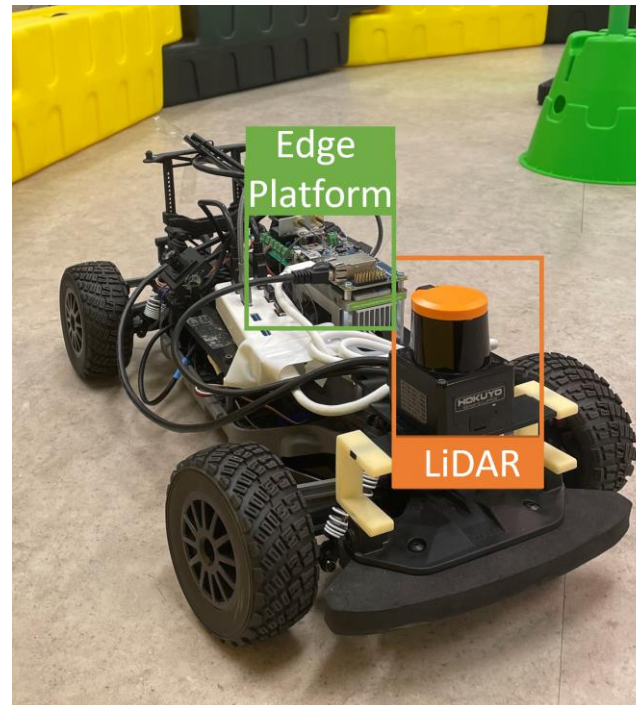
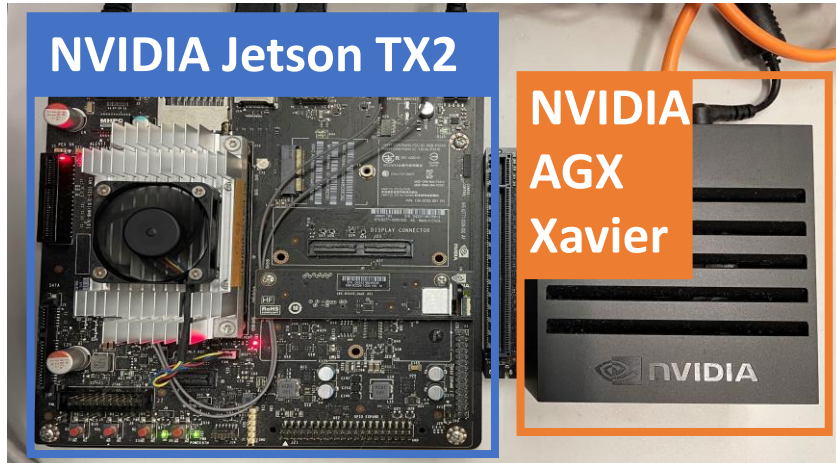
Accuracy

Dynamic Cross-Processor Scheduling



Experiment Setup

- 3 edge platforms including an autonomous driving car



- 3 types of DL tasks
- 5 DNN models
- 3 datasets

Image Classification



Image source: Hertz

Sign Recognition



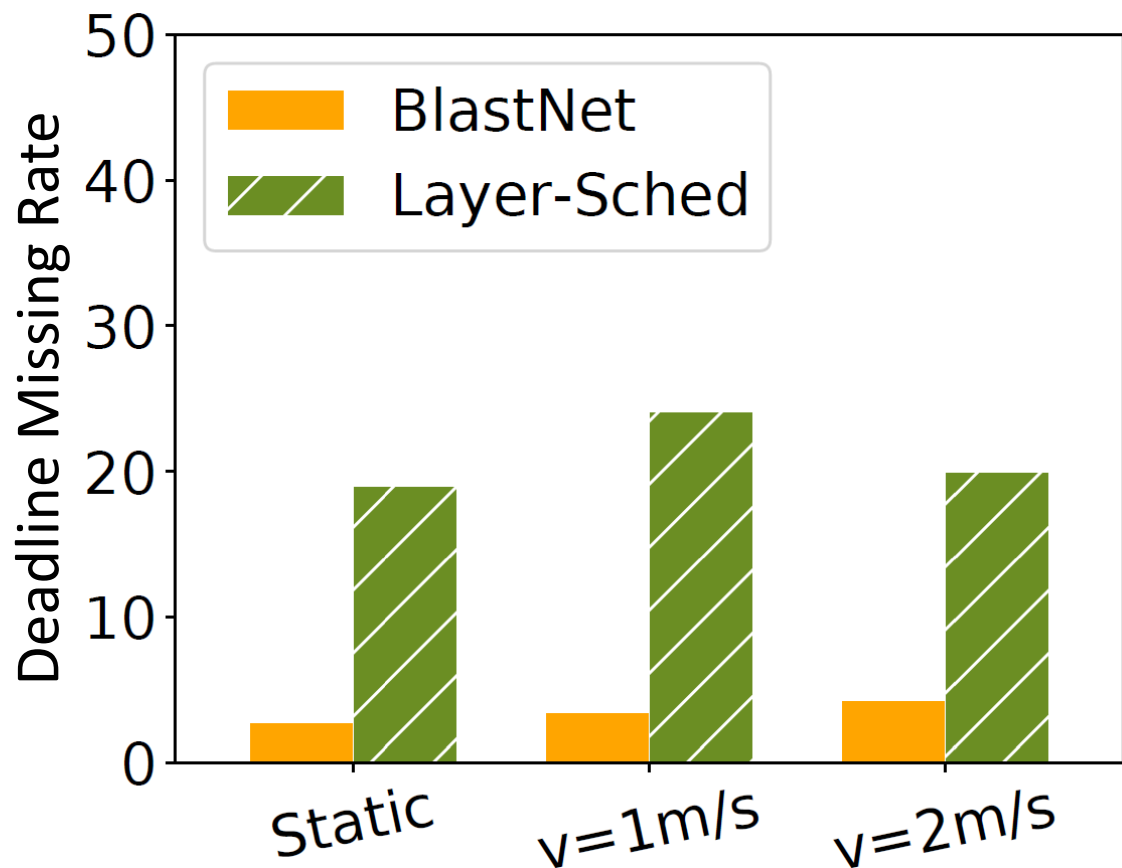
Image source: Datangxs

Object Detection

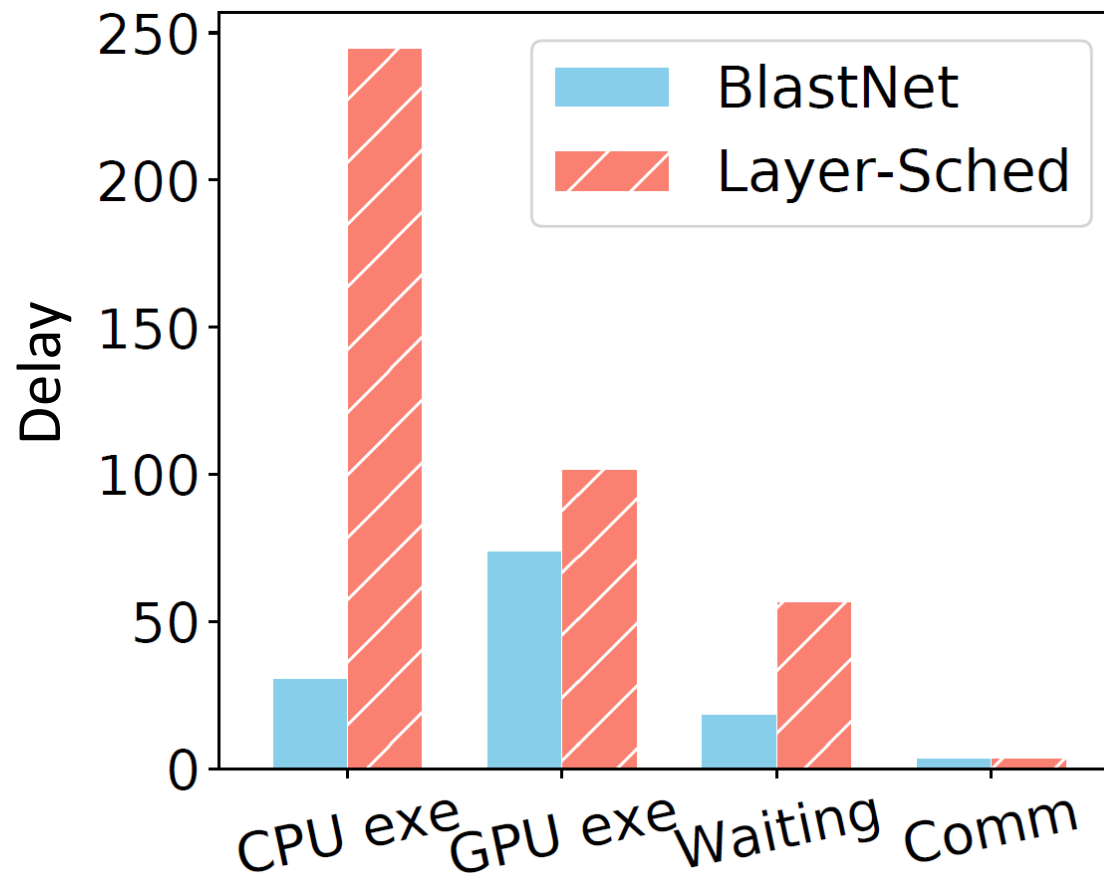


Image source: Innovtrics

Experimental results: End-to-end system evaluation



< 5% deadline missing rate



Reduced CPU/GPU execution and waiting time

Conclusion & Future Work

□ **BlastNet:** novel block-level model optimization and cross-processor scheduling

- Duo-block generation
- A dynamic block-level DNN scheduler

□ **Future Work**

- Scalability to new heterogeneous platforms
- Uncertain workloads

Thanks for Listening!

- ❖ *Check our paper:* BlastNet: Exploiting Duo-Blocks for Cross-Processor Real-Time DNN Inference
- ❖ *Authors:* Neiwen Ling, Xuan Huang, Zhihe Zhao, Guan Nan, Zhenyu Yan, Guoliang Xing
- ❖ *Website:* <http://aiot.ie.cuhk.edu.hk>

