# Fast Learning on Slow Hardware

Tor Aamodt (UBC)

Fastpath Workshop, MICRO 2022

Oct 2, 2022

# Acknowledgements
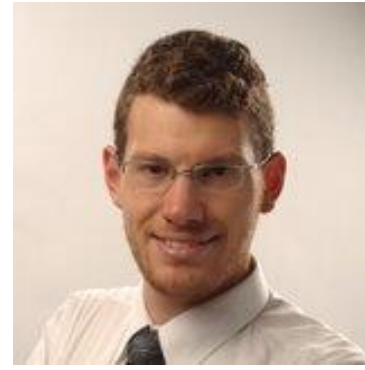


Negar Goli    Aamir Raihan    Jonathan Lew    Dave Evans
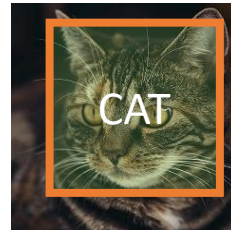
Also: Wenyi Gong, Yunpeng Liu
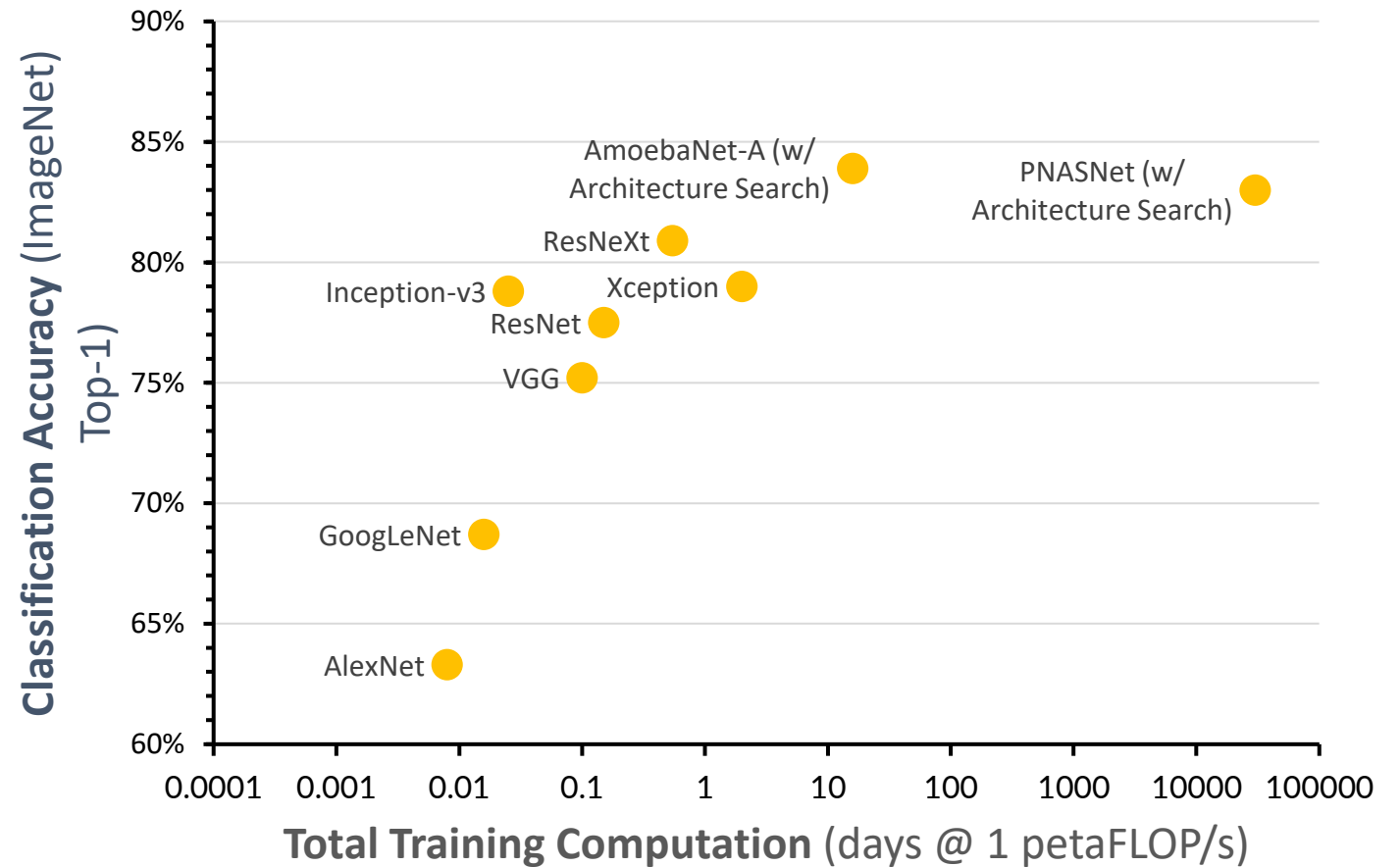
# Neural Networks are Expensive to Train

Classification

CAT

Detection

CAT

Game AI

Self Driving



Classification Accuracy (ImageNet) Top-1

- AmoebaNet-A (w/ Architecture Search)
- PNASNet (w/ Architecture Search)
- ResNeXt
- Inception-v3
- Xception
- ResNet
- VGG
- GoogLeNet
- AlexNet

**Total Training Computation** (days @ 1 petaFLOP/s)
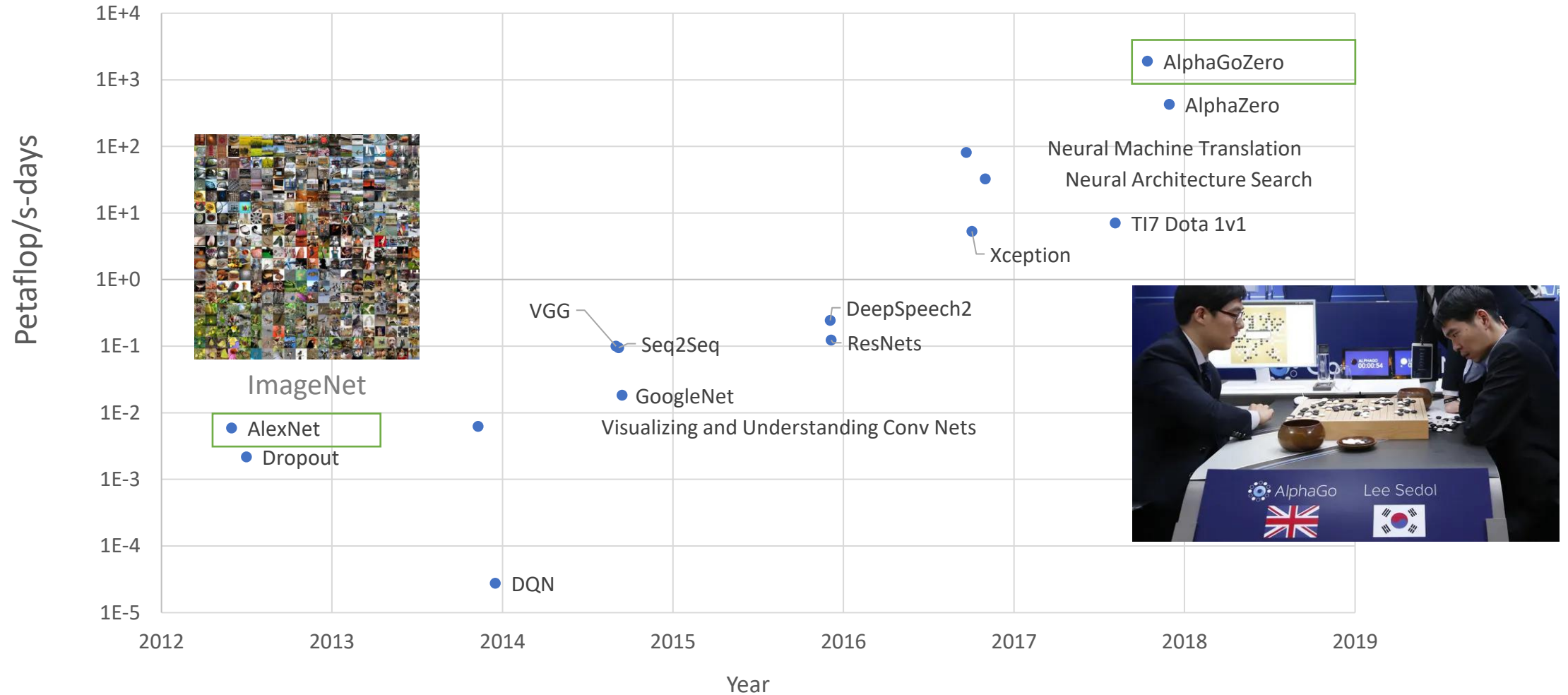
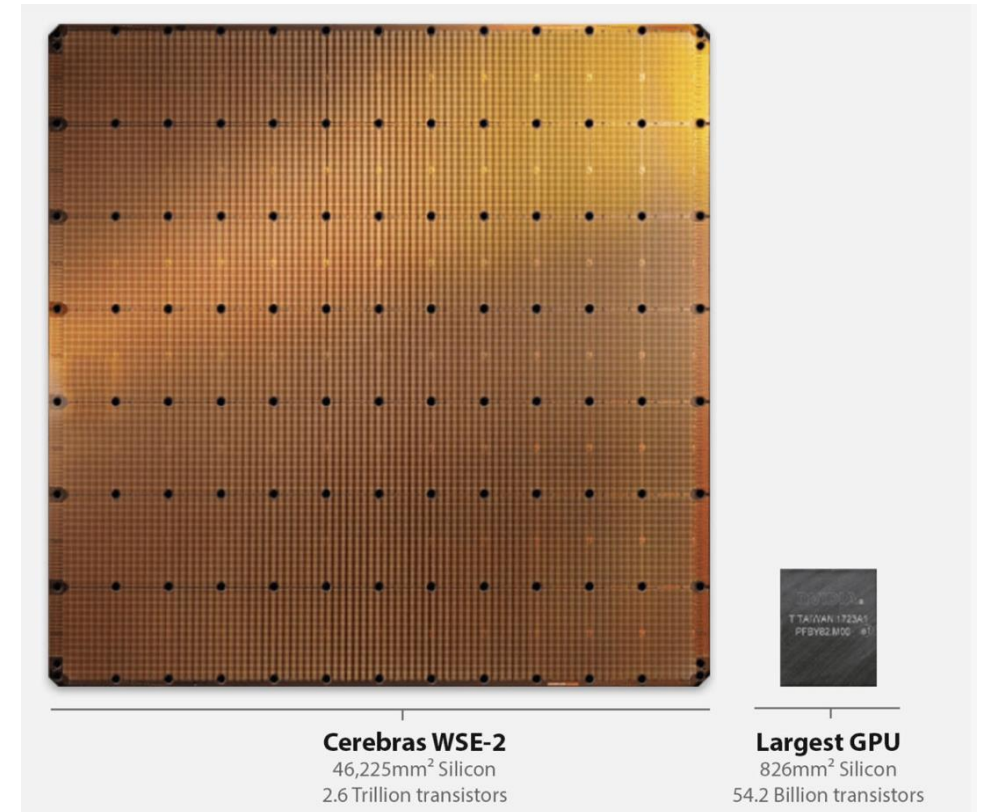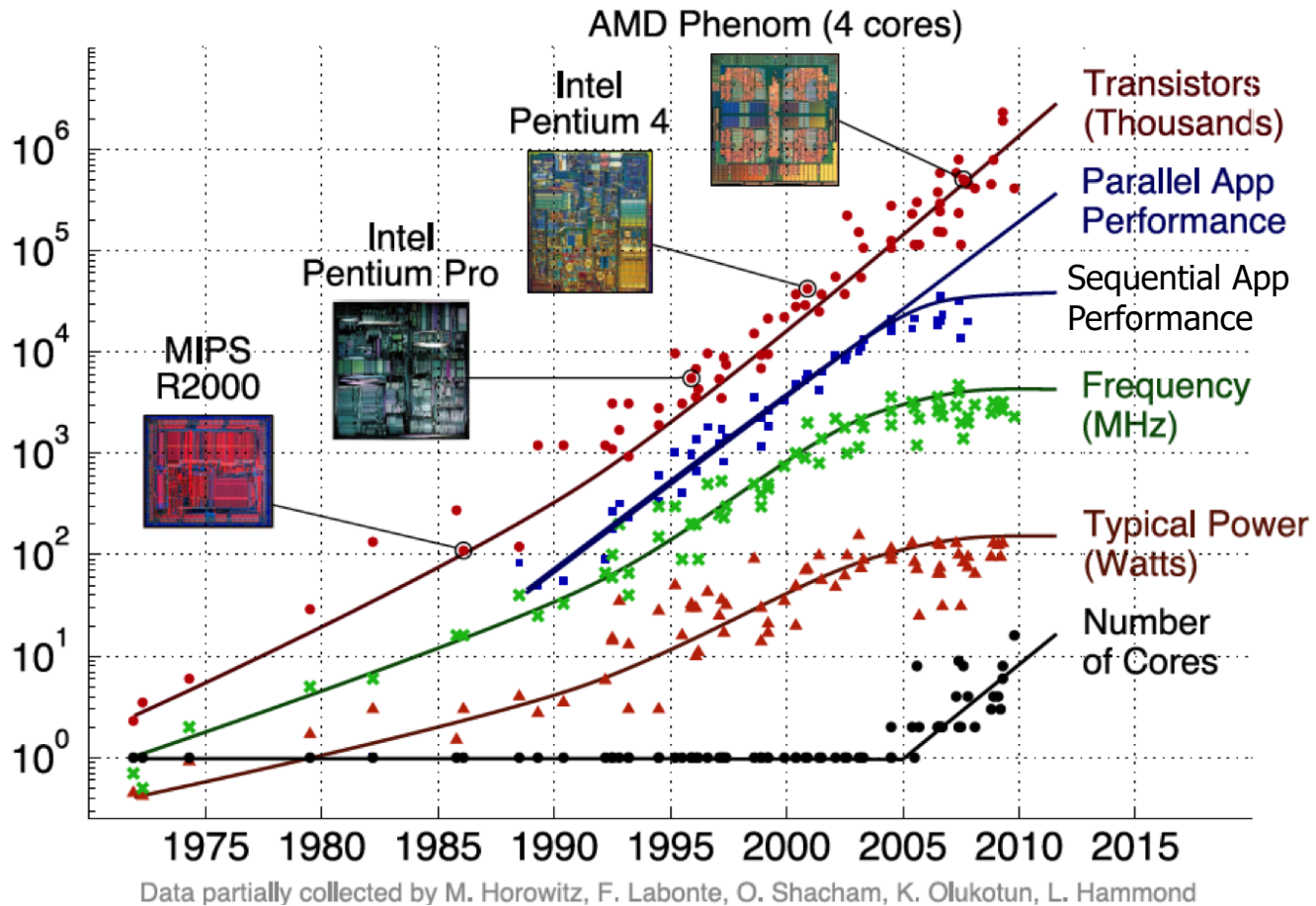# 300,000× Increase in Training Compute



Graph reproduced from openai.com/blog/ai-and-compute/

# Transistor Scaling will Reach Limits (eventually)



AMD Phenom (4 cores)

Intel Pentium 4

Intel Pentium Pro

MIPS R2000

Transistors (Thousands)

Parallel App Performance

Sequential App Performance

Frequency (MHz)

Typical Power (Watts)

Number of Cores

Data partially collected by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond



**Cerebras WSE-2**
46,225mm² Silicon
2.6 Trillion transistors

**Largest GPU**
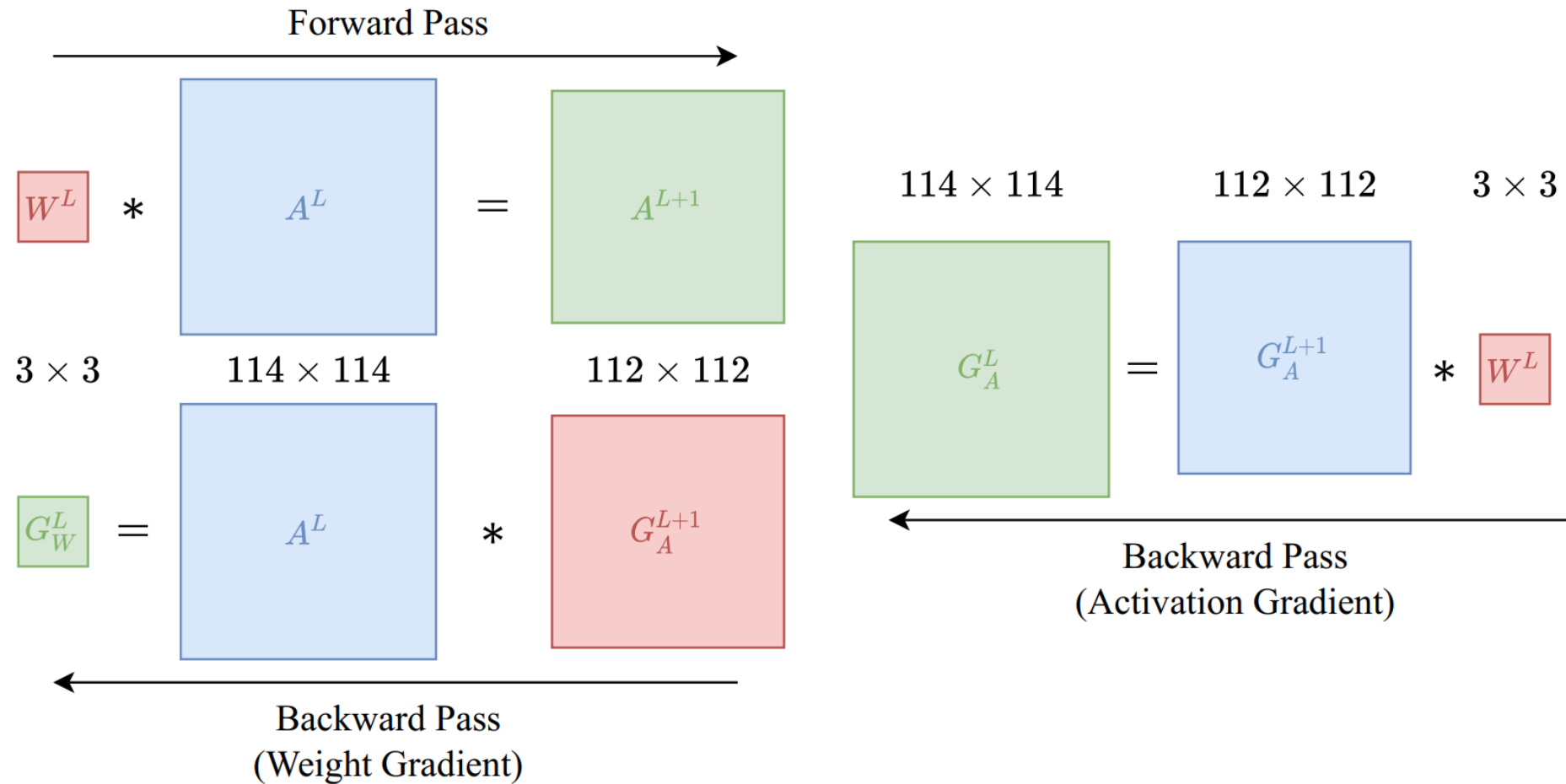826mm² Silicon
54.2 Billion transistors

[image: cerebras.ai]

# How to improve training speed?

- Exploit data parallelism
  - Each node has complete model, handles subset of all training data, accumulate gradients across nodes (size of model limited by node memory capacity)
- Model parallelism
  - Split model across compute nodes  (e.g., AlexNet, Tesla Dojo)
- Reduce computations
  - Reduce number of iterations (e.g., batch normalization, ADAM, etc…)
  - Reduce computation per iteration (e.g., stochastic depth, sparsity)
- ML hardware accelerators
  - Inference, training or both
  - Datacenter: GPU, Google TPU, Cerebras WSE, Graphcore IPU, Huawei DaVinci
  - Edge: Apple Neural Engine, Samsung NPU, Tesla FSD SoC, …

# Example: ResNet Convolution



Forward Pass

$$W^L \ * \ A^L \ = \ A^{L+1}$$

$3 \times 3$    $114 \times 114$    $112 \times 112$

$$G_W^L \ = \ A^L \ * \ G_A^{L+1}$$

Backward Pass
(Weight Gradient)

$114 \times 114$    $112 \times 112$    $3 \times 3$

$$G_A^L \ = \ G_A^{L+1} \ * \ W^L$$

Backward Pass
(Activation Gradient)

# Convolutional Neural Network Training

Activation Gradient $(G_A^{L+1})$

Activation $(A^L)$

Weight Gradient $(G_W^L)$

| 2 | -3 |
|---|----|
| 0 | 0 |

\*

| 1 | 0 | 0 |
|---|----|---|
| 0 | -1 | 2 |
| 0 | 0 | 3 |

=

| | |
|---|---|
| | |

# Convolutional Neural Network Training

$$G_A^{L+1}$$

| 2 | -3 |
|---|---|
| 0 | 0 |

$*$

$$A^L$$

| 2× 1 | -3× 0 | 0 |
|---|---|---|
| 0× 0 | 0× -1 | 2 |
| 0 | 0 | 3 |

$=$

$$G_W^L$$

| | |
|---|---|
| | |

# Convolutional Neural Network Training

$G_A^{L+1}$

$$\begin{array}{|c|c|} \hline 2 & -3 \\ \hline 0 & 0 \\ \hline \end{array}$$

$*$

$A^L$

$$\begin{array}{|c|c|c|} \hline 1 & 0 & 0 \\ \hline 0 & -1 & 2 \\ \hline 0 & 0 & 3 \\ \hline \end{array}$$

$=$

$G_W^L$

$$\begin{array}{|c|c|} \hline 2 & \\ \hline & \\ \hline \end{array}$$

# Convolutional Neural Network Training

$$G_A^{L+1}$$

| 2 | -3 |
|---|---|
| 0 | 0 |

$$*$$

$$A^L$$

| 1 | 0 | 0 |
|---|---|---|
| 0 | -1 | 2 |
| 0 | 0 | 3 |

$$=$$

$$G_W^L$$

| 2 | |
|---|---|
| | |

# Convolutional Neural Network Training

$$G_A^{L+1} \qquad A^L \qquad G_W^L$$

# Convolutional Neural Network Training

$G_A^{L+1}$        $A^L$        $G_W^L$

| 2 | -3 |
|---|----|
| 0 | 0 |

$*$

| 1 | 0 | 0 |
|---|----|---|
| 0 | -1 | 2 |
| 0 | 0 | 3 |

$=$

| 2 | 0 |
|---|---|
|   |   |

# Convolutional Neural Network Training

$$G_A^{L+1} \qquad\qquad A^L \qquad\qquad G_W^L$$

| | |
|---|---|
| 2 | -3 |
| 0 | 0 |

$*$

| | | |
|---|---|---|
| 1 | 0 | 0 |
| 0 | -1 | 2 |
| 0 | 0 | 3 |

$=$

| | |
|---|---|
| 2 | 0 |
| | |

# Convolutional Neural Network Training

$$G_A^{L+1}$$

| 2 | -3 |
|---|---|
| 0 | 0 |

**\***

$$A^L$$

| 1 | 0 | 0 |
|---|---|---|
| 2× 0 | -3× -1 | 2 |
| 0× 0 | 0× 0 | 3 |

**=**

$$G_W^L$$

| 2 | 0 |
|---|---|
|   |   |

# Convolutional Neural Network Training

$G_A^{L+1}$

$A^L$

$G_W^L$

| 2 | -3 |
|---|----|
| 0 | 0 |

$*$

| 1 | 0 | 0 |
|---|----|---|
| 0 | -1 | 2 |
| 0 | 0 | 3 |

$=$

| 2 | 0 |
|---|---|
| 3 | |

# Convolutional Neural Network Training

$$G_A^{L+1}$$

| 2 | -3 |
|---|---|
| 0 | 0 |

$*$

$$A^L$$

| 1 | 0 | 0 |
|---|---|---|
| 0 | -1 | 2 |
| 0 | 0 | 3 |

$=$

$$G_W^L$$

| 2 | 0 |
|---|---|
| 3 | |

# Convolutional Neural Network Training

$G_A^{L+1}$

$A^L$

$G_W^L$

| 2 | -3 |
|---|----|
| 0 | 0 |

$*$

| 1 | 0 | 0 |
|---|---|---|
| 0 | 2× -1 | -3× 2 |
| 0 | 0× 0 | 0× 3 |

$=$

| 2 | 0 |
|---|---|
| 3 | |

| 2 | × | -1 | = | -2 |
|---|---|----|---|----|

| -3 | × | 2 | = | -6 |
|----|---|---|---|----|

# Convolutional Neural Network Training

$$G_A^{L+1} \qquad A^L \qquad G_W^L$$

| | |
|---|---|
| 2 | -3 |
| 0 | 0 |

$*$

| | | |
|---|---|---|
| 1 | 0 | 0 |
| 0 | -1 | 2 |
| 0 | 0 | 3 |

$=$

| | |
|---|---|
| 2 | 0 |
| 3 | -8 |

# Encouraging sparsity
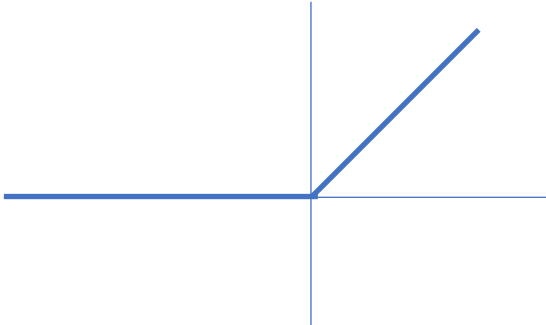
- Convolutions essentially perform many dot-products:

$$y = [10 \ 0 \ 3 \ 2] \ [1 \ 5 \ 0 \ 0]^T = 10*1 + 0*5 + 3*0 + 2*0 = 10*1 = 10$$

- Multiplications by zero can be skipped.

- Exploiting this can reduce computations and/or model size.

- Much work on this in past ~5 years.  Extensions to fully connected layers, RNNs, transformers, etc…
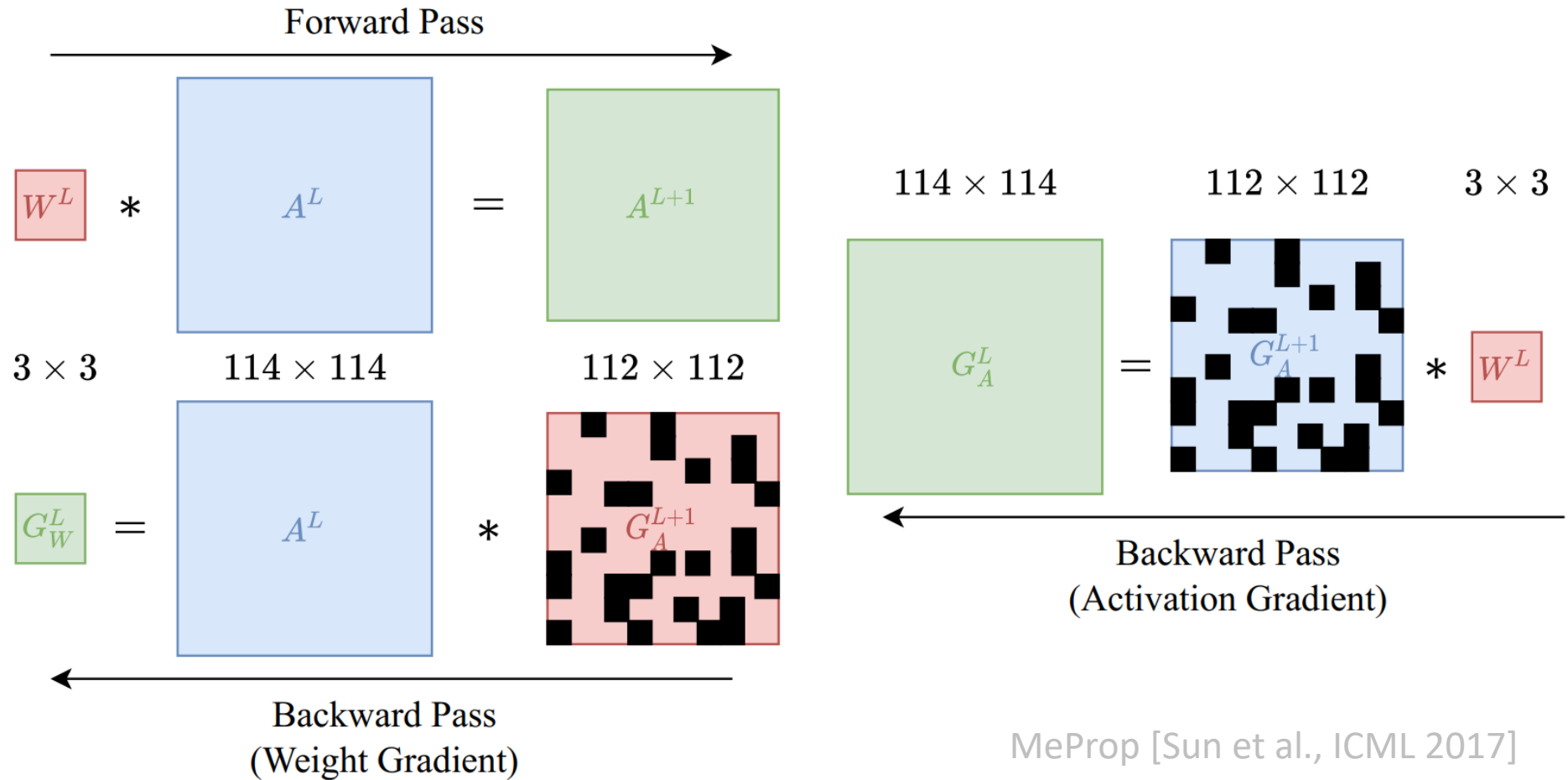
# Weight Pruning



Forward Pass

$V^L$ * $A^L$ = $A^{L+1}$

$3 \times 3$    $114 \times 114$    $112 \times 112$

$G_W^L$ = $A^L$ * $G_A^{L+1}$

Backward Pass
(Weight Gradient)

$114 \times 114$    $112 \times 112$    $3 \times 3$

$G_A^L$ = $G_A^{L+1}$ * $V^L$

Backward Pass
(Activation Gradient)

[Han et al., NeurIPS 2015]

# ReLU

**Forward Pass**

$W^L$ * $A^L$ = $A^{L+1}$

$3 \times 3$  $114 \times 114$  $112 \times 112$

$G_W^L$ = $A^L$ * $G_A^{L+1}$

**Backward Pass (Weight Gradient)**

$114 \times 114$  $112 \times 112$  $3 \times 3$

$G_A^L$ = $G_A^{L+1}$ * $W^L$

**Backward Pass (Activation Gradient)**

[Nair and Hinton, ICML 2010]
[Albericio et al. ISCA 2016]

# Sparse Gradients



Forward Pass

$W^L$ * $A^L$ = $A^{L+1}$

$3 \times 3$    $114 \times 114$    $112 \times 112$

$G_W^L$ = $A^L$ * $G_A^{L+1}$

Backward Pass
(Weight Gradient)

$114 \times 114$    $112 \times 112$    $3 \times 3$

$G_A^L$ = $G_A^{L+1}$ * $W^L$

Backward Pass
(Activation Gradient)

MeProp [Sun et al., ICML 2017]

ReSprop [Goli and Aamodt, CVPR 2020]

# Sparse Weight Activation Training (SWAT)



FORWARD PASS

BACKWARD PASS

Implication: In forward pass sparsify weights (but not activations).

Implication: In backward pass sparsify weights and activations (but not gradients)

# SWAT Algorithm (highlights)

- During each training iteration:
  - Sparse weight topology (top-k) induced, which partitions weights into **active** and **non-active** sets.
  - Forward pass:
    - use active (sparse) weights and full (dense) activations to compute layer outputs
  - Backward pass:
    - use active (sparse) weights and full (dense) gradients to compute activation gradients
    - use **sparse** activations (top-k) and dense gradients to compute **dense** weight updates

- Updating weights with dense weight gradients enables topology search (avoids "lottery ticket" problem)

# Sparse Weight Activation Training (SWAT)



[Raihan and Aamodt, NeurIPS 2020]

# Comparison of unstructured SWAT with sparse learning algorithms on the ImageNet

| Methods | Weight Sparsity (%) | Activation Sparsity (%) | Top-1 Accuracy (%) | Accuracy Change (%) | Training FLOP ⬇(%) | Inference FLOP ⬇(%) | Model Compression(x) |
|---|---|---|---|---|---|---|---|
| SET | 80 | - | 73.4 | -3.4 | 58.1 | 73.0 | 3.4 |
| | 90 | - | 71.3 | -5.5 | 63.8 | 82.1 | 5.0 |
| DSR | 80 | - | 74.1 | -2.7 | 51.6 | 59.4 | 3.4 |
| | 90 | - | 71.9 | -4.9 | 58.9 | 70.7 | 5.0 |
| SNFS | 80 | - | 74.9 | -2.1 | 45.8 | 43.3 | 5.0 |
| | 90 | - | 72.9 | -4.1 | 57.6 | 59.7 | 10.0 |
| RigL | 80 | - | 74.6 | -2.2 | 67.2 | 80.0 | 5.0 |
| | 90 | - | 72.0 | -4.8 | 74.1 | 90.0 | 10.0 |
| DST | 80.4 | - | 74.0 | -2.8 | 67.1 | 84.9 | 5.0 |
| | 90.1 | - | 72.8 | -4.0 | 75.8 | 91.3 | 10.0 |
| SWAT-U | 80 | 80 | 75.2 | -1.6 | 76.1 | 77.7 | 5.0 |
| | 90 | 90 | 72.1 | -4.7 | 85.6 | 87.4 | 10.0 |
| SWAT-U | 80 | 80 | 75.2 | -1.6 | 76.1 | 77.7 | 5.0 |
| | 90 | 90 | 72.1 | -4.7 | 85.6 | 87.4 | 10.0 |
| SWAT-U | 80 | 80 | 75.2 | -1.6 | 76.1 | 77.7 | 5.0 |
| | 90 | 90 | 72.1 | -4.7 | 85.6 | 87.4 | 10.0 |

# ReSprop: Reuse Sparsified Backpropagation



[Goli and Aamodt, CVPR 2020]

# ReSprop



Iteration i-1 — Gradient of Output

Iteration i — Gradient of Output

Iteration i — Hybrid Output Gradient (HG)
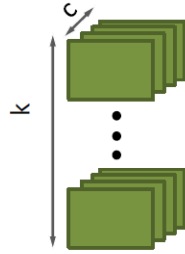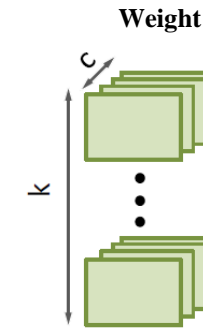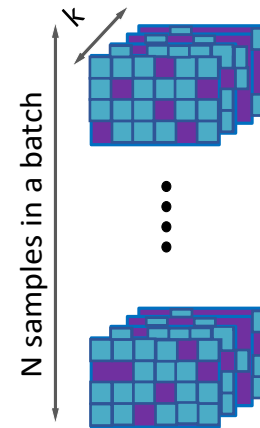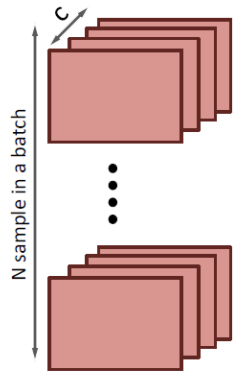
Threshold (Adaptive)

Figure 2. HG and meProp angles for different reuse percentages and sparsities, respectively. The angle is calculated by finding the average angle of all layers while training ResNet-18 on CIFAR-10 for 100 iterations (batch size=128).
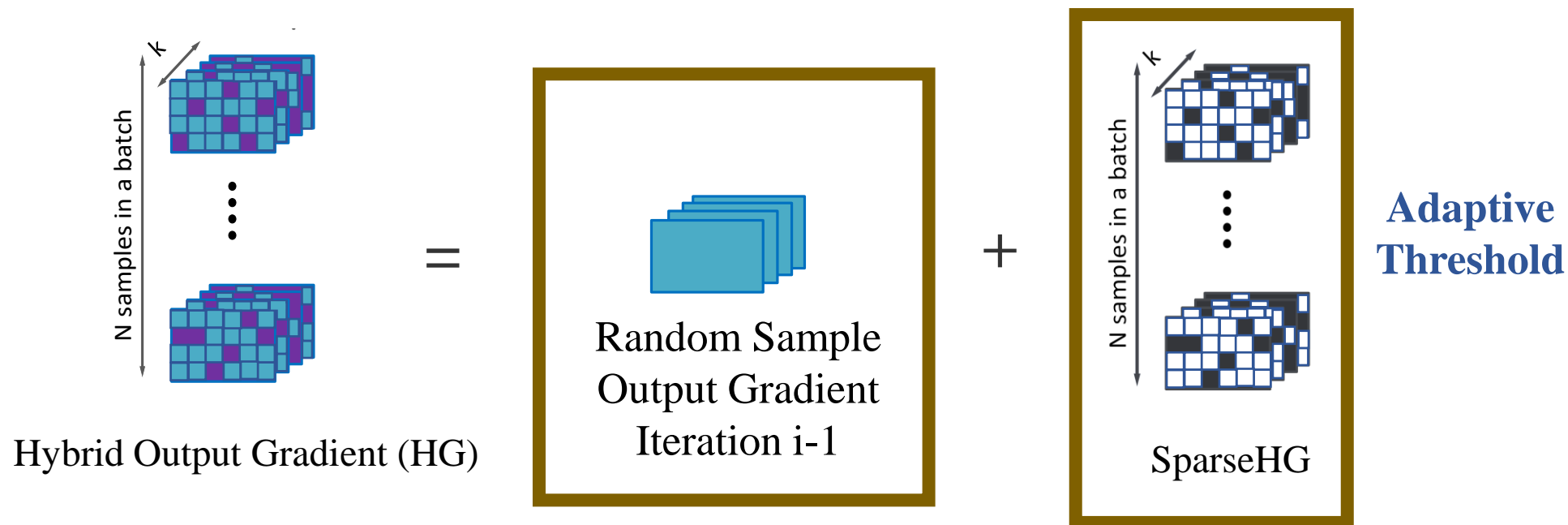
**Gradient of Weight** = Hybrid Output Gradient $\otimes$ Input

**Gradient of Input** = Hybrid Output Gradient $\otimes$ Weight

**Iteration i:**



Hybrid Output Gradient (HG) = Random Sample Output Gradient Iteration i-1 + SparseHG

**Adaptive Threshold**

1) Reuse

2) Use ==sparse== gradients for expensive convolutions in backward pass

**Fast**

| RS | Algorithm | ResNet34 | WRN-50-2 | VGG16 |
|---|---|---|---|---|
| 50% | ReSprop | 73.08 | 78.69 | 70.09 |
| | W-ReSprop | 73.21 | 78.81 | 70.41 |
| 70% | ReSprop | 67.12 | 73.34 | 68.73 |
| | W-ReSprop | 72.73 | 78.25 | 70.01 |
| **90%** | ReSprop | 63.78 | 67.72 | 60.76 |
| | **W-ReSprop** | **72.44** | **77.93** | **69.46** |
| Baseline | | 73.34 | 78.88 | 70.50 |

Table 4. Top 1 validation accuracy of ReSprop and W-ReSprop algorithms at different reuse-sparsity constraints on the ImageNet dataset.
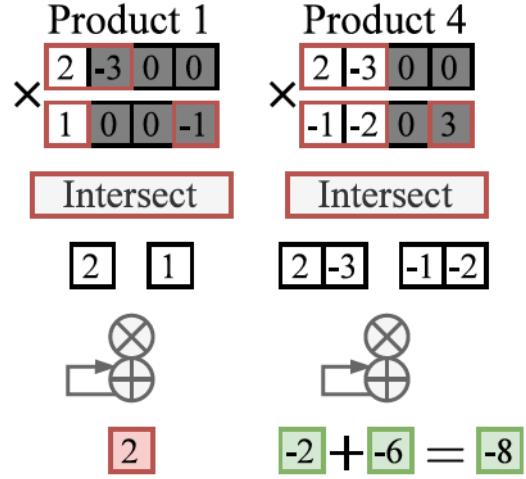
# Hardware Support for Sparsity
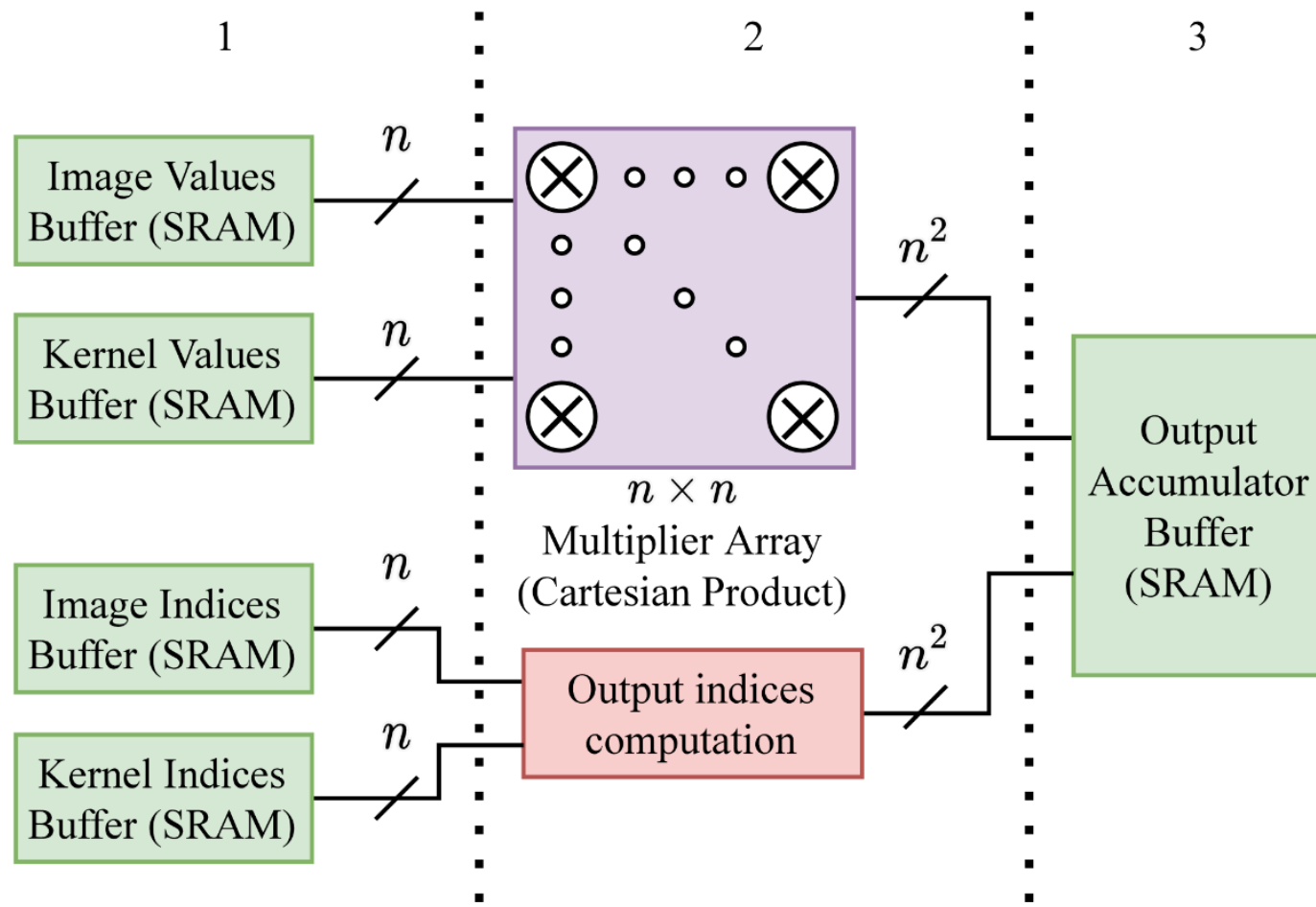


Figure 2: Convolution accelerator classes, showing zero products and Redundant Cartesian Products (RCPs). a) An example convolution of a $2 \times 2$ kernel and $3 \times 3$ image, b) Inner product/dot product, c) Intersection/streaming and, d) Outer-product.

[Lew et al., ISCA 2022]

# Exploiting Two-Sided Dynamic Sparsity: SCNN



Parashar et al. (2017)

# Convolutional Neural Network Training

$$G_A^{L+1} \qquad A^L \qquad G_W^L$$



$$
\begin{array}{|c|c|}
\hline
2 & -3 \\
\hline
0 & 0 \\
\hline
\end{array}
\quad * \quad
\begin{array}{|c|c|c|}
\hline
1 & 0 & 0 \\
\hline
0 & -1 & 2 \\
\hline
0 & 0 & 3 \\
\hline
\end{array}
\quad = \quad
\begin{array}{|c|c|}
\hline
2 & 0 \\
\hline
3 & -8 \\
\hline
\end{array}
$$

# Cartesian (Outer) Product

$$G_A^{L+1} \qquad A^L \qquad G_W^L$$

$$
\begin{array}{|c|c|}
\hline
2 & -3 \\
\hline
0 & 0 \\
\hline
\end{array}
\; * \;
\begin{array}{|c|c|c|}
\hline
1 & 0 & 0 \\
\hline
0 & -1 & 2 \\
\hline
0 & 0 & 3 \\
\hline
\end{array}
\; = \;
\begin{array}{|c|c|}
\hline
2 & 0 \\
\hline
3 & -8 \\
\hline
\end{array}
$$

$$
\begin{array}{|c|}
\hline
2 \\
\hline
-3 \\
\hline
\end{array}
\qquad
\begin{array}{|c|c|c|c|}
\hline
1 & -1 & 2 & 3 \\
\hline
\end{array}
$$

# Cartesian Product

$$G_A^{L+1} \qquad A^L \qquad = \qquad G_W^L$$



$$G_A^{L+1} * A^L = G_W^L$$

| | 1 | 0 | 0 |
|---|---|---|---|
| | 0 | -1 | 2 |
| | 0 | 0 | 3 |

$G_A^{L+1}$:

| 2 | -3 |
|---|---|
| 0 | 0 |

$G_W^L$:

| 2 | 0 |
|---|---|
| 3 | -8 |

$A^L$ / $G_A^{L+1}$ table:

| | 1 | -1 | 2 | 3 |
|---|---|---|---|---|
| 2 | 2 | -2 | 4 | 6 |
| -3 | -3 | 3 | -6 | -9 |

# Cartesian Product

$G_A^{L+1}$   $A^L$   $G_W^L$

$$\begin{bmatrix} 2 & -3 \\ 0 & 0 \end{bmatrix} * \begin{bmatrix} 1 & 0 & 0 \\ 0 & -1 & 2 \\ 0 & 0 & 3 \end{bmatrix} = \begin{bmatrix} 2 & 0 \\ 3 & -8 \end{bmatrix}$$

$A^L$

| | 1 | -1 | 2 | 3 |
|---|---|---|---|---|
| 2 | 2 | -2 | 4 | 6 |
| -3 | -3 | 3 | -6 | -9 |

$G_A^{L+1}$

# Cartesian Product

$G_A^{L+1}$ * $A^L$ = $G_W^L$

$$\begin{bmatrix} 2 & -3 \\ 0 & 0 \end{bmatrix} * \begin{bmatrix} 1 & 0 & 0 \\ 0 & -1 & 2 \\ 0 & 0 & 3 \end{bmatrix} = \begin{bmatrix} 2 & 0 \\ 3 & -8 \end{bmatrix}$$

$A^L$

$G_A^{L+1}$

| | 1 | -1 | 2 | 3 |
|---|---|---|---|---|
| 2 | 2 | -2 | 4 | 6 |
| -3 | -3 | 3 | -6 | -9 |

$$2 \times -1 = -2$$

$$-3 \times 2 = -6$$

35

# Cartesian Product

$G_A^{L+1}$     $A^L$     $G_W^L$

| 2 | -3 |
|---|----|
| 0 | 0  |

*

| 1 | 0  | 0 |
|---|----|---|
| 0 | -1 | 2 |
| 0 | 0  | 3 |

=

| 2 | 0  |
|---|----|
| 3 | -8 |

$A^L$

$G_A^{L+1}$

|   | 1 | -1 | 2 | 3 |
|---|---|----|---|---|
| 2 | 2 | -2 | 4 | 6 |
| -3 | -3 | 3 | -6 | -9 |

# The Problem: Redundant Cartesian Products

$G_A^{L+1}$

$$\begin{bmatrix} 2 & -3 \\ 0 & 0 \end{bmatrix}$$

$*$

$A^L$

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & -1 & 2 \\ 0 & 0 & 3 \end{bmatrix}$$

$=$

$G_W^L$

$$\begin{bmatrix} 2 & 0 \\ 3 & -8 \end{bmatrix}$$

$A^L$

| | 1 | -1 | 2 | 3 |
|---|---|---|---|---|

$G_A^{L+1}$

| | 1 | -1 | 2 | 3 |
|---|---|---|---|---|
| 2 | 2 | -2 | 4 | 6 |
| -3 | -3 | 3 | -6 | -9 |

**Redundant Cartesian Products (RCPs)** dominate during **training**

$$\text{Outer-product Efficiency} = \frac{H_{out} \times W_{out}}{H \times W}$$

Forward Pass →

96.52%  $W^L$  *  $A^L$  =  $A^{L+1}$

$3 \times 3$   $114 \times 114$   $112 \times 112$

$114 \times 114$   $112 \times 112$   $3 \times 3$

$G_A^L$  =  $G_A^{L+1}$  *  $W^L$

0.07%  $G_W^L$  =  $A^L$  *  $G_A^{L+1}$

← Backward Pass (Weight Gradient)

← Backward Pass (Activation Gradient)

[Lew et al., ISCA 2022]

# Redundant Cartesian Products



a)

No prod.

RCPs

b)

c)

d)

x →

y ↓

Image

s →

r ↓

Kernel

Output

No RCPs

$$\frac{y - r}{stride} < 0$$

$$\frac{x - s}{stride} < 0$$

$$\frac{y - r}{stride} \geq H_{out}$$

$$\frac{x - s}{stride} \geq W_{out}$$

[Lew et al., ISCA 2022]

# Anticipator Accelerator (ANT)

# Mapping onto a Multiplier Array

$A^L$

$G_A^{L+1}$

| 1 | -1 | 2 | 3 |

| 2 | 2 | -2 | 4 | 6 |
| -3 | -3 | 3 | -6 | -9 |

Cycle 1

| 1 | -1 |

| 2 | 2 | -2 |

# Mapping onto a Multiplier Array

$A^L$

$G_A^{L+1}$

Cycle 2

# Mapping onto a Multiplier Array: Skipping RCPs

$A^L$

$G_A^{L+1}$

Cycle 3

output coordinate = $(\frac{x - s0}{stride}, \frac{0 \cdot y - 0 \cdot r}{stride})$  $((2, 0))$

# Mapping onto a Multiplier Array: Skipping RCPs

$A^L$

$G_A^{L+1}$

Cycle 3

output coordinate = $(\frac{x \cdot s \cdot 0}{stride}, \frac{0 \cdot y \cdot 0 \cdot r}{stride})$  $((3, 0))$

# Mapping onto a Multiplier Array

$A^L$

1 | -1 | 2 | 3

$G_A^{L+1}$

2 | 2 | -2 | 4 | 6
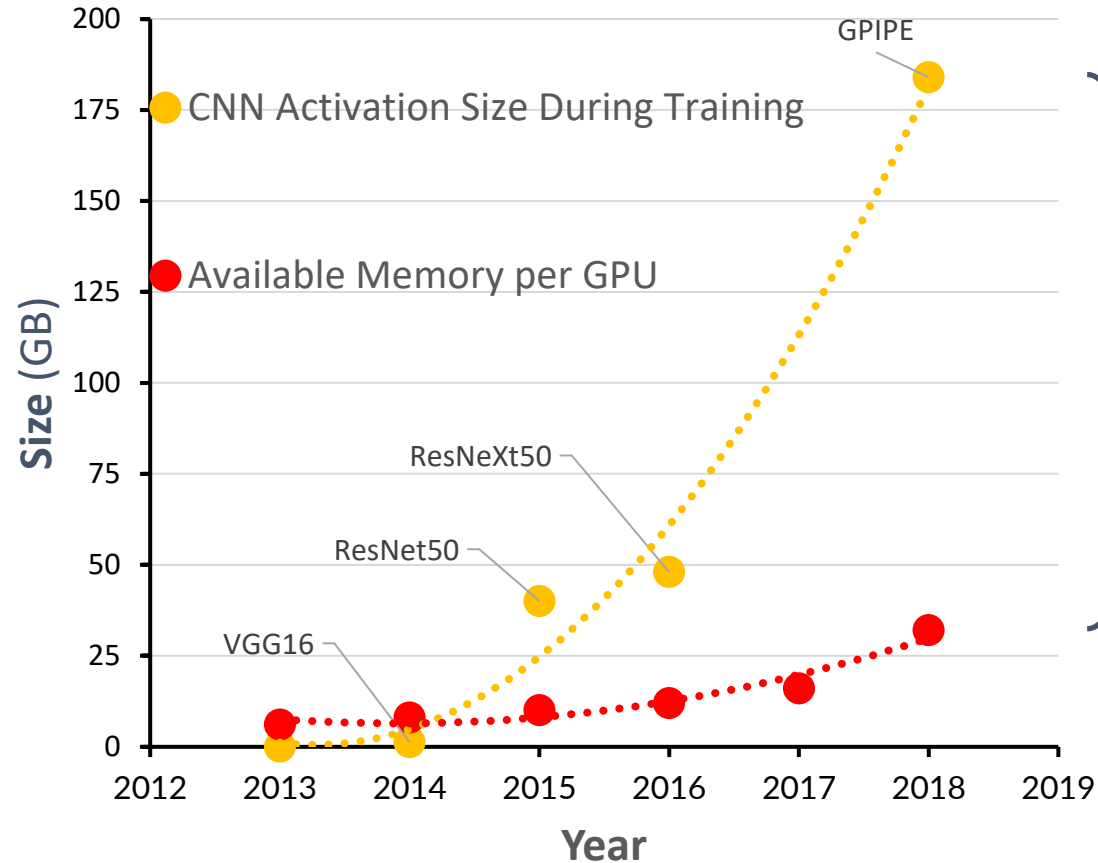
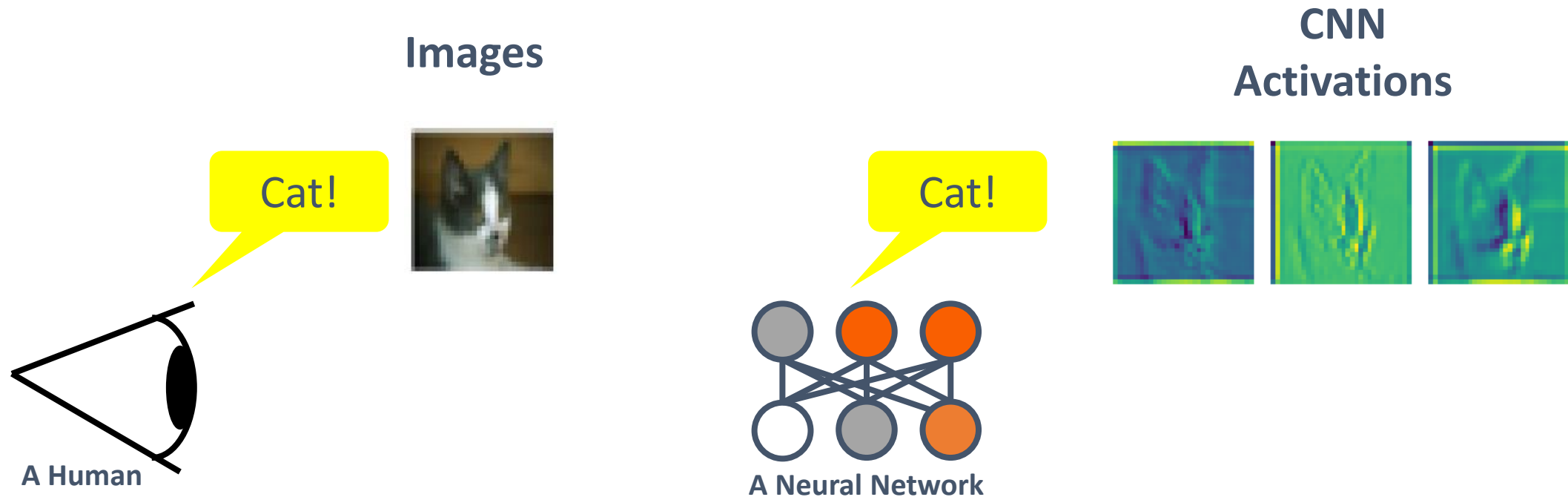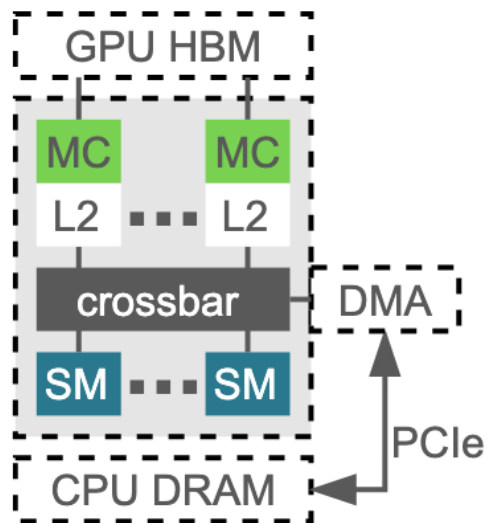-3 | -3 | 3 | -6 | -9

Cycle 3

2 | 3

-3 | -6 | -9

# Bigger Models, More Memory



[Evans et al., ISCA 2020]

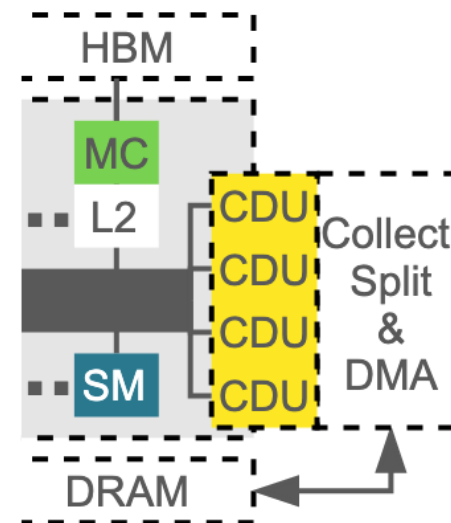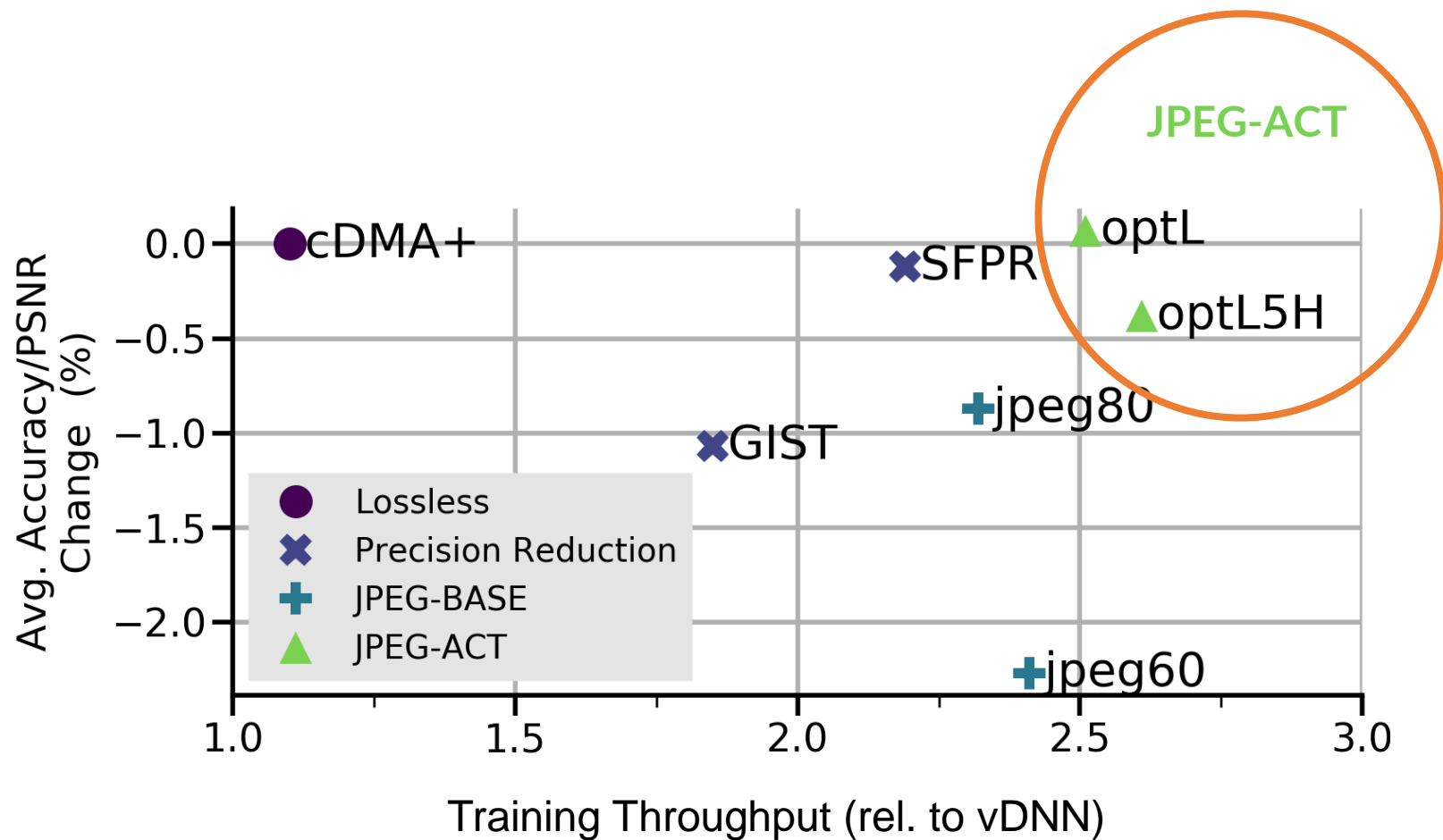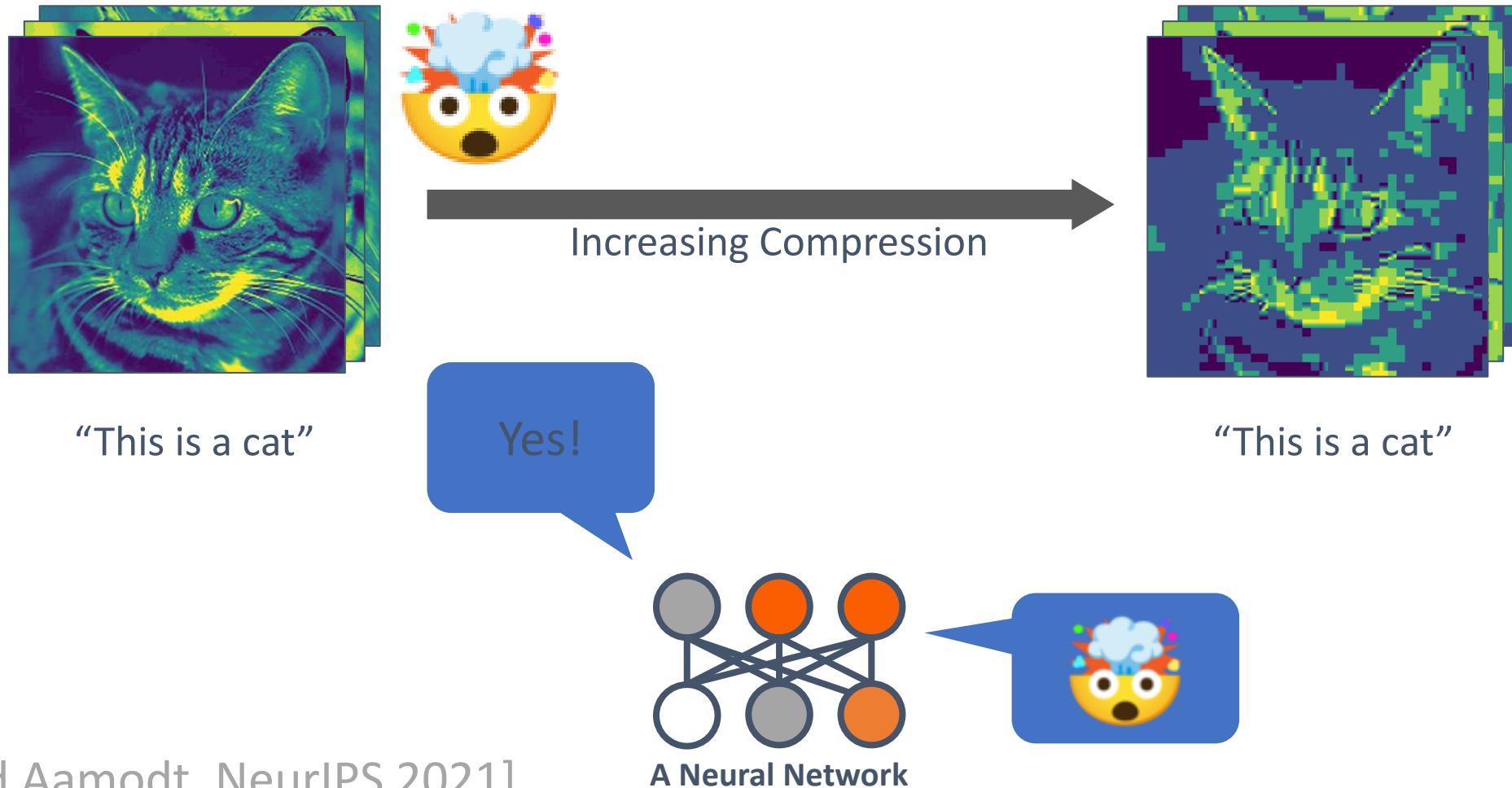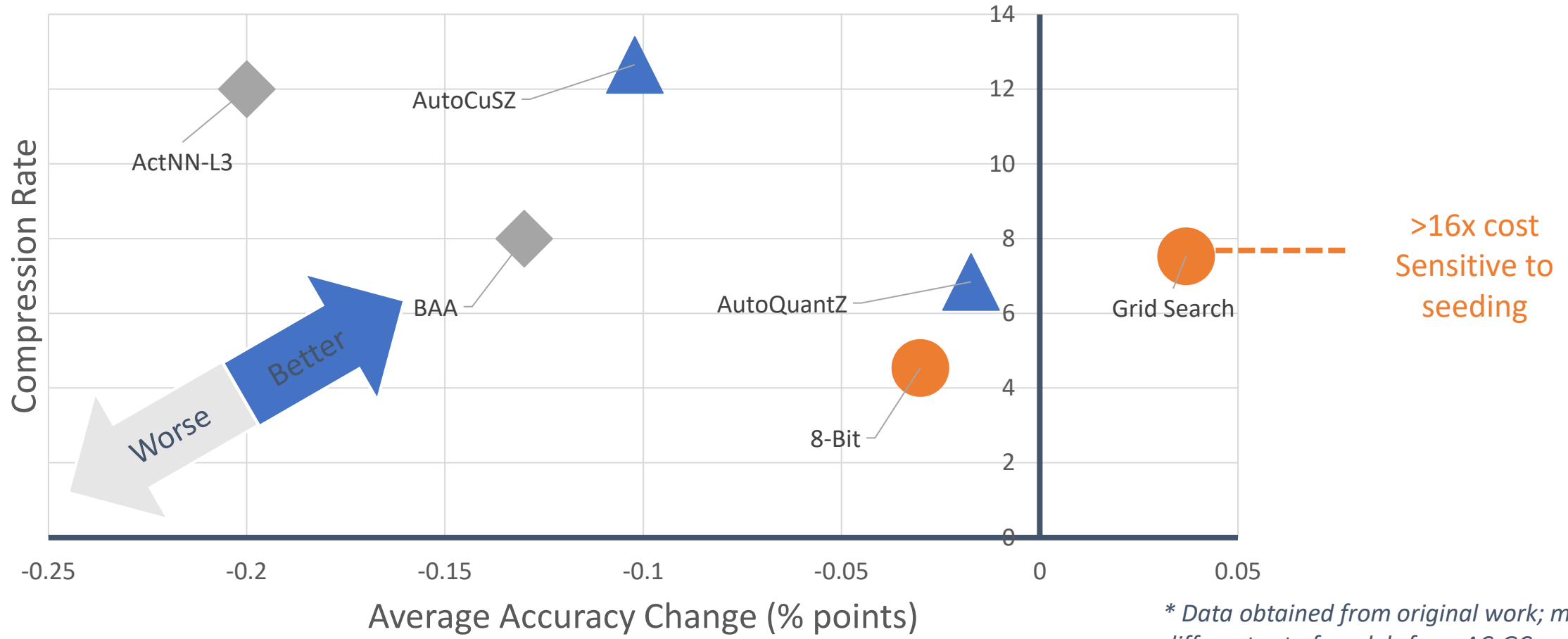# Images versus Activations

**Images**

**CNN Activations**

Cat!

Cat!

**A Human**

**A Neural Network**

[Evans et al., ISCA 2020]

Baseline (vDNN)

JPEG-ACT

# Convergence versus (Lossy) Compression

Increasing Compression

"This is a cat"

"This is a cat"

Yes!

**A Neural Network**

Average Accuracy Change (% points)

AC-GC ▲   SFPR ●   Related Works* ◆

>16x cost
Sensitive to
seeding

* Data obtained from original work; may use a
different set of models from AC-GC

(BAA) A. Chakrabarti, B. Moseley, in NeurIPS 2020
(ACTNN) J. Chen, L. Zheng, et. al, in ICML 2021

# Summary

- Obtaining greater performance for machine learning will increasingly require shifting towards specialized hardware

- ReSprop, SWAT can reduce computation demand during training by identifying computations to elide at minimal impact on accuracy leading to sparse computations

- Efficiently supporting sparse computations in hardware is a challenge, and (even more so during training since parameters are changing)

- Lossy compression can help performance (and/or increase model size) by greatly reducing memory demands.  Challenging to use during training since need to anticipate impact on validation accuracy.