# ECE 598NSG/498NSU
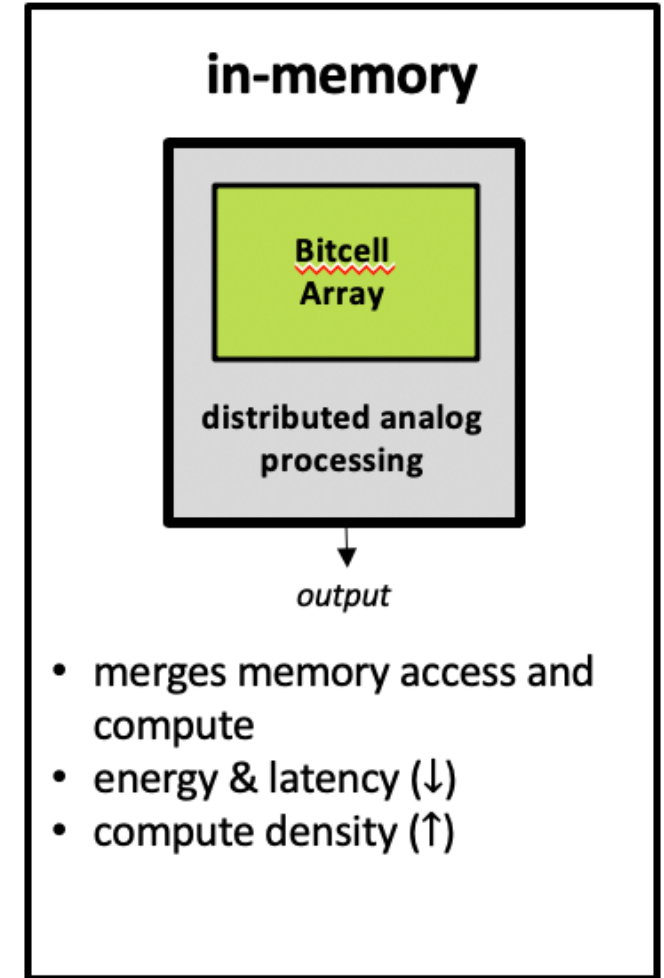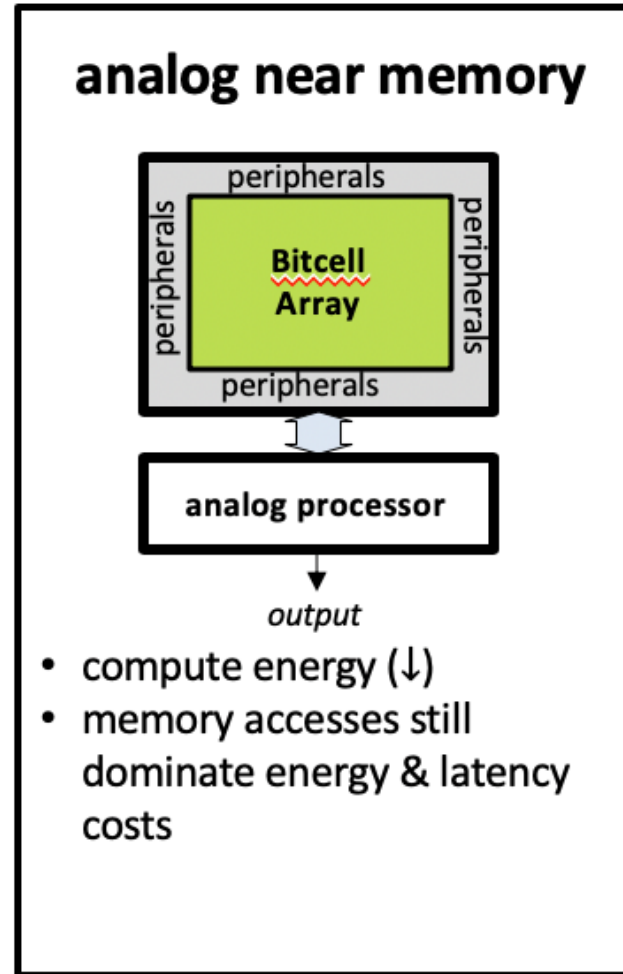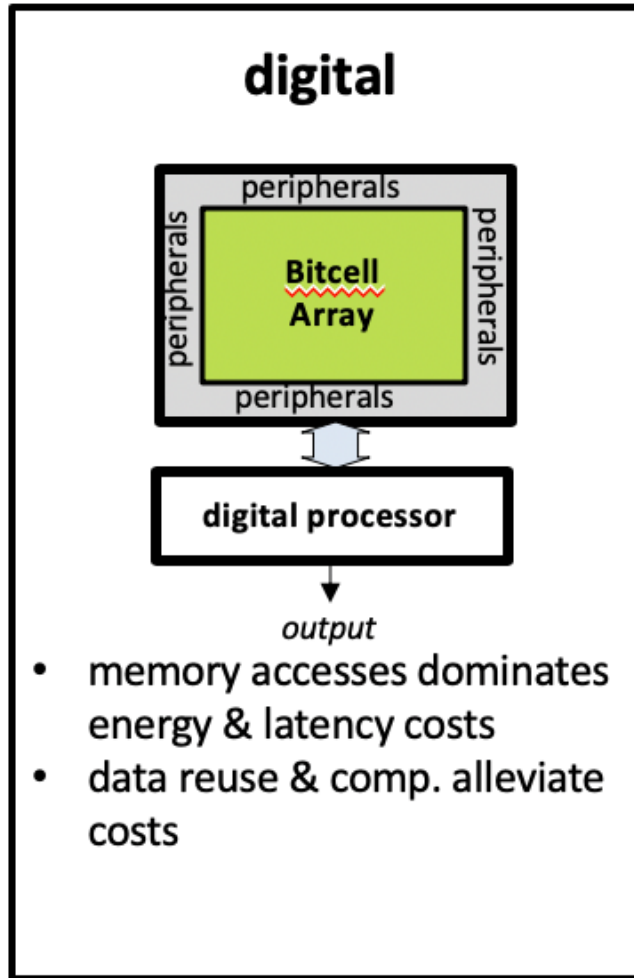# Deep Learning in Hardware
# Fall 2020

## DNN Accelerators – Overview
## Roofline, Floorline plots, Reuse

Naresh Shanbhag

Department of Electrical and Computer Engineering
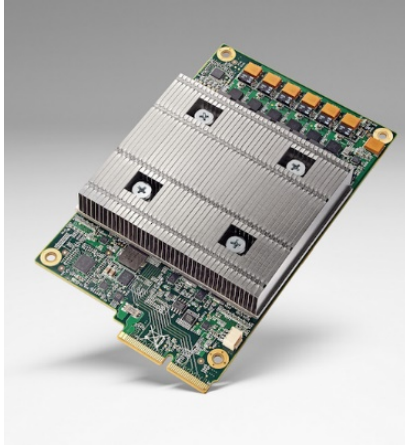
University of Illinois at Urbana-Champaign

http://shanbhag.ece.uiuc.edu

# Architectures for Inference

## digital

peripherals
peripherals
**Bitcell Array**
peripherals
peripherals

**digital processor**

*output*

- memory accesses dominates energy & latency costs
- data reuse & comp. alleviate costs

## analog near memory

peripherals
peripherals
**Bitcell Array**
peripherals
peripherals

**analog processor**

*output*

- compute energy (↓)
- memory accesses still dominate energy & latency costs

## in-memory

**Bitcell Array**

**distributed analog processing**

*output*

- merges memory access and compute
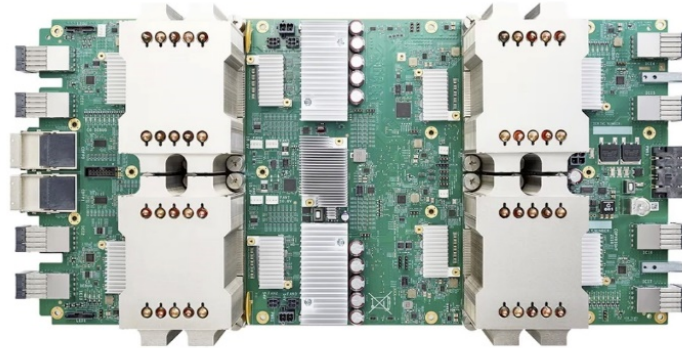- energy & latency (↓)
- compute density (↑)
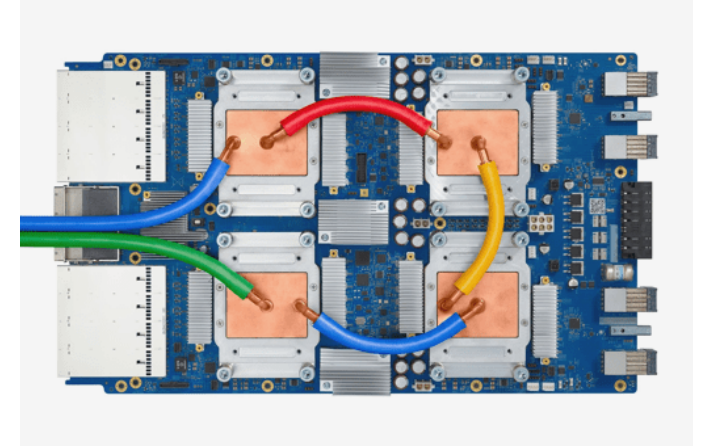
# DNN Accelerators

**TPU V1**



inference only
8b fixed-point

**TPU V2**



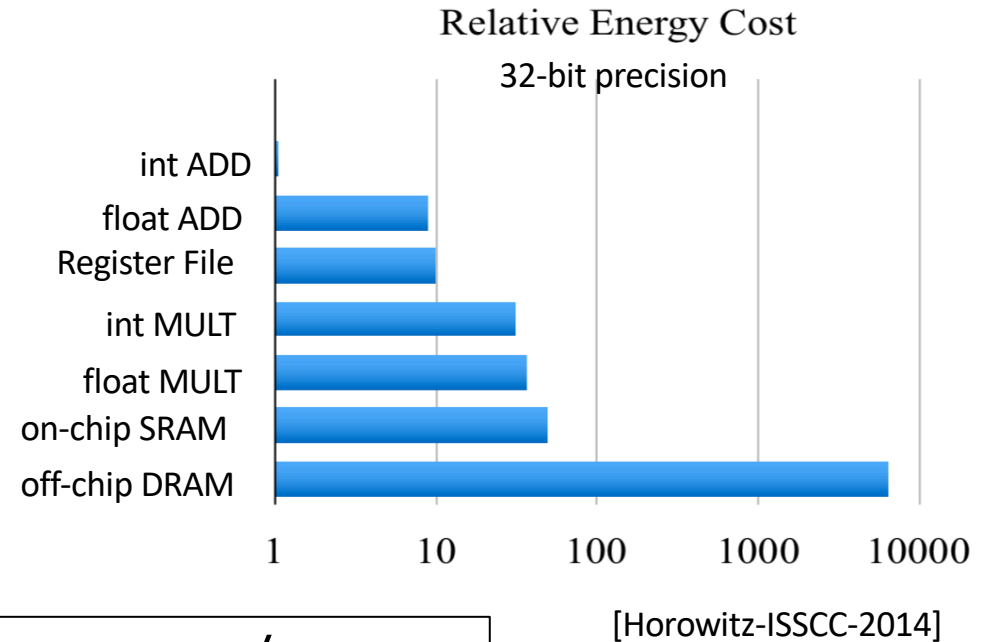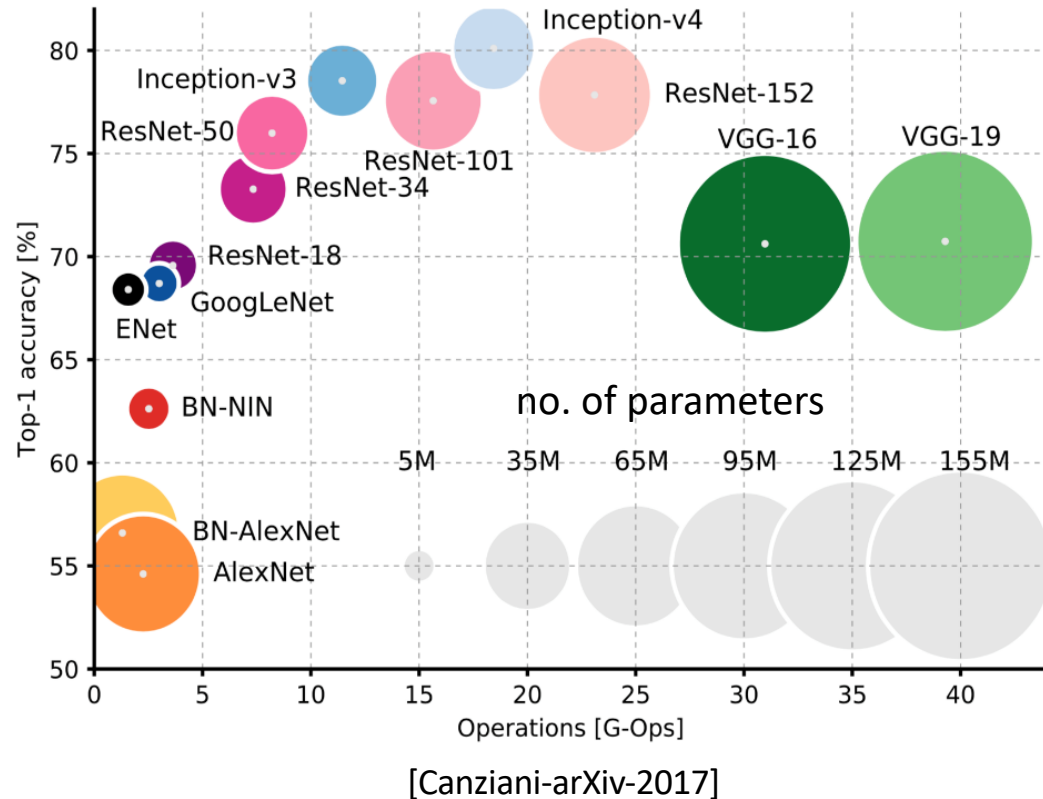inference & training
fixed-point & floating-point

**TPU V3**
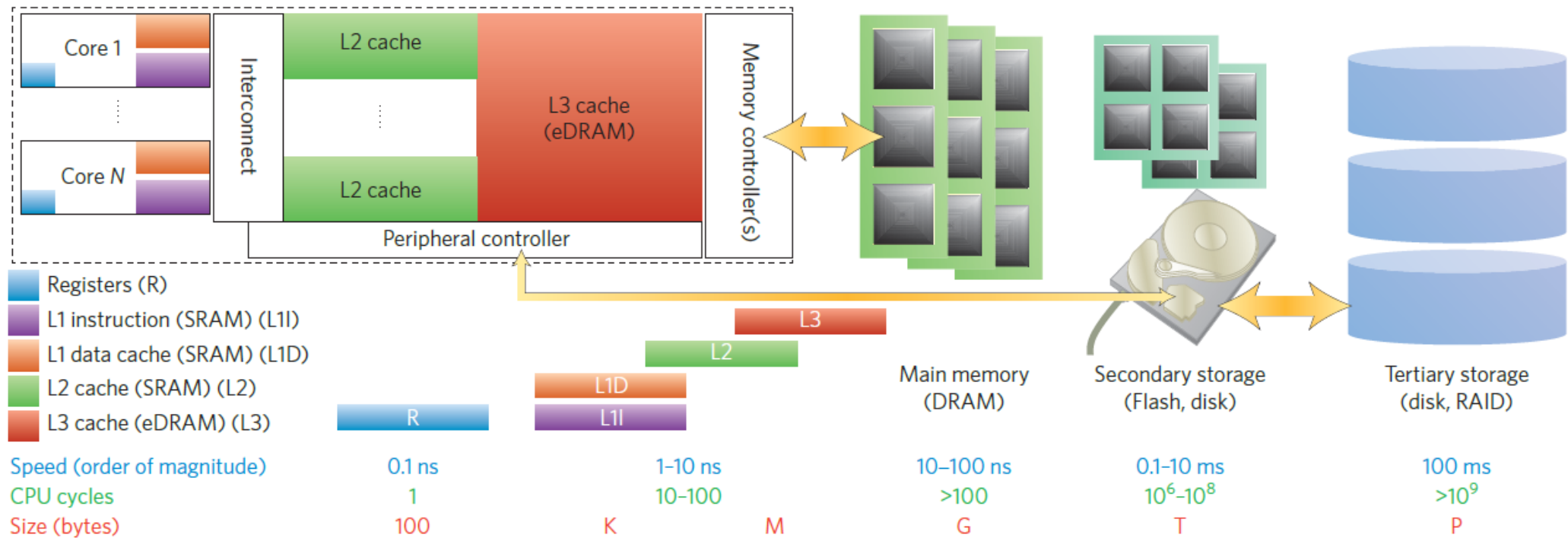


inference & training
fixed-point & floating-point

- going beyond GPUs

- accelerating inference & training in the cloud

- intersection of workload, architecture, circuits, software

# Challenge: Memory Access Energy in DNNs



[Canziani-arXiv-2017]



[Horowitz-ISSCC-2014]

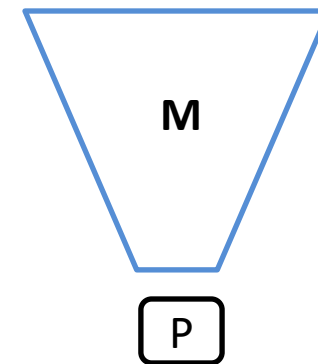$$E_R = \frac{energy/access}{energy/OP}$$

- memory access-to-compute energy consumption ratio

- SRAM accesses: $E_R = 100\times$ ; DRAM accesses: $E_R = 1000\times$
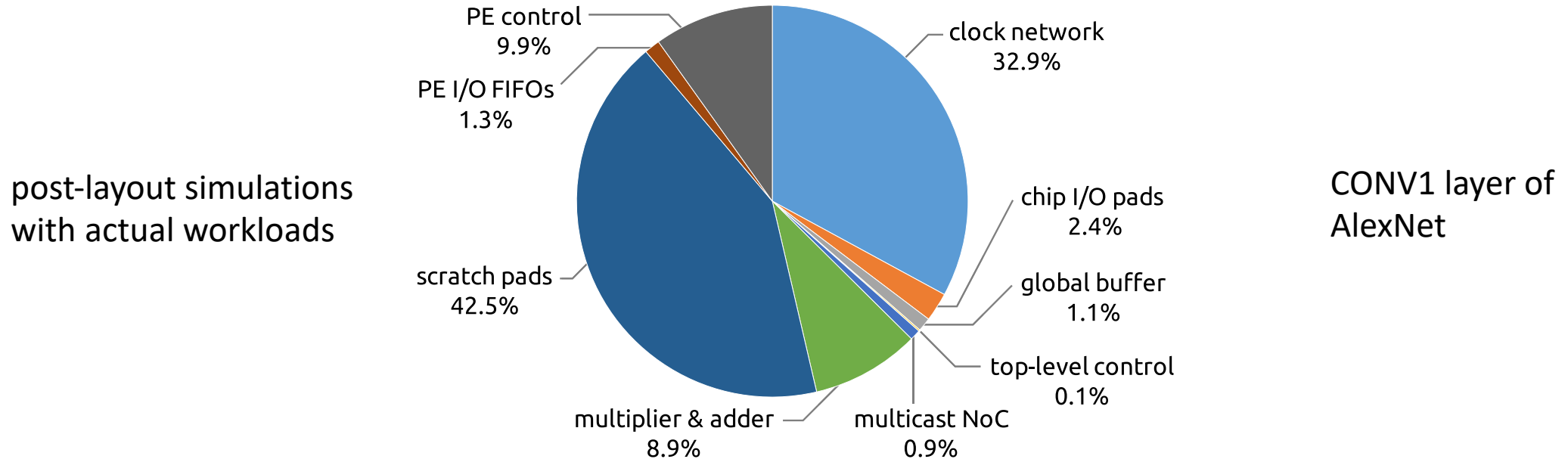
[Wong, Salahuddin, Nature Nanotechnology, 2015]

| | Registers (R) | L1 instruction (SRAM) (L1I) | L1 data cache (SRAM) (L1D) | L2 cache (SRAM) (L2) | L3 cache (eDRAM) (L3) | Main memory (DRAM) | Secondary storage (Flash, disk) | Tertiary storage (disk, RAID) |
|---|---|---|---|---|---|---|---|---|
| Speed (order of magnitude) | 0.1 ns | | 1–10 ns | | | 10–100 ns | 0.1–10 ms | 100 ms |
| CPU cycles | 1 | | 10–100 | | | >100 | $10^6$–$10^8$ | $>10^9$ |
| Size (bytes) | 100 | | K | M | | G | T | P |

the memory funnel

- deep memory hierarchy

- high density → more latency, energy

- data movement → dominant cause of energy dissipation in machine learning applications
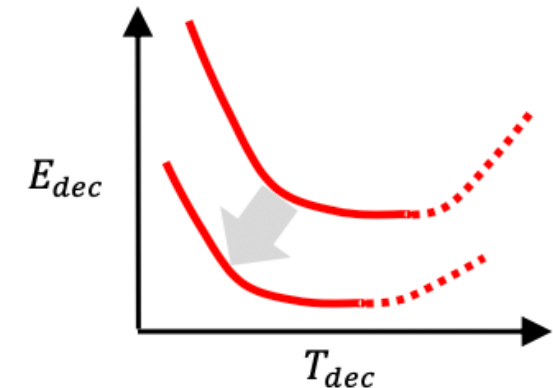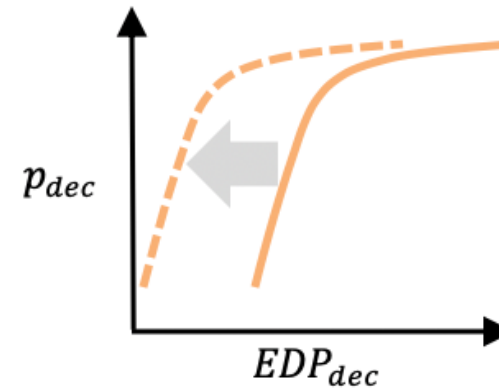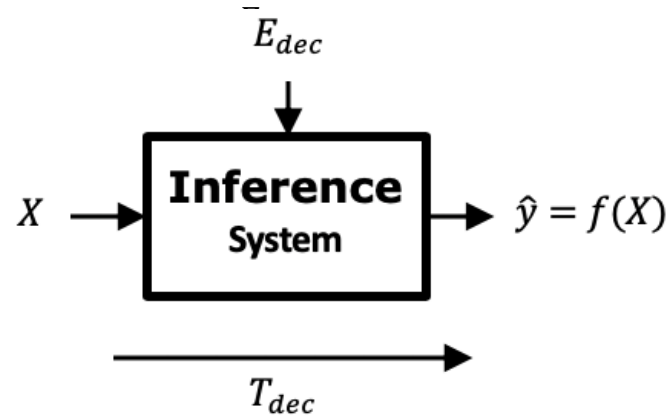
# Power Breakdown

**[Eyeriss, Y.-H. Chen, JSSC 2017]**

post-layout simulations
with actual workloads

CONV1 layer of
AlexNet



- PE control 9.9%
- PE I/O FIFOs 1.3%
- clock network 32.9%
- chip I/O pads 2.4%
- global buffer 1.1%
- top-level control 0.1%
- multicast NoC 0.9%
- multiplier & adder 8.9%
- scratch pads 42.5%

- Data movement costs 45% of total power

- Computation costs about 10%

- Clock network = 33% (hidden cost of digital realizations)

# System Metrics for Inference

**system-level energy vs. latency vs. accuracy trade-off**



- $T_{dec}$: decision latency
- $p_{dec}$: decision accuracy
- $EDP_{dec} = E_{dec} \times T_{dec}$: decision EDP
- $f_{dec}$: decision throughput

- $E_{dec} \propto 1/T_{dec} \propto$ data & network size
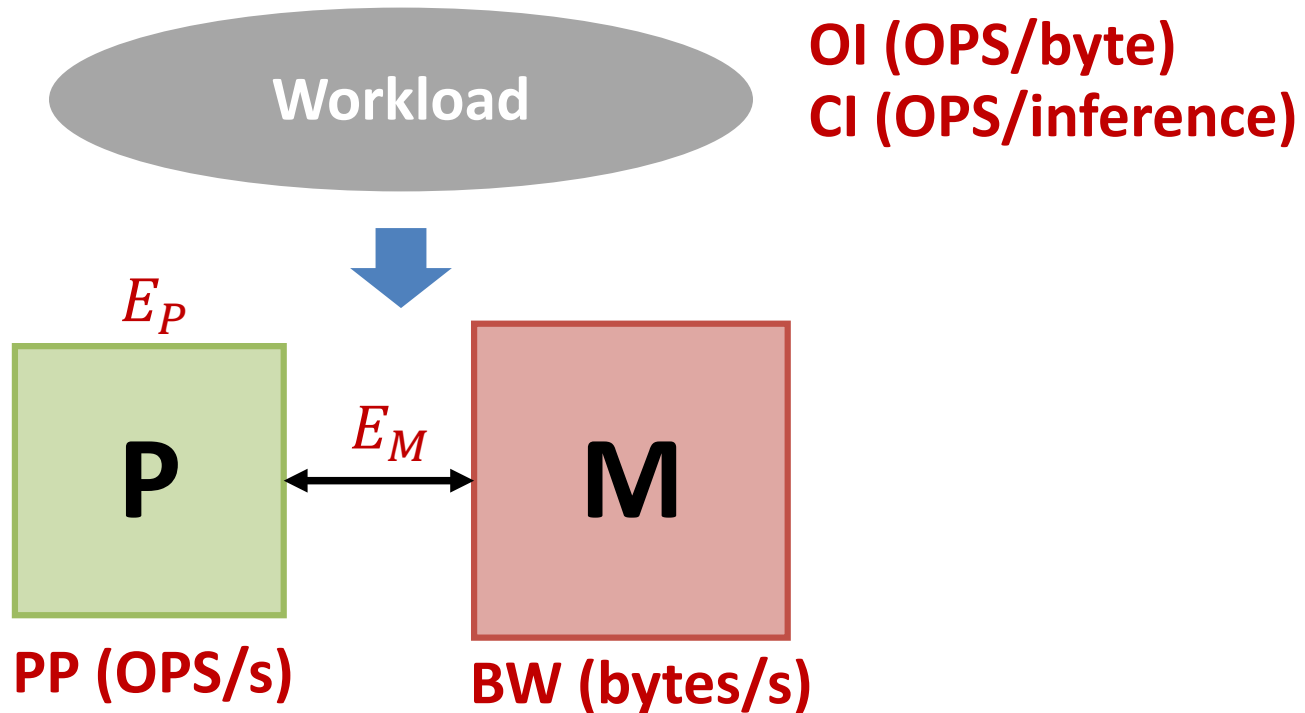- $p_{dec} \propto EDP_{dec}$

# Outline

- DNN architecture abstraction

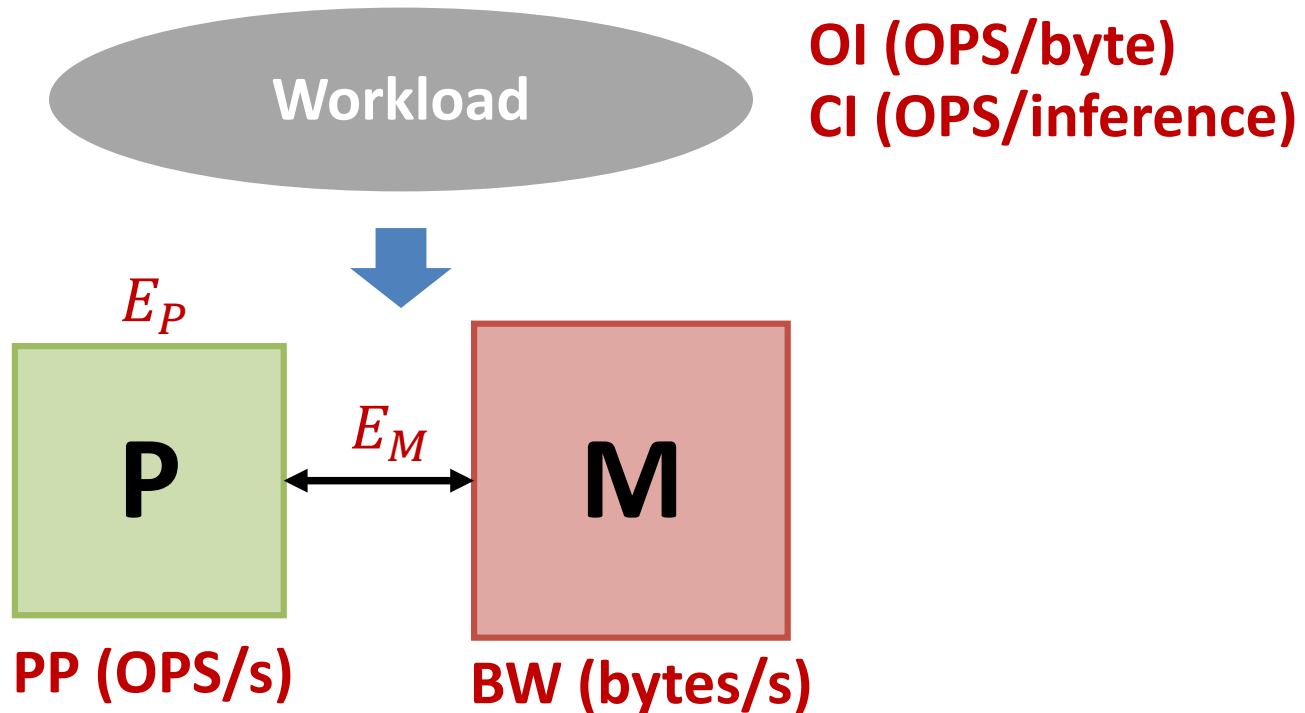- visualization via the roofline & floorline plots

- data reuse

# DNN Architectures – An Abstract View
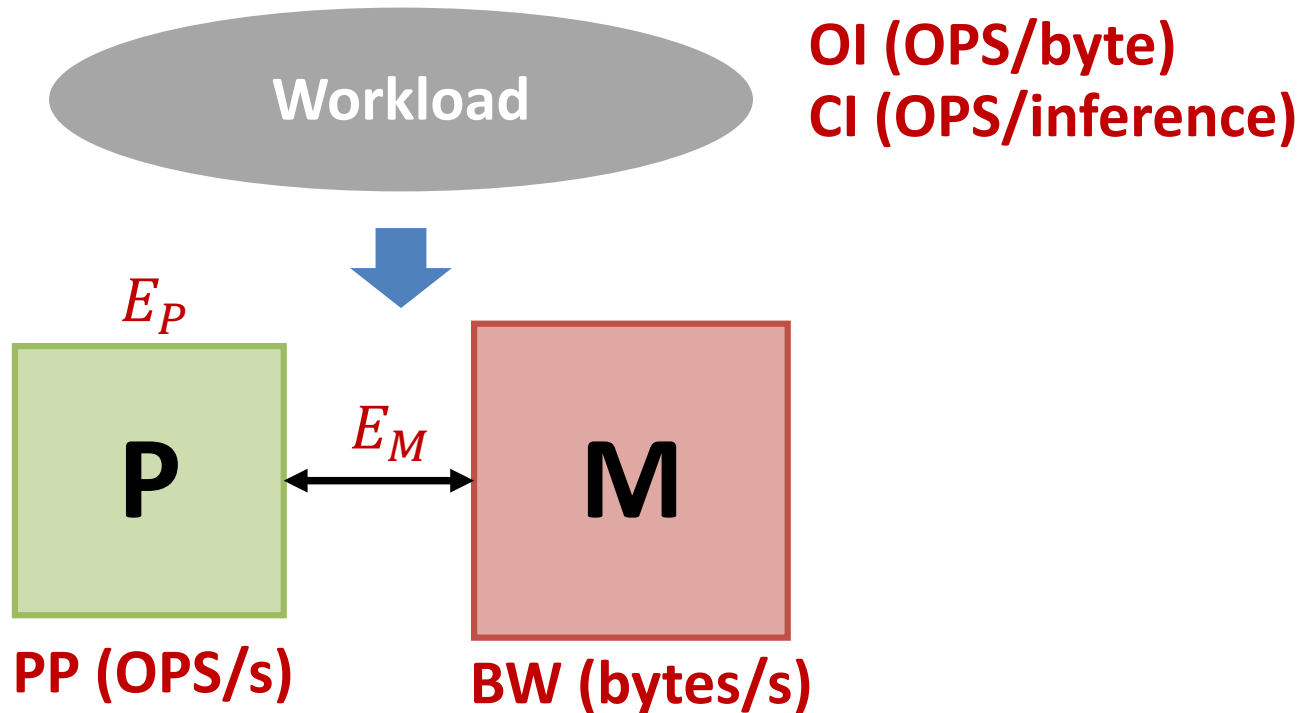
# DNN Architecture - Abstract View

**Workload**

**OI (OPS/byte)**
**CI (OPS/inference)**

$E_P$

**P**

$E_M$

**M**

**PP (OPS/s)**

**BW (bytes/s)**

- Workload – DNNs, CNNs, LSTMs,….

- Processor **P** – GPU, accelerators, many-core architectures,….

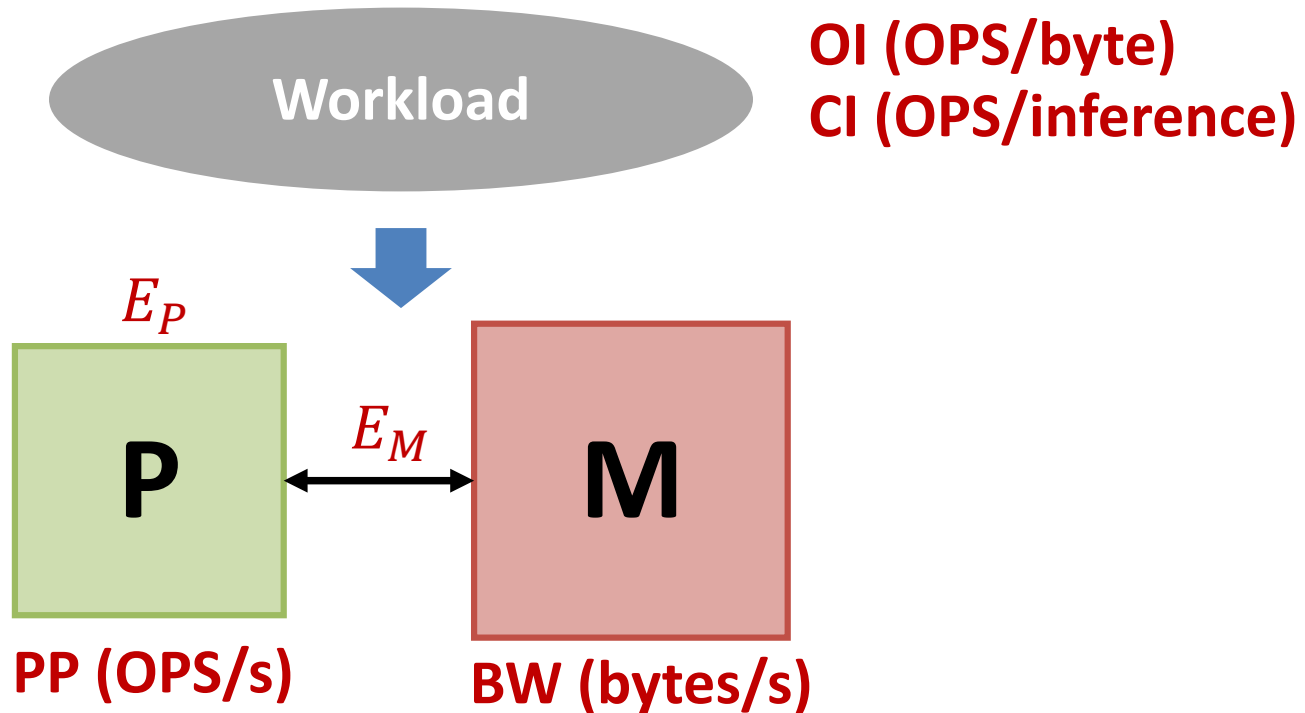- Memory **M** – register, SRAM, DRAM, flash,…..

# Workload Parameters



- operational intensity (OI): # of operations/byte (high is better) – reuse factor
- computational intensity (CI): # of operations per decision (small is better)
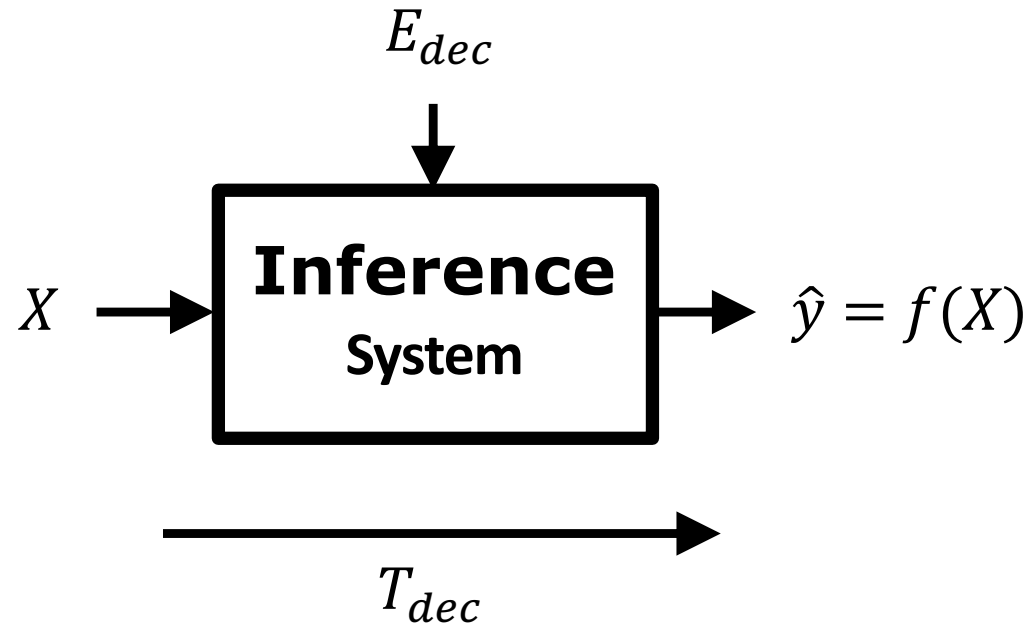
# Processor Parameters



OI (OPS/byte)
CI (OPS/inference)

PP (OPS/s)

BW (bytes/s)

- peak performance (PP): # of operations/s  (e.g., TOPS = $10^{12}$ OPS/s)
- arithmetic efficiency ($E_P$): energy per operation (e.g. $E_P = 0.2 pJ$/op)

# Memory Parameters



**OI (OPS/byte)**
**CI (OPS/inference)**

$E_P$
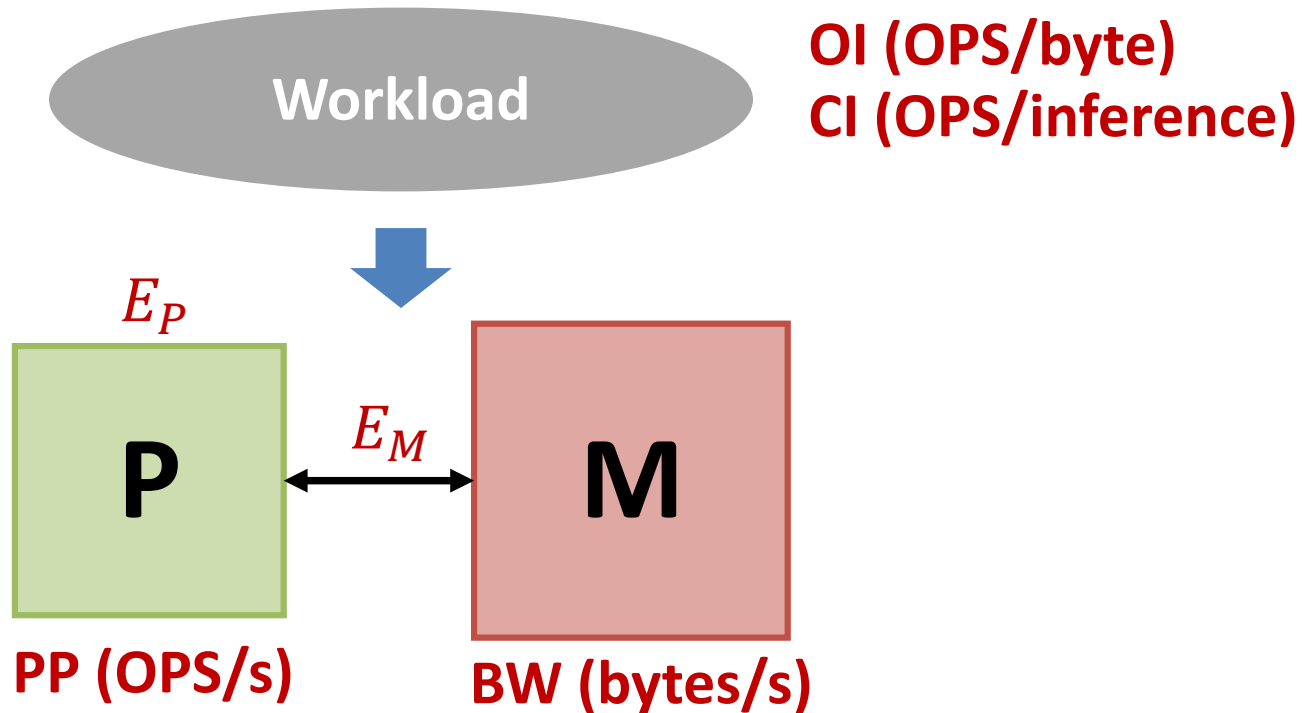
**P**

$E_M$

**M**

**PP (OPS/s)**

**BW (bytes/s)**

- read bandwidth (BW): # of bytes read/s (higher the better)

- energy/read ($E_M$): energy per byte (lower the better)

- $E_R = \dfrac{E_M}{E_P} = 100$ (SRAM); $= 500$ (DRAM); $= 1000$ (flash)
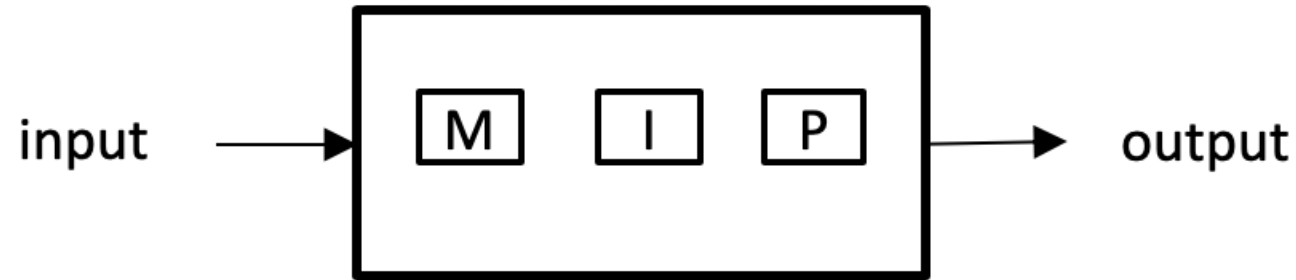
# System Level Metrics

$$E_{dec}$$

$$X \rightarrow \boxed{\begin{array}{c} \textbf{Inference} \\ \textbf{System} \end{array}} \rightarrow \hat{y} = f(X)$$

$$T_{dec}$$

- useful for system-level benchmarking
- energy/decision $(E_{dec})$
- decision throughput $(f_{dec} = \frac{1}{T_{dec}})$

# Component Level Metrics



- used for hardware benchmarking
- energy/OP ($E_P$); peak OP/s ($PP$); bytes read/s ($BW$); energy/byte ($E_M$)
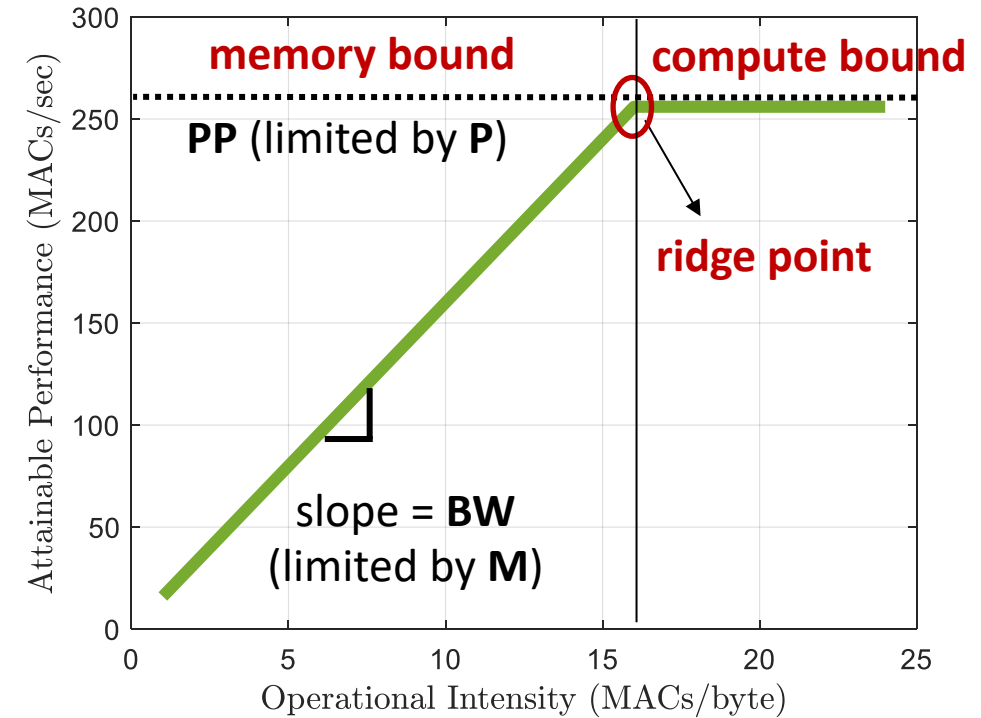
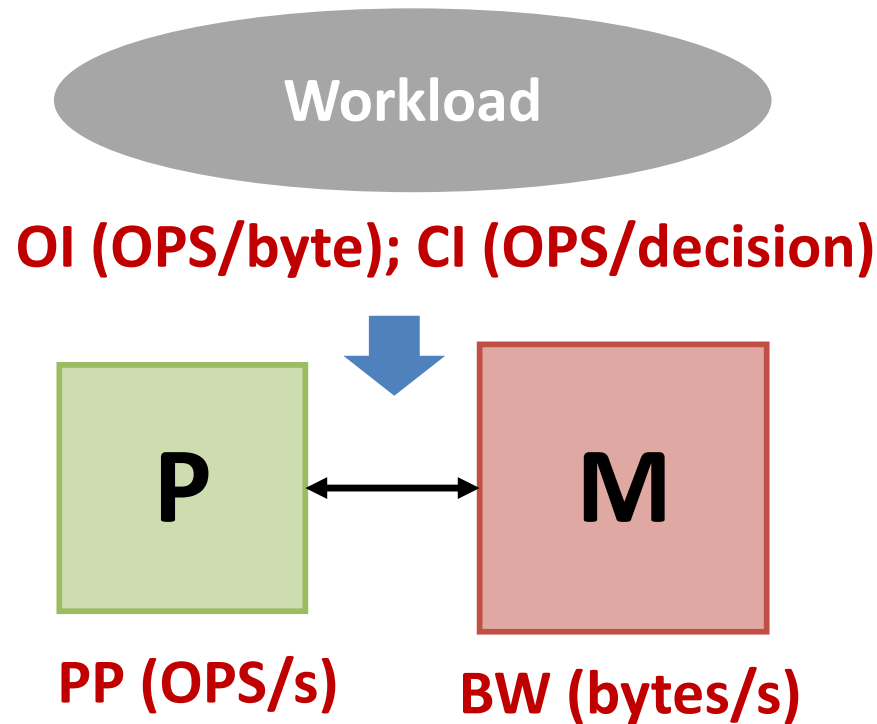# System vs. Component Level Metrics



- system-level metrics are a function of component-level metrics, workload requirements, and workload mapping onto hardware
- system-level metrics are the true metrics – hard to measure and report
- alternative – measure @ the component-level and estimate @ the system-level

# Roofline & Floorline Plots – Visual Performance & Efficiency Models

# Architectural-level Performance - Roofline Plot

Workload

OI (OPS/byte); CI (OPS/decision)

P

M

PP (OPS/s)

BW (bytes/s)



$$AP = \min(\mathrm{PP}, \mathrm{OI} \times \mathrm{BW})$$
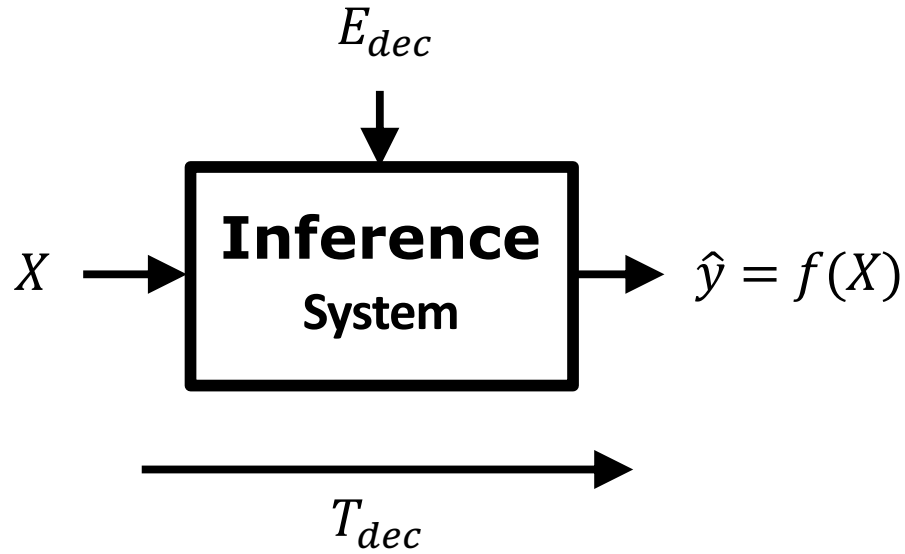
- $AP$: attainable performance of the architecture for the workload(s)
- Note: $AP \leq PP$
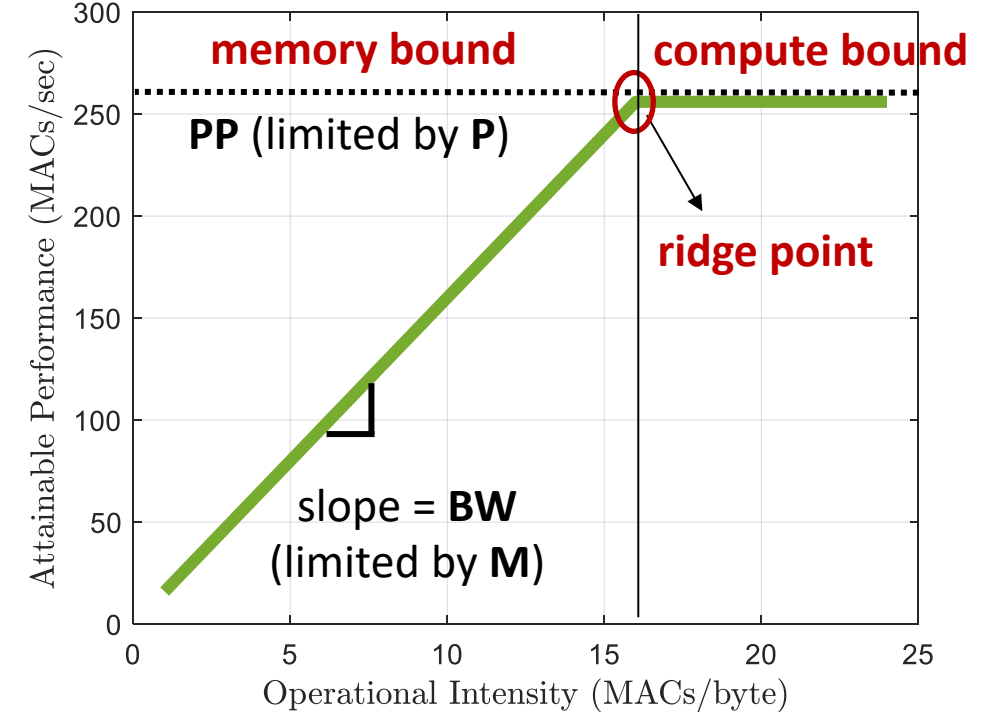
(Williams, Waterman, Patterson, Comm. of ACM '09)

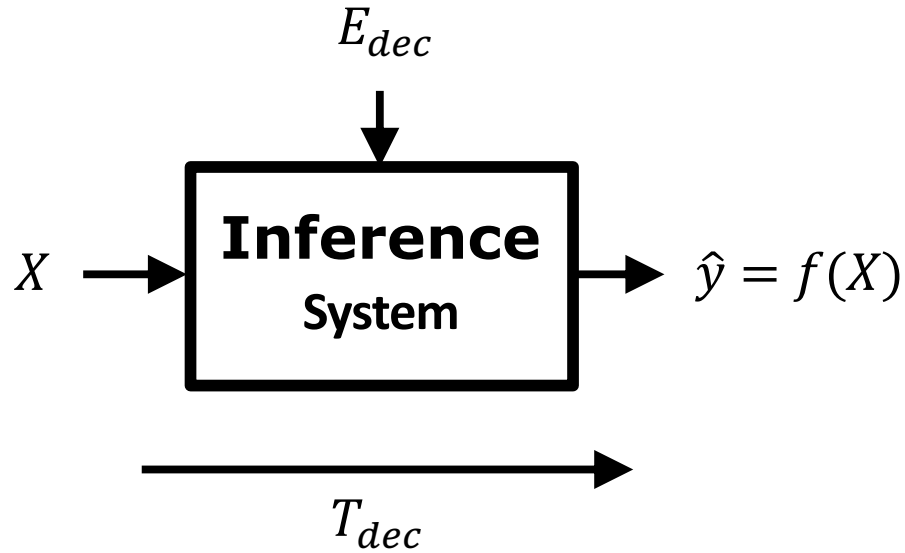# System-level Performance - Bytes/Decision - $M_{dec}$

**Roofline plot**





- Memory bandwidth (BW) $\neq$ bytes/decision ($M_{dec}$)
- Given a workload $OI$ and $CI$, bytes/decision :

$$M_{dec} = \frac{\text{bytes}}{\text{decision}} = \frac{\text{number of OPs/dec}}{\text{number of OPs/byte}} = \frac{CI}{OI}$$

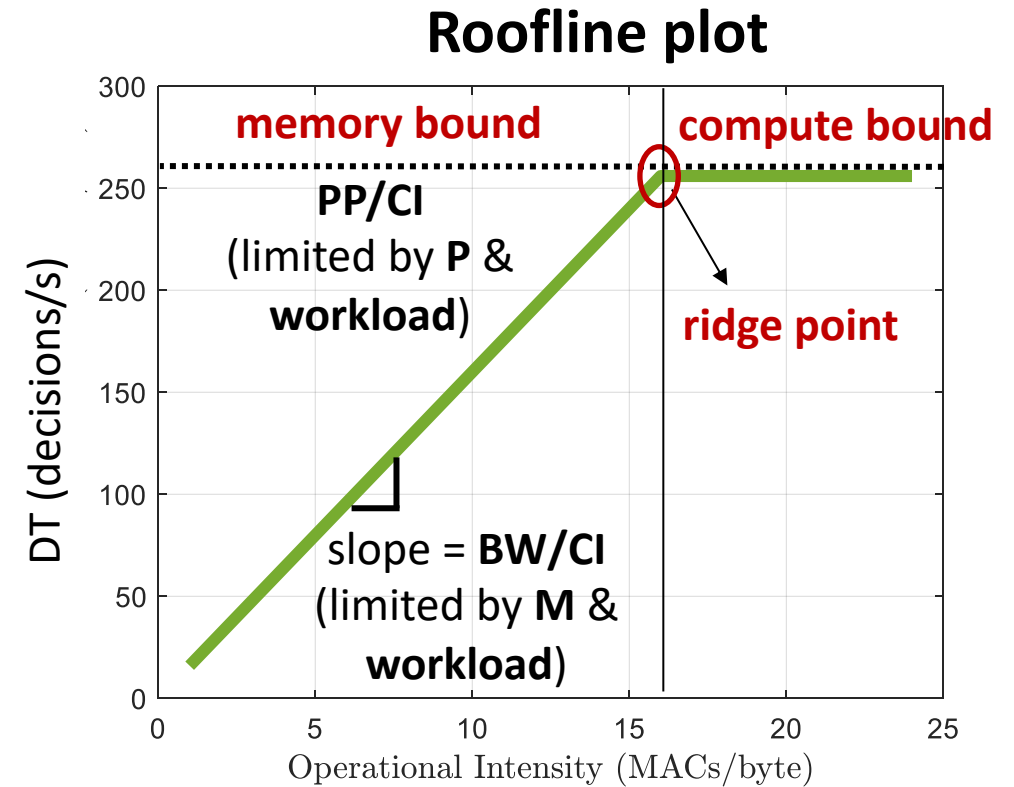# System-level Performance - Decision Throughput ($DT$)



- Attainable Performance (AP) $\neq$ Decision throughput (DT)
- Given a workload $OI$ and $CI$, decision throughput (dec/s):

$$DT = \frac{\text{decisions}}{s} = \frac{\text{achievable OPs/s}}{\# \text{ of OPs/dec}} = \frac{AP}{CI} = \min\left[\frac{PP}{CI}, OI \times \left(\frac{BW}{CI}\right)\right] = \min\left[\frac{PP}{CI}, \frac{BW}{M_{dec}}\right]$$
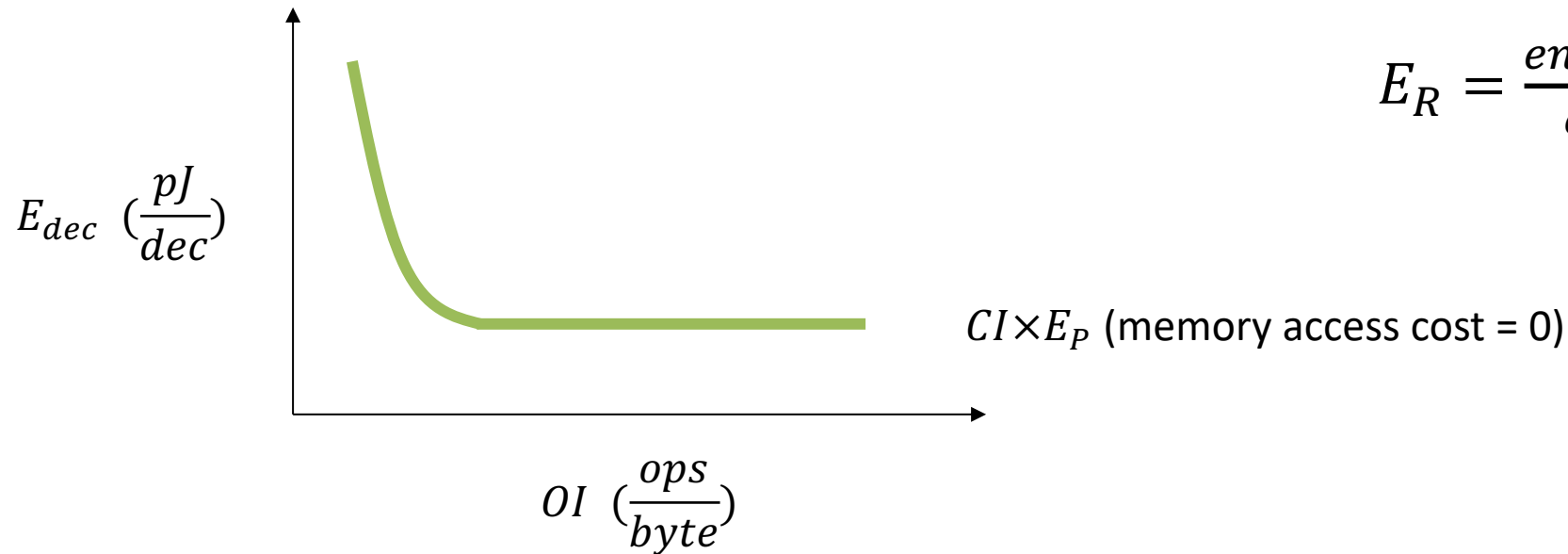
# System Energy Efficiency – Floorline Plot

(Shanbhag, Dbouk, Sakr, ECE 498NSU/ECE 598NSG 2019)

- Decision energy (pJ/dec):
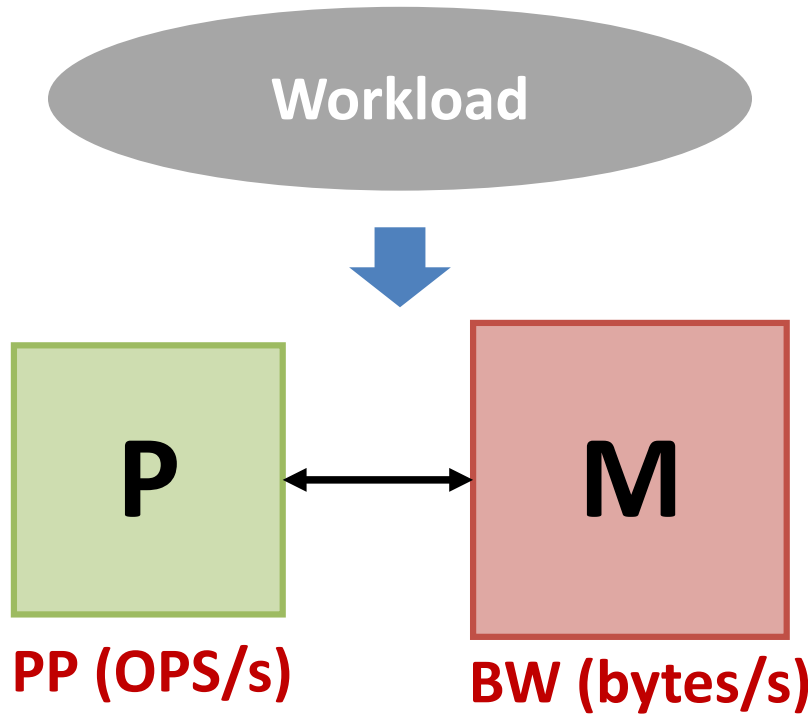
$$E_{dec} = CI{\times}E_P + M_{dec}{\times}E_M = CI{\times}E_P + \frac{CI}{OI}{\times}E_M = CI\left(E_P + \frac{E_M}{OI}\right) = CI{\times}E_P\left(1 + \frac{E_R}{OI}\right)$$

**Floorline plot**

$$E_R = \frac{energy/access}{energy/OP} \sim 10^2 \text{ (SRAM)}$$



$CI{\times}E_P$ (memory access cost = 0)

$E_{dec}$ $(\frac{pJ}{dec})$

$OI$ $(\frac{ops}{byte})$

# Maximizing Inference Throughput

**Workload**

**OI (OPS/byte)**
**CI (OPS/inference)**

**P** ↔ **M**

**PP (OPS/s)**     **BW (bytes/s)**
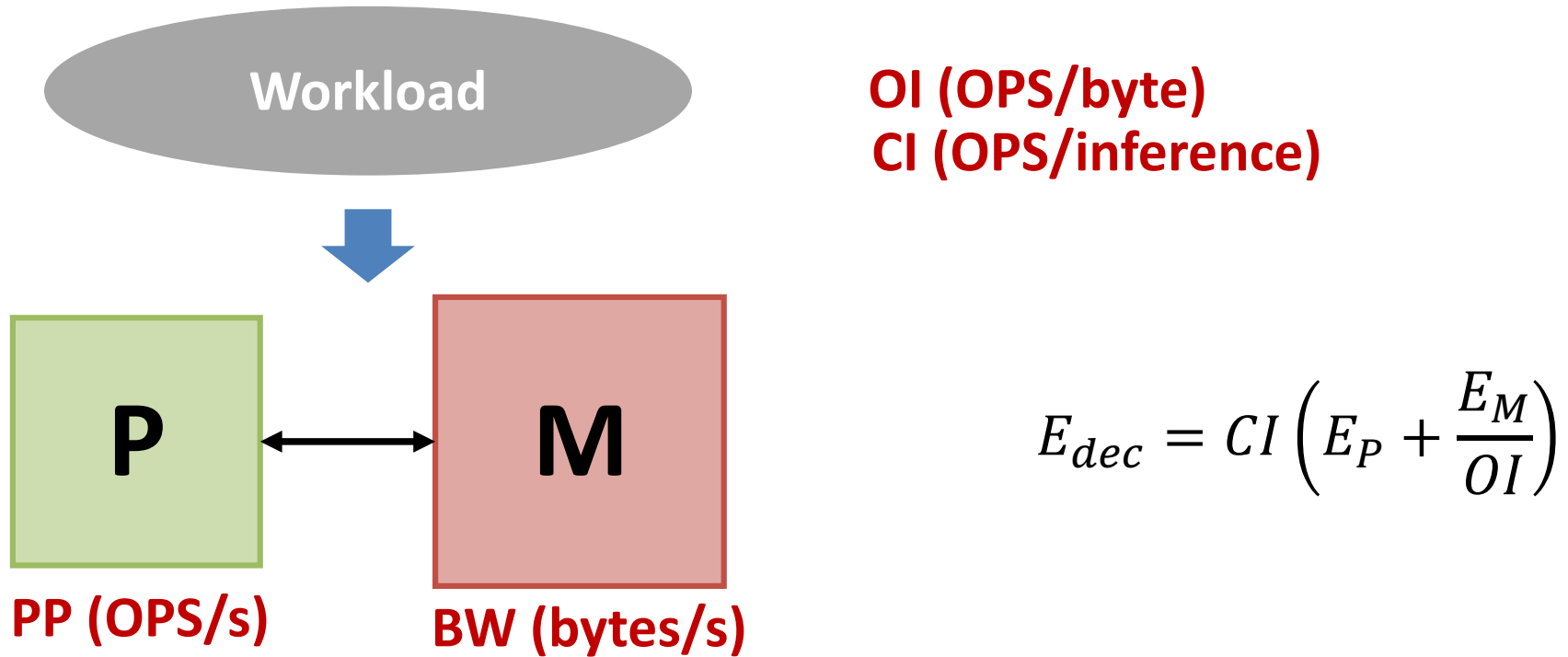
$$DT = \frac{AP}{CI} = \frac{\min(PP, OI \times BW)}{CI}$$

- Increase OI: data reuse & HW mapping → helps reduce decision energy
- Decrease CI: low-complexity DNNs, fast algorithms (FFT) → helps reduce decision energy
- Increase PP: systolic architectures & algorithm transforms
- Increase BW: new memory technologies + in-memory computing

# Maximizing Energy Efficiency

**Workload**

**OI (OPS/byte)**
**CI (OPS/inference)**

**P** &harr; **M**

**PP (OPS/s)**   **BW (bytes/s)**

$$E_{dec} = CI\left(E_P + \frac{E_M}{OI}\right)$$

- Increase OI: data reuse & HW mapping → helps increase decision throughput
- Decrease CI: low-complexity DNNs → helps increase decision throughput
- Decrease $E_P$: energy-efficient computing (e.g. near-$V_t$ computing), reduce precision
- Decrease $E_M$: new memory technologies + in-memory computing
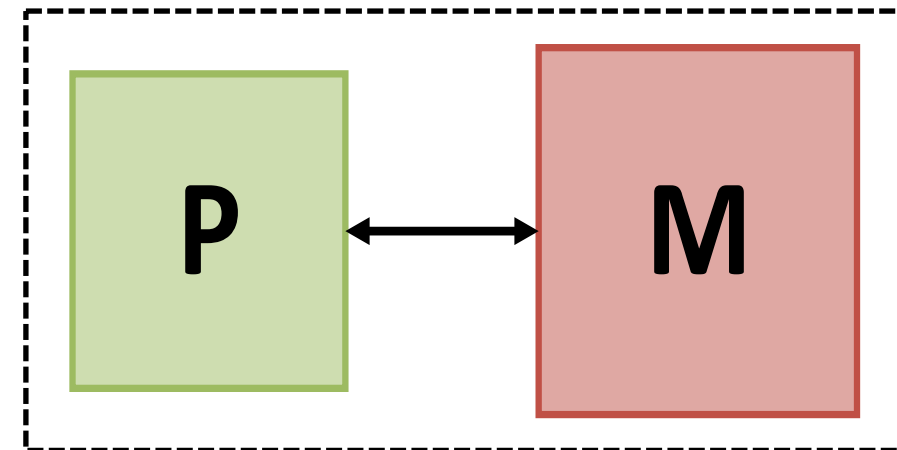
# Example

Given:

- an inference workload:
  - $CI = 1000$ OPs/dec
  - $OI = 10$ OPs/byte

- hardware platform with:
  - $PP = 256$ OPs/s
  - $BW = 16$ bytes/s
  - $E_P = 0.23$ pJ/OP
  - $E_M = 2.5$ pJ/byte

**Application**

Workload

**Hardware**
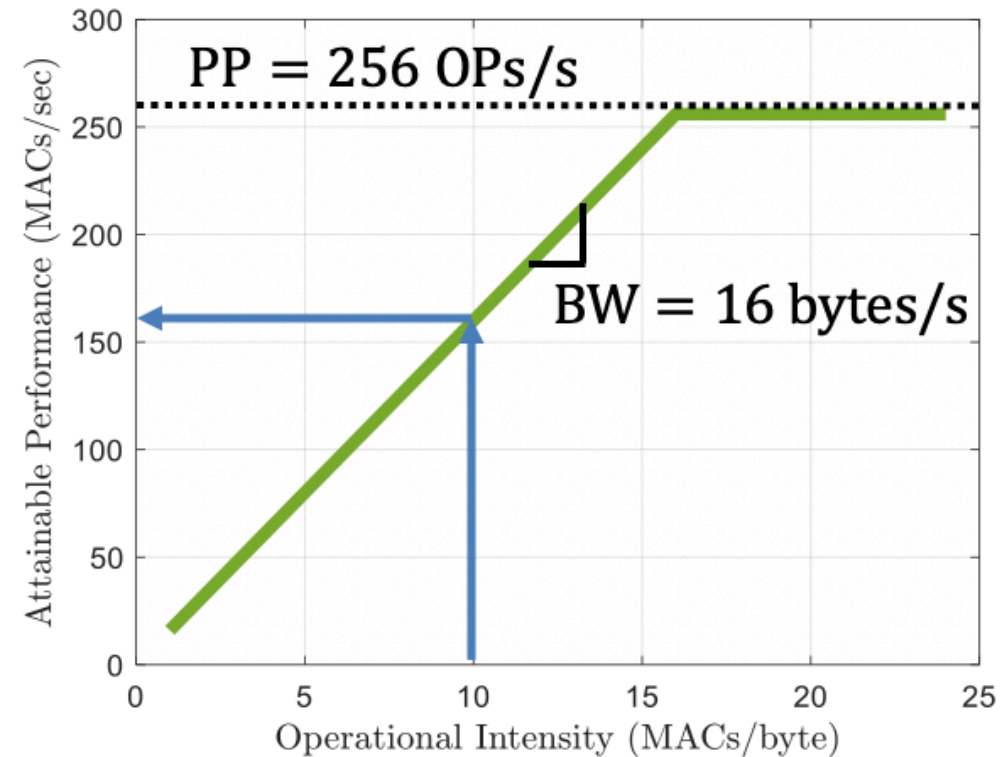
P ⟷ M

determine decision energy and latency

# Throughput - Roofline plot

- Attainable performance:

$$AP = \min(PP, OI \times BW) = 160 \text{ OPs/s}$$

- Decision throughput:

$$DT = \frac{AP}{CI} = \frac{160}{1000} = 0.16 \text{ dec/s}$$

# Energy Efficiency - Floorline plot
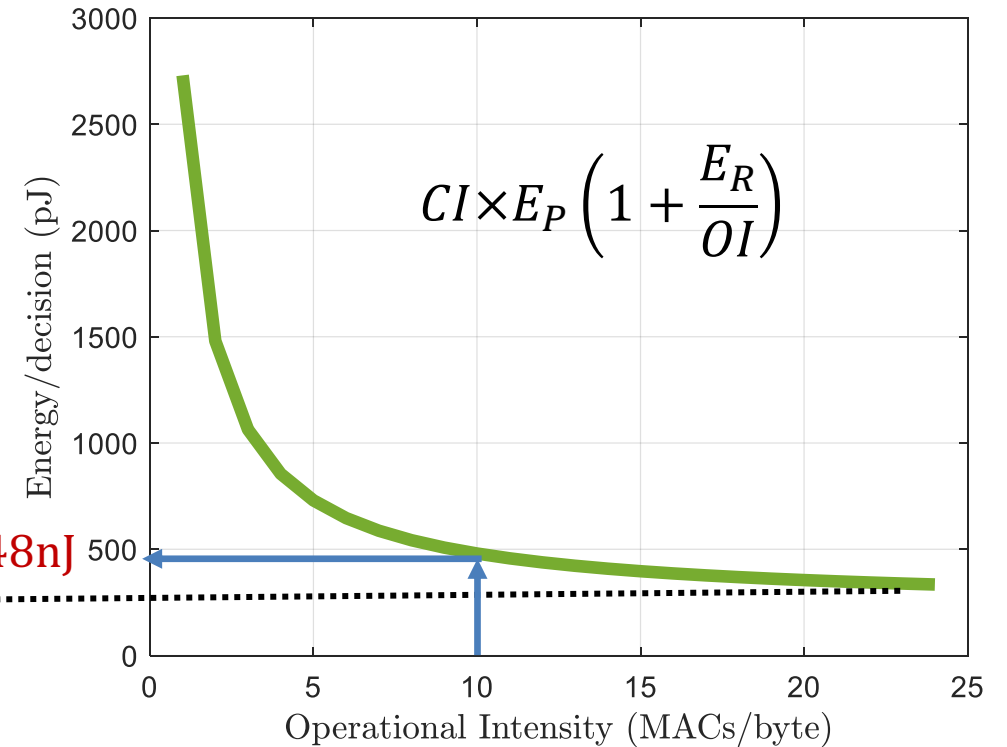
- Decision energy:

$$CI\left(E_P + \frac{E_M}{OI}\right) = 1000\left(0.23 + \frac{2.5}{10}\right)$$
$$= 0.48 \text{ nJ/dec}$$

- Floorline plot:

(lower bound) $CI \times E_P = 0.23$ nJ

$E_{dec} = 0.48$ nJ

$$E_R = \frac{E_M}{E_P} = 10.87$$

$$\text{E}_{\textbf{dec}}(\text{nJ}) = 0.23\left(1 + \frac{10.87}{\text{OI}}\right)$$
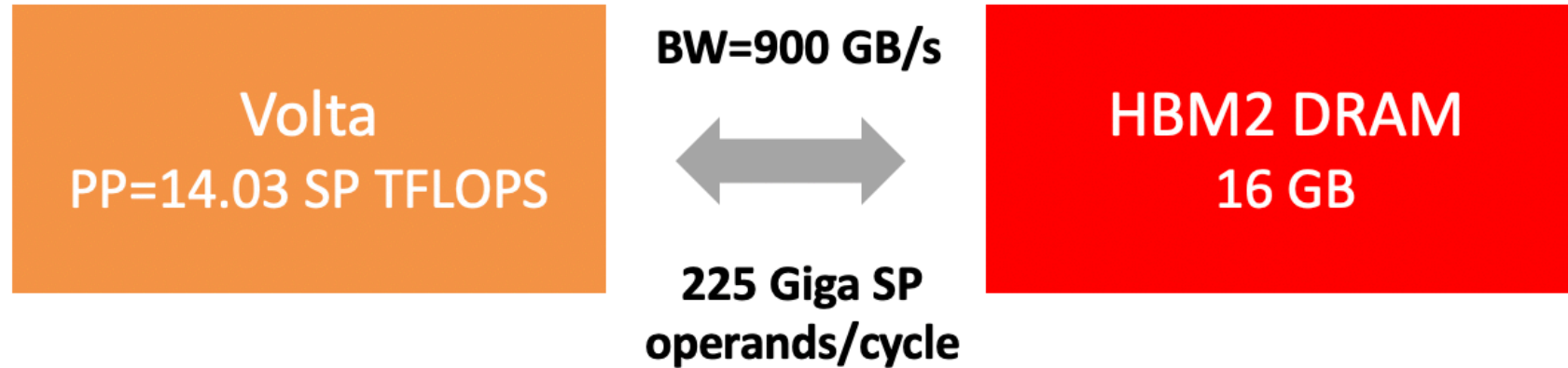
$$CI \times E_P\left(1 + \frac{E_R}{OI}\right)$$

# Component-level Energy-Delay Trade-offs

- Energy/OP ($E_P$) & Peak OP/s ($PP$) – low-power and high-throughput digital architectures and circuits (high-ceiling home)

- Bytes read/s ($BW$) & Energy/byte ($E_M$) – energy-efficient memory and interface design

- some of these are mature topics but keep evolving to meet the needs of machine learning applications

- energy consumption vs. delay (vs. robustness) trade-off

- energy vs. delay trade-off captured via **energy-delay product** $EDP$

- energy/delay vs. robustness trade-off is more subtle (future lecture)

# Roofline Plot – A Visual Performance Model

# Example – Memory Bound Scenario

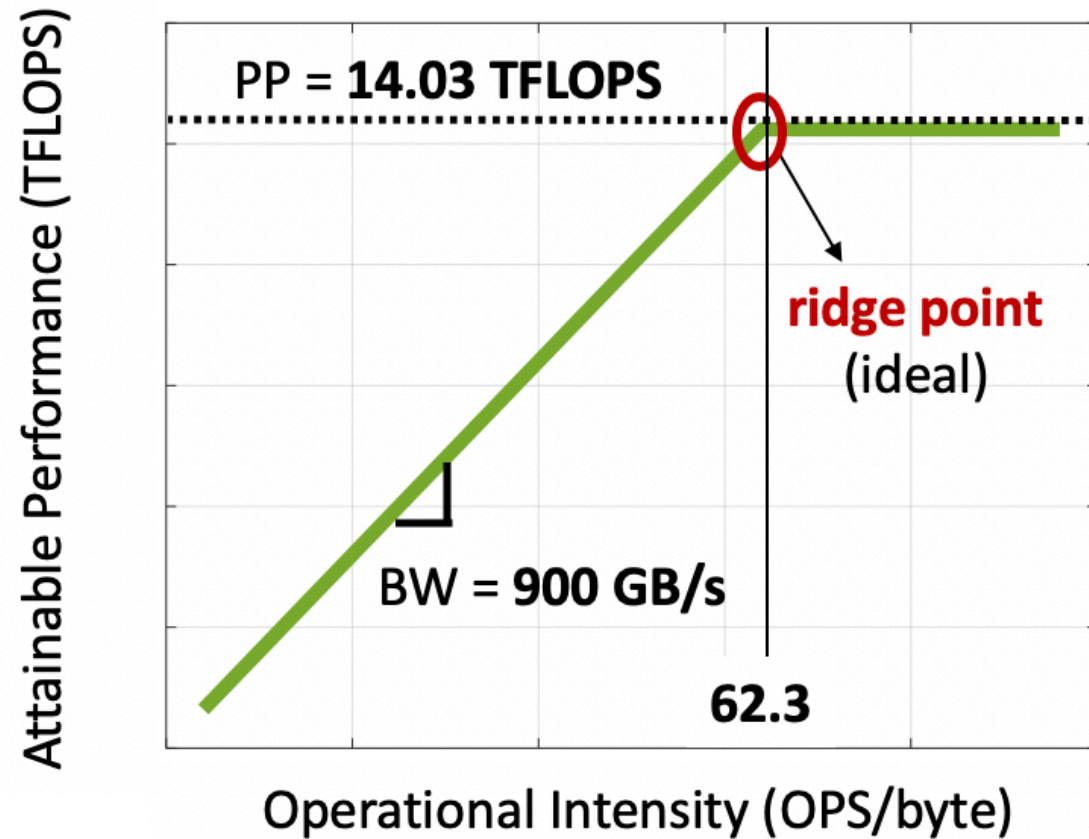[adapted from Hwu – IBM AI Compute Symposium, 2019]



Volta
PP=14.03 SP TFLOPS

BW=900 GB/s

225 Giga SP operands/cycle

HBM2 DRAM
16 GB

SP: single-precision (32-bit)

Each operand must be used **62.3 times** (reuse factor) once
fetched to achieve peak FLOPS rate (ridge point)
or
Sustain < **1.6%** of peak without data reuse

# Achievable Performance $(AP) = \min(\text{PP}, \text{OI} \times \text{BW})$

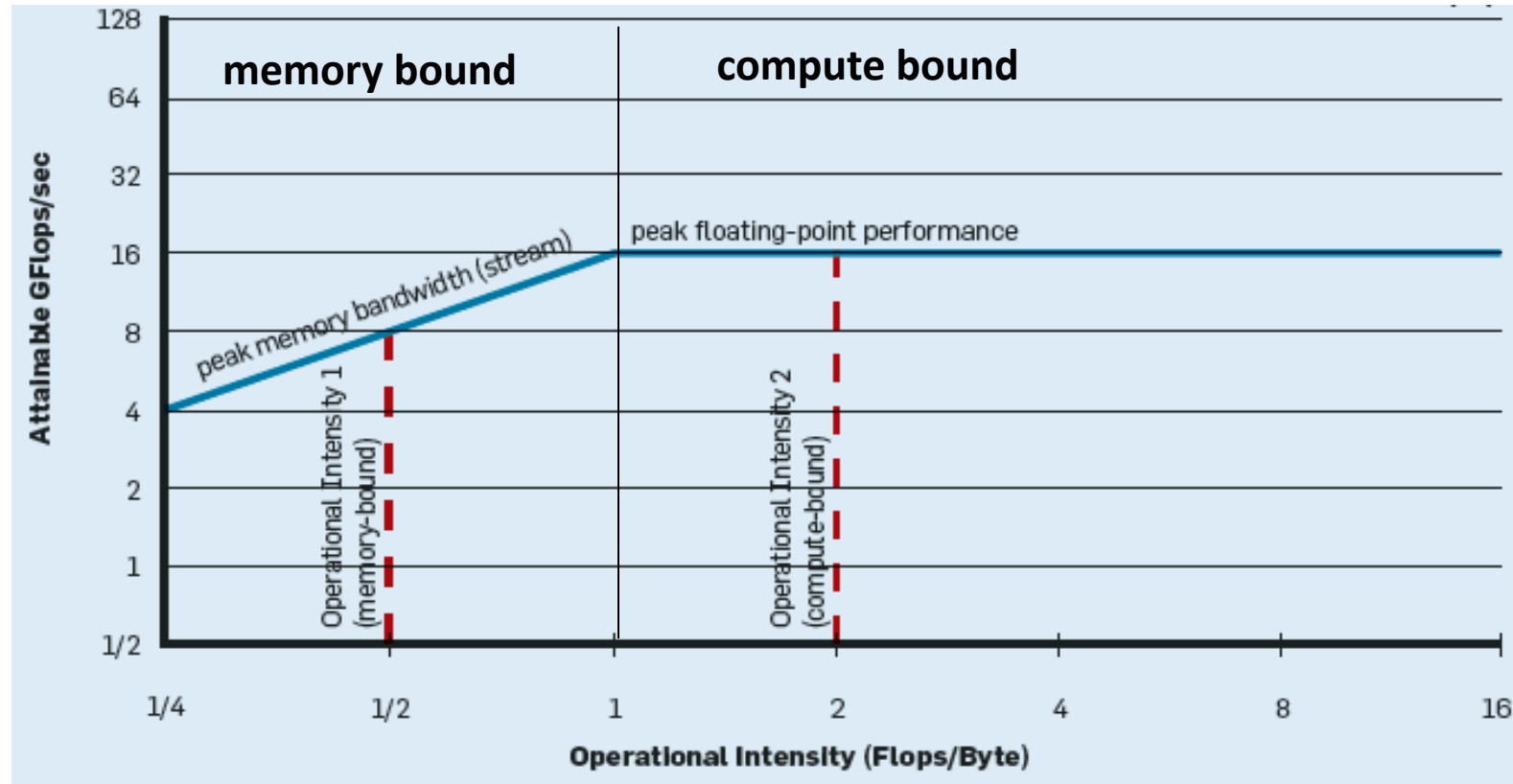# Roofline: An Insightful Visual Performance Model for Floating-Point Programs and Multicore Architectures*

Samuel Williams, Andrew Waterman, and David Patterson

Parallel Computing Laboratory, 565 Soda Hall, U.C. Berkeley, Berkeley, CA 94720-1776, 510-642-6587
samw, waterman, pattrsn@eecs.berkeley.edu

*"A model need not be perfect, just insightful."*

- Motivation: off-chip memory bandwidth will limit the performance

- Need a simple model relates processor performance to off-chip memory traffic similar to Amdhal's Law

- **Amdhal's Law**: performance gain of a parallel computer is limited by the serial portion of the workload
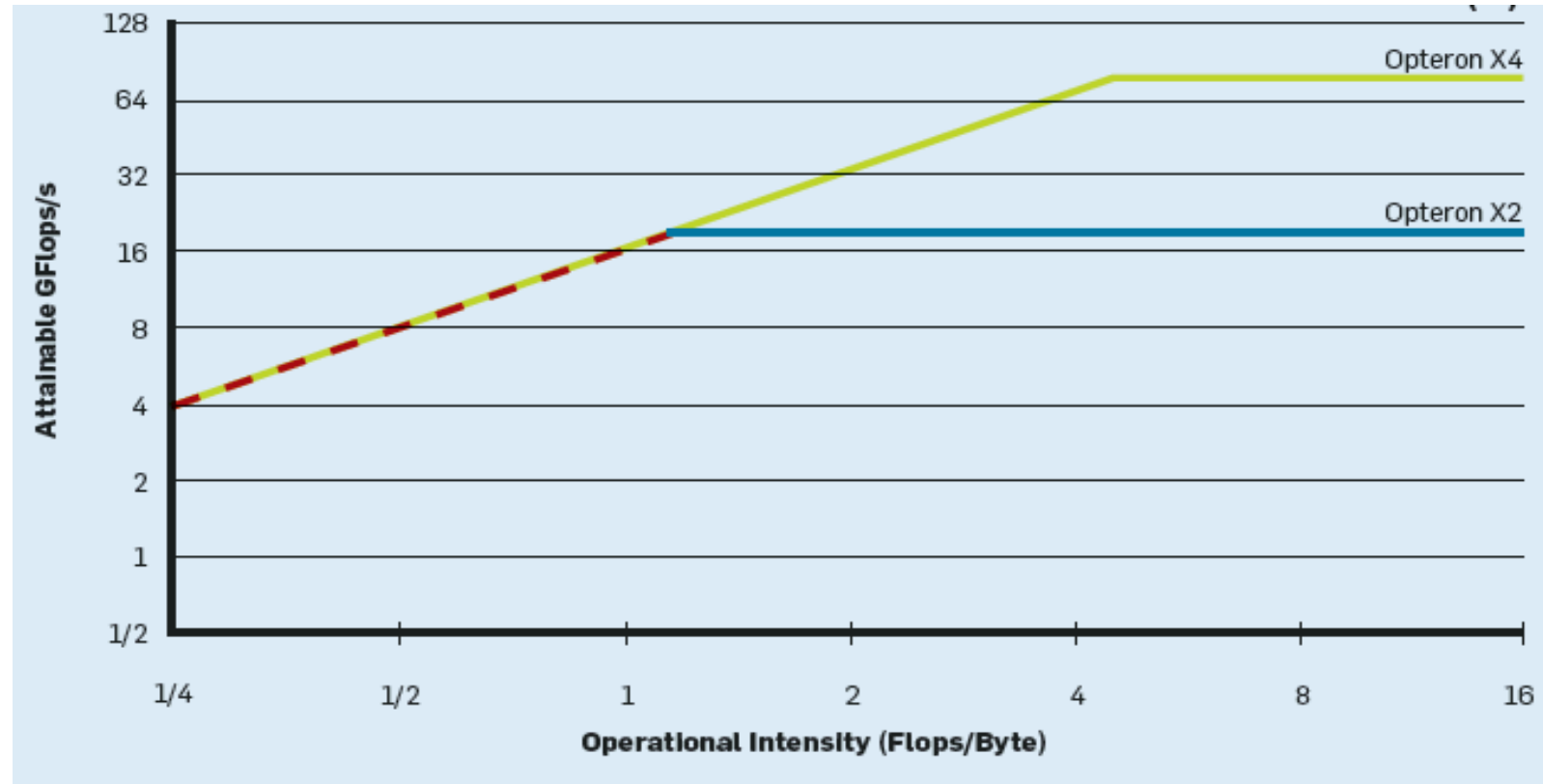
# Roofline for AMD Opteron X2



[Williams, ACM Communications, 2019]

- 2 core processor
- BW = 15 GB/s; PP = 17.6 GFlops/s

ILLINOIS
Electrical & Computer Engineering
COLLEGE OF ENGINEERING

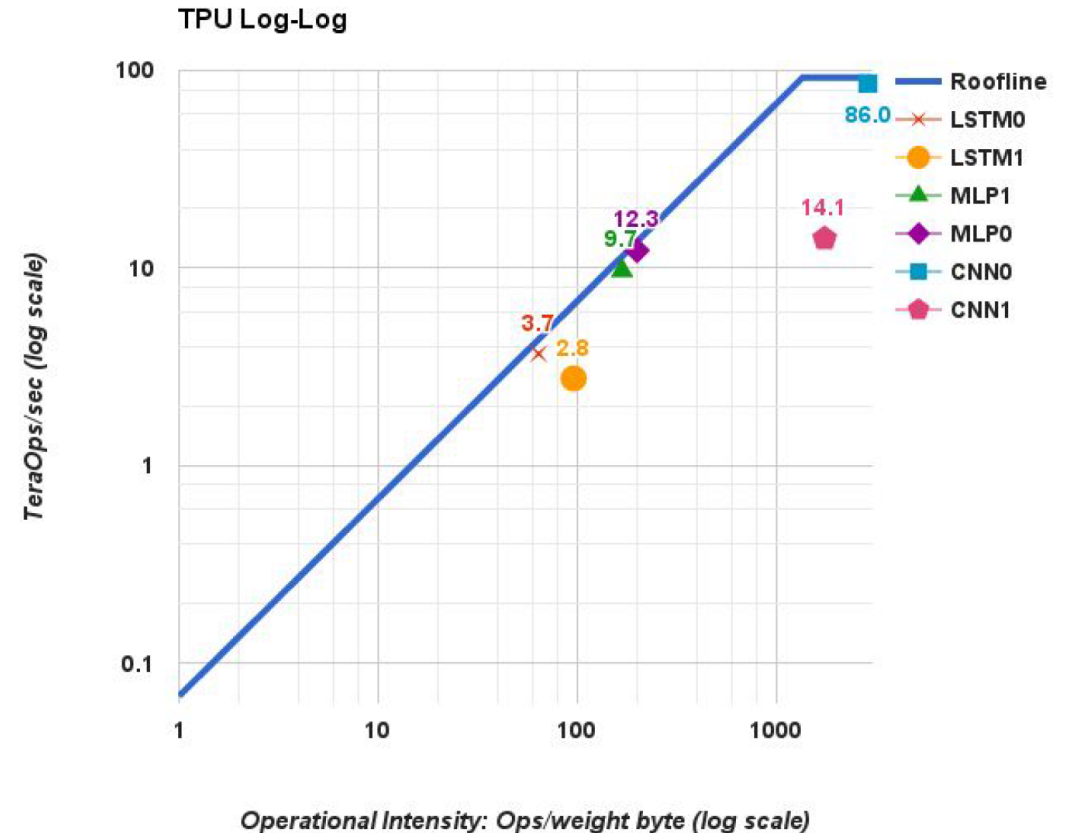2020 Notes for *Deep Learning in Hardware* - Naresh Shanbhag

# Opteron X2 vs. X4



- X2 = 2 core (2.2 GHz) X4 = 4 core (2.3 GHz)
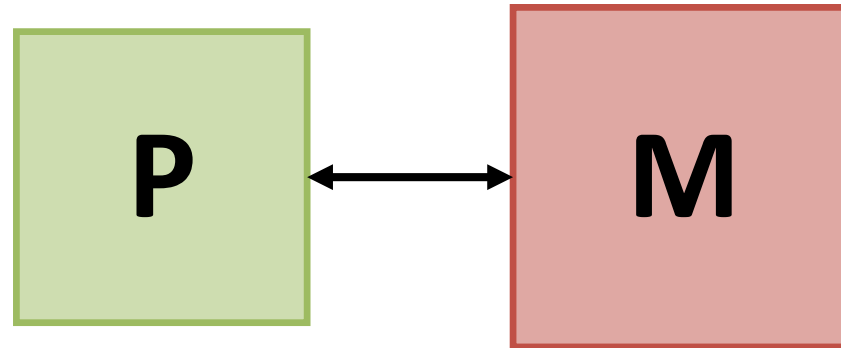- workloads with > 1 Flop/Byte OI needed for X4

# Roofline Plots – The Reality

- Roofline plot provides a ceiling on the achievable performance

- They ignore intricate data movement required by applications

- Different layer shapes might cause idle PEs => unutilized resources



TPU Log-Log

- Roofline
- LSTM0
- LSTM1
- MLP1
- MLP0
- CNN0
- CNN1

86.0
12.3
9.7
14.1
3.7
2.8

TeraOps/sec (log scale)

Operational Intensity: Ops/weight byte (log scale)

[TPU – ISCA'17]

# Dot Product Example



- Processor **P** with a peak performance (PP) of 256 operations/sec
- Memory **M** with a bandwidth (BW) of 16 bytes/sec
- Workload: perform a 1D "convolution" (correlator):

$$y_i = \sum_{j=1}^{M} w_j x_{i+j-1}$$

where $x \in \mathbb{R}^N, w \in \mathbb{R}^M \Rightarrow y \in \mathbb{R}^{N-M+1}$ (with zero-padding), and $M \in [N]$

assume 32b floating point operations; $N = 128$;

# Maximum AP

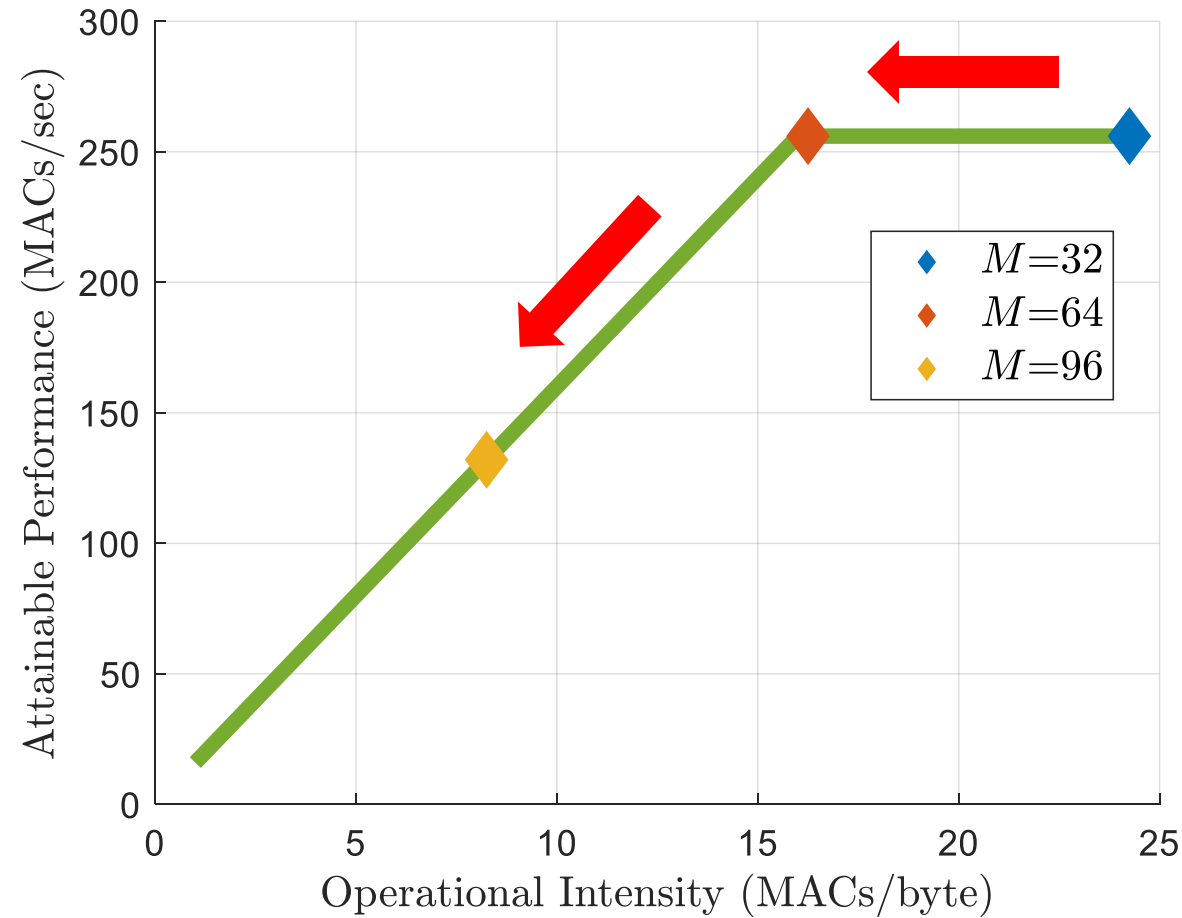- Assume we only have to fetch $w$ from **M**.

- OI for 1D convolution (DP):

**number of dot products**

$$OI = \frac{\text{number of MACs}}{\text{number of bytes read}} = \frac{(N - M + 1) \times M}{4 \times M}$$

**# MACs/DP**

$$OI = \frac{(N - M + 1)}{4}$$

**each weight value (32b) needs 4 memory reads**

- OI **reduces** with kernel size $M$
- OI **increases** with input data size $N$

# Varying Kernel Length ($M$)
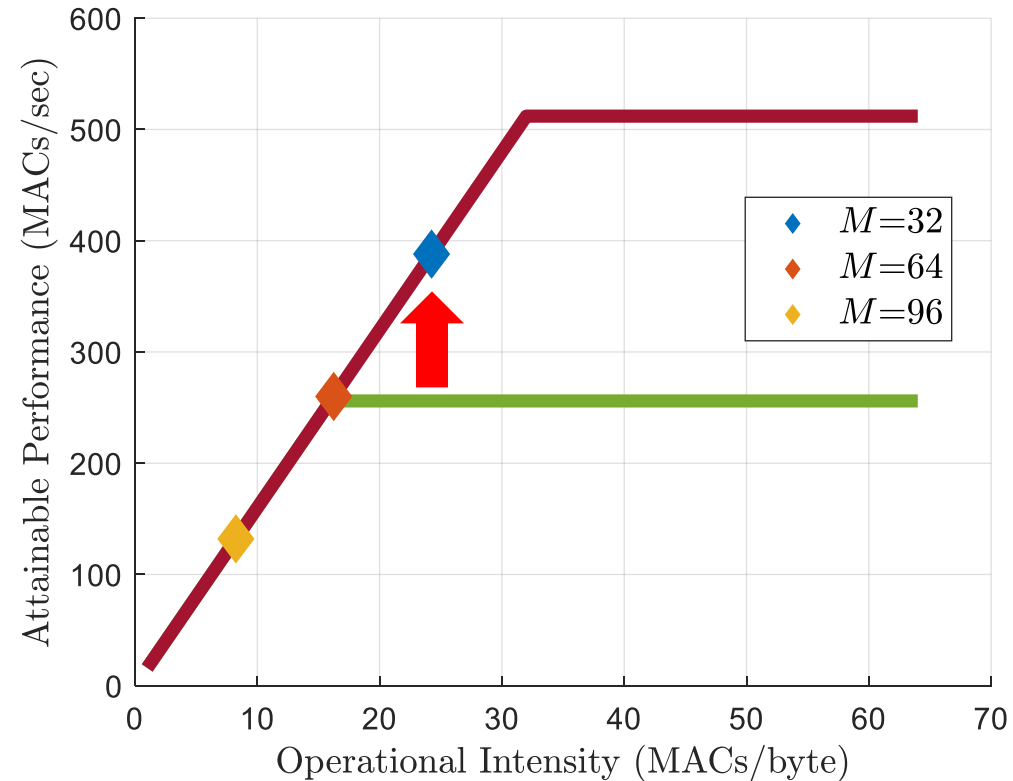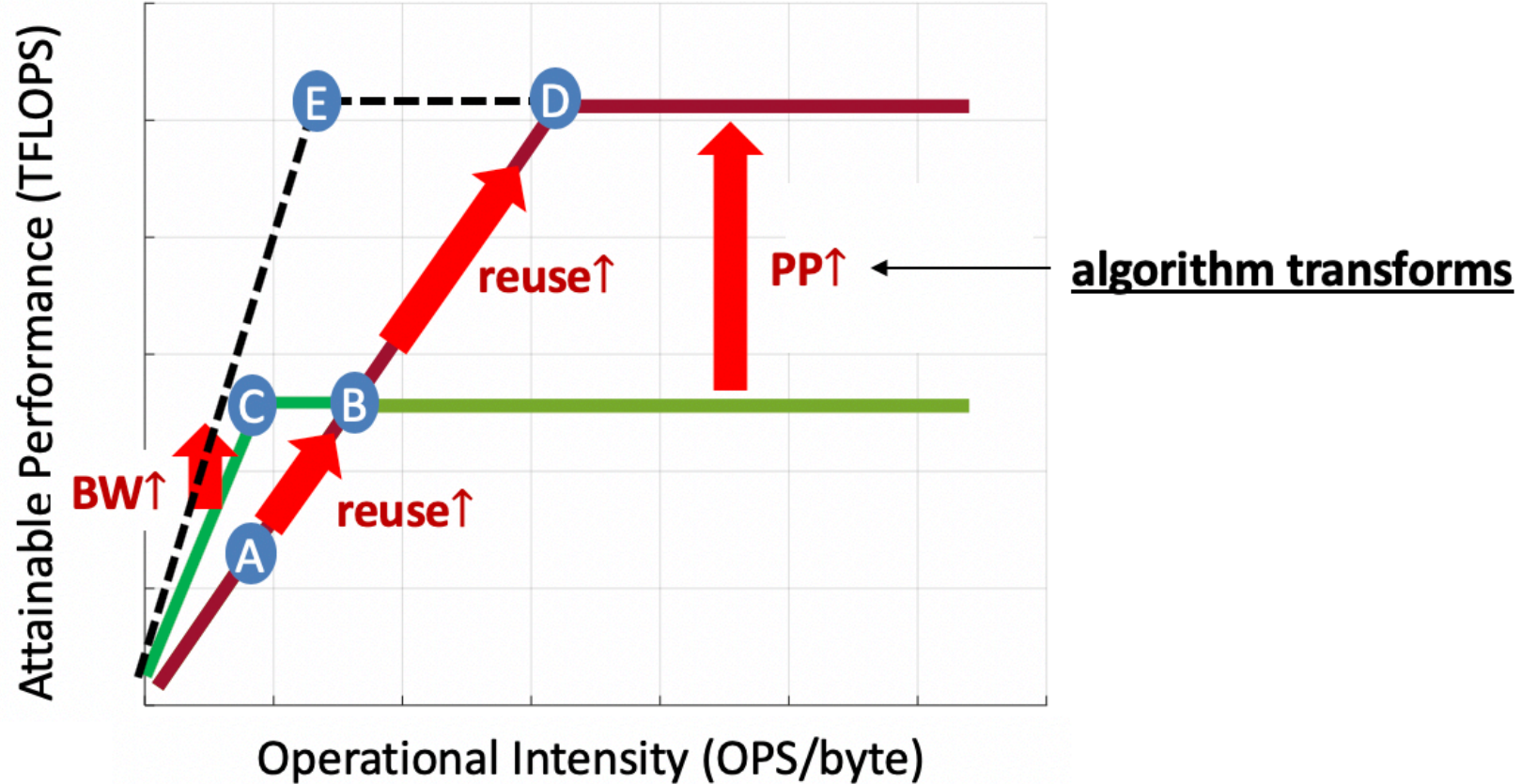
# Doubling Memory $BW$



- $M = 32$: $AP$ is *compute bound* → no benefit from doubling BW
- $M = 96$; $AP$ is *memory bound* → benefits from doubling BW

# Doubling Peak Performance ($PP$)



- $M = 32$: $AP$ increases and it becomes *memory bound*

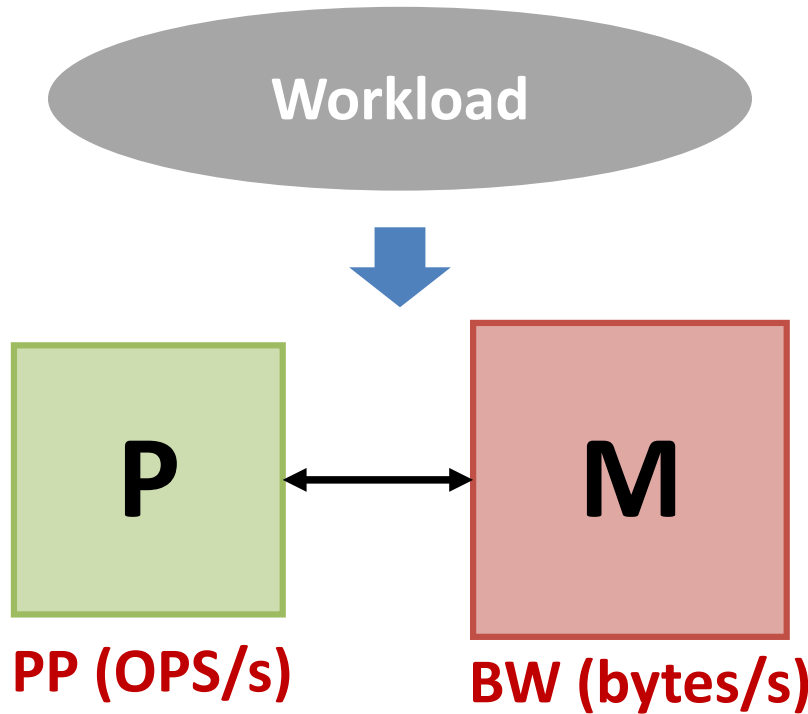- $M = 96$: $AP$ is already *memory bound,* so no change

# Methods to Improve AP



- A → B: same architecture - **enhance reuse**
- A → C: **same reuse** & processor – enhance BW
- B → D: same interface – **enhance reuse** and PP
- A → E: tunable - **enhance reuse**, BW and PP

# Data Reuse

**Workload**

**OI (OPS/byte)**
**CI (OPS/inference)**

**P**

**M**

**PP (OPS/s)**

**BW (bytes/s)**

$$DT = \frac{AP}{CI} = \min\left[\frac{PP}{CI}, OI \times \left(\frac{BW}{CI}\right)\right]$$

$$E_{dec} = CI\left(E_P + \frac{E_M}{OI}\right)$$

- **Increase OI: data reuse** & HW mapping
- Decrease CI: low-complexity DNNs
- Increase PP: systolic architectures & algorithm transforms
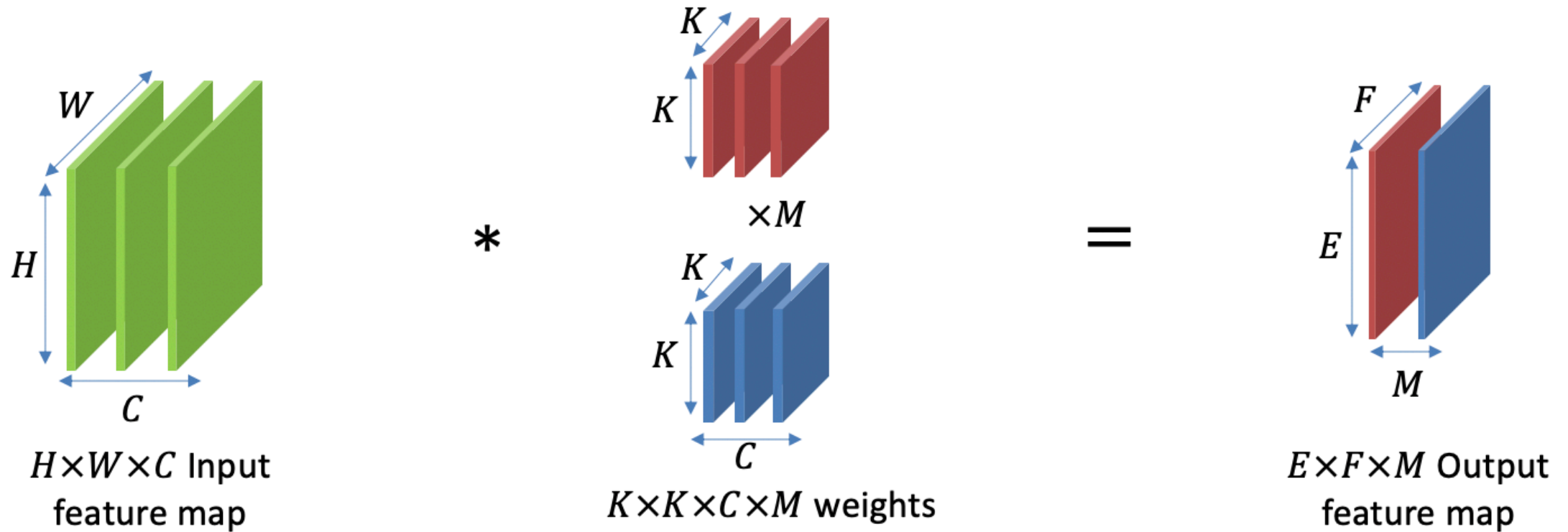- Increase BW: new memory technologies + in-memory computing

# Terminology

- **Data reuse**: # of MACs per data value **read** (i.e., MACs/data).
- For every MAC operation involved in a convolution operation, there are 3 data types:

$$\text{output partial sum} \longleftarrow \quad z \leftarrow z + w \times x \quad \longrightarrow \text{input activation}$$

with **weight** pointing up to $w$.

- Hence, three types of data reuse: input, weight, and partial sum
- Partial sum reuse = DP length = # of partial sums needed to compute one DP

# 2D Convolution



$H \times W \times C$ Input feature map

$*$

$\times M$

$K \times K \times C \times M$ weights

$=$

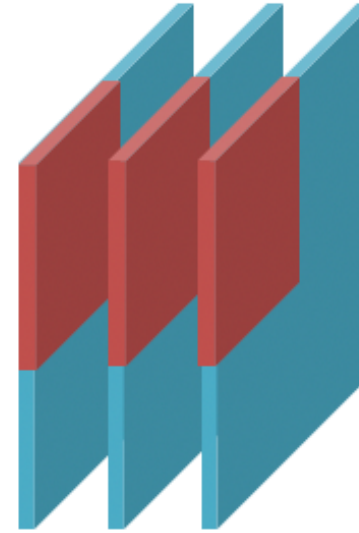$E \times F \times M$ Output feature map

# Input Reuse



- Each **input pixel** is **reused $\sim K \times K$ times** across the same weight filter
- Each **input pixel** is **reused $M$ times** due to $M$ filters processing the same input feature
- Total input reuse: $K^2 M$
- e.g., Layer 1 of VGG-16 on ImageNet $- 3 \times 3 \times 64 = 576$ MACs/byte = OI

# Weight Reuse



- Each **weight pixel** is **reused $E{\times}F$ times** across the same input feature map
- Each **weight pixel is being used $N$ times** when multiple $(N > 1)$ input batches are processed
- Total weight reuse: $E{\times}F{\times}N$
- e.g., Layer 1 of VGG-16 on ImageNet – $224{\times}224{\times}1 = 50176$ MACs/byte = OI

# Course Web Page

https://courses.grainger.illinois.edu/ece598nsg/fa2020/
https://courses.grainger.illinois.edu/ece498nsu/fa2020/

http://shanbhag.ece.uiuc.edu

ILLINOIS
Electrical & Computer Engineering
COLLEGE OF ENGINEERING