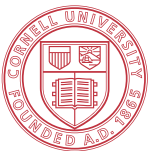




ECE 6775  
High-Level Digital Design Automation  
Fall 2023

# DNN Acceleration on FPGAs



Cornell University



# Announcements

- ▶ Lab 4 on BNN acceleration
  - Use the latest zip file from the class folder
- ▶ Start forming teams for the final project
  - 3 or 4 students per team

# Course Roadmap

## ▶ **Background**

- Introduction
- Hardware specialization
- Algorithm basics

## ▶ **High-level synthesis**

- C-based synthesis for FPGAs
- Front-end compilation
- Scheduling
- Resource sharing
- Pipelining

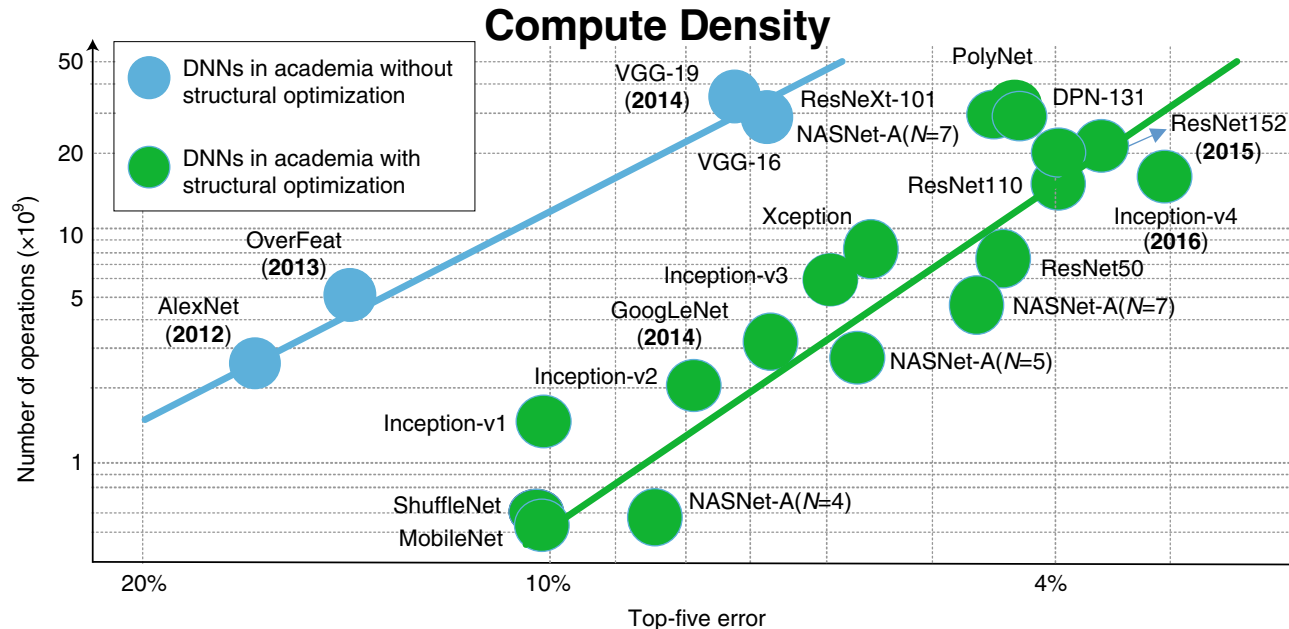
## ▶ **Advanced topics**

- Deep learning acceleration
- Domain-specific programming

# Agenda

- ▶ CNN acceleration on FPGAs
- ▶ Overview of BNN acceleration

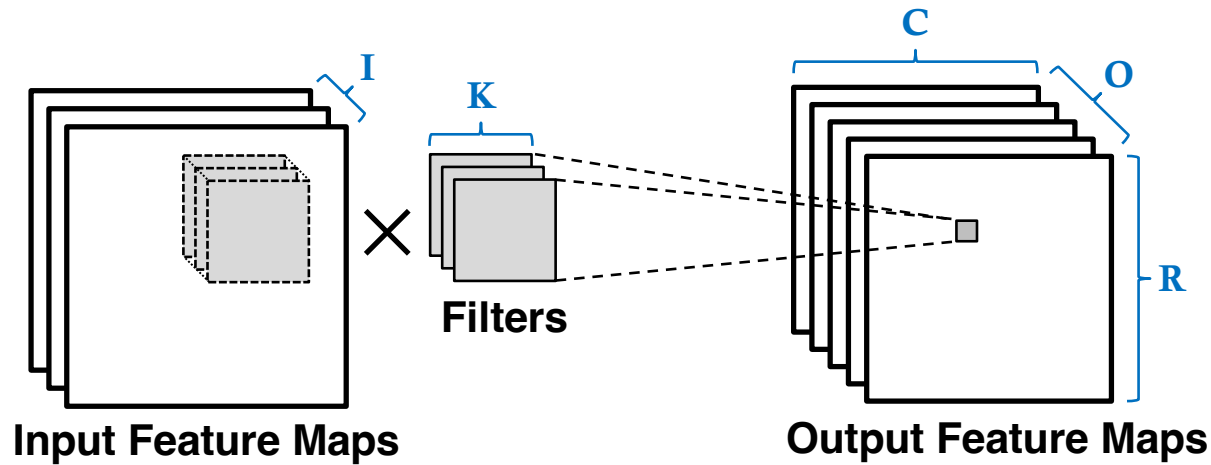
# Modern ML Models are Computationally Expensive



X. Xu, Y. Ding, S. X. Hu, M. Niemier, J. Cong, Y. Hu, and Y. Shi. **Scaling for Edge Inference of Neural Networks**. *Nature Electronics*, vol 1, Apr 2018.

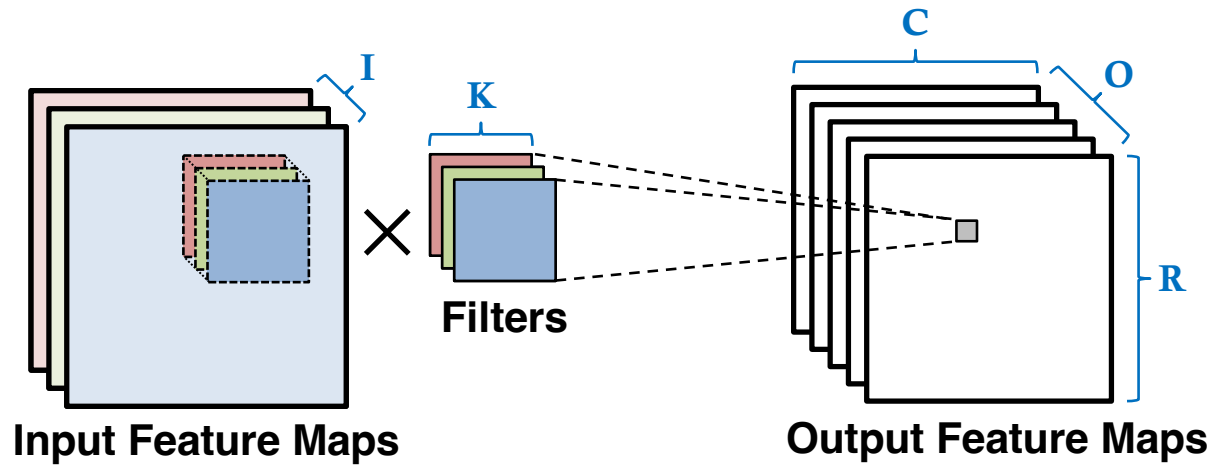
- ▶ For example, ResNet50, a 70-layer CNN model performs 7.7 billion operations required to classify one image

# Convolutional Layer



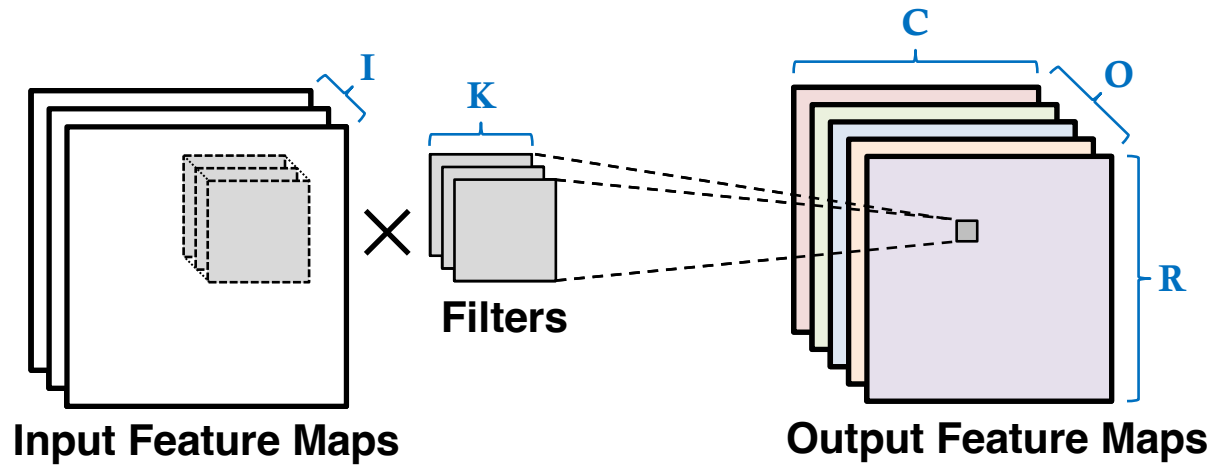
- ▶ An output pixel is connected to its neighboring region on each input feature map (fmap)
- ▶ All pixels on an output feature map use the same filter weights

# Parallelism in the Convolutional Layer



- ▶ Four main sources of parallelism
  1. **Across input feature maps (i.e., input channels)**

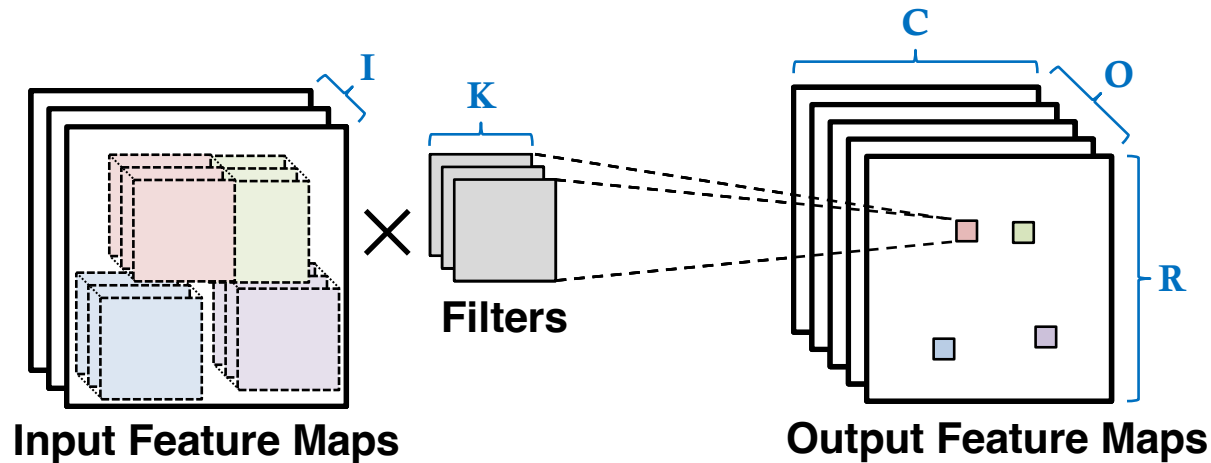
# Parallelism in the Convolutional Layer



- ▶ Four main sources of parallelism
  1. Across input feature maps (i.e., input channels)
  2. **Across output feature maps (i.e., output channels)**

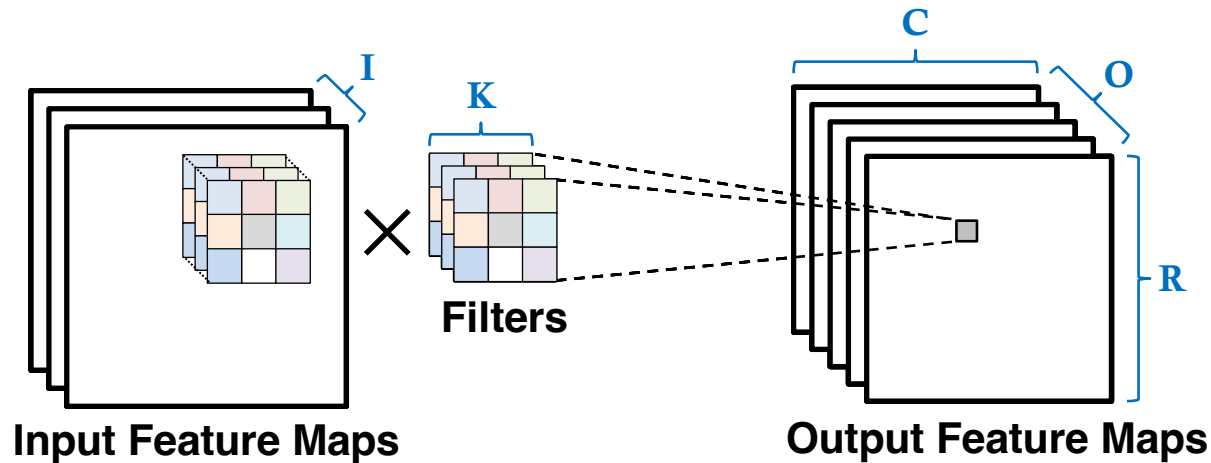


# Parallelism in the Convolutional Layer



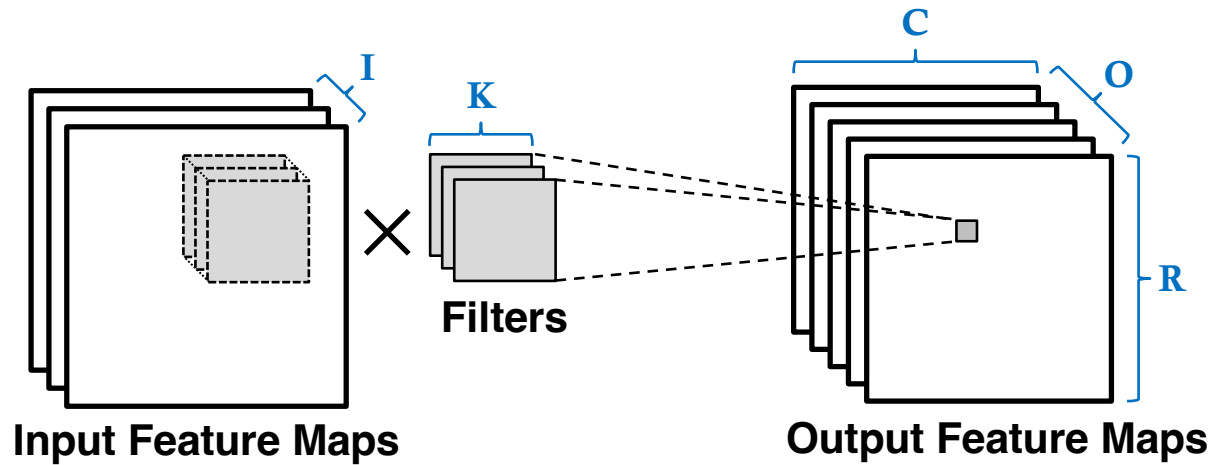
- ▶ Four main sources of parallelism
  1. Across input feature maps (i.e., input channels)
  2. Across output feature maps (i.e., output channels)
  3. **Across different output pixels (i.e., filter positions)**

# Parallelism in the Convolutional Layer



- ▶ Four main sources of parallelism
  1. Across input feature maps (i.e., input channels)
  2. Across output feature maps (i.e., output channels)
  3. Across different output pixels (i.e., filter positions)
  4. **Across filter pixels**

# Parallelism in the Code



```

1 for (row=0; row<R; row++) {
2   for (col=0; col<C; col++) {
3     for (to=0; to<O; to++) {
4       for (ti=0; ti<I; ti++) {
5         for (ki=0; ki<K; ki++) {
6           for (kj=0; kj<K; kj++) {
              output_fm[to][row][col] +=
              weights[to][ti][ki][kj]*input_fm[ti][S*row+ki][S*col+kj];
            }
          }
        }
      }
    }
  }
}

```

} Loop over output pixels  
 } Loop over output channels  
 } Loop over input channels  
 } Loop over filter pixels  
 ↓ S is the stride

# Optimizing FPGA-based Accelerator Design for Deep Convolutional Neural Networks

Cheng Zhang<sup>1</sup>, Peng Li<sup>2</sup>, Guangyu Sun<sup>1,3</sup>, Yijin Guan<sup>1</sup>, Bingjun Xiao<sup>2</sup>, Jason Cong<sup>2,3,1</sup>

<sup>1</sup>Center for Energy-Efficient Computing and Applications, Peking University

<sup>2</sup>Computer Science Department, UCLA

<sup>3</sup>PKU/UCLA Joint Research Institute in Science and Engineering

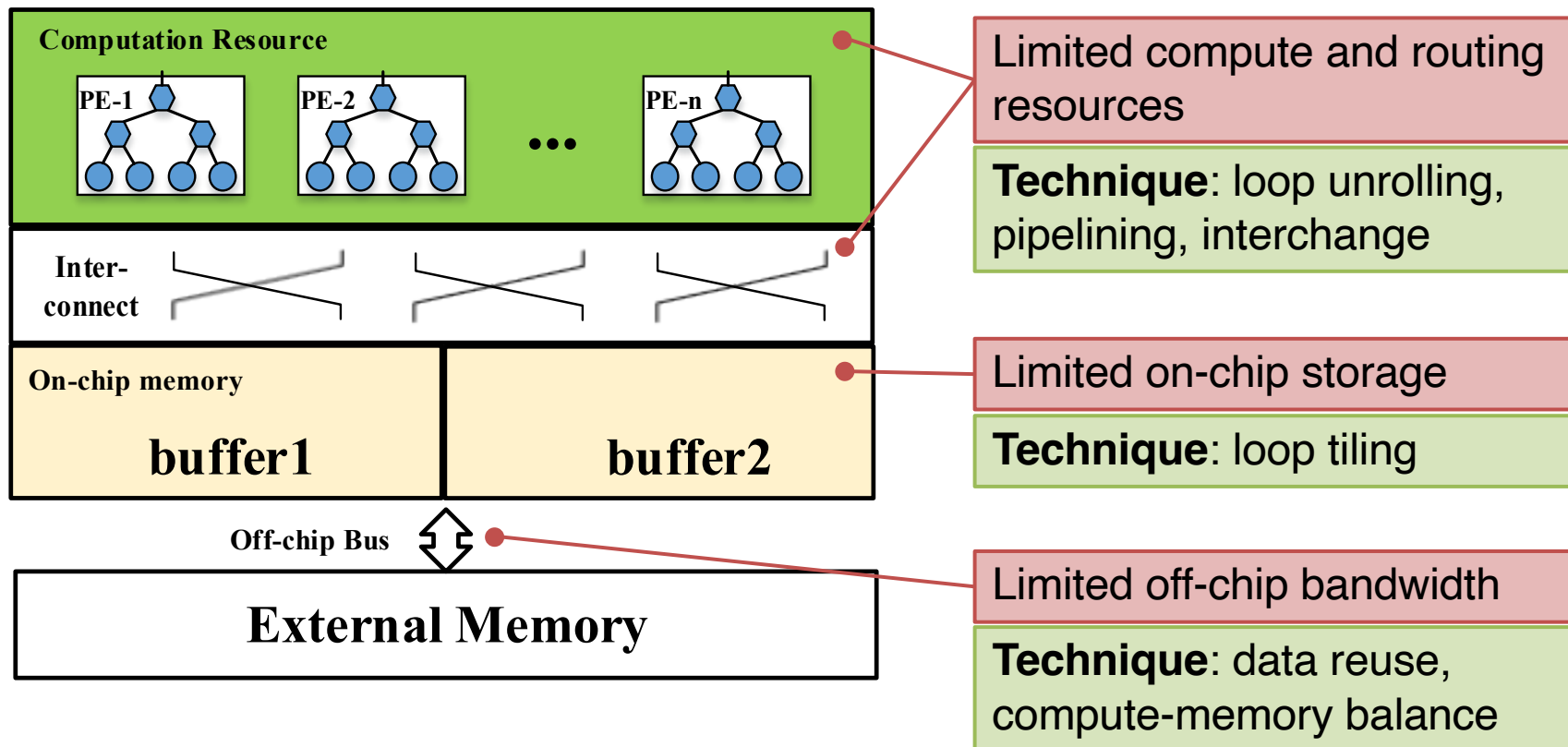
***FPGA'15, Feb 2015***

# Main Contributions

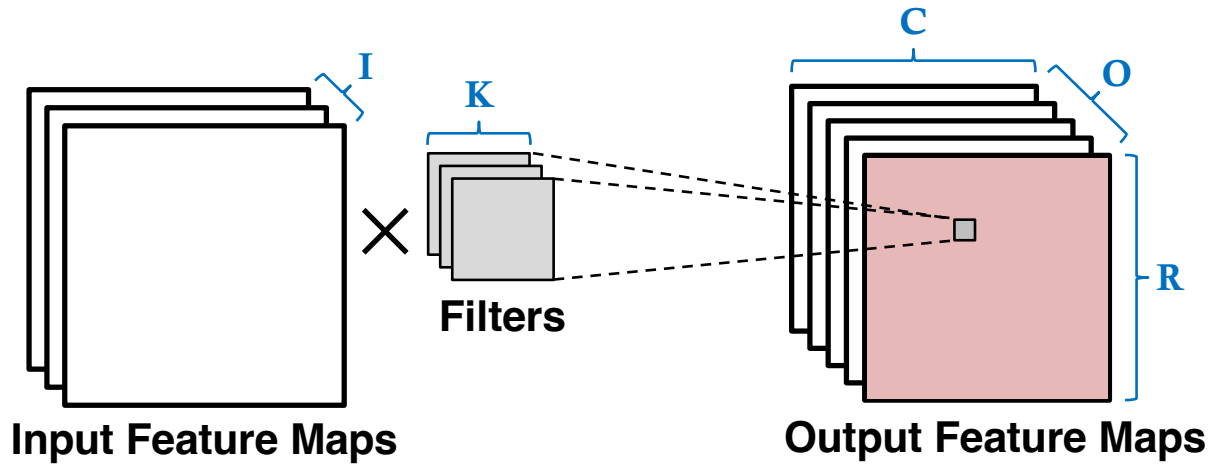
1. Analysis of the different sources of parallelism in the convolution kernel of a CNN
2. Quantitative performance modeling of the hardware design space using the Roofline method
3. Design and implementation of a CNN accelerator for FPGA using Vivado HLS, evaluated on AlexNet

# Challenges to FPGA Acceleration

- ▶ We can't just unroll all the loops due to limited FPGA resources
- ▶ Must choose the right code transformations to exploit the parallelism in a resource efficient way



# Loop Tiling

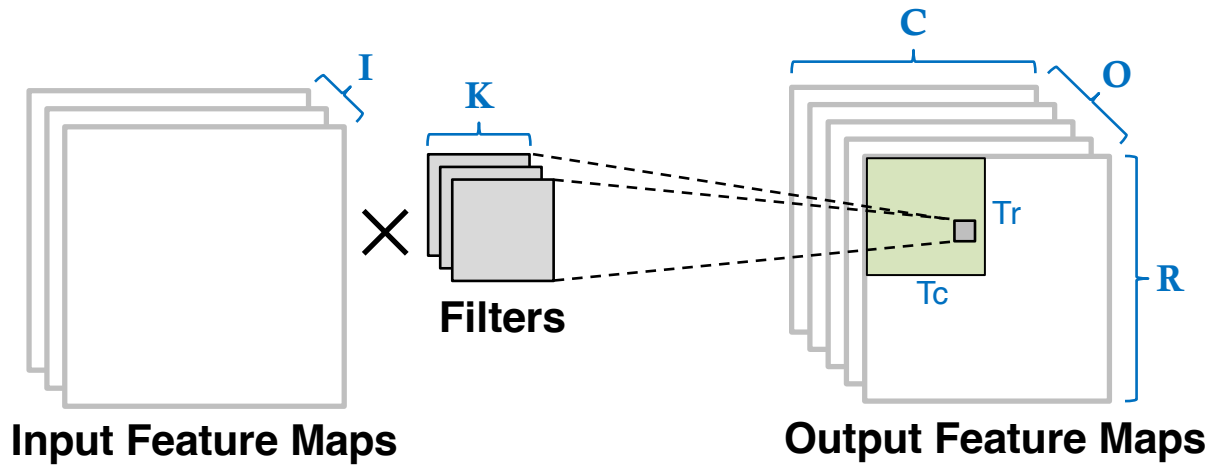


```
1 for (row=0; row<R; row++) {  
2   for (col=0; col<C; col++) {  
3     for (to=0; to<O; to++) {  
4       for (ti=0; ti<I; ti++) {  
5         for (ki=0; ki<K; ki++) {  
6           for (kj=0; kj<K; kj++) {  
               output_fm[to][row][col] +=  
               weights[to][ti][ki][kj]*input_fm[ti][S*row+ki][S*col+kj];  
           }  
         }  
       }  
     }  
   }  
}
```

} Loop over **pixels** in an output fmap

# Loop Tiling

Offloading just the inner loops requires only a small portion of the data to be stored on FPGA chip



```
1 for (row=0; row<R; row+=Tr) {  
2   for (col=0; col<C; col+=Tc) {  
3     for (to=0; to<O; to++) {  
4       for (ti=0; ti<I; ti++) {  
5         for (trr=row; trr<min(row+Tr, R); trr++) {  
6           for (tcc=col; tcc<min(col+Tc, C); tcc++) {  
7             for (ki=0; ki<K; ki++) {  
8               for (kj=0; kj<K; kj++) {  
                 output_fm[to][trr][tcc] +=  
                 weights[to][ti][ki][kj]*input_fm[ti][S*trr+ki][S*tcc+kj];  
               }  
             }  
           }  
         }  
       }  
     }  
   }  
}
```

Loop over different tiles

Loops over pixels in each tile



# Code with Loop Tiling

**Tile sizes:** rows ( $T_r$ ), columns ( $T_c$ ),  
input channels ( $T_i$ ), and output channels ( $T_o$ )

```
1 for (row=0; row<R; row+=Tr) {
2   for (col=0; col<C; col+=Tc) {
3     for (to=0; to<O; to+=To) {
4       // software: write output feature map
5       for (ti=0; ti<I; ti+=Ti) {
6         // software: write input feature map + filters
7
8         for (trr=row; trr<min(row+Tr, R); trr++) {
9           for (tcc=col; tcc<min(col+Tc, C); tcc++) {
10            for (too=to; too<min(to+To, O); too++) {
11              for (tii=ti; tii<(ti+Ti, I); tii++) {
12                for (ki=0; ki<K; ki++) {
13                  for (kj=0; kj<K; kj++) {
14                    output_fm[too][trr][tcc] +=
15                      weights[too][tii][ki][kj]*input_fm[tii][S*trr+ki][S*tcc+kj];
16                }
17              }
18            }
19          }
20        }
21      }
22      // software: read output feature map
23    }
24  }
25 }
```

**CPU Portion**

**FPGA Portion**

# Optimizing for On-Chip Performance

```
5     for (trr=row; trr<min(row+Tr, R); trr++) {
6         for (tcc=col; tcc<min(col+Tc, C); tcc++) {
7             for (too=to; too<min(to+To, O); too++) {
8                 for (tii=ti; tii<(ti+Ti, I); tii++) {
9                     for (ki=0; ki<K; ki++) {
10                        for (kj=0; kj<K; kj++) {
                            output_fm[too][trr][tcc] +=
                                weights[too][tii][ki][kj]*input_fm[tii][S*trr+ki][S*tcc+kj];
                        }
                    }
                }
            }
        }
    }
```

# Optimizing for On-Chip Performance

```
9      for (ki=0; ki<K; ki++) {      Reorder these two loops to the top level
10     for (kj=0; kj<K; kj++) {
5       for (trr=row; trr<min(row+Tr, R); trr++) {
6         for (tcc=col; tcc<min(col+Tc, C); tcc++) {
7           for (too=to; too<min(to+To, O); too++) {
8             for (tii=ti; tii<(ti+Ti, I); tii++) {
                output_fm[too][trr][tcc] +=
                weights[too][tii][ki][kj]*input_fm[tii][S*trr+ki][S*tcc+kj];
            }
        }
    }
}
}}}}}}
```

# Optimizing for On-Chip Performance

```
9     for (ki=0; ki<K; ki++) {
10        for (kj=0; kj<K; kj++) {
5           for (trr=row; trr<min(row+Tr, R); trr++) {
6              for (tcc=col; tcc<min(col+Tc, C); tcc++) {
3                 #pragma HLS pipeline
7                    for (too=to; too<min(to+To, O); too++) {
8                       for (tii=ti; tii<(ti+Ti, I); tii++) {
11                          output_fm[too][trr][tcc] +=
12                             weights[too][tii][ki][kj]*input_fm[tii][S*trr+ki][S*tcc+kj];
13                      }
14                  }
15              }
16          }
17      }
18  }
```

# Optimizing for On-Chip Performance

```
9      for (ki=0; ki<K; ki++) {
10     for (kj=0; kj<K; kj++) {
5      for (trr=row; trr<min(row+Tr, R); trr++) {
6      for (tcc=col; tcc<min(col+Tc, C); tcc++) {
      #pragma HLS pipeline
7      for (too=to; too<min(to+To, O); too++) {
      #pragma HLS unroll
8      for (tii=ti; tii<(ti+Ti, I); tii++) {
        output_fm[too][trr][tcc] +=
        weights[too][tii][ki][kj]*input_fm[tii][S*trr+ki][S*tcc+kj];
      }}}}}}
```

# Optimizing for On-Chip Performance

```
9     for (ki=0; ki<K; ki++) {
10        for (kj=0; kj<K; kj++) {
5           for (trr=row; trr<min(row+Tr, R); trr++) {
6              for (tcc=col; tcc<min(col+Tc, C); tcc++) {
3                 #pragma HLS pipeline
7                    for (too=to; too<min(to+To, O); too++) {
8                       #pragma HLS unroll
11                      for (tii=ti; tii<(ti+Ti, I); tii++) {
12                         #pragma HLS unroll
13                            output_fm[too][trr][tcc] +=
14                               weights[too][tii][ki][kj]*input_fm[tii][S*trr+ki][S*tcc+kj];
15                      }
16              }
17          }
18      }
19  }
```

Parallelize across input (Ti)  
and output (To) channels

Number of cycles to execute the above loop nest

$$\approx K \times K \times Tr \times Tc + L \approx Tr \times Tc \times K^2$$

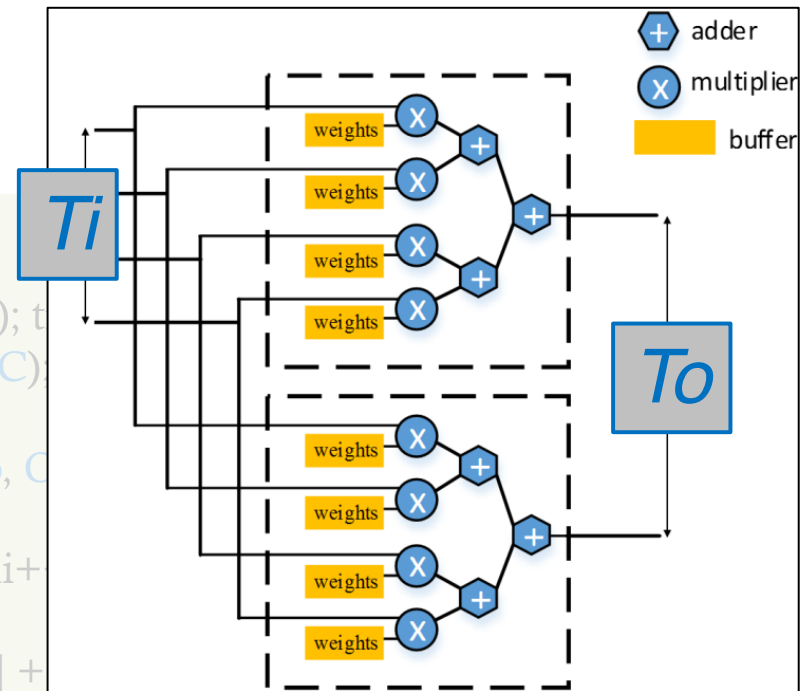
$L$  is the pipeline depth (# of pipeline stages,  $l=1$ )

# Optimizing for On-Chip Performance

```

9   for (ki=0; ki<K; ki++) {
10  for (kj=0; kj<K; kj++) {
11  for (trr=row; trr<min(row+Tr, R); trr++) {
12  for (tcc=col; tcc<min(col+Tc, C); tcc++) {
13  #pragma HLS pipeline
14  for (too=to; too<min(to+To, C); too++) {
15  #pragma HLS unroll
16  for (tii=ti; tii<(ti+Ti, I); tii++) {
17  #pragma HLS unroll
18  output_fm[too][trr][tcc] +
19  weights[too][tii][ki][kj]*input_fm[tii][S*trr+ki][S*tcc+kj];
20  }}}}}

```



## Generated Hardware

Performance and size of the accelerator determined by tile factors  $T_i$  and  $T_o$

Number of data transfers determined by  $T_r$  and  $T_c$

# Design Space Complexity

- ▶ **Challenge:** Number of available optimizations present a huge space of possible designs
  - What is the optimal loop order?
  - What tile size to use for each loop?
- ▶ Implementing and testing each design by hand will be slow and error-prone
  - Some designs will exceed the on-chip compute/memory capacity
- ▶ **Solution:** Performance modeling + automated design space exploration



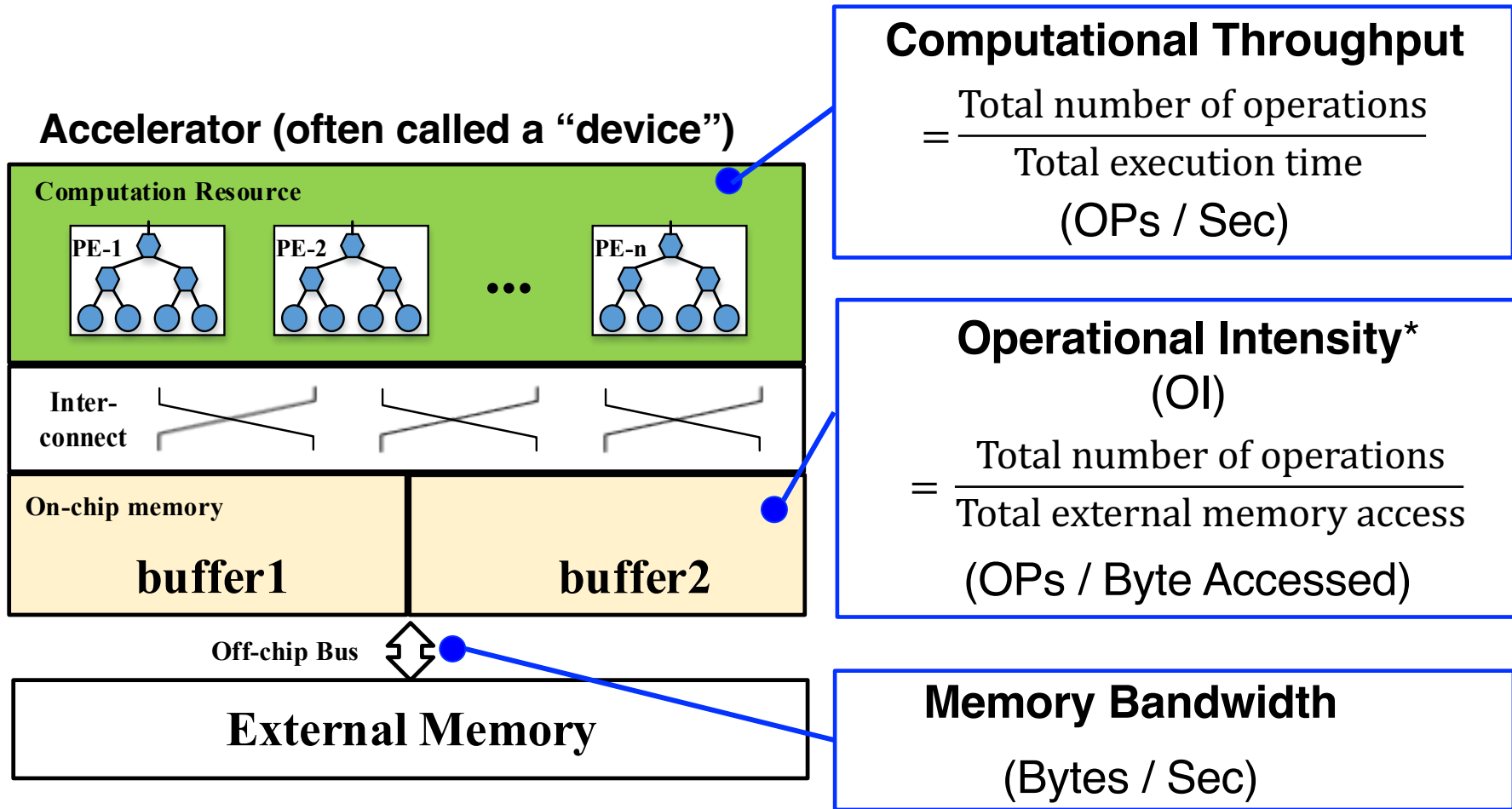
# Performance Modeling

- ▶ We calculate the following design metrics:
  - **Total number of operations (FLOP)**
    - Depends on the CNN model parameters
  - **Total external memory access (Byte)**
    - Depends on the CNN weight and activation size
  - **Total execution time (Sec)**
    - Depends on the hardware architecture (e.g., tile factors  $T_o$  and  $T_i$ )
    - Ignore resource constraints for now

# Performance Modeling

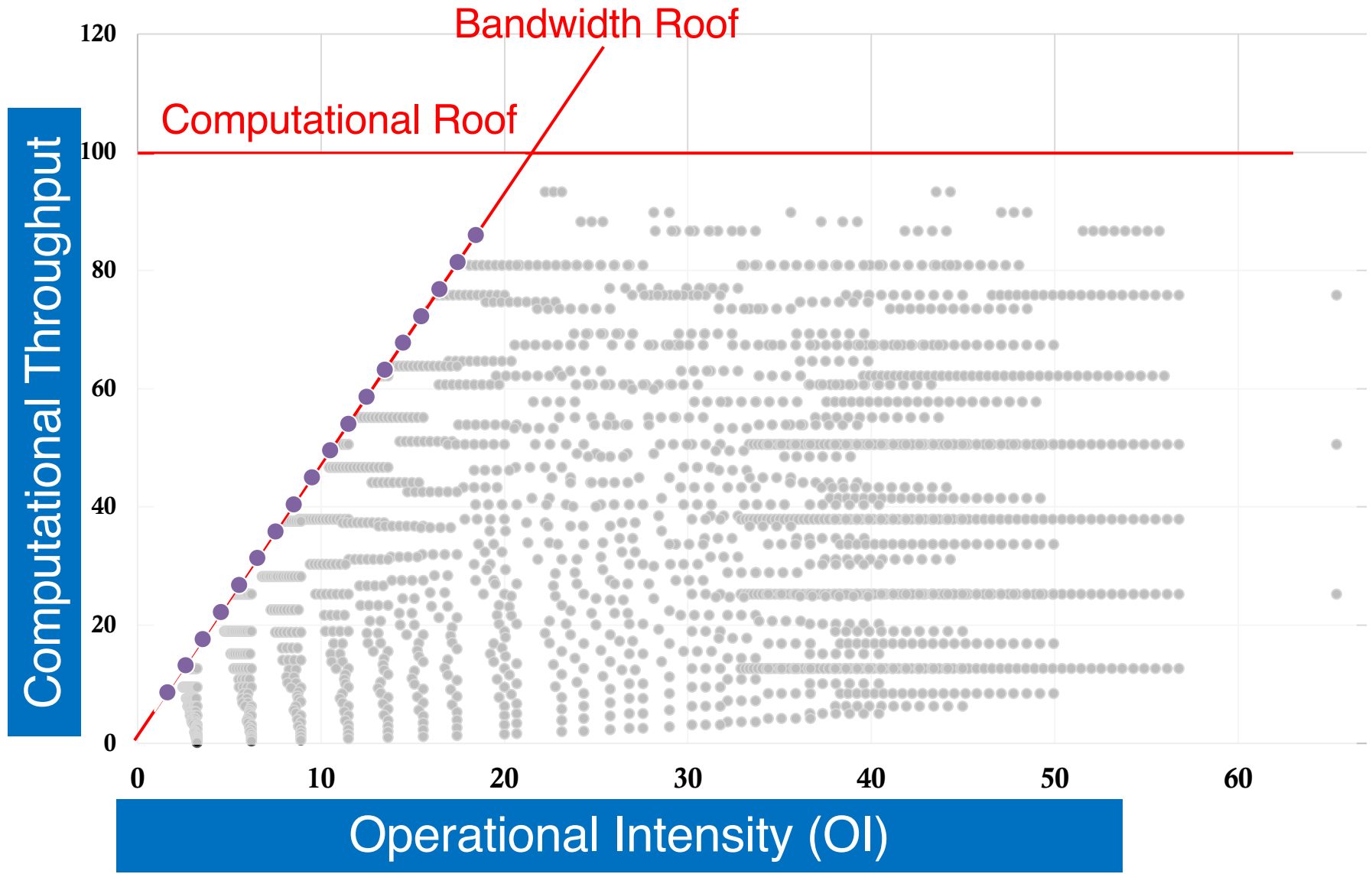
- ▶ Total operations FLOPS  $\approx 2 \times O \times I \times R \times C \times K^2$
- ▶ Execution time = Number of Cycles  $\times$  Clock Period
  - Number of cycles  $\approx \left\lceil \frac{O}{T_o} \right\rceil \times \left\lceil \frac{I}{T_i} \right\rceil \times \left\lceil \frac{R}{T_r} \right\rceil \times \left\lceil \frac{C}{T_c} \right\rceil \times (T_r \times T_c \times K^2)$   
 $\approx \left\lceil \frac{O}{T_o} \right\rceil \times \left\lceil \frac{I}{T_i} \right\rceil \times R \times C \times K^2$
- ▶ External memory accesses =  $a_i \times B_i + a_w \times B_w + a_o \times B_o$ 
  - Size of input fmap buffer:  $B_i = T_i \times (T_r + K - 1) \times (T_c + K - 1)$  with stride=1
  - Size of output fmap buffer:  $B_o = T_o \times T_r \times T_c$
  - Size of weight buffer:  $B_w = T_i \times T_o \times K^2$
  - External access times:  $a_o = \left\lceil \frac{O}{T_o} \right\rceil \times \left\lceil \frac{R}{T_r} \right\rceil \times \left\lceil \frac{C}{T_c} \right\rceil$ ,  $a_i = a_w = \left\lceil \frac{I}{T_i} \right\rceil \times a_o$ 
    - Input features and weights are reused across multiple output tiles

# Performance Modeling using Roofline

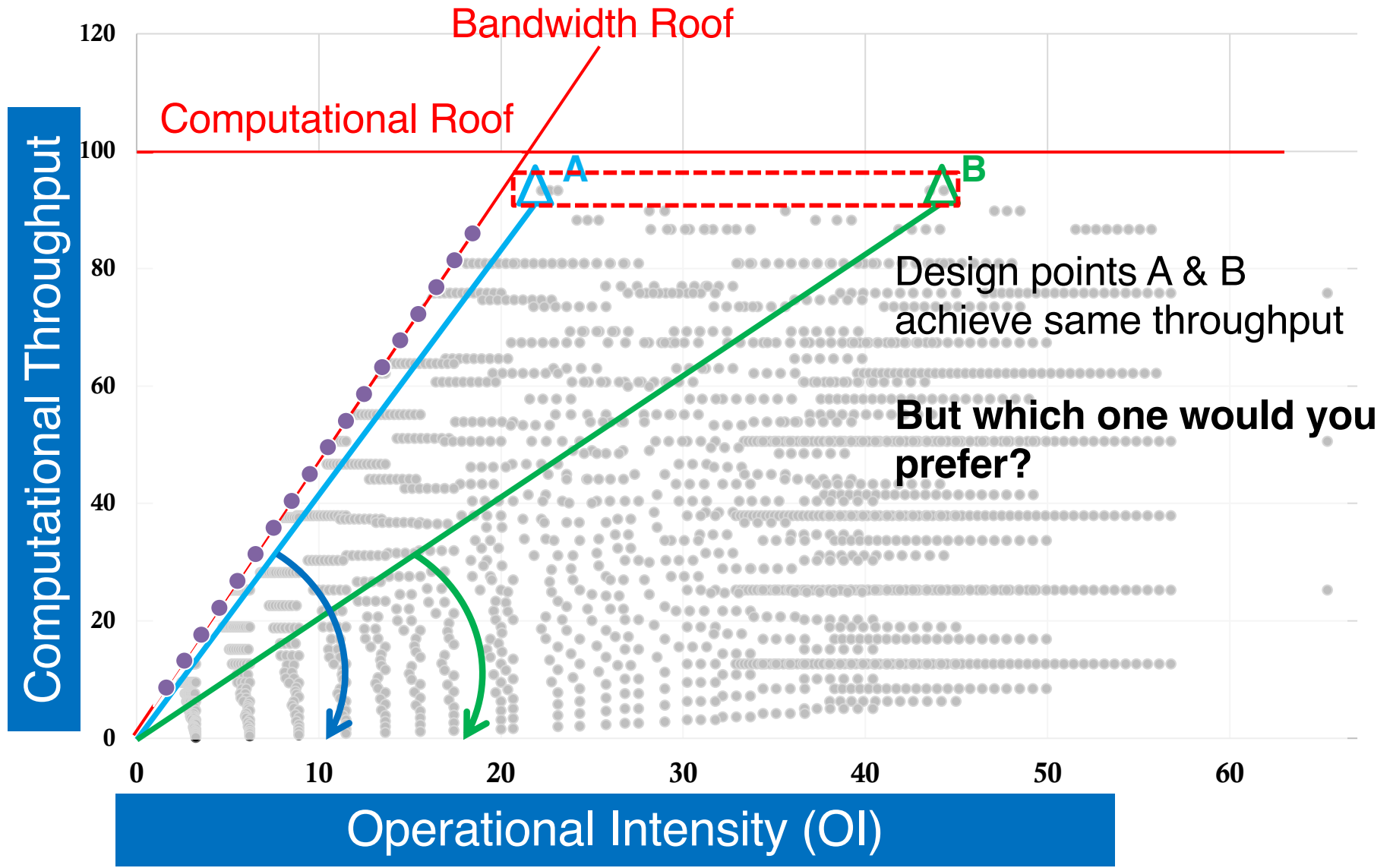


\* OI is also known as computation to communication ratio (CTC) or arithmetic intensity (AI)

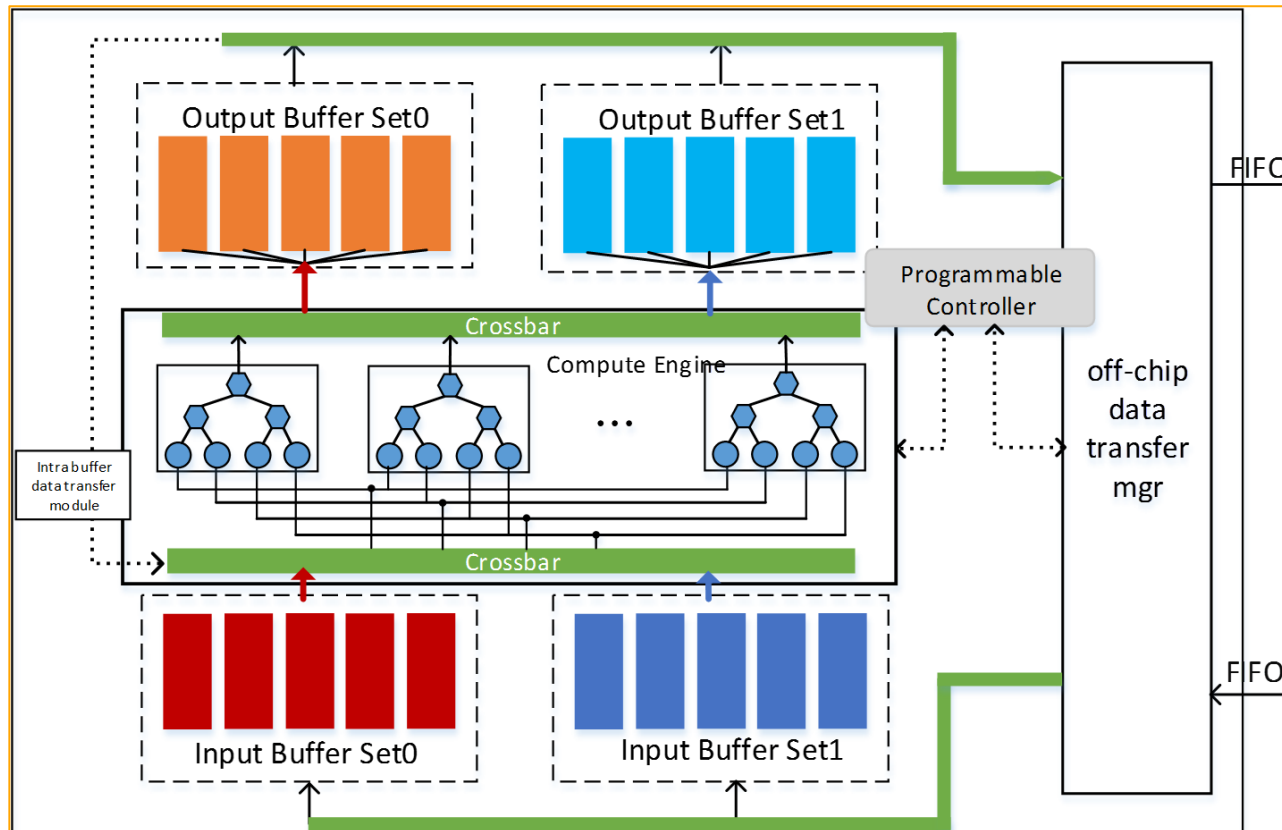
# Design Space Exploration with Roofline



# Design Space Exploration with Roofline



# Hardware Implementation



- ▶ All values in floating-point
- ▶ Only handles Conv layers, not dense or pooling
- ▶ Target FPGA: Virtex7-485t (28nm)

# **FracBNN: Accurate and FPGA-Efficient Binary Neural Networks with Fractional Activations**

Yichi Zhang<sup>1</sup>, Junhao Pan<sup>2</sup>, Xinheng Liu<sup>2</sup>, Hongzheng Chen<sup>1</sup>,  
Deming Chen<sup>2</sup>, Zhiru Zhang<sup>1</sup>

<sup>1</sup>Electrical and Computer Engineering, Cornell

<sup>2</sup>Electrical & Computer Engineering, UIUC

***FPGA'21, Feb/Mar 2021***

# CNNs with Reduced Numerical Precision

- ▶ Hardware architects widely quantization to improve the efficiency of CNN execution
  - **Motivation:** both neural nets and image/video apps naturally tolerate small amounts of noise
  - **Approach:** take a trained floating-point model and apply quantization
    - 16 or 8-bit fixed-point have been shown to be practical
- ▶ Can we go even lower by training a reduced-numerical-precision CNN from the ground up?

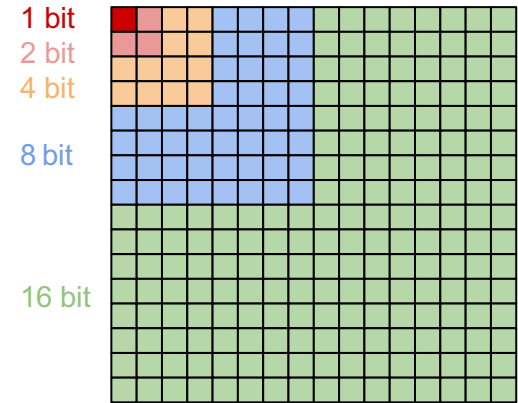


# Quantization Improves Efficiency

Reducing bitwidth by a factor of 2

- **1/4** of multiplication power cost
- **1/2** of addition power cost
- **1/2** of SRAM area/power cost

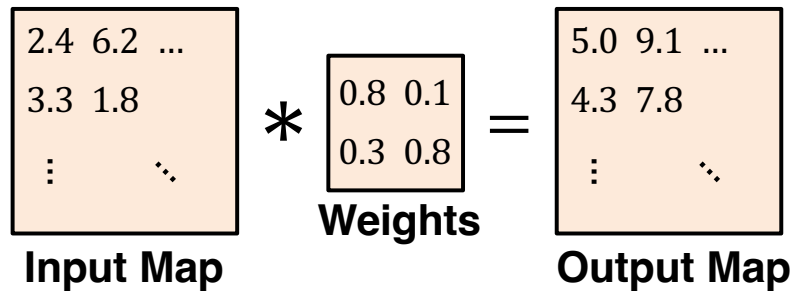
cost / multiplication



	INPUT OPERANDS	ACCUMULATOR	TOPS	X-factor vs. FFMA	SPARSE TOPS	SPARSE X-factor vs. FFMA
V100	FP32	FP32	15.7	1x	-	-
	FP16	FP32	125	8x	-	-
A100	FP32	FP32	19.5	1x	-	-
	TF32	FP32	156	8x	312	16x
	FP16	FP32	312	16x	624	32x
	BF16	FP32	312	16x	624	32x
	FP16	FP16	312	16x	624	32x
	INT8	INT32	624	32x	1248	64x
	INT4	INT32	1248	64x	2496	128x
	<b>BINARY </b>	INT32	<b>4992</b>	<b>256x</b>	-	-
IEEE FP64		19.5	1x	-	-	

# Binarization: Pushing Quantization to the Extreme

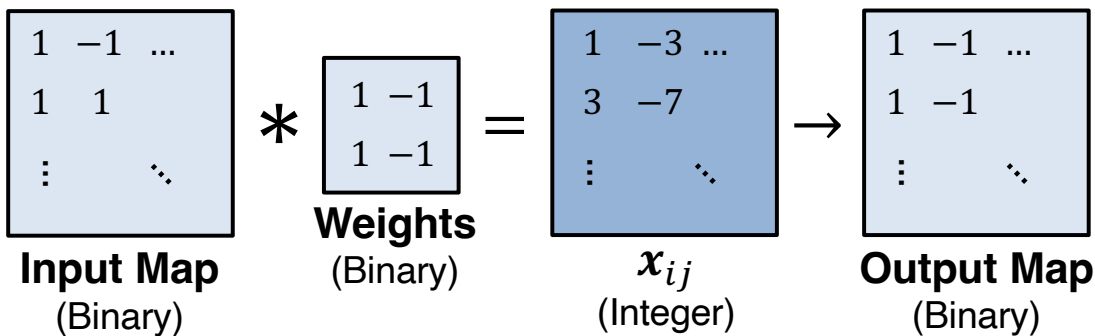
## CNN



### Key Differences

1. Inputs are binarized (-1 or +1)
2. Weights are binarized (-1 or +1)
3. MAC becomes XNOR+Popcount

## BNN



BNNs are well suited for FPGAs (rich in LUTs)

# Advantages of BNN

## 1. Floating point ops replaced with binary logic ops

$b_1$	$b_2$	$b_1 \times b_2$
-1	-1	+1
+1	-1	-1
+1	-1	-1
+1	+1	+1

$b_1$	$b_2$	$b_1 \text{ XNOR } b_2$
0	0	1
0	1	0
1	0	0
1	1	1

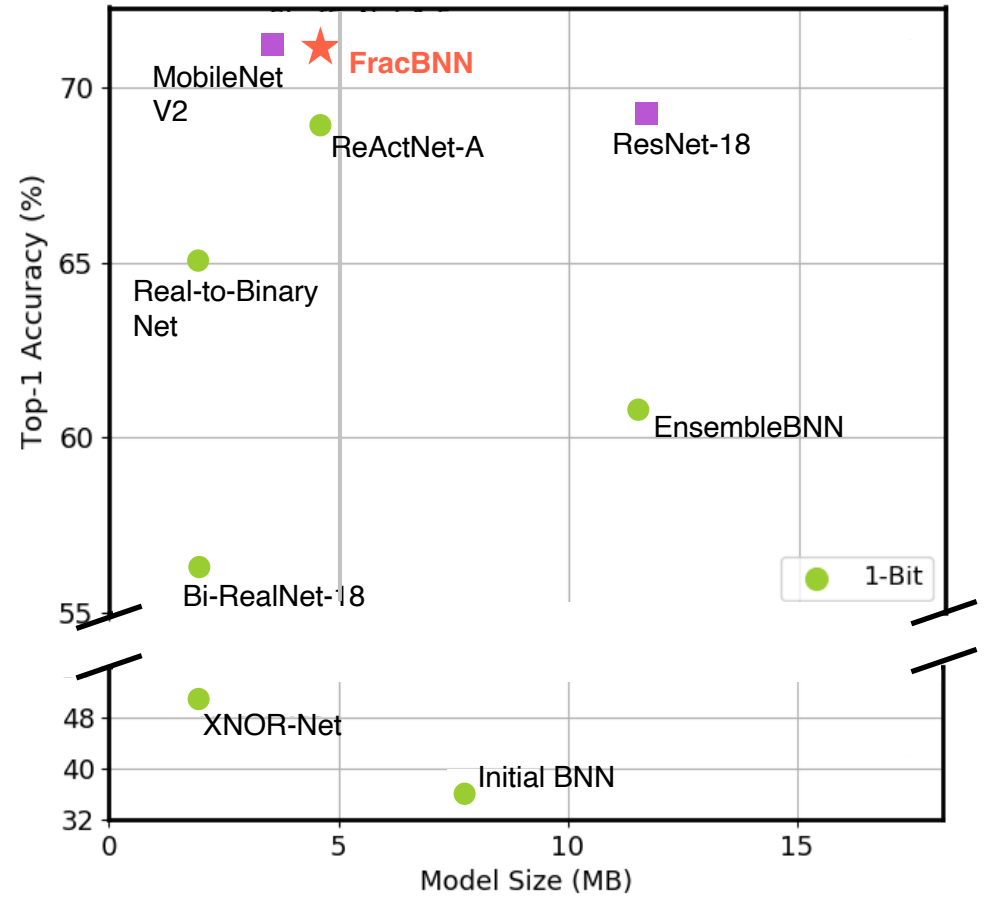
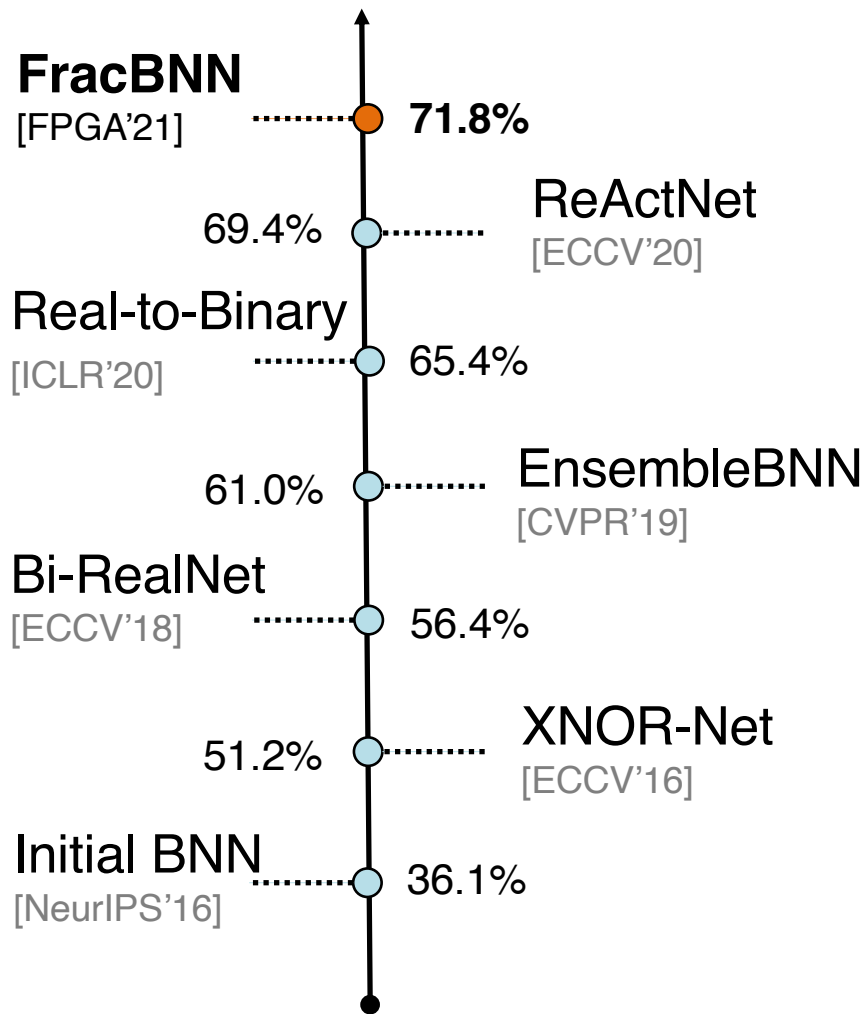
- Encode  $\{-1, +1\}$  as  $\{0,1\}$   $\rightarrow$  multiplies become XNORs
- Conv/dense layers do dot products  $\rightarrow$  XNOR and popcount
- Operations can map to LUT fabric as opposed to DSPs

## 2. Binarized weights may reduce total model size

- But note that fewer bits per weight may be offset by having more weights

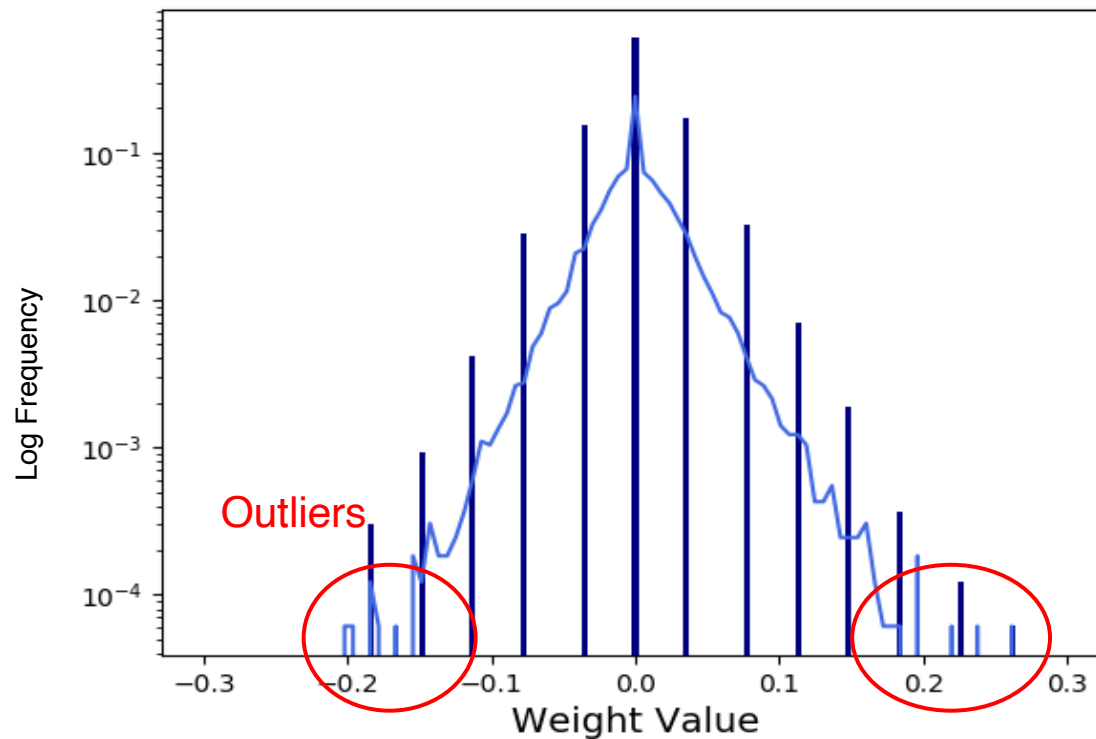
# Evolution of BNN Accuracy

ImageNet Top-1 Accuracy

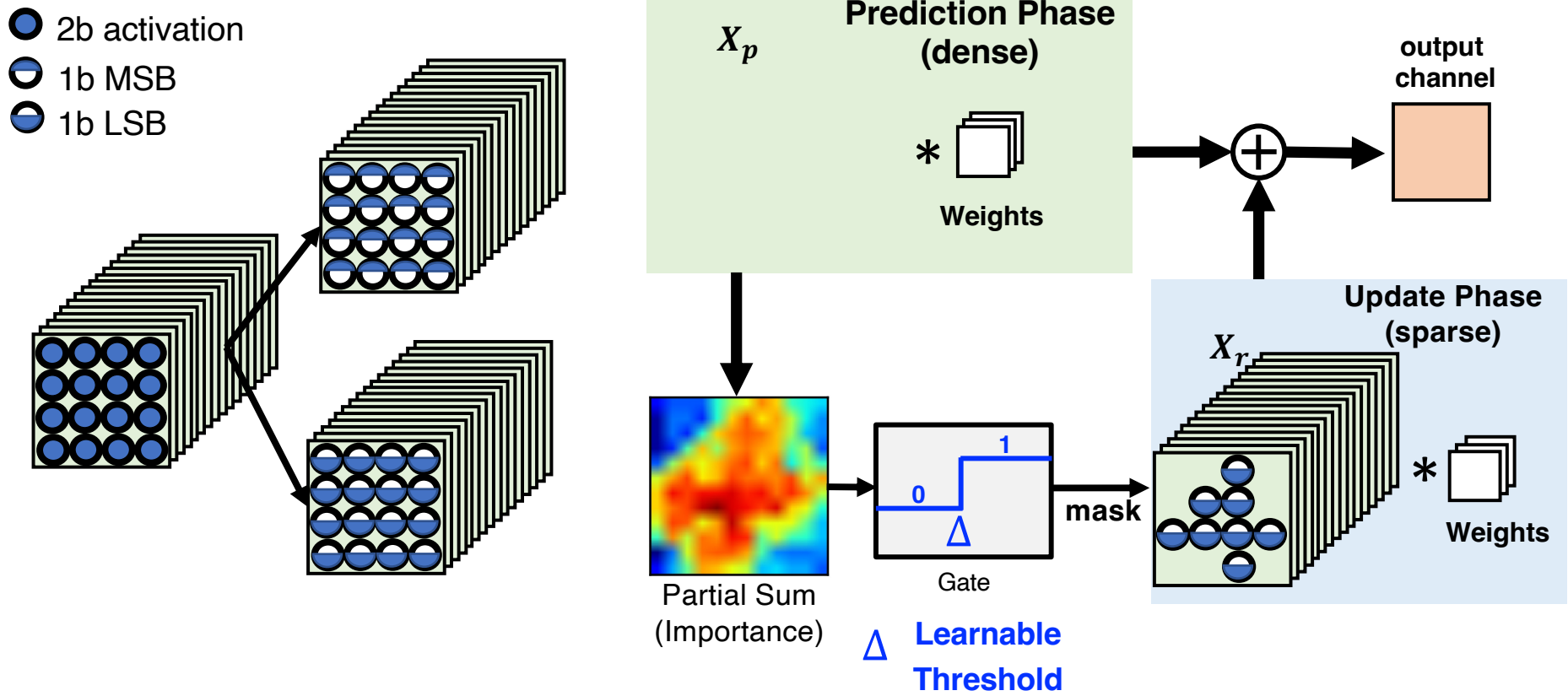


# The Outlier Problem in Quantization

- ▶ DNN weights and activations are typically distributed in a bell curve
  - Most values are close to zero, but rare *outliers* are large
  - Outliers stretch the quantization grid, resulting in poor resolution

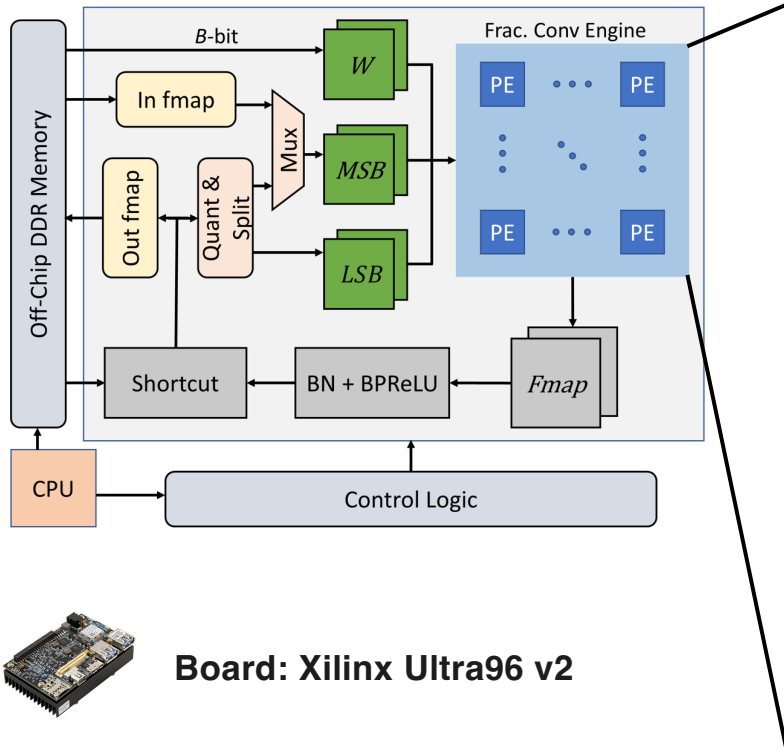


# Fractional Activation: Dual-Precision Scheme

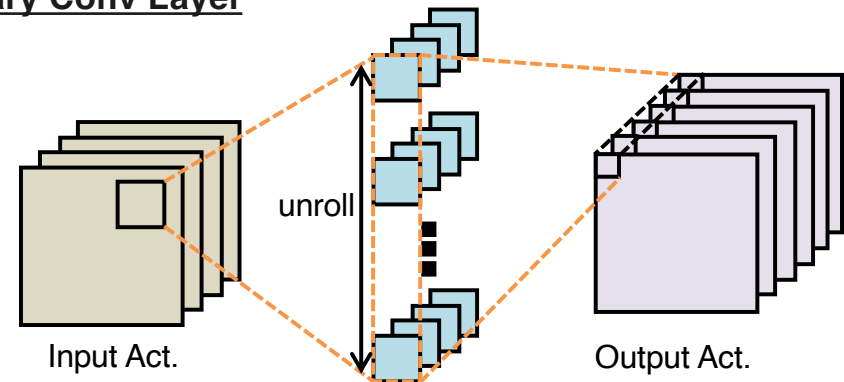


Exploits input-dependent characteristics  
Avoids overly aggressive quantization that degrades accuracy

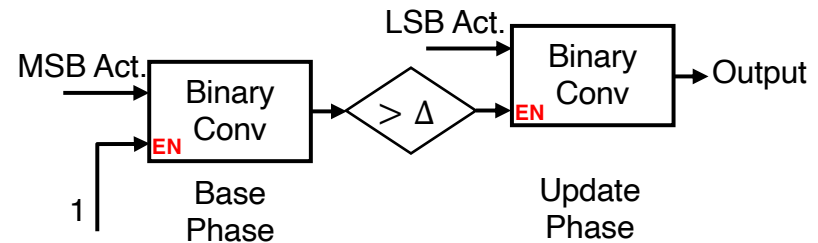
# FracBNN Accelerator Architecture on FPGA



## Binary Conv Layer



## Fractional Conv



All layers share the same conv engine using binary MACs (Implemented in HLS)

# FracBNN Hardware Evaluation

## Evaluation on CIFAR-10

	Device	Bits (W/A)	Freq. (MHz)	Top-1 (%)	Frame Rate
ReBNet	Zynq ZC702	1/1	200	86.98	2000
Initial BNN	Zynq 7Z020	1/1	143	88.8	168.4
<b>FracBNN</b>	Zynq ZU3EG	1/1.4	<b>250</b>	<b>89.1</b>	<b>2806.9</b>

## Evaluation on ImageNet

	Device	Bits (W/A)	Freq. (MHz)	Top-1/Top-5 (%)	Frame Rate
FINN-R	Zynq ZU3EG	1/2	220	50.3/-	200
Synetgy	Zynq ZU3EG	4/4	<b>250</b>	68.3/88.1	41.1
<b>FracBNN</b>	Zynq ZU3EG	1/1.4	<b>250</b>	<b>71.8/90.1</b>	48.1



## Additional Useful Resources

- ▶ Recent papers on neural networks on silicon
  - <https://github.com/fengbintu/Neural-Networks-on-Silicon>
- ▶ Tutorial on hardware architectures for DNNs
  - <http://eyeriss.mit.edu/tutorial.html>
- ▶ Landscape of neural network inference accelerators
  - <https://nicsefc.ee.tsinghua.edu.cn/projects/neural-network-accelerator>

# Acknowledgements

- ▶ This tutorial contains/adapts materials developed by
  - Ritchie Zhao (Cornell ECE PhD, now Microsoft)
  - Yichi Zhang (PhD Student, Cornell ECE)
  - Authors of the following papers
    - Optimizing FPGA-based Accelerator Design for Deep Convolutional Neural Networks (FPGA'15, PKU-UCLA)
    - FracBNN: Accurate and FPGA-Efficient Binary Neural Networks with Fractional Activations (FPGA'21, Cornell-UIUC)