



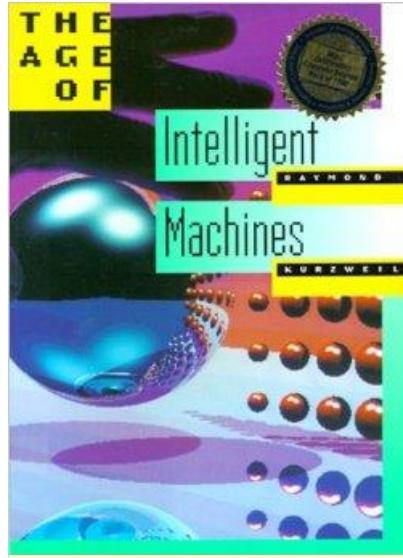
# Exploration of FPGA based Accelerators for Convolutional Neural Networks

Guangyu Sun  
Peking University

[gsun@pku.edu.cn](mailto:gsun@pku.edu.cn)



# Computer vs. Intelligence

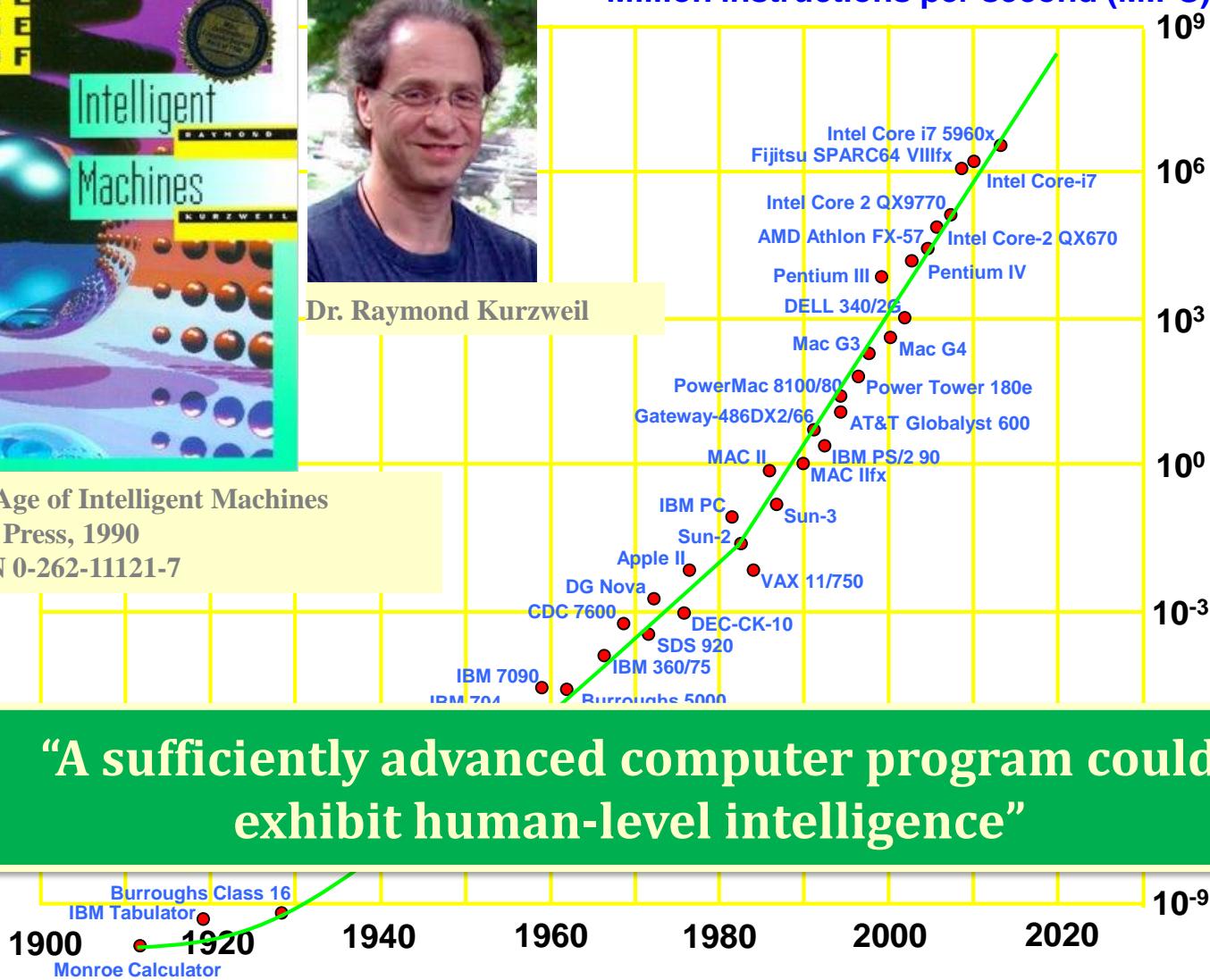


The Age of Intelligent Machines  
MIT Press, 1990  
ISBN 0-262-11121-7



Dr. Raymond Kurzweil

Million instructions per second (MIPS)



"A sufficiently advanced computer program could exhibit human-level intelligence"



Human

100G Neurons



Monkey

3G Neurons



Lizard

2M Neurons



Worm

300 Neurons



Bacteria

1 "Neuron"

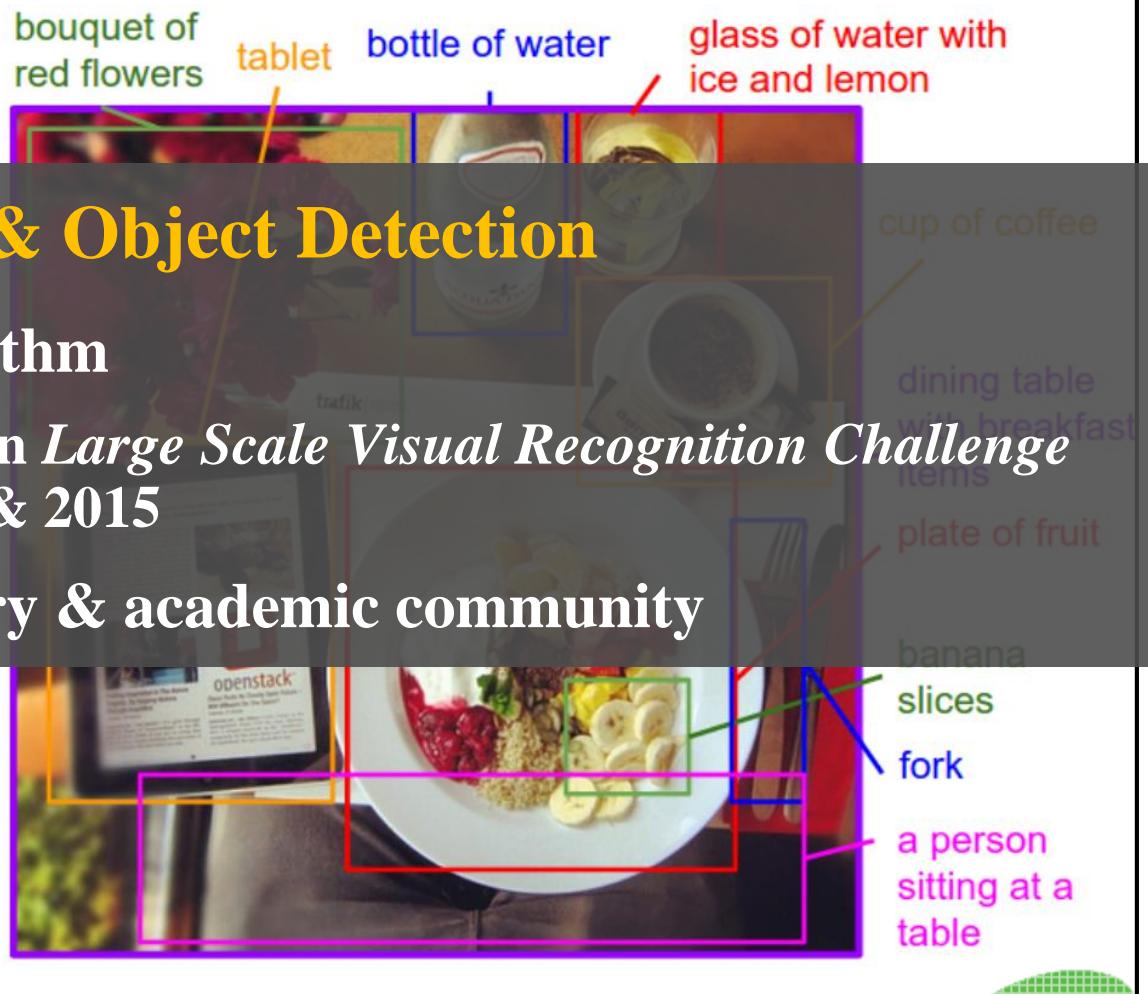


# Convolutional Neural Network

Automotive Safety



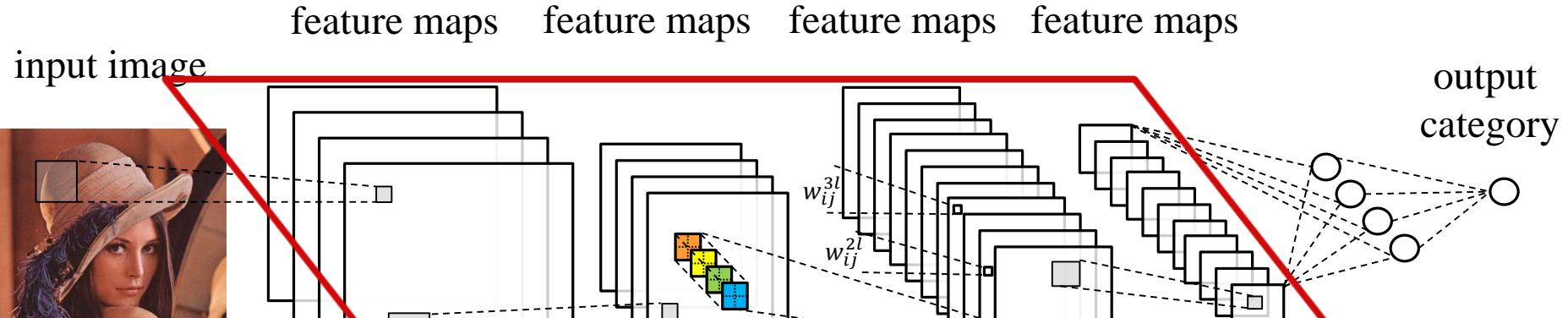
Image Descriptions @ Stanford



## Image Classification & Object Detection

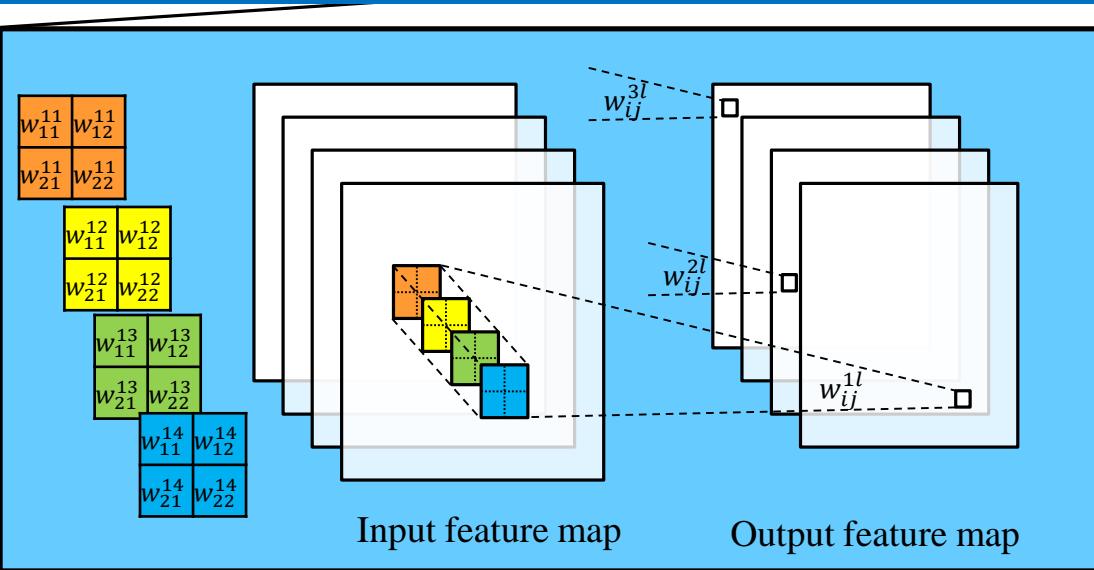
- State-of-the-art Algorithm
  - Highest Correctness in *Large Scale Visual Recognition Challenge* 2012 & 2013 & 2014 & 2015
- Widely used in industry & academic community

# Convolutional Neural Network



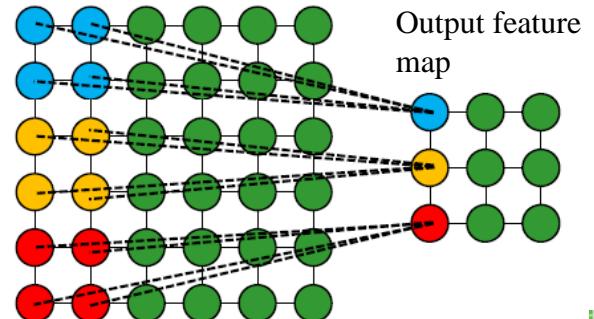
**Convolutional layers account for over 90% computation**

- [1] A. Krizhevsky, etc. Imagenet classification with deep convolutional neural networks. NIPS 2012.
- [2] J. Cong and B. Xiao. Minimizing computation in convolutional neural networks. ICANN 2014



Max-pooling is optional

Input feature map



# Related Work

---



## FPGA-based:

1. **A programmable parallel accelerator for learning and classification.** PACT 2010;
2. **A Dynamically Configurable Coprocessor for Convolutional Neural Networks.** ISCA 2010;
3. **Memory-Centric Accelerator Design for Convolutional Neural Networks.** ICCD 2013;
4. **Going deeper with embedded FPGA platform for convolutional neural network.** FPGA 2016;
5. **Throughput-optimized OpenCL-based FPGA accelerator for large-scale convolutional neural networks.** FPGA 2016;
6. **Deepburning: automatic generation of FPGA-based learning accelerators for the neural network family.** DAC 2016;

## ASIC-based:

1. **DianNao: A Small-Footprint High-Throughput Accelerator for Ubiquitous Machine-Learning.** ASPLOS 2014;
2. **ShiDianNao: Shifting Vision Processing Closer to the Sensor.** ISCA 2015;
3. **ISAAC: A Convolutional Neural Network Accelerator with In-Situ Analog Arithmetic in Crossbars.** ISCA 2016 ...



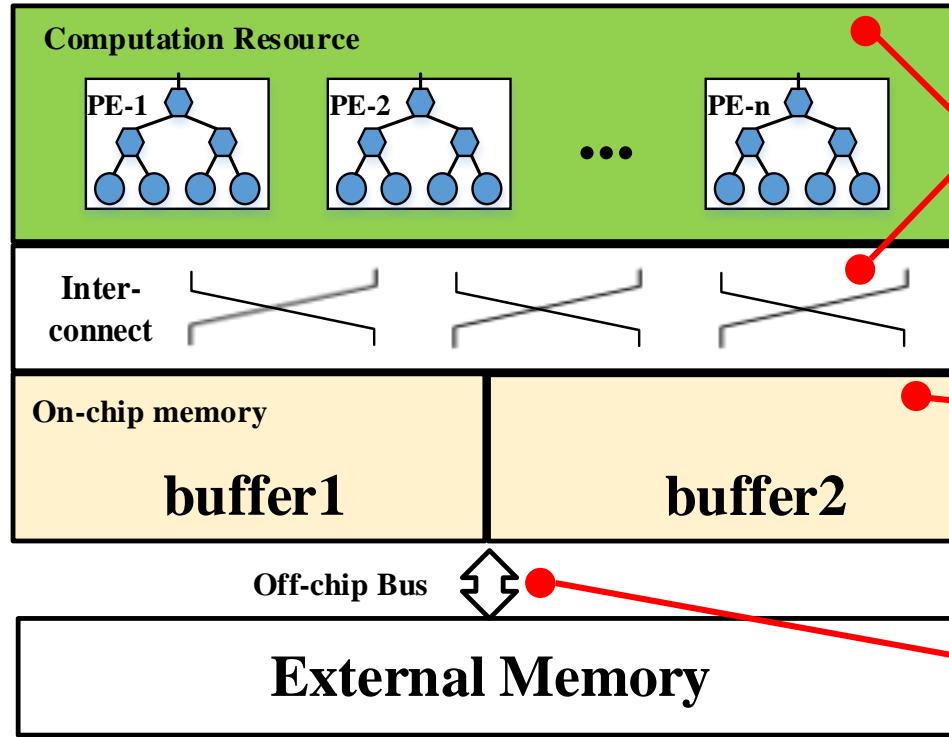
# ***Our Work***

---

- ◆ **Single-board Acceleration**
  - ◆ **Optimizing FPGA-based Accelerator Design for Deep Convolutional Neural Networks (FPGA 2015)**
- ◆ **Multi-board Acceleration**
  - ◆ **Energy-Efficient CNN Implementation on a Deeply Pipelined FPGA Cluster (ISLPED 2016)**
- ◆ **Automated Design Flow**
  - ◆ **Caffeine: Towards Uniformed Representation and Acceleration for Deep Convolutional Neural Networks (ICCAD 2016)**



# Convolution Layer Computing on FPGA



Challenge : Resource

Solution: Unroll & Pipeline



Challenge: BRAM

Solution: Loop Tiling



Challenge: Bandwidth

Solution: Data reuse



Challenge: Co-optimization,  
-- match ‘computation’ & ‘communication’

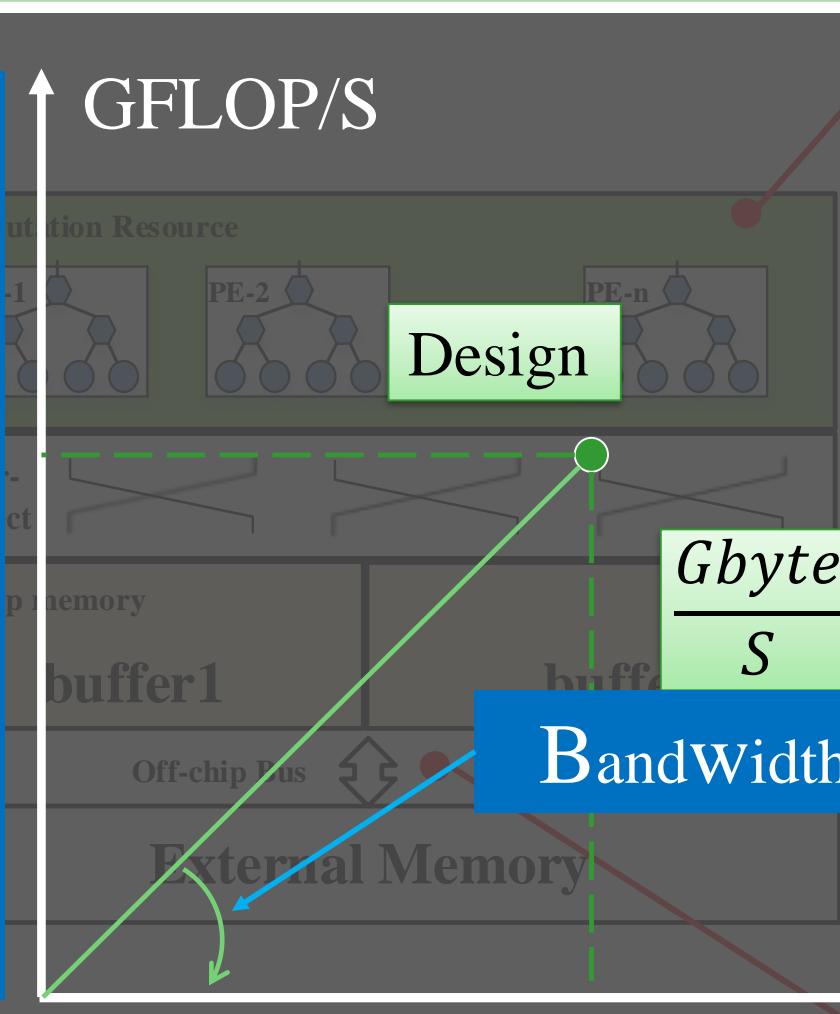
Solution: Roofline Model

Main  
Contribution



# Definitions in Roofline Model

Computational Performance



Computation To Communication (CTC) Ratio

Computational Perf.

$$= \frac{\text{total number of operations}}{\text{execution cycles} \times \text{cycle period}}$$

GFLOP/S

Computation To

$$\frac{Gbyte}{S} = \frac{GFLOP/s}{FLOP/(Mem.\ Byte\ Access)}$$

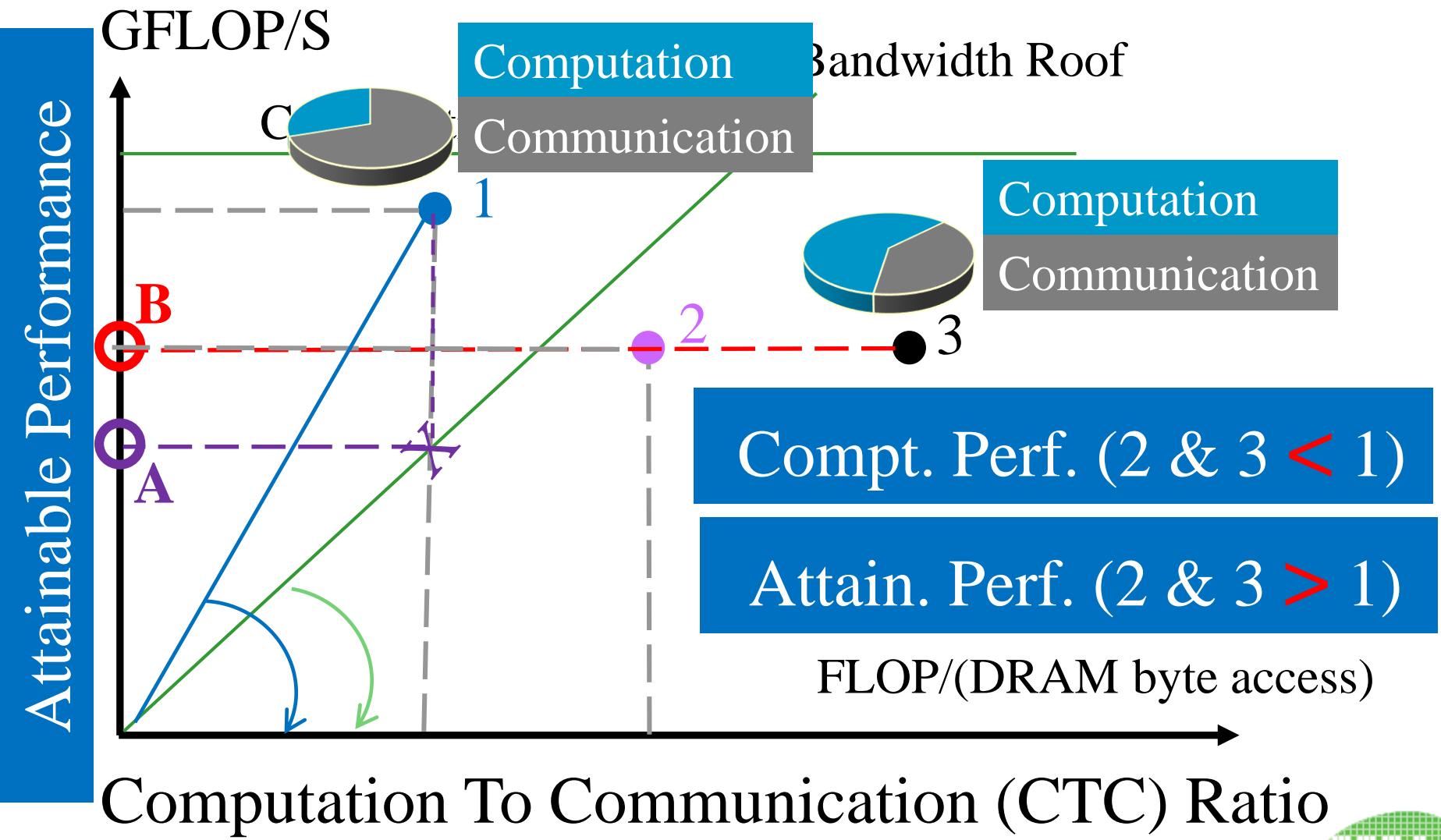
$$= \frac{\text{Total number of operations}}{\text{Amount of external data access}}$$

FLOP / (Mem. Byte Access)

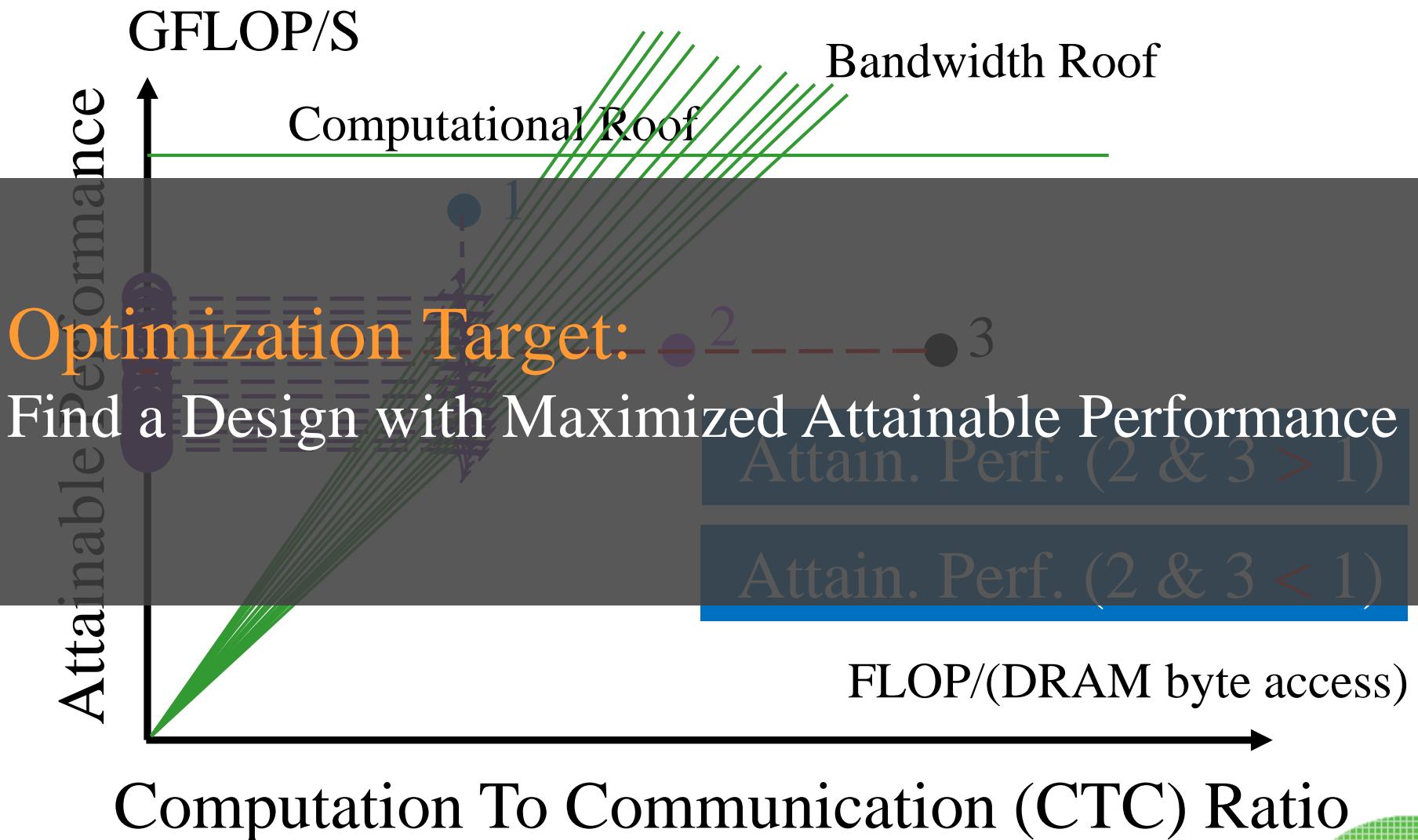
FLOP/(Memory byte access)

Data Throughput

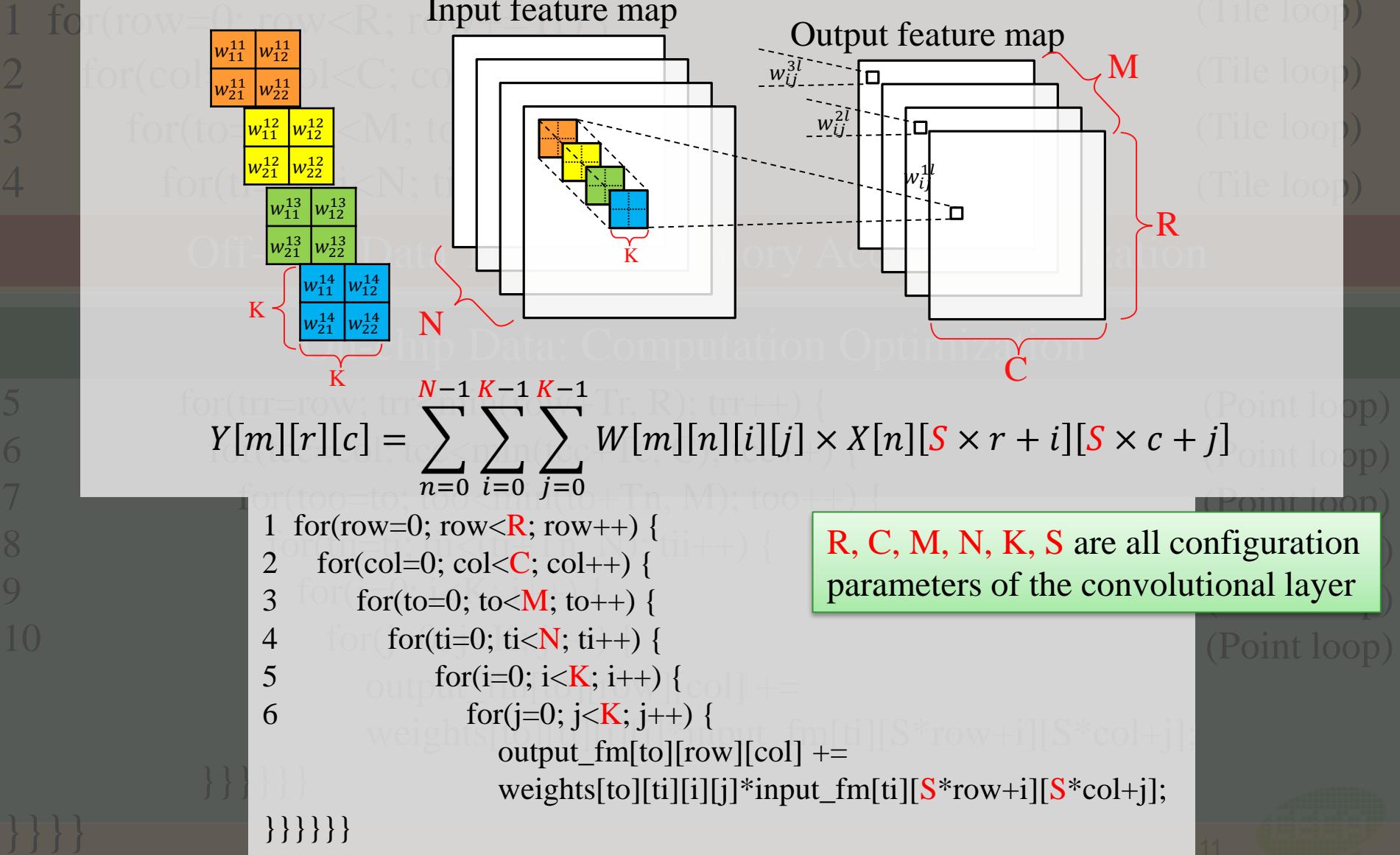
# Attainable Performance



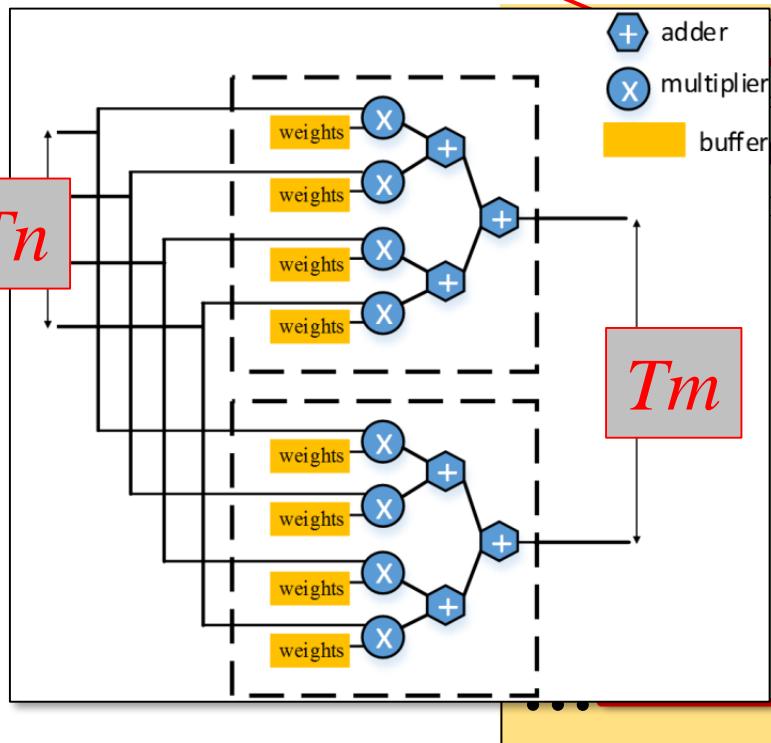
# Attainable Performance



# Loop tiling



# Computation Optimization



```
(Tile loops)
for(ti=0; i<K; i++) {
    for(j=0; j<K; j++) {
        for(trr=row; trr<min(row+Tr, R); trr++) {
            for(tcc=col; tcc<min(tcc+Tc, C); tcc++) {
                HLS pipeline
                for(too=to; too<min(to+Tm, M); too++) {
                    HLS UNROLL
                    for(tii=ti; tii<(ti+Tn, N); tii++) {
                        HLS UNROLL
                        output_fm[too][row][col] +=
                            weights[too][tii][i][j]*input_fm[tii][S*row+i][S*col+j];
                    }
                }
            }
        }
    }
}
```

$$\text{Computational Performance} = \frac{\text{total number of operations}}{\text{execution cycles}}$$

$$\text{execution cycles} = \frac{R}{Tr} \times \frac{C}{Tc} \times \left\lceil \frac{M}{Tm} \right\rceil \times \left\lceil \frac{N}{Tn} \right\rceil \times (K \times K \times Tc \times Tr + P)$$

$$\approx \left\lceil \frac{M}{Tm} \right\rceil \times \left\lceil \frac{N}{Tn} \right\rceil \times R \times C \times K \times K$$



# Memory Access Optimization

```

1 for(row=0; row<R; row+=Tr) {
2   for(col=0; col<C; col+=Tc) {
3     for(to=0; to<M; to+=Tm) {
4       for(ti=0; ti<N; ti+=Tn) {

```

load output feature map

S: foo(output\_fm(to, row, col));

store output feature map

```

      }
    }
  }
}
```

```

1 for(row=0; row<R; row+=Tr) {
2   for(col=0; col<C; col+=Tc) {
3     for(to=0; to<M; to+=Tm) {
4       for(ti=0; ti<N; ti+=Tn) {
5         for(trr=row; trr<min(row+Tr, R); trr++) {
6           for(tcc=col; tcc<min(tcc+Tc, C); tcc++) {
7             for(too=to; too<min(to+Tn, M); too++) {
8               for(tii=ti; tii<(ti+Tn, N); tii++) {
9                 for(i=0; i<K; i++) {
10                for(j=0; j<K; j++) {
11                  output_fm[to][row][col] +=
12                    weights[to][ti][i][j]*
13                    input_fm[ti][S*row+i][S*col+j];
14                }
15              }
16            }
17          }
18        }
19      }
20    }
21  }
22}

```

**Before local memory promotion**

$$\begin{aligned}
 & \text{'output\_fm' memory access} \\
 & = 2 \times \frac{R}{Tr} \times \frac{C}{Tc} \times \boxed{\frac{M}{Tm} \times \frac{N}{Tn}} \times \text{Size}_{\text{array}}
 \end{aligned}$$

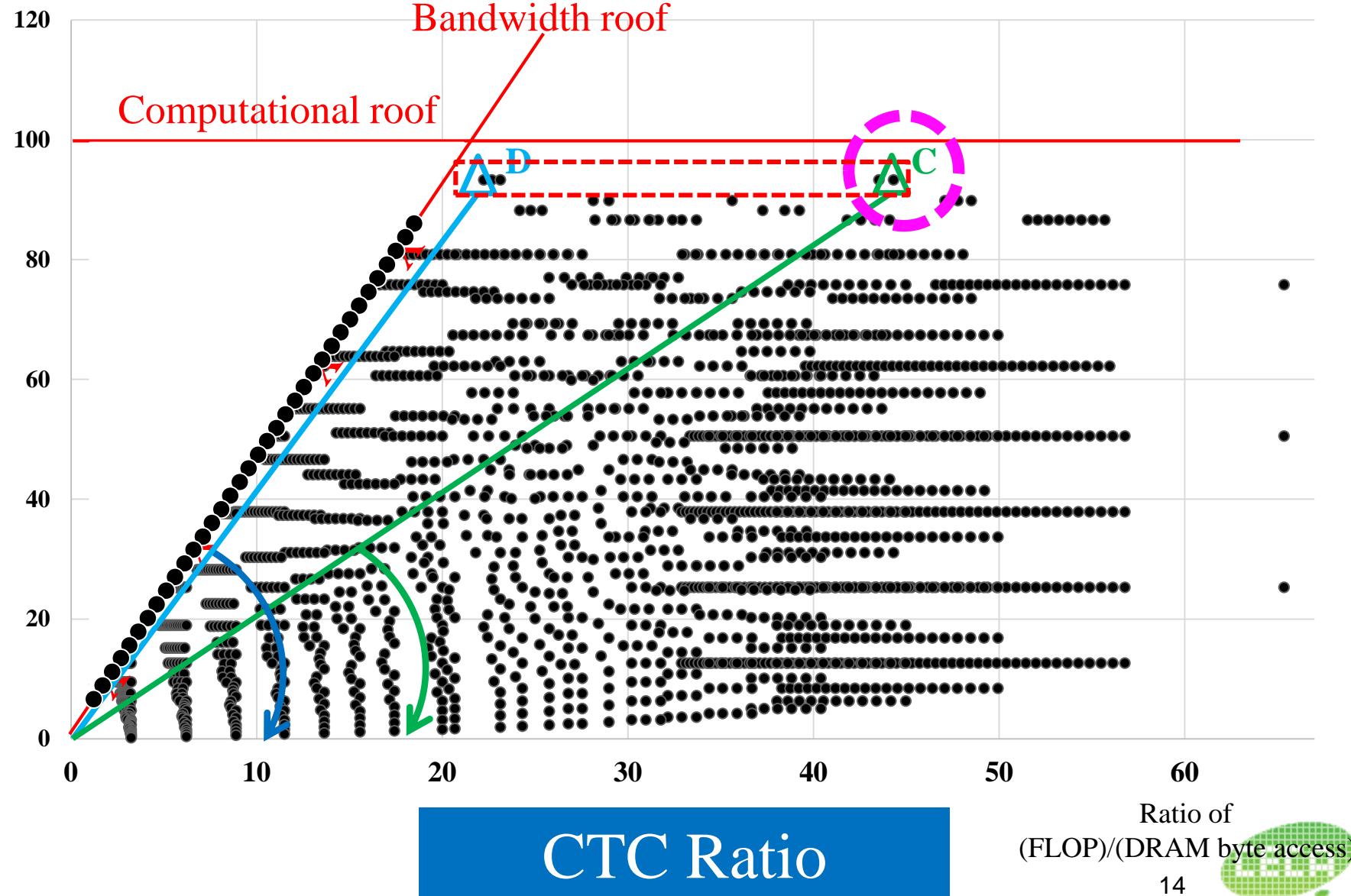
**After local memory promotion**

$$\begin{aligned}
 & \text{'output\_fm' memory access} \\
 & = 2 \times \frac{R}{Tr} \times \frac{C}{Tc} \times \frac{M}{Tm} \times \text{Size}_{\text{array}}
 \end{aligned}$$

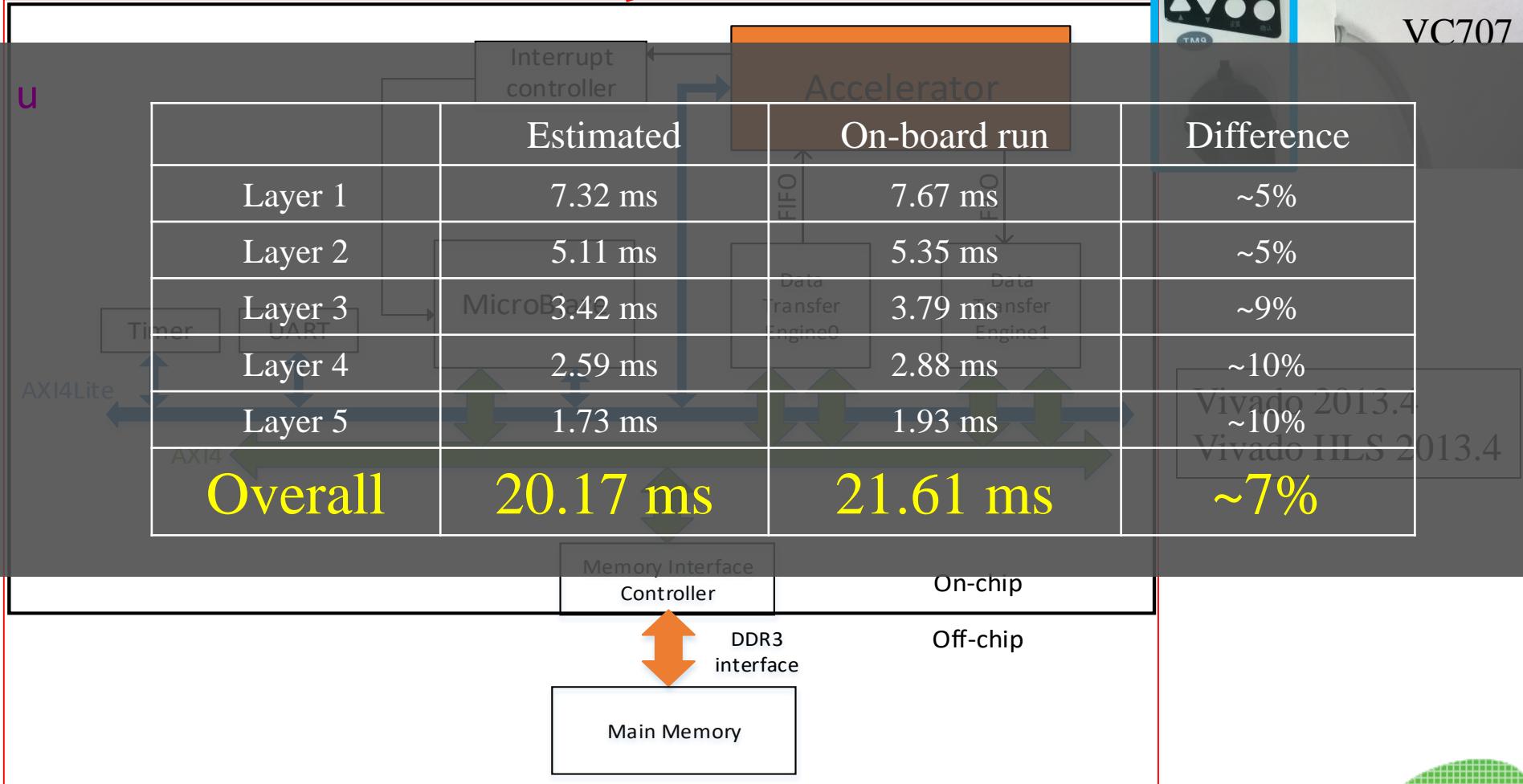
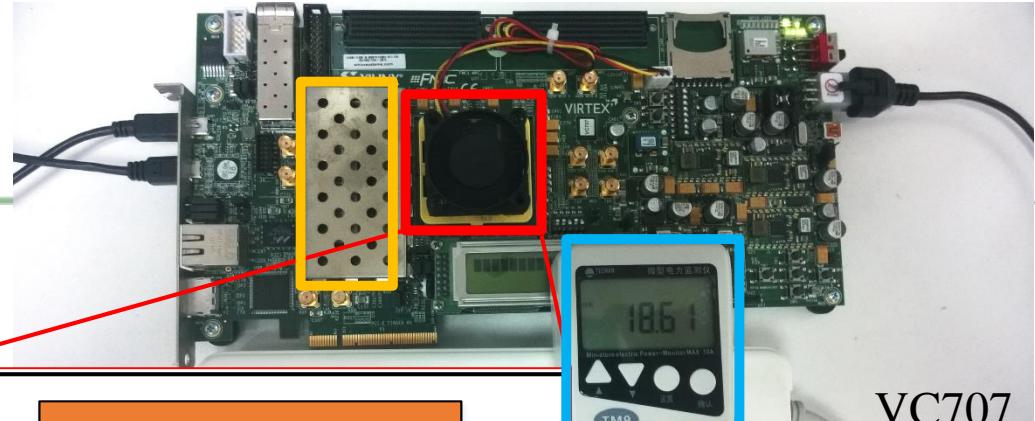
*Computation to Communication Ratio* =  $\frac{\text{Total number of operations}}{\text{Total amount of external data access}}$

# Design Space Exploration

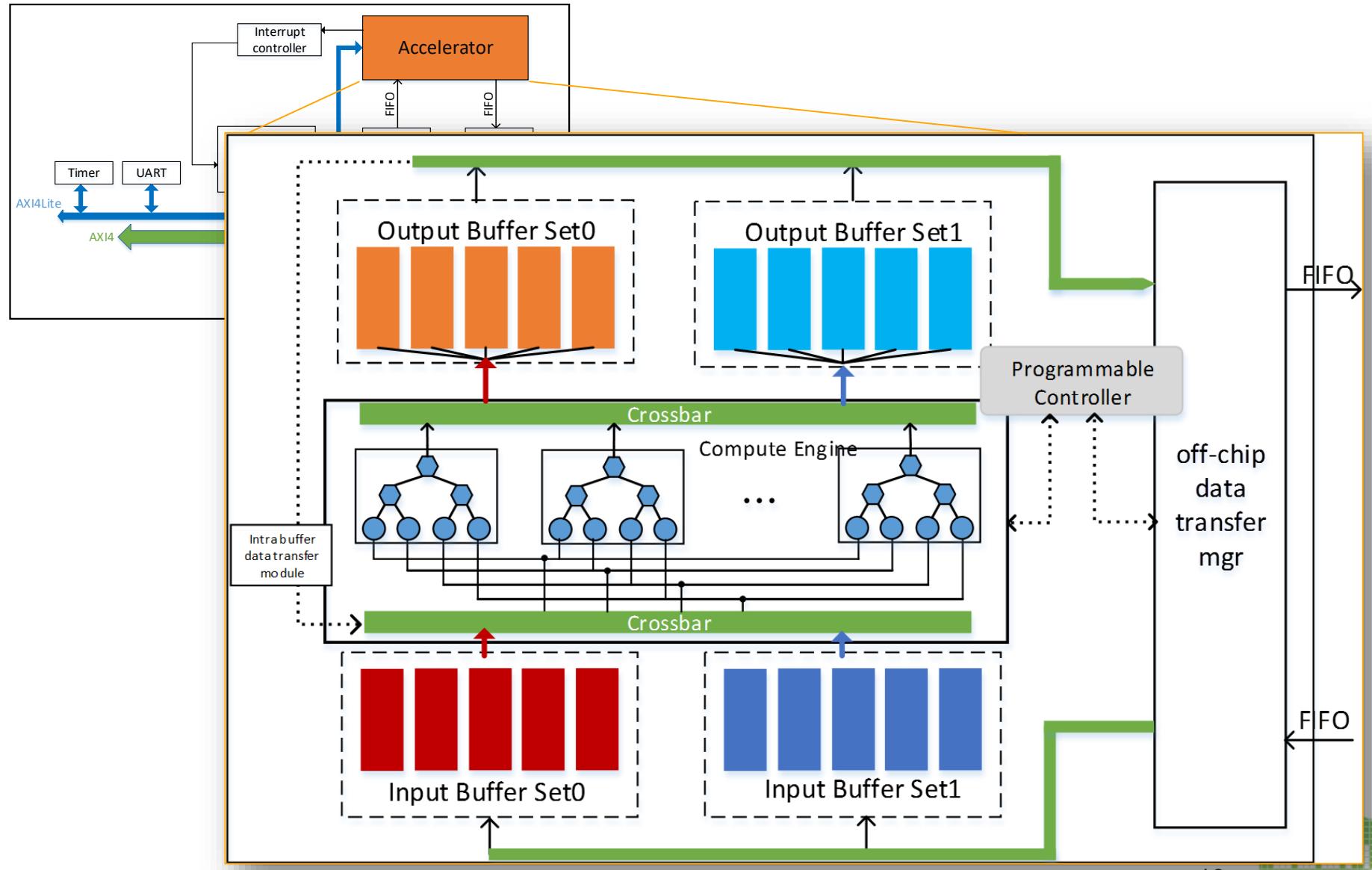
Attainable performance (GFLOPS)



# System Overview

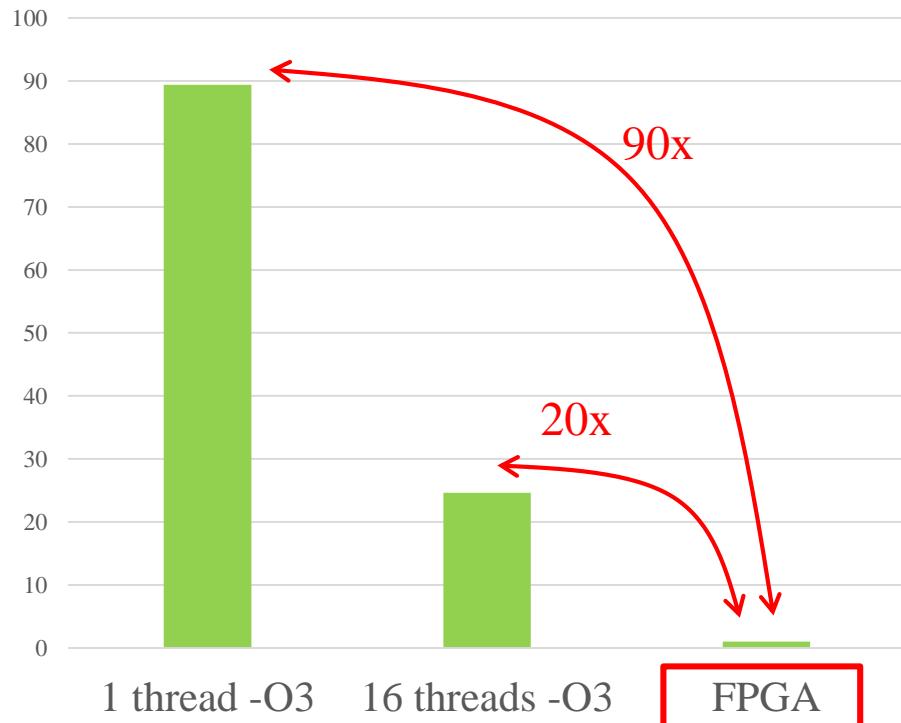


# Accelerator

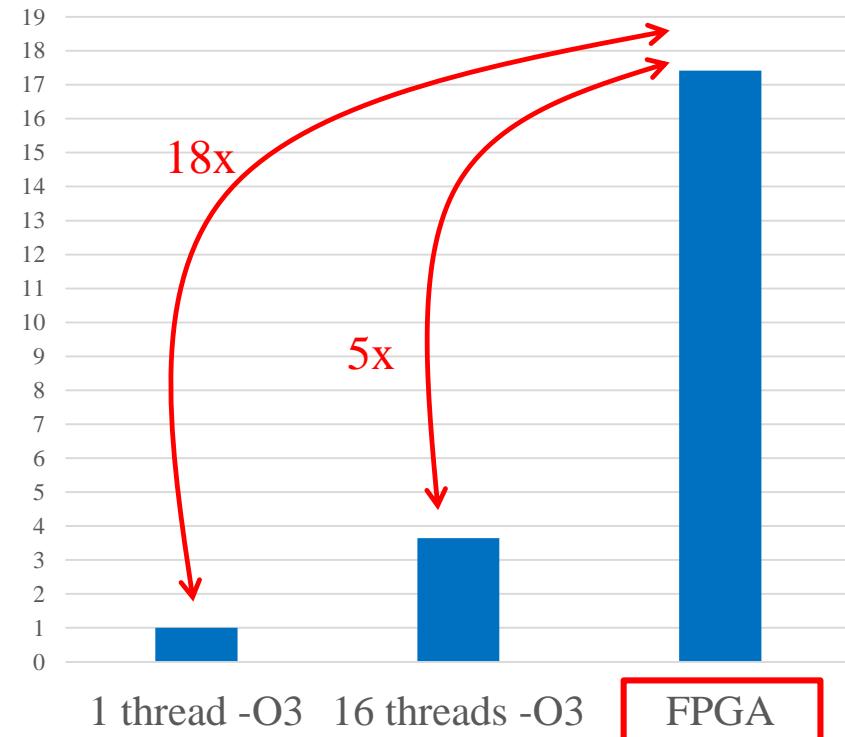


# *Experimental Results: vs. CPU*

Energy



Speedup



CPU	Xeon E5-2430 (32nm)	16 cores	2.2 GHz	gcc 4.7.2 -O3 OpenMP 3.0
FPGA	Virtex7-485t (28nm)	448 PEs	100MHz	Vivado 2013.4 Vivado HLS 2013.4

# *Mapping CNN on an FPGA-cluster*

## ◆ Problem 1

- ◆ Single FPGA board has limited number of floating-point units and memory buffers

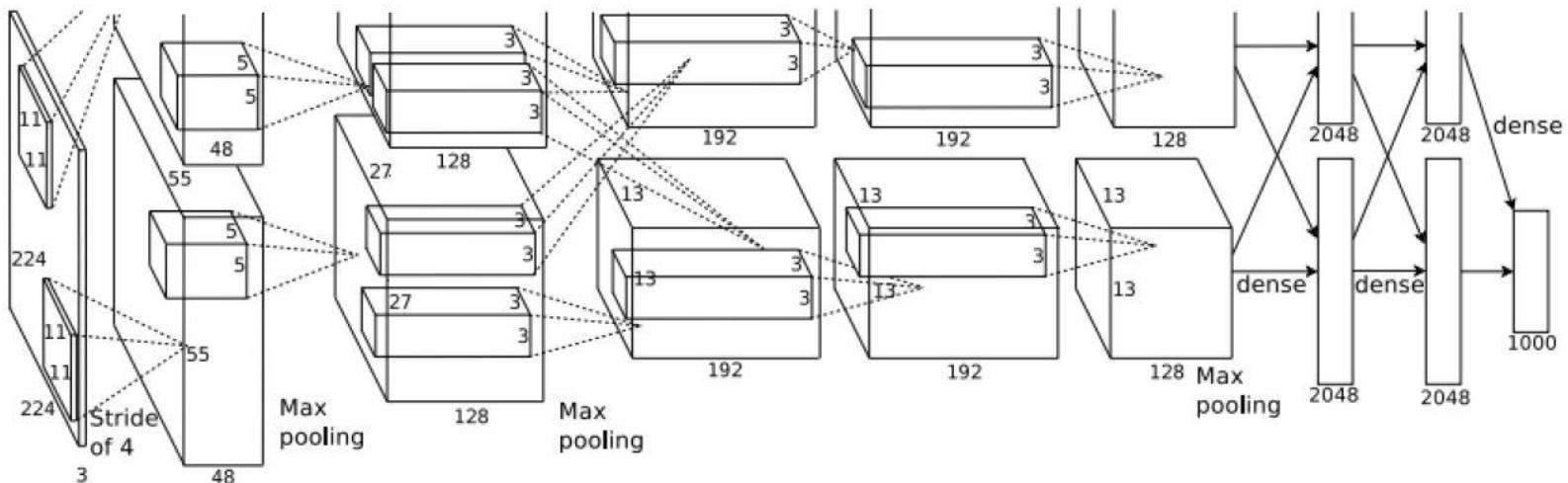
## ◆ Problem 2

- ◆ CNN computing kernels have diverse requirements for compute resource, memory bandwidth, and on-chip buffers

## ◆ Solution: Multi-board Acceleration



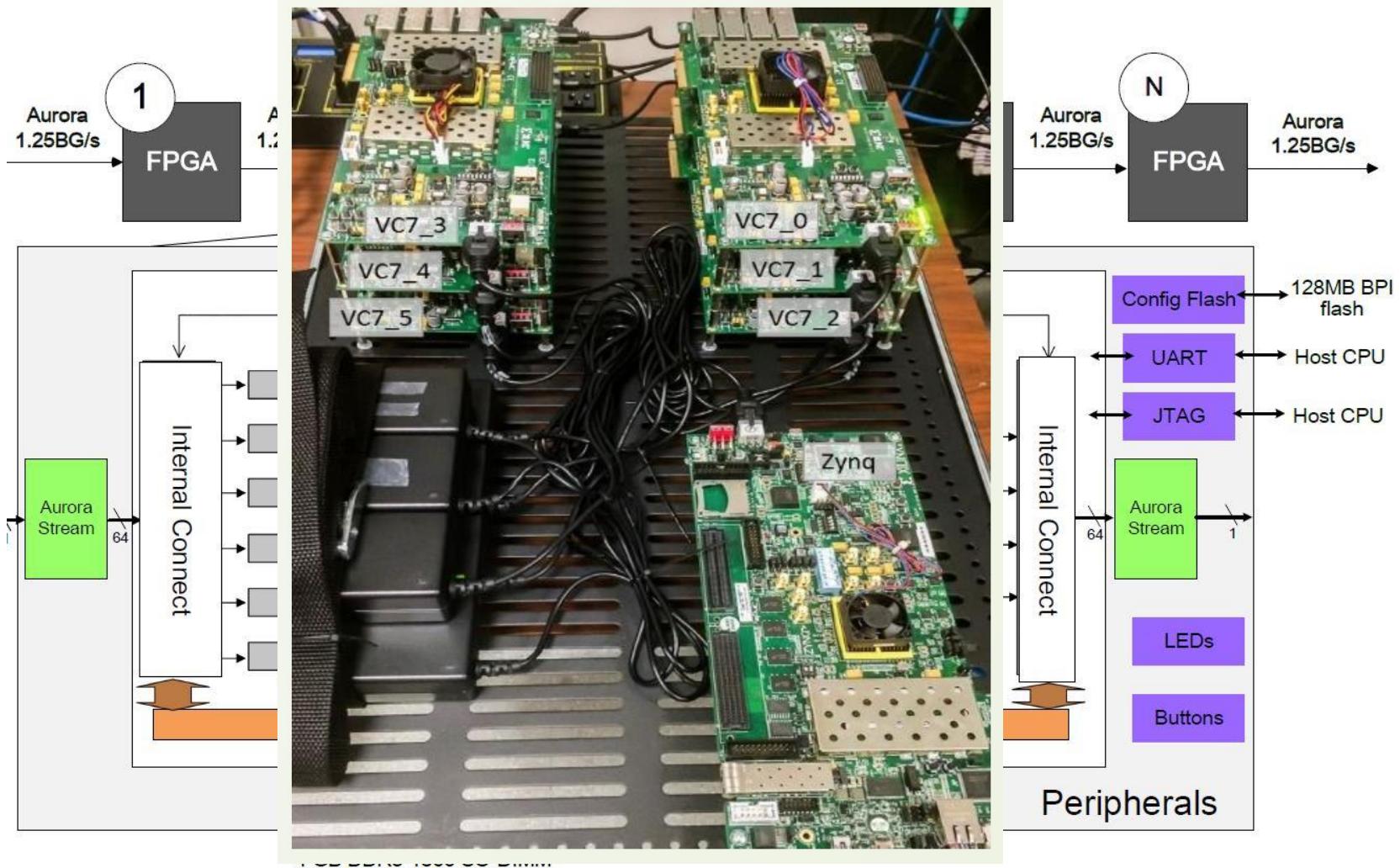
# An Example of Real-life CNN



	Convolution	LR Norm	POOL	Non Linear	Fully Connect
Total OP (10^6)	591	5.5	0.5	0.0	37.7
Percentage(%)	93.2%	0.8%	0.1%	0.0%	6.0%
	Convolution	LR Norm	POOL	Non Linear	Fully Connect
Weight Size(MB)	10	0	0	0	224
Percentage(%)	4.31%	0.0%	0.0%	0.0%	95.7%



# Overview



# Mapping Problem

## Objective 1: Minimize Latency

$$L(i, j, k) = \begin{cases} L(i, j, 1), & k = 1 \\ \min_{r=i}^{j-1} \left( \begin{array}{l} L(i, r, k-1) + \\ L(r+1, j, 1) + \\ T_{ext}(r) \end{array} \right), & k \geq 2 \end{cases}.$$

$L(i, j, 1)$  is the minimum latency of implementing layer  $i$  to layer  $j$  on one FPGA, which can be obtained using existing work.

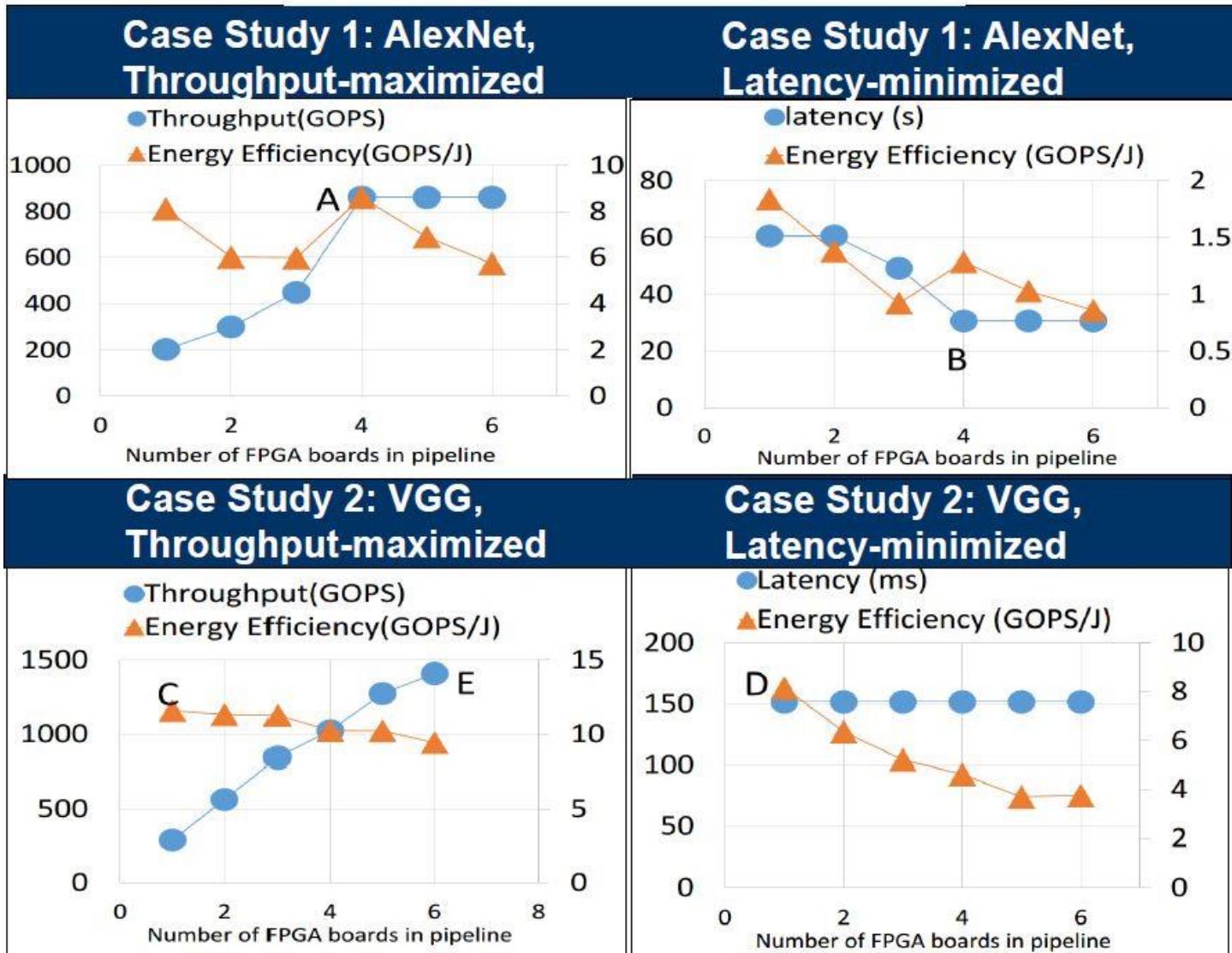
## Objective 2: Maximize Throughput

$$T(i, j, k) = \begin{cases} T(i, j, 1), & k = 1 \\ \max_{r=i}^{j-1} \min \left( \begin{array}{l} T(i, r, k-1), \\ T(r+1, j, 1), \\ \frac{1}{T_{ext}(r)} \end{array} \right), & k \geq 2 \end{cases}$$

$T(i, j, 1)$  is the maximum throughput achieved by implementing layer  $i$  to layer  $j$  on one FPGA, which can be obtained using existing work.



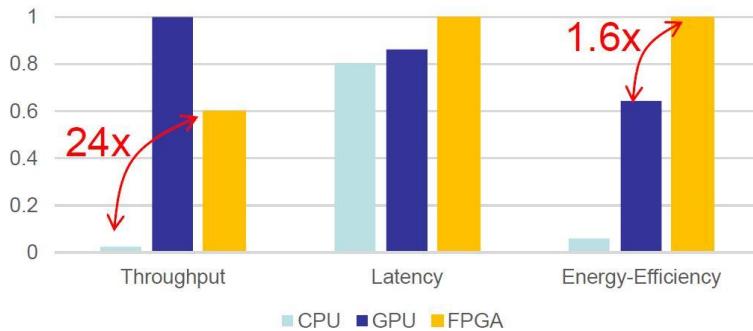
# Case Study



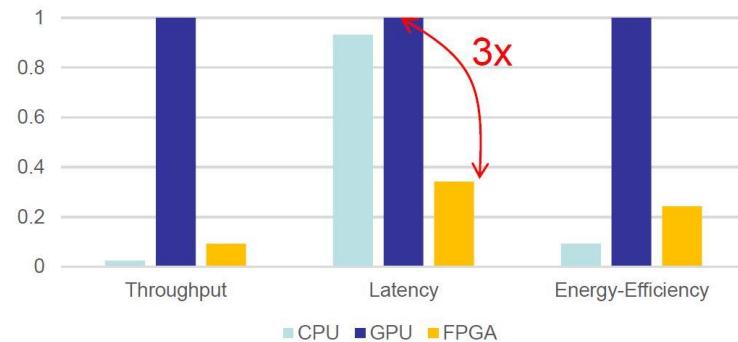
# Cross-Platform Comparison

## ◆ AlexNet

Solution A: throughput optimized

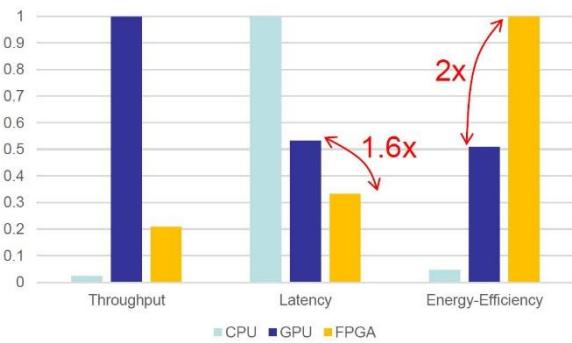


Solution B : latency optimized

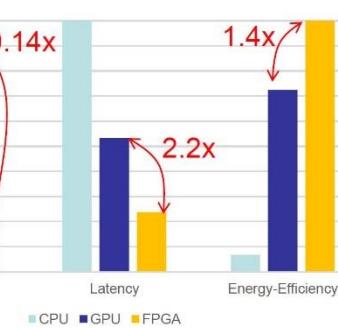


## ◆ VGG

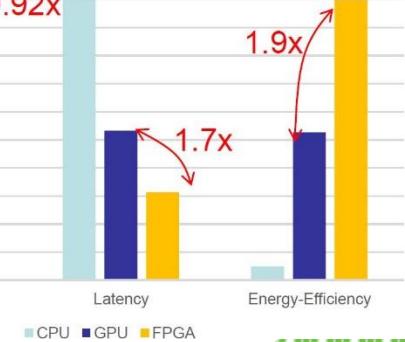
Solution C: energy-efficiency optimized



Solution D:  
latency & energy optimized



Solution E:  
throughput & energy optimized



# *Automated Design Flow*

---

## ◆ Deep CNN Framework

- Caffe, Theano, TensorFlow, etc.
- Software-based
- Support CPU/GPGPU only

## ◆ FPGA-based Accelerator

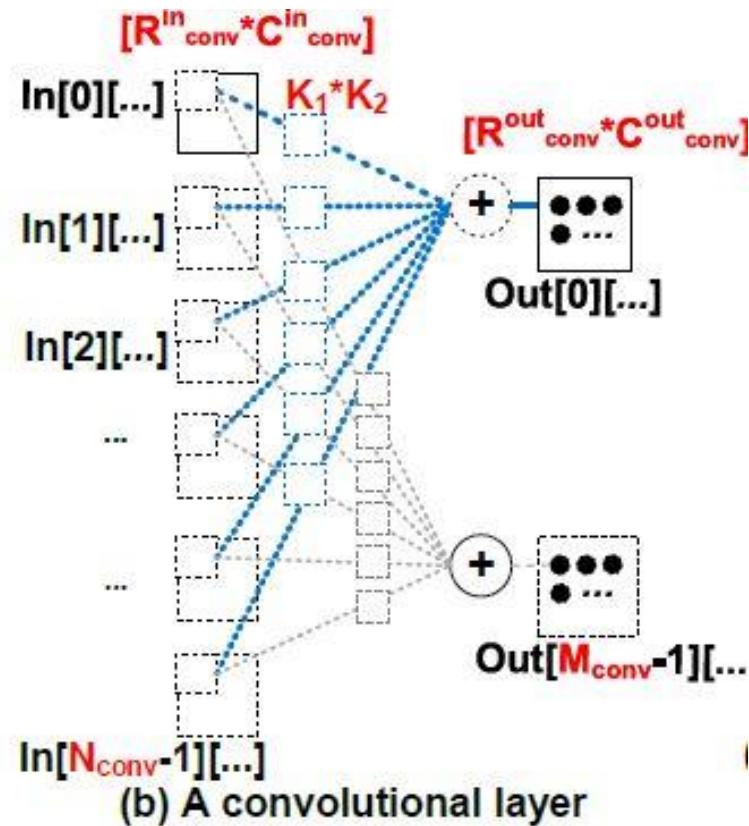
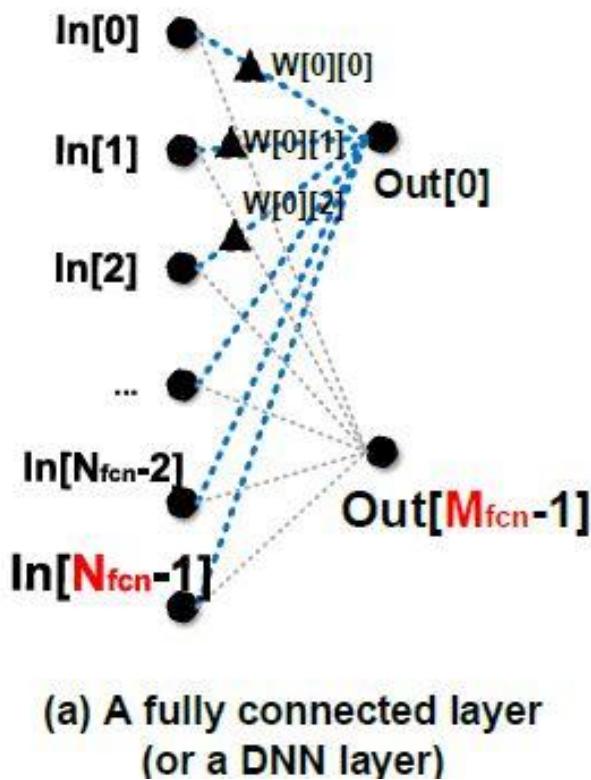
- Energy-efficiency (GOP/J)
- Performance (GOP/S)
- Flexibility

## ◆ Automated Flow

- From Model Descriptions(Caffe) to FPGA Accelerators

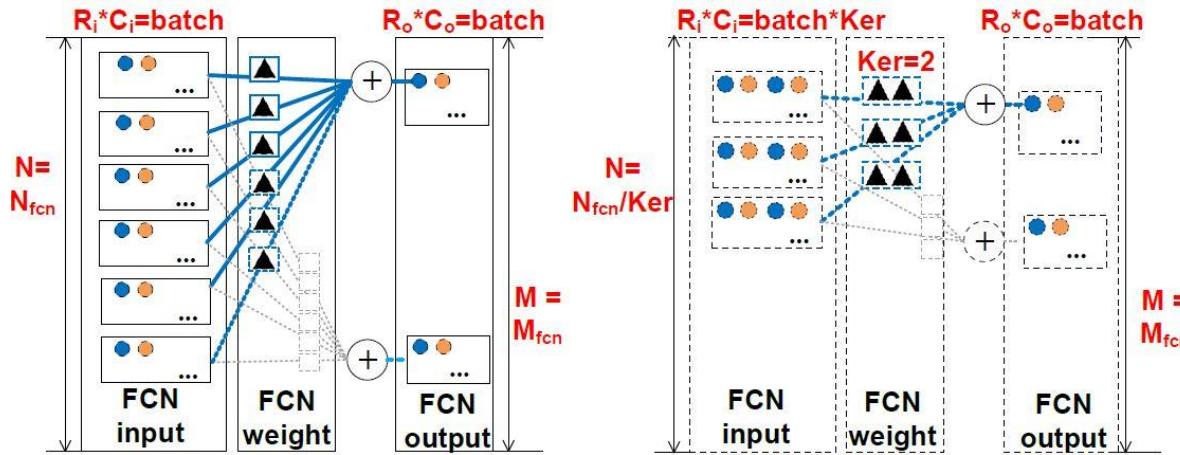


# Layer Representation (FC and Conv)

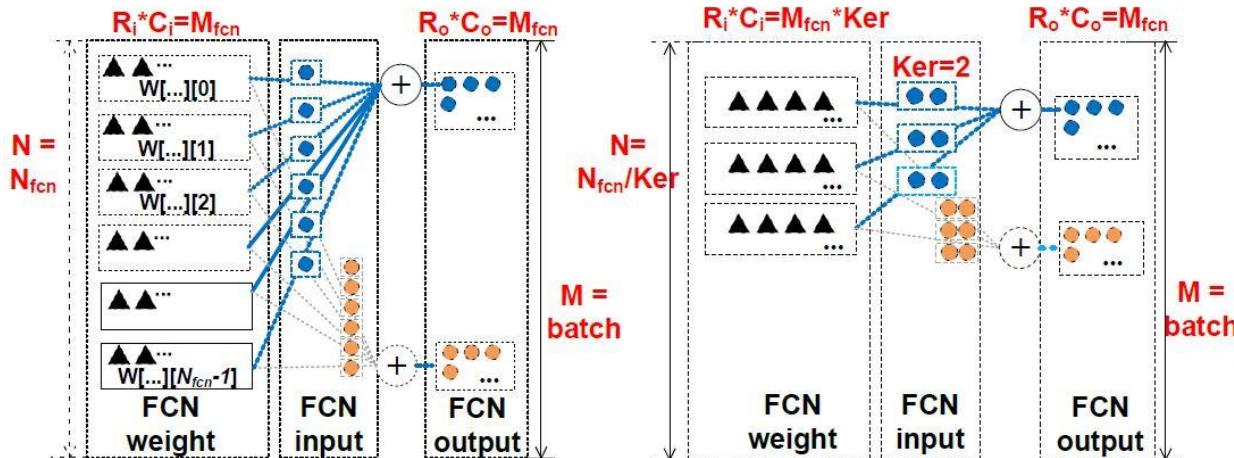


# Uniformed Representation

## Input-major Mapping



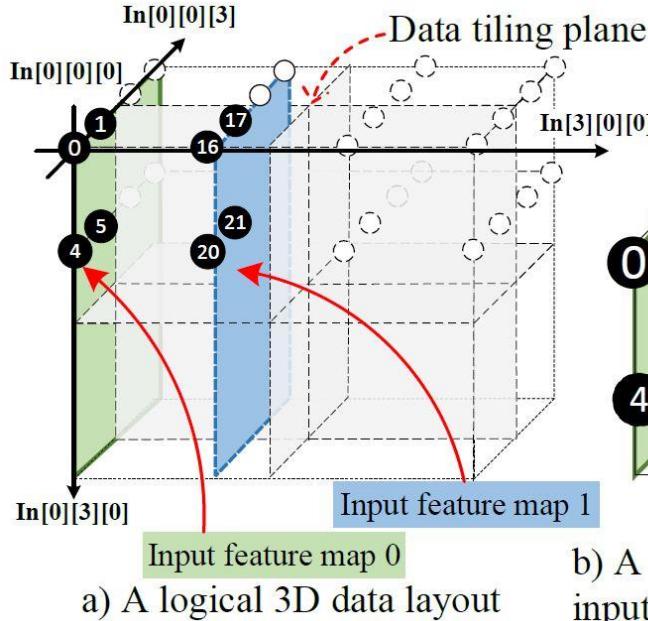
## Weight-major Mapping



For small batch sizes, Weight-major mapping is better.



# Data Layout Optimization



Data tile of input feature map 0  
is buffered in BRAM Bank 0

0	1	4	5
---	---	---	---

Data tile of input feature map 1  
is buffered in BRAM Bank 1

16	17	20	21
----	----	----	----

c) Physical data layout in on-chip buffer per BRAM bank

Data	0	1	...	4	5	...	...	16	17	...	20	21
DRAM Addr.	x0	x4	...	x10	x14	...	...	x40	x44	...	x50	x54

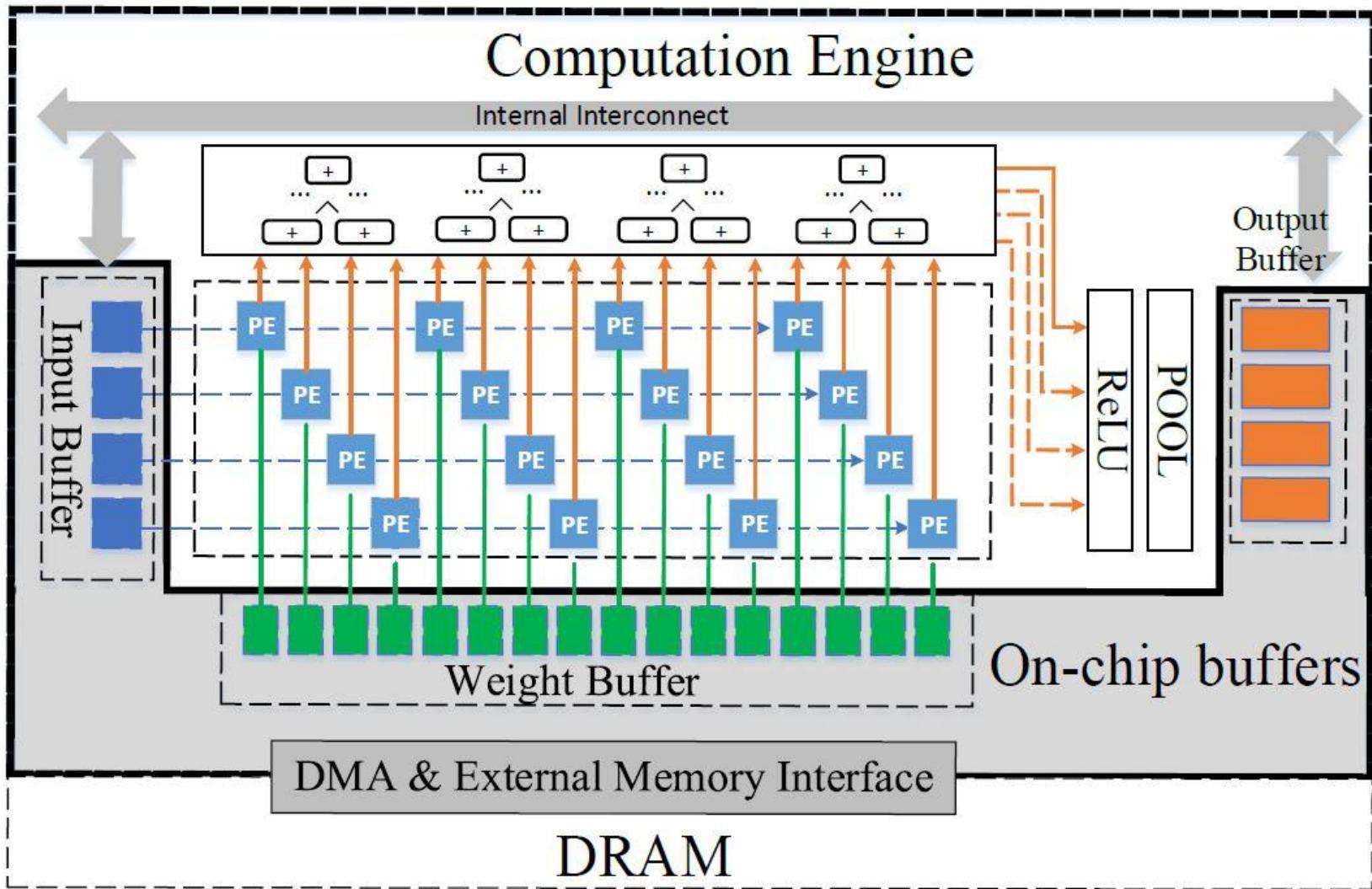
d) Row-major data layout in DRAM space

Data	0	16	1	17	4	20	5	21	...	...	...	...
DRAM Addr.	x0	x4	x8	xc	x10	x14	x18	x1c	...	...	...	...

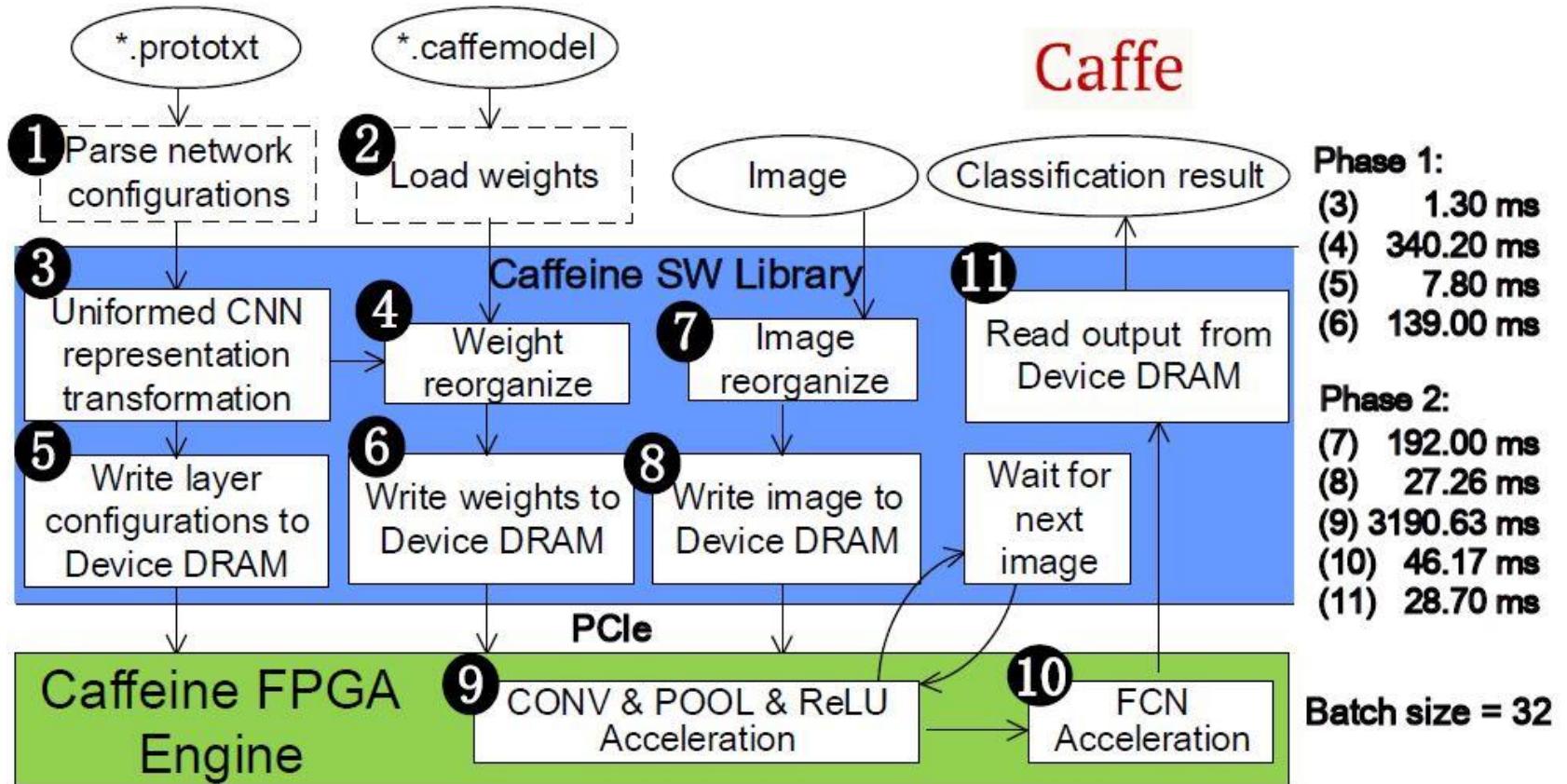
e) Proposed data layout in DRAM space



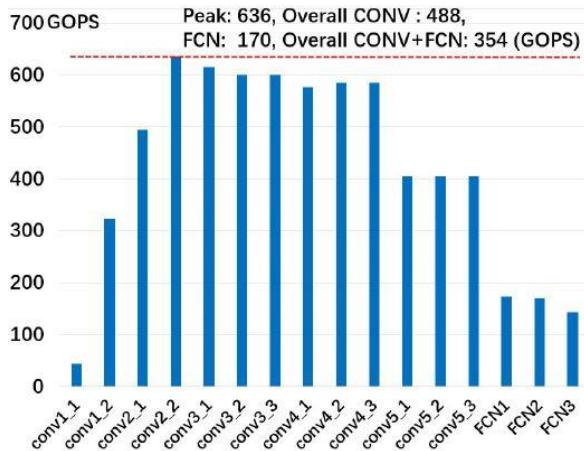
# System Overview



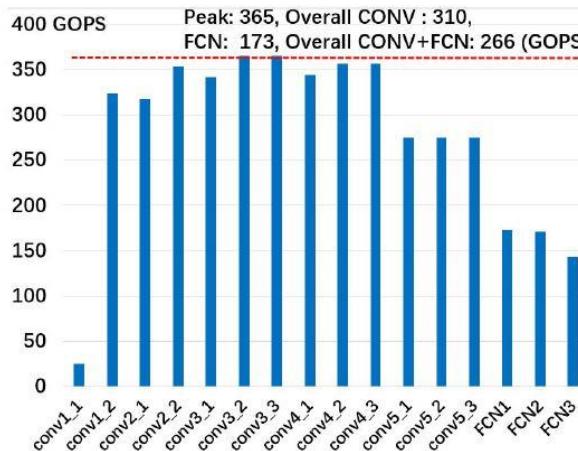
# Caffe-Caffeine integration



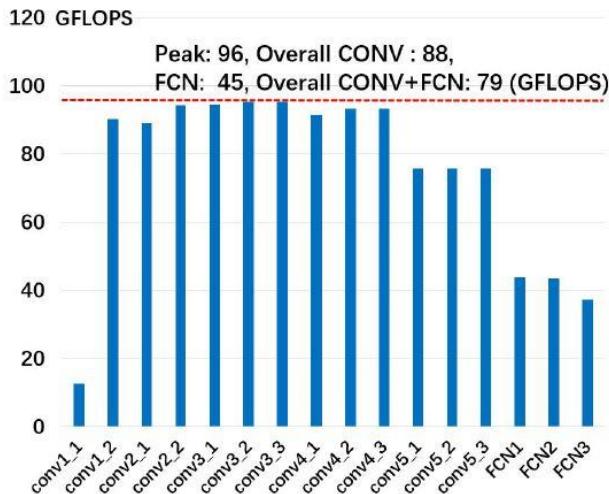
# Experimental Results



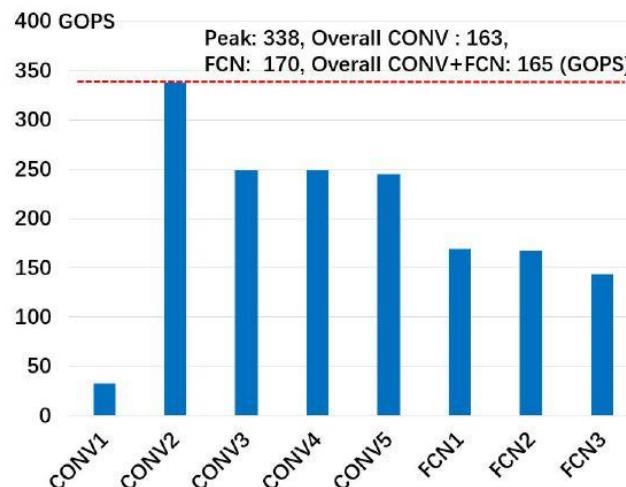
(a) VC709 VGG 16-bit fixed-point



(b) KU VGG 16-bit fixed-point



(c) KU VGG 32-bit floating-point



(d) KU AlexNet 16-bit fixed-point



# Cross-platform Comparison

platforms	CPU	CPU+GPU	CPU+FPGA
Device	E5-2609	K40	KU060
Technology	22nm	28nm	20nm
Freq.	1.9GHz	1GHz	200MHz
Power(Watt)	150	250	25
Latency (ms) per img.	733.7	15.3	101.15
Speedup	1x	48x	7.3x
J per image	110	3.8	2.5
Energy Efficiency	1x	28.7x	43.5x
			65x



# **Conclusions**

---

- „ **FPGA based CNN accelerators are promising**
- „ **Design space should be carefully explored on single/multiple FPGA boards**
- „ **Design automation flow is helpful**



# Acknowledgement

## Faculty Members:



**Jason Cong**  
(UCLA)



**Guangyu Sun**  
(PKU)

## PhD Students:



**Chen Zhang**  
(PKU)



**Yijin Guan**  
(PKU)



**Bingzhe Wu**  
(PKU)

## Undergraduate Students:



**Zhihang Yuan**  
(PKU)



**Yao Fu**  
(PKU)



**Sixue Wang**  
(PKU)



**Thanks!**

