

# **Hardware/Software Co-Exploration of Neural Architecture**

—— Published and Ongoing Works

Weiwen Jiang

wjiang2@nd.edu

*The College of Engineering  
at the University of Notre Dame*

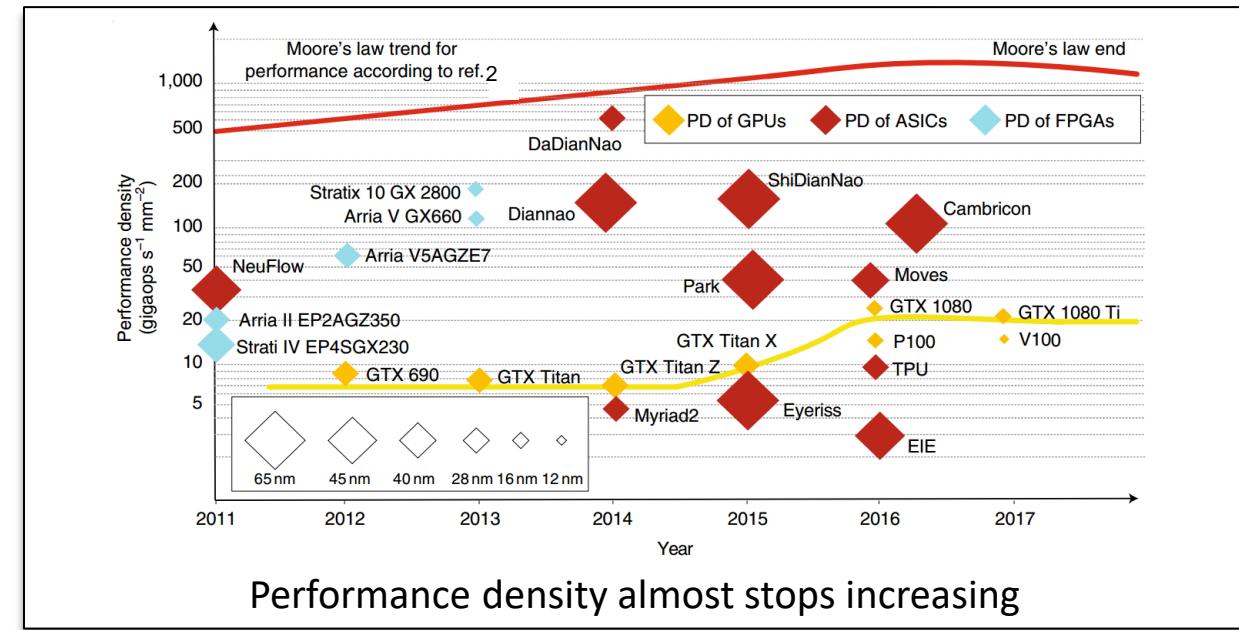
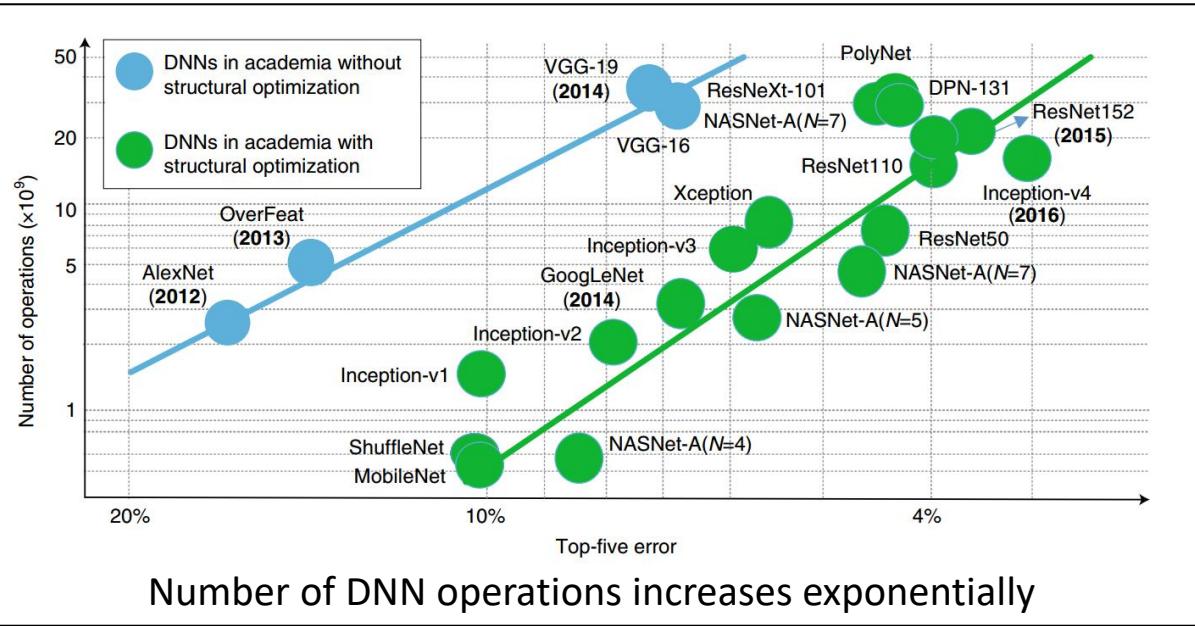


# Outline

- **Introduction**
- **Co-Exploration of Neural Architecture and FPGA Implementation**
  - DAC'19 (Best Paper Nomination)
- **Other Research Directions**
- **Conclusion**



# Computation Gap

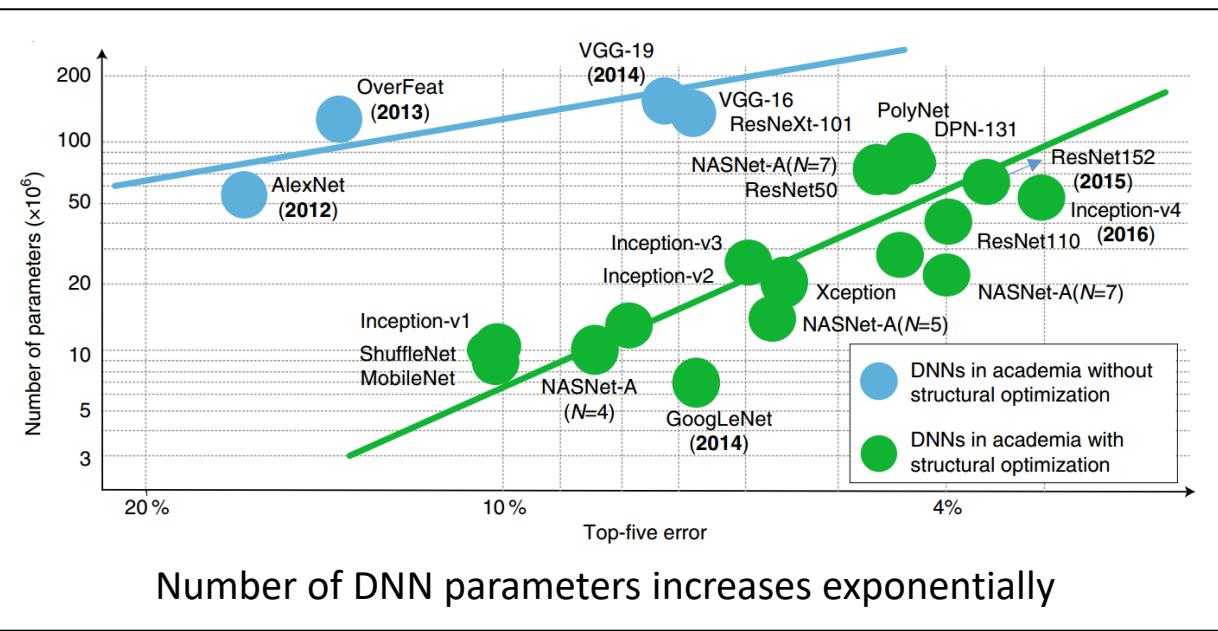


[ref.1] Xu, Xiaowei, et al. "Scaling for edge inference of deep neural networks." *Nature Electronics* 1.4 (2018): 216.

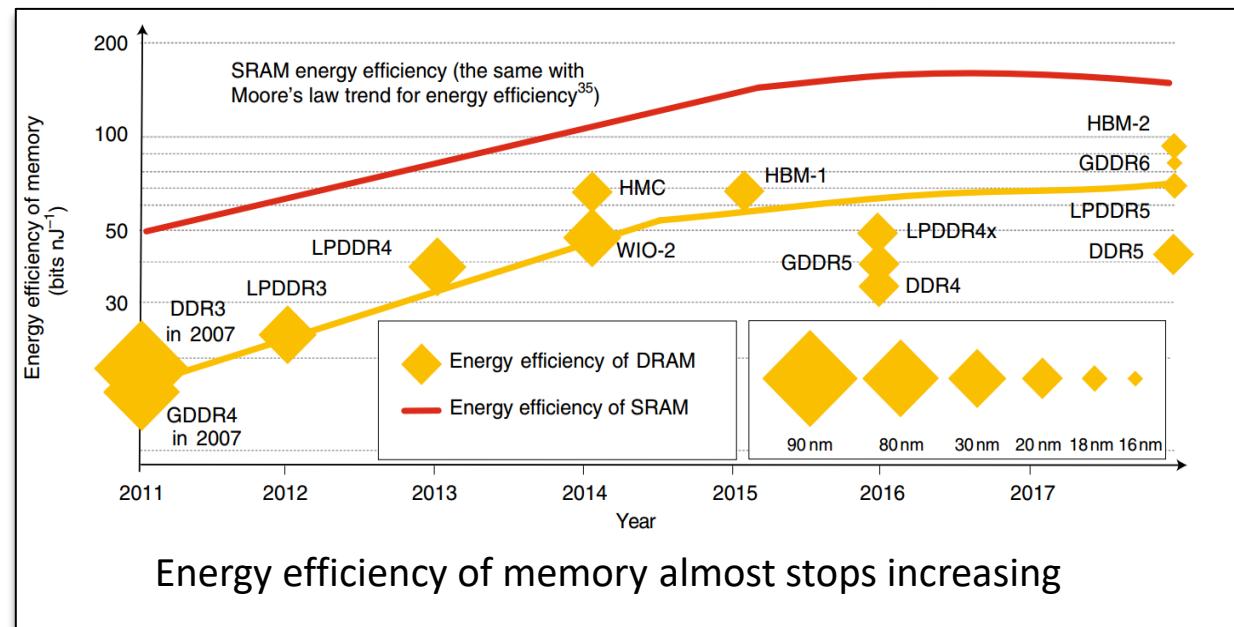
[ref.2] Sutter, H. The free lunch is over: A fundamental turn toward concurrency in software. *Dr Dobb's J.* 30, 202–210.



# Storage Gap



Number of DNN parameters increases exponentially



Energy efficiency of memory almost stops increasing

[ref.1] Xu, Xiaowei, et al. "Scaling for edge inference of deep neural networks." *Nature Electronics* 1.4 (2018): 216.

[ref.2] Sutter, H. The free lunch is over: A fundamental turn toward concurrency in software. *Dr Dobb's J.* 30, 202–210.

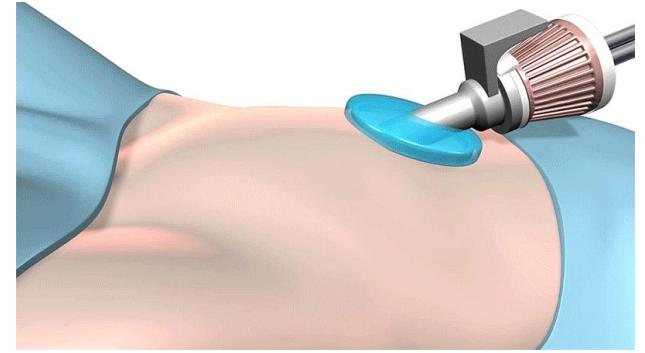
# Real-Time Requirement



Traffic Surveillance



Self-Driving Vehicle



Robots in Surgery

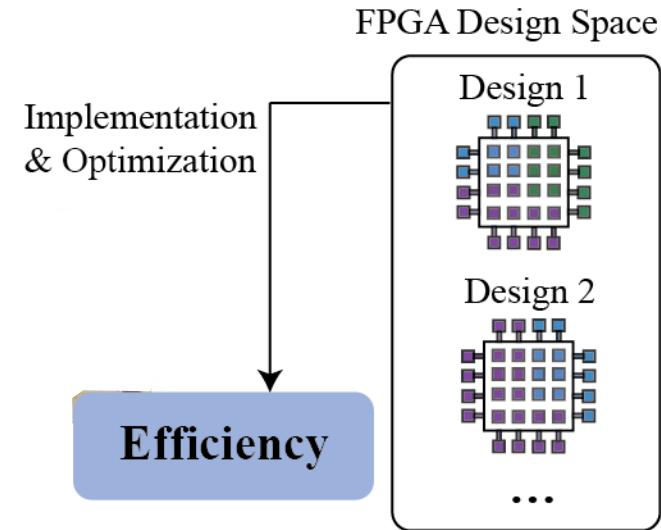
Increasing Demands on Real-Time Machine Learning Applications

[ref.1] Google's Waymo

[ref.2] da Vinci Surgical Robot at St. Francis Health Center - Robotic Gizmos

# Efforts on Both Software and Hardware Sides are Demanded

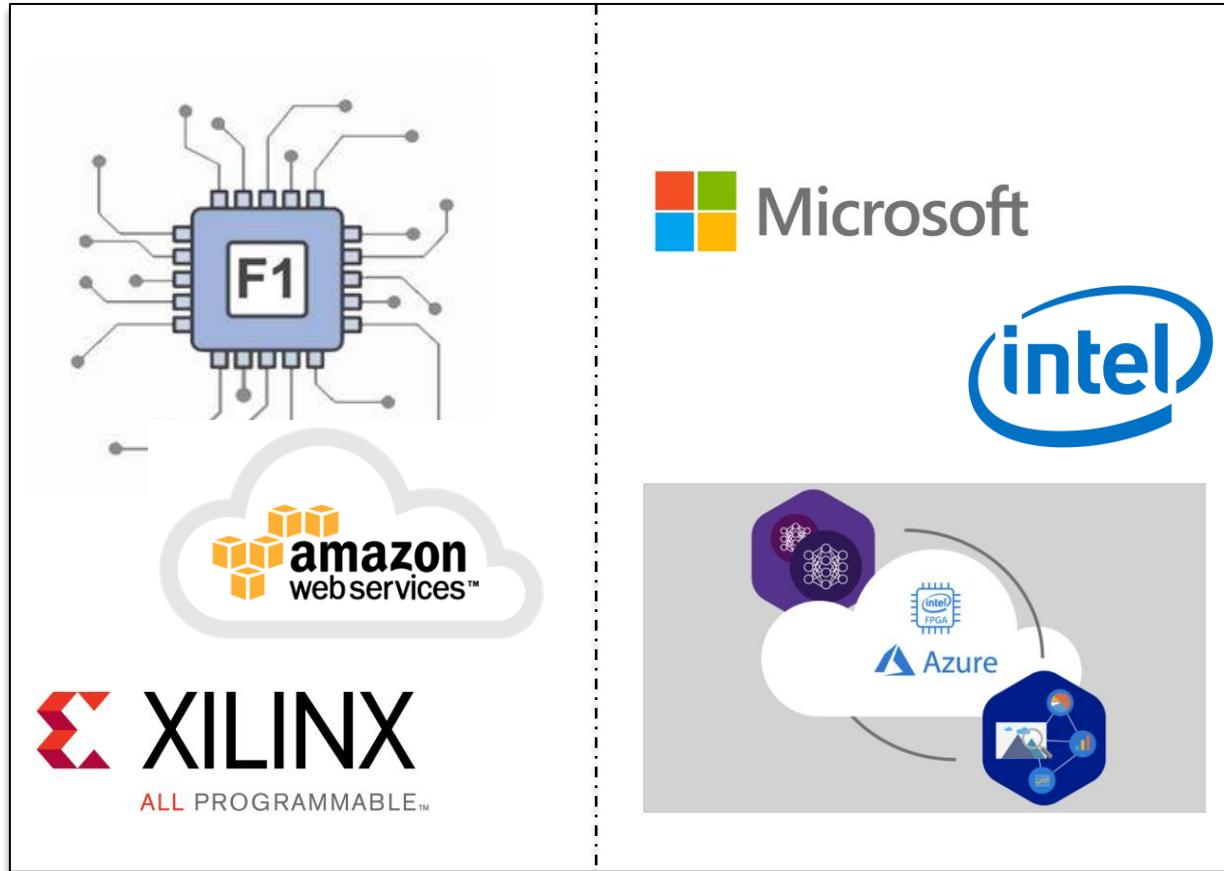
## — Hardware Implementation and Optimization



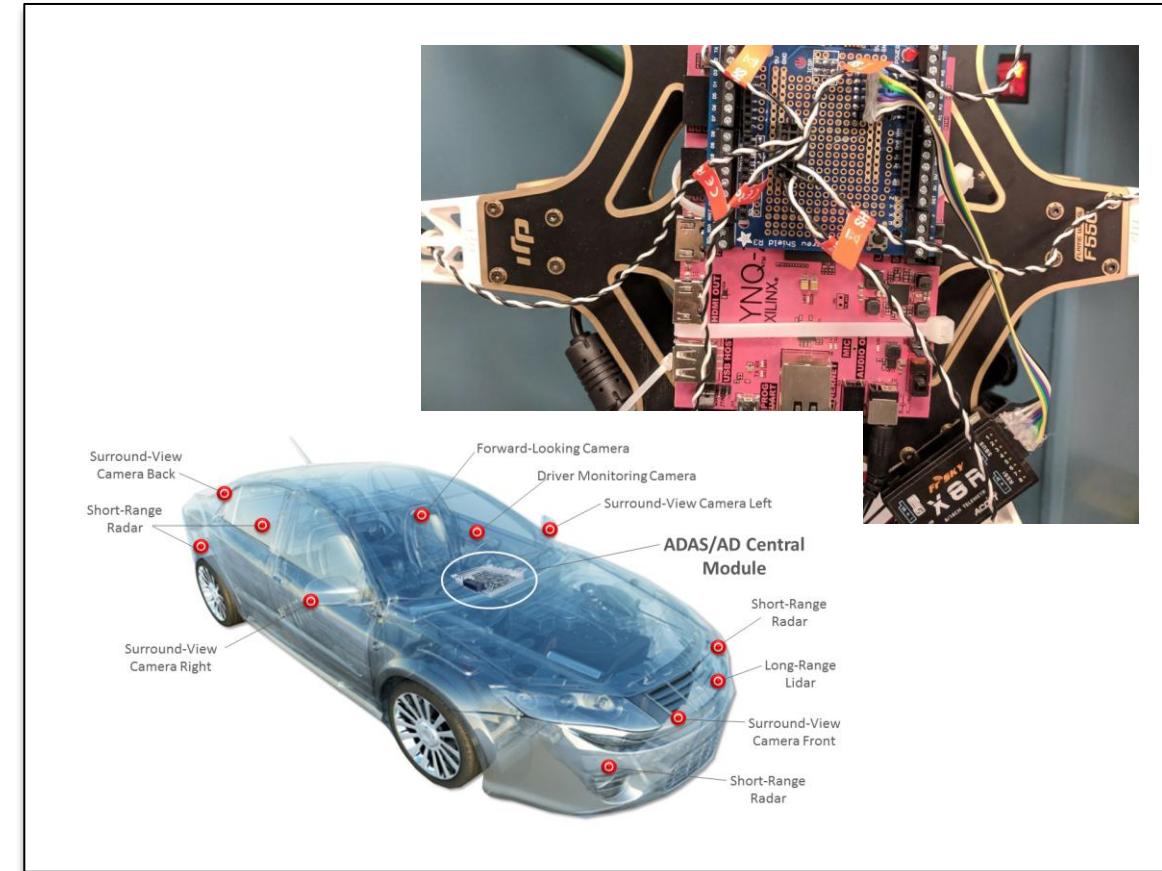
Neural Architecture Implementation on Hardware

# FPGAs in DNN Applications

## FPGA in Cloud Computing



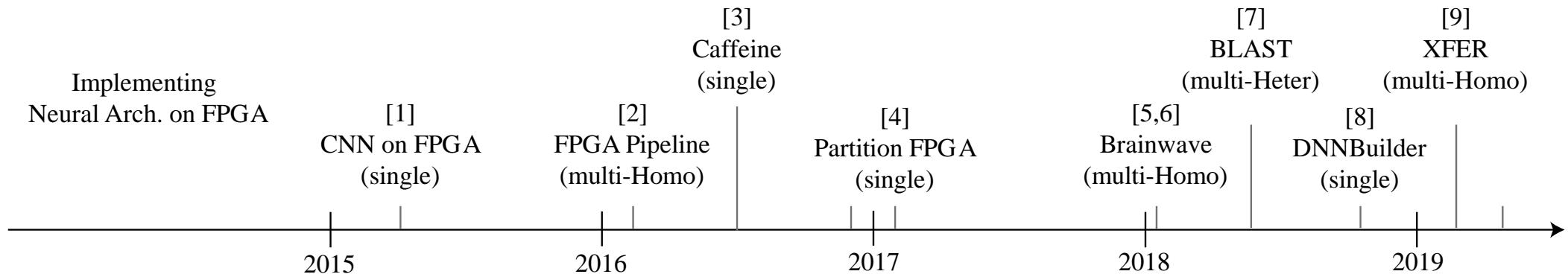
## FPGA in Edge Computing



[ref] Where Do FPGAs Stand in Auto IC Race? [https://www.eetimes.com/document.asp?doc\\_id=1333419#](https://www.eetimes.com/document.asp?doc_id=1333419#)

[ref] PYNQ on UAV <http://brennancain.com/pynqcopter-an-open-source-fpga-overlay-for-uavs/>

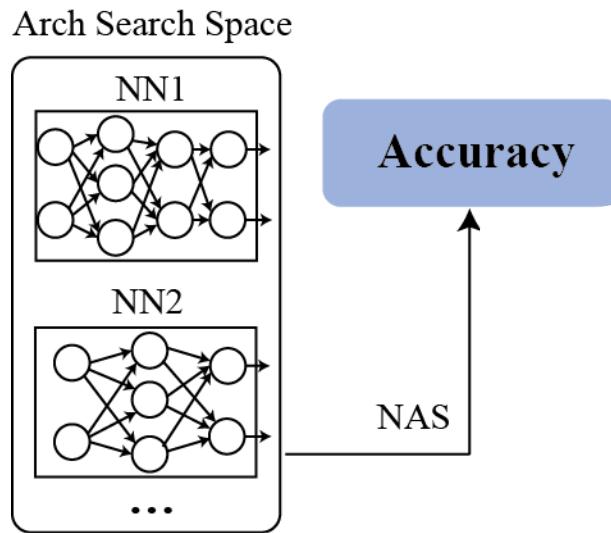
# Implementing Fixed NNs on FPGAs



- [1] Chen Zhang, et al. Optimizing fpga-based accelerator design for deep convolutional neural networks. In *FPGA*, 2015.
- [2] Chen Zhang, et al. Energy-efficient cnn implementation on a deeply pipelined fpga cluster. In *ISLPED*, 2016.
- [3] Chen Zhang, et al. Caffeine: Towards uniformed representation and acceleration for deep convolutional neural networks. *ICCAD*, 2016.
- [4] Yongming Shen, et al. Maximizing cnn accelerator efficiency through resource partitioning. In *ISCA*, 2017.
- [5] Eric Chung, et al. Serving dnns in real time at datacenter scale with project brainwave. *IEEE Micro*, 2018.
- [6] Jeremy Fowers, et al. A configurable cloud-scale dnn processor for real-time ai. In *ISCA*, 2018.
- [7] Weiwen Jiang, et al. Heterogeneous fpga-based cost-optimal design for timing-constrained cnns. *CASES – TCAD*, 2018.
- [8] Xiaofan Zhang, et al. Dnnbuilder: an automated tool for building high-performance dnn hardware accelerators for fpgas. In *ICCAD*, 2018.
- [9] Weiwen Jiang, et al. Xfer: A novel design to achieve super-linear performance on multiple fpgas for real-time ai. In *FPGA* 2019.

# Efforts on Both Software and Hardware Sides are Demanded

## — Neural Architecture Design



Neural Architecture Search

re

8

*The College of Engineering  
at the University of Notre Dame*



# Evolution of Exploring Deep Neural Architectures

## III. Hardware-Aware Neural Architecture Search

## II. Pure Automatically Explore Neural Architecture Space

## I. Human Invented

2012

AlexNet

2013

ZFNet

2014

VGGNet

2015

GoogLeNet

2018.July

2018.Nov

2018.Dec

2019

MnasNet  
from Google

proxylessNAS  
from MIT

FBNet  
from UCB &  
Facebook

Neural Architecture  
Search (RL-based)

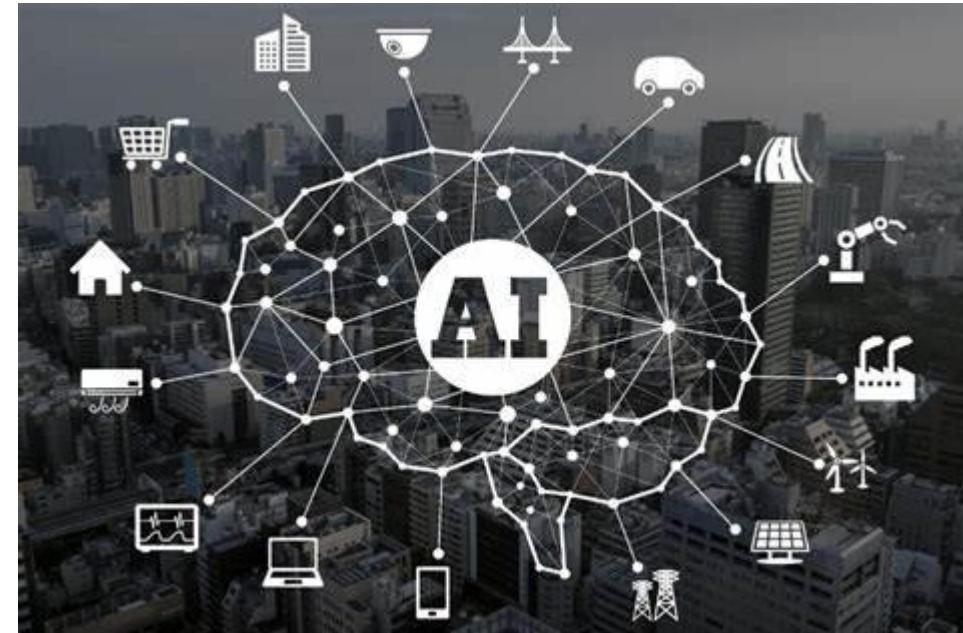
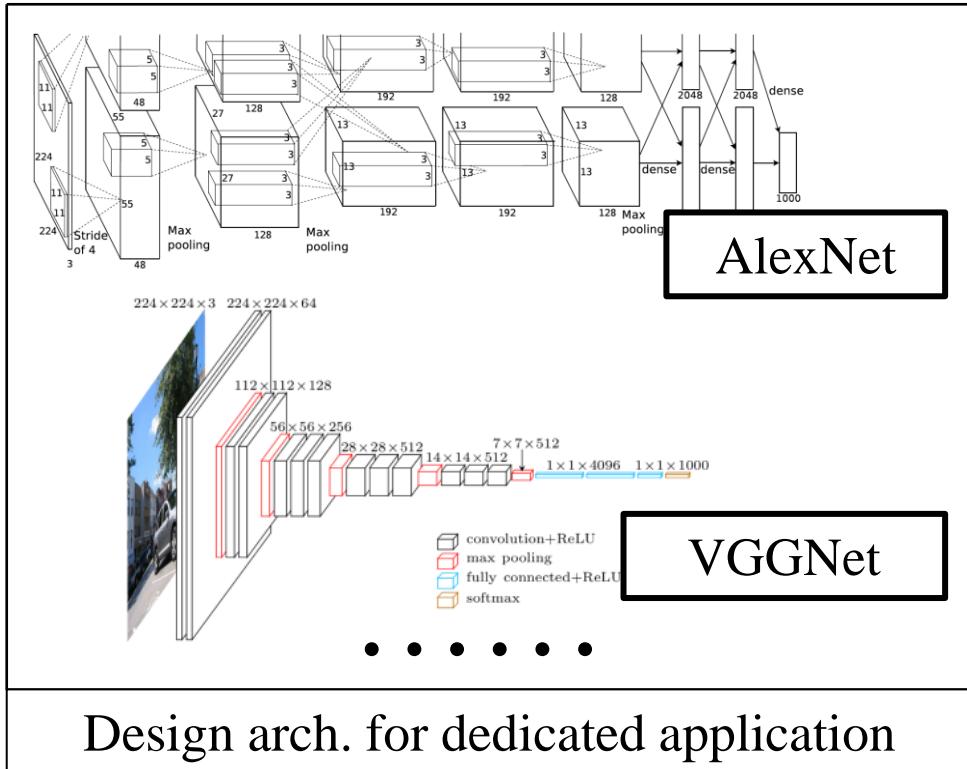
2017

2018

2019

2020

# Phase I: Human Invented

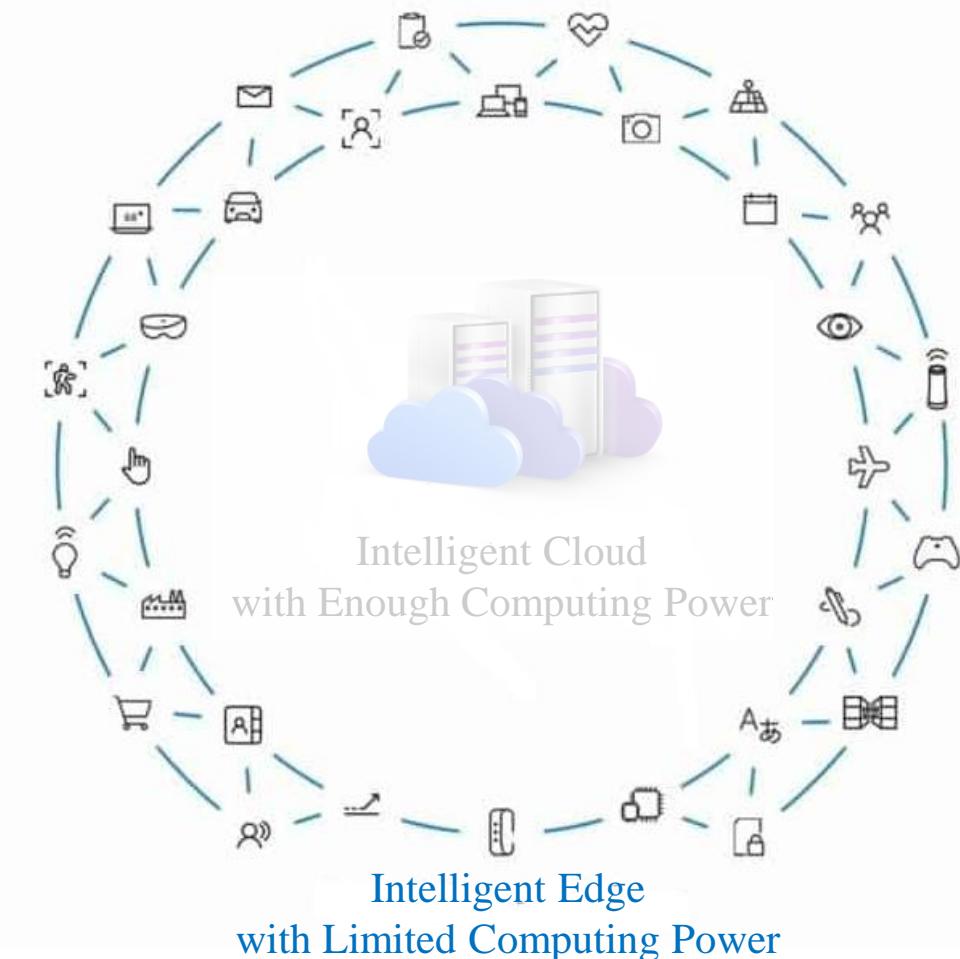
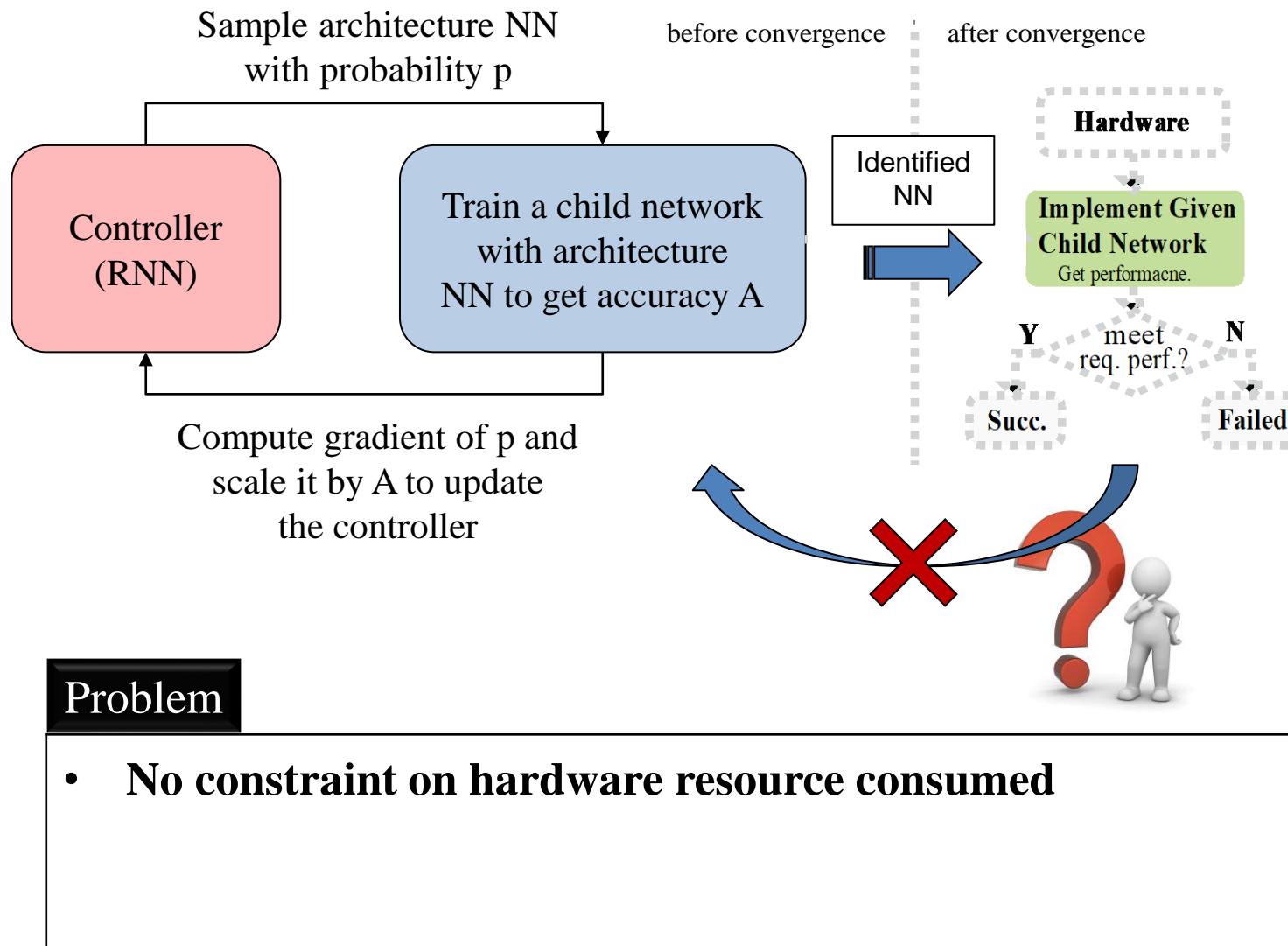


Era of AI Democratization

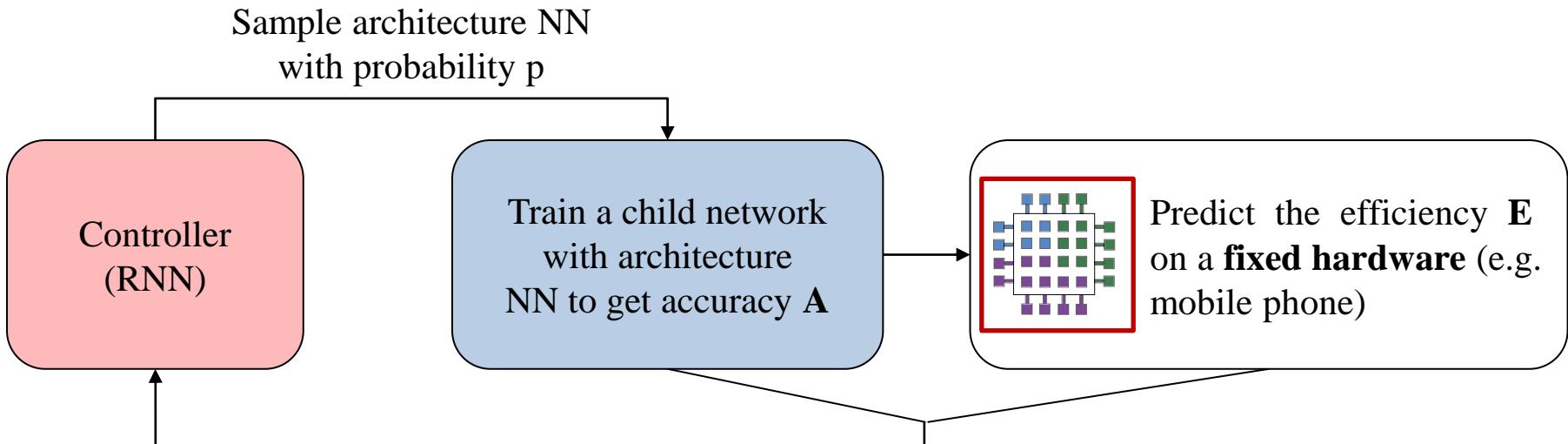
## Problem

- Domain knowledge and excessive labor
- It is impossible to manually design specific arch. for dedicated application in the era of AI democratization

## Phase II: Neural Architecture Search (NAS)



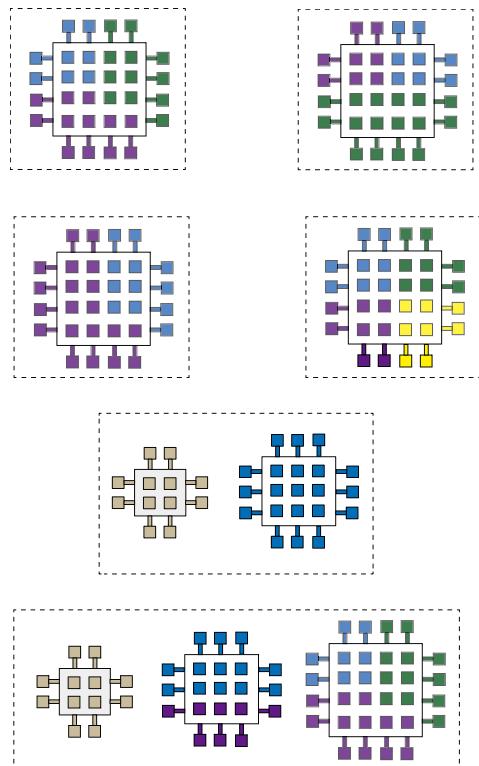
# Phase III. Hardware-Aware NAS



## Problem

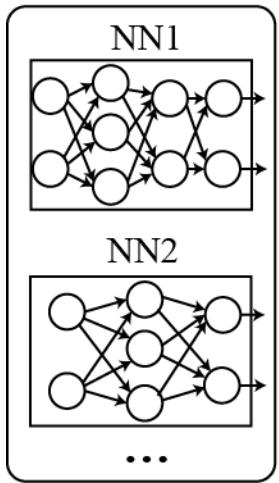
- **It works for particular fixed hardware, but not suitable for programmable hardware**

## Different Hardware Designs



# A Missing Link between Two Design Spaces

Arch Search Space

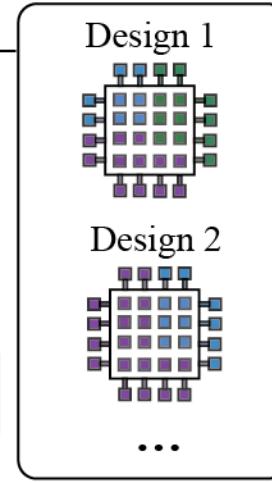


Accuracy

NAS

Neural Architecture Search

FPGA Design Space



Implementation & Optimization

Efficiency

Neural Architecture Implementation on Hardware

# Evolution of Exploring Deep Neural Architectures

## IV. Co-Explore Neural Architecture Space and Hardware Design Space

## III. Hardware-Aware Neural Architecture Search

## II. Pure Automatically Explore Neural Architecture Space

## I. Human Invented

2012

2013

2014

2015

2016

2017

2018

2019

2020

AlexNet

ZFNet

VGGNet

GoogLeNet

MnasNet  
from Google

proxylessNAS  
from MIT

FBNet  
from UCB & Facebook

DAC19 FNAS  
from ND

2018.July

2018.Nov

2018.Dec

2019

Neural Architecture  
Search (RL-based)



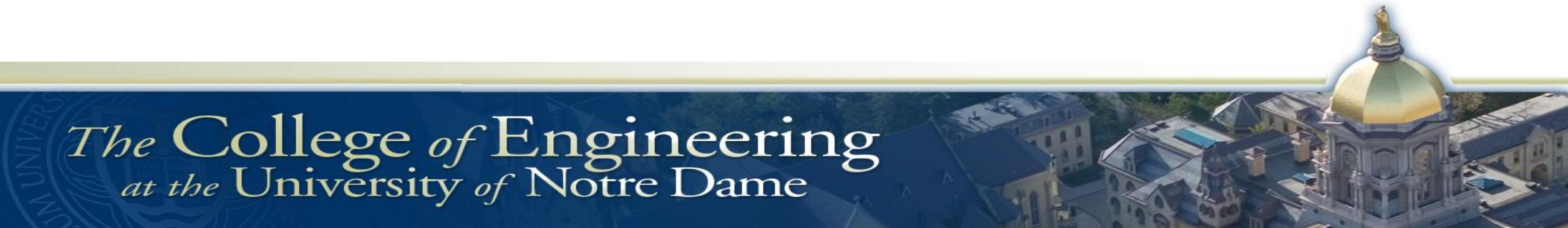
# Outline

- Introduction
- **Co-Exploration of Neural Architecture and FPGA Implementation**
  - DAC 2019 (Best Paper Nomination)
- Other Research Directions
- Conclusion

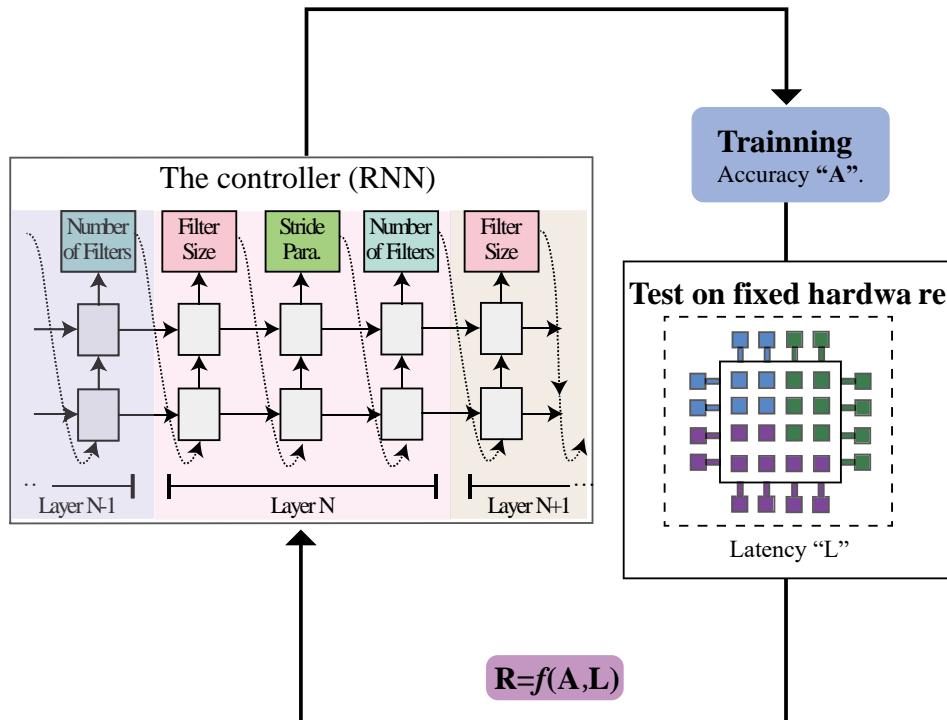


# **FNAS Framework**

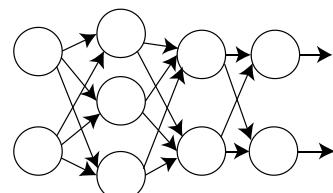
*The College of Engineering  
at the University of Notre Dame*



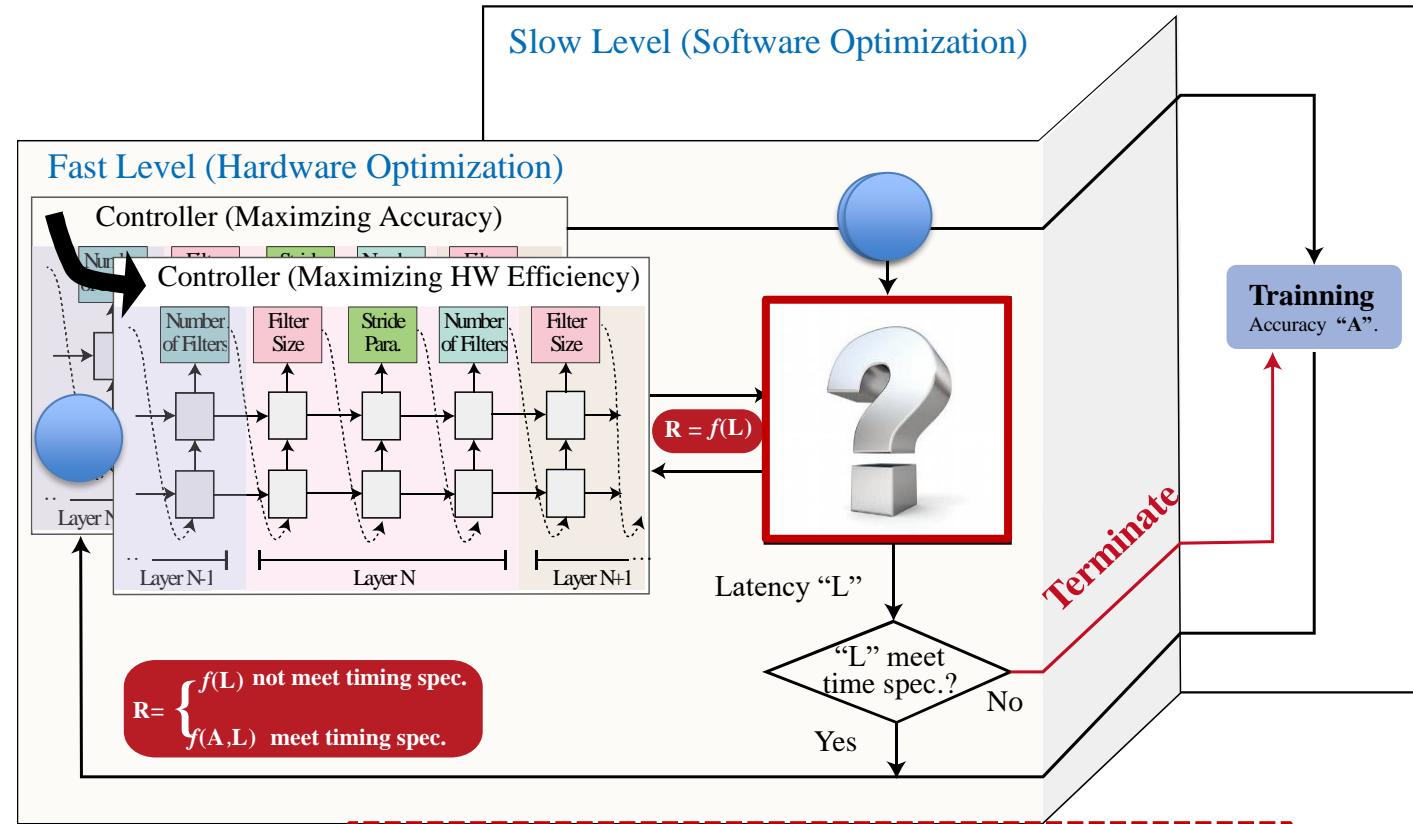
# HW-Aware NAS vs. FPGA-Implementation Aware NAS (FNAS)



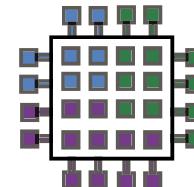
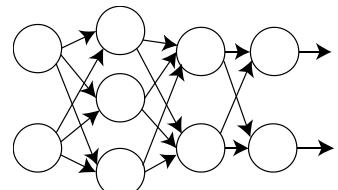
**Output:** A neural architecture



HW-Aware NAS

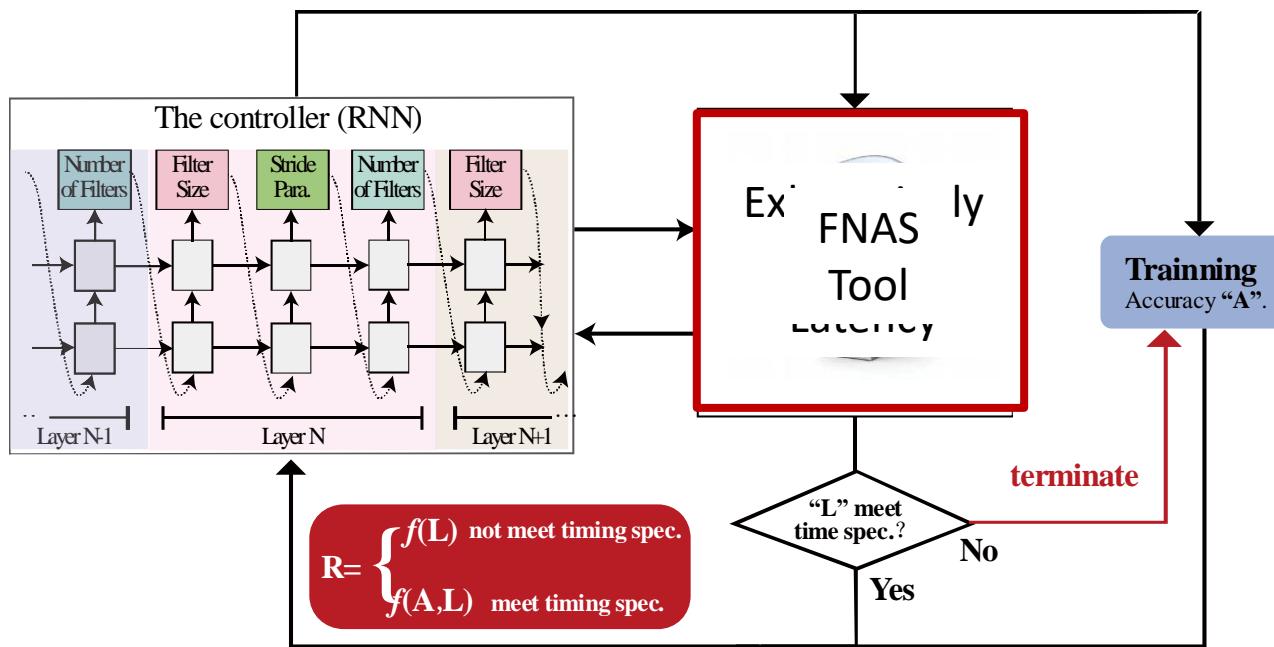


**Output:** A pair of neural architecture and hardware design



FNAS

# Solutions & Challenges



## Our Solution: FNAS tools to response to challenges

### FNAS-Design

“Design on Program Logic”

**C1**

### FNAS-GG

“Tile-based Task Graph Generator”

**C2**

### FNAS-Sched

“Scheduler on Processing System”

**C2**

### FNAS-Analyzer

Estimate Performance “ L ”

**C3**

## Naïve Solution: HW-Aware + Exhaustively Evaluate Lat.

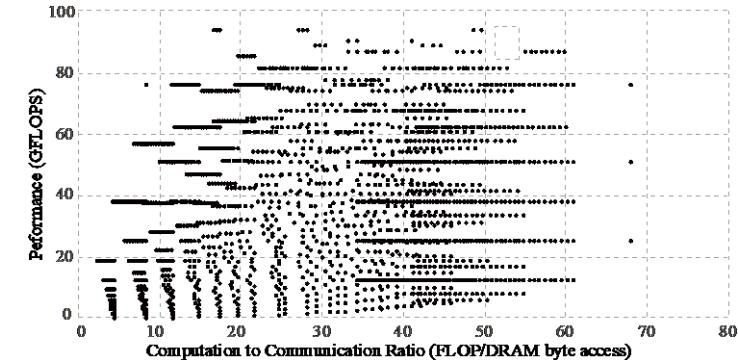


Fig1. Possible designs for Layer 5 of AlexNet on ZCU102

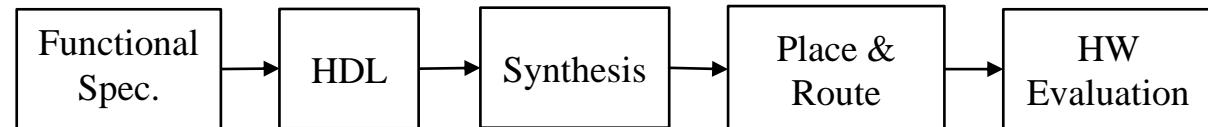


Fig2. Procedure of performance evaluation

## Challenges:

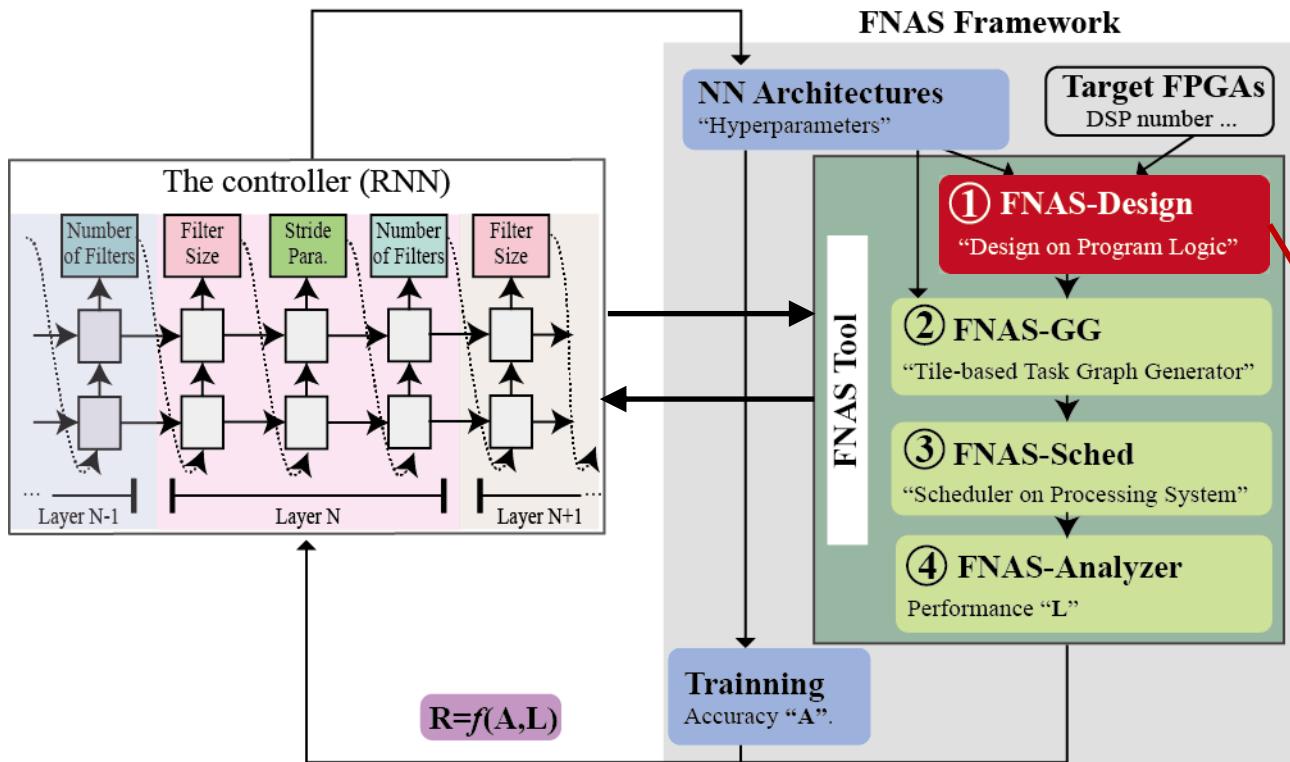
**C1: Huge design space!**

**C2: Multi-FPGA design!**

**C3: Time-consuming evaluation!**

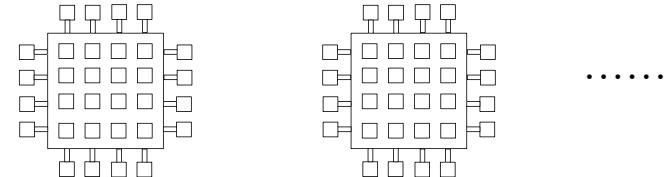
**Infeasible**

# FNAS: Design Optimization (on-chip design)



*Given:*

1. FPGAs with attributes including LUTs, DSPs, BRAM, etc.



2. A neural architecture with determined hyperparameters

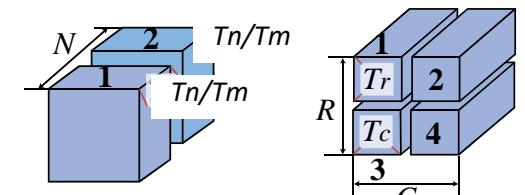
## On-chip accelerator design:

*Determine:*

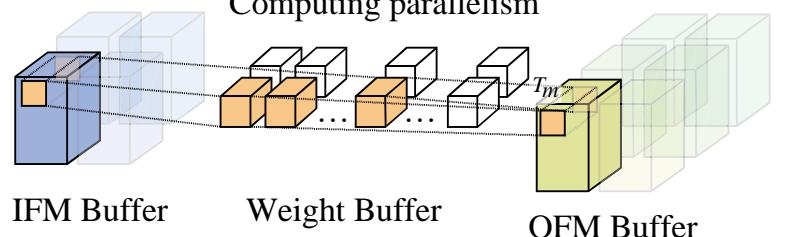
1. On-chip buffer allocation; 2. Accelerator size for computing

(note: both are determined by tiling parameters,  $T_m$ ,  $T_n$ ,  $T_r$ ,  $T_c$ )

One layer:

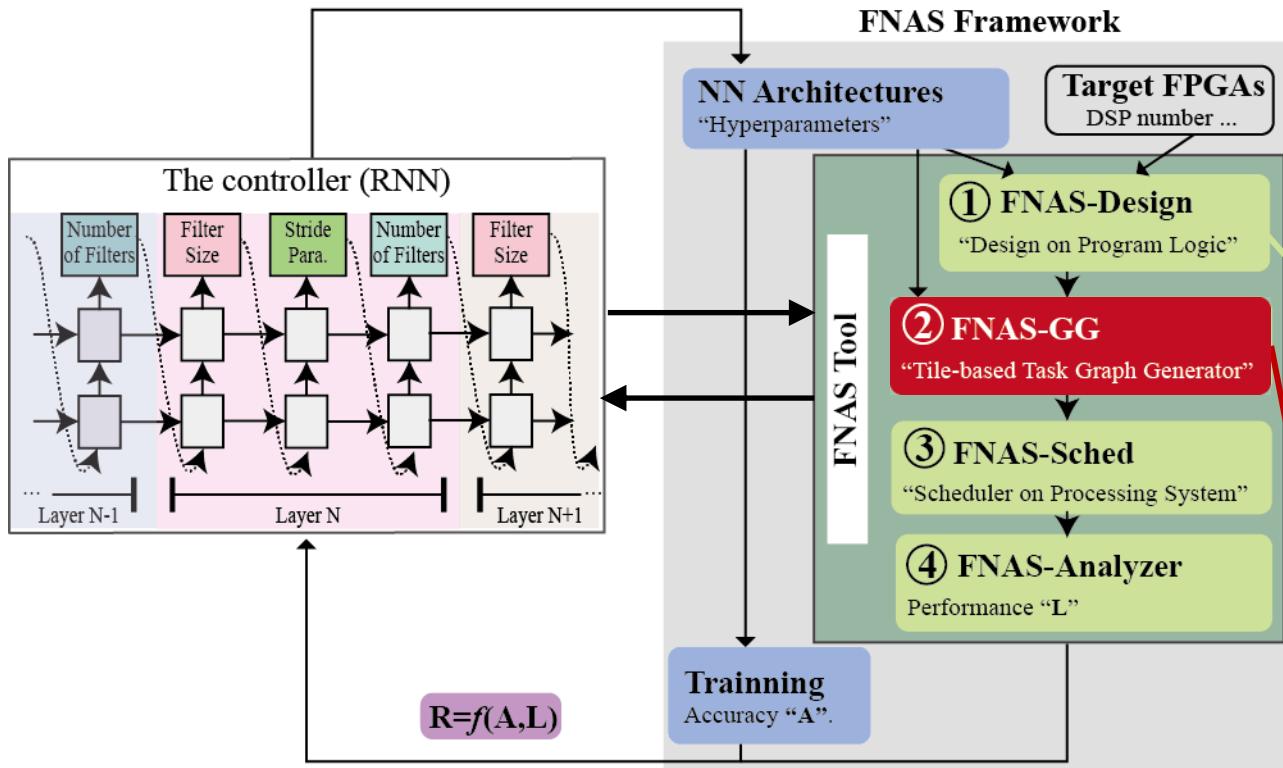


Computing parallelism



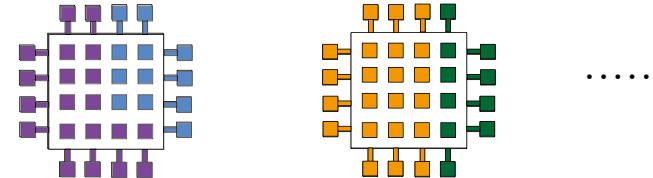
Multiple layers:

# FNAS: Graph Generator

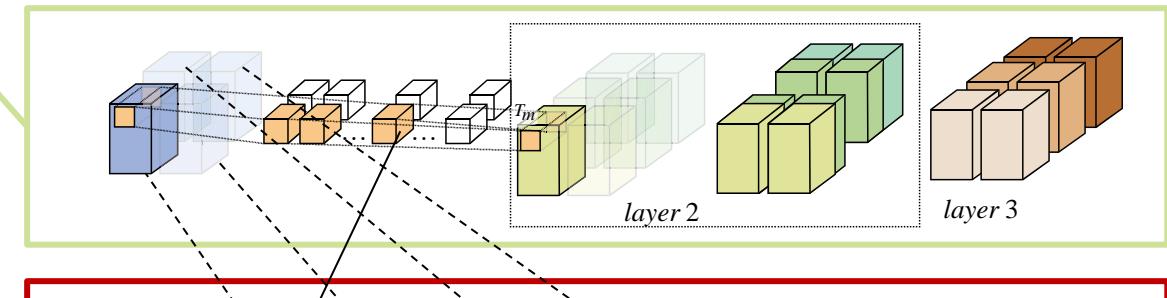


*Given :*

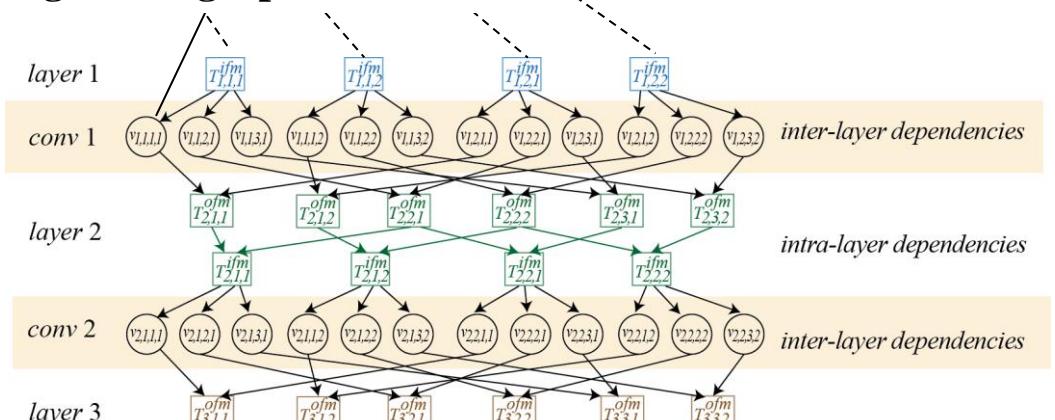
1. FPGAs with attributes including LUTs, DSPs, BRAM, etc.



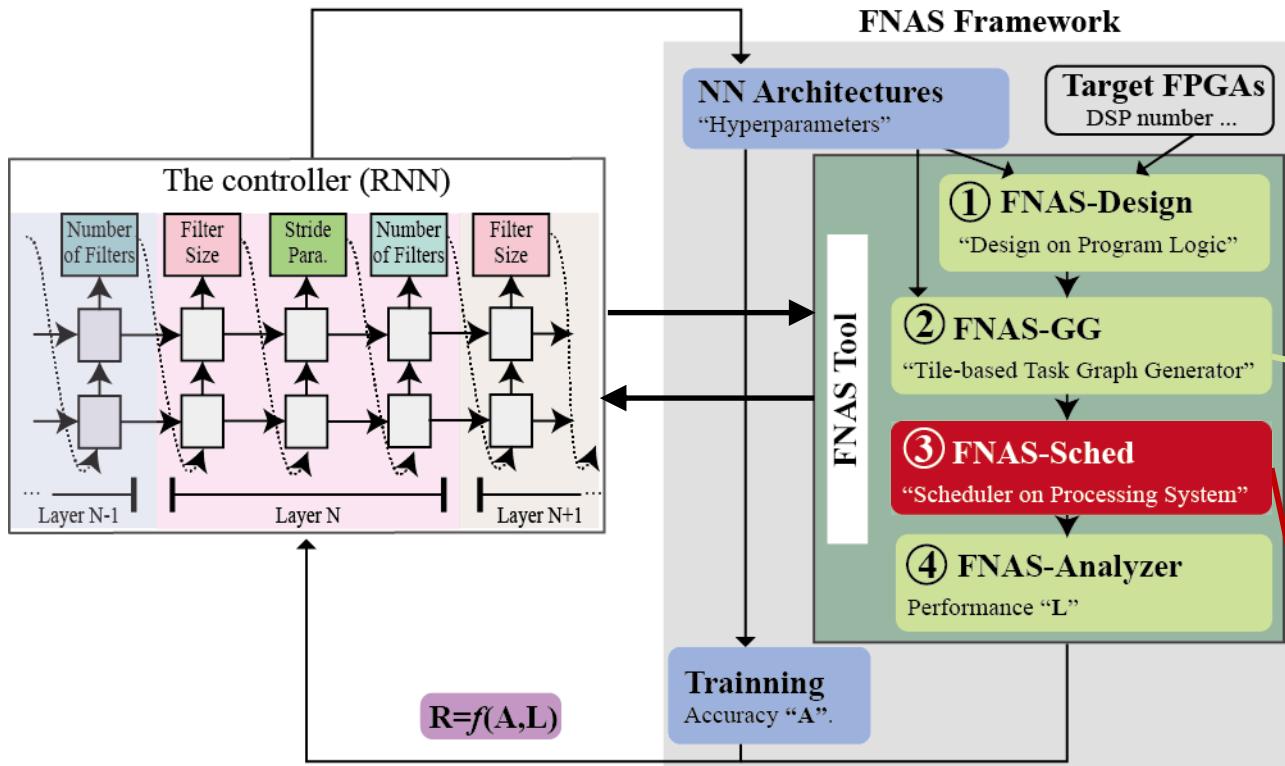
2. A neural architecture with determined hyperparameters



## High-level graph abstraction

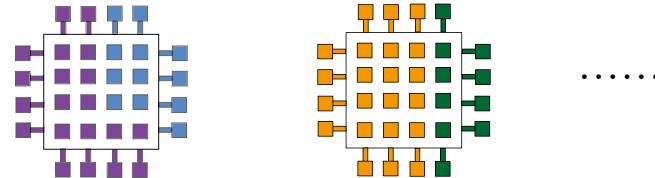


# FNAS: Schedule (off-chip design)

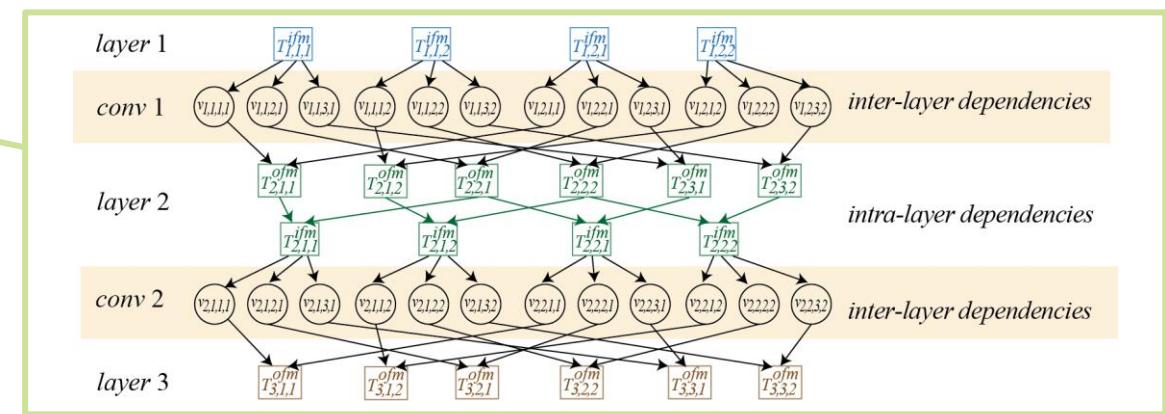


*Given :*

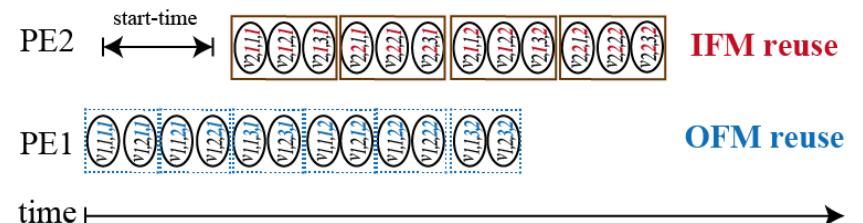
1. FPGAs with attributes including LUTs, DSPs, BRAM, etc.



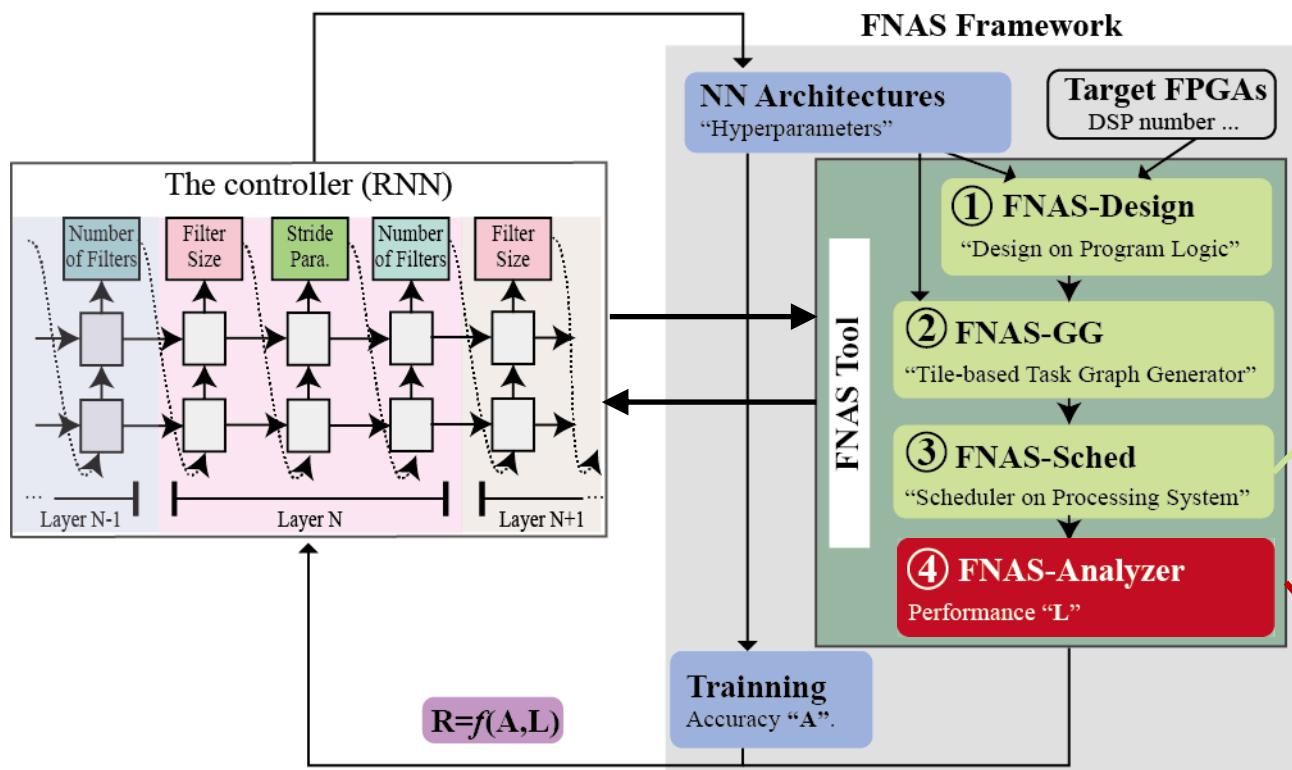
2. A neural architecture with determined hyperparameters



**Schedule of tasks in graph on multiple FPGAs**

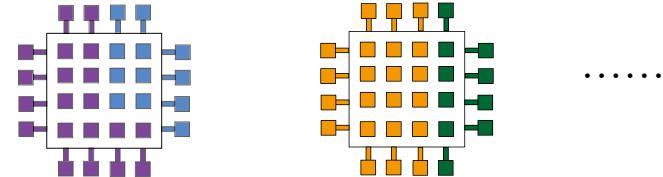


# FNAS: Analyzer



*Given :*

1. FPGAs with attributes including LUTs, DSPs, BRAM, etc.



2. A neural architecture with determined hyperparameters



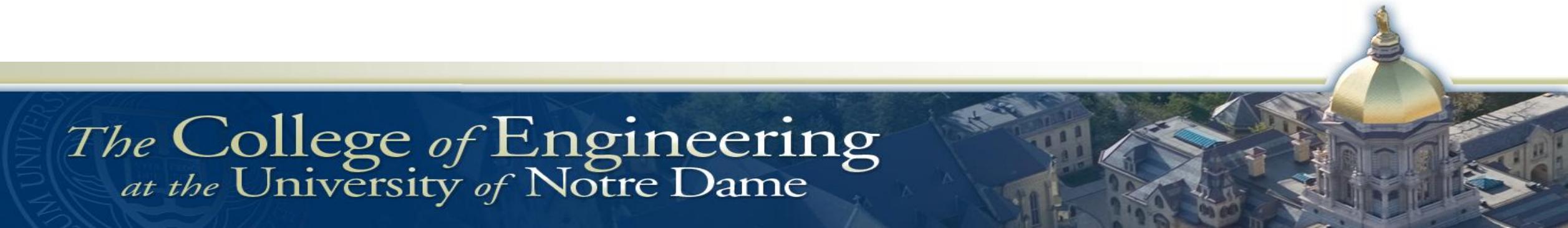
$$\text{Latency} = \text{pipeline start time} + \text{processing time}$$

*Output :*

1. A tailored FPGA Design
2. The system latency

# Experimental Results

*The College of Engineering  
at the University of Notre Dame*



# Experimental Setting

FPGAs



Xilinx 7A50T



Xilinx 7Z020

Datasets

0 0 0 0 0 0 0 0 0 0 0 0 0 0 0  
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
2 2 2 2 2 2 2 2 2 2 2 2 2 2 2  
3 3 3 3 3 3 3 3 3 3 3 3 3 3 3  
4 4 4 4 4 4 4 4 4 4 4 4 4 4 4  
5 5 5 5 5 5 5 5 5 5 5 5 5 5 5  
6 6 6 6 6 6 6 6 6 6 6 6 6 6 6  
7 7 7 7 7 7 7 7 7 7 7 7 7 7 7  
8 8 8 8 8 8 8 8 8 8 8 8 8 8 8  
9 9 9 9 9 9 9 9 9 9 9 9 9 9 9

MNIST



CIFAR-10



ImageNet

NAS Search Space      Layer Num.

up to 5

up to 10

up to 15

NAS Search Space      Filter Size

[5, 7, 14]

[1, 3, 5, 7]

[1, 3, 5, 7]

NAS Search Space      Filter Num.

[9, 18, 36]

[24, 36, 48, 64]

[16, 32, 64, 128]

HW Search Space      Channel Tiling Para. (Tm,Tn); Row Tiling Para. (Tr); Col Tiling Para. (Tc); Schedule

Timing Spec. (ms)

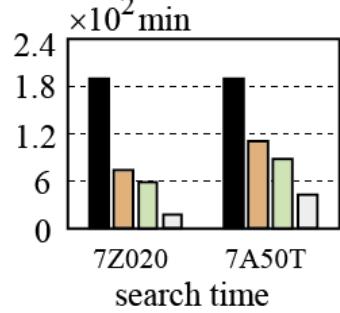
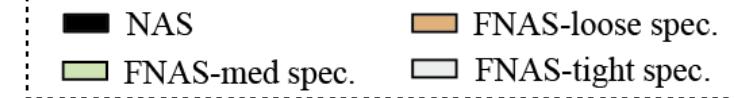
[2, 5, 10, 20]

[1.5, 2, 2.5, 10]

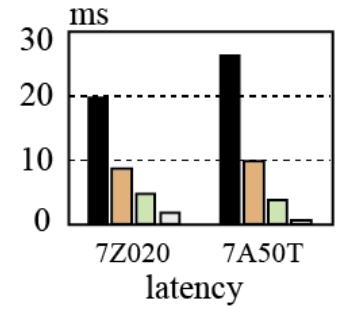
[2.5, 5, 7.5, 10]

# Experimental Results

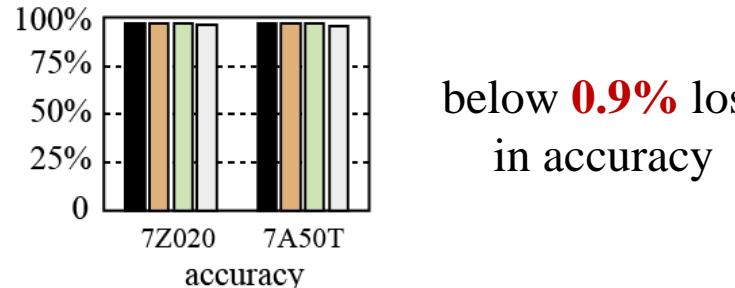
## Different Hardware (CIFAR-10)



up to **11.13X** reduction  
in search time



up to **7.81X** reduction  
in inference latency



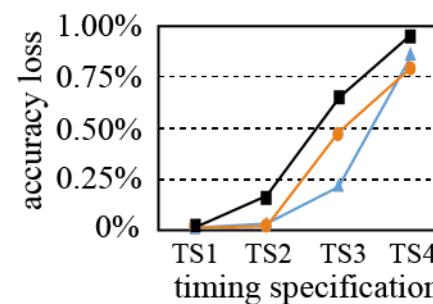
below **0.9%** loss  
in accuracy

## Different Datasets (7Z020)



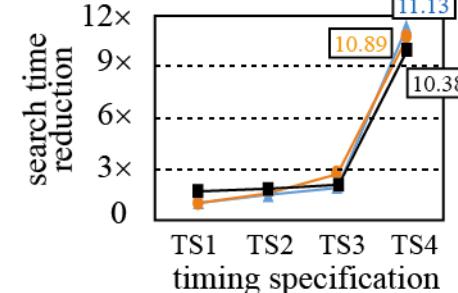
tightness of timing specification

TS1   TS2   TS3   TS4  
loose → tight



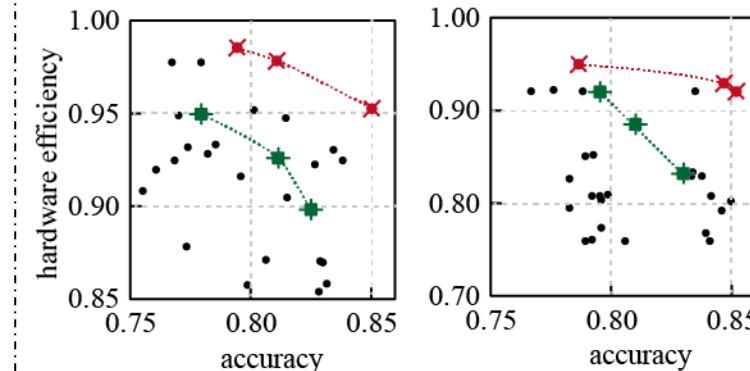
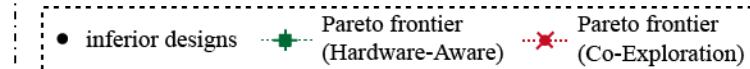
Baseline: NAS

below **1%** loss  
in accuracy



up to **10X** reduction  
in inference latency

## Compare to HW-Aware NAS (CIFAR-10 + 7Z020)



FNAS can significantly  
**push forward**  
the Pareto frontiers between  
**accuracy and efficiency**  
tradeoff

# Experimental Results: Superior to Existing Approaches

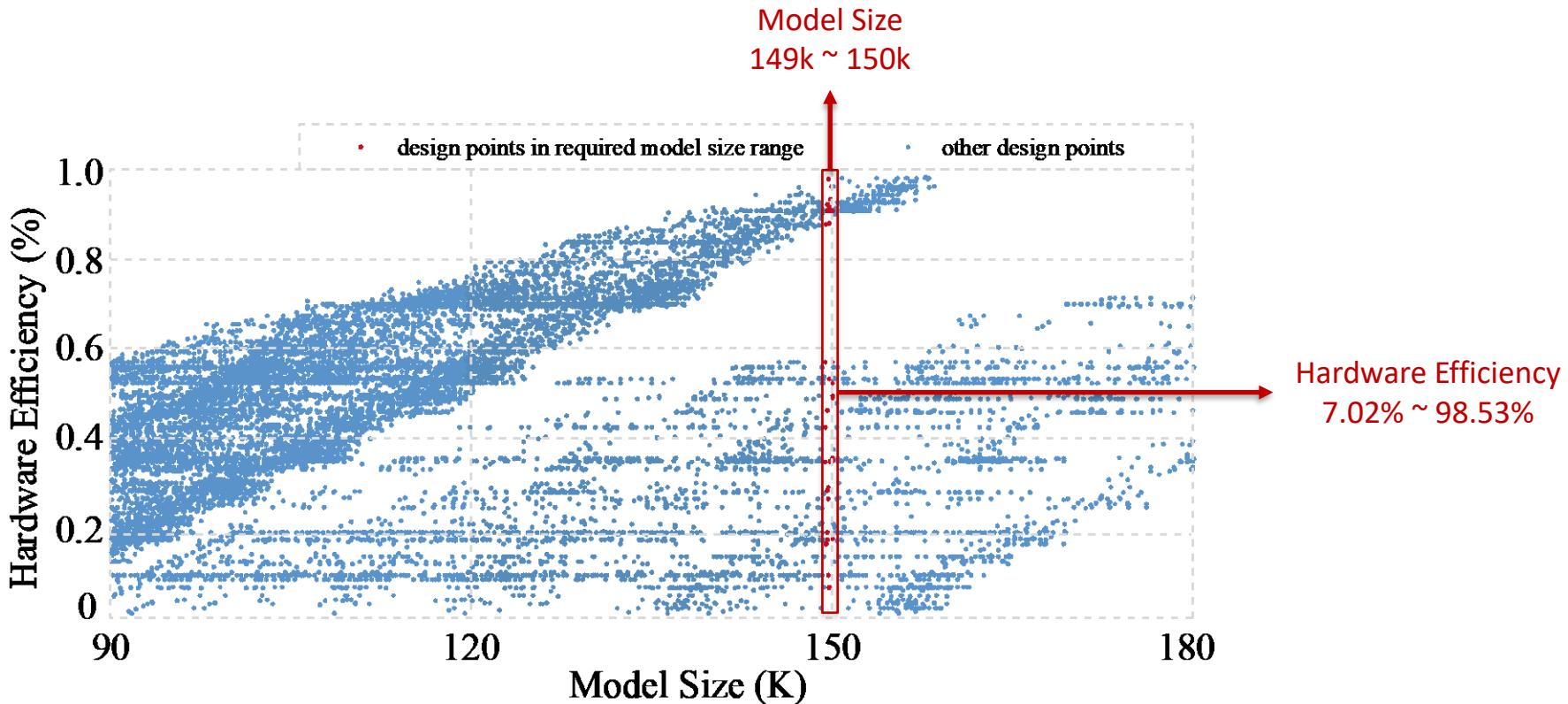
## Optimizing Hardware Efficiency

Comparison the proposed Co-Exploration with Hardware-Aware NAS and Heuristic Sequential Optimization

Dataset	Models	Depth	Parameters	Accuracy (Top1)	Accuracy (Top5)	Pipeline Eff.	FPS	Energy Eff. GOPS/W
CIFAR-10	Hardware-Aware NAS	13	0.53M	84.53%	–	73.27%	16.2	0.84
	Sequential Optimization	13	0.53M	84.53%	–	92.20%	29.7	1.36
	Co-Exploration (OptHW)	10	0.29M	80.18%	–	99.69%	35.5	2.55
	Co-Exploration (OptSW)	14	0.61M	85.19%	–	92.15%	35.5	1.91
ImageNet	Hardware-Aware NAS	15	0.44M	68.40%	89.84%	81.07%	6.8	0.34
	Sequential Optimization	15	0.44M	68.40%	89.84%	86.75%	10.4	0.46
	Co-Exploration (OptHW)	17	0.54M	68.00%	89.60%	96.15%	12.1	1.01
	Co-Exploration (OptSW)	15	0.48M	70.24%	90.53%	93.89%	10.5	0.74

## Optimizing Network Accuracy

# Experimental Results: Importance of Co-Exploration



In the design space:

Models with similar model sizes may have distinct hardware efficiency

=> **Cannot restrict model size** to guarantee **hardware efficiency**

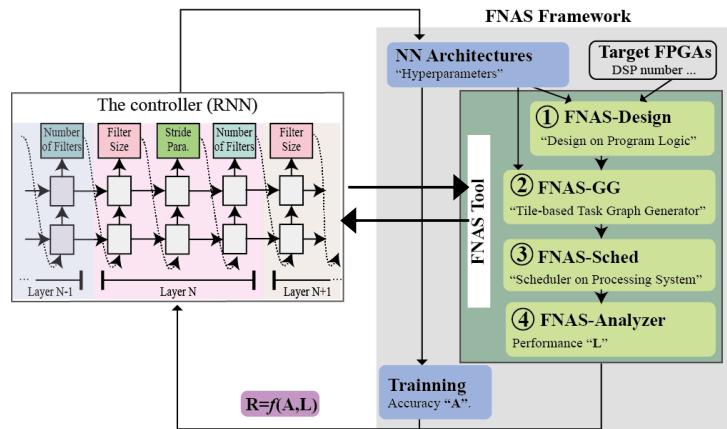
# Outline

- Introduction
- Co-Exploration of Neural Architecture and FPGA Implementation
  - DAC 2019 (Best Paper Nomination)
- Other Research Directions
- Conclusion



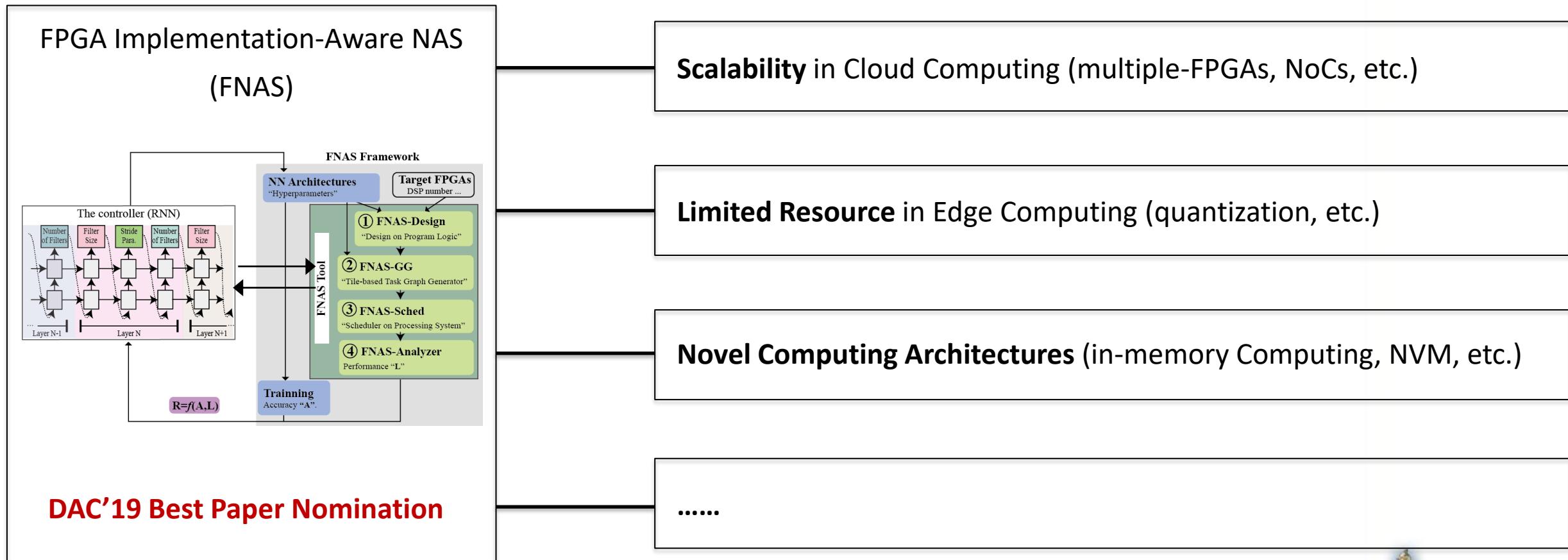
# Our Work on Co-Exploring Neural Architecture and Hardware

## FPGA Implementation-Aware NAS (FNAS)



DAC'19 Best Paper Nomination

# Our Work on Co-Exploring Neural Architecture and Hardware

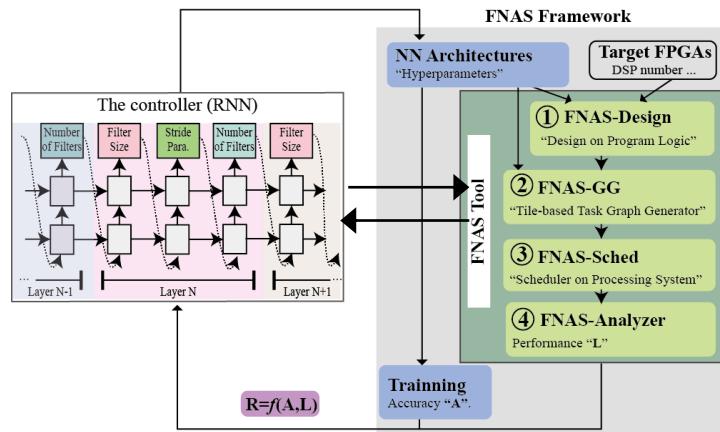


30



# Design Neural Networks on Scalable Architecture

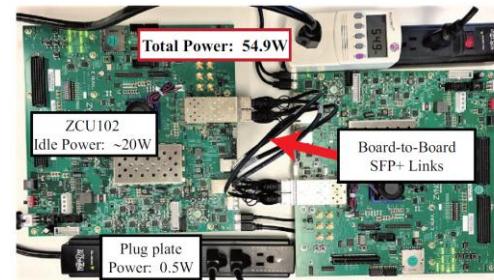
## FPGA Implementation-Aware NAS (FNAS)



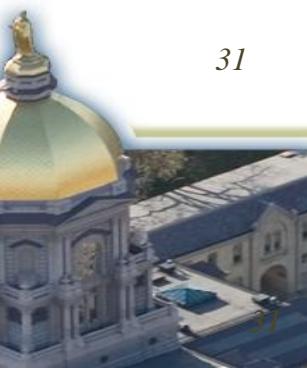
DAC'19 Best Paper Nomination

## Scalability in Cloud Computing (multiple-FPGAs, NoCs, etc.)

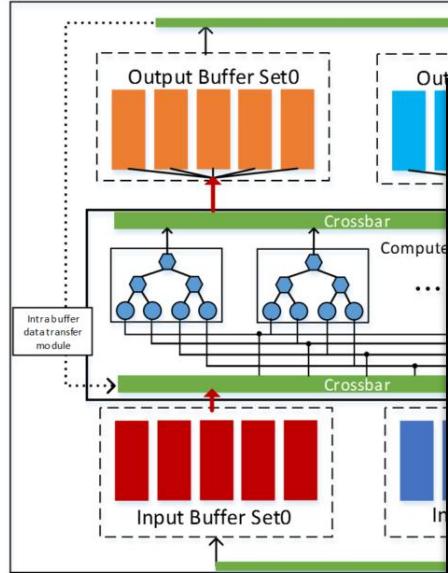
Super-Linear Speedup across  
Multiple FPGAs (XFER)



**CODES+ISSS'19**  
**Best Paper Nomination**



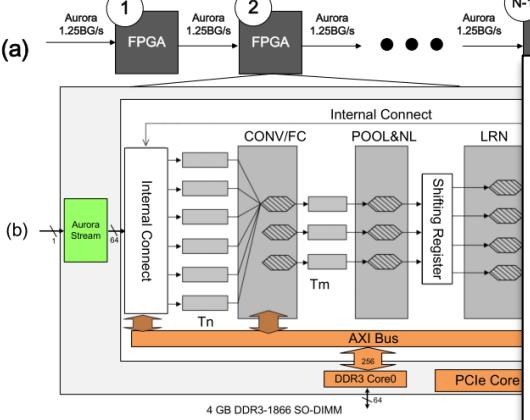
# Existing Implementations



Limited Resources

[Ref] C. Zhang et al. Optimizing neural networks. In Proc. of

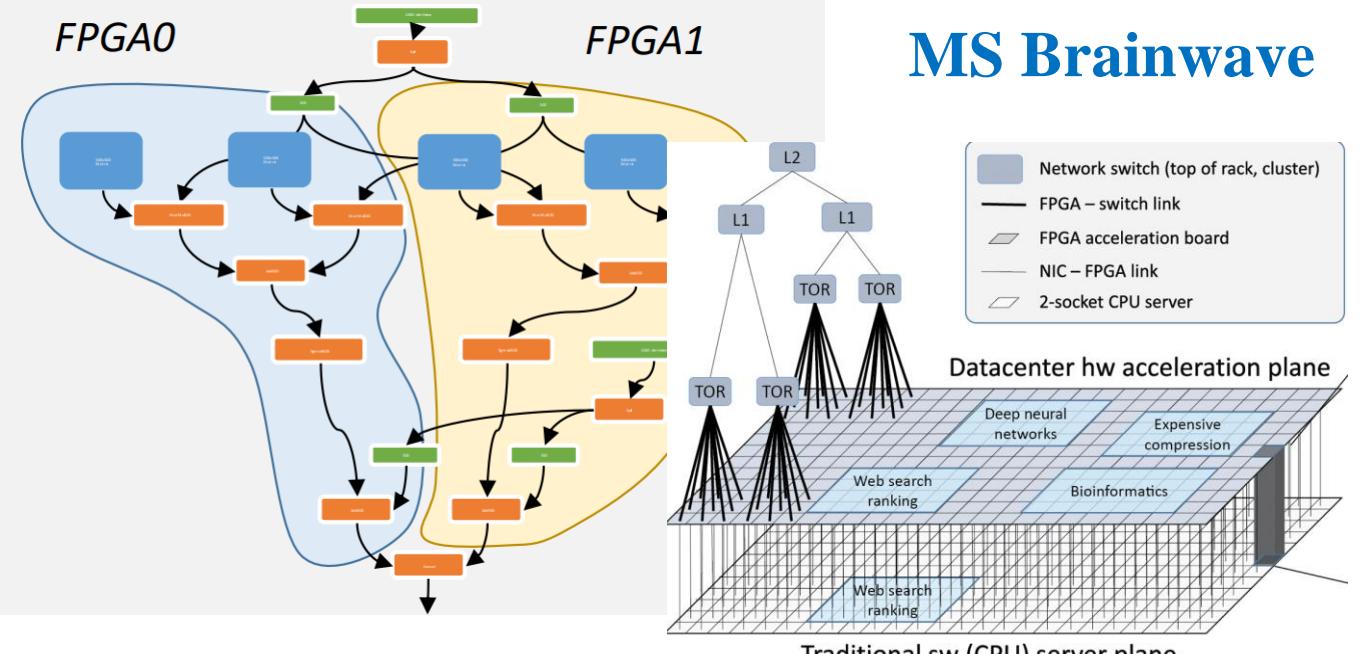
## Single-FPGA



Pipeline → High

[Ref] C. Zhang et al. Energy-efficient cluster. In Proc. of ISLPED, pages

## Pipelined FPGAs

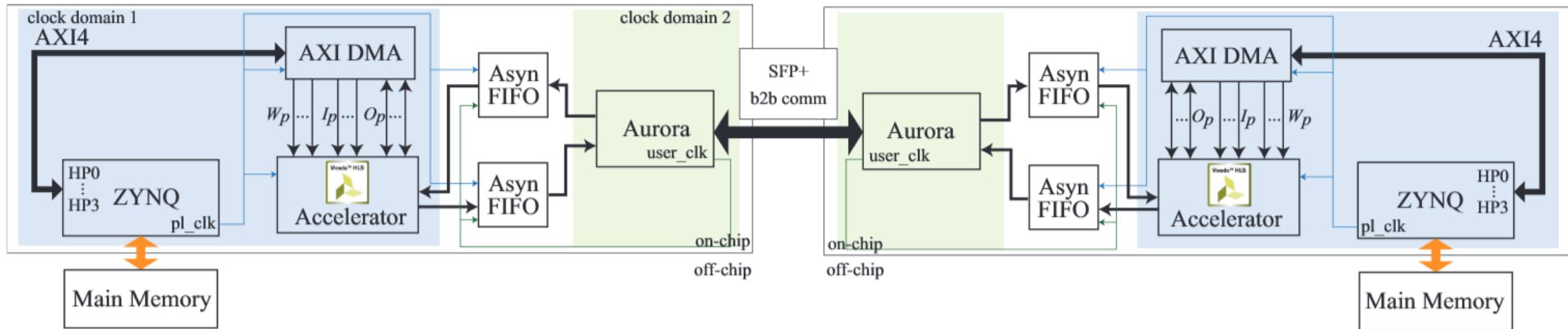


**Lock weights to multiple FPGAs → Specific for RNN  
Not applicable for sharing weights (e.g. CNN)**

[Ref] Eric Chung, et al. Serving DNNs in Real Time at Datacenter Scale with Project Brainwave, In IEEE Micro, 2018

# XFER: Architecture and Motivation

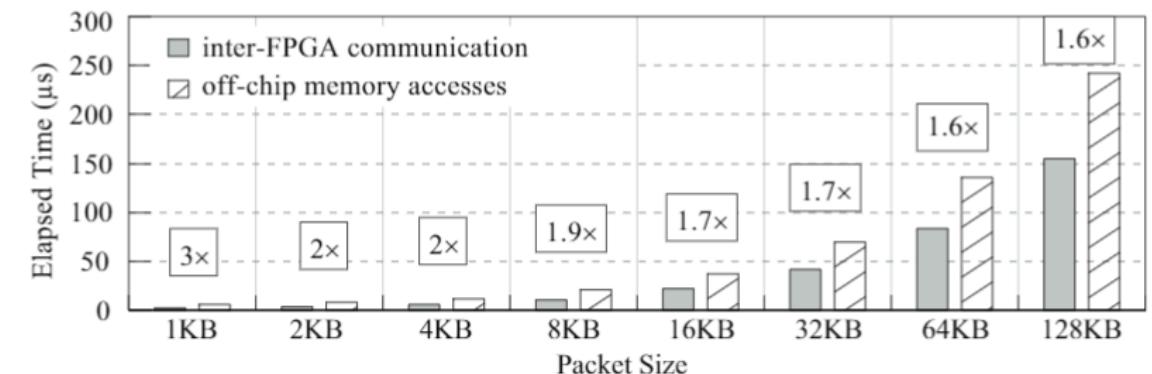
FPGA 1



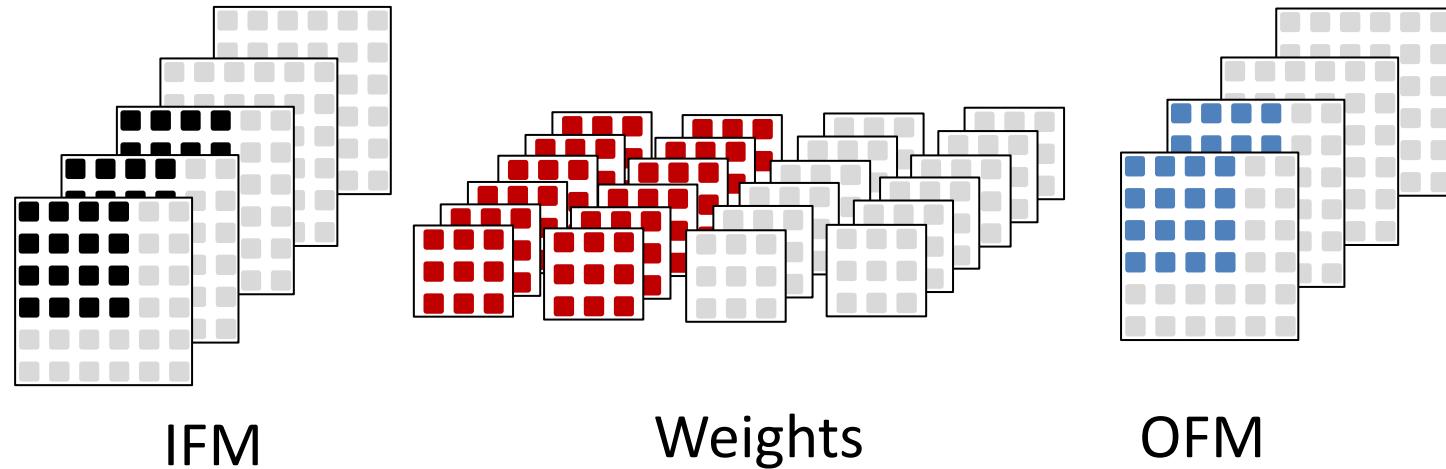
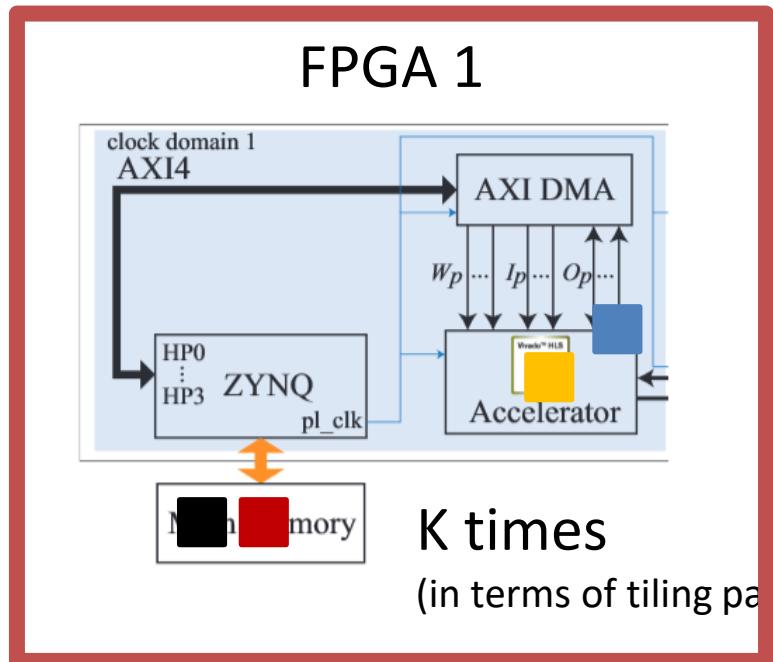
## Feature

- Different **clock domains** for computation (**low**) and communication (**high**)
- Communication **not go through Off-Chip Memory**, but **directly switch between on-chip buffers**

## Results & Motivations

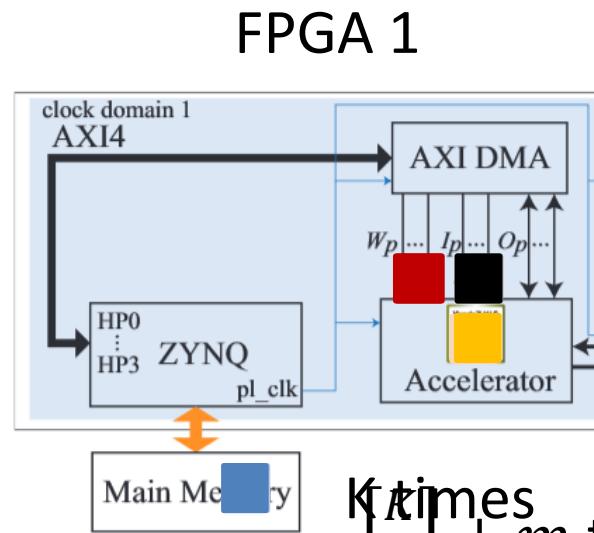


# Performance on Single FPGA bounded by Limited Resource



NN	Operation	Cycles	Note
AlexNet Layer 5	Comm_IFM	2,612	Performance Dominated by <b>Comm_Weights</b> Latency is 5,658
	Comm_Weights	<b>5,658</b>	
	Comm_OFM	368	
	Computation	3,326	

# Double Computation Resource cannot Double Performance

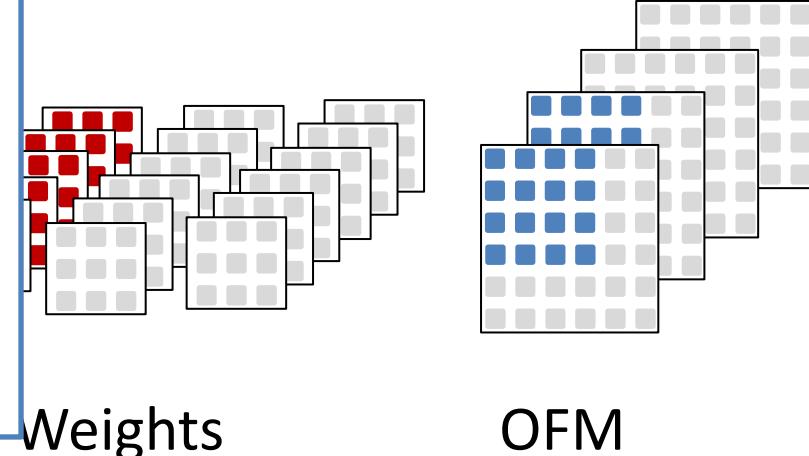


K times  
 $\frac{1}{2} + m \text{ times}$

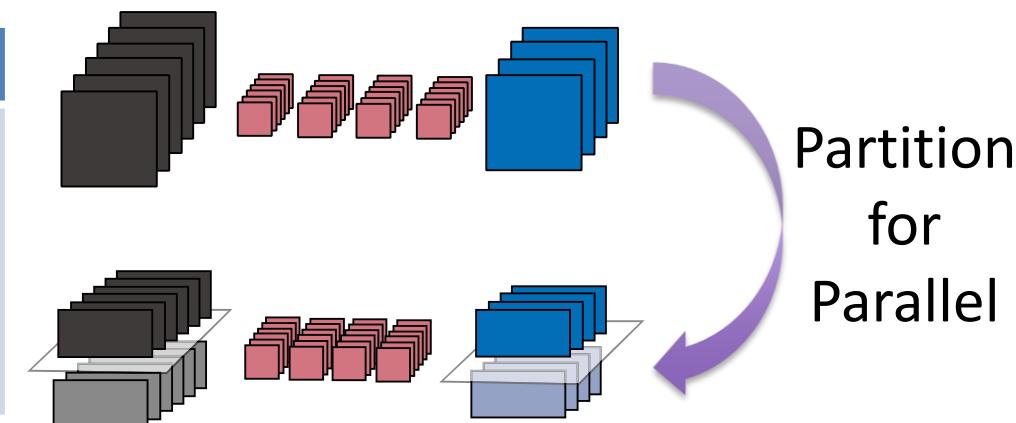
(in terms of tiling parameters)  
(K is in terms of tiling parameters)  
(m is in terms of stride and padding)

## Observation:

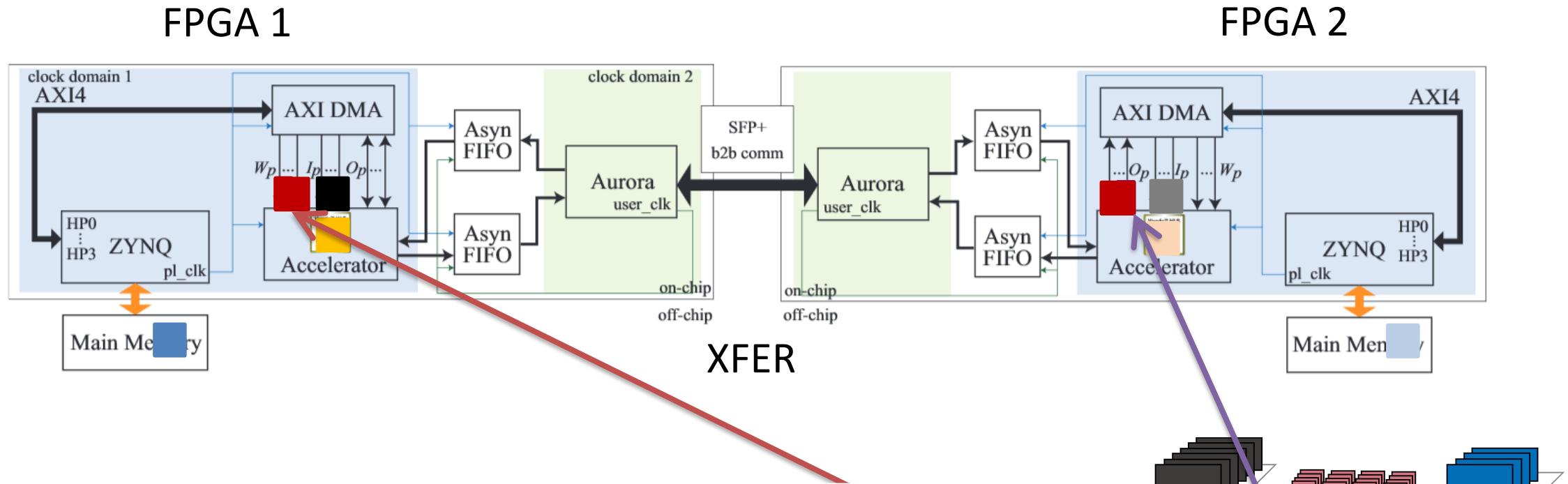
1. Loading IFM, OFM, Computation are conducted **independently** on two FPGAs
2. Performance is still dominated by **loading weights**, and weights are **duplicated** on two FPGAs



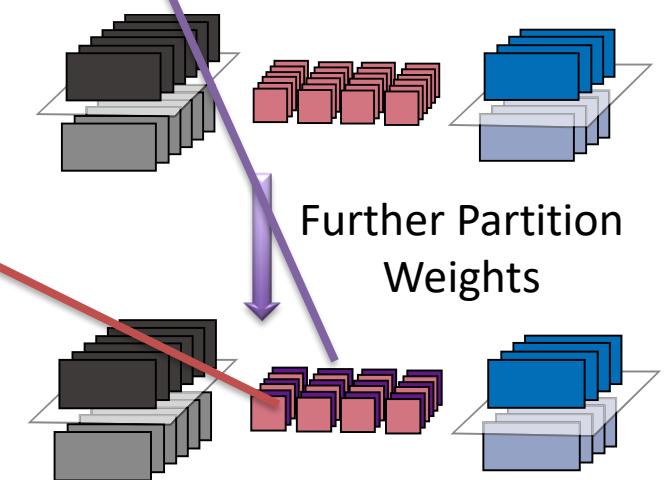
NN	Operation	Single	w/o XFER	Note
AlexNet Layer 5	Comm_IFM	2,612	1,412	<b>1.91X Speedup</b>
	Comm_Weights	<b>5,658</b>	<b>2,953</b>	
	Comm_OFM	368	234	
	Computation	3,326	1,782	



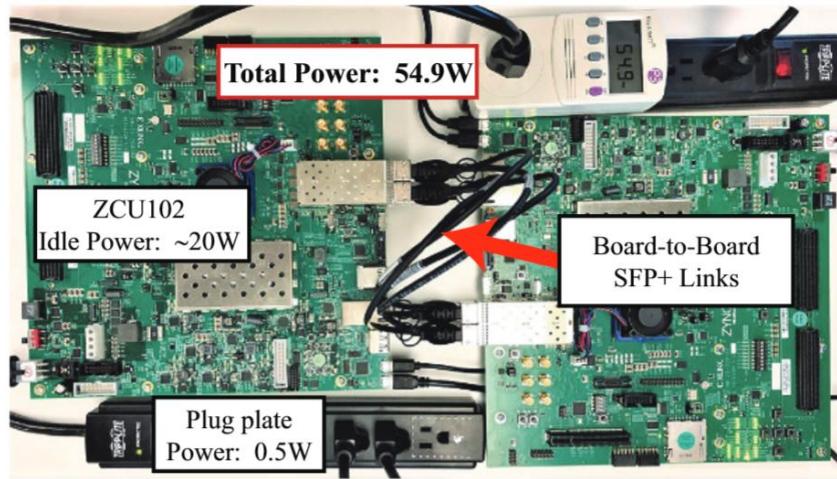
# XFER: Achieve Super-Linear Speedup by Transferring Accesses from Off-Chip Memory to Inter-FPGA Links



NN	Operation	Single	w/o XFER	w/ XFER	Note
AlexNet Layer 5	Comm_IFM	2,612	1,412	1,412	Bottleneck Moves to <b>Comp.</b> <b>3.18X</b> <b>Speedup</b>
	Comm_Weights	<b>5,658</b>	<b>2,953</b>	1,695	
	Comm_OFM	368	234	234	
	Computation	3,326	1,782	<b>1,782</b>	



# Experimental Results



**Testbed:**  
Xilinx ZUC102 FPGAs  
connected by SFP+ Links

Comparison results of XFER with comparisons to GPUs and the existing FPGA designs

Design	mGPU		GPU		FPGA15		ISCA17		ISLPED16		XFER			
Precision	32bits float		32bits float		32bits float		32bits float		16bits fixed		32bits float		16bits fixed	
Device	Jetson TX2		Titan X		VX485T		VX485T		4×VX690t		2×ZCU102		2×ZCU102	
Freq (MHz)	1300MHz		1139MHz		100MHz		100MHz		150MHz		100MHz		200MHz	
Power (Watt)	16.00		162.00		18.61		-		126.00		52.40		54.40	
DSP Uti.	-		-		80%		80%		-		90.79%		55.87%	
BRAM Uti.	-		-		49.71%		43.25%		-		72.92%		92.43%	
Overall Perf.	Lat.	Thr.	Lat.	Thr.	Lat.	Thr.	Lat.	Thr.	Lat.	Thr.	Lat.	Thr.	Lat.	Thr.
	ms	GOPS	ms	GOPS	ms	GOPS	ms	GOPS	ms	GOPS	ms	GOPS	ms	GOPS
	11.1 - 13.2	110.75	5.1 - 6.4	235.55	21.62	69.09	60.13	85.47	30.6	128.8	10.13	149.54	2.27	679.04
E.-E. (GOPS/W)	6.88		1.45		3.71		-		1.02		2.85		12.48	

**Lowest Latency among all competitors**

# Experimental Results

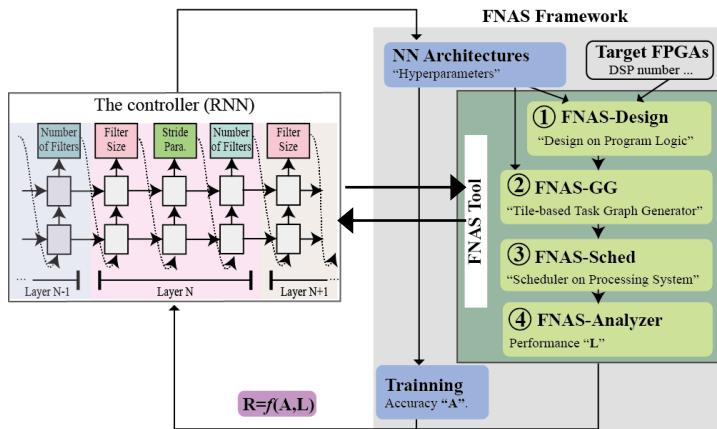
Comparison results of **XFER** with SOTA single FPGA design

Design	32bits float				16bits fixed			
	<i>FPGA15</i>	<b>Super-LIP</b>	<i>FPGA15</i>	<b>Super-LIP</b>				
$\langle T_m, T_n \rangle$	$\langle 64, 7 \rangle$	$\langle 64, 7 \rangle$	$\langle 64, 24 \rangle$	$\langle 128, 10 \rangle$				
Power (W)	25.70 (1 FPGA)	<b>52.40</b> (2 FPGAs)	26.00 (1 FPGA)	<b>54.40</b> (2 FPGAs)				
Perf.	Lat. ms	Thr. GOPs	Lat. ms	Thr. GOPs	Lat. ms	Thr. GOPs	Lat. ms	Thr. GOPs
conv1	7.36	28.6	<b>3.66</b>	<b>57.6</b>	3.74	56.5	<b>0.94</b>	224.5
conv2	5.20	86.1	<b>2.55</b>	<b>175.5</b>	1.48	302.6	<b>0.48</b>	933.1
conv3	4.50	66.4	<b>1.73</b>	<b>172.7</b>	1.20	249.6	<b>0.33</b>	906.2
conv4	3.41	65.7	<b>1.31</b>	<b>171.0</b>	0.89	252.6	<b>0.35</b>	640.8
conv5	2.28	66.0	<b>0.88</b>	<b>170.9</b>	0.59	251.7	<b>0.17</b>	879.5
overall	22.75	66.6	<b>10.13</b>	<b>149.5</b>	7.90	195.1	<b>2.27</b>	679.0
Perf. Impr.	1.00×		<b>2.25×</b>	1.00×		<b>3.48×</b>		
E.-E. (GOPs/W)	2.59		<b>2.85</b>	7.51		<b>12.48</b>		
E.-E. Impr.	-		<b>9.21%</b>	-		<b>39.86%</b>		

Achieve Super-Linear Speedup for both 32-bits and 16-bits implementations

# Design Neural Networks on Scalable Architecture

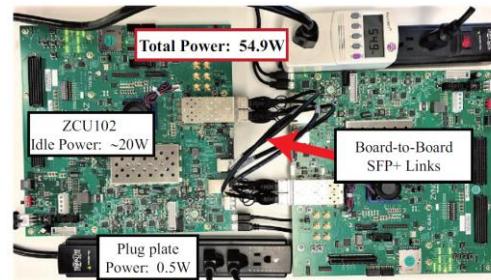
## FPGA Implementation-Aware NAS (FNAS)



DAC'19 Best Paper Nomination

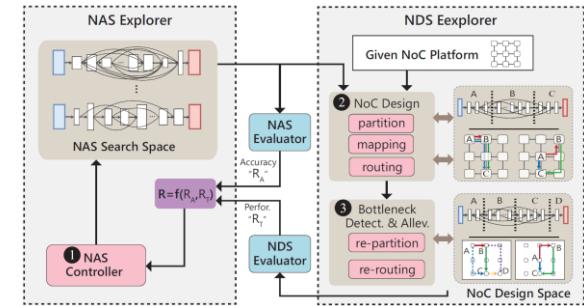
## Scalability in Cloud Computing (multiple-FPGAs, NoCs, etc.)

### Super-Linear Speedup across Multiple FPGAs (XFER)



CODES+ISSS'19  
Best Paper Nomination

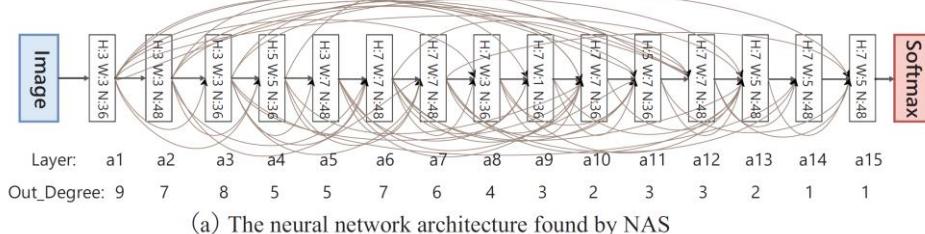
### Co-Explore Network-on-Chip Design and Neural Architectures (NANDS)



ASP-DAC'20  
Best Paper Nomination

# NANDS: Co-Explore NoC Design and Neural Architectures

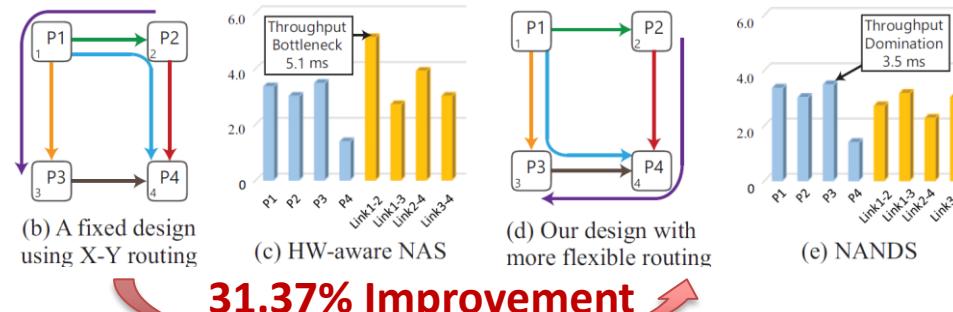
## Why NoC? Complicated Structure from NAS



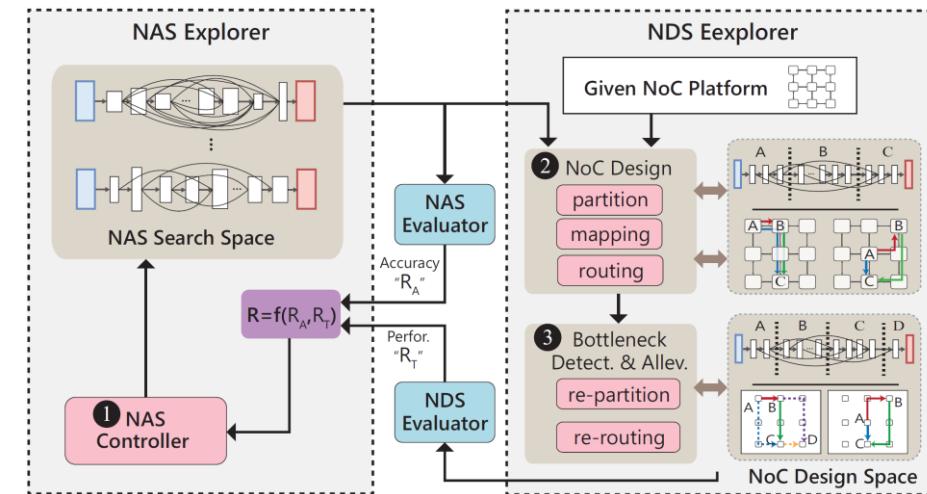
Operation Time	Platforms	Single Processing Element		4 Processing Elements	
		Bus Interconnection	2-D Mesh NoC	Bus Interconnection	2-D Mesh NoC
Computation (ms)		12.4		3.4	3.4
Data transmission (ms)	—		14.7		6.2

(b) The timing performance of network implementations on different platforms

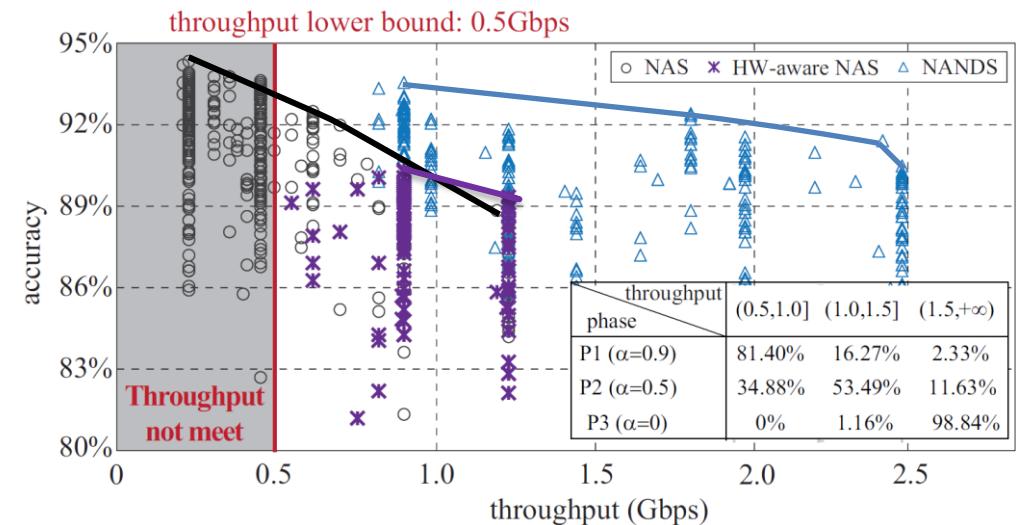
## Fixed design (X-Y routing) leads lower performance



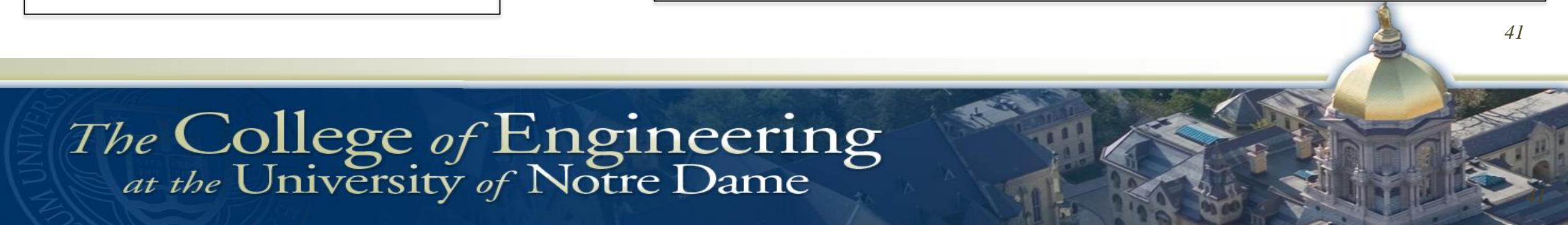
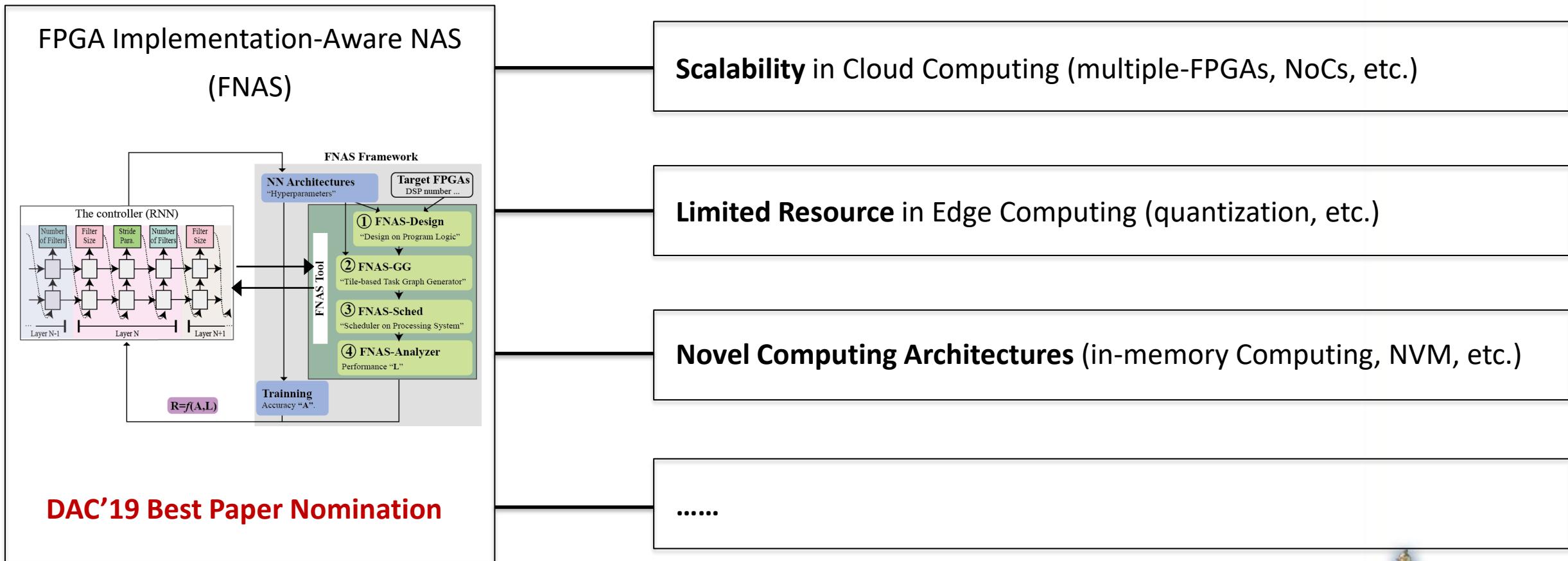
## Co-Exploration Framework



## Results

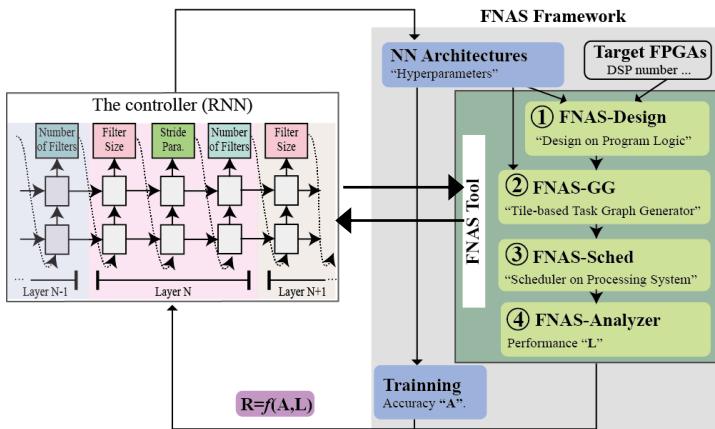


# Design Neural Networks for Resource-Limited Device



# Design Neural Networks for Resource-Limited Device

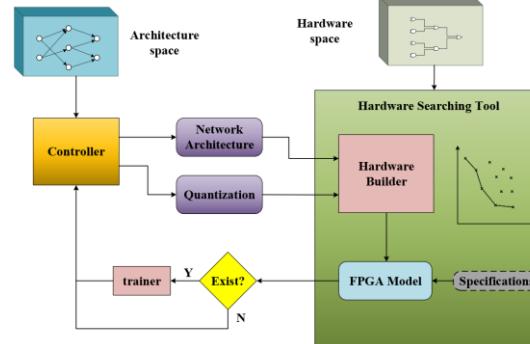
## FPGA Implementation-Aware NAS (FNAS)



DAC'19 Best Paper Nomination

## Limited Resource in Edge Computing (quantization, etc.)

### Co-Explore Quantization and Neural Architectures (QuanNAS)



ICCAD'19



*The College of Engineering  
at the University of Notre Dame*

# QuanNAS: Co-Explore Quantization and Neural Architectures

## Search Space

Parameter	Symbol	Value
# filters	$N$	(24, 36, 48, 64)
filter height	$F_h$	(1, 3, 5, 7)
filter width	$F_w$	(1, 3, 5, 7)
stride height	$S_h$	(1, 2, 3)
stride width	$S_w$	(1, 2, 3)
pooling size	$P_s$	(1, 2)
activation integer bits	$A_i$	(0, 1, 2, 3)
activation fractional bits	$A_f$	(0, 1, 2, 3, 4, 5, 6)
weight integer bits	$W_i$	(0, 1, 2, 3)
weight fractional bits	$W_f$	(0, 1, 2, 3, 4, 5, 6)

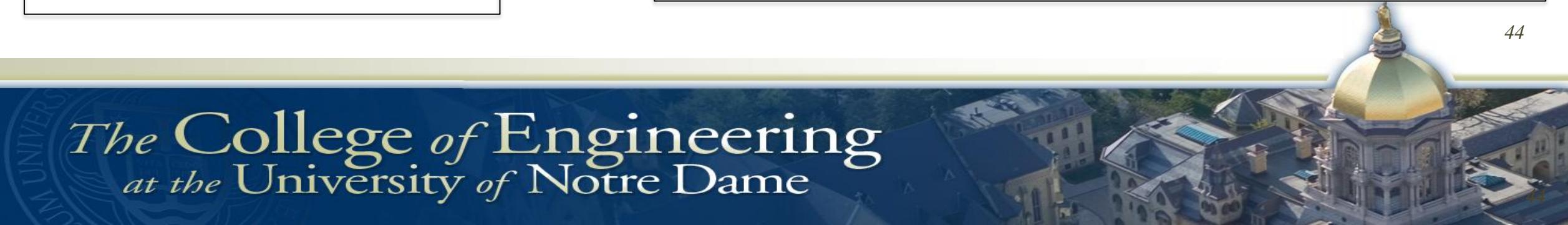
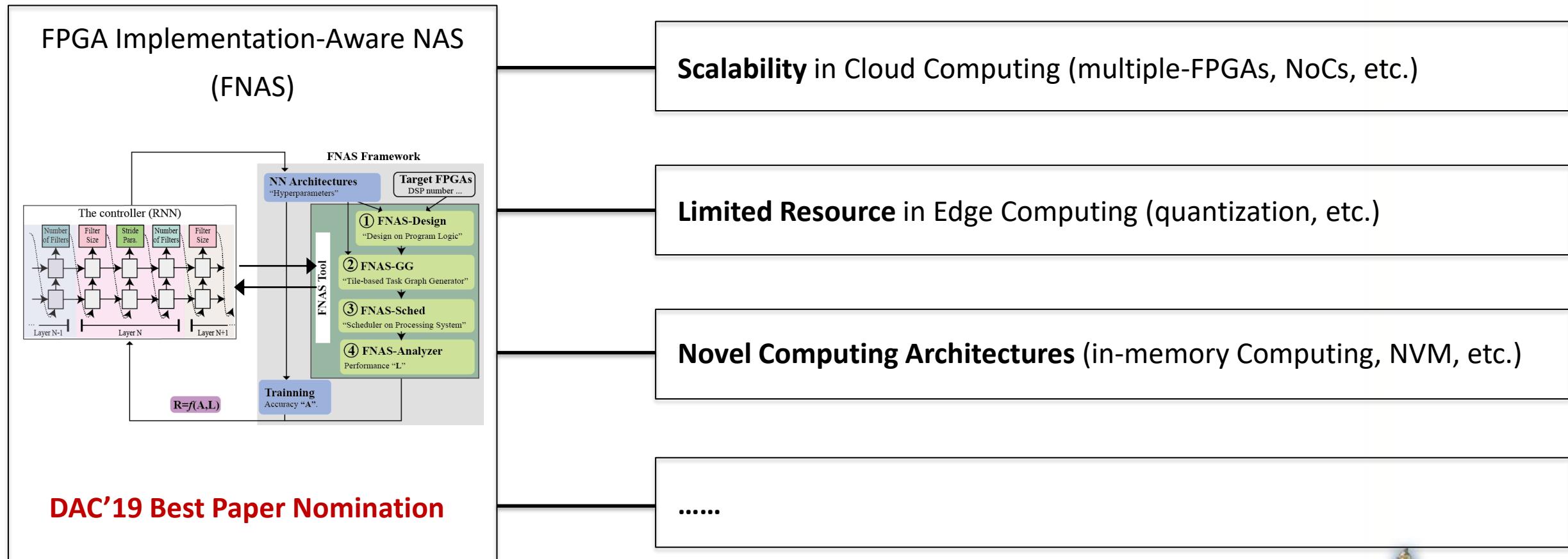
Besides filter size and shape, **we co-explore bit width of weights and activation data, for FPGAs and ASICs.**

## Results

Design	rL	rT	Acc	Acc	#LUTs	FPS	para size (kbits)
			w/o Quan	w/ Quan			
A <sub>1</sub> -d <sub>1</sub>	100,000	500	87.76%	80.23%	99,871	556	1,867
A <sub>1</sub> -d <sub>2</sub>	100,000	1000	87.76%	25.79%	99,848	1157	1,189
B <sub>1</sub> -d <sub>1</sub>	100,000	500	89.71%	87.64%	96,904	512	3,463
B <sub>1</sub> -d <sub>2</sub>	100,000	1000	89.71%	64.35%	98,752	1020	2,784
B <sub>1</sub> -d <sub>3</sub>	300,000	2000	89.71%	50.93%	285,441	2083	2,835
D	30,000	1000	83.65%	82.98%	29,904	1293	457
E	100,000	1000	86.99%	82.76%	94,496	1042	1,923
F	300,000	2000	87.03%	84.92%	299,860	2089	1,217

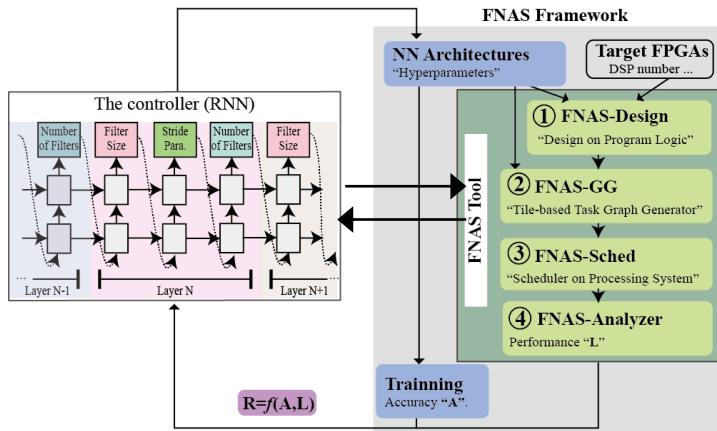
- A,B are first explored by NAS, then to be quantized. In this way, the **accuracy loss** by quantization is huge.
- D,E,F are explored by our Co-Exploration framework, which output **better accuracy** for quantized network with **smaller parameter size**

# Design Neural Networks for Novel Computing Architectures



# Design Neural Networks for Novel Computing Architectures

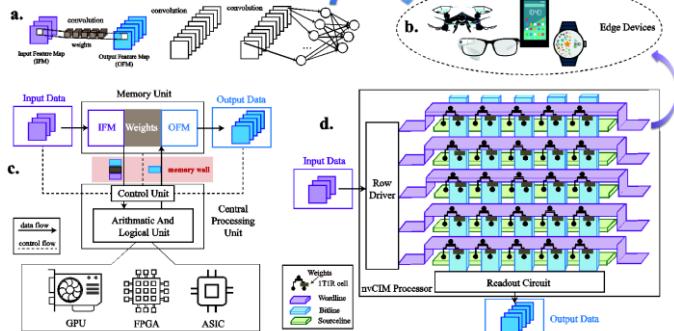
## FPGA Implementation-Aware NAS (FNAS)



DAC'19 Best Paper Nomination

## Novel Computing Architectures (in-memory Computing, NVM, etc.)

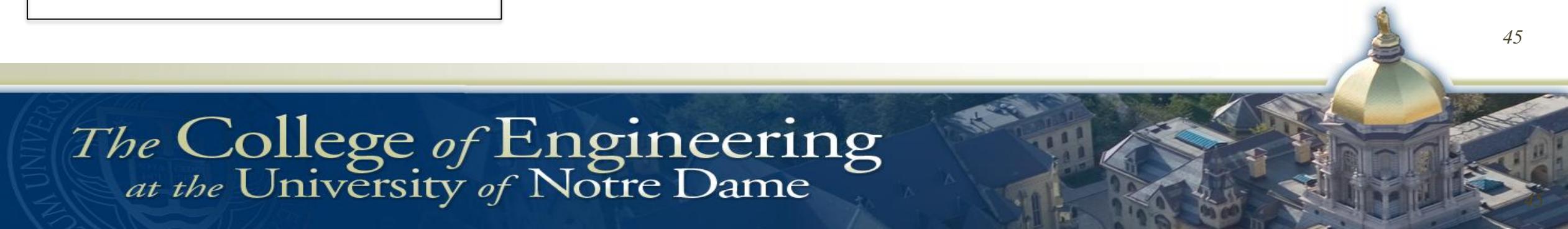
### Integrating Memristors and CMOS for Better AI



Co-Explore Neural Architectures and Devices

Submitted to TC

Nature Electronics'19

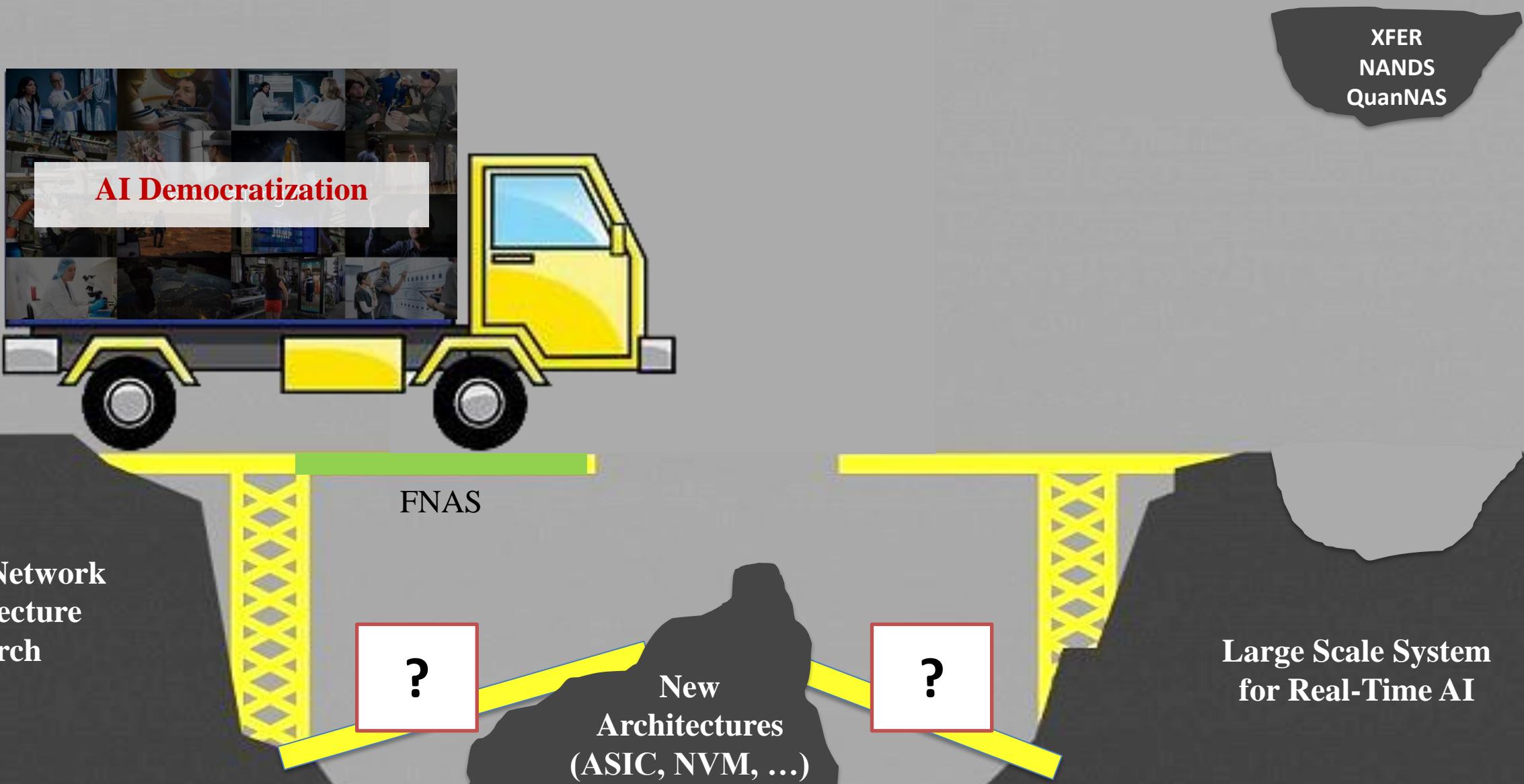


# Outline

- Introduction
- Co-Explore Neural Architecture and FPGA Implementation
  - DAC 2019 (Best Paper Nomination)
- Research Directions
- Conclusion



# Conclusion and Future Work



# Reference

- [1] Weiwen Jiang, Xinyi Zhang, Edwin H.-M. Sha, Qingfeng Zhuge, Lei Yang, Yiyu Shi and Jingtong Hu, "Accuracy vs. Efficiency: Achieving Both through FPGA-Implementation Aware Neural Architecture Search," in Proc. of **DAC 2019**. (**Nominated for Best Paper Award**)
- [2] Weiwen Jiang, Edwin Sha, Xinyi Zhang, Lei Yang, Qingfeng Zhuge, Yiyu Shi and Jingtong Hu, "Achieving Super-Linear Speedup across Multi-FPGA for Real-Time DNN Inference," **CODES+ISSS 2019** and **ACM TECS** (**Nominated for Best Paper Award**)
- [3] Lei Yang, Weiwen Jiang, Weichen Liu, Edwin Sha, Yiyu Shi and Jingtong Hu, "Co-Exploring Neural Architecture and Network-on-Chip Design for Real-Time Artificial Intelligence," **ASP-DAC 2020** (**Nominated for Best Paper Award**)
- [4] Qing Lu, Weiwen Jiang, Xiaowei Xu, Yiyu Shi and Jingtong Hu, "On Neural Architecture Search for Resource-Constrained Hardware Platforms," in Proc. of **ICCAD 2019** (Invited Paper)
- [5] Weiwen Jiang, Bike Xie, Chun-Chen Liu and Yiyu Shi, "Integrating Memristors and CMOS for Better AI," **Nature Electronics**, September 2019

# Thank You!

*The College of Engineering  
at the University of Notre Dame*

