



Algorithm-Hardware Co-Design for Edge DNN Deployment and ML for Hardware Design

Qijing Jenny Huang
University of California, Berkeley

Outline

- Motivation
- Codesign for Image Classification
 - Synetgy: Shift-based DiracDeltaNet [FPGA'19]
- Codesign for Object Detection
 - Deformable Convolution [EMC2'19]
- Reinforcement Learning for HLS
 - Autophase: Compiler Phase-Ordering [MLSys'20]
- Conclusion

Embedded Computer Vision

Applications



Robots

Drones

Autonomous
Vehicles

Security
cameras

Mobile
phones

CV Kernels/Tasks

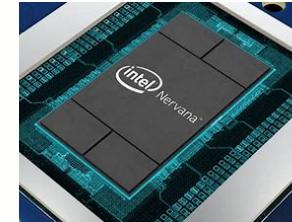


Image
Classification

Object Detection

Semantic
Segmentation

Embedded Platforms



CPU

GPU

FPGA

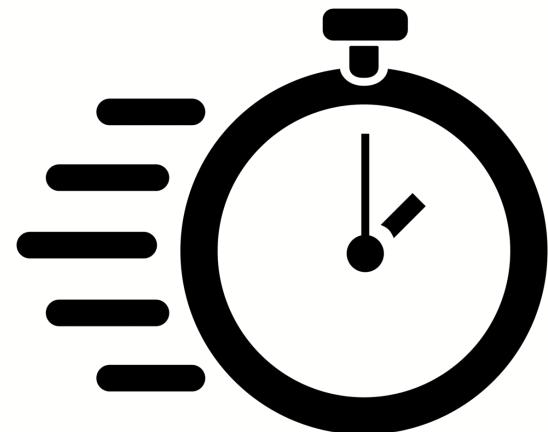
ASIC

Goals for Embedded CV

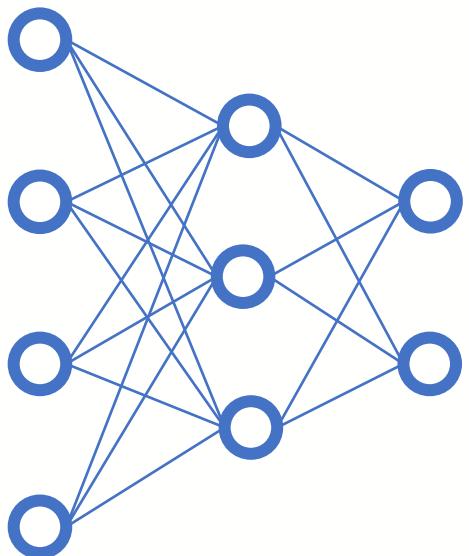
Accuracy



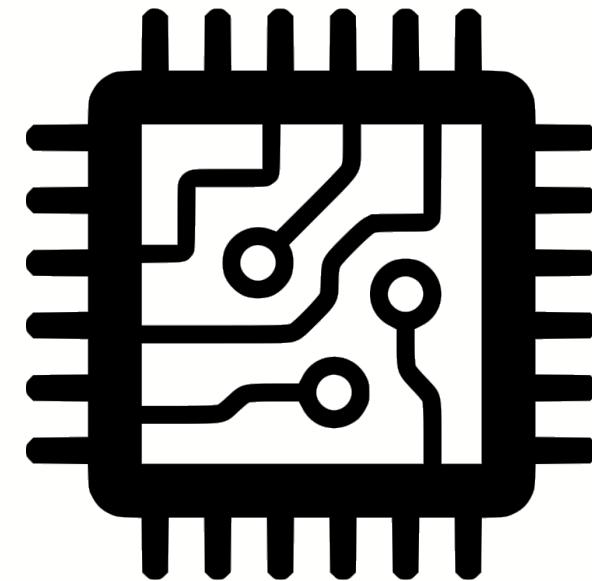
Efficiency



How to improve accuracy and efficiency?



Design better ConvNet



Design better hardware

Motivation

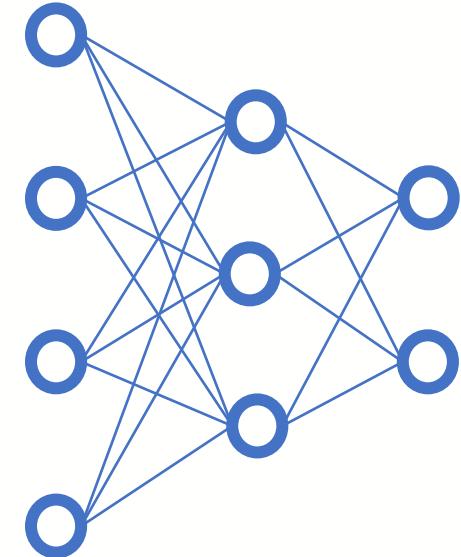
- Why codesign algorithm and hardware?
 - **Inefficient Model Designs** – many CV tasks use large inefficient models and operations solely optimized for accuracy
 - **Limited Hardware Resources** – embedded devices have limited compute resources and a strict energy and power budget
 - **Real-time Requirements** – accelerators must guarantee response within certain time constraints
- Goals: codesign **algorithms** and **accelerators** that *satisfy embedded system constraints and fall on the pareto optimal curve of the accuracy-latency tradeoff.*

Outline

- Motivation
- Codesign for Image Classification
 - Synetgy: Shift-based DiracDeltaNet [FPGA'19]
- Codesign for Object Detection
 - Synetgy: Algorithm-hardware Co-design for ConvNet Accelerators on Embedded FPGAs**
- Reinforcement Learning for ML
 - Autophase: Compiler Phase-Ordering [MLSys'20]
- Conclusion

ConvNet Design Strategies

- Strategy 1: Use efficient models
- Strategy 2: Simplify operators
- Strategy 3: Quantize

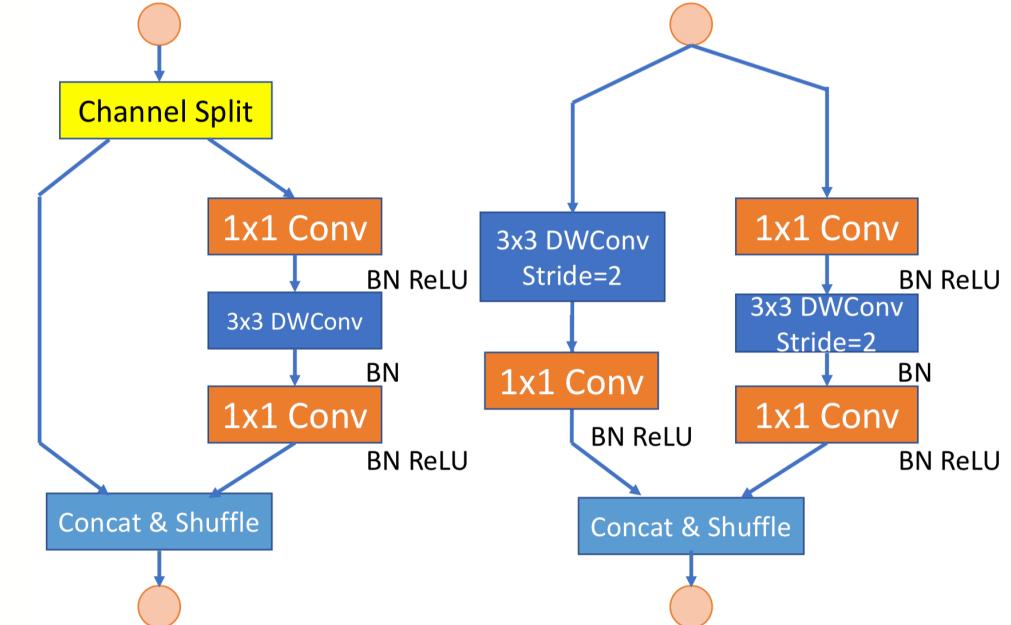


Design better ConvNet

CNN Strategy 1: Use efficient models

- **ShuffleNetV2-1.0x** [1] as our starting point
- Compared to VGG16:
 - 65x fewer OPs
 - 48x fewer parameters
 - Near equal accuracy on ImageNet

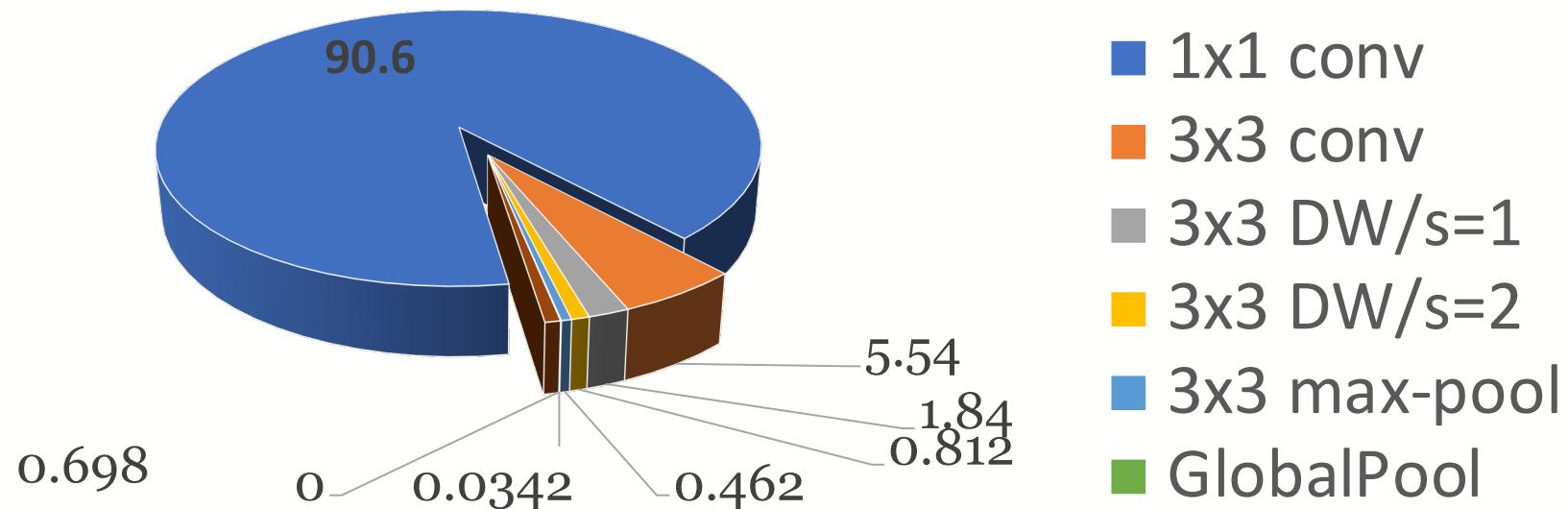
	MACs	#Params	Top-1 Acc
ShuffleNetV2-1.0x	146M	2.3M	69.4%
VGG16	15.3G	138M	71.5%



ShuffleNetV2 building blocks

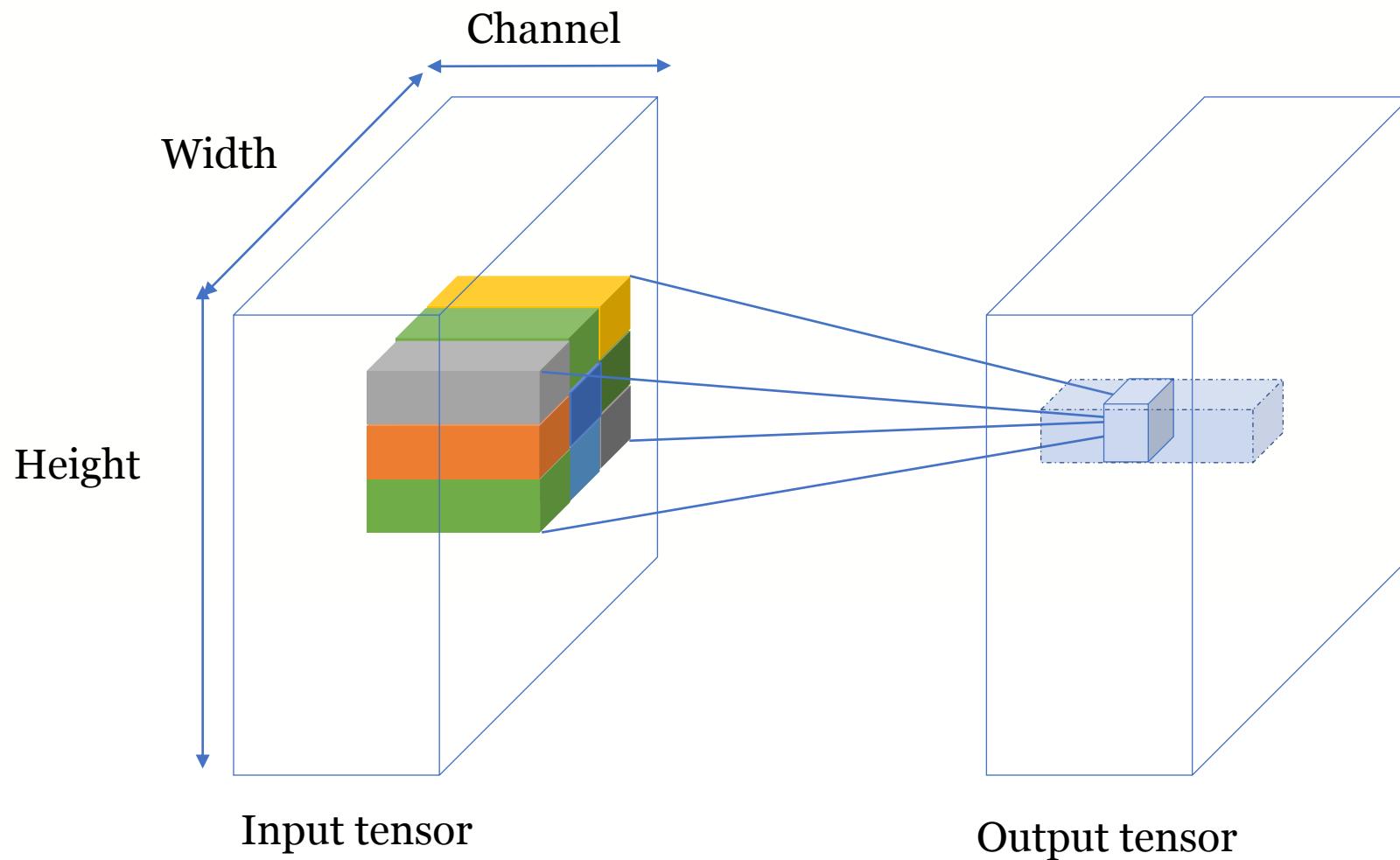
CNN Strategy 2: Simplify operators

- Can we reduce the number of operator types?
- Can we make the operation more hw-friendly?



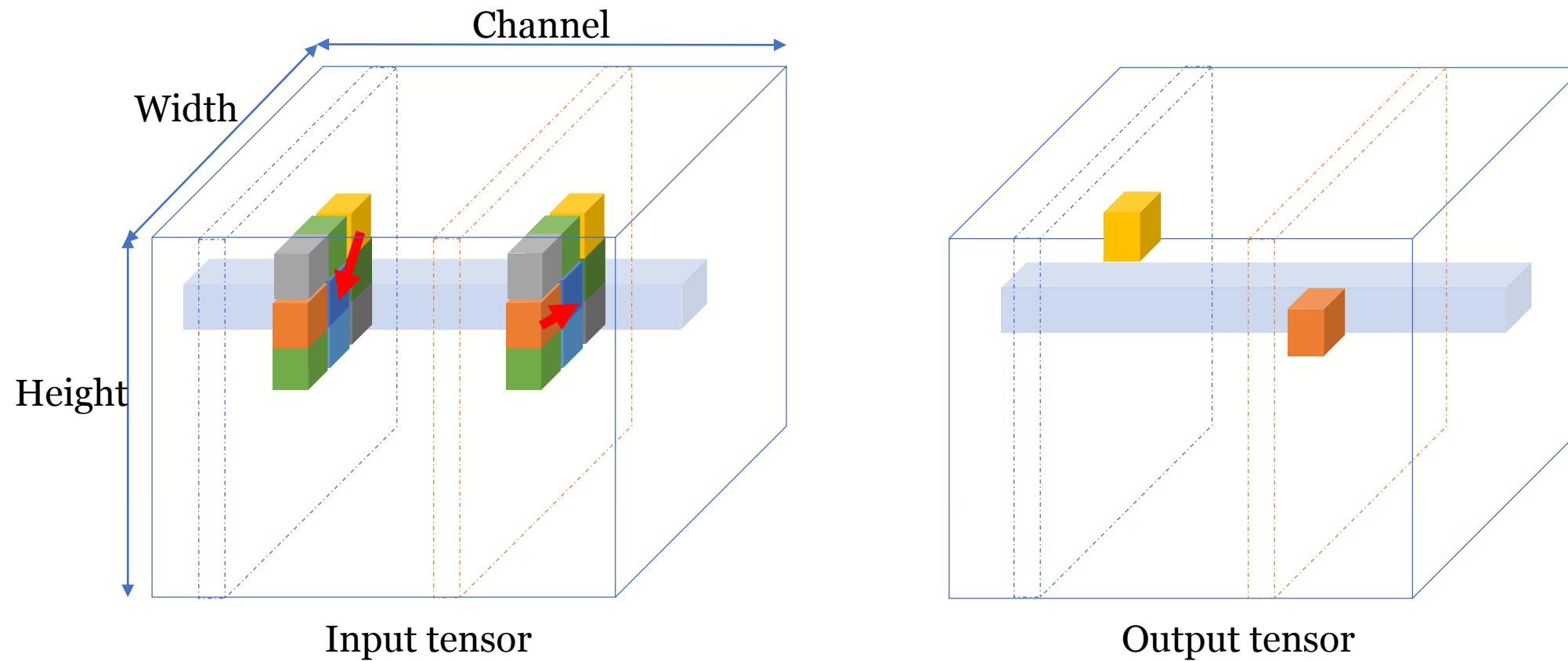
CNN Strategy 2: Simplify operators

- Can we replace 3x3 convolutions?



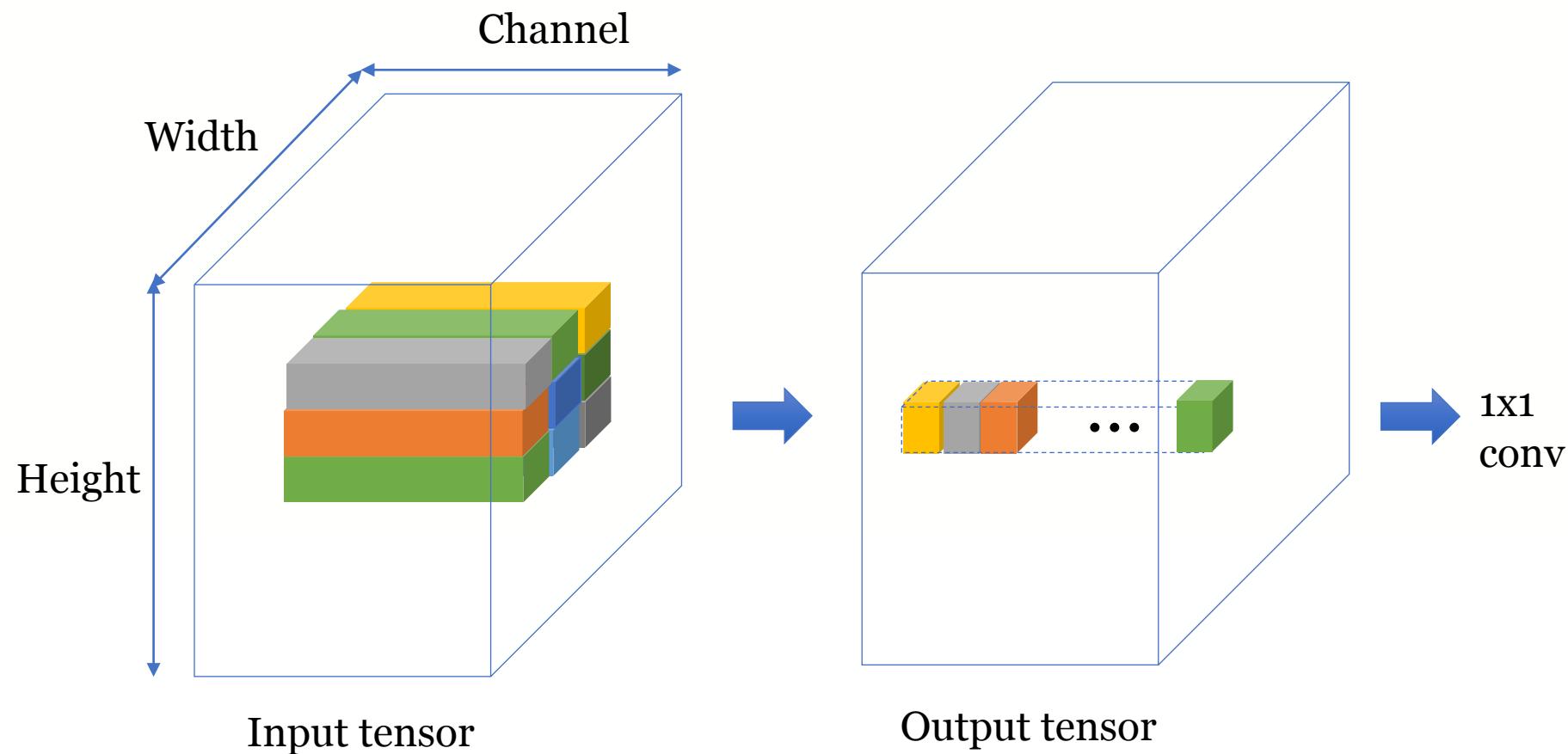
The Shift Operation

- The shift operation moves a neighboring pixel to the center position



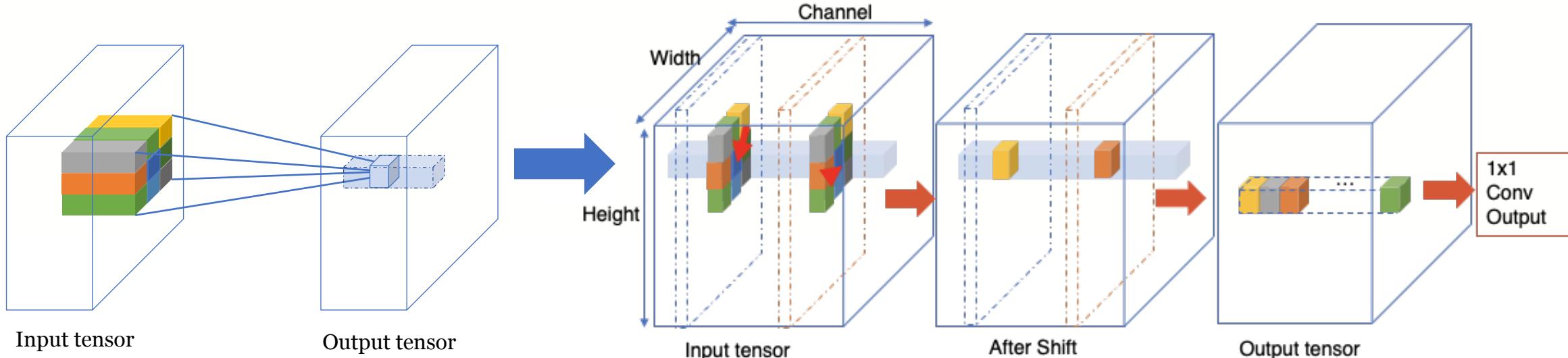
The Shift Operation

- 1x1 conv aggregates spatial information along the channel dimension

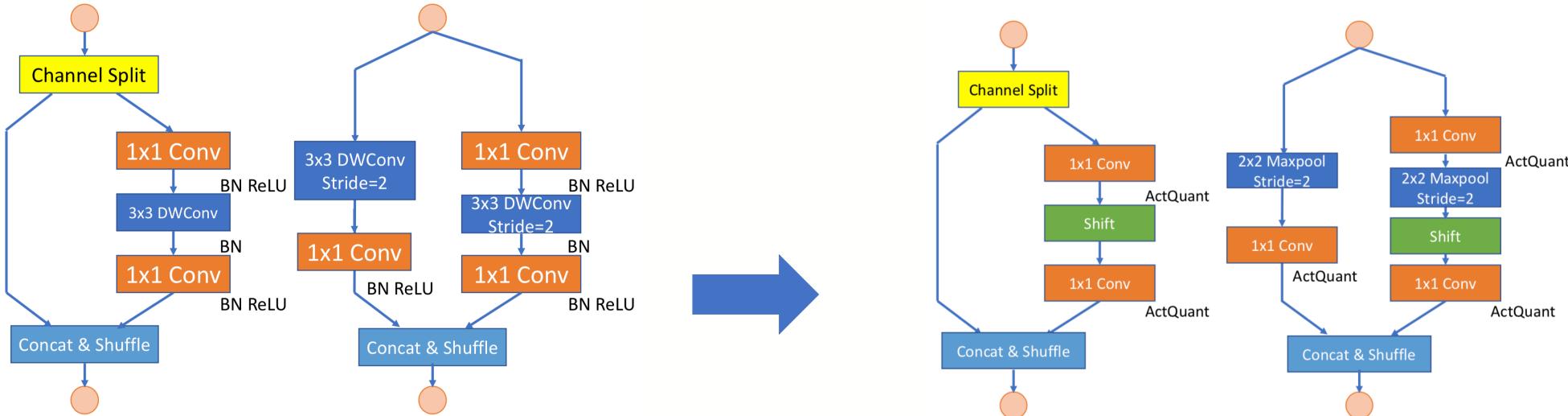


Replace 3 x 3 Conv

- **3x3 conv:** ➡
 - Aggregates neighboring pixels
 - Mixes channel info
- **shift:** Re-aligns pixels
- **1x1 conv:** Mixes channel info



ShuffleNetV2 -> DiracDeltaNet



- Accuracy (full precision): **69.4%**
- Operators involved:
 - 1x1 convolution
 - 3x3 convolution
 - 3x3 DW convolution
 - 3x3 max pooling
 - Channel split/shuffle(concat)

- Accuracy (full precision): **69.7%**
- Operators involved:
 - 1x1 convolution
 - 2x2 max pooling
 - Channel split/shuffle/shift/concat

CNN Strategy 3: Quantize

- Quantization has been mostly demonstrated on large networks. Is it effective on the small ones like DiracDeltaNet?
- We used existing quantization methods:
 - DoReFaNet [1] method for weights
 - Modified PACT [2] method for activations
- We achieved 4-bit weight and 4-bit activation precision with competitive accuracy

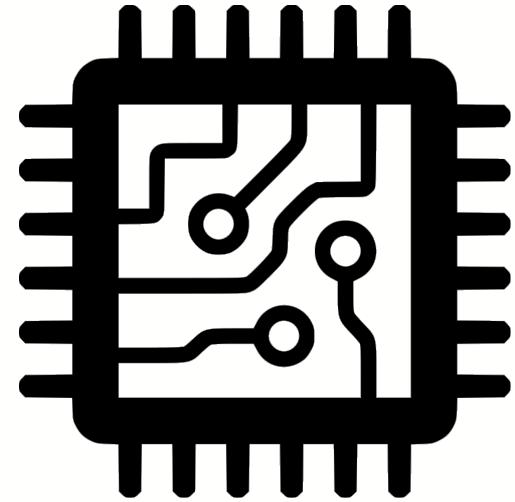
	Network	Pruning	Precision	Top-1 Acc
[3]	VGG16	Yes	8-8b	67.72%
Ours	DiracDeltaNet	No	4-4b	67.52%

[1] Zhou, S., Wu, Y., Ni, Z., Zhou, X., Wen, H. and Zou, Y. DoReFa-Net: Training Low Bitwidth Convolutional Neural Networks with Low Bitwidth Gradients.

[2] Choi, J., Wang, Z., Venkataramani, S., I-Jen Chuang, P., Srinivasan, V. and Gopalakrishnan, K. PACT: Parameterized Clipping Activation for Quantized Neural Networks.

Hardware Design Strategies

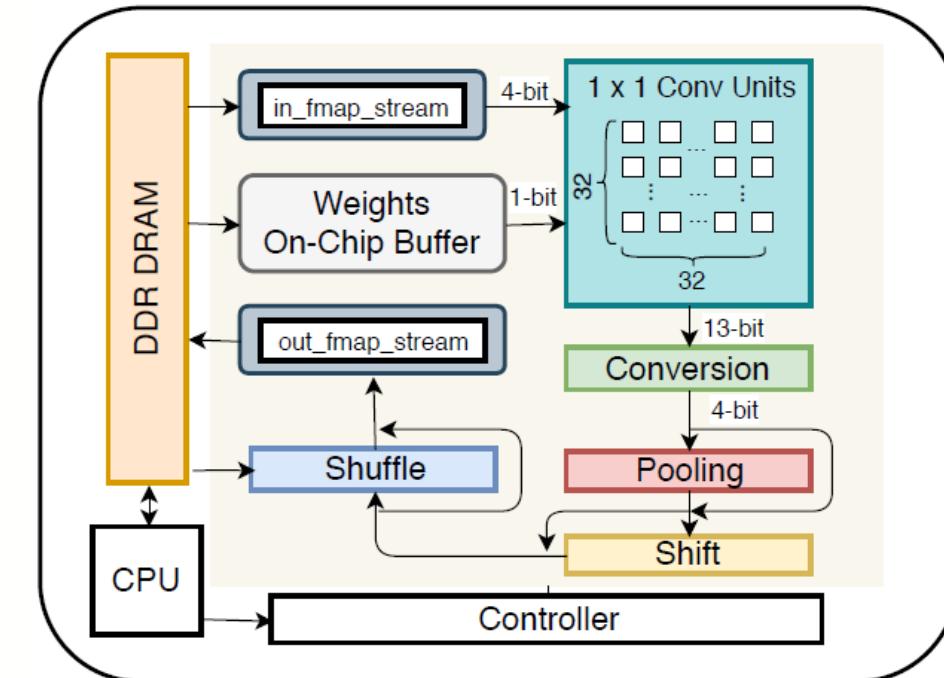
- Strategy 1: Specialize conv engine
- Strategy 2: Use dataflow architecture
- Strategy 3: Merge layers



Design better hardware

HW Strategy 1: Specialize conv engine

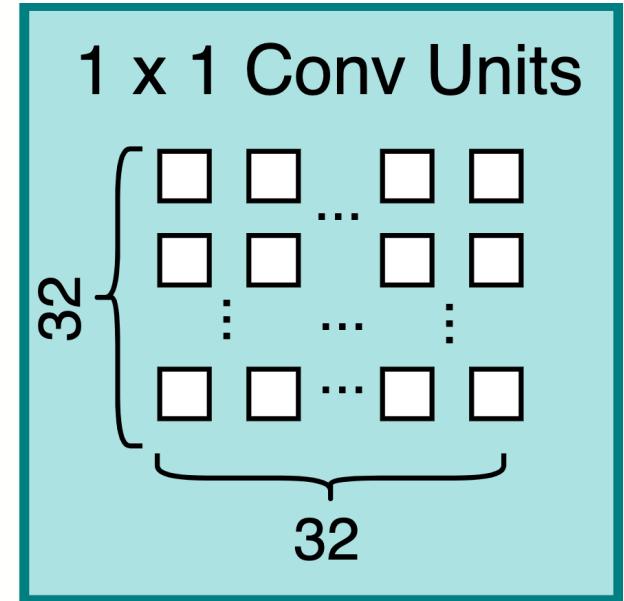
- HW engine supports:
 - 1x1 conv
 - 2x2 max-pooling
 - shift
 - shuffle
- Layer-based design
- Implemented with Vivado HLS and PYNQ



Architecture Diagram

HW Strategy 1: Specialize conv engine

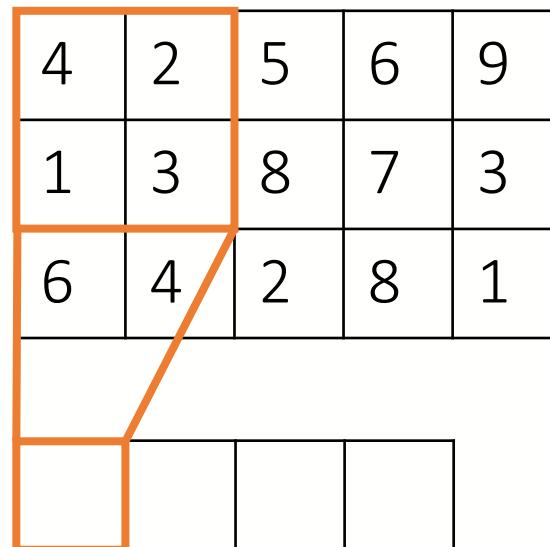
- 1x1 Conv Unit:
 - Supports matrix-vector multiplication
 - 4-bit inputs
 - 4-bit weights
 - 17-bit partial sums
 - Buffers weights and partial sums on-chip
 - Performs 32×32 MACs per iteration
 - Each input gets reused *output channel size* times



Strategy 2: Use dataflow architecture

- 1x1 conv
 - No line-buffer
- shift
 - 3x3 sliding window, $lI=1$
 - 2x2 max-pooling
 - 2x2 sliding window, $lI=2$

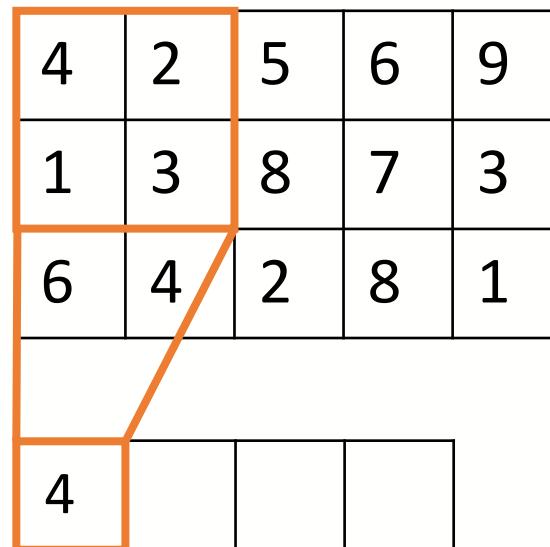
**2x2 max-pooling
example:**



Strategy 2: Use dataflow architecture

- 1x1 conv
 - No line-buffer
- shift
 - 3x3 sliding window, $lI=1$
 - 2x2 max-pooling
 - 2x2 sliding window, $lI=2$

2x2 max-pooling
example:



Strategy 2: Use dataflow architecture

- 1x1 conv
 - No line-buffer
- shift
 - 3x3 sliding window, $lI=1$
 - 2x2 max-pooling
 - 2x2 sliding window, $lI=2$

2x2 max-pooling
example:

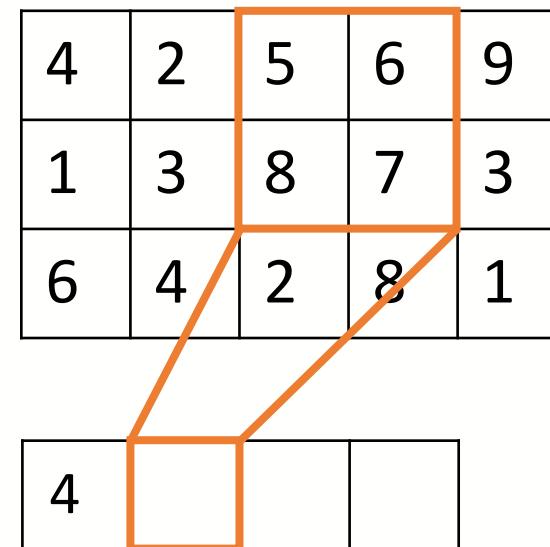
4	2	5	6	9
1	3	8	7	3
6	4	2	8	1

4			
---	--	--	--

Strategy 2: Use dataflow architecture

- 1x1 conv
 - No line-buffer
- shift
 - 3x3 sliding window, $l=1$
- 2x2 max-pooling
 - 2x2 sliding window, $l=2$

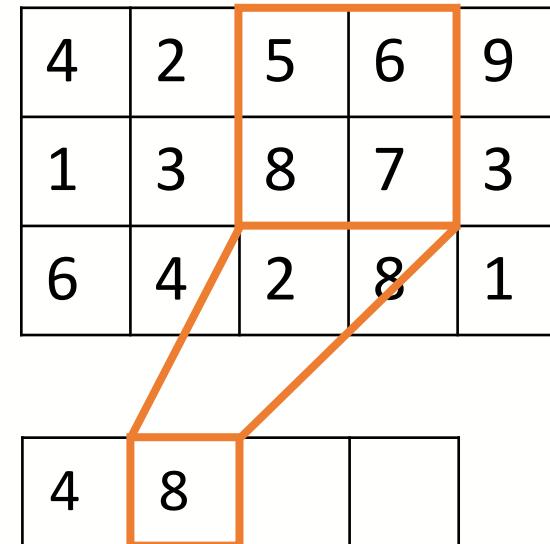
2x2 max-pooling
example:



Strategy 2: Use dataflow architecture

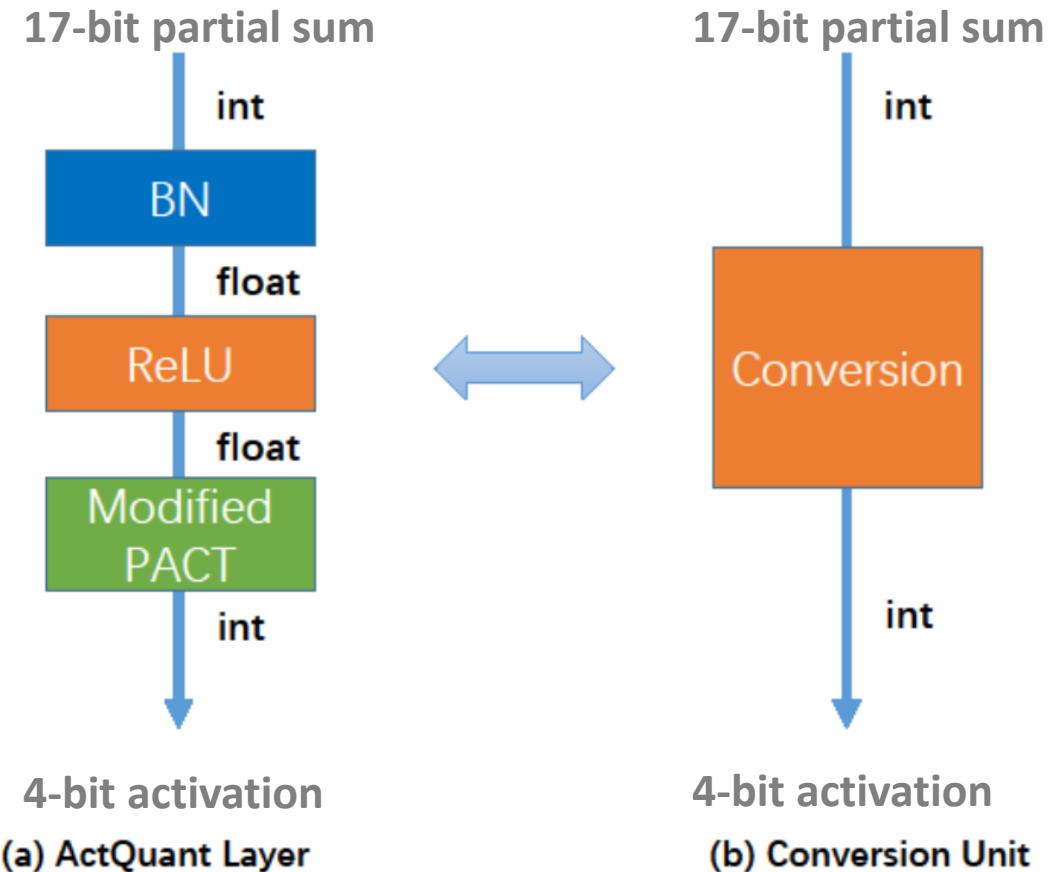
- 1x1 conv
 - No line-buffer
- shift
 - 3x3 sliding window, $l=1$
- 2x2 max-pooling
 - 2x2 sliding window, $l=2$

2x2 max-pooling
example:



Strategy 3: Merge layers

- Conversion unit includes:
 - Batch Norm
 - ReLU
 - Modified PACT
- It performs 17-bit to 4-bit conversion
- It is implemented with comparators



Comparison with Previous Work

					Energy /
Algorithm-hardware co-design can achieve both high accuracy (67.5% top-1) and good efficiency (66 FPS) for image classification					
[3]	Stratix V	3.8	66.55%	3.165	0.023
Ours	Zynq ZU3EG	66.3	67.52%	4.4b	0.083

- Equal top-1 accuracy
- 11.6x higher framerate
- 6.3x more power efficient

[1] Guo, K., Han, S., Yao, S., Wang, Y., Xie, Y. and Yang, H. Software-Hardware Codesign for Efficient Neural Network Acceleration. IEEE Micro, 37.

[2] Qiu, J., Wang, J., Yao, S., Guo, K., Li, B., Zhou, E., Yu, J., Tang, T., Xu, N., Song, S., Wang, Y. and Yang, H. Going Deeper with Embedded FPGA Platform for Convolutional Neural Network. FPGA 2016.

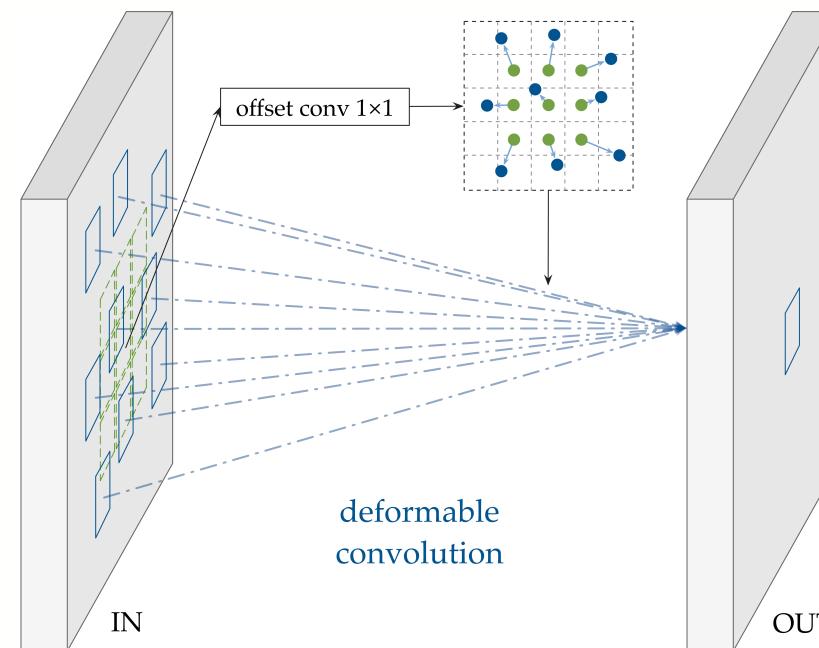
[3] Suda, N., Chandra, V., Dasika, G., Mohanty, A., Ma, Y., Vrudhula, S.B.K., Seo, J.S. and Cao, Y. "Throughput-Optimized OpenCL-based FPGA Accelerator for Large-Scale Convolutional Neural Networks". FPGA 2016.

Outline

- Motivation
- Codesign for Image Classification
 - Shift-based DiracDeltaNet
- Other Co-design Examples
 - Algorithm-Hardware Co-design for Deformable Convolution**
- Reinforcement Learning for HLS
 - Compiler Phase-Ordering
- Conclusion

Deformable Convolution

- **Deformable Convolution** is an input-adaptive dynamic operation that samples inputs from variable spatial locations



1. Generate offsets
2. Sample from input feature map

Deformable Convolution

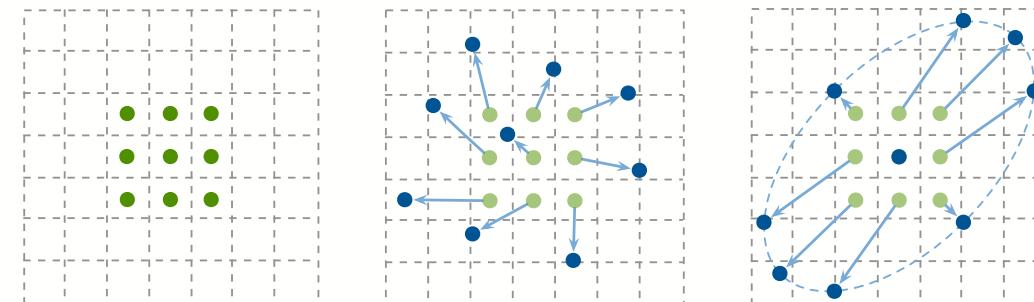
- **Deformable Convolution** is an input-adaptive dynamic operation that samples inputs from variable spatial locations
- Its sampling locations vary with:
 - Different input images
 - Different output pixel locations



Sampling Locations (in red) for Different Output Pixels (in green)

Deformable Convolution

- **Deformable Convolution** is an input-adaptive dynamic operation that samples inputs from variable spatial locations
- Its sampling locations vary with:
 - Different input images
 - Different output pixel locations
- It captures the spatial variance of objects with different:
 - Scales
 - Aspect Ratios
 - Rotation Angles



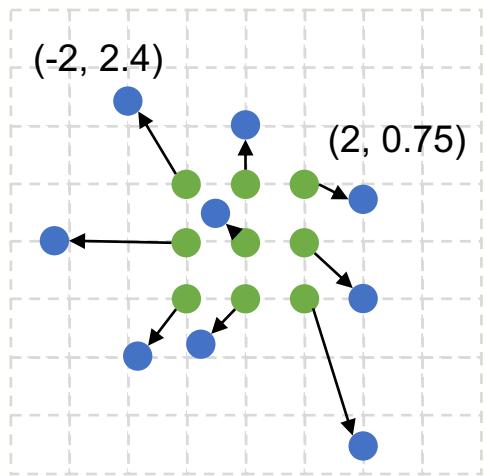
Variable Receptive Fields

Deformable Convolution

- **Deformable Convolution** is an input-adaptive dynamic operation that samples inputs from variable spatial locations
- Its sampling locations vary with:
 - Different input images
 - Different output pixel locations
- It captures the spatial variance of objects with different:
 - Scales
 - Aspect Ratios
 - Rotation Angles
- Challenges:
 - Additional compute and memory requirements for offset generation
 - Irregular input-dependent memory access patterns
 - Not friendly for dataflows that leverage the spatial reuse

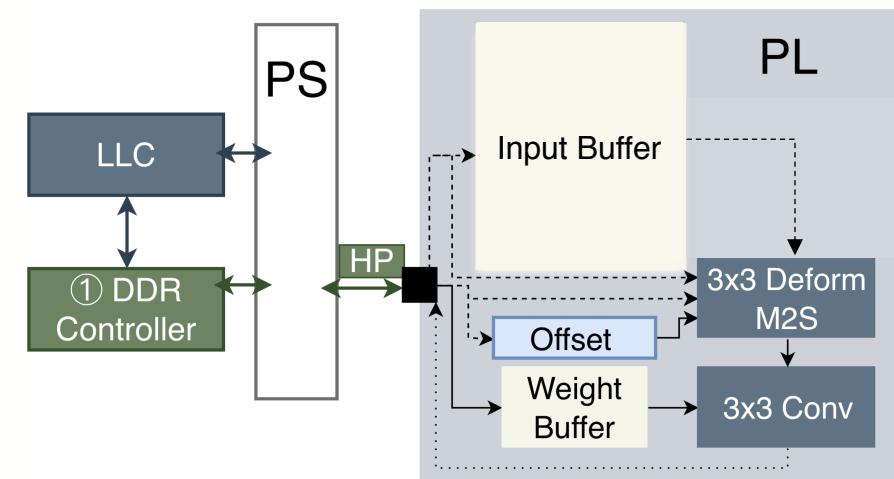
Operation Codesign

Algorithm Modification:



0. Original Deformable

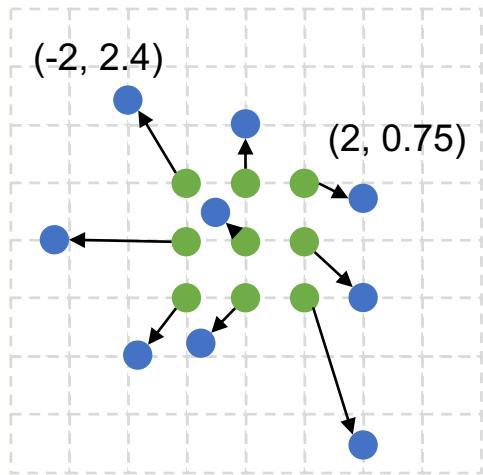
Hardware Optimization:



- Preloads weights to on-chip buffer
- Loads input and offsets directly from DRAM

Operation Codesign

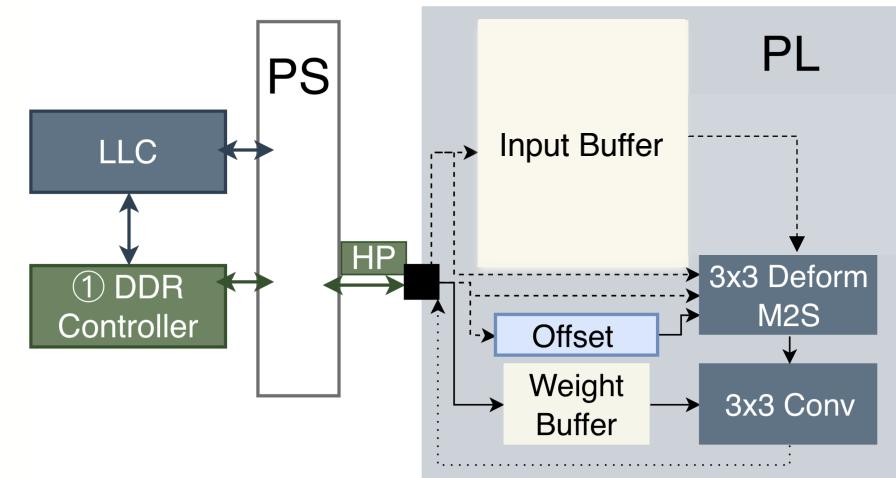
Algorithm Modification:



1. Depthwise Deformable

Accuracy¹(AP): 42.9

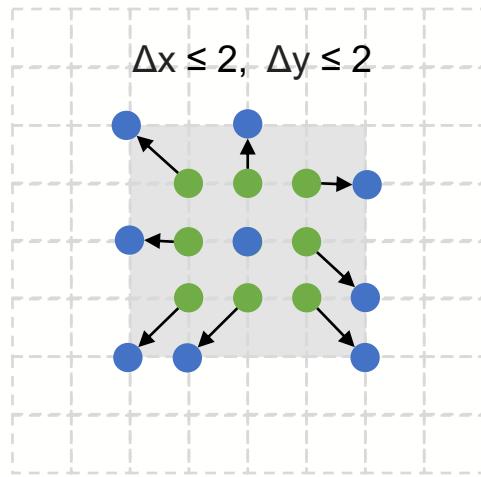
Hardware Optimization:



- Reduce the total MACs

Operation Codesign

Algorithm Modification:

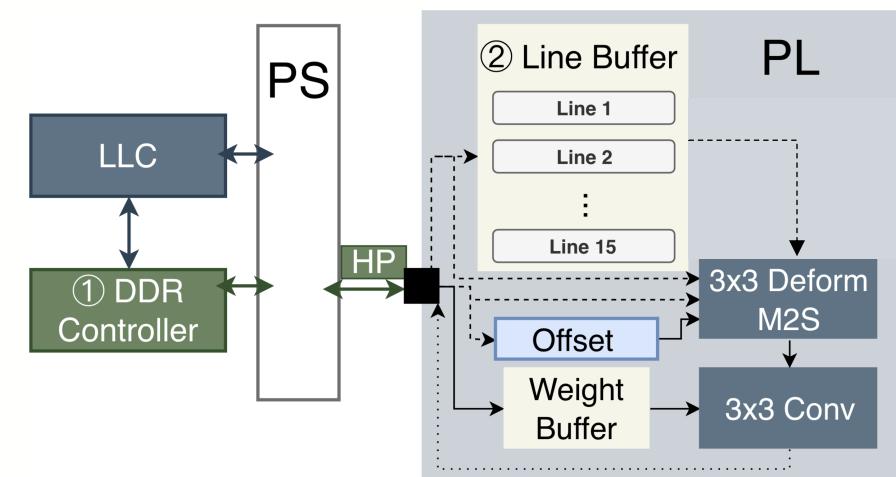


2. Bounded Range

Accuracy¹(AP): 41.0

↓ 1.9

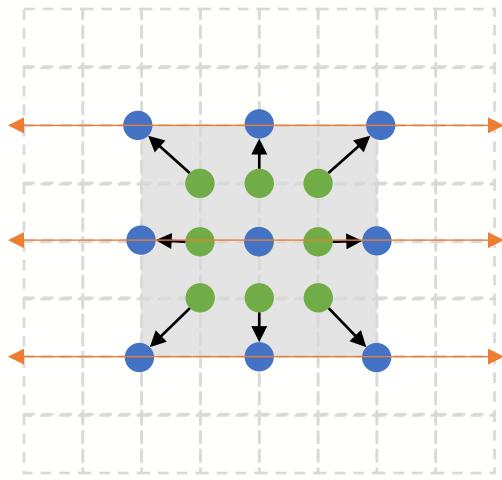
Hardware Optimization:



- **Buffers inputs in the on-chip line buffer to allow spatial reuse**

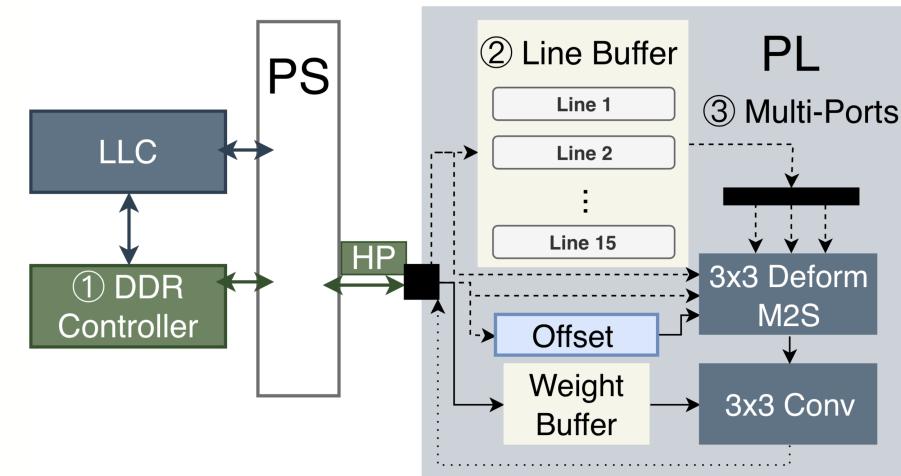
Operation Codesign

Algorithm Modification:



Accuracy¹(AP): 41.1 ↑ 0.1

Hardware Optimization:



- Improves on-chip memory bandwidth

Ablation Study

Object Detection Accuracies

Operation	Depthwise	Bound	Square	VOC			COCO			
				AP	AP50	AP75	AP	AP50	AP75	APs
3 × 3				39.2	60.8	41.2	21.4	36.5	21.5	7.3
3 × 3	✓			39.1	60.9	40.9	19.8	34.3	19.7	6.3
5 × 5	✓			40.6	62.4	42.6	21.3	36.4	21.3	6.7
7 × 7	✓			41.9	63.8	43.8	21.7	37.2	21.5	6.9
9 × 9	✓			42.3	64.8	44.3				
deform	✓			42.9	64.4	45.7	23.0	38.4	23.3	6.9
deform	✓	✓		41.0	63.0	42.9	21.3	36.4	21.1	7.2
deform	✓	✓	✓	41.1	63.1	43.7	21.5	36.8	21.5	6.5

5x less
compute

- 3x3 deformable convolution is more efficient than large convolution kernels
- The accuracy of the codesigned deformable convolution (1.8 AP difference on Pascal VoC and 1.5 AP difference on COCO)

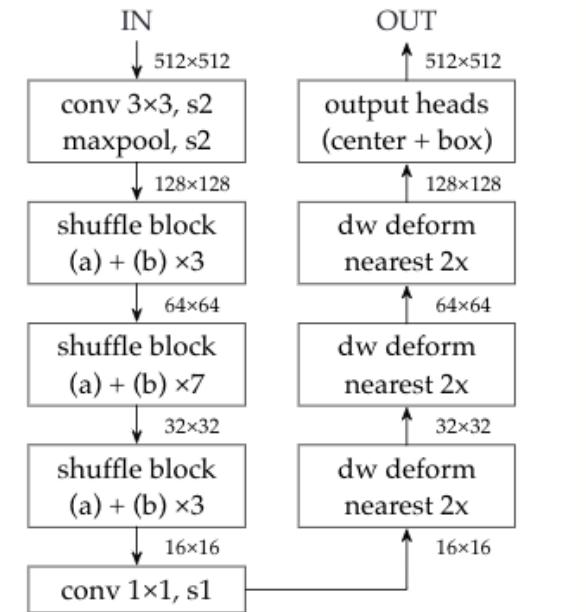
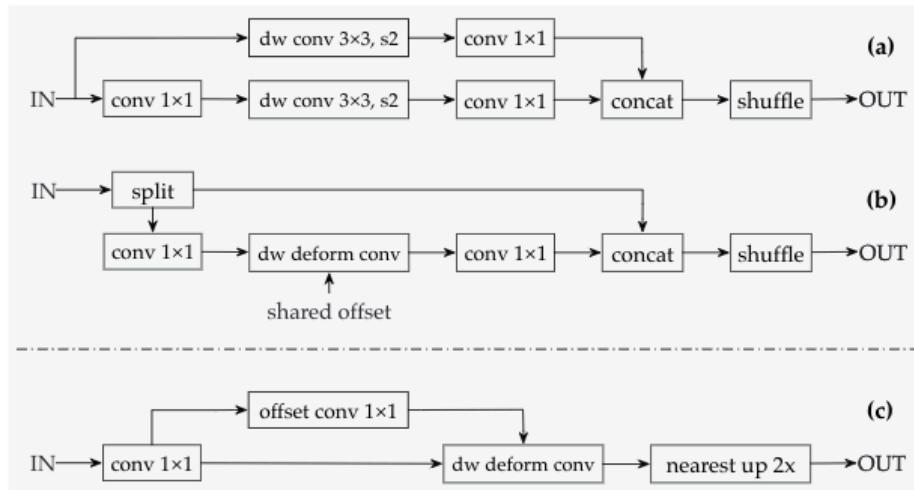
Ablation Study

Hardware Performance

Operation	Original	Deformable	Bound (buffered)	Square (multi-ported)	Without LLC		With LLC	
					Latency (ms)	GOPs	Latency (ms)	GOPs
Full 3×3 Conv	✓	✓ ✓ ✓	✓ ✓ ✓	✓	43.1	112.0	41.6	116.2
					59.0	81.8	42.7	113.1
					43.4	111.5	41.8	115.5
					43.4	111.5	41.8	115.6
Depthwise 3×3 Conv	✓	✓ ✓ ✓	✓ ✓ ✓	✓	1.9	9.7	2.0	9.6
					20.5	0.9	17.8	1.1
					3.0	6.2	3.4	5.5
					2.1	9.2	2.3	8.2

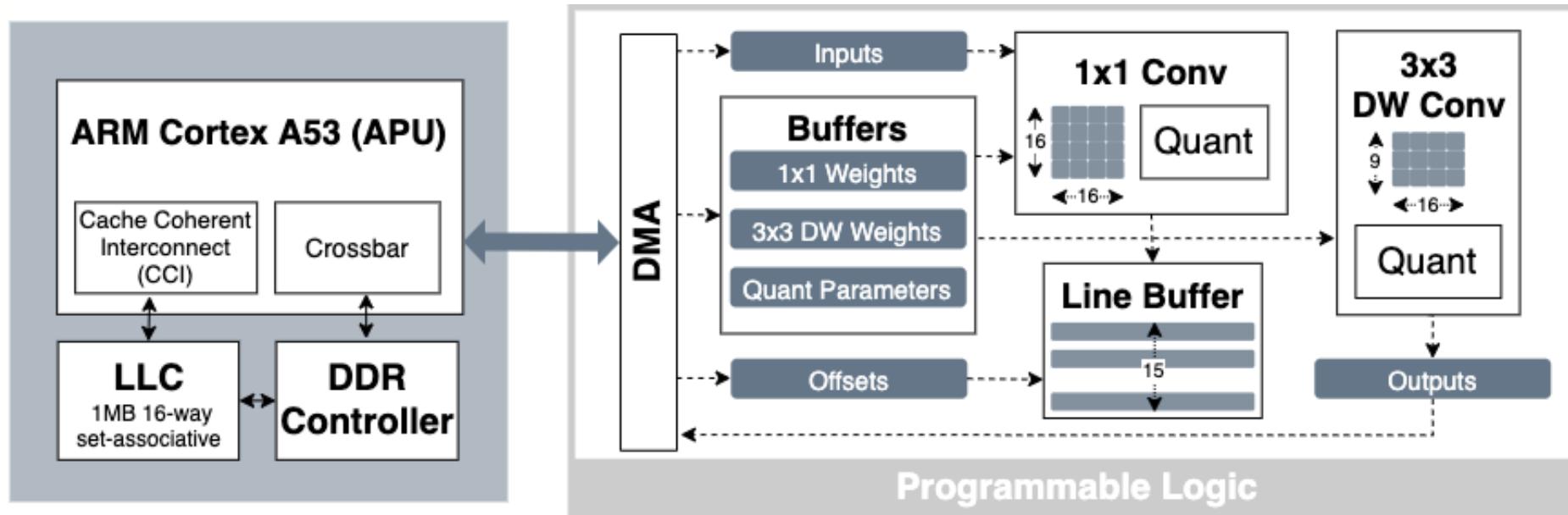
- Our algorithm-hardware co-design methodology for the deformable convolution achieves a **1.36×** and **9.76×** speedup respectively for the **full** deformable convolution and **depthwise** deformable convolution on FPGA

DNN Architecture



- Simple building blocks to reduce the hardware complexity
- Anchor-free detection system to reduce the postprocessing overhead for Non Maximum Suppression (NMS)

Hardware Accelerator



- Specialized engine for efficiency
- Weight buffers and line buffers for maximum data reuse

Results

Object Detection Accuracies

Detector	Weights	Activations	Model Size	MACs	AP50
Tiny-YOLO	32-bit	32-bit	60.5 MB	3.49 G	57.1
CoDeformNet 1x	32-bit	32-bit	6.06 MB	1.14 G	64.6

- **Hardware Acceleration** achieves both high accuracy (55.1 mAP) and real-time inference (26.9 FPS) for object detection on embedded FPGAs

	Platform	Framerate (fps)	Test Dataset	Precision	Accuracy
Finn-R [2] [26]	Zynq XCZU3EG (Ultra96)	16	VOC07	1-3b	mAP(50.1)
Ours (256 × 256)	Zynq XCZU3EG (Ultra96)	26.9	VOC07	4-8b	mAP(55.1)

- Our design achieves both higher framerate and higher accuracy

Outline

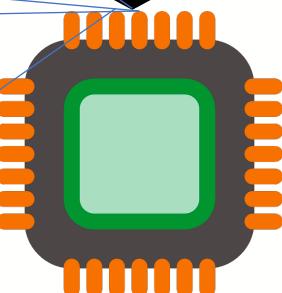
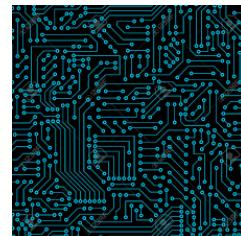
- Motivation
- Codesign for Image Classification
 - Shift-based DiracDeltaNet
- Codesign for AutoPhase: Juggling HLS Phase Orderings in Random Forests with Deep Reinforcement Learning
- Reinforcement Learning for HLS
 - Compiler Phase-Ordering
- Conclusion

Motivation

```
#include <stdio.h>
int main (void) {
    printf ("Hello, World!\n");
    return 0;
}
```



```
01001001
11011101
01001111
01001001
11011101
```



For years, compiler optimizations rely on **heuristics** and **hand engineering**

Compiler Phase Ordering

clang program.c -flag1 -flag2 ...

Which passes?
In which order?
More than 2^{247} possibilities ...

NP-Hard

Compiler Phase Ordering

Example:

Normalizing a Vector

```
for (int i=0; i<n; i++) {  
    out[i] = in[i] / mag(in, n);  
}
```

Pass #1: -licm (loop-invariant code motion)

```
double precompute = mag(in, n);  
for (int i=0; i<n; i++) {  
    out[i] = in[i] / precompute;  
}
```

Pass #2: -inline

```
double precompute, sum = 0;  
for (int i=0; i<n; i++) {  
    sum += A[i] * A[i];  
}  
precompute = sqrt(sum);  
for (int i=0; i<n; i++) {  
    out[i] = in[i] / precompute;  
}
```

$\Theta(n)$

Compiler Phase Ordering

Example:

Normalizing a Vector

```
for (int i=0; i<n; i++) {  
    out[i] = in[i] / mag(in, n);
```

}

Pass #1: -inline

```
for (int i=0; i<n; i++) {  
    double sum = 0;  
    for (int j=0; j<n; j++) {  
        sum += A[j] * A[j];  
    }  
    out[i] = in[i] / sqrt(sum);
```

Pass #2: -licm (loop-invariant code motion)

```
double sum;  
for (int i=0; i<n; i++) {  
    sum = 0;  
    for (int j=0; j<n; j++) {  
        sum += A[j] * A[j];  
    }  
    out[i] = in[i] / sqrt(sum);}
```

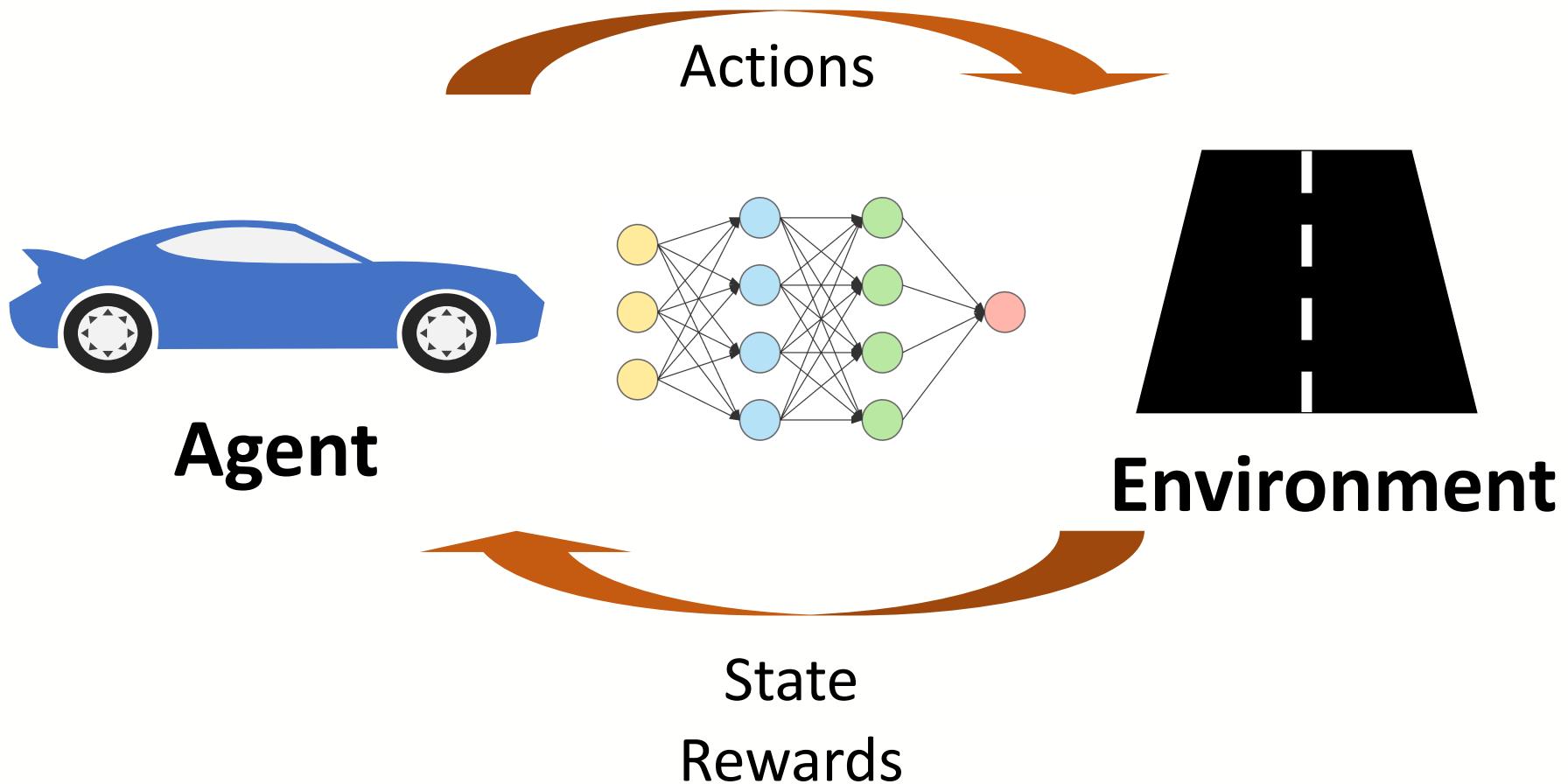


$\Theta(n)$



$\Theta(n^2)$

Deep Reinforcement Learning



Deep Reinforcement Learning

Decision(**Action**): Gas

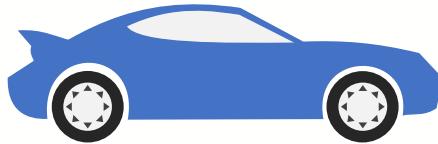


Decision(**Action**): Brake

Deep Reinforcement Learning

Decision (**Action**):

- Gas
- Brake
- Steering



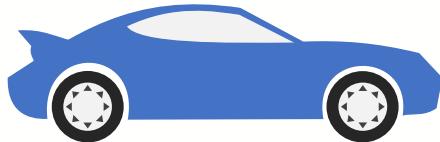
Fee (**Reward**): \$100

GPS Location (**State**):
37.87631055, -122.2388

Deep Reinforcement Learning

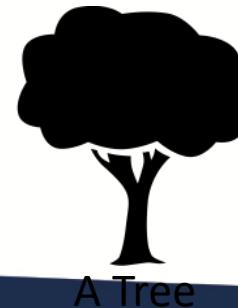
Decision (**Action**):

- Gas
- Brake
- Steering



Fee (**Reward**): \$100

GPS Location (**State**):
37.87631055, -122.2388



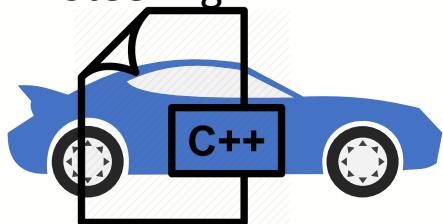
Fee (**Reward**): -\$500

GPS Location (**State**):
-30.43131384, 73.1693

Deep Reinforcement Learning

Decision (**Action**):

- Action:** Optimization Pass #
- Gas (e.g. -licm, -inline, -sroa, mem2reg)
- Brake
- Steering



Source Code



Fee (**Reward**): \$100

GPS Location (**State**):
37.87631055, -122.2388

Reward:

Cycle Count Improvement

State: a) Program Features
b) Applied Passes



Fee (**Reward**): -\$500

GPS Location (**State**):
-30.43131384, 73.1693

Generating the States

1. Program Features:

- 56 static features (# of instructions, # of loops, ...)
- An LLVM analysis pass is built

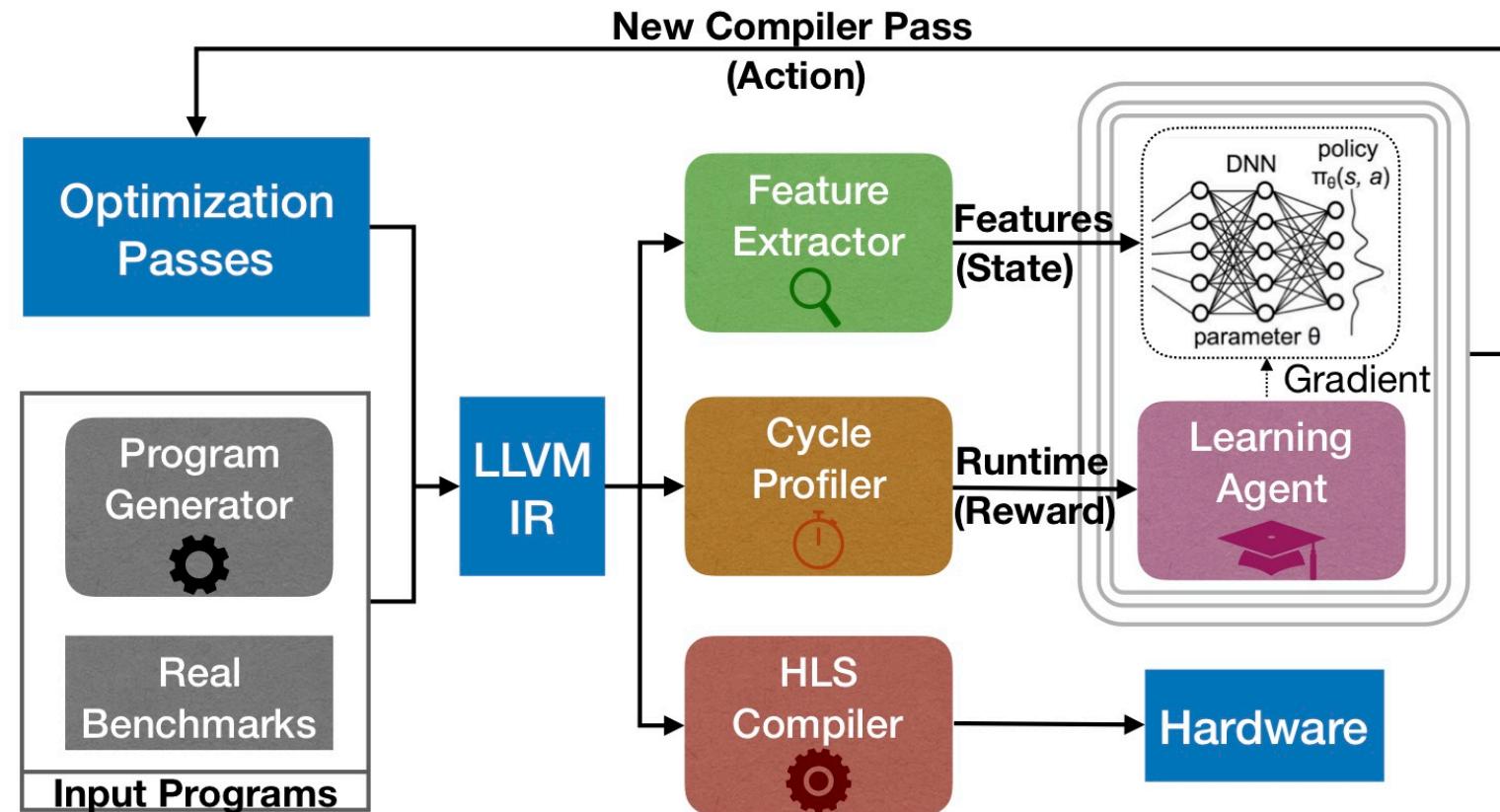
2. Histogram of Applied Passes:

- 45 optimization passes
- Once pass i is applied, $Histogram[i]++$
(**Histogram**: a vector of length 45, i : index of a pass)

Generating the Rewards

- Reward definition:
 - cycles before a pass is applied**
 - cycles after a pass is applied
- LegUp estimates the circuit cycles from C
 - Leveraging SW run information
 - 20× faster than RTL simulation
 - 0.48% error rate on our benchmarks

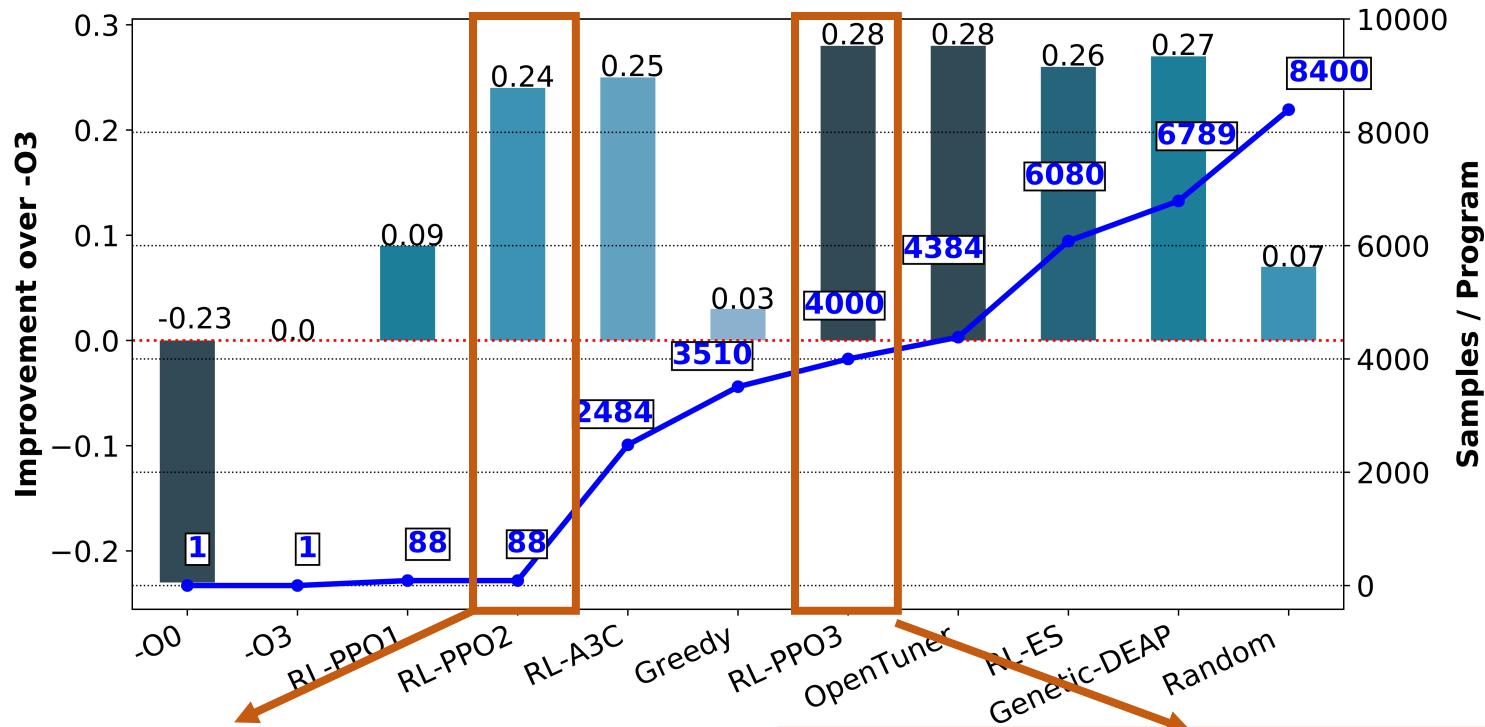
AutoPhase Framework



Evaluation

- Benchmarks:
 - 9 benchmarks from CHStone
- Algorithms:
 - Random Search
 - Genetic Algorithms
 - OpenTuner – Multiarmed Bandit/PSO
 - Greedy Algorithms/Beam Search
 - Deep Reinforcement Learning
 - Proximal Policy Optimization (PPO)
 - Asynchronous Actor-Critic Agents (A3C)
 - Evolutionary Strategies (ES)

Performance (Execution Time)



- Uses only **Histogram of Applied Passes**
- Compiles and runs **once** per trajectory
- Fewer samples required
- 28% improvement over -O3
- Uses **Program Features and Histogram of Applied Passes**

RL Generalization

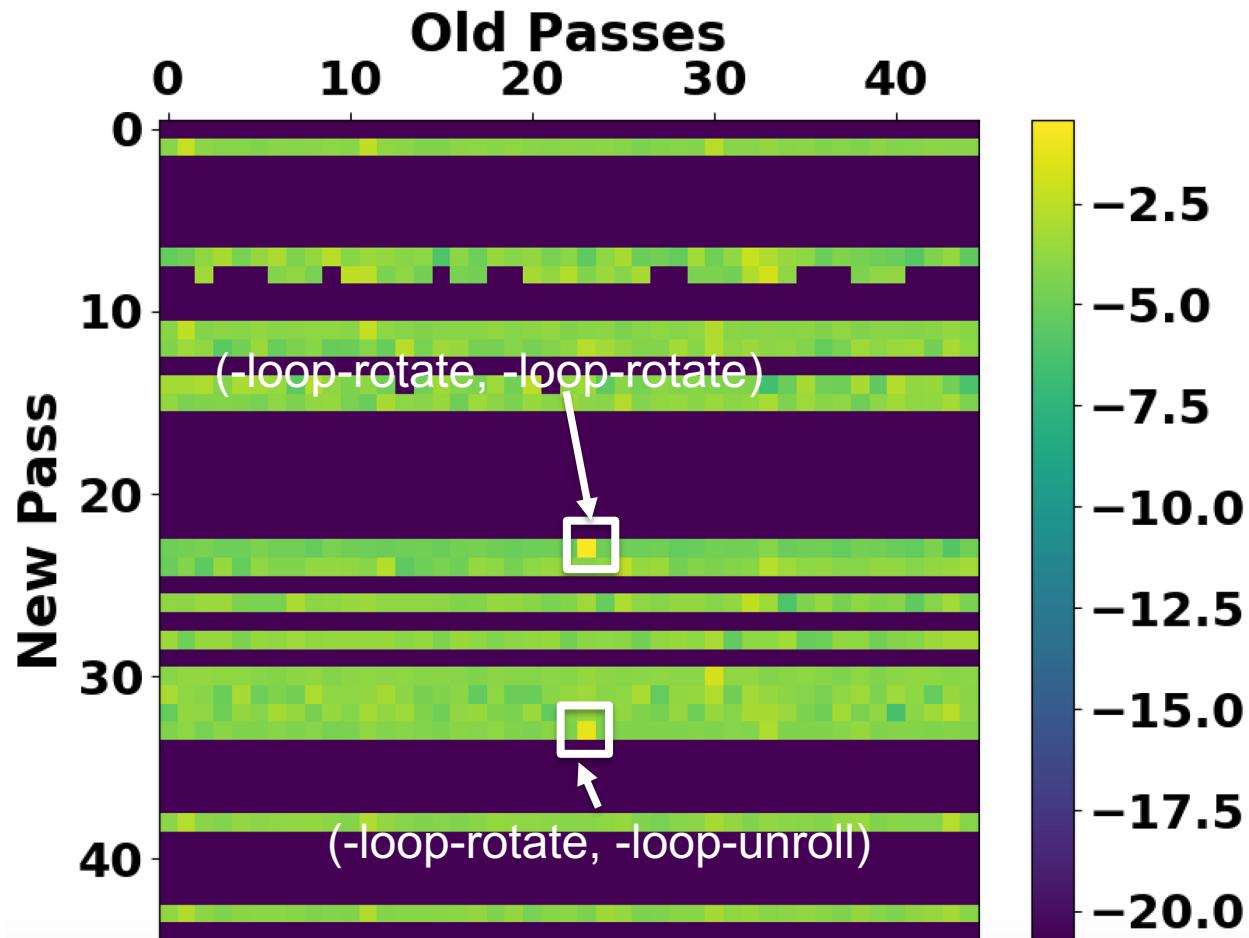
- Can we train a network and use it to predict the optimal sequence of any program in a single compilation?
 - More than 45^{45} possible phase orderings
 - Limited observations
 - Limited data

Very Difficult!

~15,000 Random HLS
Programs are Generated



Importance Analysis



Normalization Schemes

1. Logarithm of program features and rewards:

- Reduces the magnitude of features/rewards
- The neural network learns to correlate the products of features instead of a linear combination of them
 - $w_1 \log(o_{f_1}) + w_2 \log(o_{f_2}) = \log(o_{f_1}^{w_1} \cdot o_{f_2}^{w_2})$

2. Dividing the program feature values by the total number of instructions

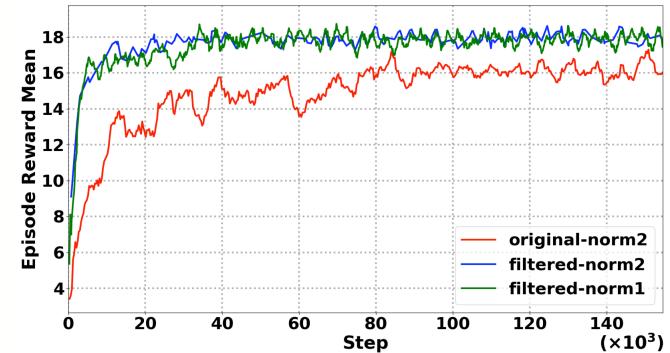
Generalization Results

Train on:

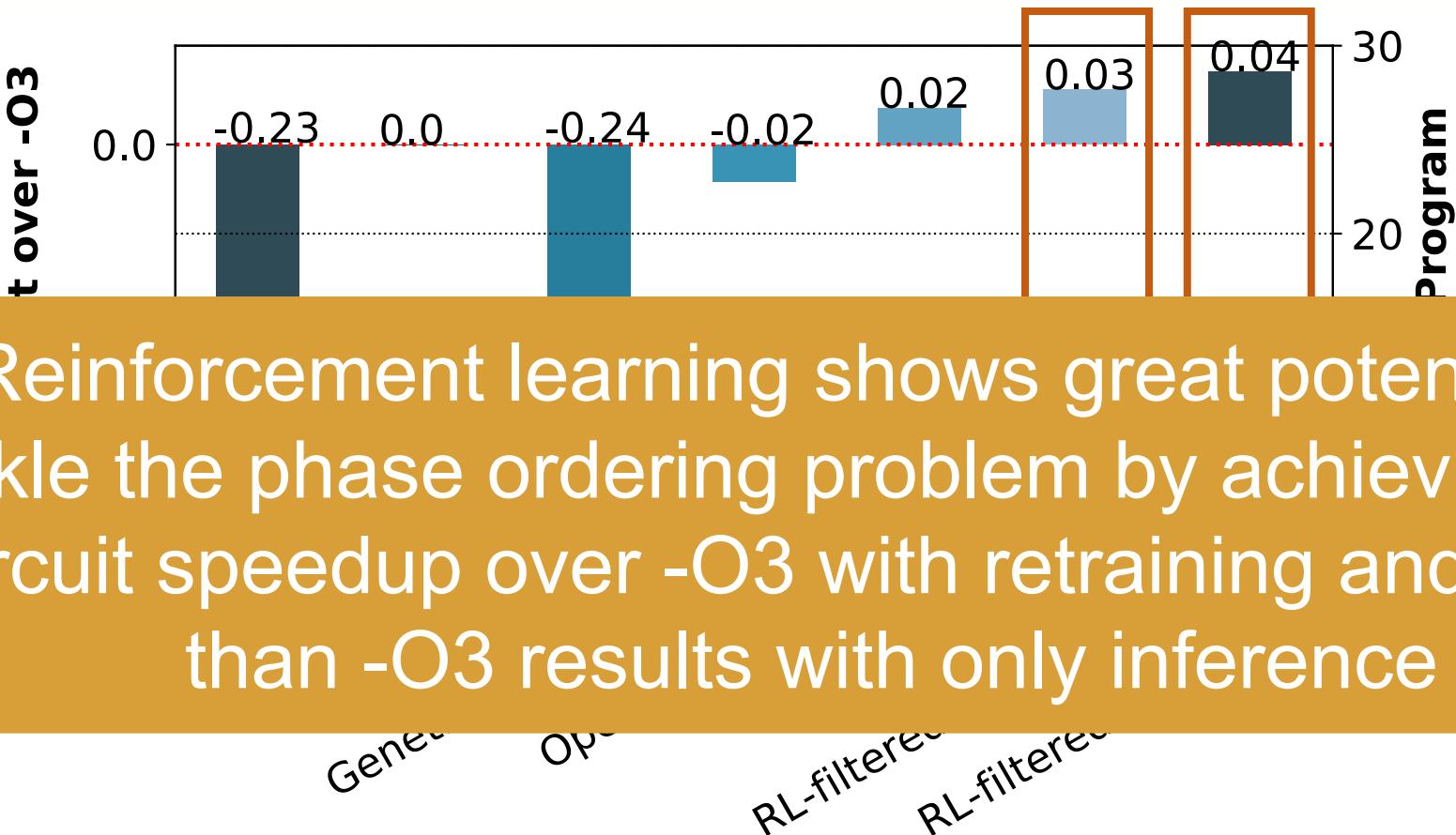
- 100 random programs

Inference-only on:

- 13,000 randomly-generated programs
 - 6% better than –O3
- CHStone benchmarks
 - different from the randomly-generated programs



Generalization Results



Reinforcement learning shows great potential to tackle the phase ordering problem by achieving 28% circuit speedup over -O3 with retraining and better than -O3 results with only inference

- ~4% improvement over -O3 on CHStone

Outline

- Motivation
- Codesign for Image Classification
 - Shift-based DiracDeltaNet
- Codesign for Object Detection
 - Deformable Convolution
- Reinforcement Learning for HLS
 - Compiler Phase-Ordering
- Conclusion

Conclusion

- *Algorithm-hardware co-design* can achieve both high accuracy and good efficiency for embedded CV applications
- *Reinforcement learning* is a promising approach for addressing the hardware design problem with intractable search space

Questions?

Email: qijing.huang@berkeley.edu

Special Thanks

Fellow Researchers:

- Yifan Yang
- Bichen Wu
- Tianjun Zhang
- Liang Ma
- Giulio Gambardella
- Michaela Blott
- Luciano Lavagno
- Kees Vissers
- Prof. Kurt Keutzer
- Dequan Wang
- Zhen Dong
- Yizhao Gao
- Ameer Haj-Ali
- William Moses
- John Xiang
- Prof. Krste Asanovic
- Prof. Ion Stoica
- Prof. John Wawrzynek

Thanks!