



香港中文大學  
The Chinese University of Hong Kong

# Deep Neural Network Design Automation

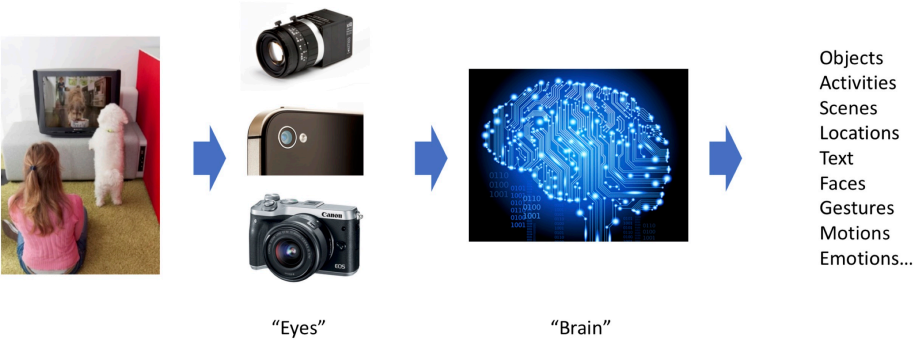
**Bei Yu**

Department of Computer Science and Engineering  
The Chinese University of Hong Kong

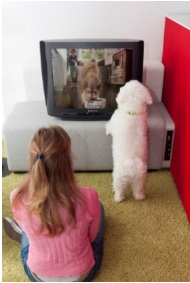


# Computer Vision

- ▶ Humans use their **eyes** and their **brains** to visually sense the world.
- ▶ Computers use their **cameras** and **computation** to visually sense the world

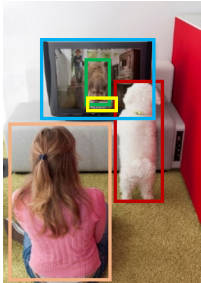


# Few More Core Problems



Classification

Image



Detection

Region



Segmentation

Pixel



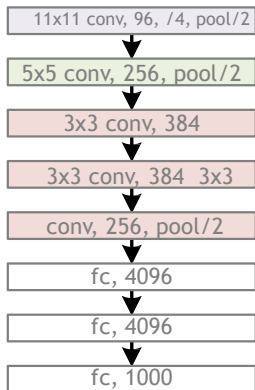
Sequence

Video



# Revolution of Depth

AlexNet, 8  
layers  
(ILSVRC 2012)



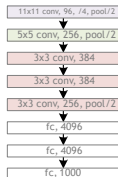
Slide Credit: He et al. (MSRA)



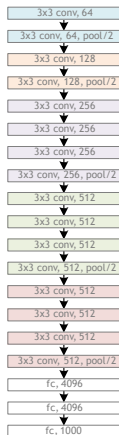


# Revolution of Depth

AlexNet, 8  
layers  
(ILSVRC 2012)



VGG, 19  
layers  
(ILSVRC  
2014)



GoogleNet, 22  
layers  
(ILSVRC 2014)



Slide Credit: He et al. (MSRA)



# Revolution of Depth

AlexNet, 8  
layers  
(ILSVRC 2012)



VGG, 19  
layers  
(ILSVRC  
2014)



ResNet, 152  
layers  
(ILSVRC 2015)



Slide Credit: He et al. (MSRA)



# Some Recent Classification Architectures

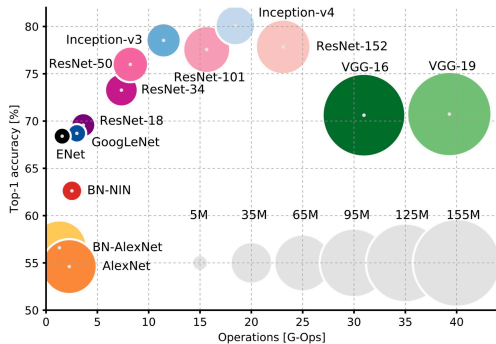
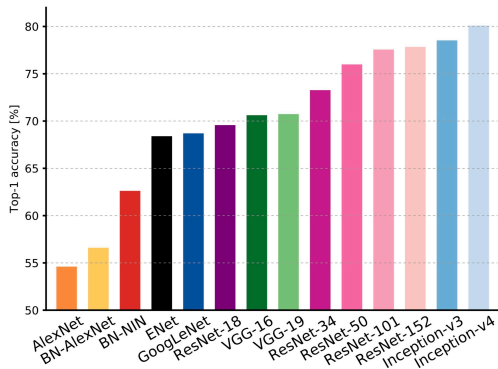
- ▶ AlexNet (Krizhevsky, Sutskever, and E. Hinton 2012) 233MB
- ▶ Network in Network (Lin, Chen, and Yan 2013) 29MB
- ▶ VGG (Simonyan and Zisserman 2015) 549MB
- ▶ GoogleNet (Szegedy, Liu, et al. 2015) 51MB
- ▶ ResNet (K. He et al. 2016) 215MB
- ▶ Inception-ResNet (Szegedy, Vanhoucke, et al. 2016)
- ▶ DenseNet (Huang et al. 2017)
- ▶ Xception (Chollet 2017)
- ▶ MobileNetV2 (Sandler et al. 2018)
- ▶ ShuffleNet (Zhang, Zhou, et al. 2018)



# Some Recent Classification Architectures

- ▶ AlexNet (Krizhevsky, Sutskever, and E. Hinton 2012) 233MB
- ▶ Network in Network (Lin, Chen, and Yan 2013) 29MB
- ▶ VGG (Simonyan and Zisserman 2015) 549MB
- ▶ GoogleNet (Szegedy, Liu, et al. 2015) 51MB
- ▶ ResNet (K. He et al. 2016) 215MB
- ▶ Inception-ResNet (Szegedy, Vanhoucke, et al. 2016) 23MB
- ▶ DenseNet (Huang et al. 2017) 80MB
- ▶ Xception (Chollet 2017) 22MB
- ▶ MobileNetV2 (Sandler et al. 2018) 14MB
- ▶ ShuffleNet (Zhang, Zhou, et al. 2018) 22MB





1

<sup>1</sup>Alfredo Canziani, Adam Paszke, and Eugenio Culurciello (2017). “An analysis of deep neural network models for practical applications”. In: *arXiv preprint*.





1

<sup>1</sup>Fei-Fei Li & Justin Johnson & Serena Yeung, Stanford cs231n.



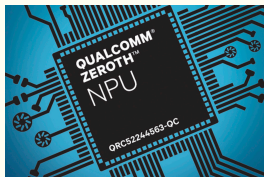
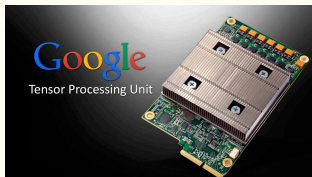
# When Machine Learning Meets Hardware

Convolution layer is one of the most expensive layers

- ▶ Computation pattern
- ▶ **Emerging** challenges

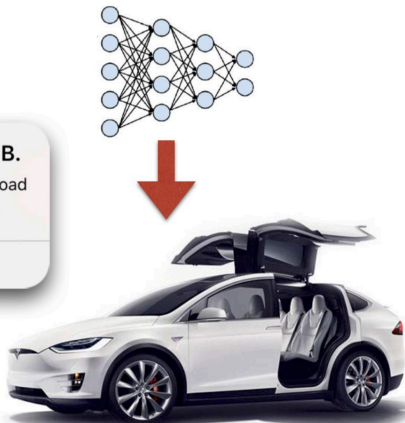
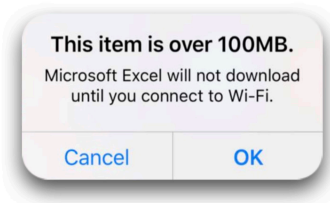
More and more end-point devices with limited memory

- ▶ Cameras
- ▶ Smartphone
- ▶ Autonomous driving



# 1st Challenge: Model Size

Hard to distribute large models through over-the-air update





## 2nd Challenge: Energy Efficiency



AlphaGo: 1920 CPUs and 280 GPUs,  
**\$3000 electric bill** per game



on mobile: **drains battery**  
on data-center: **increases TCO**



# Outline

Algorithmic Level

Architecture Level

Compilation Level

Hardware Implementation Level

Physical Synthesis Level



# Outline

Algorithmic Level

Architecture Level

Compilation Level

Hardware Implementation Level

Physical Synthesis Level



# Algorithm EX1: Object Detection [ECCV'20]

## Dive Deeper Into Box for Object Detection

Ran Chen<sup>1</sup>, Yong Liu<sup>2</sup>, Mengdan Zhang<sup>2</sup>, Shu Liu<sup>3</sup>,  
Bei Yu<sup>1</sup>, and Yu-Wing Tai<sup>2</sup>

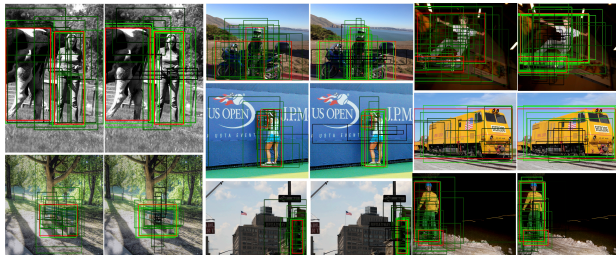
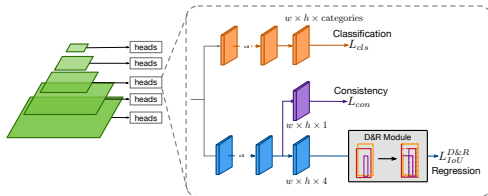
<sup>1</sup> The Chinese University of Hong Kong  
{rchen, byu}@cse.cuhk.edu.hk

<sup>2</sup> Tencent YouTu Lab

{ly.chaos, zhangmengdanrz, yuwingtai}@gmail.com

<sup>3</sup> SmartMore

sliu@smartmore.com



# Algorithm EX2: Semantic Segmentation [ECCV'20]

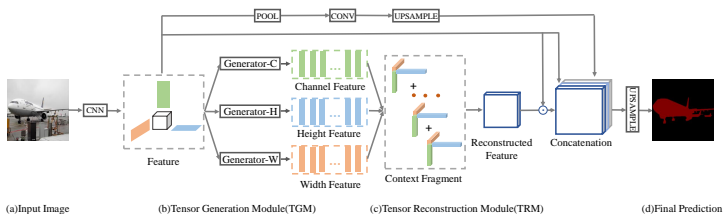
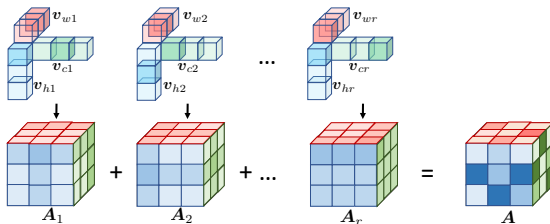
## Tensor Low-Rank Reconstruction for Semantic Segmentation

Wanli Chen<sup>1</sup>, Xinge Zhu<sup>1</sup>, Ruoqi Sun<sup>2</sup>, Junjun He<sup>2</sup>, Ruiyu Li<sup>3</sup>,  
Xiaoyong Shen<sup>3</sup>, and Bei Yu<sup>1</sup>

<sup>1</sup> The Chinese University of Hong Kong  
{wlchen,byu}@cse.cuhk.edu.hk, zx018@ie.cuhk.edu.hk

<sup>2</sup> Shanghai Jiao Tong University  
{hejunjun,ruoqisun7}@sjtu.edu.cn

<sup>3</sup> SmartMore  
{ryli,xiaoyong}@smartmore.com



# Outline

Algorithmic Level

**Architecture Level**

Compilation Level

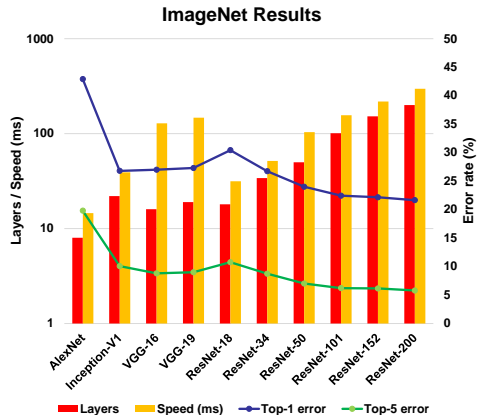
Hardware Implementation Level

Physical Synthesis Level

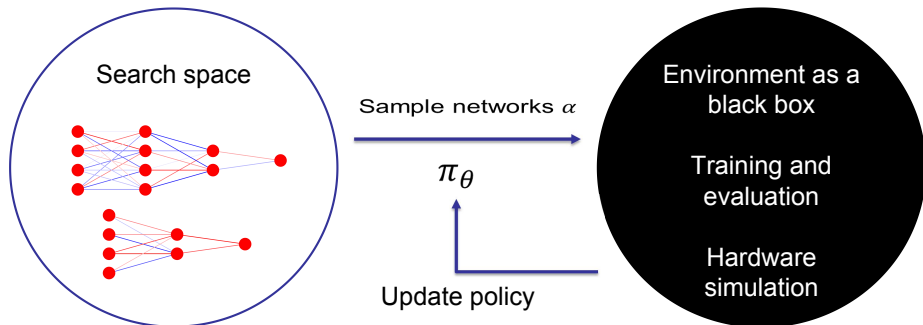


# Neural Architecture Search (NAS)

- ▶ Designing neural architecture is extremely challenging.
- ▶ Mechanism of neural networks is not well interpreted.
- ▶ Can we advance AI/ML using artificial intelligence instead of human intelligence?



# Neural Architecture Search (NAS)





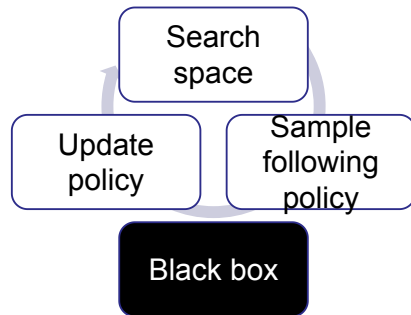
# Neural Architecture Search (NAS)

## Black box Optimization

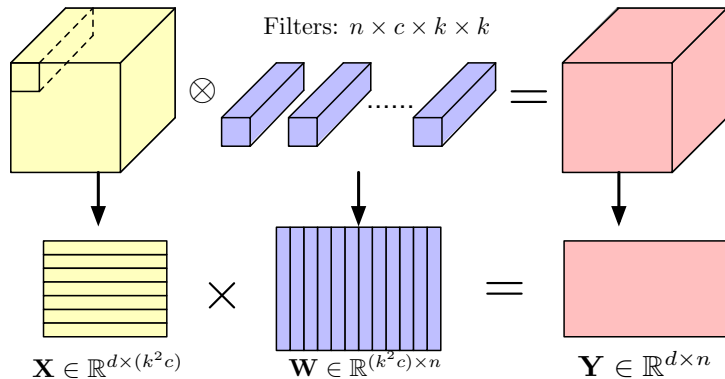
- ▶ Find the optimal network configuration to maximize the performance.
- ▶ Huge search space: *e.g.*  $1.28 \times 10^{54}$  settings.

## Available methods

- ▶ Reinforcement learning.
- ▶ Evolutionary algorithm.
- ▶ Differentiable architecture search.



# Im2col (Image2Column) Convolution



- ▶ Transform convolution to **matrix multiplication**
- ▶ **Unified** calculation for both convolution and fully-connected layers



# Compression Approach: Sparsity<sup>1, 2</sup>

$$\mathbf{X} \in \mathbb{R}^{d \times (k^2 c)} \times \mathbf{S} \in \mathbb{R}^{(k^2 c) \times n} = \mathbf{Y} \in \mathbb{R}^{d \times n}$$

## Sparse DNN

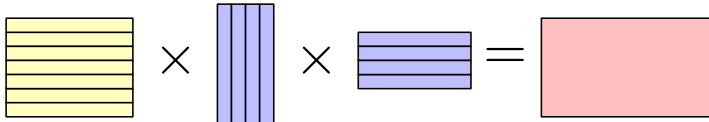
- ▶ *Sparsification*: weight pruning;
- ▶ *Compression*: compressed sparse format for storage;
- ▶ *Potential acceleration*: sparse matrix multiplication algorithm.

<sup>1</sup>Wei Wen et al. (2016). “Learning structured sparsity in deep neural networks”. In: *Proc. NIPS*, pp. 2074–2082.

<sup>2</sup>Yihui He, Xiangyu Zhang, and Jian Sun (2017). “Channel Pruning for Accelerating Very Deep Neural Networks”. In: *Proc. ICCV*.



# Compression Approach: Low-Rank<sup>1, 2</sup>


$$\mathbf{X} \in \mathbb{R}^{d \times (k^2 c)} \quad \mathbf{U} \in \mathbb{R}^{(k^2 c) \times r} \quad \mathbf{V} \in \mathbb{R}^{r \times n} \quad \mathbf{Y} \in \mathbb{R}^{d \times n}$$

## Low-rank DNN

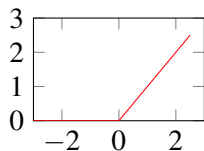
- ▶ *Low-rank approximation*: matrix decomposition or tensor decomposition.
- ▶ *Compression and acceleration*: less storage required and less FLOP in computation.

<sup>1</sup>Xiangyu Zhang, Jianhua Zou, et al. (2015). “Efficient and accurate approximations of nonlinear convolutional networks”. In: *Proc. CVPR*, pp. 1984–1992.

<sup>2</sup>Xiyu Yu et al. (2017). “On compressing deep models by low rank and sparse decomposition”. In: *Proc. CVPR*, pp. 7370–7379.



# Non-linearity Approximation<sup>1</sup>



ReLU

- ▶ Activation unit: ReLU
- ▶ Error more sensitive to positive response;
- ▶ Enlarge the solution space.

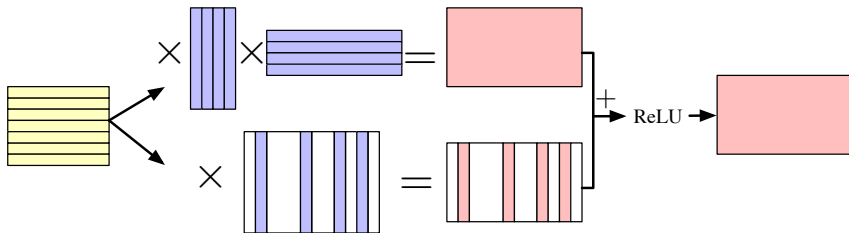
$$\min_{\mathbf{W}} \sum_{i=1}^N \|\mathbf{W}\mathbf{X}_i - \mathbf{Y}_i\|_F \rightarrow \min_{\mathbf{W}} \sum_{i=1}^N \|r(\mathbf{W}\mathbf{X}_i) - \mathbf{Y}_i\|_F$$

- ▶  $\mathbf{X}$ : input feature map
- ▶  $\mathbf{Y}$ : output feature map

<sup>1</sup>Xiangyu Zhang, Jianhua Zou, et al. (2015). "Efficient and accurate approximations of nonlinear convolutional networks". In: *Proc. CVPR*, pp. 1984–1992.



# Our Idea: Unified Structure<sup>1</sup> (Best Student Paper Award)



- ▶ Simultaneous low-rank approximation and network sparsification;
- ▶ Non-linearity is taken into account.
- ▶ Acceleration is achieved with structured sparsity.

<sup>1</sup>Yuzhe Ma et al. (2019). "A Unified Approximation Framework for Non-Linear Deep Neural Networks". In: *Proc. IGTAL*.



# Outline

Algorithmic Level

Architecture Level

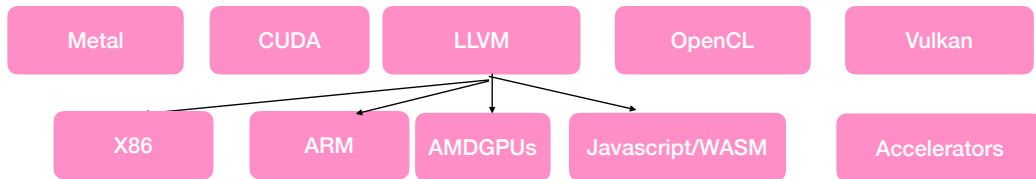
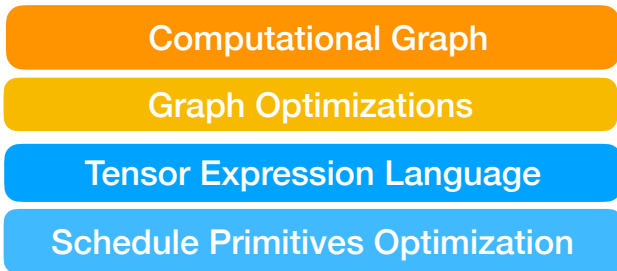
**Compilation Level**

Hardware Implementation Level

Physical Synthesis Level



# Deep Learning Compiler - TVM<sup>1</sup>



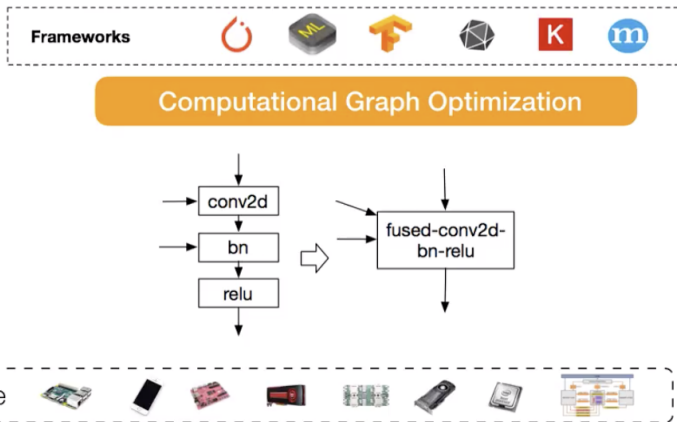
<sup>1</sup>Some materials are from CSE 599W: Systems for ML <http://dlsys.cs.washington.edu/>





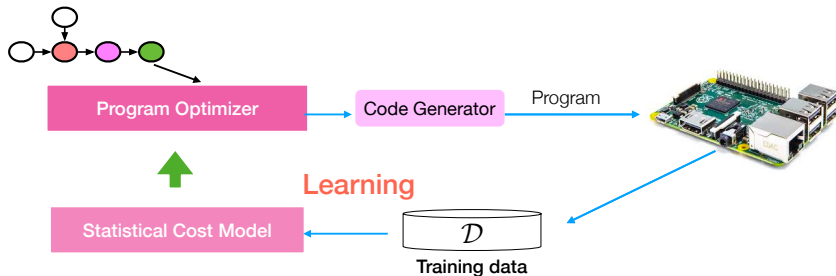
# TVM: End to End Optimization

## Computational Graph Optimization: Operator Fusion



# TVM: End to End Optimization

## Layer-wise Optimization: Autotuning



### Tuning algorithms:

- ▶ Active learning.
- ▶ Transfer learning.
- ▶ Reinforcement learning.



# TVM Domain Specific Language

## Decoupling scheduling and algorithms.

- ▶ Specify the algorithm.
- ▶ Specify the schedule.

```
A = t.placeholder((1024, 1024))
B = t.placeholder((1024, 1024))
k = t.reduce_axis((0, 1024))
C = t.compute((1024, 1024), lambda y, x:
              t.sum(A[k, y] * B[k, x], axis=k))
s = t.create_schedule(C.op)
```



```
for y in range(1024):
    for x in range(1024):
        C[y][x] = 0
        for k in range(1024):
            C[y][x] += A[k][y] * B[k][x]
```



# TVM/VTA: Full Stack Open Source System



High-level Optimizations

Tensor Program Search Space

ML-based Optimizer

VTA Runtime & JIT Compiler

VTA Hardware/Software Interface (ISA)

VTA MicroArchitecture

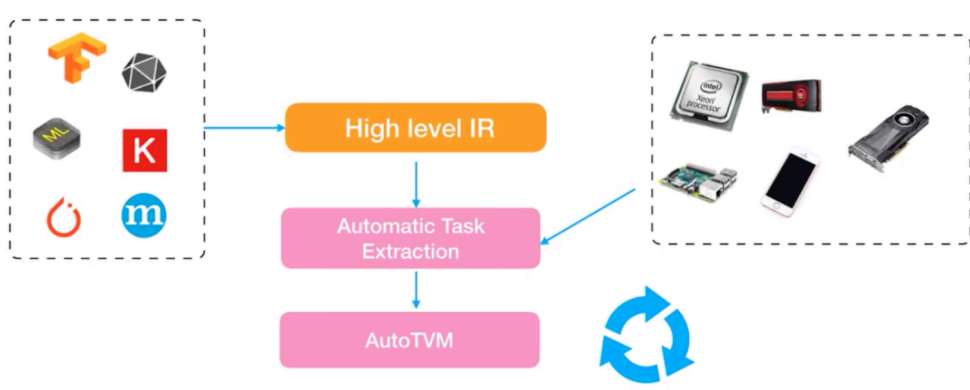
VTA Simulator



- ▶ JIT compile accelerator micro code.
- ▶ Support heterogenous devices, 10x better than CPU on the same board.
- ▶ Move hardware complexity to software.



# TVM: End-to-End Integration



# TVM Domain Specific Language + Loop Tiling

- ▶ Optimize data locality
- ▶ Minimize memory conflict
- ▶ Optimize for device cache

## + Loop Tiling

```
yo, xo, ko, yi, xi, ki = s[C].tile(y, x, k, 8, 8, 8)
```

```
for yo in range(128):  
    for xo in range(128):  
        C[yo*8:yo*8+8][xo*8:xo*8+8] = 0  
        for ko in range(128):  
            for yi in range(8):  
                for xi in range(8):  
                    for ki in range(8):  
                        C[yo*8+yi][xo*8+xi] +=  
                            A[ko*8+ki][yo*8+yi] * B[ko*8+ki][xo*8+xi]
```



# TVM Domain Specific Language + New Features

- ▶ Allow read and write to special memory scope
- ▶ Allow hook into hardware instructions
- ▶ Allow optimize for pipeline parallelism via reordering

## + Cache Data on Accelerator Special Buffer

```
CL = s.cache_write(C, vdl.a.acc_buffer)
AL = s.cache_read(A, vdl.a.inp_buffer)
# additional schedule steps omitted ...
```

## + Map to Accelerator Tensor Instructions

```
s[CL].tensorize(yi, vdl.a.gemm8x8)
```

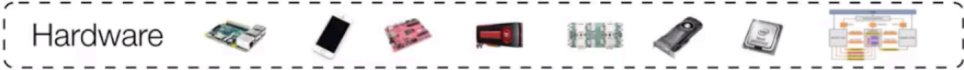
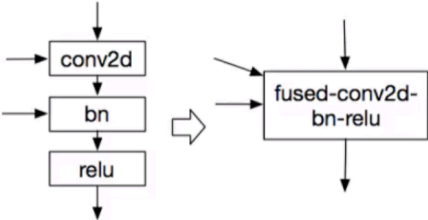
```
inp_buffer AL[8][8], BL[8][8]
acc_buffer CL[8][8]
for yo in range(128):
    for xo in range(128):
        vdl.a.fill_zero(CL)
        for ko in range(128):
            vdl.a.dma_copy2d(AL, A[ko*8:ko*8+8][yo*8:yo*8+8])
            vdl.a.dma_copy2d(BL, B[ko*8:ko*8+8][xo*8:xo*8+8])
            vdl.a.fused_gemm8x8_add(CL, AL, BL)
            vdl.a.dma_copy2d(C[yo*8:yo*8+8,xo*8:xo*8+8], CL)
```



# TVM: End to End Optimization

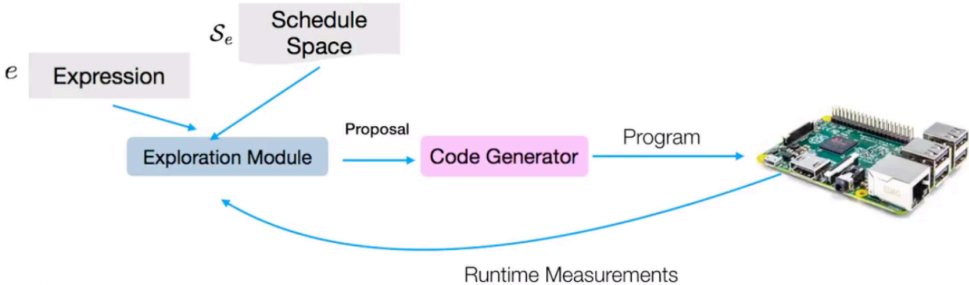


Computational Graph Optimization



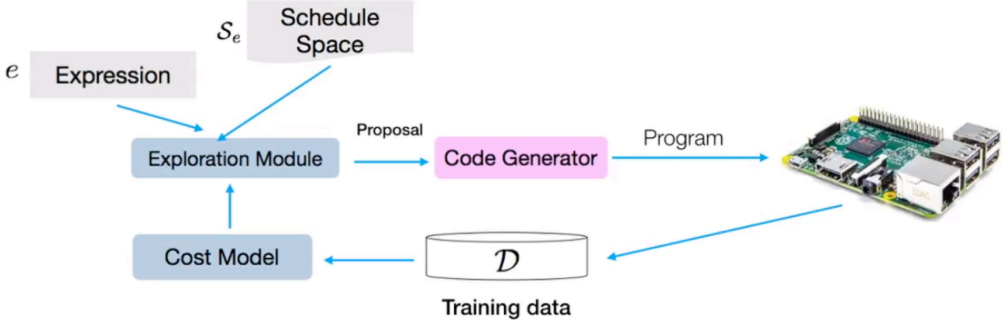


# TVM: Blackbox Autotuning



High experiment cost

# TVM: Statistical Cost Model based Approach



Learn from historical data

# Outline

Algorithmic Level

Architecture Level

Compilation Level

**Hardware Implementation Level**

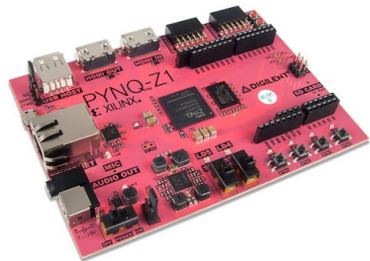
Physical Synthesis Level



# Object Detection on FPGA

- ▶ <http://www.pynq.io/>
- ▶ DAC-2018 System Design Contest
- ▶ Detailed Info:  
[www.cse.cuhk.edu.hk/~byu/2018-DAC-HDC/](http://www.cse.cuhk.edu.hk/~byu/2018-DAC-HDC/)

PYNQ™



Sponsored by:



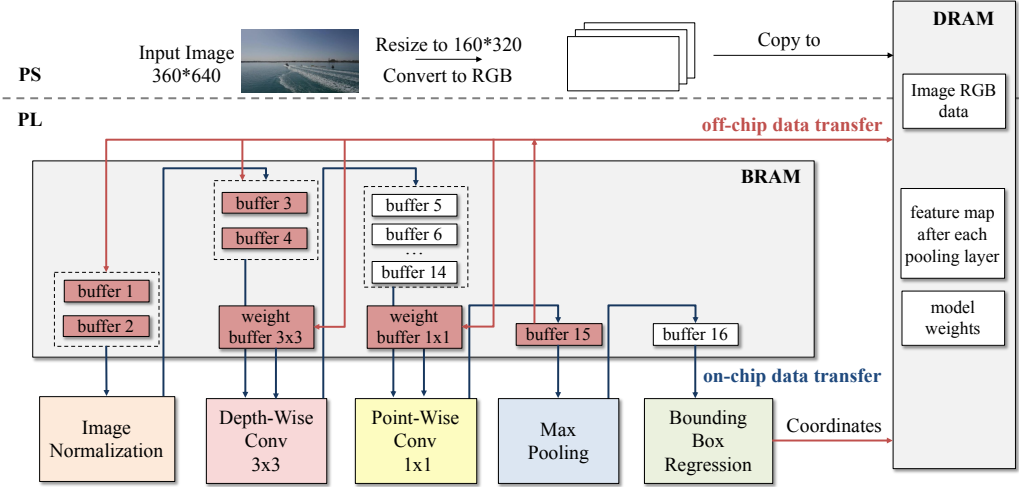
nVIDIA®



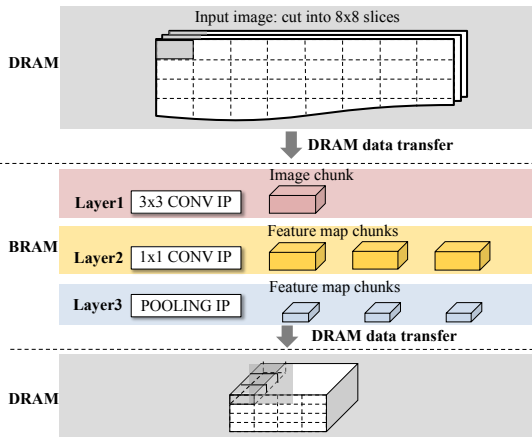
XILINX®



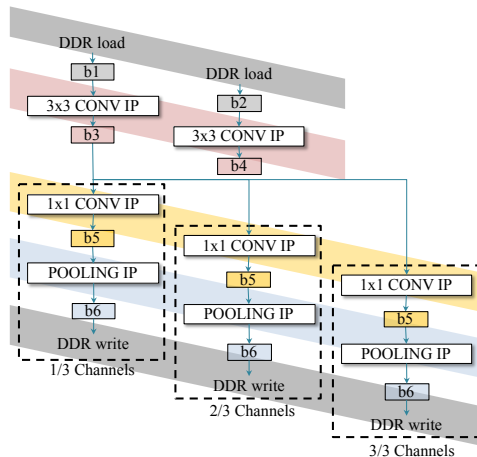
# Overall System Diagram



# Image partition, fine grained buffer scheduling, IP pipeline



**DRAM data transfer pattern**



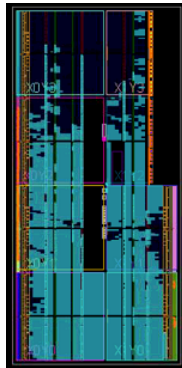
**Fine Grained IP/Buffer Scheduling**



# FPGA Design Flow

1. DNN Design in C/C++ (example)
2. Generate RTL (example) by tool Vivado HLS

```
CONV_1x1.v
CONV_3x3_group.v
Loop_0_proc.v
Relu.v
buffer_copy_from_axi.v
clear_buf.v
compute_bounding_box.v
compute_engine_16.v
copy_to_DDR_pool3.v
copy_to_DDR_pool6.v
copy_to_DDR_pool9.v
dataFlow_in_loop.v
dataFlow_parent_loop_1.v
exp_generic_Float_s.v
exp_generic_FloatqK.v
exp_generic_FloatqK_ram.dat
exp_generic_Floatrdcl.v
exp_generic_Floatrdl_ram.dat
fifo_w2_d1_A.v
fifo_w8_d2_A.v
load_bias_from_axi.v
load_image_chunk_bkb.v
load_image_chunk_bkb_ram.dat
load_image_chunk_nor.v
load_pool3_from_axi.v
load_pool6_from_axi.v
load_weight_20_from_s.v
load_weight_30_from_s.v
load_weights.v
max_pooling.v
mobilenet.v
mobilenet_AXILiteS_axi.v
mobilenet_FM_buf1Rg6.v
mobilenet_FM_buf1Rg6_ram.dat
mobilenet_FM_buf1D0.v
mobilenet_FM_buf1D0_ram.dat
mobilenet_FM_buf3bak.v
mobilenet_FM_buf3bak_ram.dat
mobilenet_FM_buf_bTr.v
mobilenet_FM_buf_bTr_ram.dat
mobilenet_IMG_m_axi.v
mobilenet_INPUT_r_m_axi.v
mobilenet_OUTPUT_r_m_axi.v
mobilenet_ama_addrqcg.v
mobilenet_ama_addrcud.v
mobilenet_ama_addrfe.v
mobilenet_ama_addr0g.v
mobilenet_ama_addrg8.v
mobilenet_ama_addrIM.v
mobilenet_ama_addrIM.v
mobilenet_ama_addrb6.v
mobilenet_ama_addrng.v
mobilenet_ama_addrpcA.v
mobilenet_ap_fpext_0_no_dsp_32_lp.tcl
mobilenet_ap_uiitofr_4_no_dsp_32_lp.tcl
mobilenet_bias_buf1hA.v
mobilenet_bias_buf1hA_ram.dat
mobilenet_fpext_3vdy.v
mobilenet_mac_mulFv1.v
mobilenet_mac_mulhbi.v
mobilenet_mac_mulsc4.v
mobilenet_mac_mulxds.v
mobilenet_nu1_mullbs.v
mobilenet_nu1_multde.v
mobilenet_max_164_jbK.v
mobilenet_sdqv_17wd1.v
mobilenet_uitofp_udo.v
mobilenet_weight_0gC.v
mobilenet_weight_DgC_ram.dat
mobilenet_weight_yd2.v
mobilenet_weight_yd2_ram.dat
ny_exp_fix.v
set_bias.v
```



3. Generate bitstream by tool Vivado IDE
4. Load bitstream to FPGA board



# Outline

Algorithmic Level

Architecture Level

Compilation Level

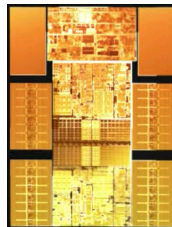
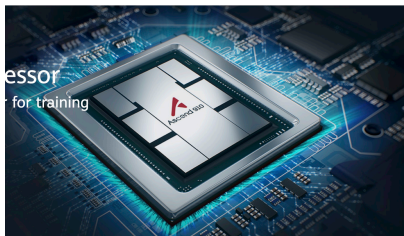
Hardware Implementation Level

**Physical Synthesis Level**

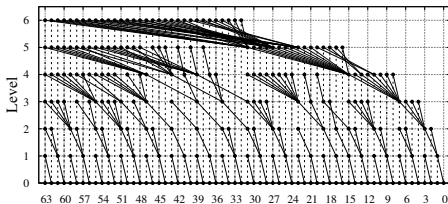




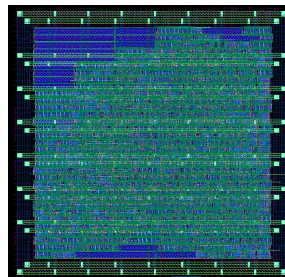
# Huawei Ascend 910



**Adder** is one of the most important component!



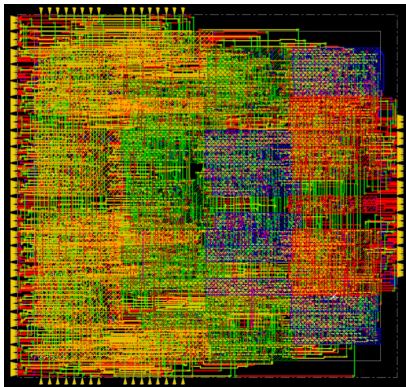
(a) Logic perspective



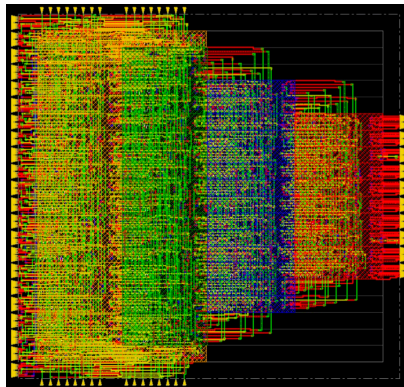
(b) Physical synthesis perspective



# EDA Challenges: How to Design an AI Chip Component?



(c) Current EDA tool output



(d) Manual design

C/C++ Programming skills are heavily required – welcome students with ICPC background

