# UNIVERSITY OF TORONTO

# Energy-Efficient Machine Learning Acceleration

**Jorge Albericio, Patrick Judd, Alberto Delmas, Mostafa Mahmoud, Milos Nikolic, Zissis Poulos, Sayeh Sharify, Kevin Siu, Dylan Stuart**

**Isak Edo, Omar Awad, Ali Hadi Zadeh**

*Tayler Hetherington, Tor Aamodt, Natalie Enright Jerger*
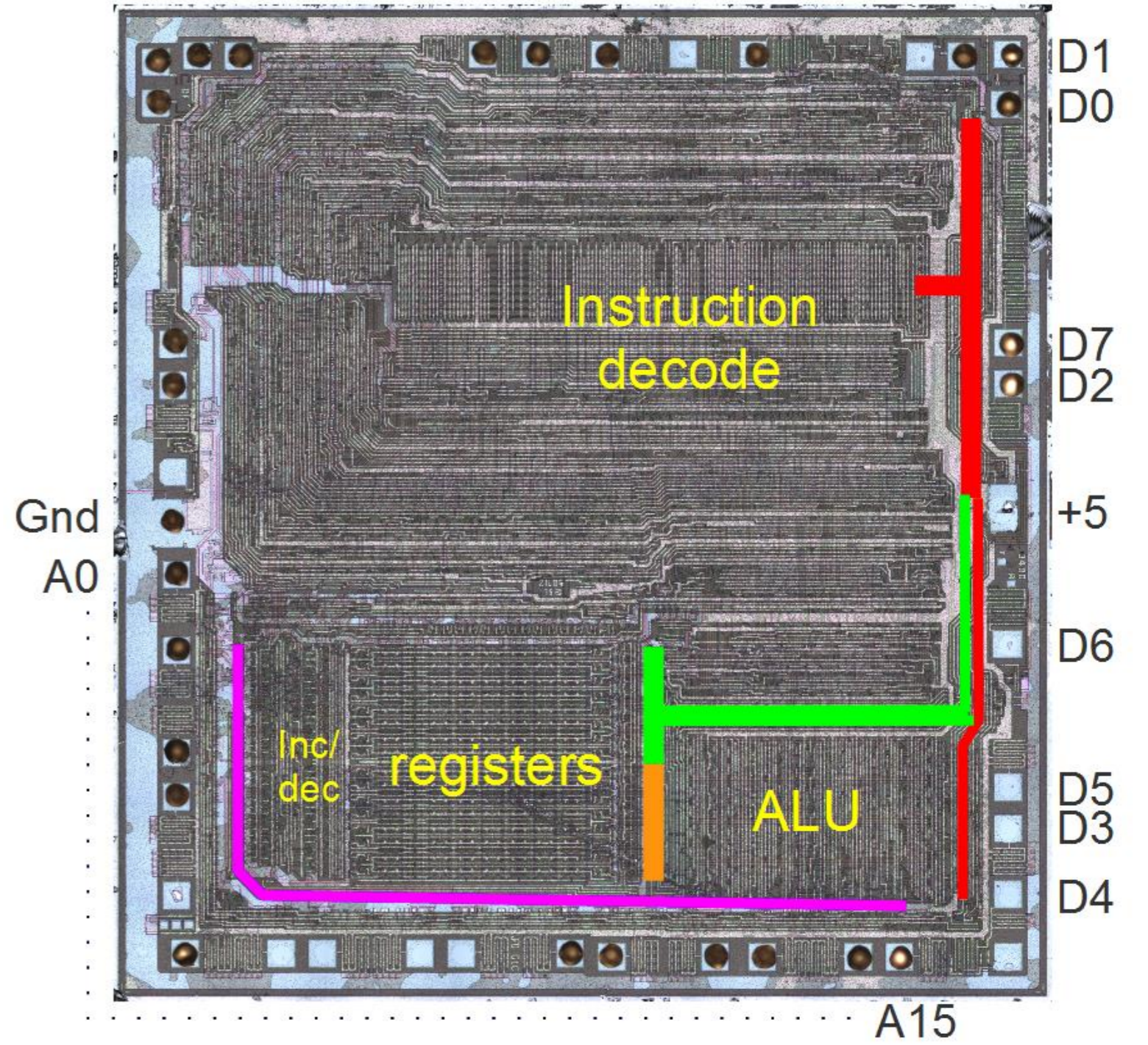
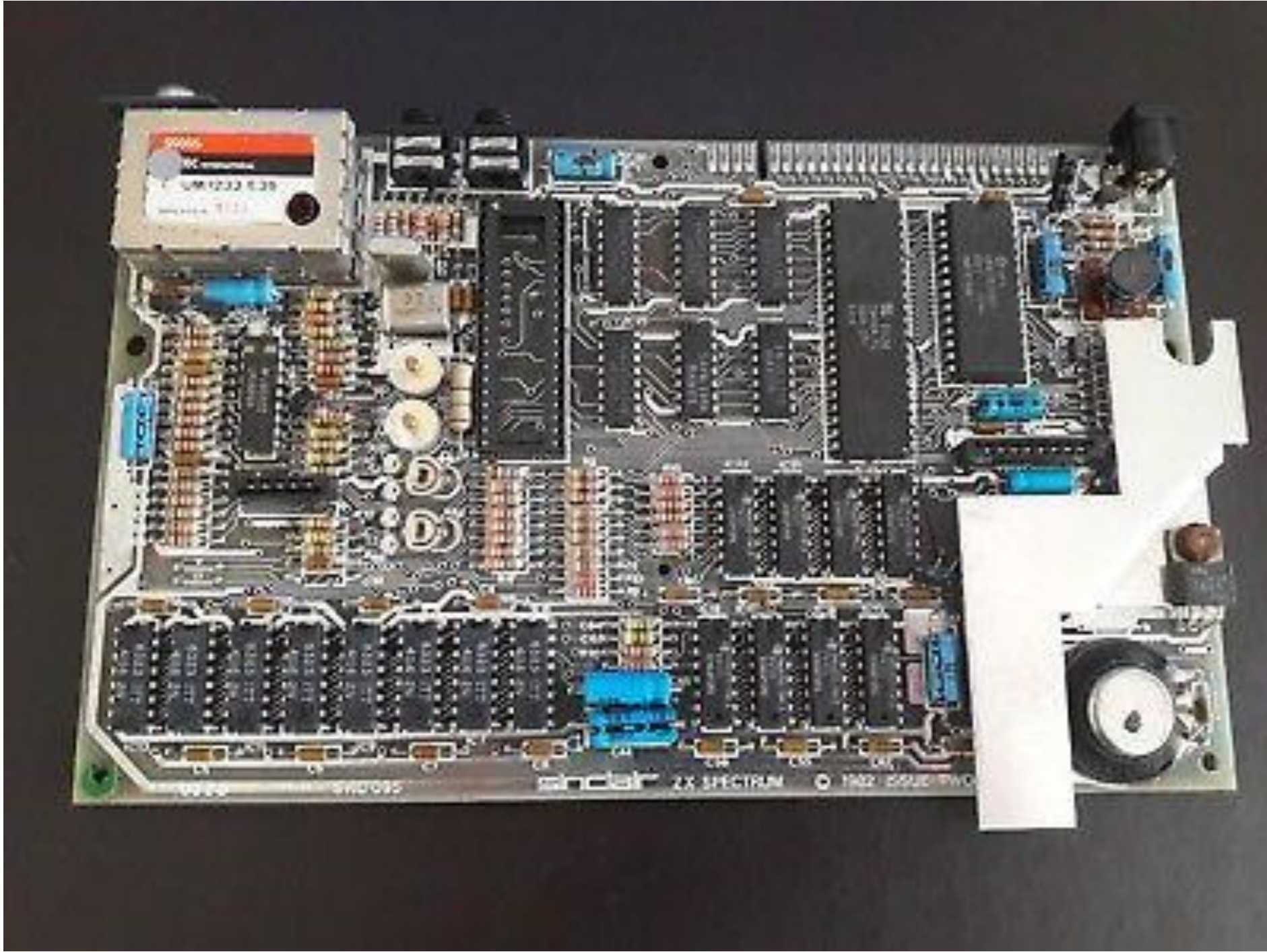**Andreas Moshovos**

**moshovos@ece.utoronto.ca**

```
10>LET a=10
20 PRINT a
```

**Computing Hardware**
We build tools
Used by "everyone" for "everything"
Science, medicine, commerce, …

# www.eecg.toronto.edu/~moshovos/CUDA08

www.eecg.toronto.edu/~moshovos/CUDA08/doku.php?id=project_presentations_reports_source_code

## ECE1742S: Programming Massively Parallel Multiprocessors Using CUDA

Trace: start → project_presentations_reports_source_code

## Convolutional Neural Networks for Object Classification

Alex Krizhevsky

Department of Computer Science, University of Toronto

akrizhevsky at gmail.com

**Abstract** I implemented a convolutional neural network with one layer of convolution. I tested it on the CIFAR-10 dataset, which consists of 6000 32×32 colour images in each of 10 classes. The convolutional net does well on the classification task and takes roughly 140x less time to train than a CPU implementation.

- Presentation: (pdf).
- Report: (pdf).
- Source code: (zip).

# Our Goal:

Hardware Acceleration of ML

## Why?

Hope to… Further Enable ML:

*Innovation*

*Applications*

# Computing Devices do…

**Data "Transformation"** → A + B

**Data Movement** → A = B

# Off-Chip

# On-Chip

**Memory**

**Compute Cores**

**On-Chip Memory**

## On- vs. Off-Chip
Energy: ~100x
Latency: ~50x

## Compute/Watt is the primary design constraint

## **Layers**



*LHC*
**maybe**

## Tons of **Out += A x W**
For other types of networks too

$$Out_0 \mathrel{+}= A_0 \times W_0$$

$$Out_0 \mathrel{+}= A_1 \times W_1$$

$$Out_0 \mathrel{+}= A_2 \times W_2$$

$$Out_0 \mathrel{+}= A_3 \times W_3$$

$$Out_0 \mathrel{+}= A_4 \times W_4$$

$$Out_0 \mathrel{+}= A_0 \times W_0$$

$$Out_0 \mathrel{+}= A_1 \times W_1$$

## Lots of Parallelism

$$Out_0 \mathrel{+}= A_3 \times W_3$$

$$Out_0 \mathrel{+}= A_4 \times W_4$$

$Out_0 += A_0 \times W_0$ $Out_1 += A_0 \times W_0$ $Out_0 += A_0 \times W_0$ $Out_1 += A_0 \times W_0$ $Out_0 += A_0 \times W_0$ $Out_1 += A_0 \times W_0$ $Out_0 += A_0 \times W_0$ $Out_1 += A_0 \times$

$Out_0 += A_1 \times W_1$ $Out_1 += A_1 \times W_1$ $Out_0 += A_1 \times W_1$ $Out_1 += A_1 \times W_1$ $Out_0 += A_1 \times W_1$ $Out_1 += A_1 \times W_1$ $Out_0 += A_1 \times W_1$ $Out_1 += A_1 \times$

$Out_0 += A_2 \times W_2$ $Out_1 += A_2 \times W_2$ $Out_0 += A_2 \times W_2$ $Out_1 += A_2 \times W_2$ $Out_0 += A_2 \times W_2$ $Out_1 += A_2 \times W_2$ $Out_0 += A_2 \times W_2$ $Out_1 += A_2 \times$

$Out_0 += A_3 \times W_3$ $Out_1 += A_3 \times W_3$ $Out_0 += A_3 \times W_3$ $Out_1 += A_3 \times W_3$ $Out_0 += A_3 \times W_3$ $Out_1 += A_3 \times W_3$ $Out_0 += A_3 \times W_3$ $Out_1 += A_3 \times$

$Out_0 += A_4 \times W_4$ $Out_1 += A_4 \times W_4$ $Out_0 += A_4 \times W_4$ $Out_1 += A_4 \times W_4$ $Out_0 += A_4 \times W_4$ $Out_1 += A_4 \times W_4$ $Out_0 += A_4 \times W_4$ $Out_1 += A_4 \times$

$Out_0 += A_0 \times W_0$ $Out_1 += A_0 \times W_0$ $Out_0 += A_0 \times W_0$ $Out_1 += A_0 \times W_0$ $Out_0 += A_0 \times W_0$ $Out_1 += A_0 \times W_0$ $Out_0 += A_0 \times W_0$ $Out_1 += A_0 \times$

$Out_0 += A_1 \times W_1$ $Out_1 += A_1 \times W_1$ $Out_0 += A_1 \times W_1$ $Out_1 += A_1 \times W_1$ $Out_0 += A_1 \times W_1$ $Out_1 += A_1 \times W_1$ $Out_0 += A_1 \times W_1$ $Out_1 += A_1 \times$

$Out_0 += A_2 \times W_2$ $Out_1 += A_2 \times W_2$ $Out_0 += A_2 \times W_2$ $Out_1 += A_2 \times W_2$ $Out_0 += A_2 \times W_2$ $Out_1 += A_2 \times W_2$ $Out_0 += A_2 \times W_2$ $Out_1 += A_2 \times$

$Out_0 += A_3 \times W_3$ $Out_1 += A_3 \times W_3$ $Out_0 += A_3 \times W_3$ $Out_1 += A_3 \times W_3$ $Out_0 += A_3 \times W_3$ $Out_1 += A_3 \times W_3$ $Out_0 += A_3 \times W_3$ $Out_1 += A_3 \times$

$Out_0 += A_4 \times W_4$ $Out_1 += A_4 \times W_4$ $Out_0 += A_4 \times W_4$ $Out_1 += A_4 \times W_4$ $Out_0 += A_4 \times W_4$ $Out_1 += A_4 \times W_4$ $Out_0 += A_4 \times W_4$ $Out_1 += A_4 \times$

$Out_0 += A_0 \times W_0$ $Out_1 += A_0 \times W_0$ $Out_0 += A_0 \times W_0$ $Out_1 += A_0 \times W_0$ $Out_0 += A_0 \times W_0$ $Out_1 += A_0 \times W_0$ $Out_0 += A_0 \times W_0$ $Out_1 += A_0 \times$

$Out_0 += A_1 \times W_1$ $Out_1 += A_1 \times W_1$ $Out_0 += A_1 \times W_1$ $Out_1 += A_1 \times W_1$ $Out_0 += A_1 \times W_1$ $Out_1 += A_1 \times W_1$ $Out_0 += A_1 \times W_1$ $Out_1 += A_1 \times$

$Out_0 += A_2 \times W_2$ $Out_1 += A_2 \times W_2$ $Out_0 += A_2 \times W_2$ $Out_1 += A_2 \times W_2$ $Out_0 += A_2 \times W_2$ $Out_1 += A_2 \times W_2$ $Out_0 += A_2 \times W_2$ $Out_1 += A_2 \times$

$Out_0 += A_3 \times W_3$ $Out_1 += A_3 \times W_3$ $Out_0 += A_3 \times W_3$ $Out_1 += A_3 \times W_3$ $Out_0 += A_3 \times W_3$ $Out_1 += A_3 \times W_3$ $Out_0 += A_3 \times W_3$ $Out_1 += A_3 \times$

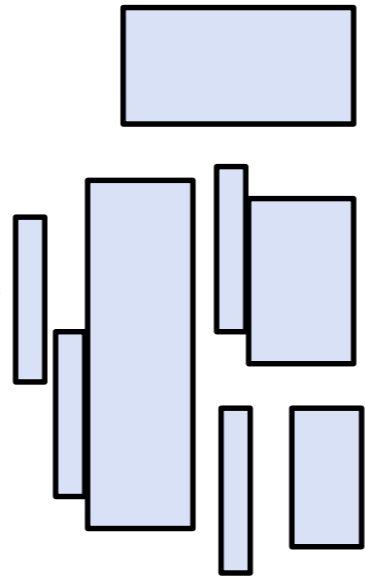$Out_0 += A_4 \times W_4$ $Out_1 += A_4 \times W_4$ $Out_0 += A_4 \times W_4$ $Out_1 += A_4 \times W_4$ $Out_0 += A_4 \times W_4$ $Out_1 += A_4 \times W_4$ $Out_0 += A_4 \times W_4$ $Out_1 += A_4 \times$

Calculate the *same* output but … do less work

$Out_0 \mathrel{+}= A_0 \times W_0$ $Out_1 \mathrel{+}= A_0 \times W_0$ $Out_0 \mathrel{+}= A_0 \times W_0$ $Out_1 \mathrel{+}= A_0 \times W_0$ $Out_0 \mathrel{+}= A_0 \times W_0$ $Out_1 \mathrel{+}= A_0 \times W_0$ $Out_0 \mathrel{+}= A_0 \times W_0$ $Out_1 \mathrel{+}= A_0 \times$

$Out_0 \mathrel{+}= A_1 \times W_1$ $Out_1 \mathrel{+}= A_1 \times W_1$ $Out_0 \mathrel{+}= A_1 \times W_1$ $Out_1 \mathrel{+}= A_1 \times W_1$ $Out_0 \mathrel{+}= A_1 \times W_1$ $Out_1 \mathrel{+}= A_1 \times W_1$ $Out_0 \mathrel{+}= A_1 \times W_1$ $Out_1 \mathrel{+}= A_1 \times$

$Out_0 \mathrel{+}= A_2 \times W_2$ $Out_1 \mathrel{+}= A_2 \times W_2$ $Out_0 \mathrel{+}= A_2 \times W_2$ $Out_1 \mathrel{+}= A_2 \times W_2$ $Out_0 \mathrel{+}= A_2 \times W_2$ $Out_1 \mathrel{+}= A_2 \times W_2$ $Out_0 \mathrel{+}= A_2 \times W_2$ $Out_1 \mathrel{+}= A_2 \times$

$Out_0 \mathrel{+}= A_3 \times W_3$ $Out_1 \mathrel{+}= A_3 \times W_3$ $Out_0 \mathrel{+}= A_3 \times W_3$ $Out_1 \mathrel{+}= A_3 \times W_3$ $Out_0 \mathrel{+}= A_3 \times W_3$ $Out_1 \mathrel{+}= A_3 \times W_3$ $Out_0 \mathrel{+}= A_3 \times W_3$ $Out_1 \mathrel{+}= A_3 \times$

$Out_0 \mathrel{+}= A_4 \times W_4$ $Out_1 \mathrel{+}= A_4 \times W_4$ $Out_0 \mathrel{+}= A_4 \times W_4$ $Out_1 \mathrel{+}= A_4 \times W_4$ $Out_0 \mathrel{+}= A_4 \times W_4$ $Out_1 \mathrel{+}= A_4 \times W_4$ $Out_0 \mathrel{+}= A_4 \times W_4$ $Out_1 \mathrel{+}= A_4 \times$

$Out_0 \mathrel{+}= A_0 \times W_0$ $Out_1 \mathrel{+}= A_0 \times W_0$ $Out_0 \mathrel{+}= A_0 \times W_0$ $Out_1 \mathrel{+}= A_0 \times W_0$ $Out_0 \mathrel{+}= A_0 \times W_0$ $Out_1 \mathrel{+}= A_0 \times W_0$ $Out_0 \mathrel{+}= A_0 \times W_0$ $Out_1 \mathrel{+}= A_0 \times$

$Out_0 \mathrel{+}= A_1 \times W_1$ $Out_1 \mathrel{+}= A_1 \times W_1$ $Out_0 \mathrel{+}= A_1 \times W_1$ $Out_1 \mathrel{+}= A_1 \times W_1$ $Out_0 \mathrel{+}= A_1 \times W_1$ $Out_1 \mathrel{+}= A_1 \times W_1$ $Out_0 \mathrel{+}= A_1 \times W_1$ $Out_1 \mathrel{+}= A_1 \times$

$Out_0 \mathrel{+}= A_2 \times W_2$ $Out_1 \mathrel{+}= A_2 \times W_2$ $Out_0 \mathrel{+}= A_2 \times W_2$ $Out_1 \mathrel{+}= A_2 \times W_2$ $Out_0 \mathrel{+}= A_2 \times W_2$ $Out_1 \mathrel{+}= A_2 \times W_2$ $Out_0 \mathrel{+}= A_2 \times W_2$ $Out_1 \mathrel{+}= A_2 \times$

$Out_0 \mathrel{+}= A_3 \times W_3$ $Out_1 \mathrel{+}= A_3 \times W_3$ $Out_0 \mathrel{+}= A_3 \times W_3$ $Out_1 \mathrel{+}= A_3 \times W_3$ $Out_0 \mathrel{+}= A_3 \times W_3$ $Out_1 \mathrel{+}= A_3 \times W_3$ $Out_0 \mathrel{+}= A_3 \times W_3$ $Out_1 \mathrel{+}= A_3 \times$

$Out_0 \mathrel{+}= A_4 \times W_4$ $Out_1 \mathrel{+}= A_4 \times W_4$ $Out_0 \mathrel{+}= A_4 \times W_4$ $Out_1 \mathrel{+}= A_4 \times W_4$ $Out_0 \mathrel{+}= A_4 \times W_4$ $Out_1 \mathrel{+}= A_4 \times W_4$ $Out_0 \mathrel{+}= A_4 \times W_4$ $Out_1 \mathrel{+}= A_4 \times$

$Out_0 \mathrel{+}= A_0 \times W_0$ $Out_1 \mathrel{+}= A_0 \times W_0$ $Out_0 \mathrel{+}= A_0 \times W_0$ $Out_1 \mathrel{+}= A_0 \times W_0$ $Out_0 \mathrel{+}= A_0 \times W_0$ $Out_1 \mathrel{+}= A_0 \times W_0$ $Out_0 \mathrel{+}= A_0 \times W_0$ $Out_1 \mathrel{+}= A_0 \times$

$Out_0 \mathrel{+}= A_1 \times W_1$ $Out_1 \mathrel{+}= A_1 \times W_1$ $Out_0 \mathrel{+}= A_1 \times W_1$ $Out_1 \mathrel{+}= A_1 \times W_1$ $Out_0 \mathrel{+}= A_1 \times W_1$ $Out_1 \mathrel{+}= A_1 \times W_1$ $Out_0 \mathrel{+}= A_1 \times W_1$ $Out_1 \mathrel{+}= A_1 \times$

$Out_0 \mathrel{+}= A_2 \times W_2$ $Out_1 \mathrel{+}= A_2 \times W_2$ $Out_0 \mathrel{+}= A_2 \times W_2$ $Out_1 \mathrel{+}= A_2 \times W_2$ $Out_0 \mathrel{+}= A_2 \times W_2$ $Out_1 \mathrel{+}= A_2 \times W_2$ $Out_0 \mathrel{+}= A_2 \times W_2$ $Out_1 \mathrel{+}= A_2 \times$

$Out_0 \mathrel{+}= A_3 \times W_3$ $Out_1 \mathrel{+}= A_3 \times W_3$ $Out_0 \mathrel{+}= A_3 \times W_3$ $Out_1 \mathrel{+}= A_3 \times W_3$ $Out_0 \mathrel{+}= A_3 \times W_3$ $Out_1 \mathrel{+}= A_3 \times W_3$ $Out_0 \mathrel{+}= A_3 \times W_3$ $Out_1 \mathrel{+}= A_3 \times$

$Out_0 \mathrel{+}= A_4 \times W_4$ $Out_1 \mathrel{+}= A_4 \times W_4$ $Out_0 \mathrel{+}= A_4 \times W_4$ $Out_1 \mathrel{+}= A_4 \times W_4$ $Out_0 \mathrel{+}= A_4 \times W_4$ $Out_1 \mathrel{+}= A_4 \times W_4$ $Out_0 \mathrel{+}= A_4 \times W_4$ $Out_1 \mathrel{+}= A_4 \times$

$$Out_0 \mathrel{+}= A_0 \times W_0 \quad Out_1 \mathrel{+}= A_0 \times W_0 \quad Out_0 \mathrel{+}= A_0 \times W_0 \quad Out_1 \mathrel{+}= A_0 \times W_0 \quad Out_0 \mathrel{+}= A_0 \times W_0 \quad Out_1 \mathrel{+}= A_0 \times W_0 \quad Out_0 \mathrel{+}= A_0 \times W_0 \quad Out_1 \mathrel{+}= A_0 \times$$

$$Out_0 \mathrel{+}= A_1 \times W_1 \quad Out_1 \mathrel{+}= A_1 \times W_1 \quad Out_0 \mathrel{+}= A_1 \times W_1 \quad \cdots \quad Out_1 \mathrel{+}= A_1 \times W_1 \quad Out_0 \mathrel{+}= A_1 \times W_1 \quad Out_1 \mathrel{+}= A_1 \times$$

$$Out_0 \mathrel{+}= A_2 \times W_2 \quad Out_1 \mathrel{+}= A_2 \times W_2 \quad Out_0 \mathrel{+}= A_2 \times \cdots \mathrel{+}= A_2 \times W_2 \quad Out_0 \mathrel{+}= A_2 \times W_2 \quad Out_1 \mathrel{+}= A_2 \times$$

$$Out_0 \mathrel{+}= A_3 \times W_3 \quad Out_1 \mathrel{+}= A_3 \times W_3 \quad Out_0 \mathrel{+}= \cdots \mathrel{+}= A_3 \times W_3 \quad Out_0 \mathrel{+}= A_3 \times W_3 \quad Out_1 \mathrel{+}= A_3 \times$$

$$Out_0 \mathrel{+}= A_4 \times W_4 \quad Out_1 \mathrel{+}= A_4 \times W_4 \quad Out_0 \mathrel{+}= \cdots \mathrel{+}= A_4 \times W_4 \quad Out_0 \mathrel{+}= A_4 \times W_4 \quad Out_1 \mathrel{+}= A_4 \times$$

$$Out_0 \mathrel{+}= A_0 \times W_0 \quad Out_1 \mathrel{+}= A_0 \times W_0 \quad Out_0 \mathrel{+}= \cdots \mathrel{+}= A_0 \times W_0 \quad Out_0 \mathrel{+}= A_0 \times W_0 \quad Out_1 \mathrel{+}= A_0 \times$$

$$Out_0 \mathrel{+}= A_1 \times W_1 \quad Out_1 \mathrel{+}= A_1 \times W_1 \quad Out_0 \mathrel{+}= \cdots A_1 \times W_1 \quad Out_0 \mathrel{+}= A_1 \times W_1 \quad Out_1 \mathrel{+}= A_1 \times$$

$$Out_0 \mathrel{+}= A_2 \times W_2 \quad Out_1 \mathrel{+}= A_2 \times W_2 \quad Out_0 \mathrel{+}= A_2 \cdots Out_1 \mathrel{+}= A_2 \times W_2 \quad Out_0 \mathrel{+}= A_2 \times W_2 \quad Out_1 \mathrel{+}= A_2 \times$$

$$Out_0 \mathrel{+}= A_3 \times W_3 \quad Out_1 \mathrel{+}= A_3 \times W_3 \quad Out_0 \mathrel{+}= A_3 \times W_3 \cdots W_3 \quad Out_1 \mathrel{+}= A_3 \times W_3 \quad Out_0 \mathrel{+}= A_3 \times W_3 \quad Out_1 \mathrel{+}= A_3 \times$$

$$Out_0 \mathrel{+}= A_4 \times W_4 \quad Out_1 \mathrel{+}= A_4 \times W_4 \quad Out_0 \mathrel{+}= A_4 \times W_4 \cdots W_4 \quad Out_1 \mathrel{+}= A_4 \times W_4 \quad Out_0 \mathrel{+}= A_4 \times W_4 \quad Out_1 \mathrel{+}= A_4 \times$$

$$Out_0 \mathrel{+}= A_0 \times W_0 \quad Out_1 \mathrel{+}= A_0 \times W_0 \quad Out_0 \mathrel{+}= A_0 \times \cdots W_0 \quad Out_1 \mathrel{+}= A_0 \times W_0 \quad Out_0 \mathrel{+}= A_0 \times W_0 \quad Out_1 \mathrel{+}= A_0 \times$$

$$Out_0 \mathrel{+}= A_1 \times W_1 \quad Out_1 \mathrel{+}= A_1 \times W_1 \quad Out_0 \mathrel{+}= A_1 \times W \cdots W_1 \quad Out_1 \mathrel{+}= A_1 \times W_1 \quad Out_0 \mathrel{+}= A_1 \times W_1 \quad Out_1 \mathrel{+}= A_1 \times$$

$$Out_0 \mathrel{+}= A_2 \times W_2 \quad Out_1 \mathrel{+}= A_2 \times W_2 \quad Out_0 \mathrel{+}= A_2 \times W_2 \cdots A_2 \times W_2 \quad Out_1 \mathrel{+}= A_2 \times W_2 \quad Out_0 \mathrel{+}= A_2 \times W_2 \quad Out_1 \mathrel{+}= A_2 \times$$

$$Out_0 \mathrel{+}= A_3 \times W_3 \quad Out_1 \mathrel{+}= A_3 \times W_3 \quad Out_0 \mathrel{+}= A_3 \times W_3 \quad Out_1 \mathrel{+}= \cdots A_3 \times W_3 \quad Out_1 \mathrel{+}= A_3 \times W_3 \quad Out_0 \mathrel{+}= A_3 \times W_3 \quad Out_1 \mathrel{+}= A_3 \times$$

$$Out_0 \mathrel{+}= A_4 \times W_4 \quad Out_1 \mathrel{+}= A_4 \times W_4 \quad Out_0 \mathrel{+}= A_4 \times W_4 \quad Out_1 \mathrel{+}= A_4 \times W_4 \quad Out_0 \mathrel{+}= A_4 \times W_4 \quad Out_1 \mathrel{+}= A_4 \times W_4 \quad Out_0 \mathrel{+}= A_4 \times W_4 \quad Out_1 \mathrel{+}= A_4 \times$$

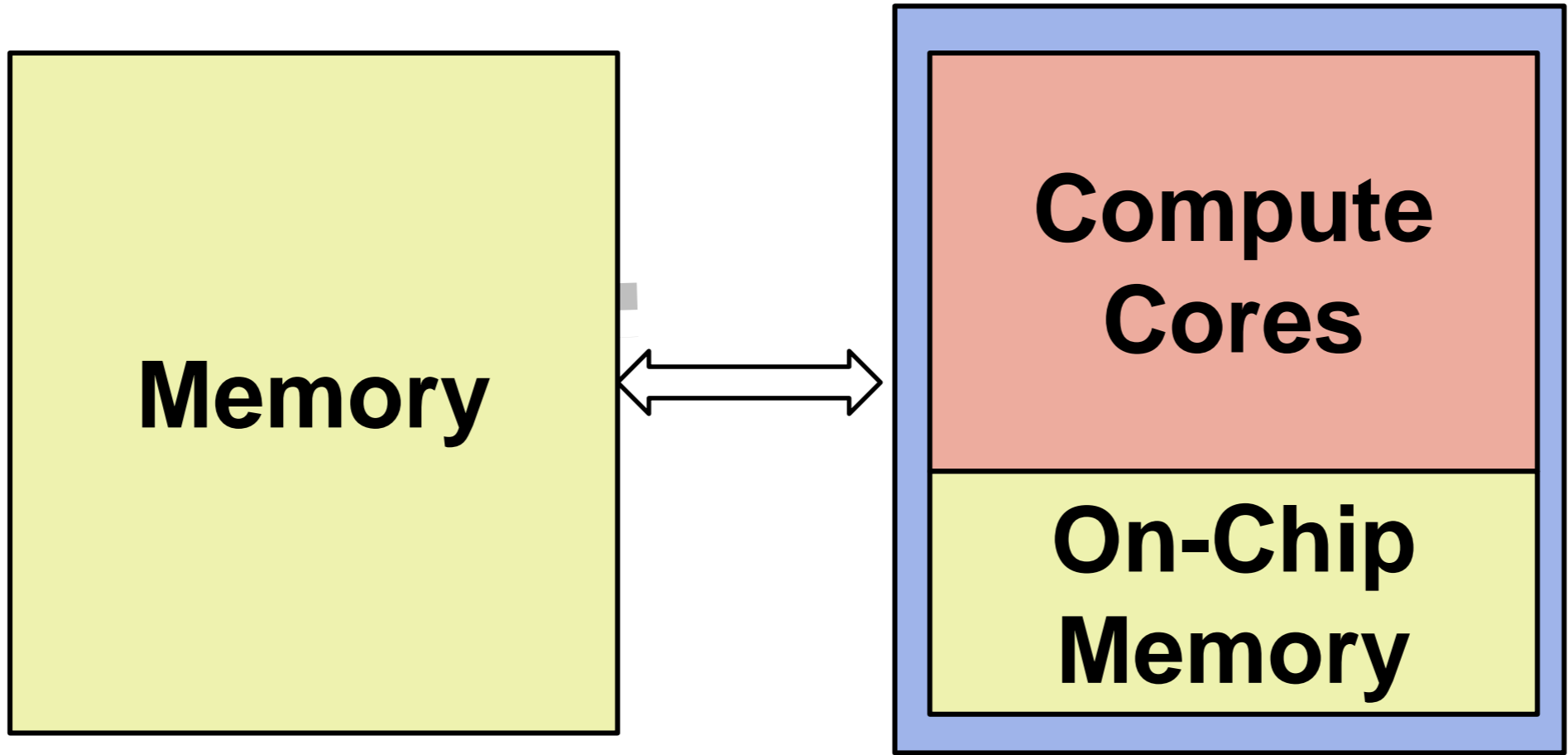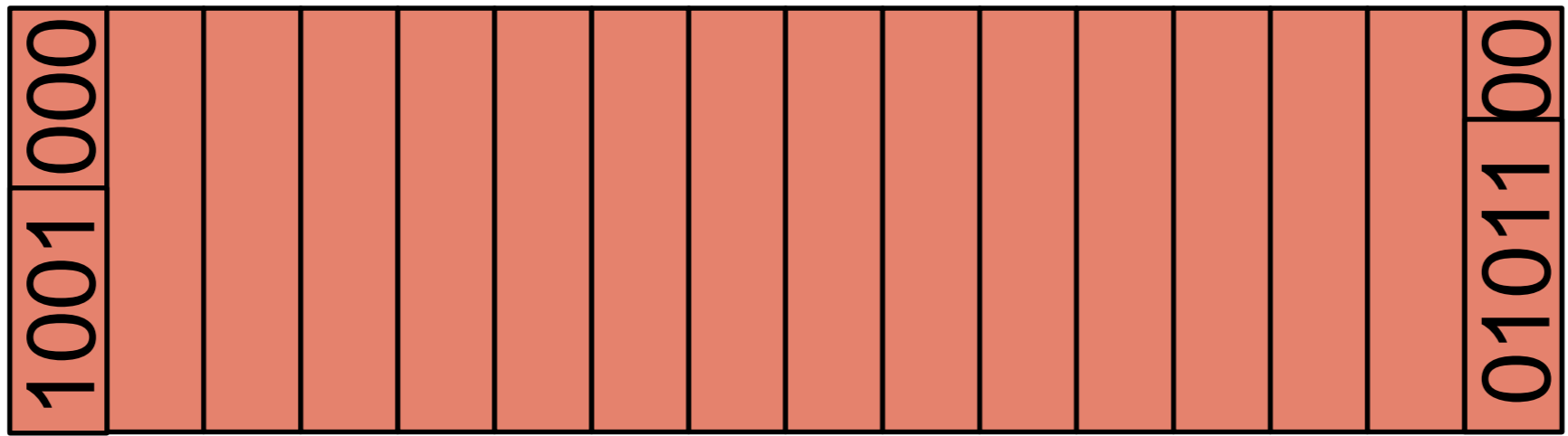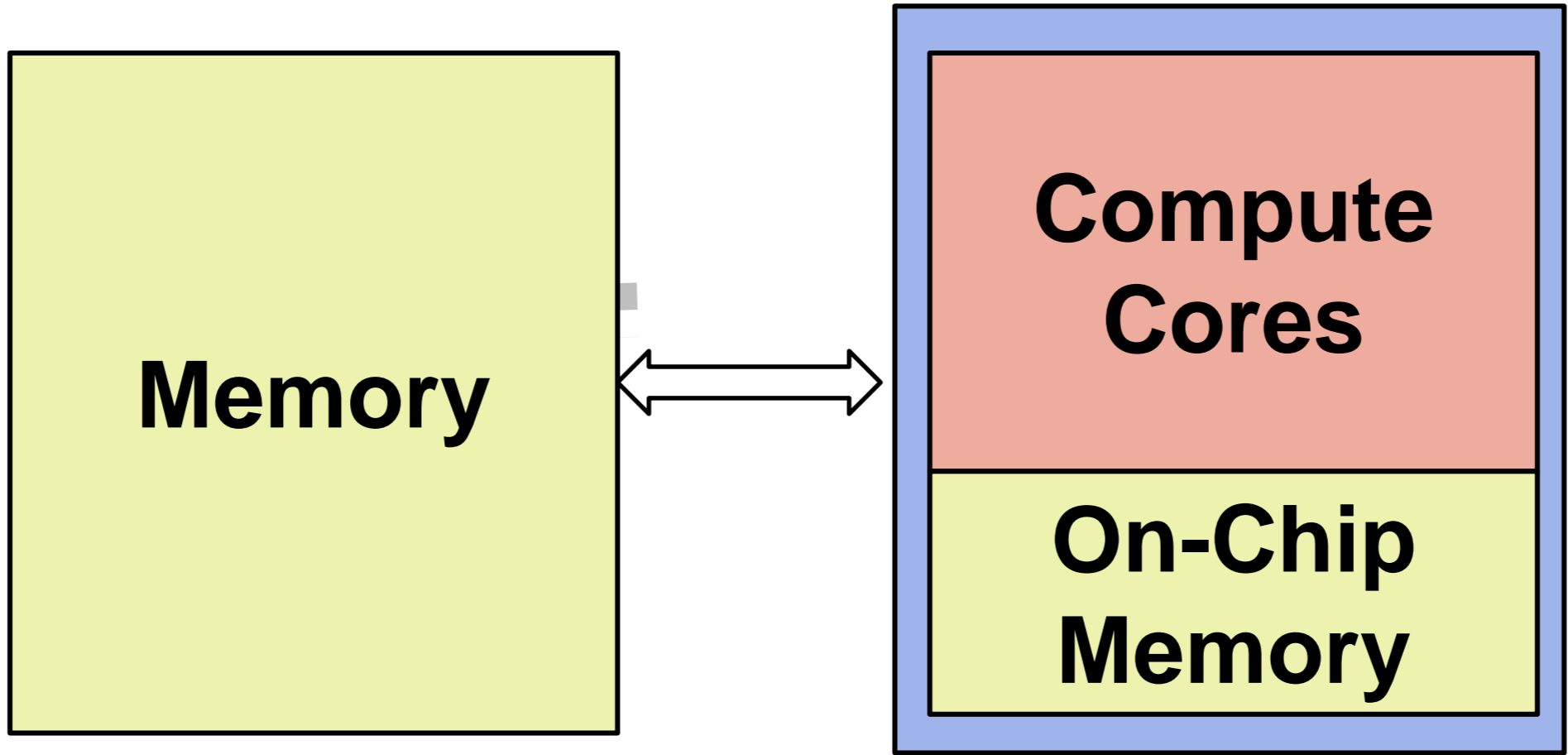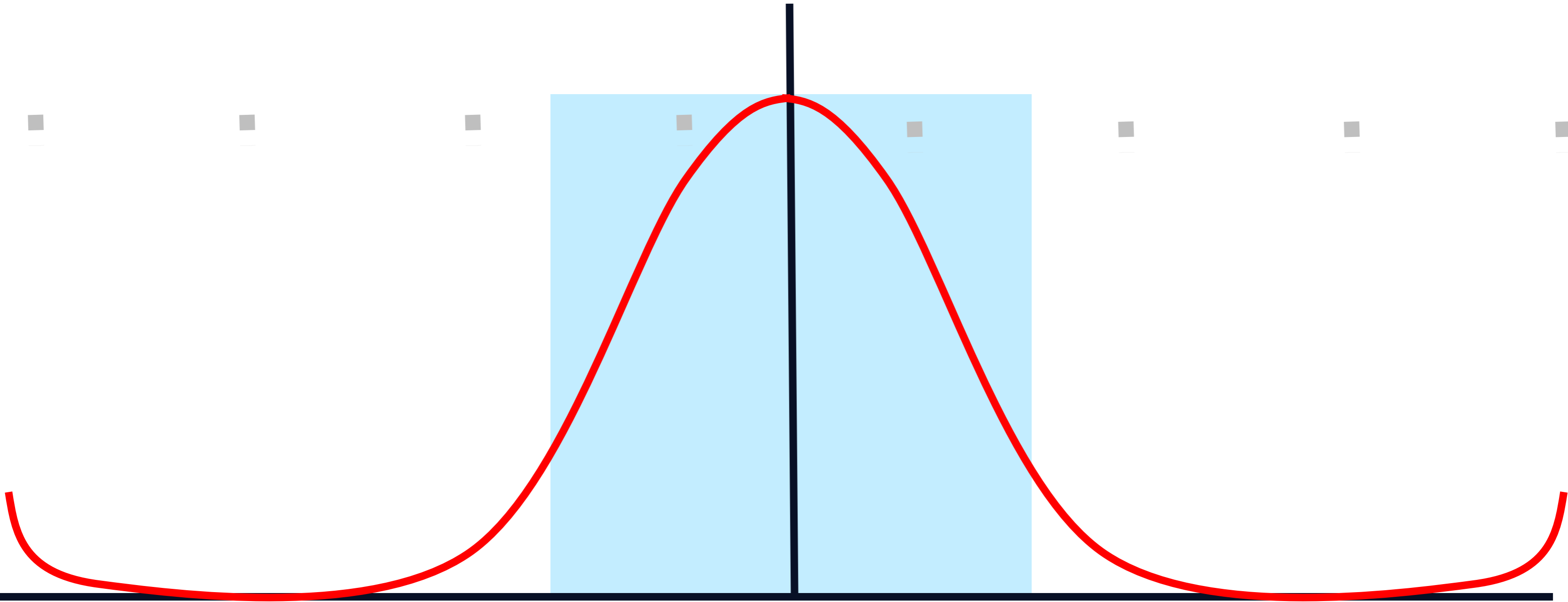$Out_0 += A_0 \times W_0$ $Out_1 += A_0 \times W_0$ $Out_0 += A_0 \times W_0$ $Out_1 += A_0 \times W_0$ $Out_0 += A_0 \times W_0$ $Out_1 += A_0 \times W_0$ $Out_1 += A_0 \times$
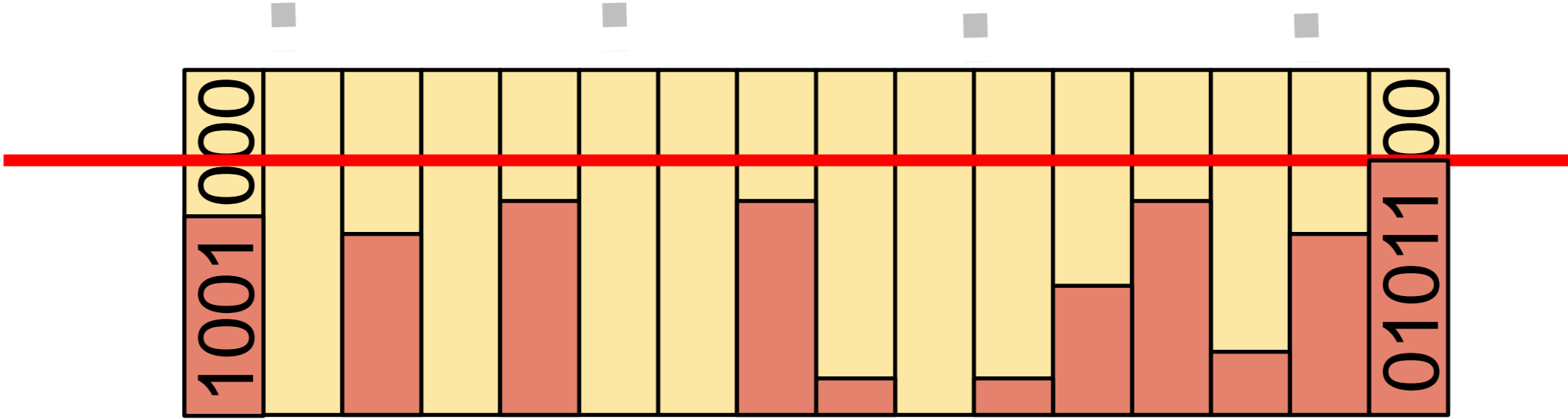
$Out_0 += A_4 \times W_4$ $Out_1 += A_1 \times W_1$ $Out_1 += A_1 \times W_1$

$Out_1 += A_3 \times W_3$

$Out_1 += A_4 \times W_4$ $Out_1 += A_4 \times$

$Out_0 += A_4 \times W_4$

$Out_1 += A_0 \times W_0$ $Out_0 += A_0 \times W_0$ $2$ $Out_0 += A_0 \times W_0$ $Out_1 += A_0 \times W_0$ $Out_1 += A_0 \times$

$Out_1 += A_1 \times W_1$ $Out_0 += A_1 \times W_1$ $Out_1 += A_3 \times W_3$ $Out_0 += A_1 \times W_1$ $Out_1 += A_1 \times$

$Out_0 += A_2 \times W_2$ $Out_1 += A_4 \times W_4$

$Out_0 += A_1 \times W_1$

$Out_0 += A_3 \times W_3$

$Out_0 += A_4 \times W_4$

$Out_1 += A_0 \times W_0$ $Out_0 += A_0 \times W_0$ $Out_1 += A_0 \times W_0$ $Out_0 += A_0 \times W_0$ $Out_1 += A_0 \times$

$Out_0 += A_1 \times W_1$

$Out_1 += A_4 \times$

$Out_1 += A_4 \times W_4$

$Out_0 += A_4 \times W_4$ $Out_1 += A_4 \times W_4$

$Out_0 \mathrel{+}= A_0 \times W_0$

$Out_0 \mathrel{+}= A_4 \times W_4$

$Out_1 \mathrel{+}= A_0 \times W_0$

$Out_0 \mathrel{+}= A_0 \times W_0$

$Out_1 \mathrel{+}= A_0 \times W_0$

$Out_1 \mathrel{+}= A_1 \times W_1$

$Out_0 \mathrel{+}= A_0 \times W_0$

$Out_1 \mathrel{+}= A_0 \times W_0$

$Out_1 \mathrel{+}= A_1 \times W_1$

$Out_1 \mathrel{+}= A_0 \times$

$Out_1 \mathrel{+}= A_3 \times W_3$

$Out_1 \mathrel{+}= A_4 \times W_4$

$Out_1 \mathrel{+}= A_4 \times$

$Out_0 \mathrel{+}= A_4 \times W_4$

$Out_1 \mathrel{+}= A_0 \times W_0$

$Out_1 \mathrel{+}= A_1 \times W_1$

$Out_0 \mathrel{+}= A_0 \times W_0$

$Out_0 \mathrel{+}= A_1 \times W_1$

2

$Out_1 \mathrel{+}= A_3 \times W_3$

$Out_0 \mathrel{+}= A_0 \times W_0$

$Out_0 \mathrel{+}= A_1 \times W_1$

$Out_1 \mathrel{+}= A_0 \times W_0$

$Out_1 \mathrel{+}= A_0 \times$

$Out_1 \mathrel{+}= A_1 \times$

$Out_0 \mathrel{+}= A_2 \times W_2$

$Out_1 \mathrel{+}= A_4 \times W_4$

$Out_0 \mathrel{+}= A_1 \times W_1$

$Out_0 \mathrel{+}= A_3 \times W_3$

$Out_0 \mathrel{+}= A_4 \times W_4$

$Out_1 \mathrel{+}= A_0 \times W_0$

$Out_0 \mathrel{+}= A_0 \times W_0$

$Out_1 \mathrel{+}= A_0 \times W_0$

$Out_0 \mathrel{+}= A_0 \times W_0$

$Out_1 \mathrel{+}= A_0 \times$

$Out_0 \mathrel{+}= A_1 \times W_1$

$Out_1 \mathrel{+}= A_4 \times$

$Out_0 \mathrel{+}= A_4 \times W_4$

$Out_1 \mathrel{+}= A_4 \times W_4$

$Out_1 \mathrel{+}= A_4 \times W_4$

"Compiler"

Runtime

Hardware

- $Out\ +=\ A_0\ x\ W_0$

$$0 \ x \ W$$

$$\sim 0 \ x \ W$$

- **Out += $A_0$ x $W_0$**

**0 x W**

**~50%**

**0 x W**

- **At Various Granularities: e.g., Layer or finer**

# A x W

# ~60%-

- At Various Granularities e.g., layer or finer

```
W 0001 0100
A 0010 1010
───────────────
  0000 0000
 0010 1010
 0000 0000
0010 1010
0000 0000
0010 1010
0000 0000
0000 0000
```

W 0001 0100

A 0010 1010

1 1 1

1 1 1

1 1 1

# A x W

# 75%-99%

**Memory** | **Compute Cores** | **On-Chip Memory**

**On- vs. Off-Chip**
Energy: ~100x
Latency: ~50x

$$Out_0 += A_0 \times W_0 \quad Out_1 += A_0 \times W_0 \quad Out_0 += A_0 \times W_0 \quad Out_1 += A_0 \times W_0 \quad Out_0 += A_0 \times W_0 \quad Out_1 += A_0 \times W_0 \quad Out_0 += A_0 \times W_0 \quad Out_1 += A_0 \times$$

$$Out_0 += A_1 \times W_1 \quad Out_1 += A_1 \times W_1 \quad Out_0 += A_1 \times W_1 \quad Out_1 += A_1 \times W_1 \quad Out_0 += A_1 \times W_1 \quad Out_1 += A_1 \times W_1 \quad Out_0 += A_1 \times W_1 \quad Out_1 += A_1 \times$$

$$Out_0 += A_2 \times W_2 \quad Out_1 += A_2 \times W_2 \quad Out_0 += A_2 \times W_2 \quad Out_1 += A_2 \times W_2 \quad Out_0 += A_2 \times W_2 \quad Out_1 += A_2 \times W_2 \quad Out_0 += A_2 \times W_2 \quad Out_1 += A_2 \times$$

$$Out_0 += A_3 \times W_3 \quad Out_1 += A_3 \times W_3 \quad Out_0 += A_3 \times W_3 \quad Out_1 += A_3 \times W_3 \quad Out_0 += A_3 \times W_3 \quad Out_1 += A_3 \times W_3 \quad Out_0 += A_3 \times W_3 \quad Out_1 += A_3 \times$$

$$Out_0 += A_4 \times W_4 \quad Out_1 += A_4 \times W_4 \quad Out_0 += A_4 \times W_4 \quad Out_1 += A_4 \times W_4 \quad Out_0 += A_4 \times W_4 \quad Out_1 += A_4 \times W_4 \quad Out_0 += A_4 \times W_4 \quad Out_1 += A_4 \times$$

$$Out_0 += A_0 \times W_0 \quad Out_1 += A_0 \times W_0 \quad Out_0 += A_0 \times W_0 \quad Out_1 += A_0 \times W_0 \quad Out_0 += A_0 \times W_0 \quad Out_1 += A_0 \times W_0 \quad Out_0 += A_0 \times W_0 \quad Out_1 += A_0 \times$$

$$Out_0 += A_1 \times W_1 \quad Out_1 += A_1 \times W_1 \quad Out_0 += A_1 \times W_1 \quad Out_1 += A_1 \times W_1 \quad Out_0 += A_1 \times W_1 \quad Out_1 += A_1 \times W_1 \quad Out_0 += A_1 \times W_1 \quad Out_1 += A_1 \times$$

$$Out_0 += A_2 \times W_2 \quad Out_1 += A_2 \times W_2 \quad Out_0 += A_2 \times W_2 \quad Out_1 += A_2 \times W_2 \quad Out_0 += A_2 \times W_2 \quad Out_1 += A_2 \times W_2 \quad Out_0 += A_2 \times W_2 \quad Out_1 += A_2 \times$$

$$Out_0 += A_3 \times W_3 \quad Out_1 += A_3 \times W_3 \quad Out_0 += A_3 \times W_3 \quad Out_1 += A_3 \times W_3 \quad Out_0 += A_3 \times W_3 \quad Out_1 += A_3 \times W_3 \quad Out_0 += A_3 \times W_3 \quad Out_1 += A_3 \times$$

$$Out_0 += A_4 \times W_4 \quad Out_1 += A_4 \times W_4 \quad Out_0 += A_4 \times W_4 \quad Out_1 += A_4 \times W_4 \quad Out_0 += A_4 \times W_4 \quad Out_1 += A_4 \times W_4 \quad Out_0 += A_4 \times W_4 \quad Out_1 += A_4 \times$$

$$Out_0 += A_0 \times W_0 \quad Out_1 += A_0 \times W_0 \quad Out_0 += A_0 \times W_0 \quad Out_1 += A_0 \times W_0 \quad Out_0 += A_0 \times W_0 \quad Out_1 += A_0 \times W_0 \quad Out_0 += A_0 \times W_0 \quad Out_1 += A_0 \times$$

$$Out_0 += A_1 \times W_1 \quad Out_1 += A_1 \times W_1 \quad Out_0 += A_1 \times W_1 \quad Out_1 += A_1 \times W_1 \quad Out_0 += A_1 \times W_1 \quad Out_1 += A_1 \times W_1 \quad Out_0 += A_1 \times W_1 \quad Out_1 += A_1 \times$$

$$Out_0 += A_2 \times W_2 \quad Out_1 += A_2 \times W_2 \quad Out_0 += A_2 \times W_2 \quad Out_1 += A_2 \times W_2 \quad Out_0 += A_2 \times W_2 \quad Out_1 += A_2 \times W_2 \quad Out_0 += A_2 \times W_2 \quad Out_1 += A_2 \times$$

$$Out_0 += A_3 \times W_3 \quad Out_1 += A_3 \times W_3 \quad Out_0 += A_3 \times W_3 \quad Out_1 += A_3 \times W_3 \quad Out_0 += A_3 \times W_3 \quad Out_1 += A_3 \times W_3 \quad Out_0 += A_3 \times W_3 \quad Out_1 += A_3 \times$$

$$Out_0 += A_4 \times W_4 \quad Out_1 += A_4 \times W_4 \quad Out_0 += A_4 \times W_4 \quad Out_1 += A_4 \times W_4 \quad Out_0 += A_4 \times W_4 \quad Out_1 += A_4 \times W_4 \quad Out_0 += A_4 \times W_4 \quad Out_1 += A_4 \times$$

$Out_0 \mathrel{+}= A_0 \times W_0$ $Out_1 \mathrel{+}= A_0 \times W_0$ $Out_0 \mathrel{+}= A_0 \times W_0$ $Out_1 \mathrel{+}= A_0 \times W_0$ $Out_0 \mathrel{+}= A_0 \times W_0$ $Out_1 \mathrel{+}= A_0 \times W_0$ $Out_0 \mathrel{+}= A_0 \times W_0$ $Out_1 \mathrel{+}= A_0 \times W$

$Out_0 \mathrel{+}= A_0 \times W_0$ $Out_1 \mathrel{+}= A_0 \times W_0$ $Out_0 \mathrel{+}= A_0 \times W_0$ $Out_1 \mathrel{+}= A_0 \times W_0$ $Out_0 \mathrel{+}= A_0 \times W_0$ $Out_1 \mathrel{+}= A_0 \times W_0$ $Out_0 \mathrel{+}= A_0 \times W_0$ $Out_1 \mathrel{+}= A_0 \times W$

**Memory**

**Compute Cores**

**On-Chip Memory**

$Out_0 += A_0 \times W_0$ $Out_1 += A_0 \times W_0$ $Out_0 += A_0 \times W_0$ $Out_1 += A_0 \times W_0$ $Out_0 += A_0 \times W_0$ $Out_1 += A_0 \times W_0$ $Out_0 += A_0 \times W_0$ $Out_1 += A_0 \times W$

**Memory**

**Compute Cores**

**On-Chip Memory**

1001 000

01011 00

$Out_0 += A_0 \times W_0$ $Out_1 += A_0 \times W_0$ $Out_0 += A_0 \times W_0$ $Out_1 += A_0 \times W_0$ $Out_0 += A_0 \times W_0$ $Out_1 += A_0 \times W_0$ $Out_0 += A_0 \times W_0$ $Out_1 += A_0 \times W$



*DPRed: Making Typical Activation and Weight Values Matter In Deep Learning Computing*, Delmas et al.,
https://arxiv.org/abs/1804.06732

$Out_0 += A_0 \times W_0$ $Out_1 += A_0 \times W_0$ $Out_0 += A_0 \times W_0$ $Out_1 += A_0 \times W_0$ $Out_0 += A_0 \times W_0$ $Out_1 += A_0 \times W_0$ $Out_0 += A_0 \times W_0$ $Out_1 += A_0 \times W$



1001 000

01011 00

# Reduces traffic to 50% -- 25% of the original

Precisions are much lower

| GoogleNet | 6.19-5.94-5.74-6.77-6.91-6.77-6.86-6.77 |
|-----------|------------------------------------------|
|           | -6.92-6.31-5.96-6.31-6.00-6.31-6.55-5.33 |
|           | -5.33-5.33-5.33-5.33-5.48-6.74-6.33-6.74 |
|           | -6.51-6.74-7.07-6.35-6.17-6.35-5.88-6.35 |
|           | -6.56-5.07-4.69-5.07-4.82-5.07-5.31-5.53 |
|           | -4.89-5.53-5.70-5.53-5.86-7.88-7.62-7.88 |
|           | -8.07-7.88-8.31-4.97-3.85-4.97-3.61-4.97-5.36 |

# Also: Computation Speed:

# Ineffectual Computations Explained

```
W 0001 0100
A 0010 1010
─────────────
  0000 0000
  0010 1010
  0000 0000
  0010 1010
  0000 0000
  0010 1010
  0000 0000
  0000 0000
```

# Ineffectual Computations Explained

```
W 0001 0100
A 0010 1010
_____

0010 1010
0000 0000
0010 1010
0000 0000
0010 1010
```

```
W 0001 0100
A 0010 1010
```

```
0010 1010
```

```
0010 1010
```

```
0010 1010
```

```
W 0001 0100
A 0010 1010
```

# 75%-95%

**ImageNet ≈ 230 x 230**



**Fixed Resolution**

pooling

pooling

Fully Connected

Dog
Cat
Lion

**Different resolution (typically lower)**

**ImageNet ≈ 230 x 230**



**Fixed Resolution**

pooling

**Different resolution (typically lower)**

pooling

Fully Connected

Dog
Cat
Lion

**More Channels**

**Any resolution**

**Same resolution**

**Per-pixel Prediction**

Out1   Out2

Out1 **+= $A_0$ x $W_0$**

Out1 **+= $A_1$ x $W_1$**

Out1 **+= $A_2$ x $W_2$**

Out1 **+= $A_3$ x $W_3$**

**A**

**W**

Out2 **+= $A'_0$ x $W_0$**

Out2 **+= $A'_1$ x $W_1$**

Out2 **+= $A'_2$ x $W_2$**

Out2 **+= $A'_3$ x $W_3$**

**Out1** $+= A_0 \times W_0$

**Out2** $+= A'_0 \times W_0$

**Out1** **+= A$_0$ x W$_0$ = 0110 x W$_0$**

Cost = 2

**Out2 += A'$_0$ x W$_0$ = 0111 x W$_0$**

**Cost = 3**

**Out1** **+= A$_0$ x W$_0$ = 0110 x W$_0$**

Cost = 2

**Out2 += A'$_0$ x W$_0$ = 0111 x W$_0$**

**Out2 += Out1 + (A'$_0$ - A$_0$) x W$_0$**

**Out1** $+= A_0$ x $W_0$ = **0110** x $W_0$

Cost = 2

**Out2** $+= A'_0$ x $W_0$ = **0111** x $W_0$

**Out2** $+=$ **Out1** $+ (A'_0 - A_0)$ x $W_0$

**Out2** $+=$ **Out1** $+ ($**0111** - **0110**$)$ x $W_0$

**Out1** **+=** $A_0$ **x** $W_0$ **= 0110 x** $W_0$

Cost = 2

**Out2** **+=** $A'_0$ **x** $W_0$ **= 0111 x** $W_0$

**Out2** **+=** **Out1** **+** $(A'_0 - A_0)$ **x** $W_0$

**Out2** **+=** **Out1** **+(0111- 0110) x** $W_0$

**Out2** **+=** **Out1** **+** **0001 x** $W_0$

Cost = 1

# Pruning and Sparsity

**Do as you are told?**

$$Out\ {+}{=}\ A_0 \times W_0$$

$$Out\ {+}{=}\ A_1 \times W_1$$

$$Out\ {+}{=}\ A_2 \times W_2$$

$$Out\ {+}{=}\ A_3 \times W_3$$

$$Out\ {+}{=}\ A_4 \times W_4$$

**Do as you are told?**

$$\text{Out} \mathrel{+}= A_0 \text{ x } W_0$$

$$\text{Out} \mathrel{+}= A_1 \text{ x } 0$$

$$\text{Out} \mathrel{+}= A_2 \text{ x } W_2$$

$$\text{Out} \mathrel{+}= A_3 \text{ x } W_3$$

$$\text{Out} \mathrel{+}= A_4 \text{ x } W_4$$

**Do as you are told?**

$$\text{Out} \mathrel{+}= A_0 \times W_0$$

$$\text{Out} \mathrel{+}= A_2 \times W_2$$

$$\text{Out} \mathrel{+}= A_3 \times W_3$$

$$\text{Out} \mathrel{+}= A_4 \times W_4$$

**Do as you are told?**

$$Out \mathrel{+}= A_0 \times W_0$$

$$Out \mathrel{+}= A_2 \times W_2$$

$$Out \mathrel{+}= A_3 \times W_3$$

$$Out \mathrel{+}= A_4 \times W_4$$

# Removing Zero Weights

**Prior efforts**

⭐ Designed for

reward

effort

**Cambricon-X**

# Removing Zero Weights

# Bit-Tactical

- **We build tools**
- **Used by "everyone" for "everything"**
- **Science, medicine, commerce, …**

**Provide Immediate Benefits**

**Reward but do not Require Model Optimization**

**Path to Innovation**

**Intelligent Sensing**

System Software

Hardware

**Sanja Fidler**

**Chris Pal**

**Raquel Urtasun**

**Yoshua Bengio**