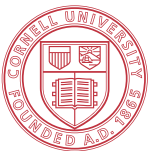




ECE 6775
High-Level Digital Design Automation
Fall 2025

More Hardware Specialization



Cornell University



Announcements

- ▶ In-class tutorial: Thu 9/4, led by Niansong (TA)
 - Bring your laptop
 - Complete the HLS tool setup beforehand
- ▶ TA office hours: posted on Ed
 - Wed: 4:30-5:30pm @ Rhodes 312
 - Fri: 4:30-5:30pm @ Rhodes 380
- ▶ Lab 1 will be released soon
- ▶ Paper reading next Tuesday 9/9
 - A. Boutros and V. Betz, “[FPGA Architecture: Principles and Progression](#)”, IEEE CAS-M 2021

Review: How to Improve Energy Efficiency

Do less work

- **Amortize overhead** in control and data supply across multiple instructions

Do even less work

- **Use smaller (or simpler) data** => cheaper operations, lower storage & communication costs
- **Move data locally and directly**
 - Store data nearby in simpler memory
 - Wire compute units for direct communication when possible

Common HW Specialization Techniques

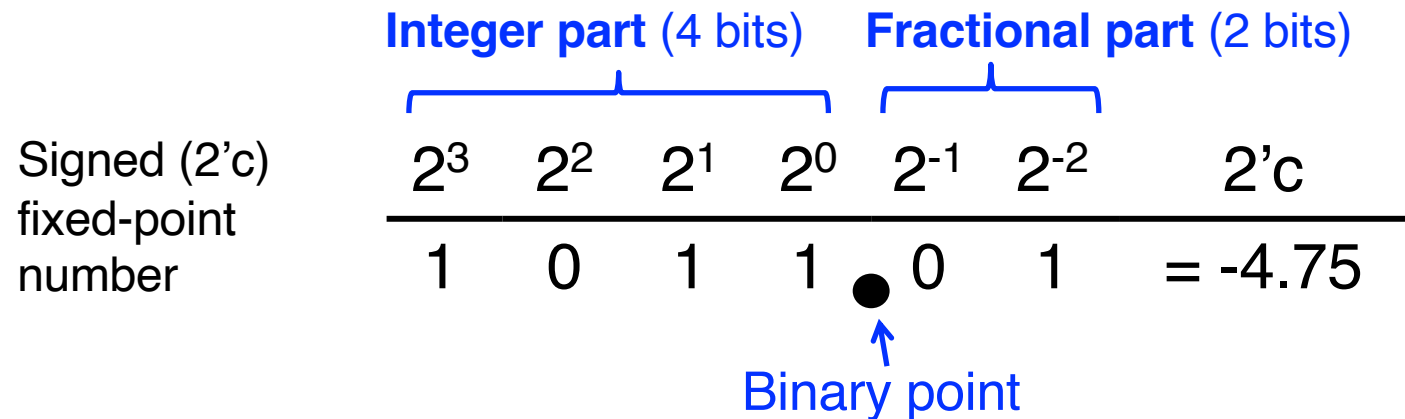
- Lab 1 — **Custom Compute Units:** Use complex instructions to amortize overhead (e.g., SIMD, “ASIC”-in-an-instruction)
- Lab 2 — **Custom Numeric Types:** Balance accuracy and efficiency with data types that use smaller bit widths or simpler arithmetic
- Lab 3 — **Custom Memory Hierarchy:** Exploit data access patterns to reduce energy per memory operation
- Lab 4 — **Custom Communication Architecture:** Tailor on-chip networks to data movement patterns

Agenda

- ▶ Common hardware specialization techniques
 - Custom numeric: a deeper look at fixed-point types
 - Custom memory hierarchy: a case study
 - Custom communication architecture: systolic arrays
- ▶ Roofline-based performance modeling
 - Operational intensity (OI) analysis

Review: Fixed-Point Representation

- ▶ The positional binary encoding can also represent fractional values, by using a **fixed** position of the binary point and place values with negative exponents
 - (-) Less convenient to use in software, compared to floating point
 - (+) Much more efficient in hardware



Overflow and Quantization Issues

- ▶ When working with fixed-point types, quantization and overflow issues often arise due to the limited precision and range of representation
- ▶ **Overflow** occurs when a value exceeds the maximum representable range for the fixed-point format
 - In our class, overflow mainly concerns the integer part of the fixed-point number
- ▶ **Quantization** is needed when converting real numbers from a higher-precision format (like float) to fixed point

Common Modes for Handling Overflow (1)

- ▶ **Wrapping** or wraparound: The value wraps around within the range using modulo arithmetic
 - Efficient in hardware, as it involves simply dropping the MSB(s) of the original number

-2^3	2^2	2^1	2^0	2^{-1}	2^{-2}	2^c
1	0	1	1	0	1	= -4.75


Dropping MSB
when integer width
is reduced

						2^c
0	1	1	0	1		= ??

Wrapping can cause a negative number to become positive, or a positive to negative

Common Modes for Handling Overflow (2)

- ▶ **Saturation:** The closest representable value (either maximum or minimum) is used, preventing wraparound
 - Implementing saturation in hardware requires additional logic to check for overflows and apply the clamping

	-2^3	2^2	2^1	2^0	2^{-1}	2^{-2}	2^c	
	1	0	1	1	0	1	= -4.75	
Saturation								
	?	?	?	?	?	?	= ??	2^c

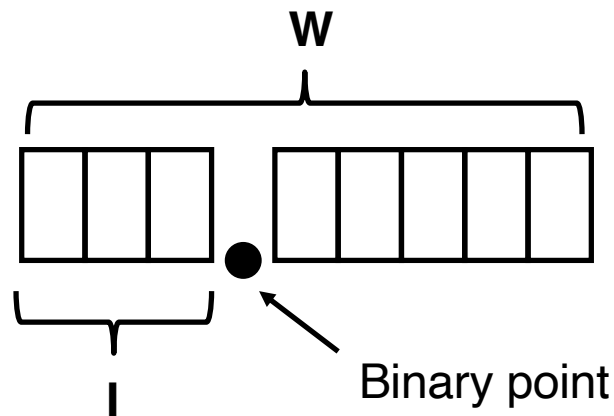
What is the saturated result for the example above?

Common Quantization Modes

- ▶ **Truncation:** Cuts off the excess bits that don't fit in the target precision
 - Efficient in hardware, as it involves simply dropping the LSB(s) of the original number
- ▶ **Rounding:** Rounds the value to a representable fixed-point number, potentially reducing quantization error compared to truncation
 - **Round to nearest:** The most common method, where the value is rounded to the nearest fixed-point representation
 - **Round toward zero:** Rounds towards zero, effectively truncating the fractional part
 - **Round toward infinity:** Rounds away from zero, towards positive or negative infinity

Fixed-Point Types in Vivado HLS

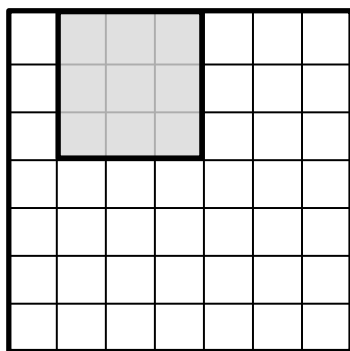
- ▶ **ap_fixed** is a templated C++ data type used for representing fixed-point numbers
 - Signed: **ap_fixed**; Unsigned: **ap_ufixed**
 - Template parameters `ap_(u)fixed<W, I, Q, O>`
 - W: total bitwidth
 - I: integer bitwidth
 - Q: quantization mode (optional, default is AP_TRN)
 - O: overflow mode (optional, default is AP_WRAP)



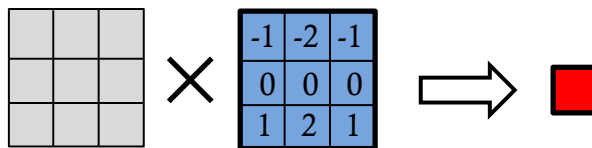
Custom Memory Hierarchy: Case Study on Convolution

- ▶ **Convolution** is pervasive in image/video processing and ML – performed over overlapping windows (aka stencils)

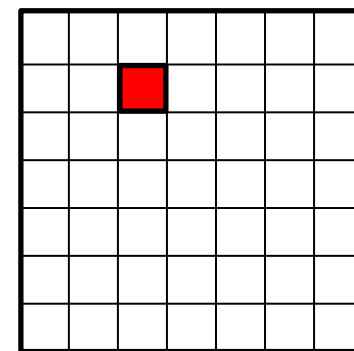
$$(Img \otimes f)_{\left[n+\frac{k-1}{2}, m+\frac{k-1}{2}\right]} = \sum_{i=0}^{k-1} \sum_{j=0}^{k-1} Img_{[n+i][m+j]} \cdot f_{[i,j]}$$



Input image
frame



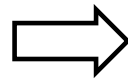
A K-by-K **dot product** is performed
for each output pixel (K=3 here)



Output image
frame

An Application of Convolution: Edge Detection

- ▶ Identifies discontinuities in an image where brightness (or image intensity) changes sharply
 - Very useful for feature extractions in computer vision



Sobel operator
 $G = (G_X, G_Y)$

$$G_X = \begin{bmatrix} -1 & -2 & -1 \\ 0 & 0 & 0 \\ 1 & 2 & 1 \end{bmatrix}$$

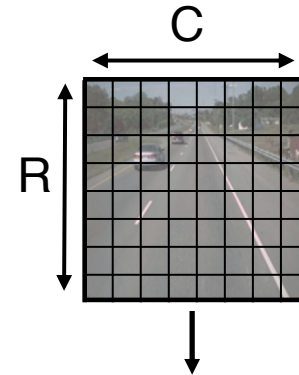
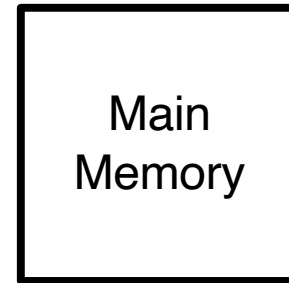
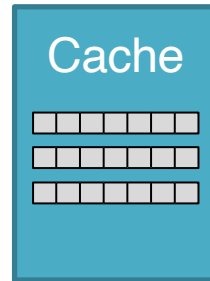
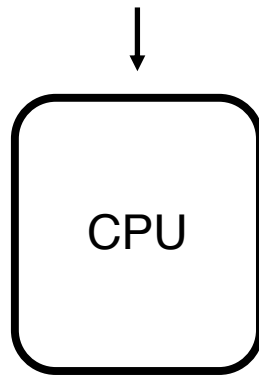
$$G_Y = \begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix}$$



Figures: Pilho Kim, GaTech

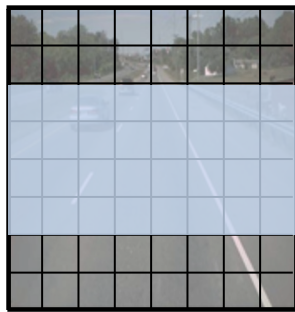
CPU Implementation of a 3x3 Convolution

```
for (r = 1; r < R; r++)  
  for (c = 1; c < C; c++)  
    for (i = 0; i < K; i++)  
      for (j = 0; j < K; j++)  
        out[r][c] += img[r+i-1][c+j-1] * f[i][j];
```

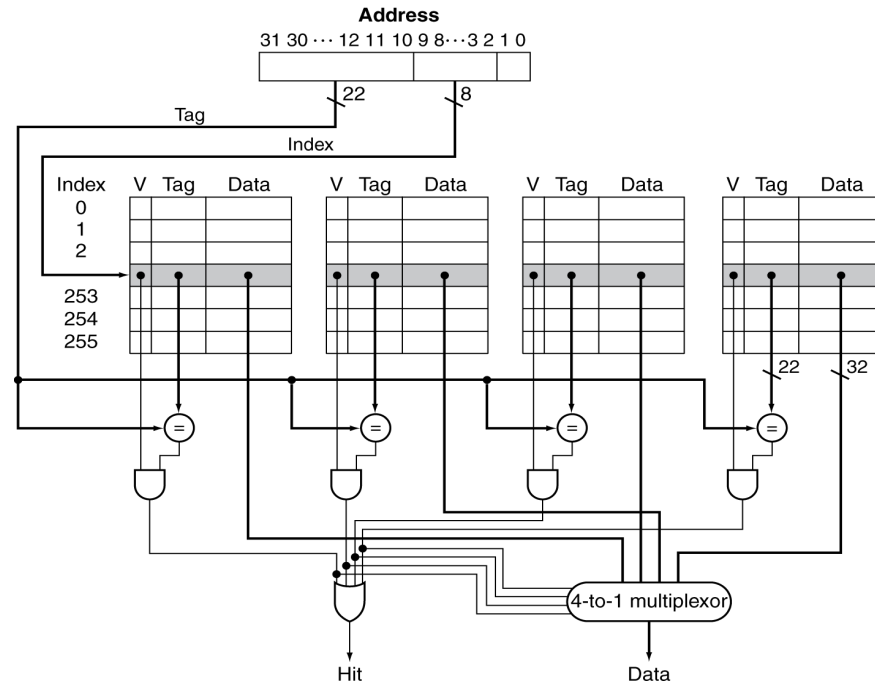


General-Purpose Cache for Convolution

- ▶ A general-purpose cache can significantly reduce external memory accesses, but it is costly and incurs high energy overhead

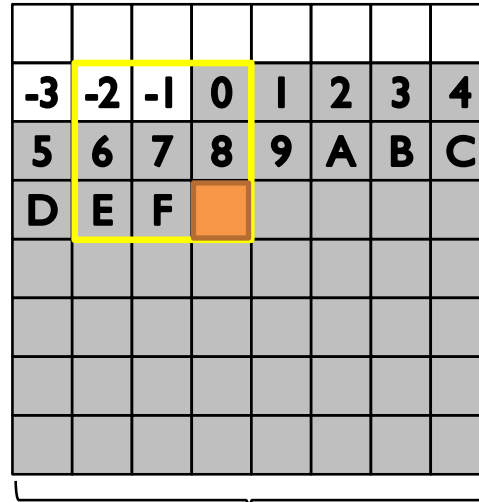


A subset of image data stored in cached



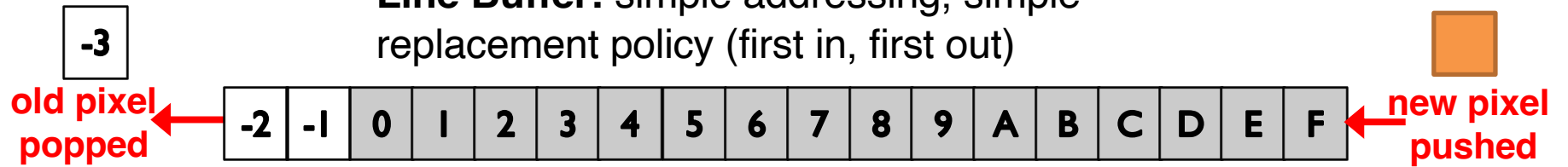
Line Buffer: Customized “Cache” for Convolution

- ▶ **“Cache” the input pixels in a line buffer:** Each time we move the KxK window (in yellow) to the right and push in a new pixel (in orange) to the specialized “cache”



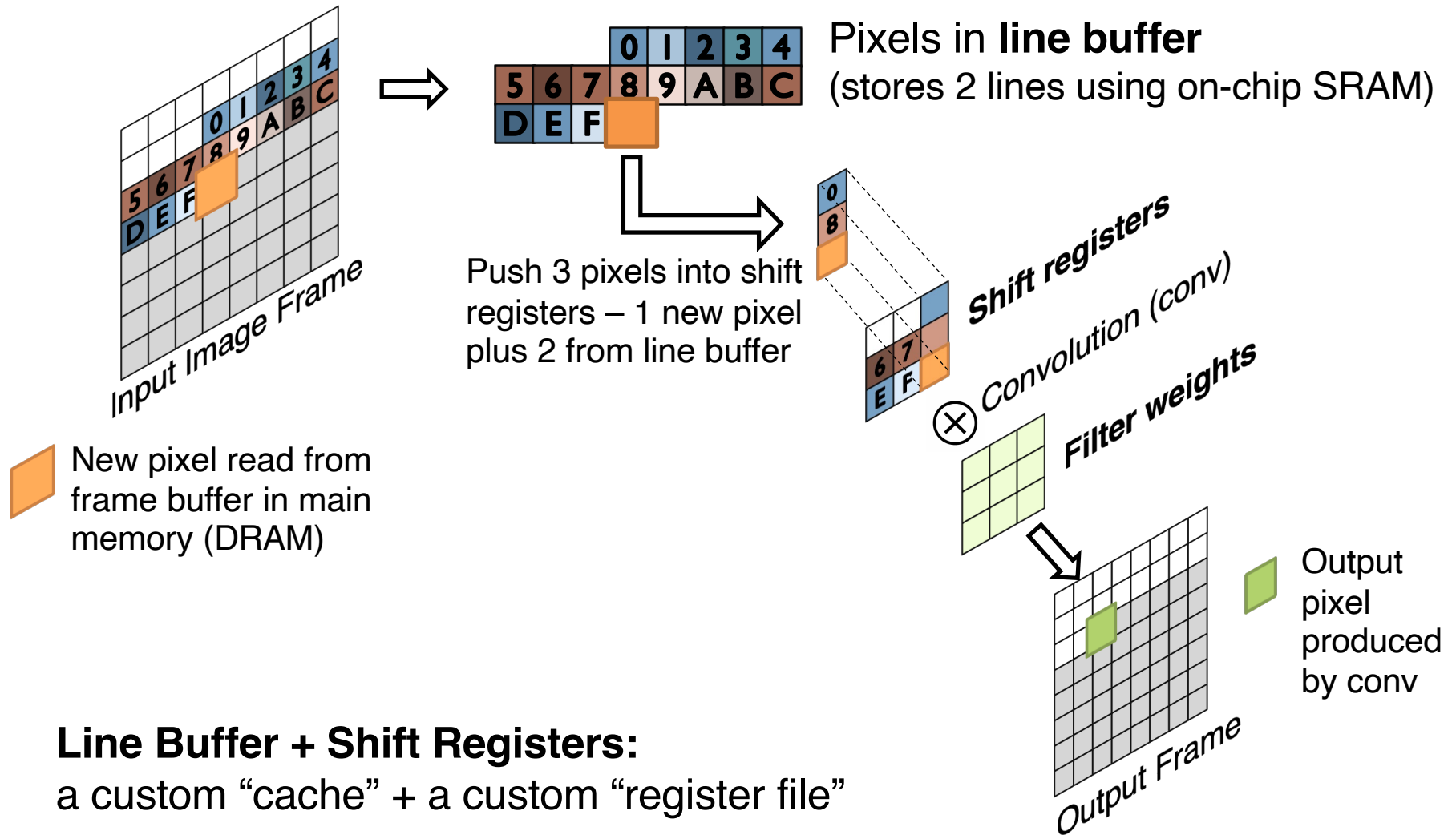
C (number of columns or width)

Line Buffer: simple addressing, simple replacement policy (first in, first out)



$(K-1)*C+K$ pixels in flight, $K=3$ and $C=8$ here

A More Complete Picture of the Custom On-Chip Memory Hierarchy



Line Buffer + Shift Registers:
 a custom “cache” + a custom “register file”

Custom Communication Architecture: Systolic Arrays as an Example

- ▶ An array of processing elements (PEs) that process data in a systolic manner using nearest-neighbor communication

Systolic Arrays (for VLSI)

H. T. Kung[†] and Charles E. Leiserson[†]

*And now I see with eye serene
The very pulse of the machine.
--William Wordsworth*

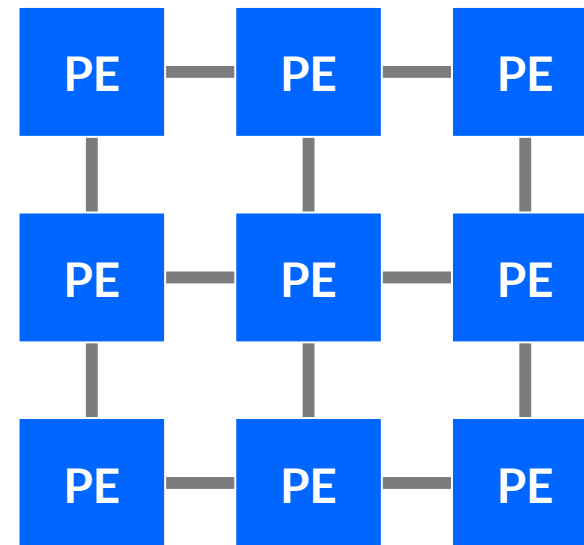
Abstract

A **systolic** system is a network of processors which rhythmically compute and pass data through the system. Physiologists use the word "systole" to refer to the rhythmically recurrent contraction of the heart and arteries which pulses blood through the body. In a **systolic** computing system, the function of a processor is analogous to that of the heart. Every processor regularly pumps data in and out, each time performing some short computation, so that a regular flow of data is kept up in the network.

Many basic matrix computations can be pipelined elegantly and efficiently on **systolic** networks having an array structure. As an example, hexagonally connected processors can optimally perform matrix multiplication. Surprisingly, a similar **systolic** array can compute the LU-decomposition of a matrix. These **systolic** arrays enjoy simple and regular communication paths, and almost all processors used in the networks are identical. As a result, special purpose hardware devices based on **systolic** arrays can be built inexpensively using the **VLSI** technology.

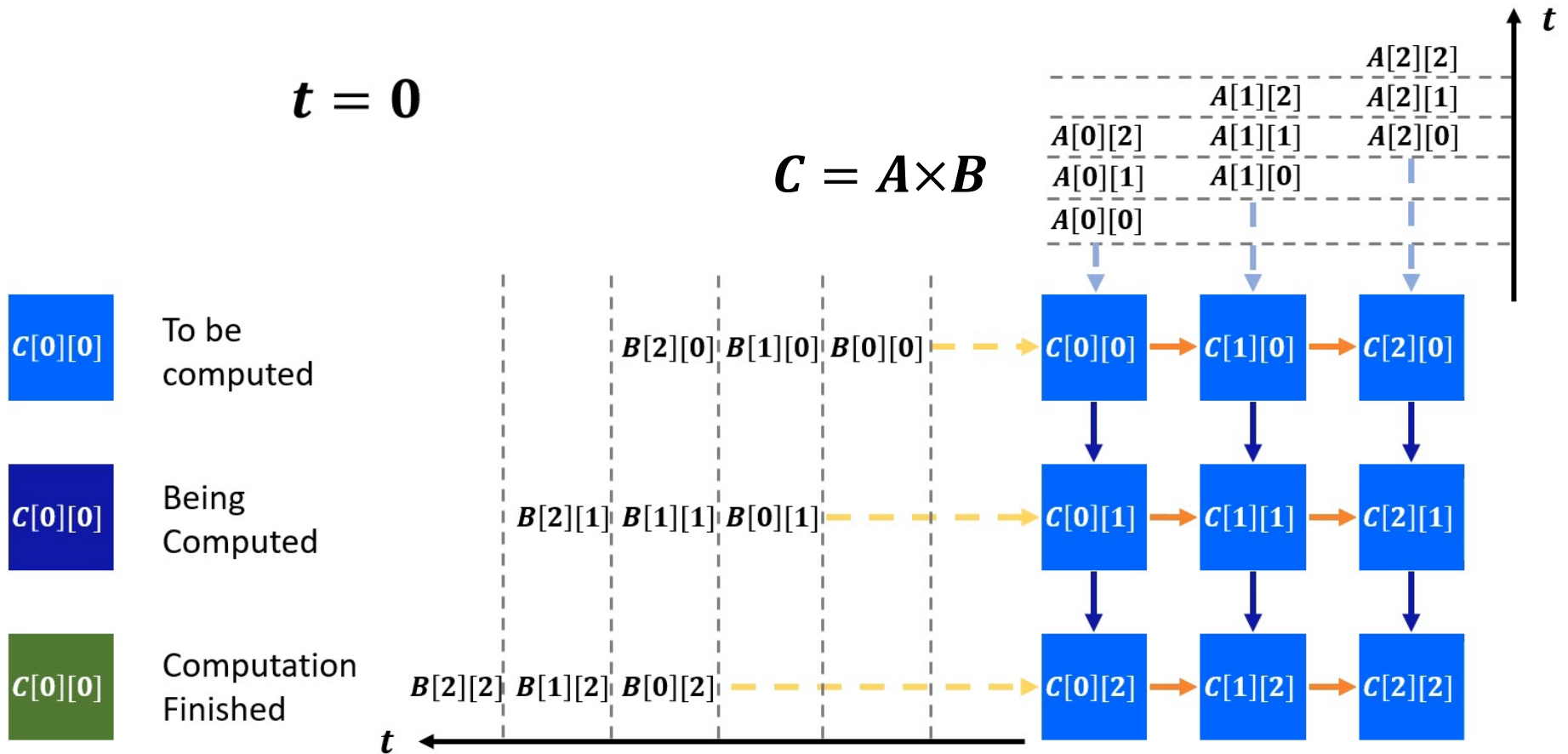
1. Introduction

Developments in microelectronics have revolutionized computer design. Integrated circuit technology has increased the number and complexity of components that can fit on a chip or a printed circuit board. Component density has been doubling every one-to-two years and already, a multiplier can fit on a very large scale integrated



- + Simple & regular design
- + Massive parallelism
- + **Short nearest-neighbor interconnection**
- + Balancing compute with I/O

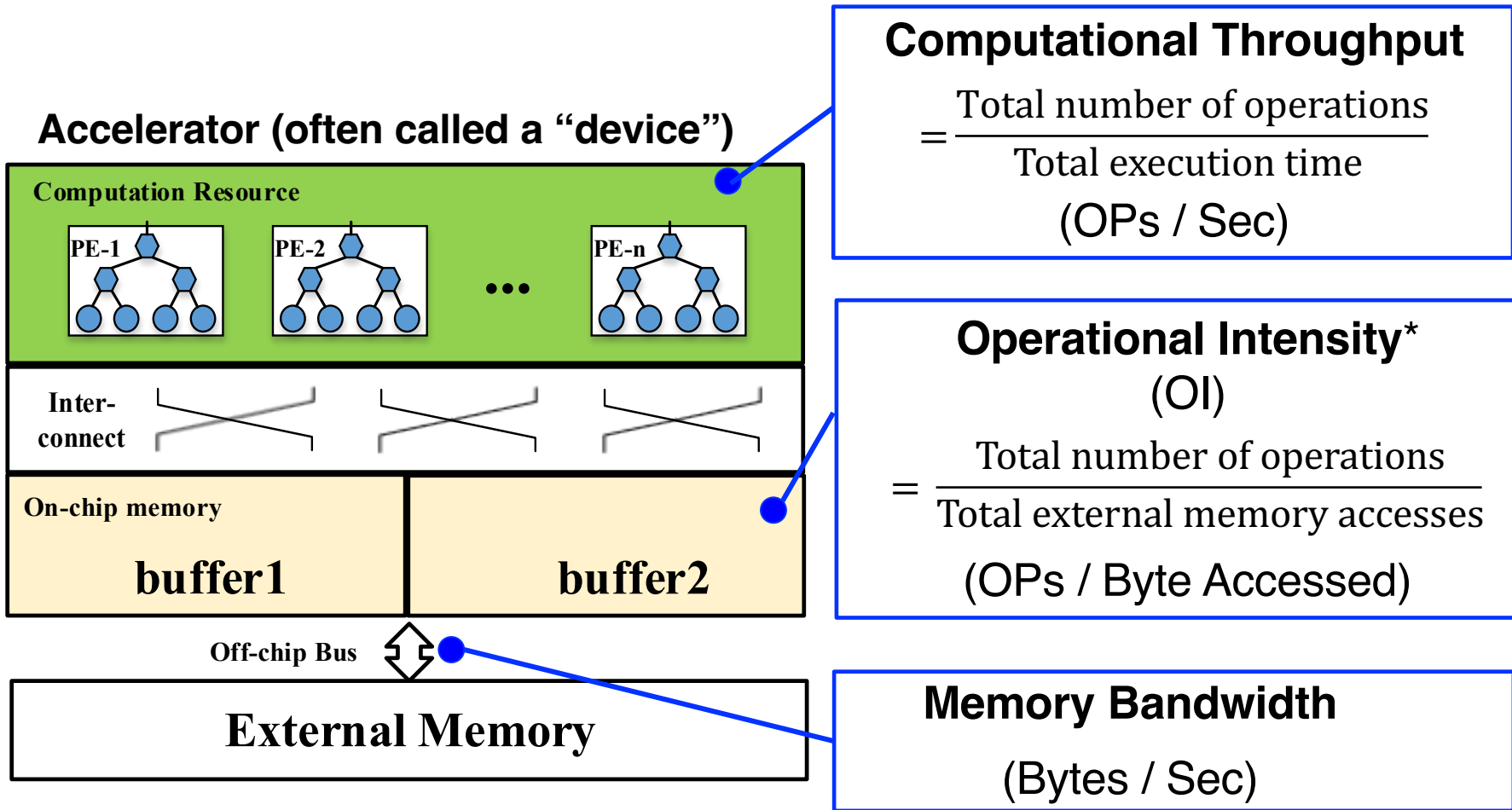
Matrix Multiplication (MM) on a Systolic Array



Accelerator Performance Modeling

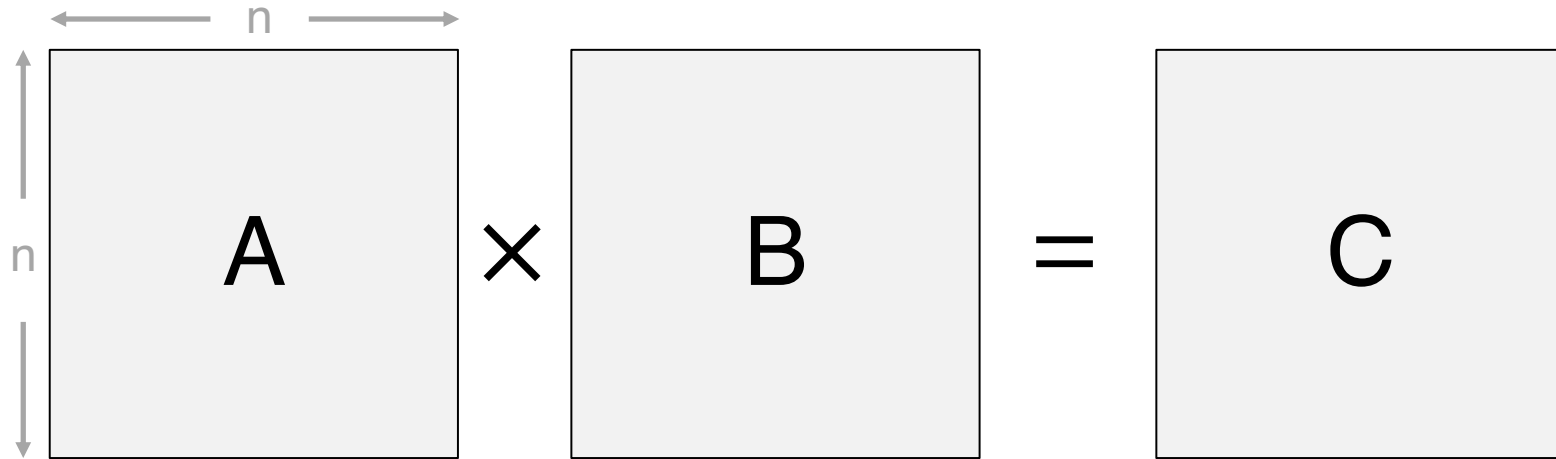
- ▶ How do we model the performance of an application running on a compute device (e.g., an accelerator)?
 - What's the maximum (peak) throughput of the device?
 - What's the actual attainable throughput?
 - Where're the bottlenecks?
- ▶ It requires characterization of both application and hardware

Key Metrics to Consider



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

Estimating OI: Matrix Matrix Multiplication



A, B, and C use single-precision floats (32b per element, i.e., 4 bytes each)

Number of operations: $2n^3$ (multiply & add counted separately)

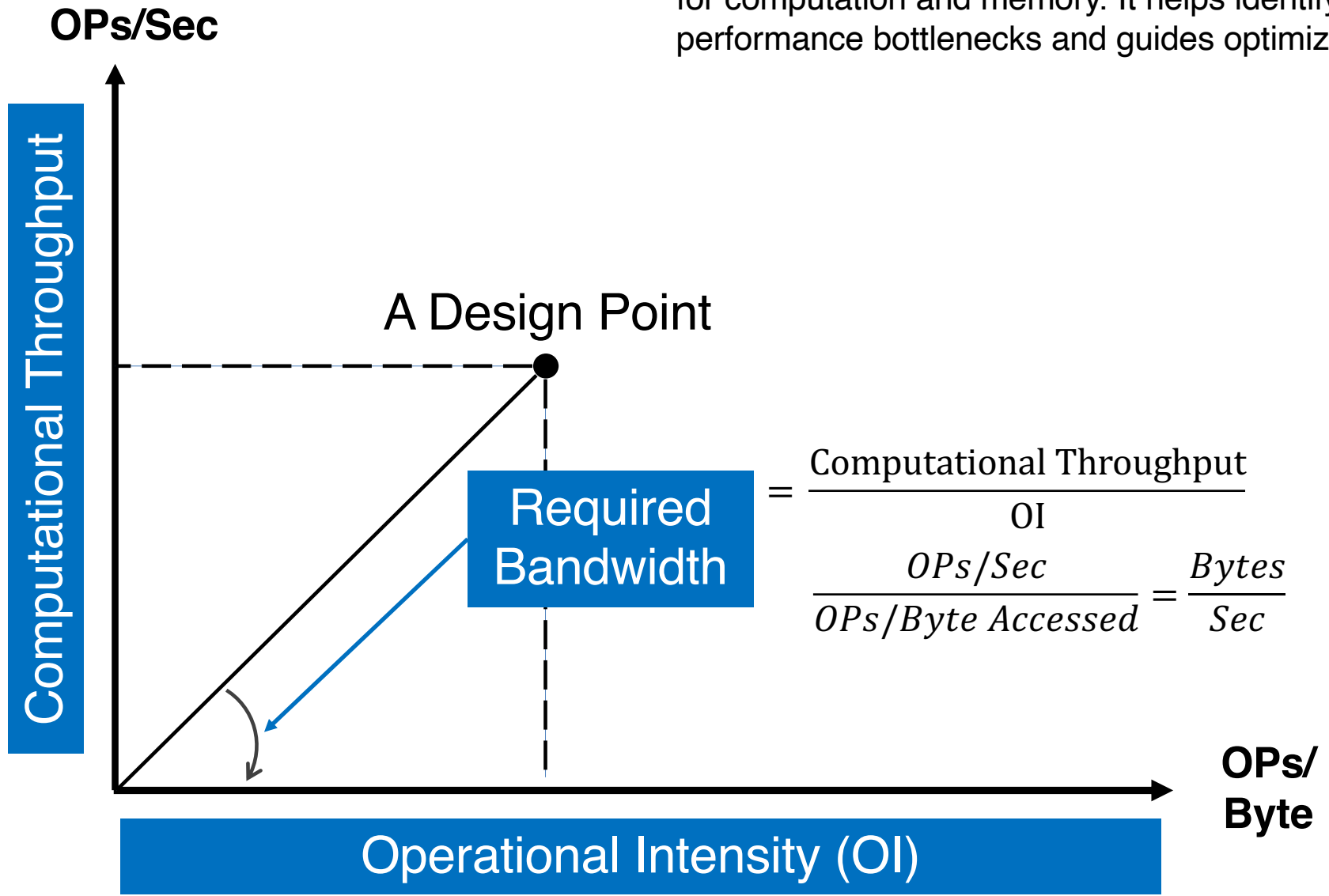
External mem accesses: $3n^2 \times 4$ bytes (reads or writes)

(assuming perfect data reuse)

$$OI = \frac{\# \text{ of operations}}{\# \text{ of bytes read/written}} = \frac{2n^3}{12n^2} = \frac{n}{6}$$

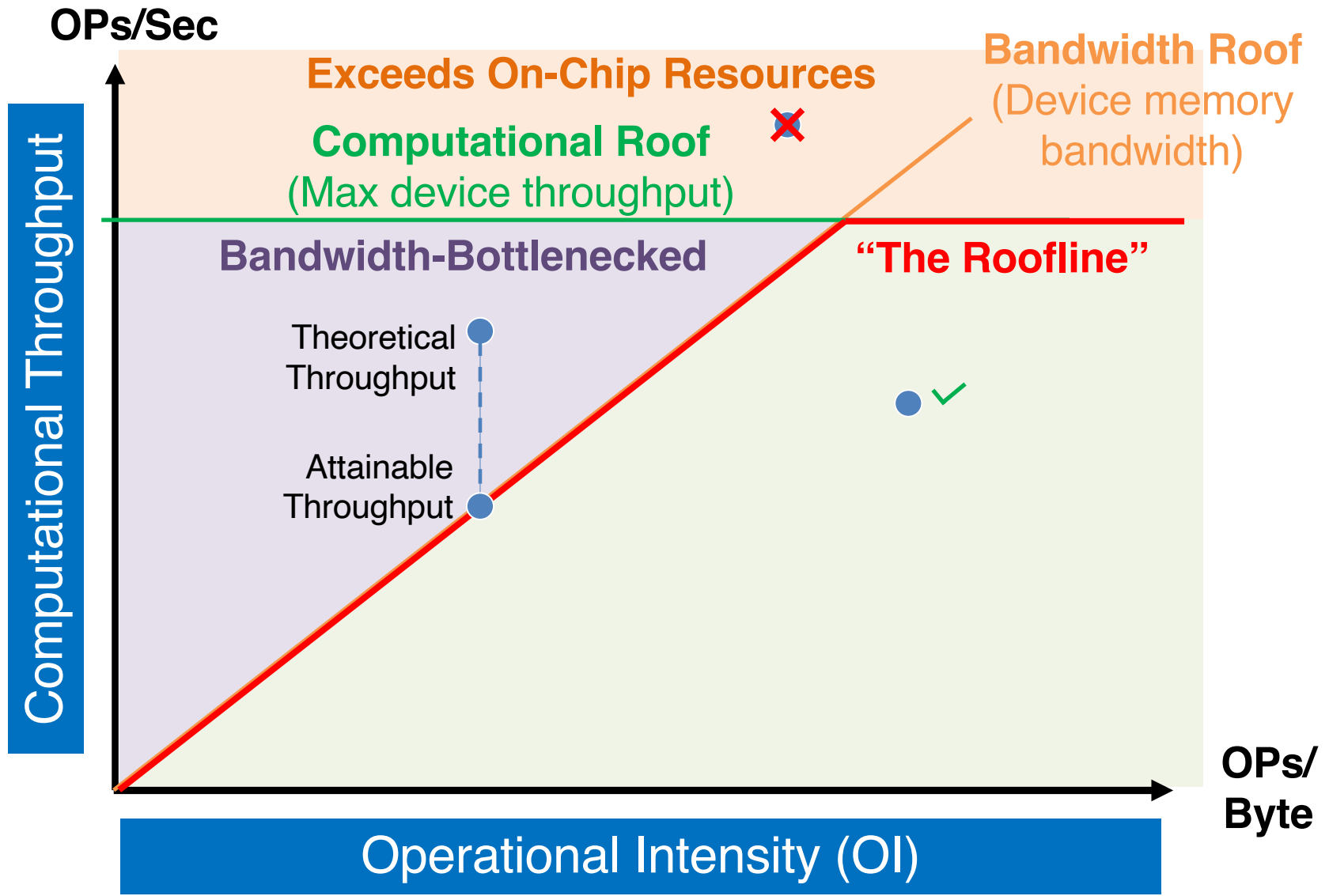
Roofline Model [1]

The Roofline Model visualizes a system's performance limits, showing how an application's throughput compares to the theoretical maximums for computation and memory. It helps identify performance bottlenecks and guides optimization.



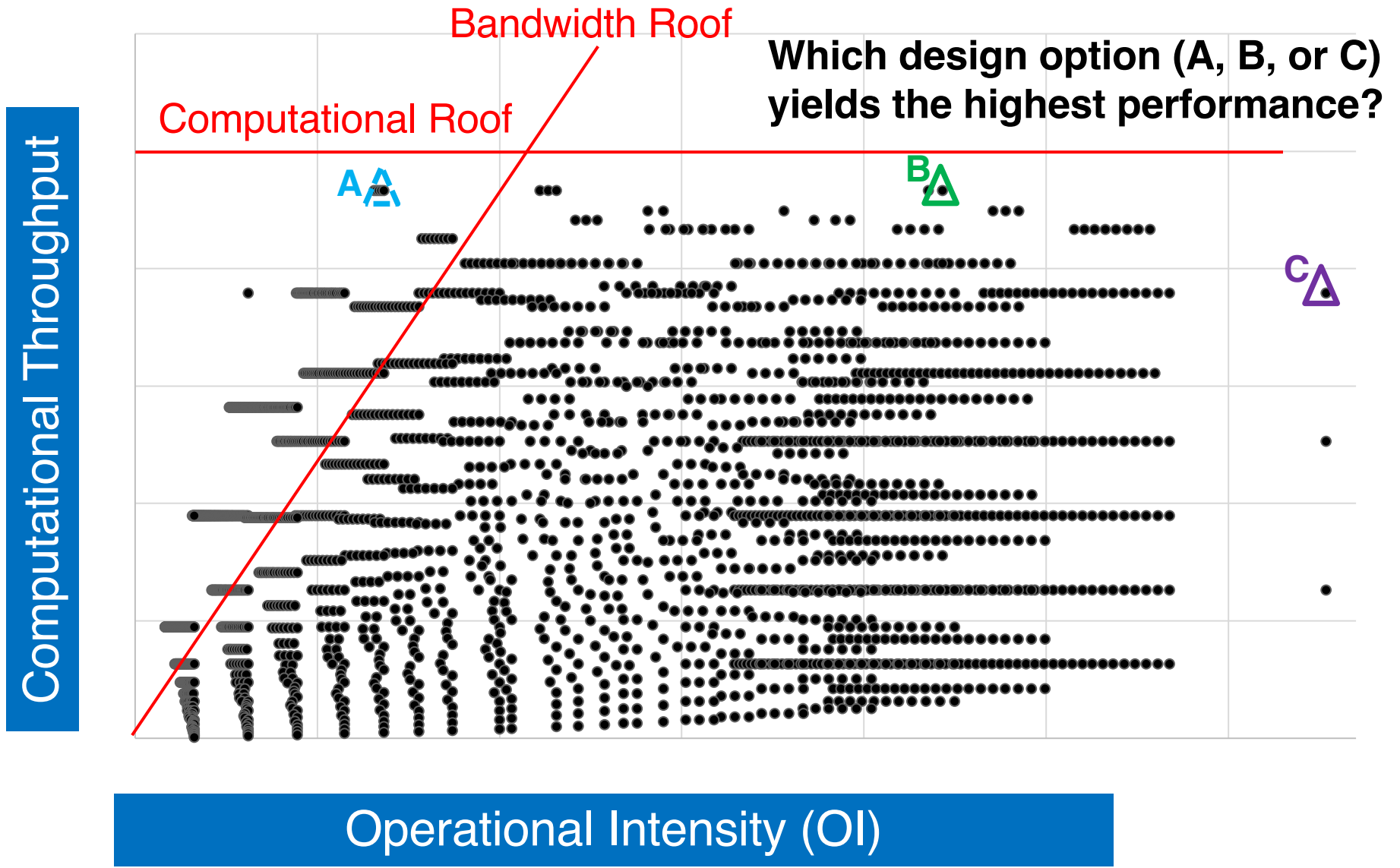
[1] S. Williams, A. Waterman, and D. Patterson, Roofline: an insightful visual performance model for multicore architectures, CACM, 2009.

Roofline Model [1]

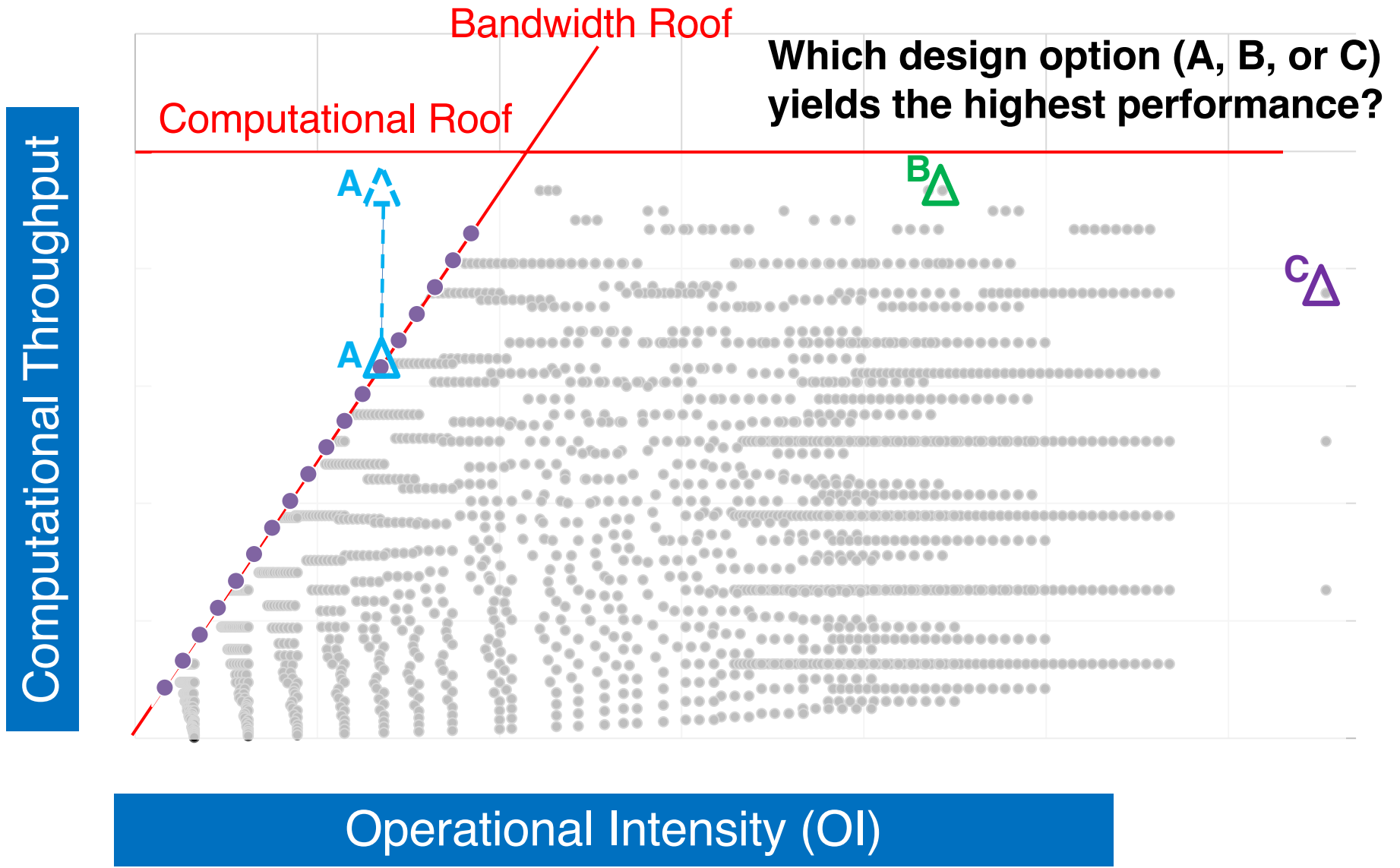


[1] S. Williams, A. Waterman, and D. Patterson, Roofline: an insightful visual performance model for multicore architectures, CACM, 2009.

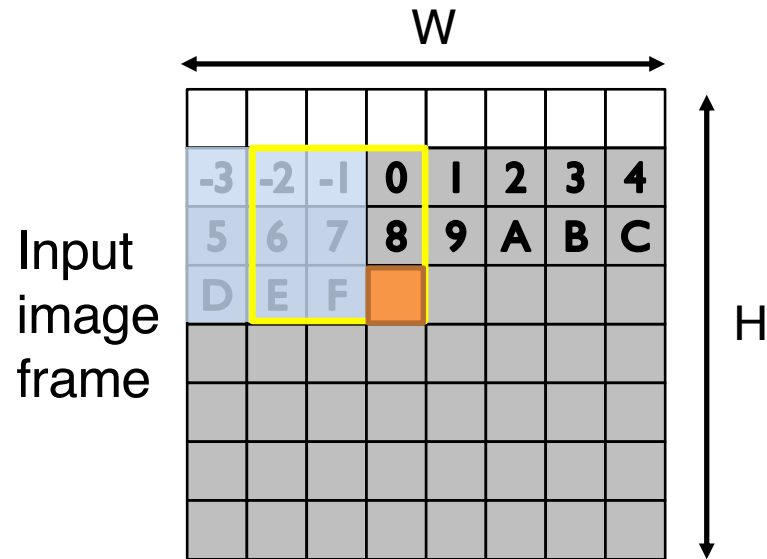
Design Space Exploration with Roofline



Design Space Exploration with Roofline

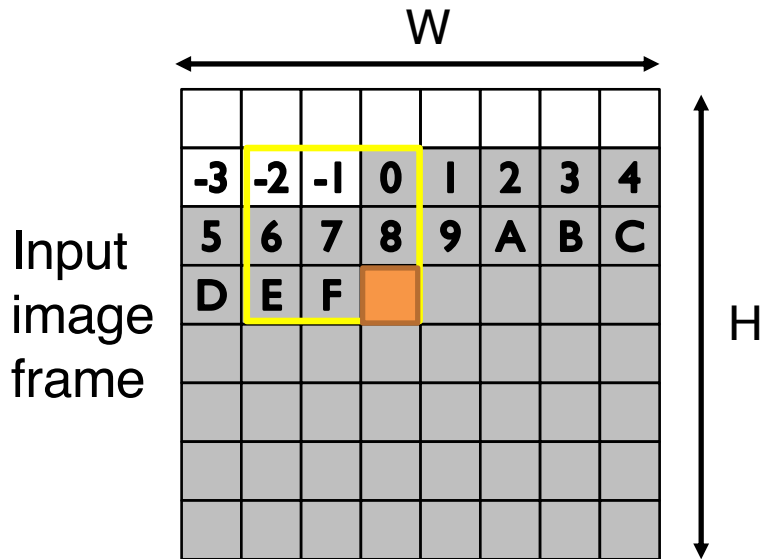


Exercise: OI Analysis of 2D Convolution



Estimate the OI for the 2D convolution kernel both without and with the use of a line buffer

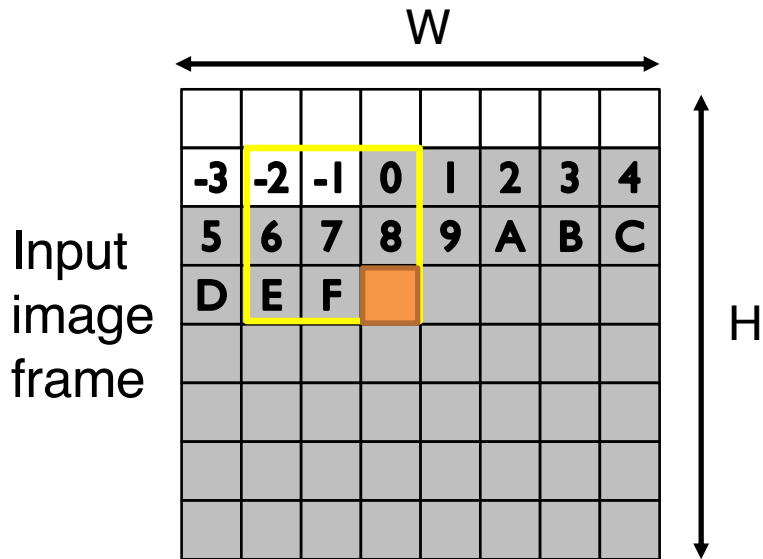
2D Convolution: OI without Line Buffer



```
for (r = 1; r < H; r++)
  for (c = 1; c < W; c++)
    for (i = 0; i < 3; i++)
      for (j = 0; j < 3; j++)
        out[r][c] += img[r+i-1][c+j-1] * f[i][j];
```

- ▶ OI without line buffer, i.e., no data reuse
 - Number of operations = $W \cdot H \cdot 9 \cdot 2$
(1 multiply + 1 add per pixel)
 - External mem accesses = $W \cdot H \cdot 9$ bytes (reads) + $W \cdot H$ (writes)
(assuming 1 byte per pixel in grayscale)

2D Convolution: OI with Line Buffer



```
for (r = 1; r < H; r++)  
  for (c = 1; c < W; c++)  
    for (i = 0; i < 3; i++)  
      for (j = 0; j < 3; j++)  
        out[r][c] += img[r+i-1][c+j-1] * f[i][j];
```

- ▶ OI without line buffer, i.e., no data reuse
 - Number of operations = $W \cdot H \cdot 9 \cdot 2$
(1 multiply + 1 add per pixel)
 - External mem accesses = $W \cdot H \cdot 2$ bytes (reads & writes)
(assuming 1 byte per pixel in grayscale)

Summary

- ▶ End of Dennard scaling leads to increasing **hardware specialization** to sustain improvement in hardware performance and energy efficiency
- ▶ Special-purpose hardware **accelerators** commonly leverage customized (1) compute units, (2) numeric types, (3) memory hierarchy, and (4) communication architecture
- ▶ **Roofline modeling** is a useful tool for first-order analysis of the accelerator performance

Next Lecture

- ▶ TA-led Tutorial on C-Based HLS
 - **Bring your laptop!**

Acknowledgements

- ▶ These slides contain/adapt materials developed by
 - Prof. Jason Cong (UCLA)