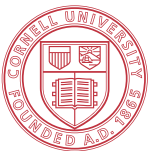


ECE 6775  
High-Level Digital Design Automation  
Fall 2025

# Deep Learning Acceleration on FPGAs



Cornell University



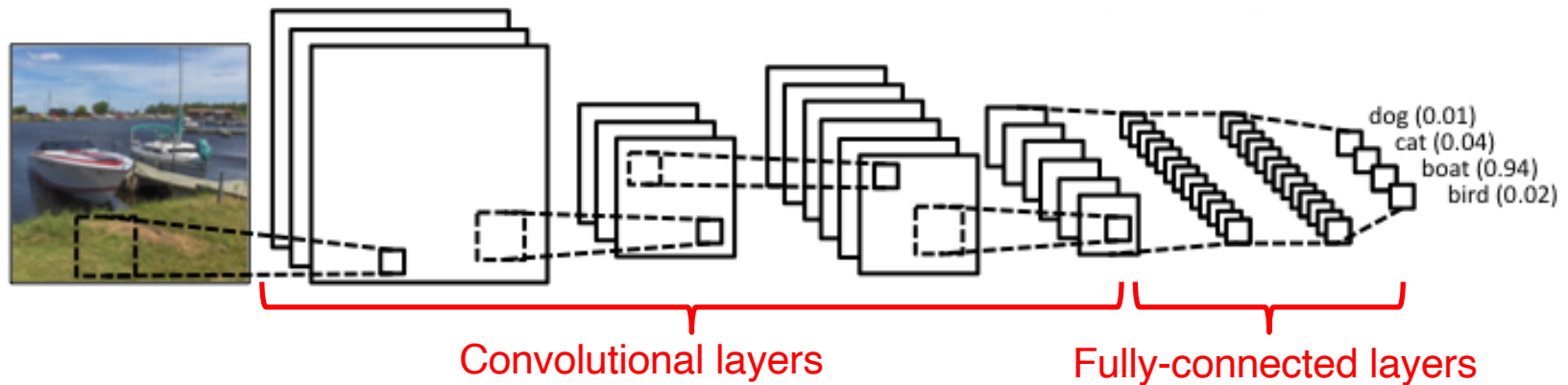
# Announcements

- ▶ **HW 2 hard deadline:** tonight at 11:59pm
- ▶ **Lab 4 teams finalized on CMS: get started early!**
- ▶ **In-class prelim on Tuesday 10/28**
  - **Coverage: Lectures 01–11, 13, 14**
    - Lec12 NN tutorial excluded
    - Lec15 includes useful discussion on URE and pipelining
    - Lec16 (today) reviews the roofline model
  - We'll do a more in-depth prelim review in the next lecture

# Agenda

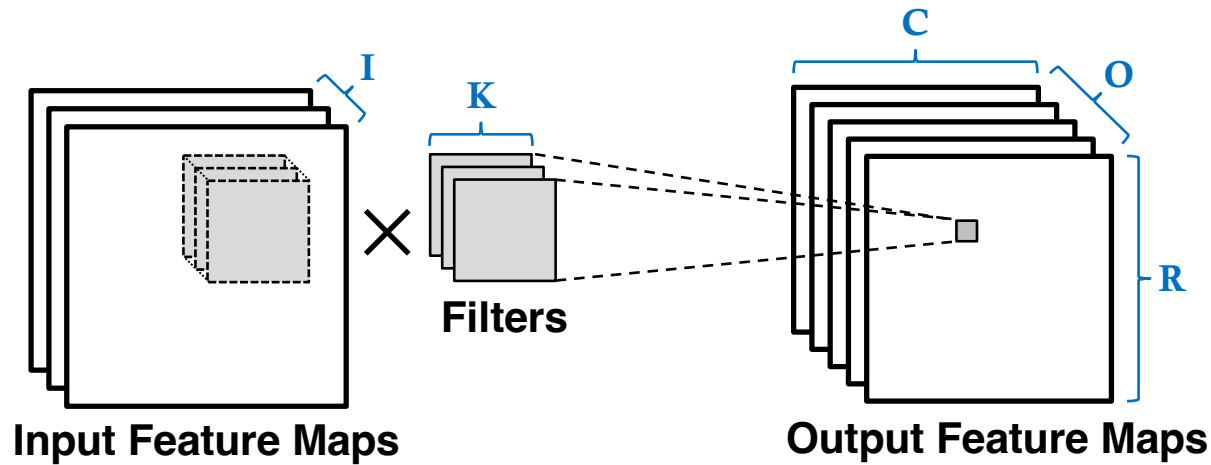
- ▶ CNN acceleration on FPGAs
- ▶ Potential of FPGA-based LLM inference

# Recap: Convolutional Neural Network (CNN)



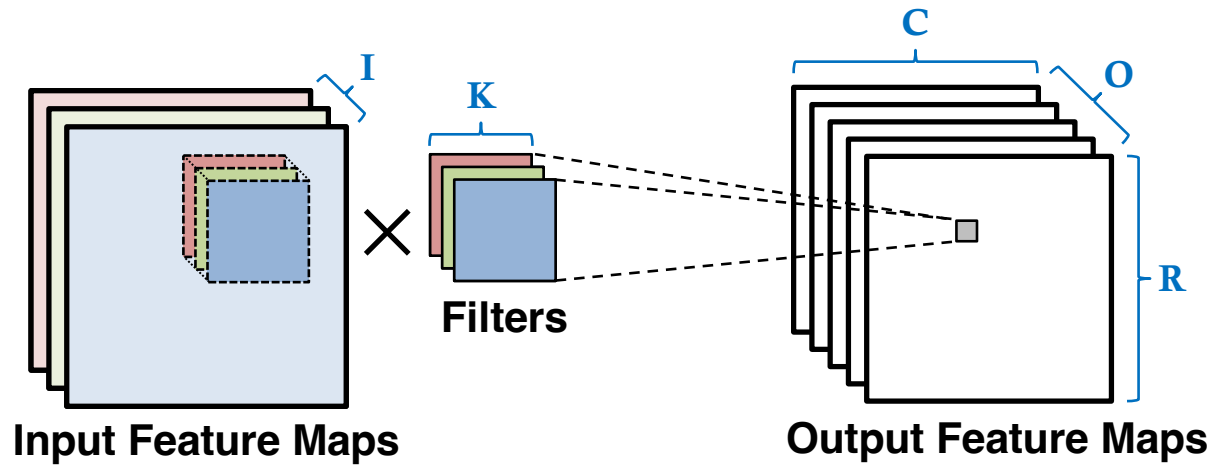
- ▶ Front: convolutional layers learn visual features
- ▶ Back: fully-connected layers perform classification using the visual features

# Convolutional Layer



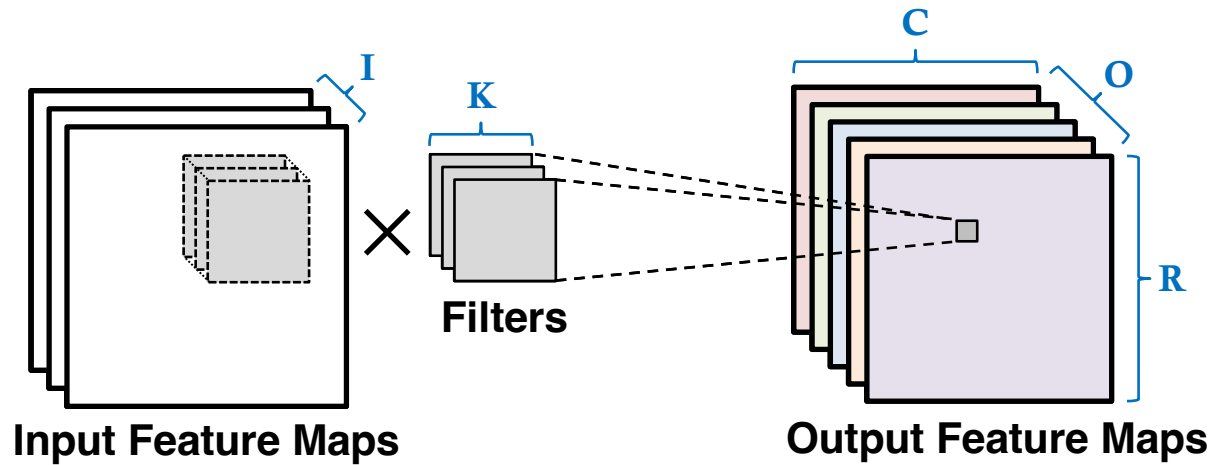
- ▶ An output pixel is connected to its neighboring region on each input feature map (fmap)
- ▶ All pixels on an output feature map use the same filter weights

# Parallelism in the Convolutional Layer



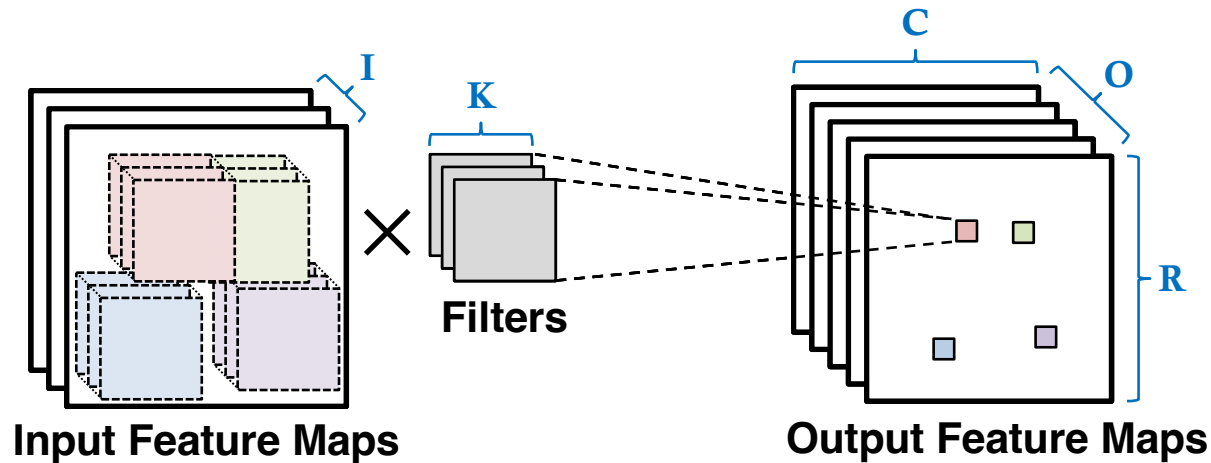
- ▶ Four main sources of parallelism
  1. **Across input feature maps (i.e., input channels)**

# Parallelism in the Convolutional Layer



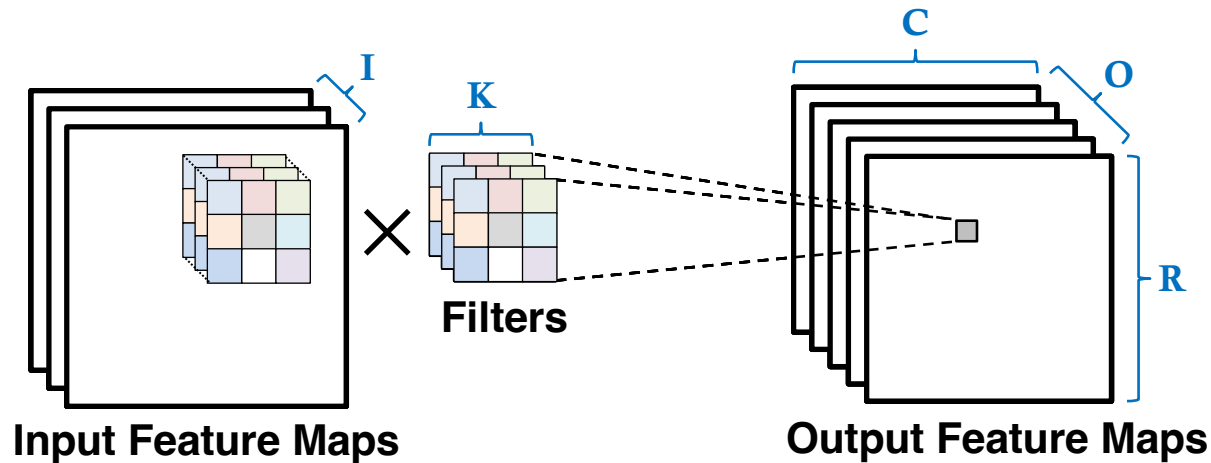
- ▶ Four main sources of parallelism
  1. Across input feature maps (i.e., input channels)
  2. **Across output feature maps (i.e., output channels)**

# Parallelism in the Convolutional Layer



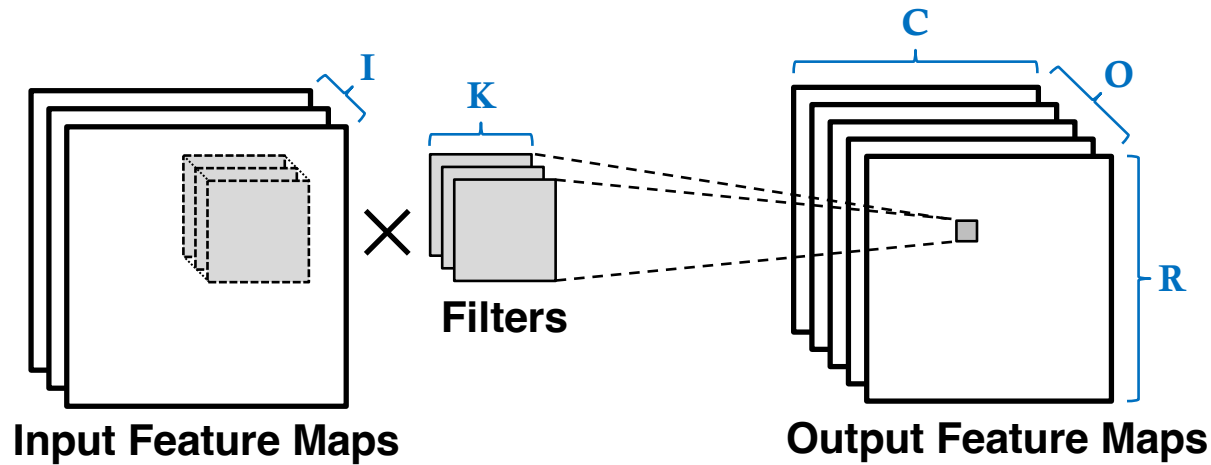
- ▶ Four main sources of parallelism
  1. Across input feature maps (i.e., input channels)
  2. Across output feature maps (i.e., output channels)
  3. **Across different output pixels (i.e., filter positions)**

# 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> Peking University

<sup>2</sup> UCLA

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

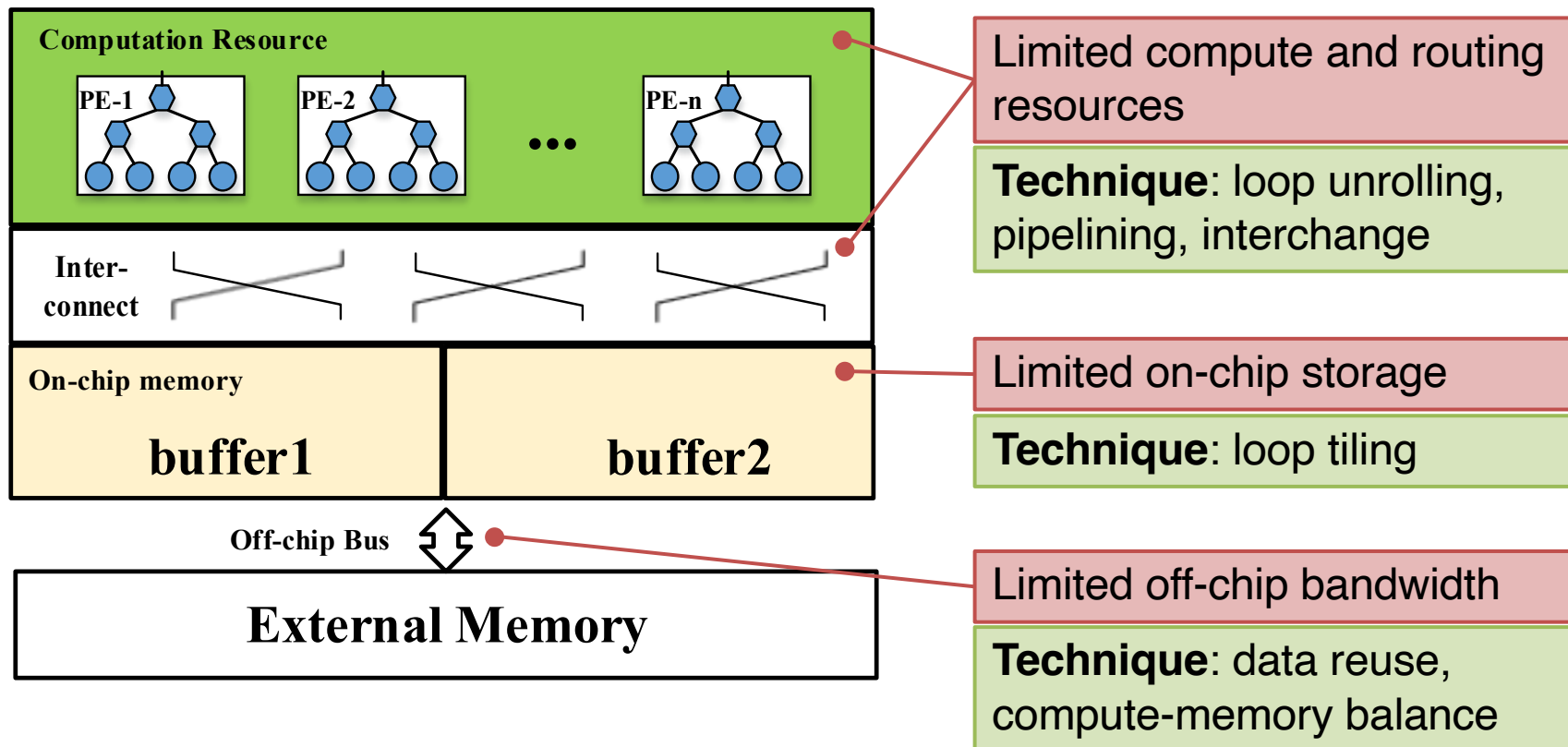
***FPGA 2015***

# Main Contributions

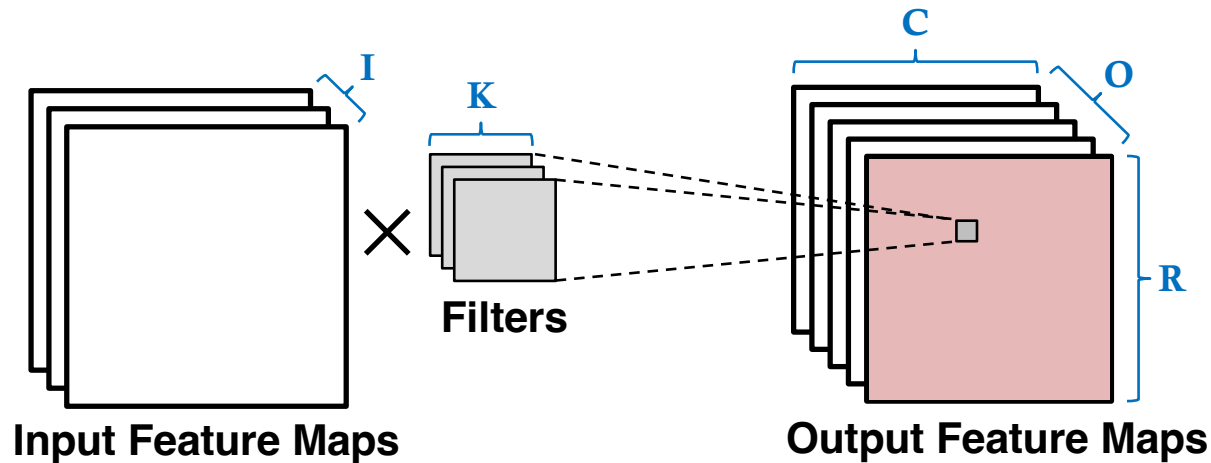
1. Analysis of the different sources of parallelism in the convolution kernel of a CNN
2. Quantitative performance modeling of the hardware design space using the Roofline method
3. Design and implementation of a CNN accelerator for FPGA using Vivado HLS, evaluated on AlexNet

# 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



# Motivation for Loop Tiling



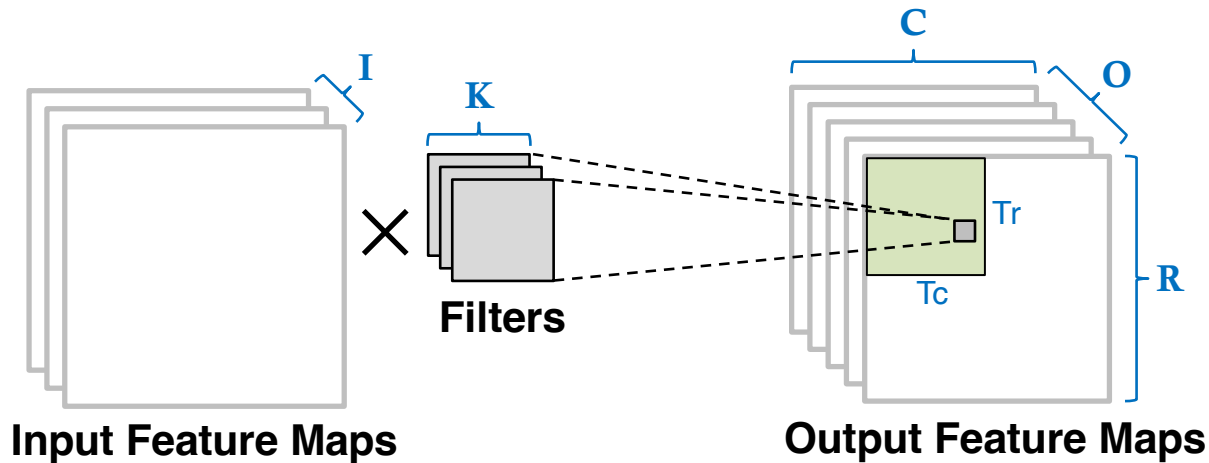
```
1 for (row=0; row<R; row++) {  
2   for (col=0; col<C; col++) {  
3     for (to=0; to<O; to++) {  
4       for (ti=0; ti<I; ti++) {  
5         for (ki=0; ki<K; ki++) {  
6           for (kj=0; kj<K; kj++) {  
               output_fm[to][row][col] +=  
               weights[to][ti][ki][kj]*input_fm[ti][S*row+ki][S*col+kj];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

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

If we offload these loops, we must store entire output feature maps on-chip

# 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<min(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      }
23      // software: read output feature map
24    }
25  }
26 }
```

**CPU Portion**

**FPGA Portion  
(Accelerator)**

# Optimizing for Accelerator 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<min(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];
                        }
                    }
                }
            }
        }
    }
```

**FPGA Portion**  
**(Accelerator)**

# Optimizing for Accelerator 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<min(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 Accelerator 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++) {
8                   for (tii=ti; tii<min(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 Accelerator 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<min(ti+Ti, I); tii++) {
12                             output_fm[too][trr][tcc] +=
13                                weights[too][tii][ki][kj]*input_fm[tii][S*trr+ki][S*tcc+kj];
14                          }
15                    }
16              }
17           }
18        }
19     }
```

# Optimizing for Accelerator 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<min(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

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

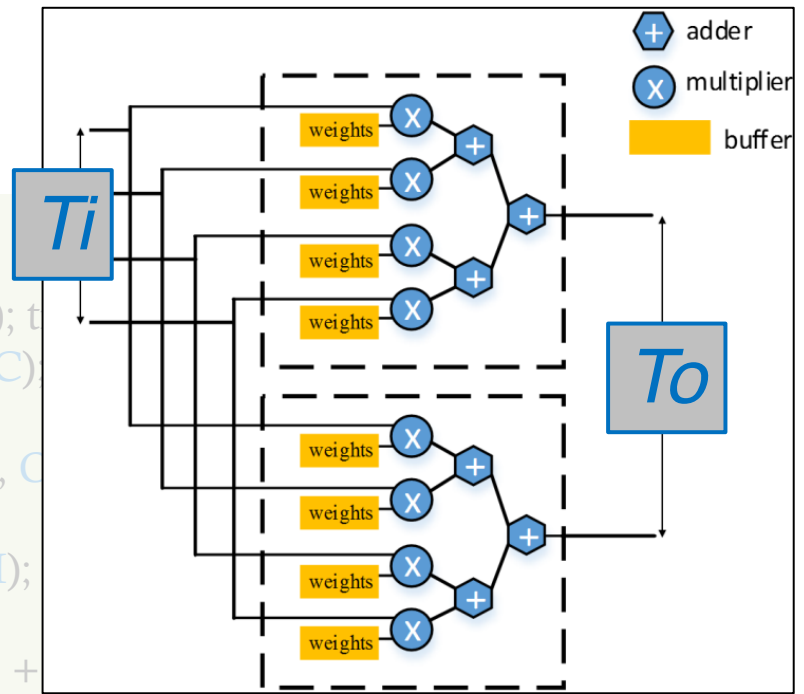
$L$  is the pipeline depth (i.e., # of pipeline stages), assume  $ll=1$

# Optimizing for Accelerator 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<min(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  }}}}}

```

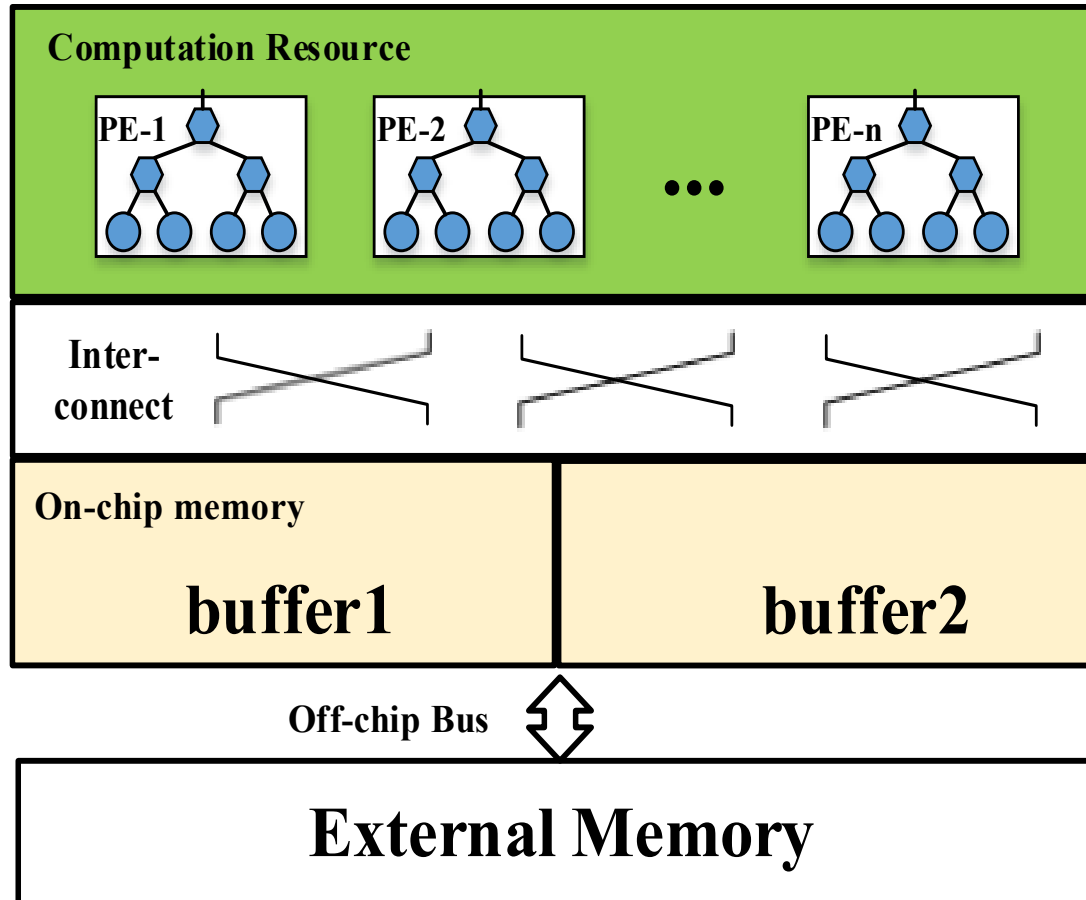


## Generated Hardware

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

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

# Overall Accelerator Architecture



# 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

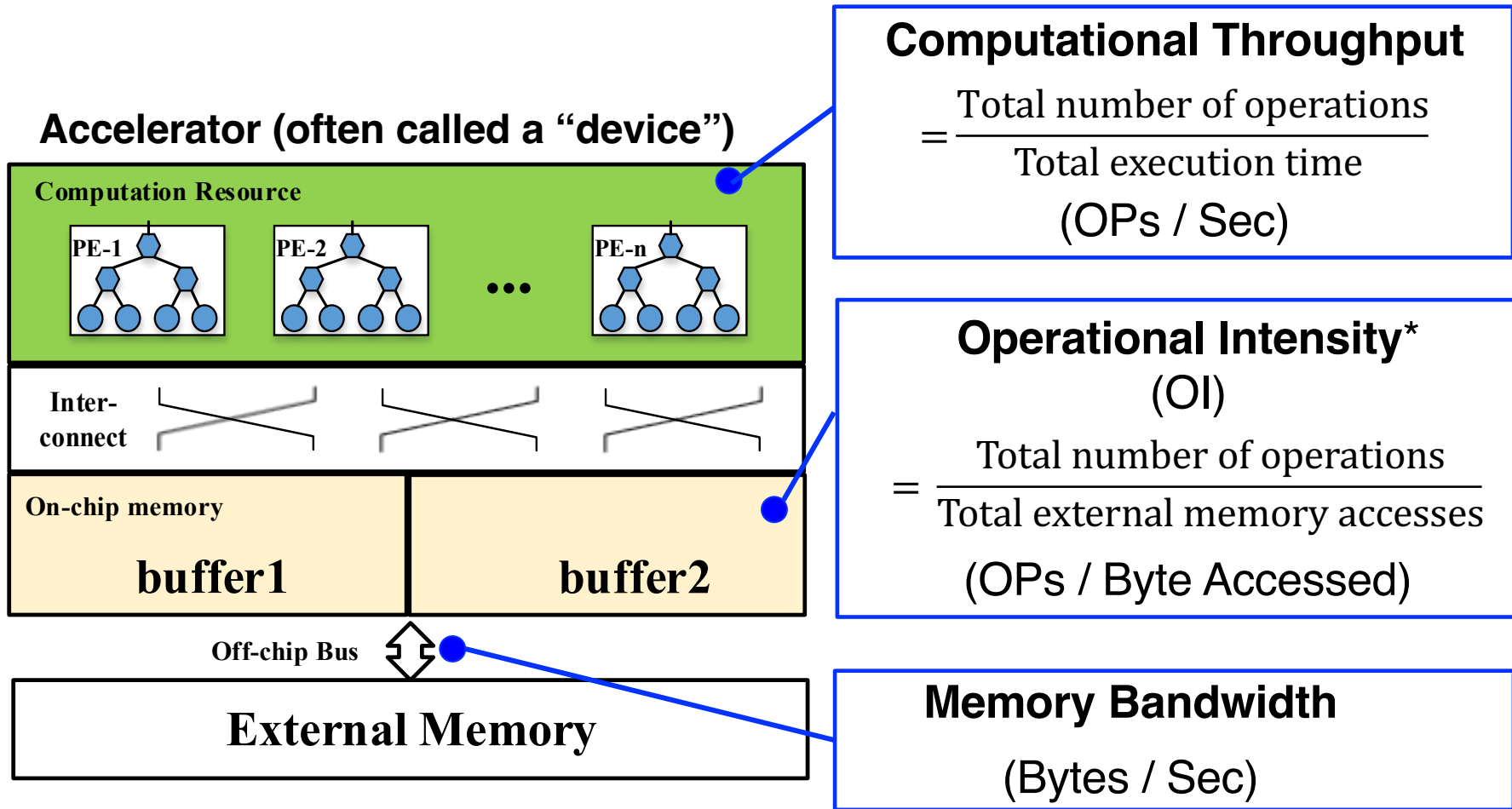
# Performance Modeling

- ▶ Three design metrics are estimated:
  - **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  $T_i$ )

# 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 =  $a_o \times B_o + a_i \times B_i + a_w \times B_w$ 
  - Size of output fmap buffer:  $B_o = T_o \times T_r \times T_c$
  - Size of input fmap buffer:  $B_i = T_i \times (T_r + K - 1) \times (T_c + K - 1)$  with stride=1
  - Size of weight buffer:  $B_w = T_i \times T_o \times K^2$
  - External access times:  $a_o = \left\lceil \frac{O}{T_o} \right\rceil \times \left\lceil \frac{R}{T_r} \right\rceil \times \left\lceil \frac{C}{T_c} \right\rceil$ ,  $a_i = a_w = \left\lceil \frac{I}{T_i} \right\rceil \times a_o$ 
    - $a_o$ ,  $a_i$ , and  $a_w$  denote the number of transactions between external memory and the output, input, and weight buffers
    - Input features and weights are reused across multiple output tiles

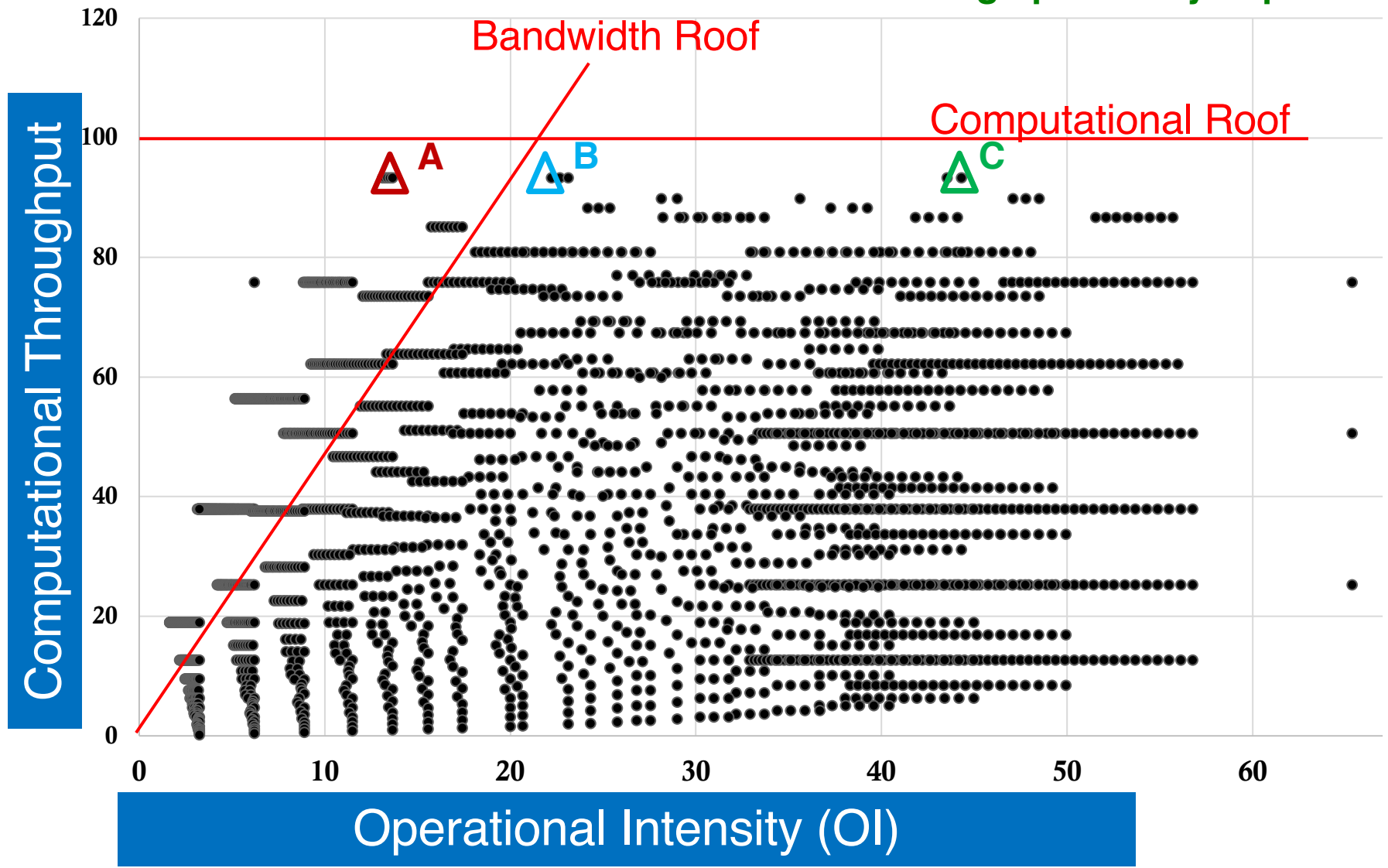
# Revisiting the Roofline Model



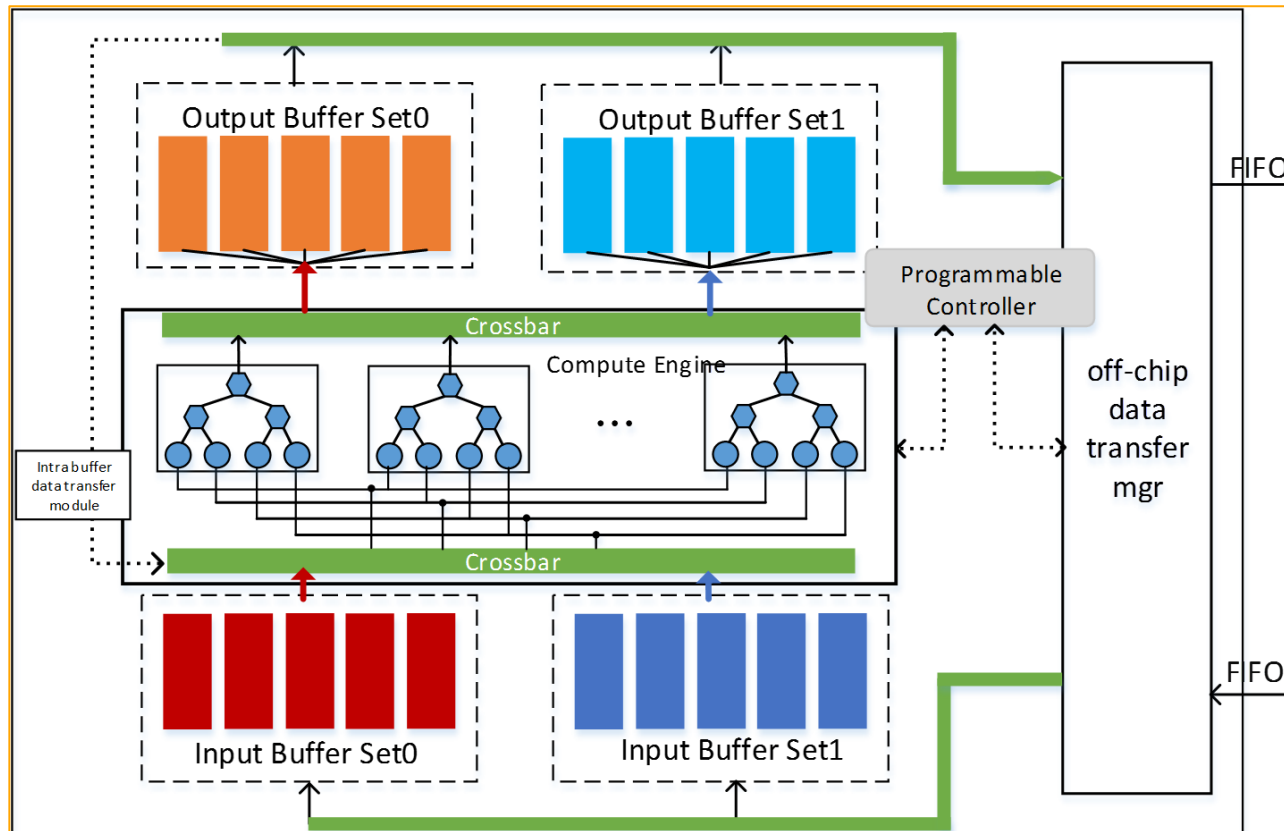
\* **OI** is also known as computation to communication ratio (CTC) or arithmetic intensity (AI)

# Design Space Exploration with Roofline

Which design point do you prefer?



# FPGA Implementation



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

# Understanding the Potential of FPGA-Based Spatial Acceleration for Large Language Model Inference

Hongzheng Chen<sup>1</sup>, Jiahao Zhang<sup>2,1</sup>, Yixiao Du<sup>1</sup>, Shaojie Xiang<sup>1</sup>, Zichao Yue<sup>1</sup>, Niansong Zhang<sup>1</sup>, Yaohui Cai<sup>1</sup>, Zhiru Zhang<sup>1</sup>







<sup>1</sup> Cornell University

<sup>2</sup> Tsinghua University

***ACM TRETS 2024 (FCCM Journal Track)***

# The Era of Large Language Models (LLMs)



Devices	Cost	Energy
 A100 GPU (7nm)	 ~\$20K	 ~250W
 U280 FPGA (16nm)	 ~\$8K	 ~30W

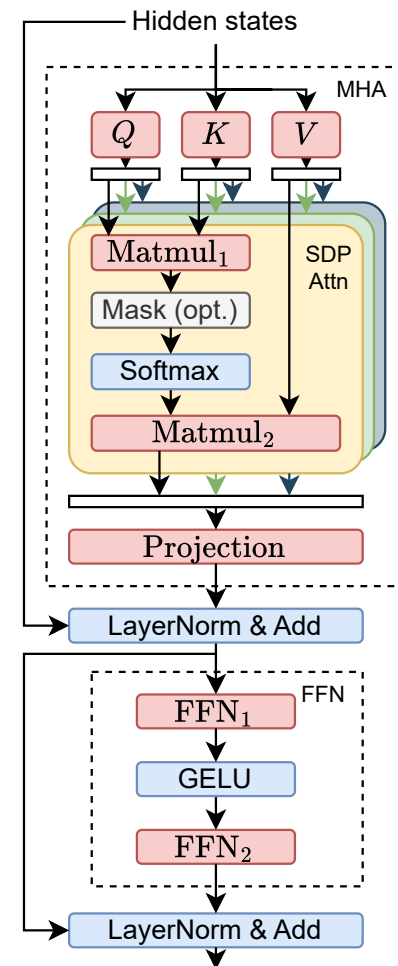
Is FPGA suitable for efficient LLM inference?

# Transformer GPU Execution Breakdown

- ▶ LLMs aren't just about matrix multiplications (MatMul)
- ▶ Non-linear & element-wise operators also play a significant role
  - Low compute-to-memory ratio (namely, low OI)
  - High kernel launch overheads

Operator Class	Representative Op	% FLOP	% Run Time
Tensor contraction	MM, MV	<b>99.80</b>	61.0
Stat. normalization	softmax, layernorm	0.17	<b>25.5</b>
Element-wise	bias, dropout	0.03	<b>13.5</b>

Proportions for operator classes in the BERT model, implemented in PyTorch, profiled with an NVIDIA V100 GPU  
 (MM: matrix-matrix multiply; MV: matrix-vector multiply)



# Two-Stage LLM Inference: OI Analysis

## Stage 1: Prefill

Takes in user prompts and generates the first token

### MM

$$\mathbf{A}_{m \times k} \times \mathbf{B}_{k \times n} = \mathbf{C}_{m \times n}$$

$$\frac{\text{compute}}{\text{memory footprint}} = \frac{mkn}{mk + kn + mn}$$

MM OI: Cubic / Quadratic

**Compute bound:** performance dictated by the compute roof

## Stage 2: Decode

Processes the previously generated token to produce new tokens one at a time in an auto-regressive manner

### MV

$$\mathbf{A}_{m \times k} \times \mathbf{b}_k = \mathbf{c}_m$$

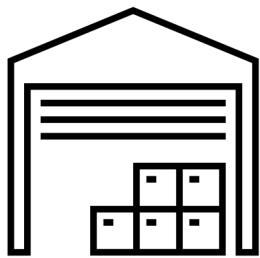
$$\frac{\text{compute}}{\text{memory footprint}} = \frac{mk}{mk + k + m}$$

MV OI: Quadratic / Quadratic

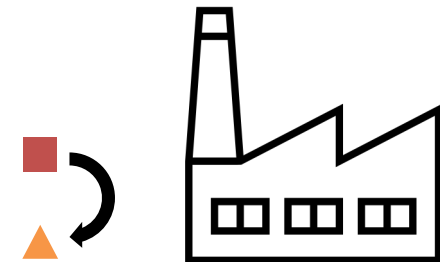
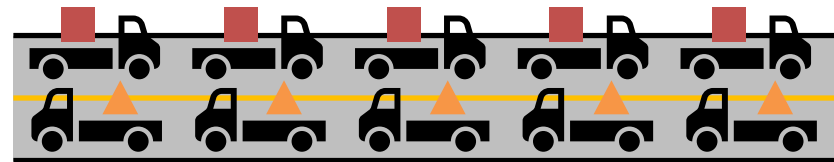
**Memory bound:** performance dictated by the bandwidth roof

# Typical GPU Execution Cycle for a Single Operator

- ▶ **Analogy:** Compute is like a factory, requiring instructions (kernel launch overhead) and a steady data supply (memory bandwidth) to run efficiently
- ▶ If compute efficiency grows faster than data supply, it limits the system's ability to operate at peak performance



Warehouse  
(memory)



Factory  
(compute units)

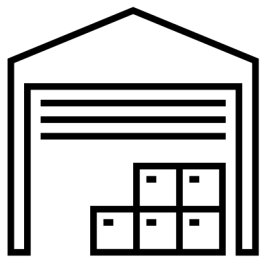
# Typical GPU Execution Cycle for Multiple Operators

- ▶ Moving data to and from GPU compute units incurs a high memory bandwidth cost
- ▶ For memory-bound operations, more time is spent moving data than performing computation

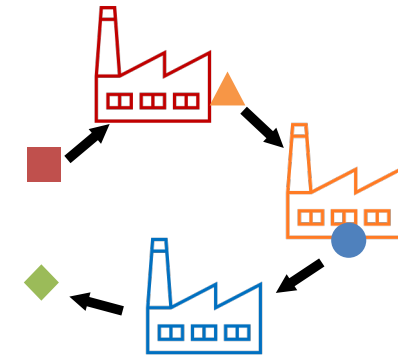


# The Potential of Dataflow Execution

- ▶ Model-specific dataflow accelerator
  - Pass intermediate results directly to the next operator
  - Reduced memory traffic leads to higher performance and higher energy efficiency



Warehouse  
(memory)



Factory  
(specialized  
compute units)

# Contributions of This Work

- ▶ An analytical model to estimate performance and resource/bandwidth usage for model-specific dataflow accelerator on FPGAs

- **Compute resource:**  $M$  is compute power in MACs/cycle and  $C$  is the # of layers per FPGA

$$\sum M_i C < M_{\text{tot}}, i \in \{q, k, v, a_1, a_2, p, f_1, f_2\}$$

- **Memory capacity:**  $S$  is buffer size

$$S_{\text{param}} C < \text{DRAM}_{\text{tot}},$$

$$\sum S_i C < \text{SRAM}_{\text{tot}}, i \in \{\text{tile, KV, FIFO}\}$$

- **Memory port:**  $s$  is tensor size and  $b$  is bitwidth

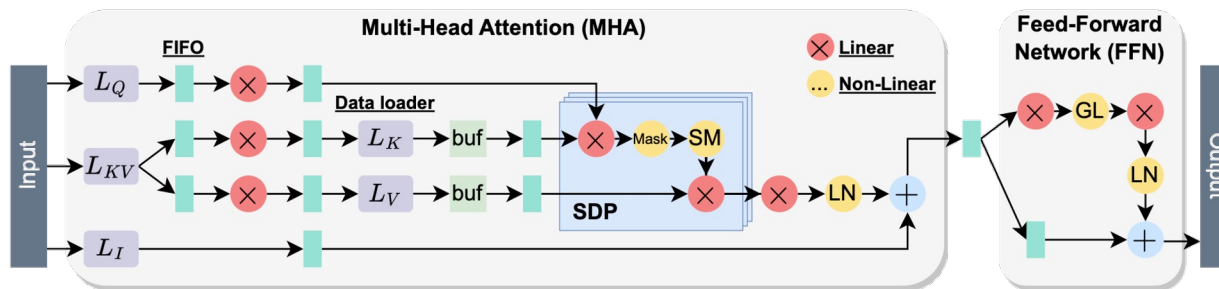
$$R_i = \left\lceil \frac{s_i b_{\text{BRAM}}}{M_i / r_i \times S_{\text{BRAM}}} \right\rceil \times \frac{M_i / r_i}{k}$$

$$\sum_i C R_i + 2C(R_{a_1} + R_{a_2}) < \text{SRAM}_{\text{tot}}, i \in \{q, k, v, p, f_1, f_2\}$$

- **Memory bandwidth:**  $B$  is bandwidth

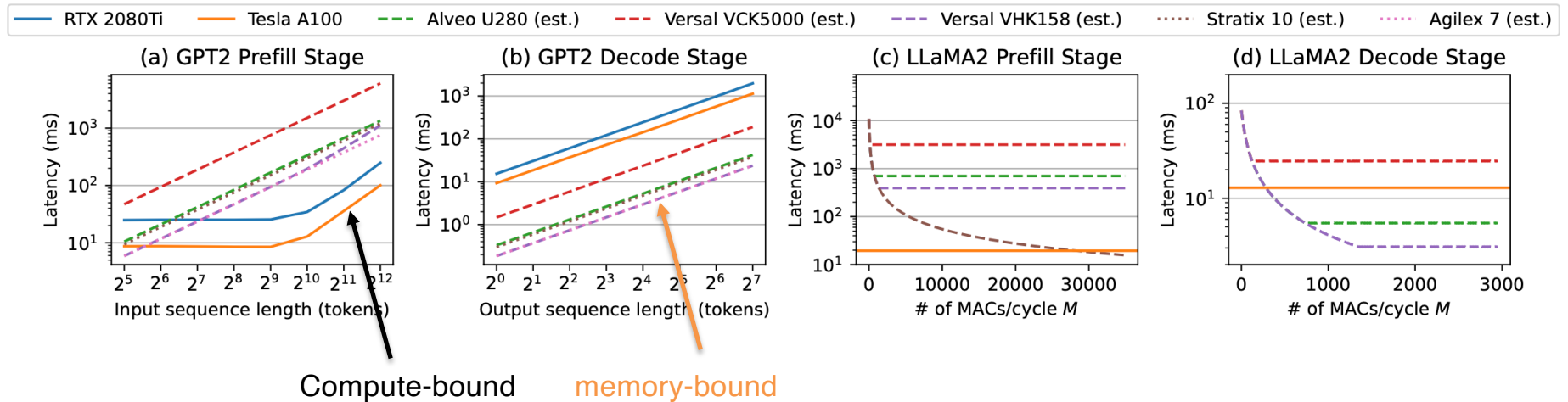
$$\sum_i C B_i < B_{\text{tot}}, i \in \{q, k, v, p, f_1, f_2\}$$

- ▶ An optimized HLS kernel library and compose accelerator designs for LLM on FPGAs



# Estimation on Different Models and Devices

- ▶ Is FPGA capable of efficient LLM inference?
  - Latency estimation of GPT2 and LLaMA2 on different FPGAs
  - GPU results are obtained through actual profiling



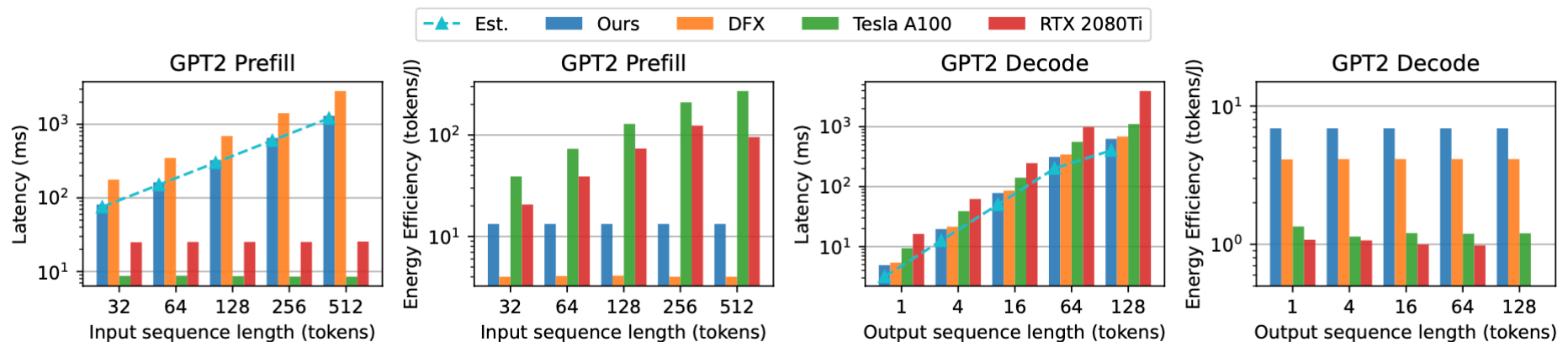
Existing FPGAs are inferior in the compute-intensive prefill stage, but can outperform GPUs in the memory-intensive decode stage

# FPGA Implementation of an GPT2 Accelerator

- ▶ GPT2: single-batch, low-latency settings, adjust input/output sequence length
  - Implemented on AMD Alveo U280, W8A8 quantization, 250MHz
  - 92% estimation accuracy in the prefill stage using the analytical framework
  - 2.2x speedup in prefill over DFX (an FPGA-based temporal arch.)
  - **1.7x speedup in decode stage, 5.4x more energy-efficient vs A100**

**Alveo U280**<sub>(ours)</sub> *16nm*  
Max M: 1289  
Bandwidth: **498 GB/s**

**A100 GPU** *7nm*  
Max M: 16415 (temporal)  
Bandwidth: **1935 GB/s**



# Acknowledgements

- ▶ This tutorial contains/adapts materials developed by
  - Ritchie Zhao (Cornell ECE PhD, now NVIDIA)
  - Authors of the following papers
    - Optimizing FPGA-based Accelerator Design for Deep Convolutional Neural Networks (FPGA'15, PKU-UCLA)
    - Understanding the Potential of FPGA-Based Spatial Acceleration for Large Language Model Inference (ACM TRETS, FCCM'24 Journal Track)

## Next Lecture

- ▶ Prelim Review