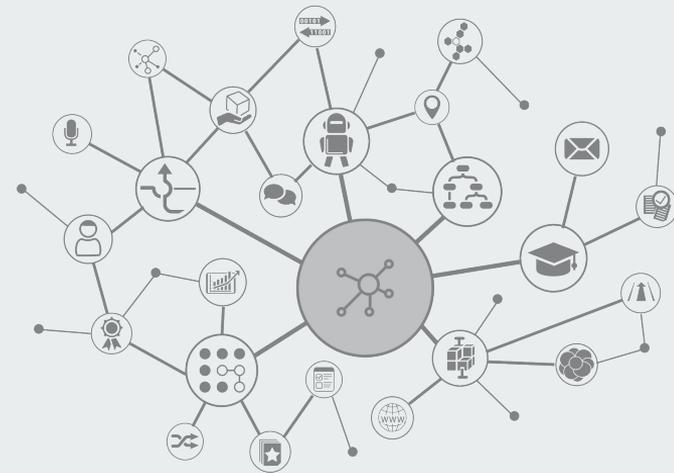




Deep Learning

Time and Memory Efficiency

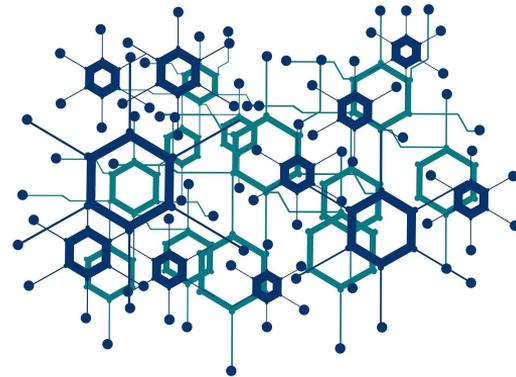


Part 1 : JAEYOUNG CHUN
Part 2 : ADIR MORGAN

APRIL 29, 2018
WEIZMANN INSTITUTE OF SCIENCE
DEPARTMENT OF COMPUTER SCIENCE



Agenda



PART 1

1. Motivation
2. Classical Algorithms / Build Blocks
3. Hardware

PART 2

4. Recent Works

Motivation

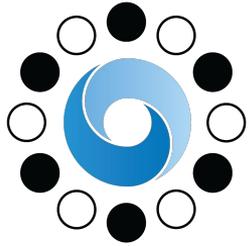
Why Seeking Efficiency?





Facebook Under Fire: How Privacy Crisis Could Change Big Data Forever

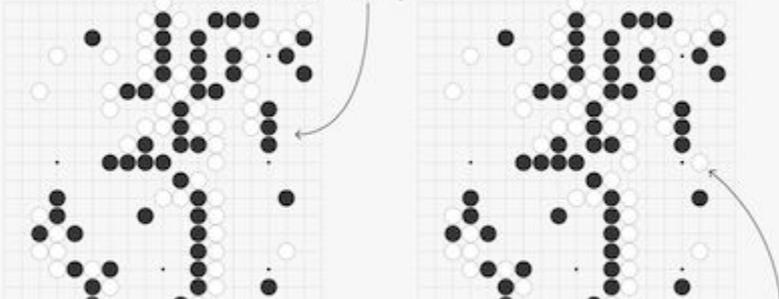
AlphaGo



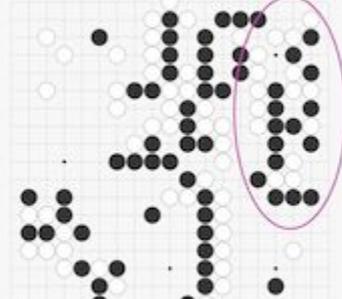
1,920 CPUs
280 GPUs
\$3,000 Electricity Bill
per game

AlphaGo vs Lee Sedol, Game 1

Lee, playing black, is in total control of this region

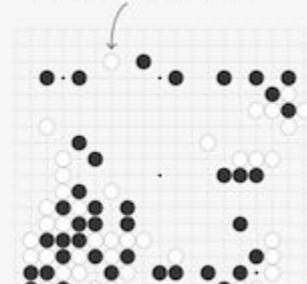


And ends up taking a significant amount of territory from Lee

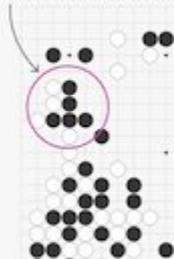


Lee Sedol vs AlphaGo, Game 2

Lee, playing white, invades a region controlled by AlphaGo



AlphaGo defies convention by the invasion altogether, instead up its defenses in another area



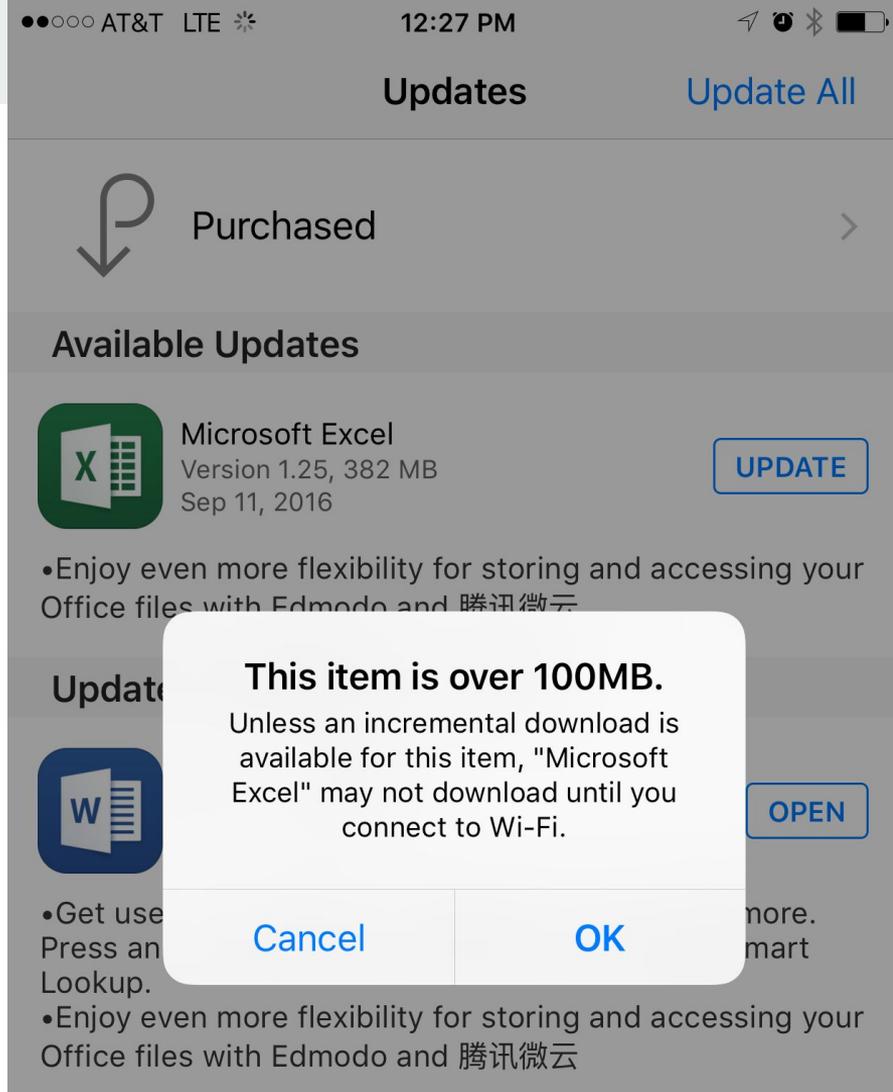


AppStore Download Restriction

To download an app over 100MB onto your mobile device,

you must connect to WiFi.

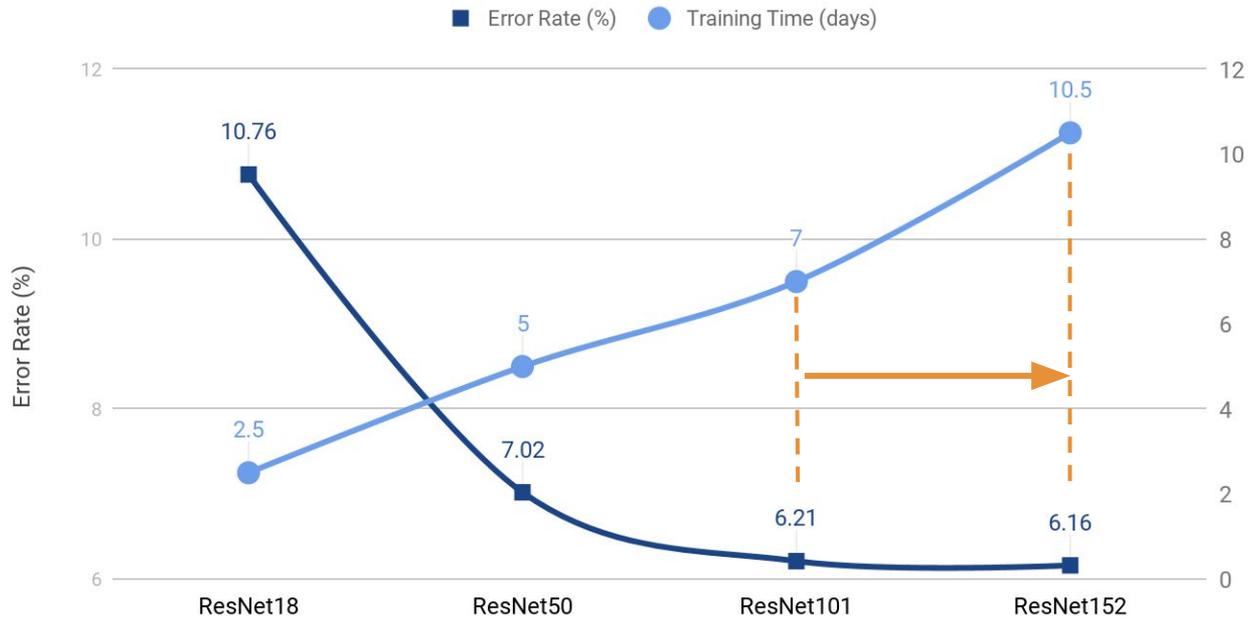
Putting this in perspective, VGG-16 has 130 million parameters (520MB).



Model Size



Speed of Training

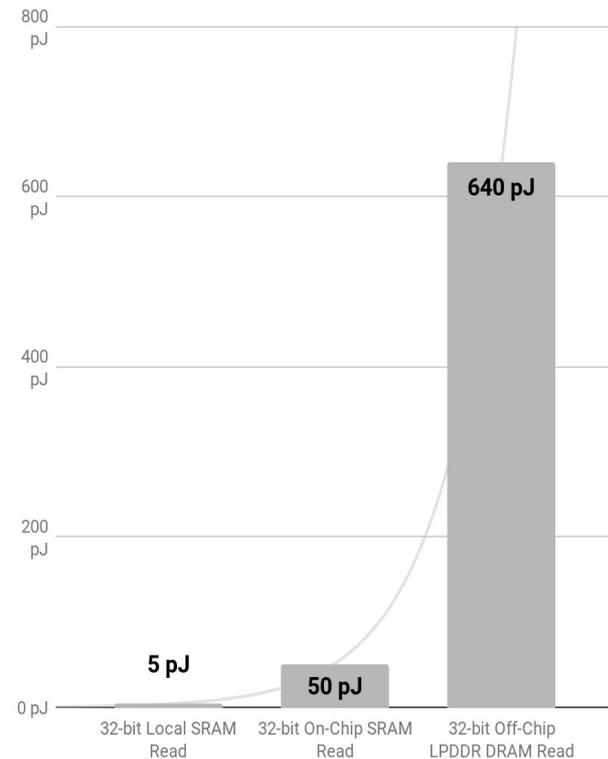
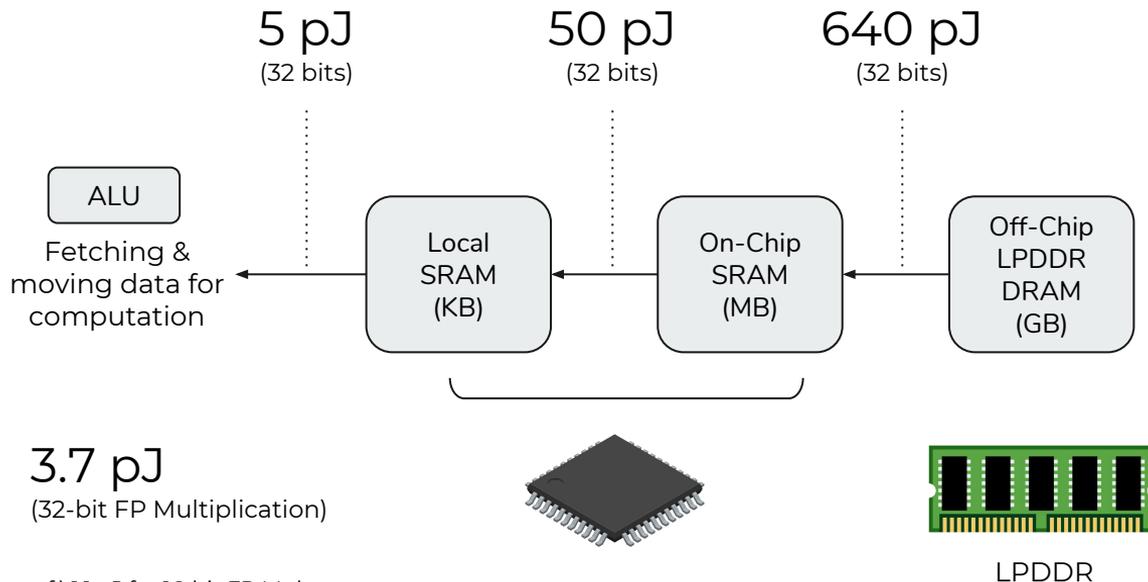


0.05 ↓
ERR %

3.5 ↑
DAYS

Energy Consumption

Cost of data movement is much more huge.
When compared, arithmetic ops is more like a noise.



c.f.) 1.1 pJ for 16-bit FP Mult
0.2 pJ for 8-bit Mult



Modules

 Efficient Inference

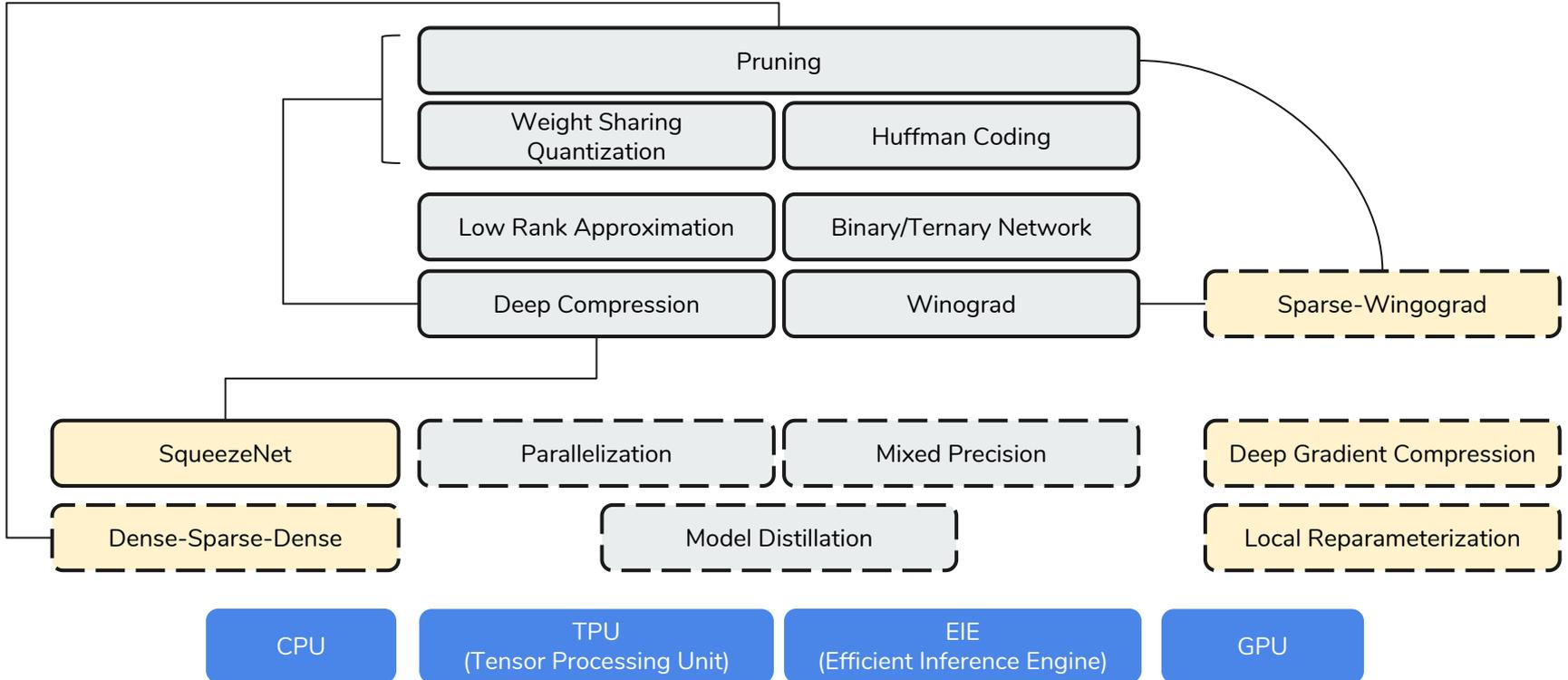
 Efficient Training

(*) no clear cut

 Algos: Classic/Building Block

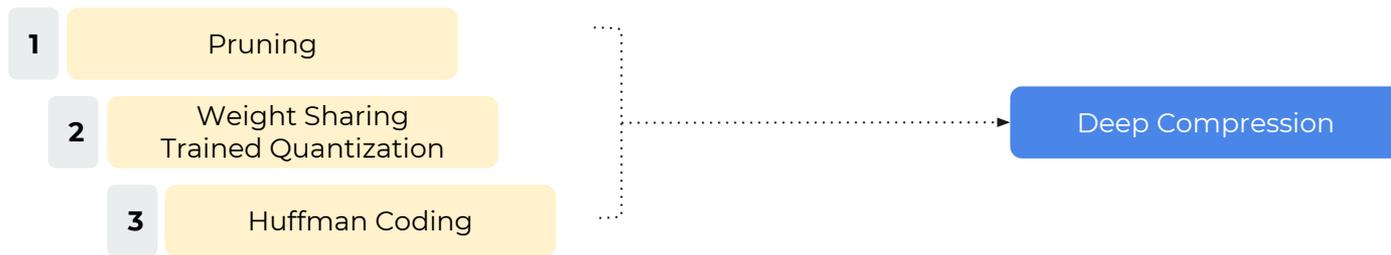
 Algos: Relatively New

 Hardware



Something to Keep in Mind...

- Losing any accuracy?
- Multiple methods interfering each other?



Understanding the underlying concepts
& getting insights, and intuitions



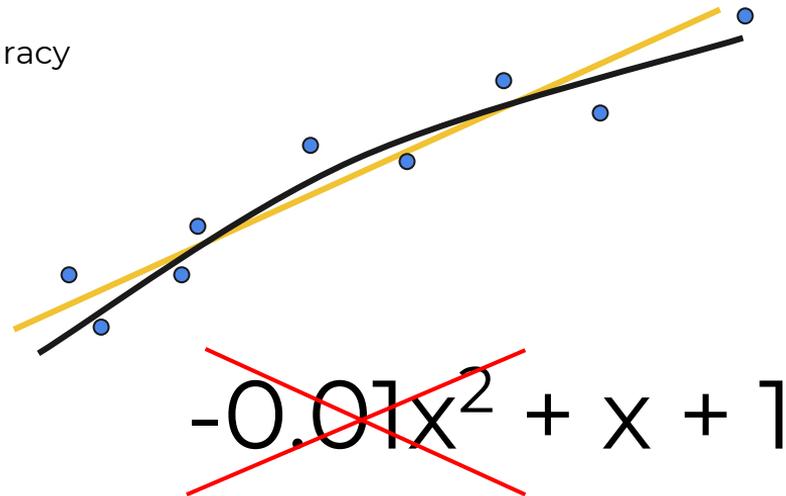
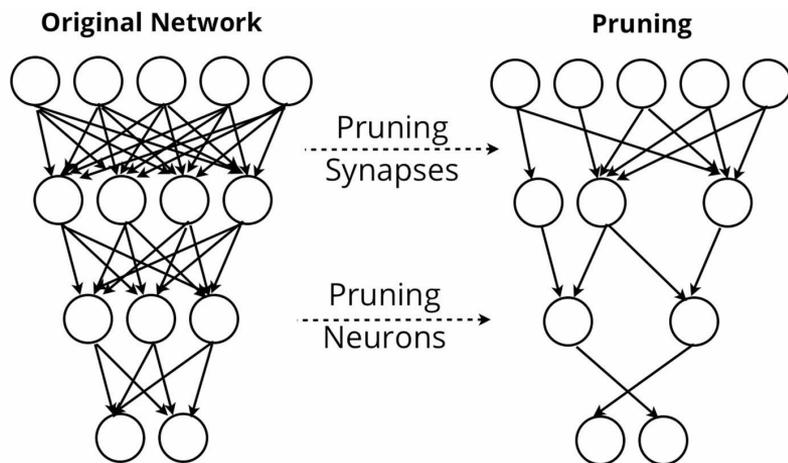
Implementation and
mathematical details

Algorithms for Efficient Inference



Pruning

Less number of parameters with almost no loss of accuracy



not just to reduce the network complexity,
but also to avoid overfitting

Pruning

Train Connectivity

Learn the connectivity via normal network training, as you would normally do.

1

Prune Connections

Prune the small-weight connections from the network.
(below a certain threshold)

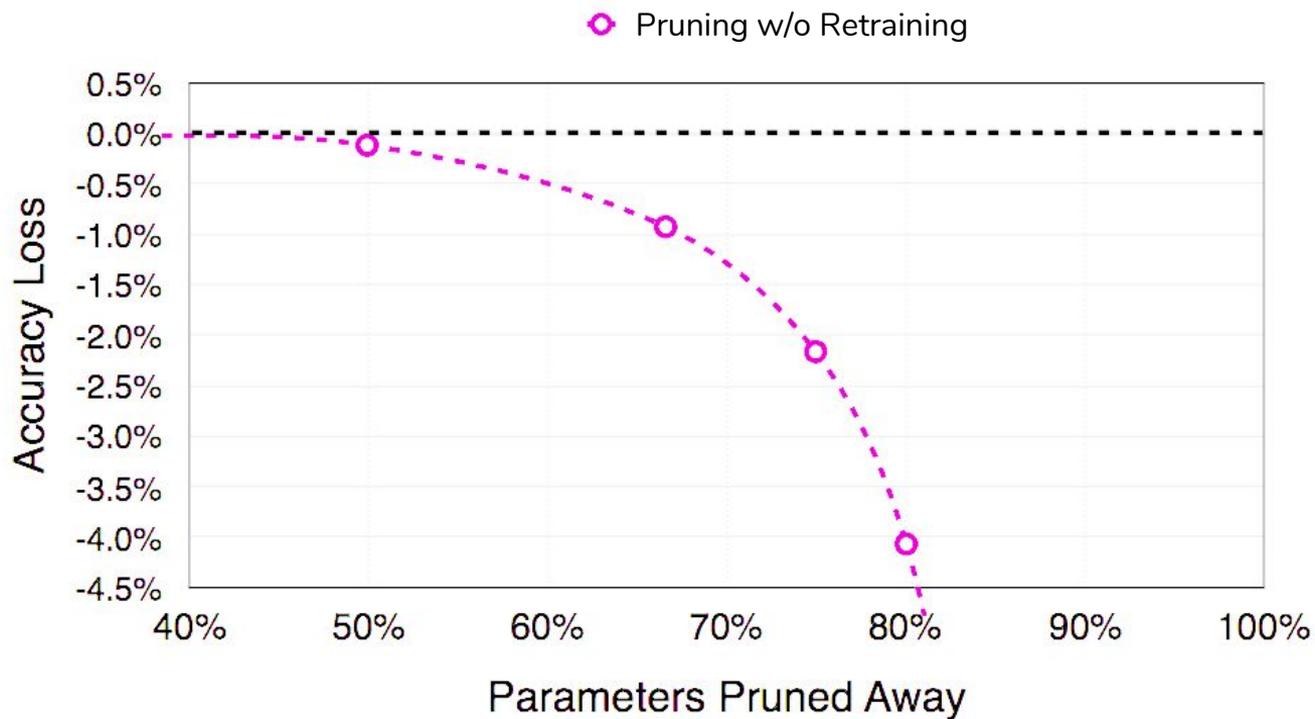
2

Retrain Weights

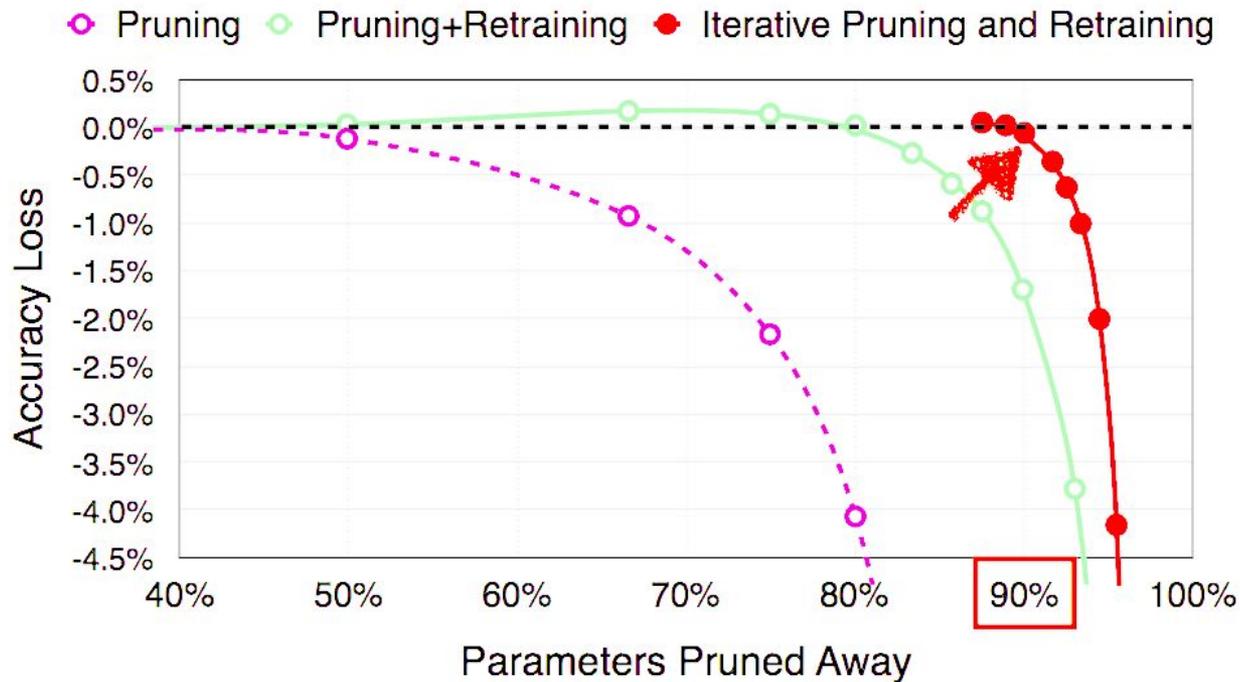
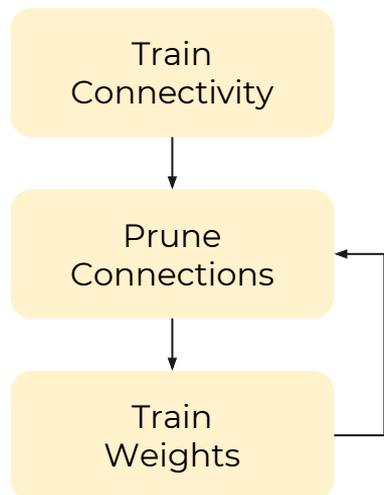
Retrain the network to learn the final weights for the remaining sparse connections.

3

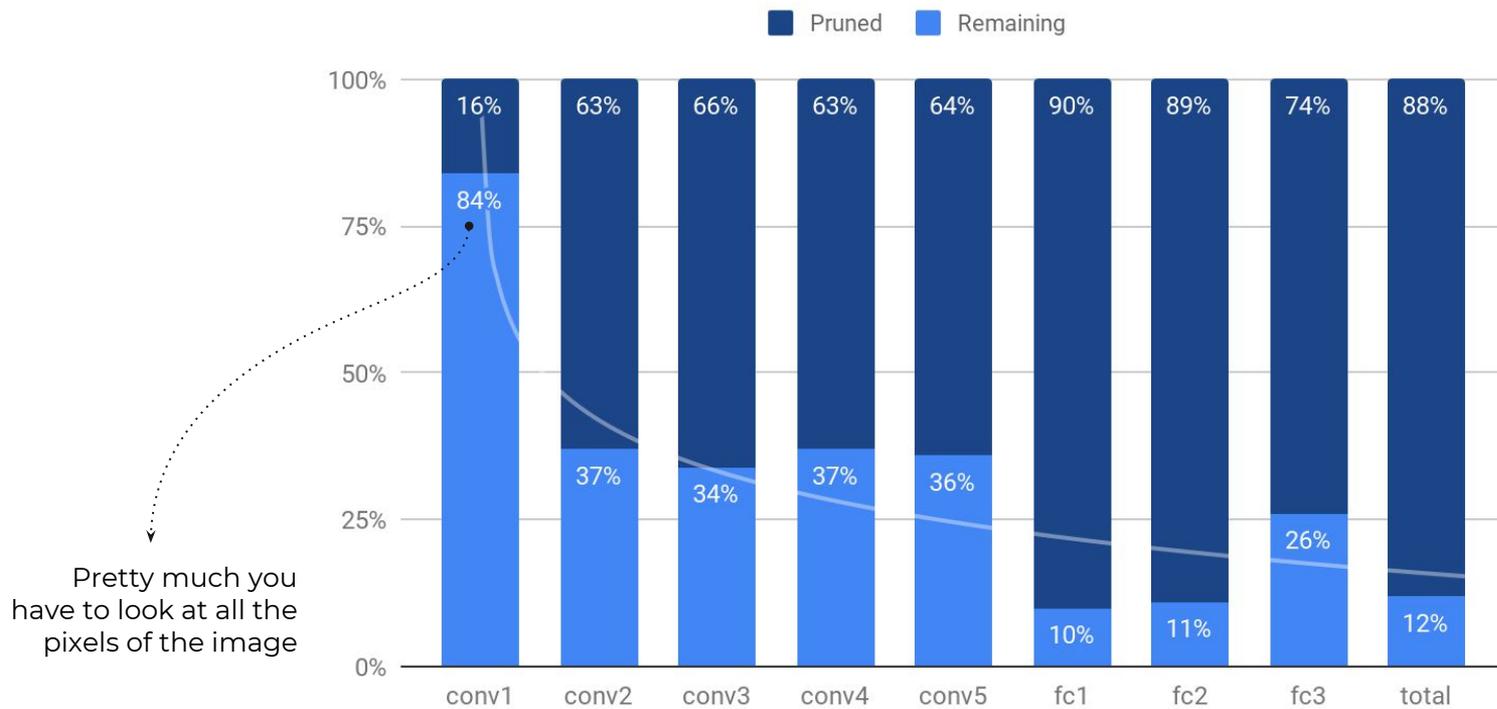
Pruning



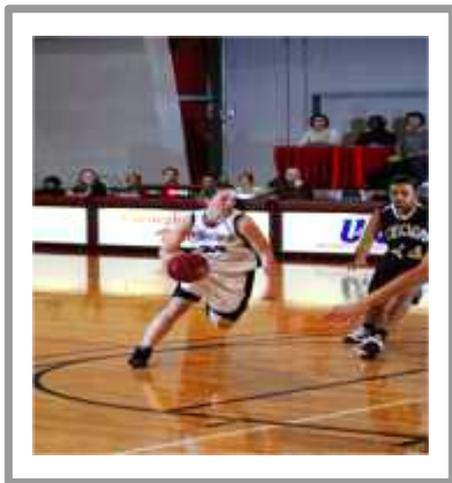
Pruning



Pruning: AlexNet



Pruning + Image Captioning



ORIGINAL

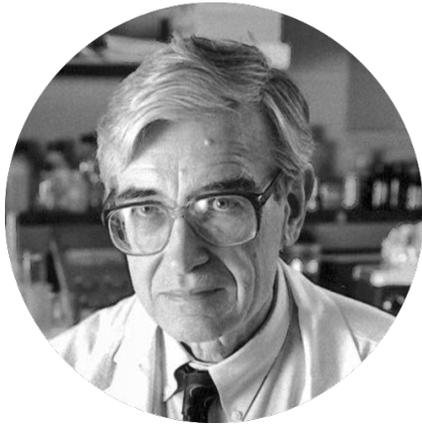
a basketball player in a white uniform
is a playing with a **ball**.

PRUNED 90%

a basketball player in a white uniform
is a playing with a **basketball**.

If our brain loses 90% of neurons,
can we still describe this image
with this high accuracy?

Pruning in Human?

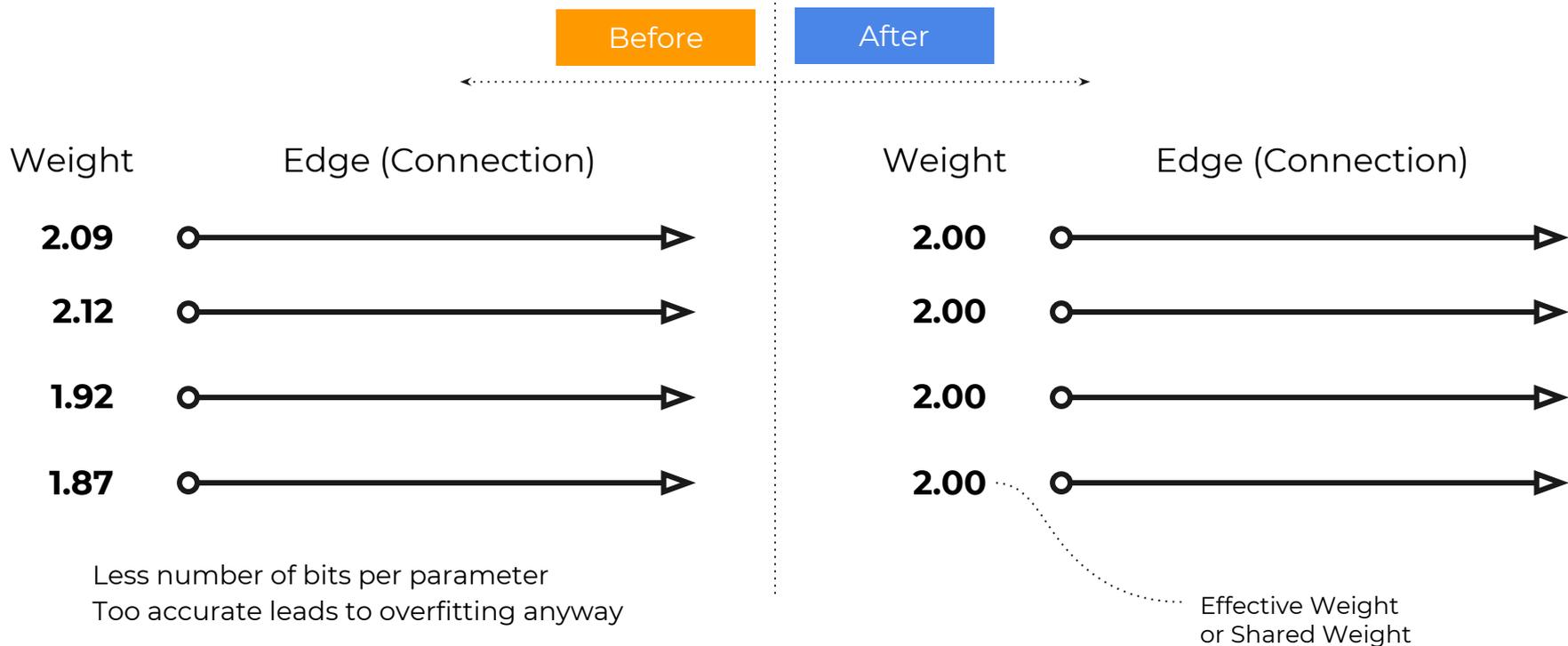


Peter Huttenlocher
(1931-2013)

500 trillion synapses
1,000 trillion synapses
50 trillion synapses



Weight Sharing



Weight Sharing

Weights are not shared across layers. The shared weights approximate the original network because the method determines weight sharing after a network is fully trained.

K-Means Clustering (with K=4)

2.09	-0.98	1.48	0.09
0.05	-0.14	-1.08	2.12
-0.91	1.92	0	-1.03
1.87	0	1.53	1.49

4x4=16 numbers



Color (Cluster)	Original Weights	Effective Weight (Centroid)
Purple	[2.09, 2.12, 1.92, 1.87]	2.00
Green	[1.48, 1.53, 1.49]	1.50
Orange	[0.09, 0.05, -0.14, 0, 0]	0.00
Yellow	[-0.98, -1.08, -0.91, -1.03]	-1.00

Only 4 numbers

Weight Sharing



Also called "Codebook"

Weight Sharing / Quantization

Color	Effective Weight	Index (Int)	Index (Binary)
Purple	2.00	0	00
Green	1.50	1	01
Orange	0.00	2	10
Yellow	-1.00	3	11

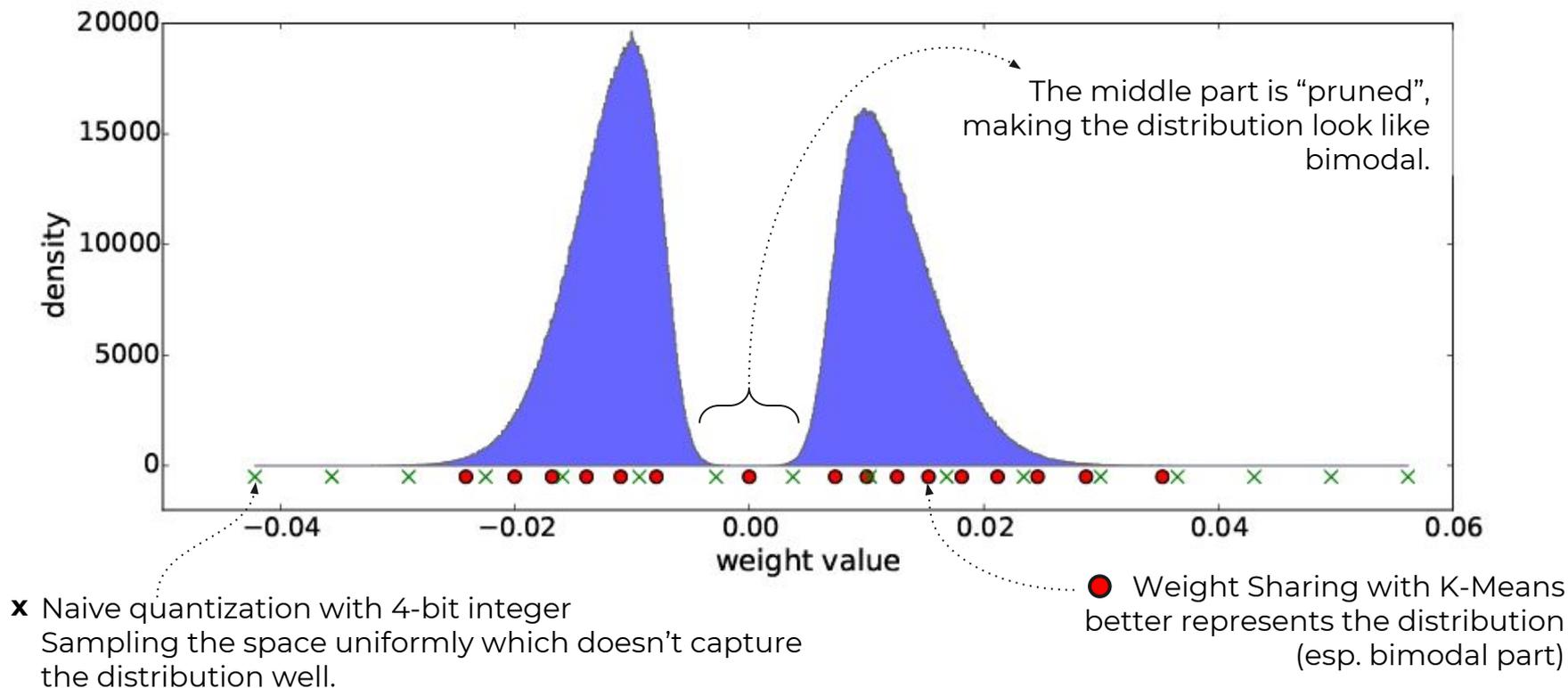


00	11	01	10
10	10	11	00
11	00	10	11
00	10	01	01

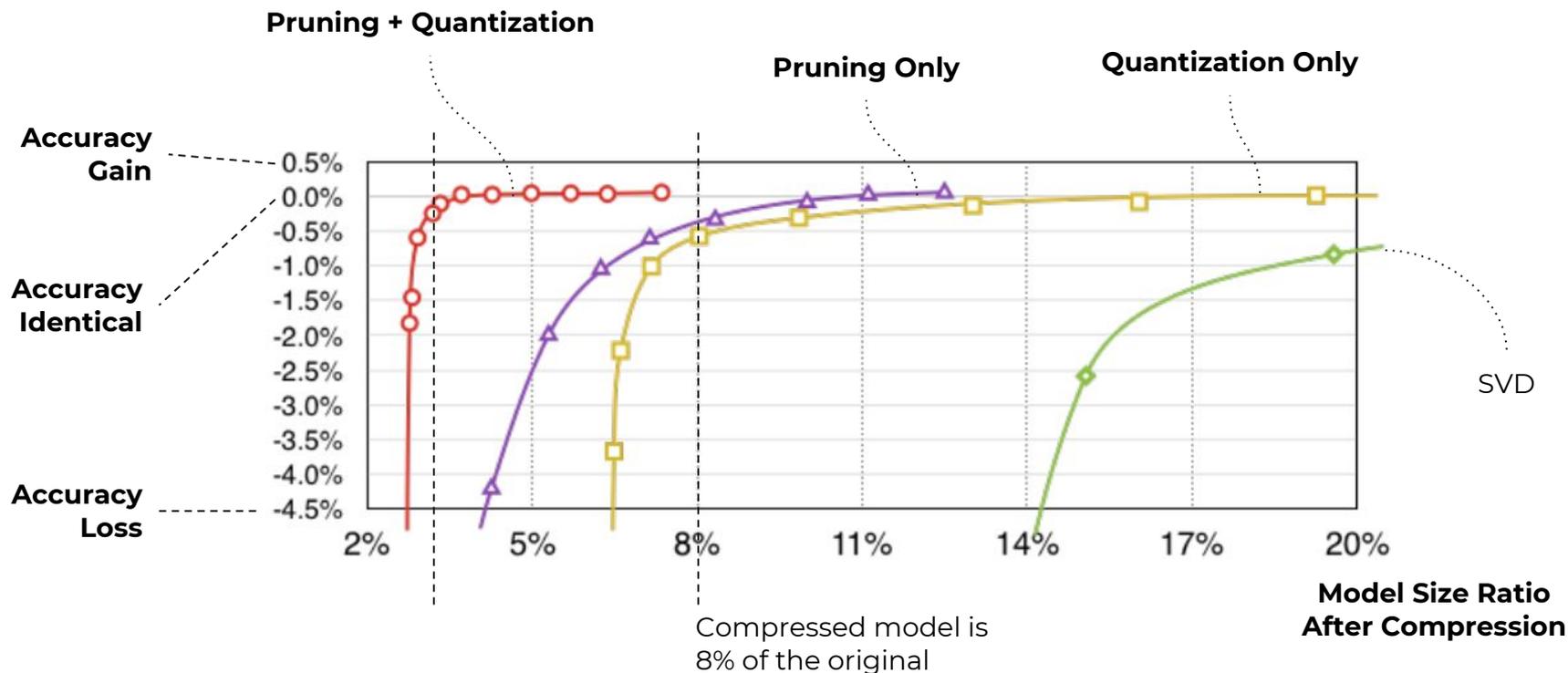
Lookup

No more 32-bit FP and only 2 bits

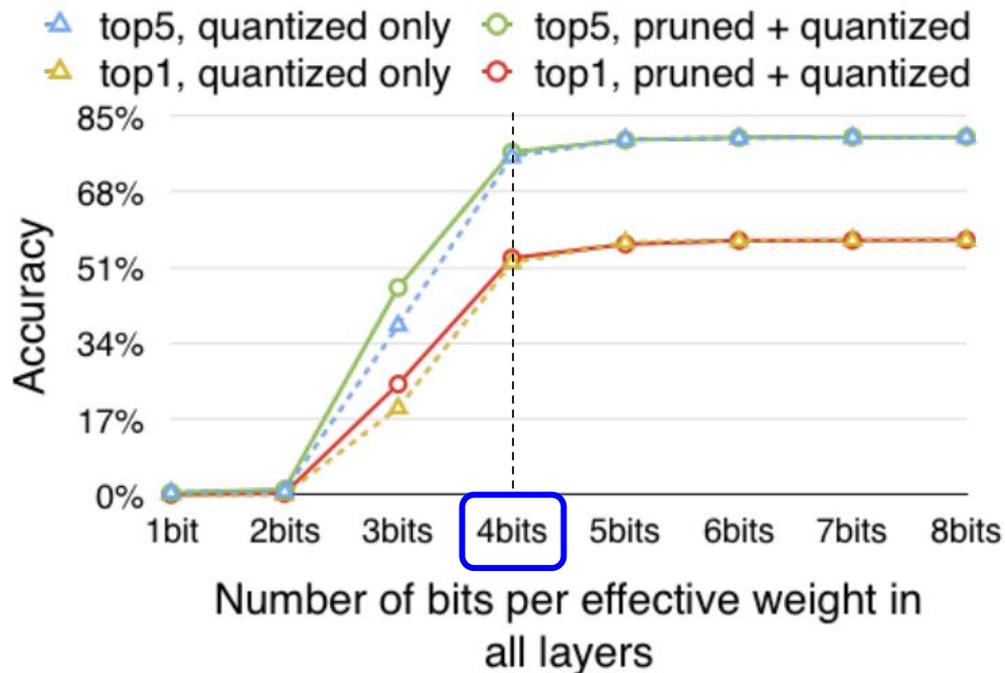
Trained Quantization: Weight Distribution



Pruning and/or Quantization: Accuracy



Pruning and/or Quantization: Accuracy



Huffman Coding

$$2 \text{ bits} * (24 + 5 + 4 + 3) = 72 \text{ bits}$$



00	11	00	01	10	00
00	00	00	00	00	00
10	10	00	11	00	00
11	00	00	10	11	00
00	10	00	01	01	00
00	00	00	00	00	00

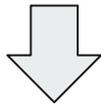
6x6=36 numbers

Color	Effective Weight	Index (Int)	Index (Binary)	Count	%
Purple	2.00	0	00	24	66.67
Orange	0.00	2	10	5	13.89
Yellow	-1.00	3	11	4	11.11
Green	1.50	1	01	3	8.33

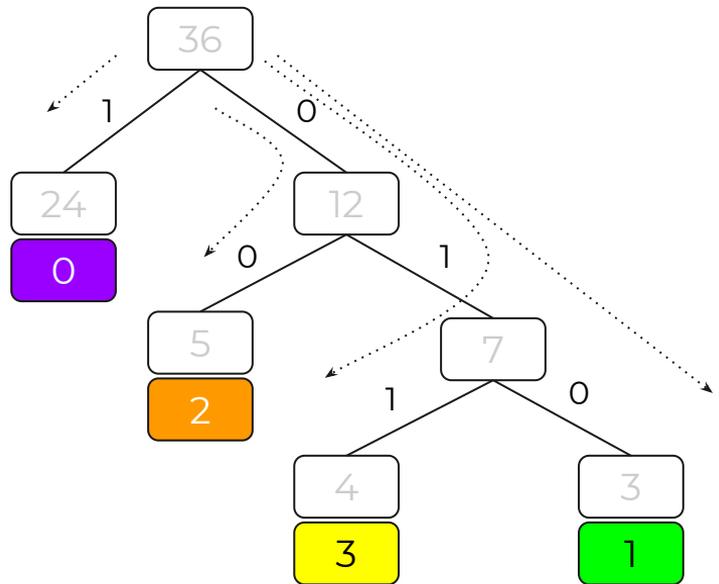
Frequent weights → use **less** bits to represent
Infrequent weights → use **more** bits to represent

Huffman Coding

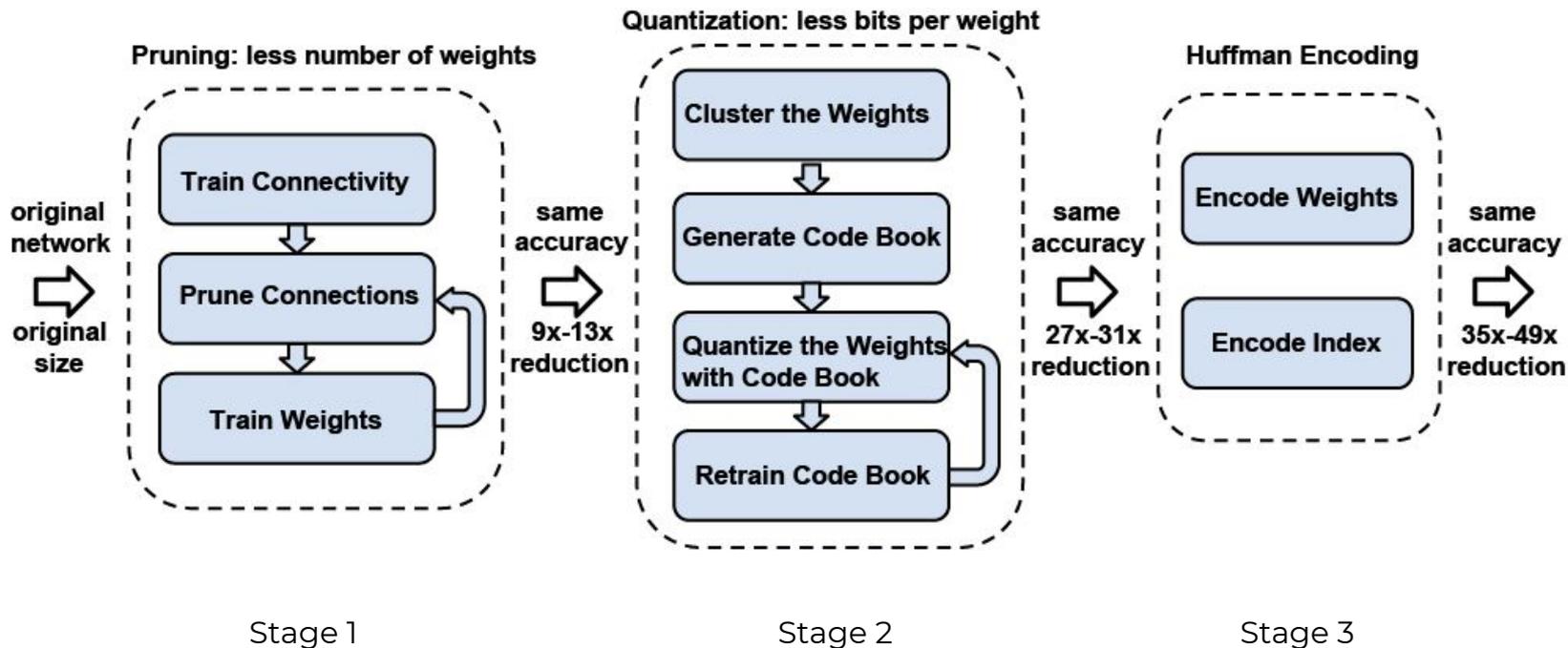
Color	Effective Weight	Index (Int)	Huffman Code	Count	%
Purple	2.00	0	1	24	66.67
Orange	0.00	2	00	5	13.89
Yellow	-1.00	3	011	4	11.11
Green	1.50	1	010	3	8.33



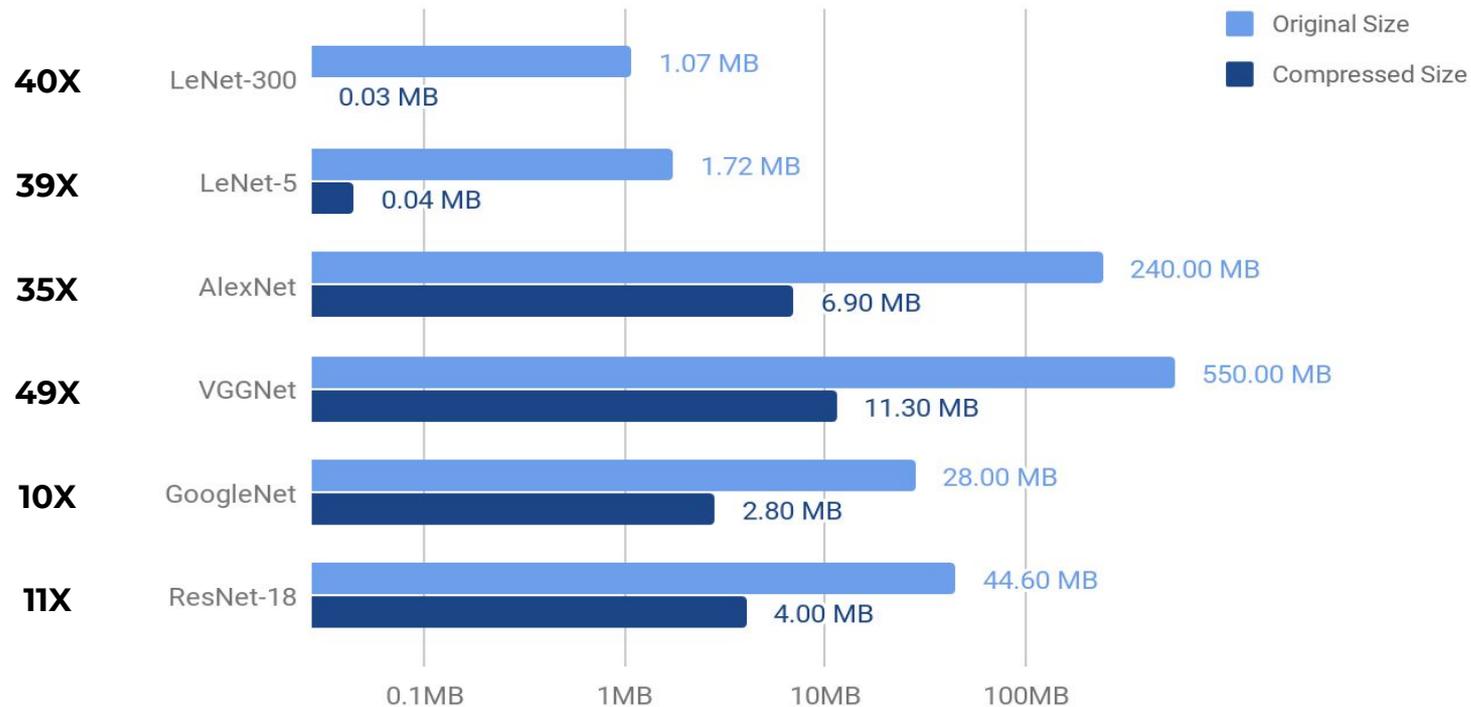
$$(1 \text{ bit} * 24) + (2 \text{ bits} * 5) + (3 \text{ bits} * 4) + (3 \text{ bits} * 3) = 55 \text{ bits}$$



Deep Compression

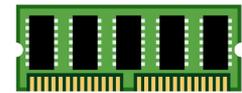
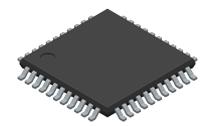
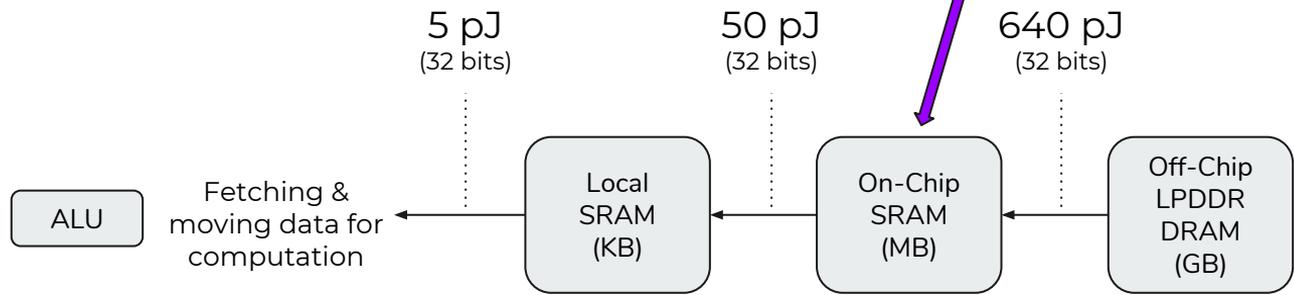
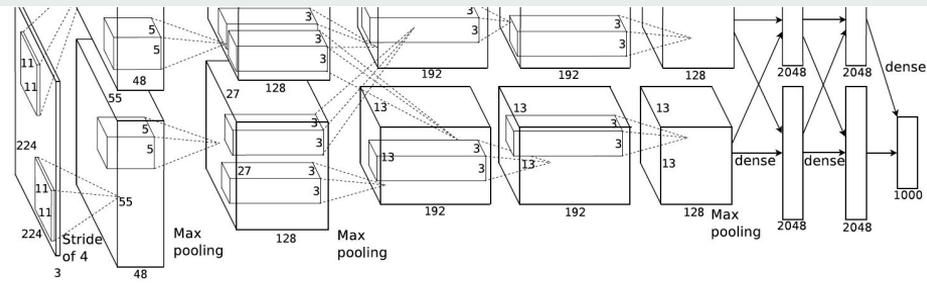


Compression Ratio (w/o accuracy loss)



Deep Compression

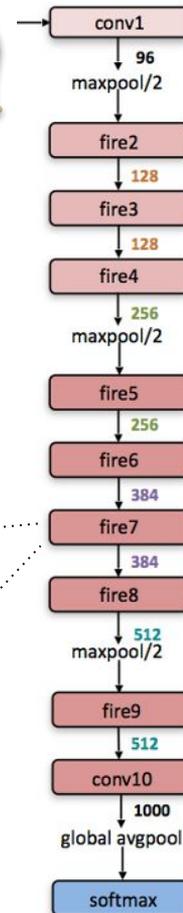
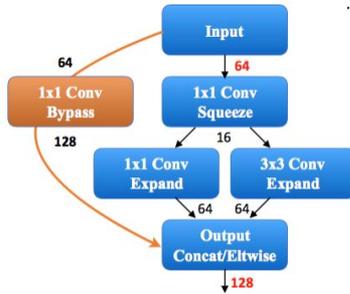
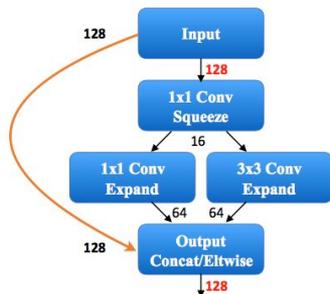
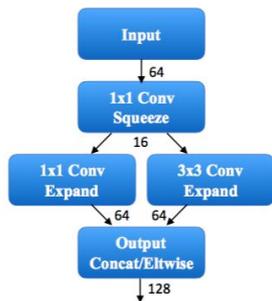
Large DNNs such as AlexNet, VGGNet are **fully fit** into **on-chip SRAM**.



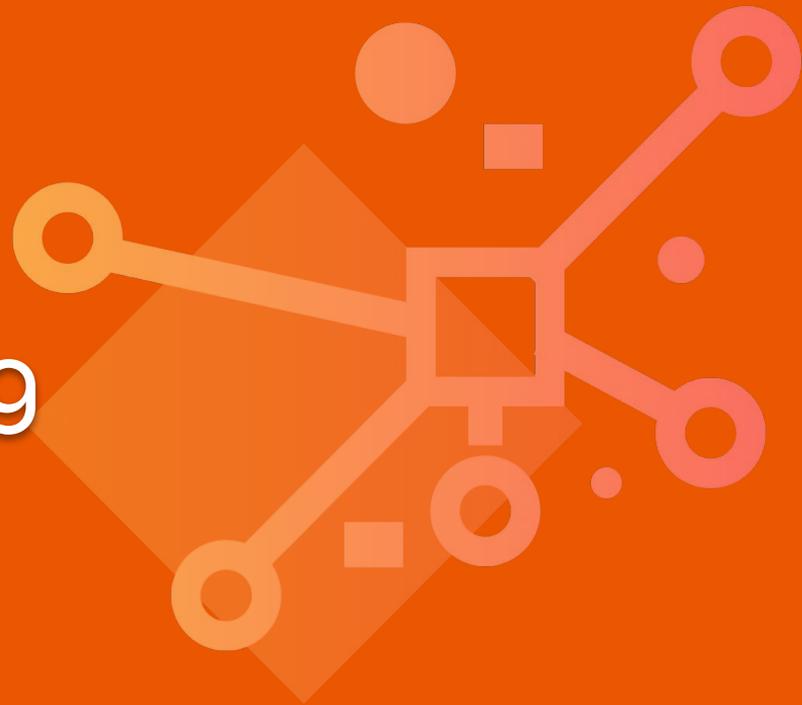
LPDDR

SqueezeNet + Deep Compression

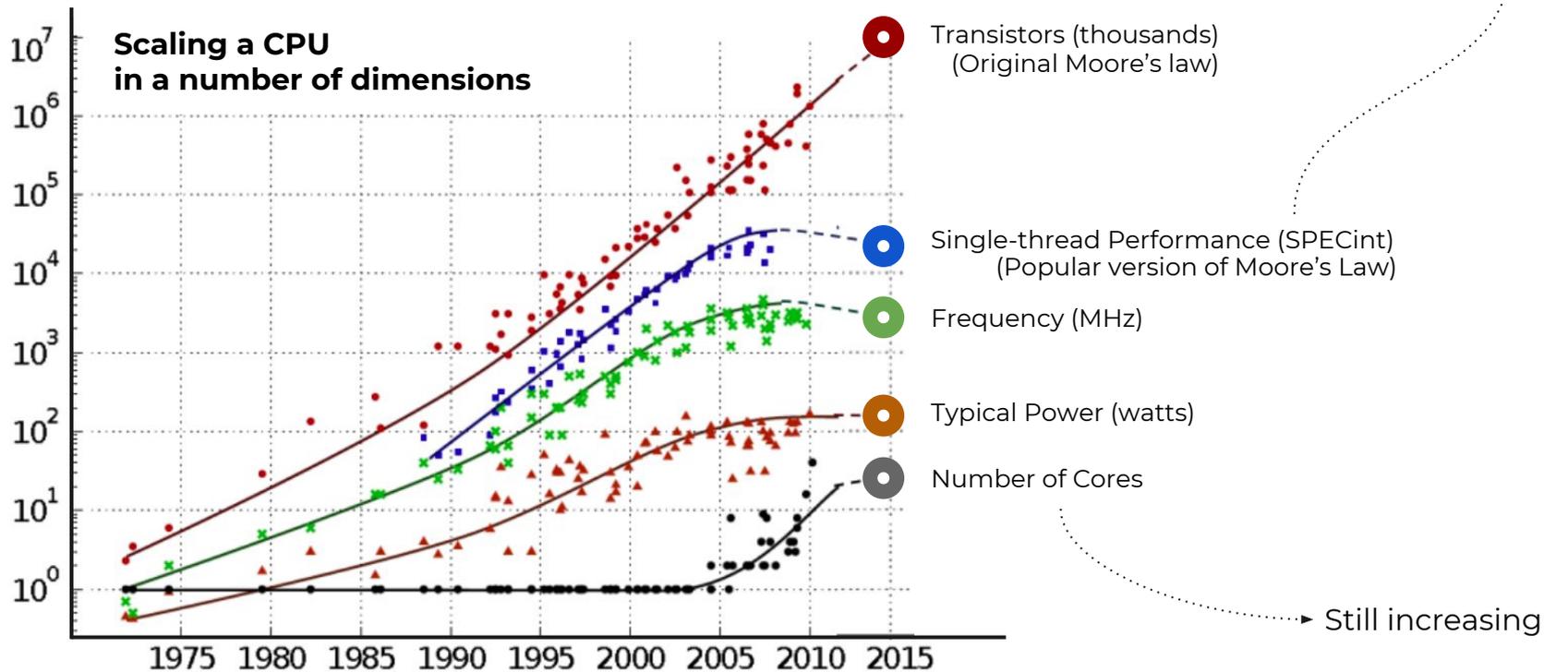
Can we apply Deep Compression to **already compact model** such as SqueezeNet?



Algorithms for Efficient Training

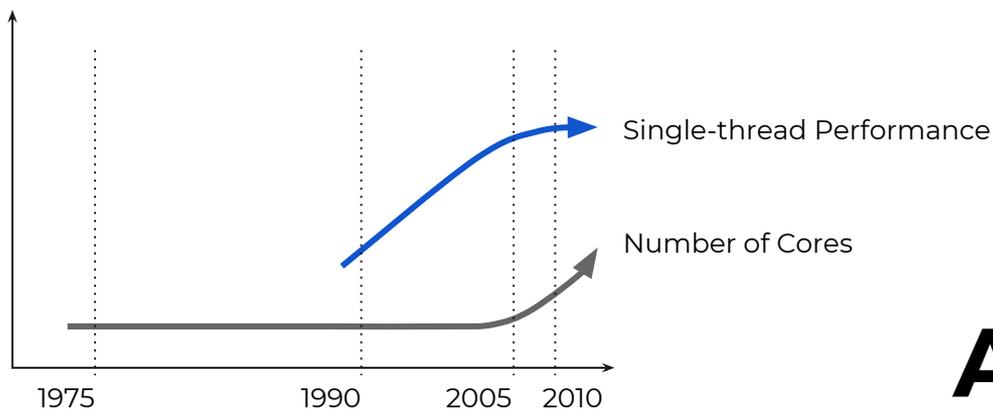


Moore's Law



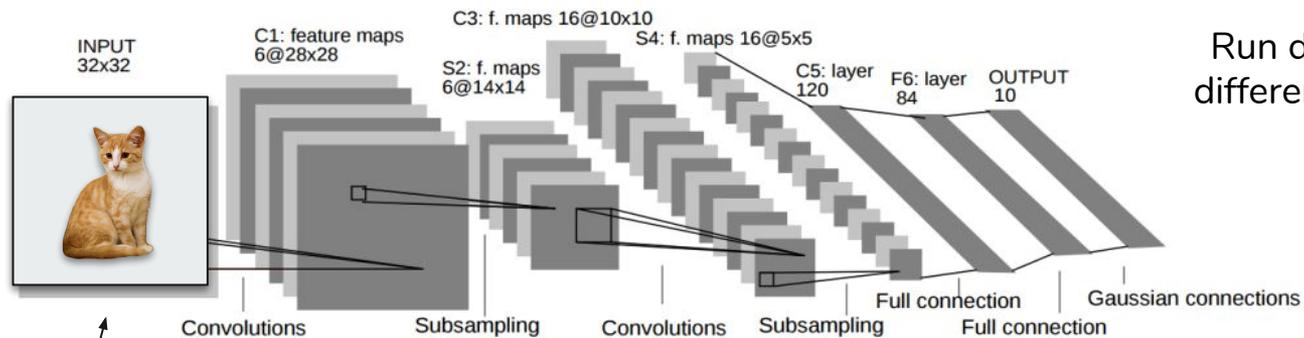
Without Moore's Law (popular version)

Q: How are we going to continue to scale the performance we need to build a better DNN?



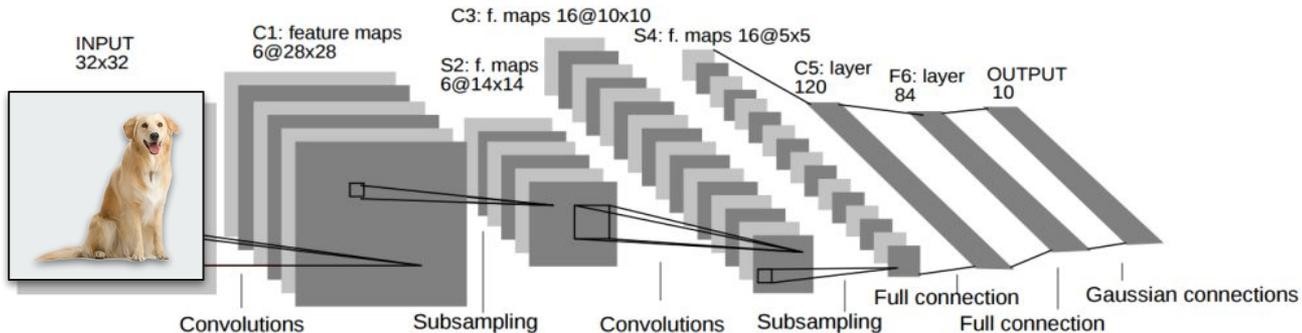
A: Use multiple processors in parallel

Data Parallelism



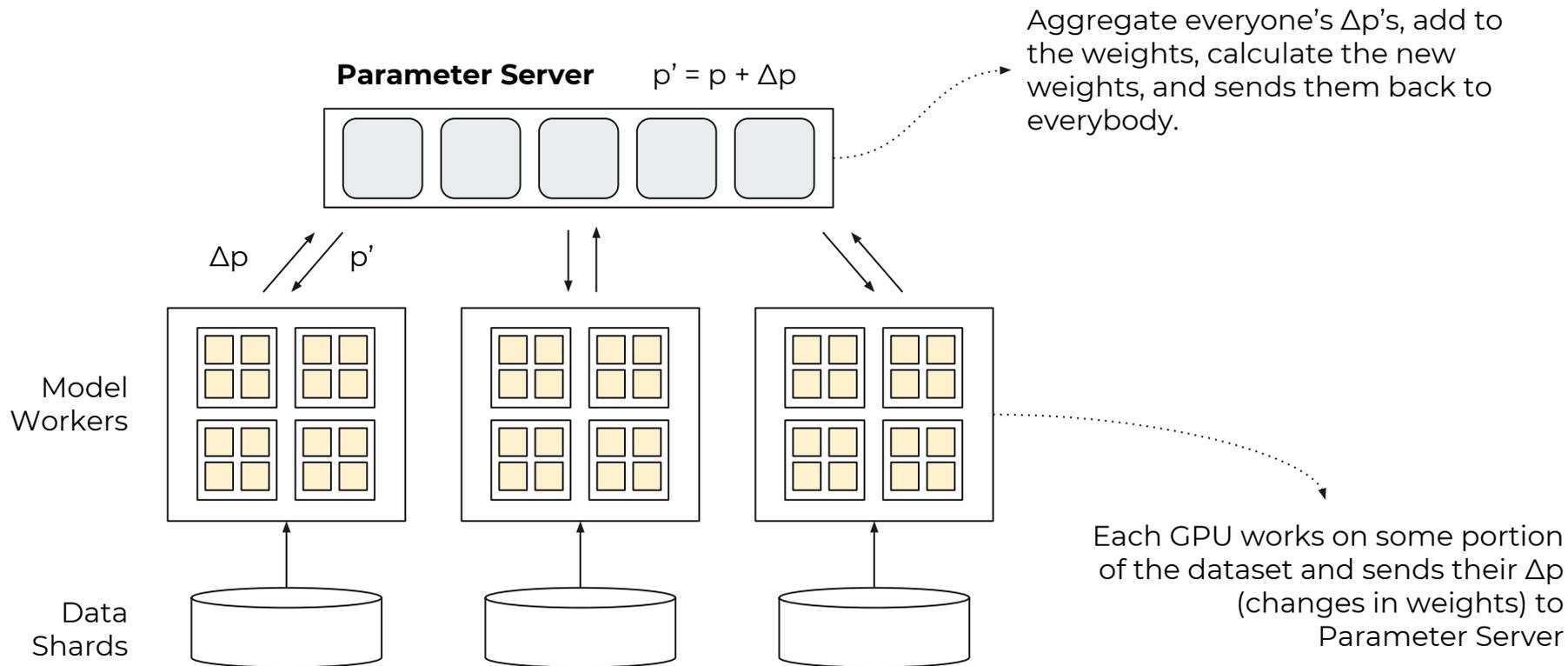
Run different images on different GPUs in parallel

CIFAR-10

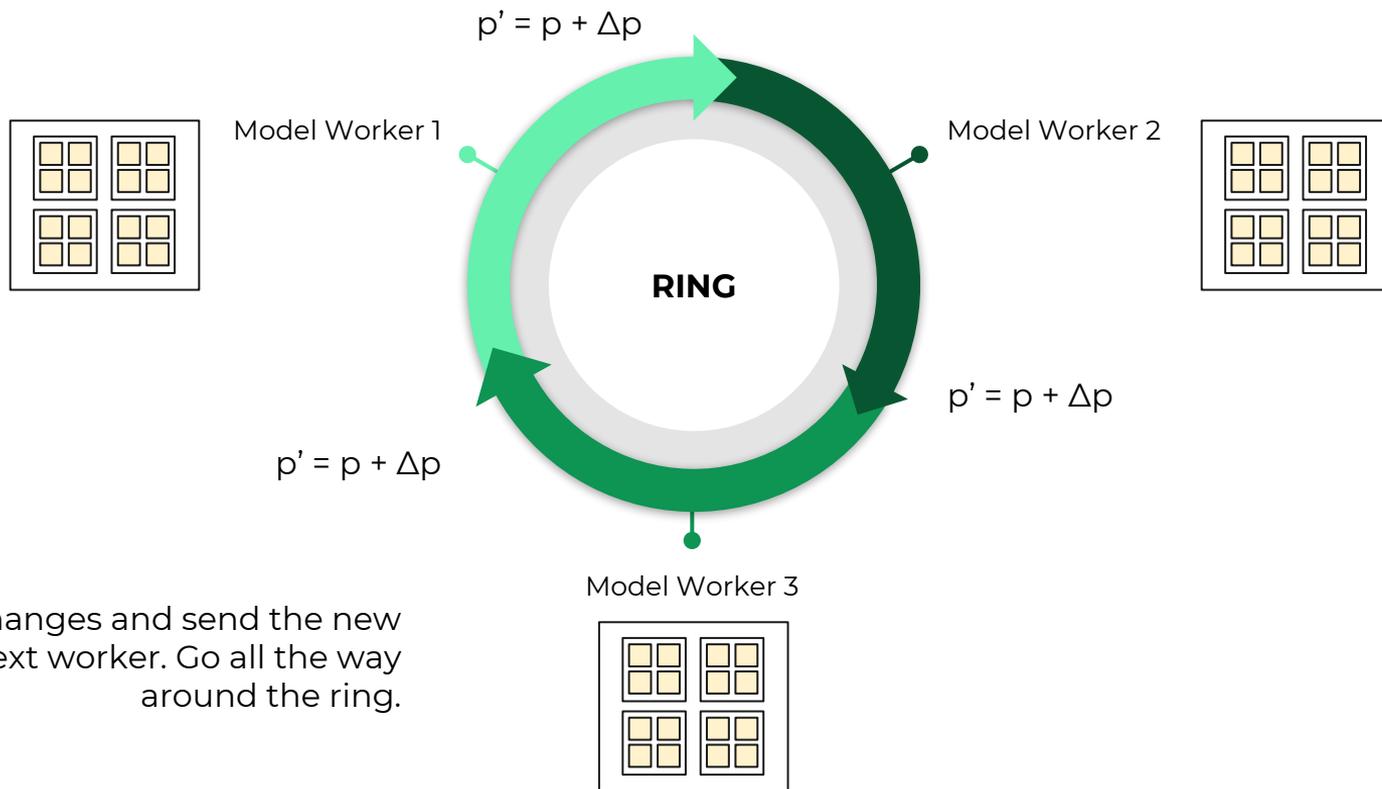


Weight updates must be coordinated between workers.

Data Parallelism



Data Parallelism



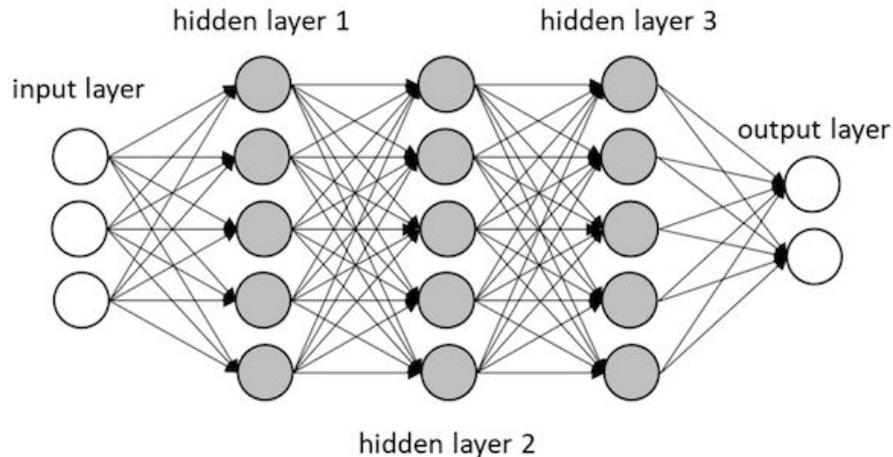
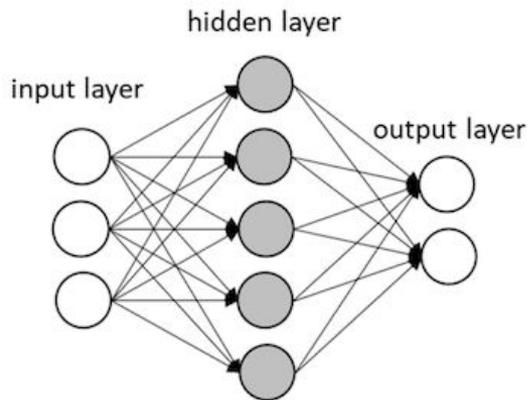
Add my weight changes and send the new weight to the next worker. Go all the way around the ring.

Hyper-Parameter Parallelism

Try many alternative networks in parallel

- Different number of layers
- Different size of convolutional kernels
- Different number of neurons per layer
- ...

—————→ Search in the parameter space

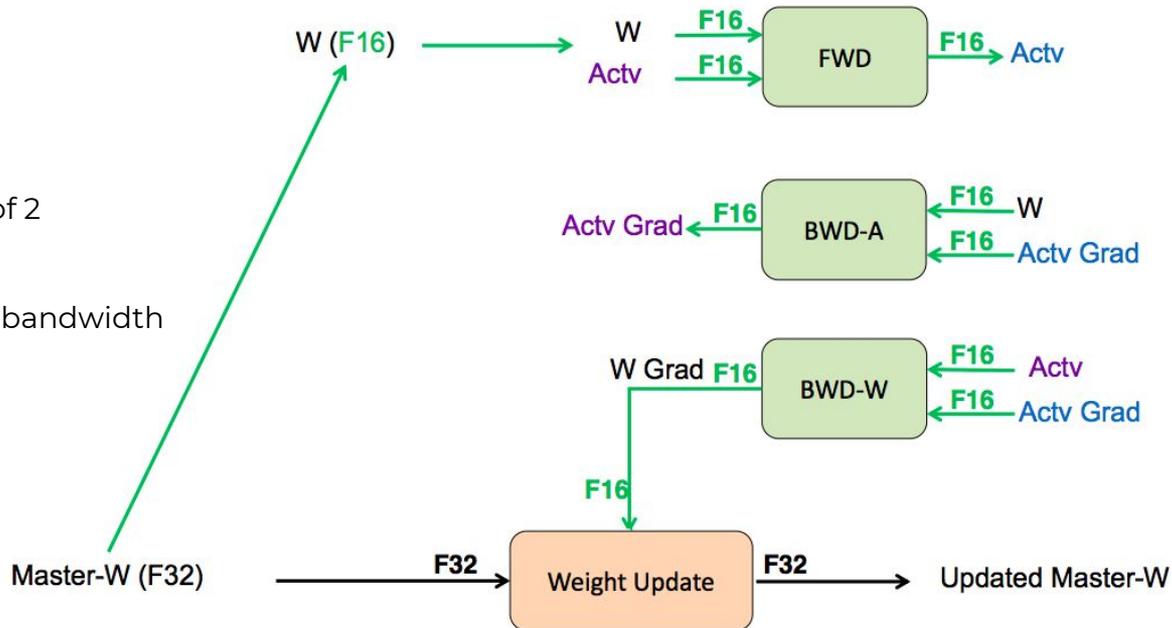


Mixed Precision Training

More precision than required \rightarrow reduce precision

Save by factor of 2

- Storage
- Memory bandwidth



Mixed Precision Training: Comparison

AlexNet

Mode	Top1 accuracy, %	Top5 accuracy, %
Fp32	58.62	81.25
Mixed precision training	58.12	80.71

Inception V3

Mode	Top1 accuracy, %	Top5 accuracy, %
Fp32	71.75	90.52
Mixed precision training	71.17	90.10

ResNet-50

Mode	Top1 accuracy, %	Top5 accuracy, %
Fp32	73.85	91.44
Mixed precision training	73.6	91.11

Hardware





Quick Overview of Hardware Side, But Why?

01

We write algorithms and software that runs on hardware.

02

Some of the algorithms we reviewed are actually used in the hardware design.

Google TPU

- Tensor Processing Unit
- Compared to GPU & CPU
 - 15x to 30x faster
 - 30x to 80x better energy efficiency
- Only internal use
 - e.g. AlphaGo, Street View, Photos, ...

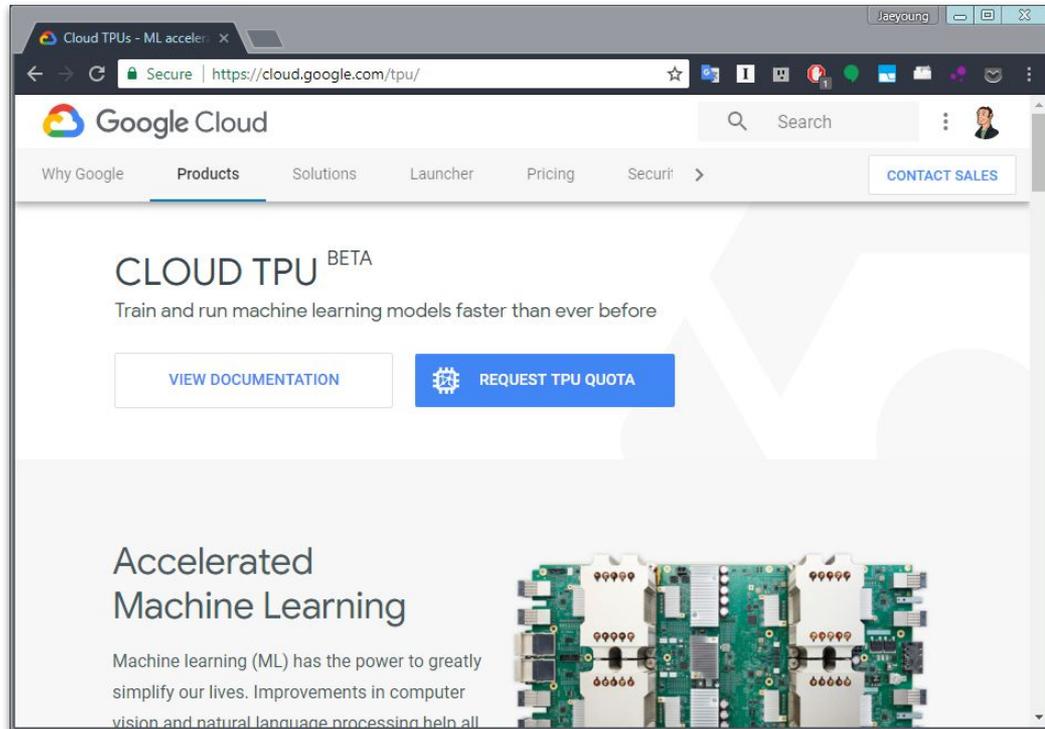


Can be inserted into a SATA hard disk slot for easy/fast deployment to existing server infrastructure



Google Cloud TPU (2nd gen)

- Mid Feb 2018
Cloud TPU^{BETA} announced.
- Supports for inference as well as training



Google Cloud TPU (2nd gen)

The screenshot displays the Google Cloud Platform console interface for creating a Cloud TPU. The browser address bar shows the URL `https://console.cloud.google.com/c...`. The page title is "Create a Cloud TPU".

Left Sidebar (Navigation):

- Compute Engine
- VM instances
- Instance groups
- Instance templates
- Disks
- Snapshots
- Images
- TPUs** (highlighted)
- Committed use discounts

Main Form Fields:

- Name:** jaeyoung-node-1
- Zone:** us-central1-c
- TPU type:** tpu-v2
- TensorFlow version:** 1.8
- Network:** default

The Secret of Google TPU

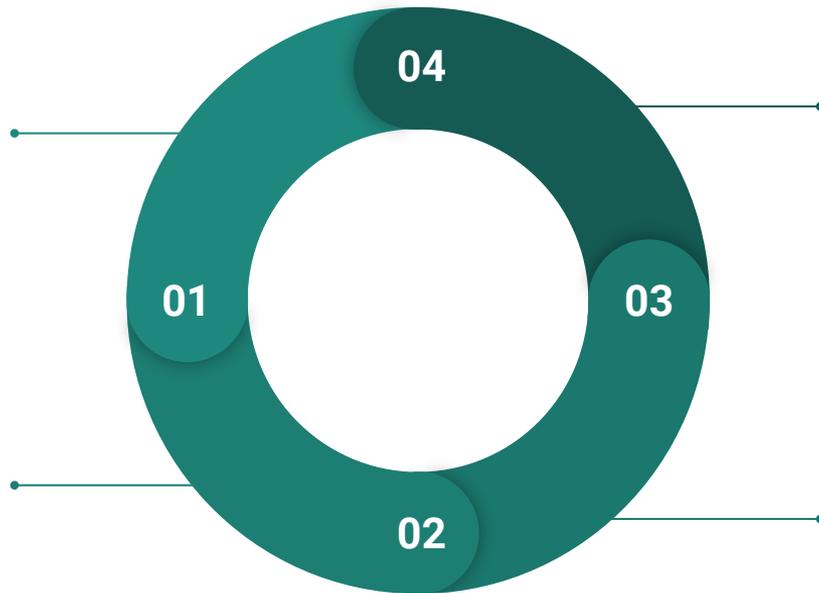


Quantization

Optimization technique that uses an 8-bit integer to approximate an arbitrary value between a preset minimum and a maximum value

CISC Instruction Set

High-level instructions specifically designed for neural network inference.



Matrix Multiply Unit

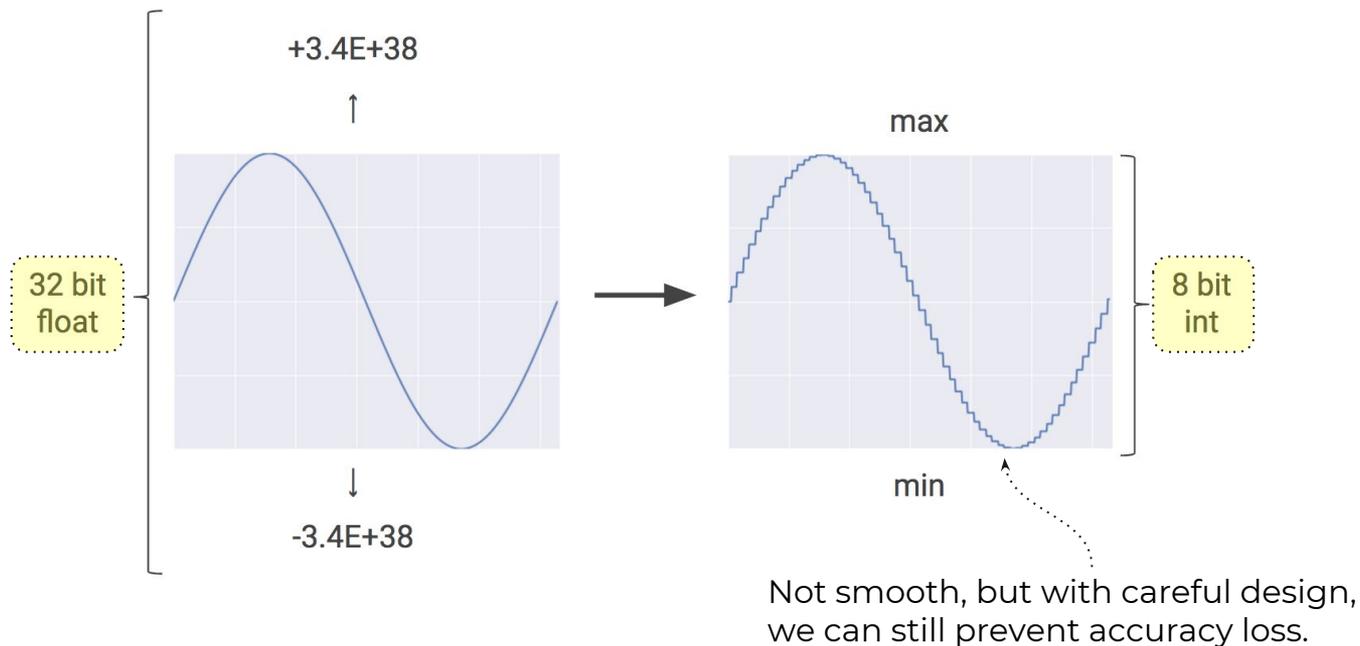
Processes hundreds of thousands of operations (= matrix operation) in a single clock cycle.

Minimal Design

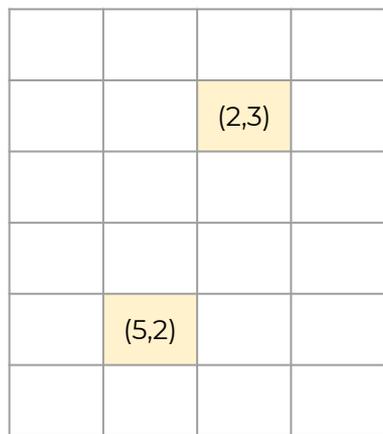
Optimized for neural network inference only. In the TPU, the control logic is minimal and takes under 2% of the die.

Google TPU: Quantization

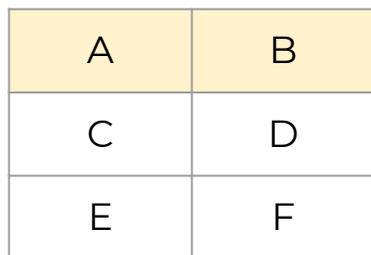
We have already seen the power of Quantization when discussing Deep Compression.



Google TPU: CISC Instruction Set



6x4 Memory



CPU Register



Reduced vs. Complex Instruction Set Computer

RISC

LOAD A, 2:3
LOAD B, 5:2
PROD A, B
STORE 2:3, A

Low-level simple instructions
that are commonly used

e.g.
load, store, add, multiply

CISC

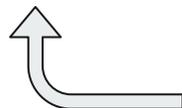
MULT 2:3, 5:2

High-level instructions
that perform complex operations

e.g.
compute
multiply-and-add
many times

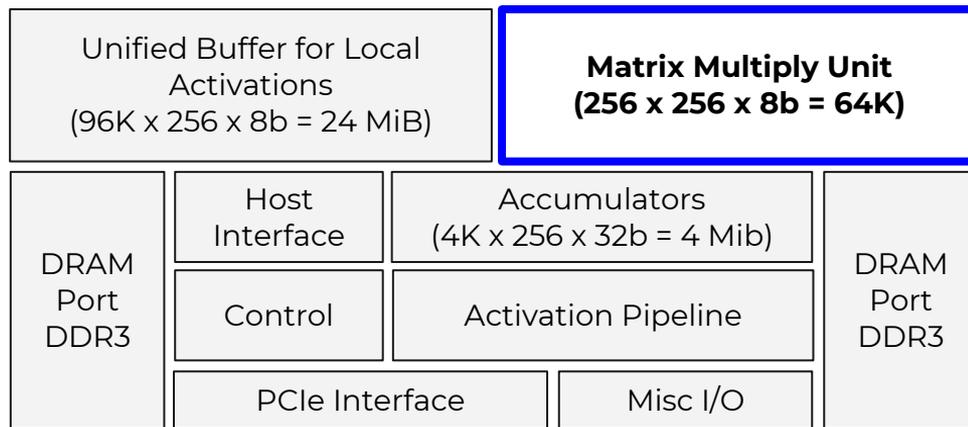
Google TPU: CISC Instruction Set

TPU Instruction	Function
Read_Host_Memory	Read data from memory
Read_Weights	Read weights from memory
MatrixMultiply / Convolve	Multiply or convolve with the data and weights, accumulate the results
Activate	Apply activation functions
Write_Host_Memory	Write result to memory



High-level instructions specifically designed for neural network inference

Google TPU: Matrix Multiplier Unit

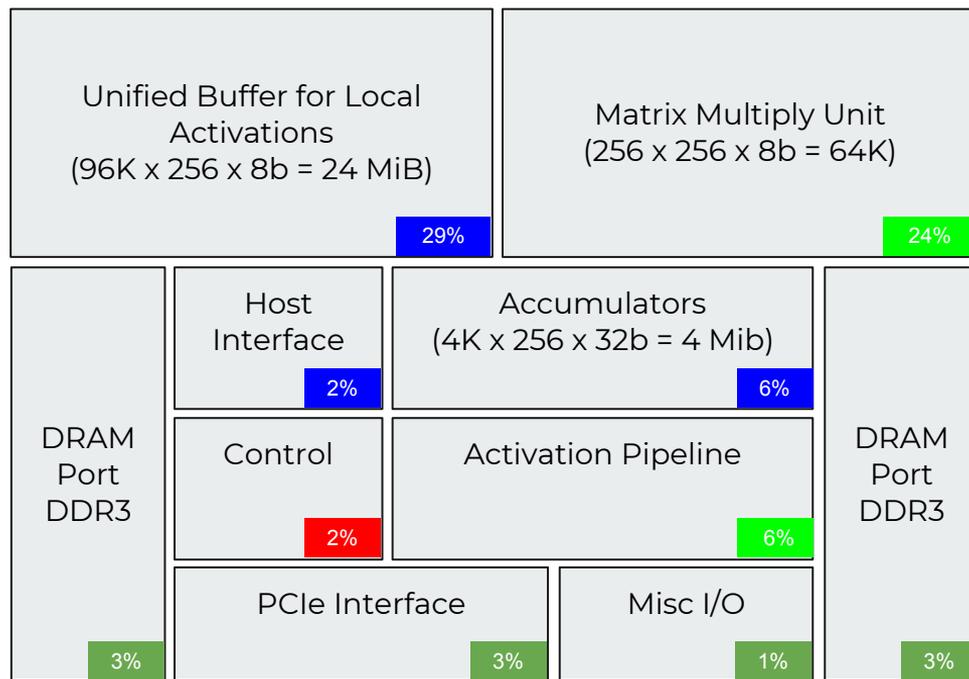
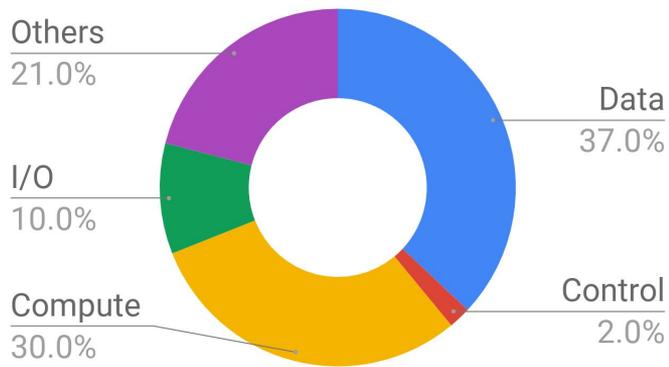


65,536
(multiply-and-add operations per cycle)

700
(TPU clock in MHz)

46×10^{12}
(65,536 × 700M operations per sec)

Google TPU: Minimal Design



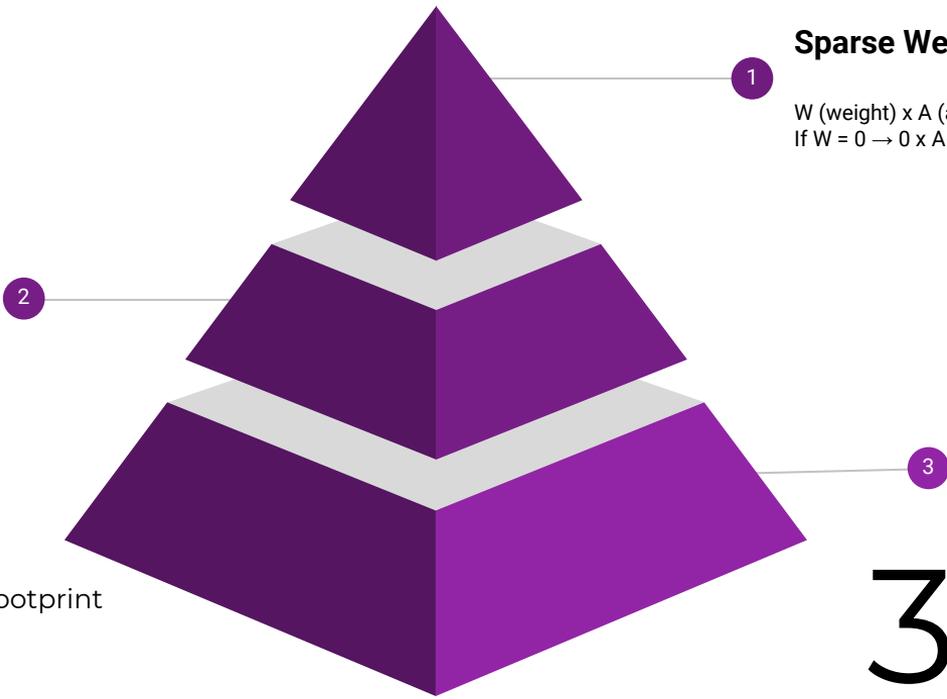
EIE (Efficient Inference Engine)

Weight Sharing

With K-means clustering, e.g. the blues (2.09, 2.12., 1.92, 1.87) are treated as 2.0 instead (i.e. effective weights)

2.09	-0.98	1.48	0.09
0.05	-0.14	-1.08	2.12
-0.91	1.92	0	-1.03
1.87	0	1.53	1.49

8x less memory footprint



Sparse Weight

W (weight) \times A (activation)
If $W = 0 \rightarrow 0 \times A = 0$

10x less computation

5x less memory footprint

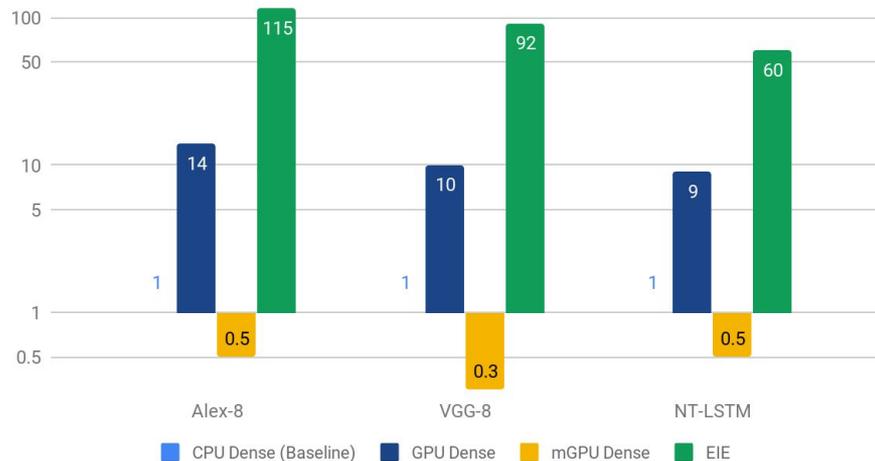
Sparse Activation

W (weight) \times A (activation)
If $A = 0 \rightarrow W \times 0 = 0$

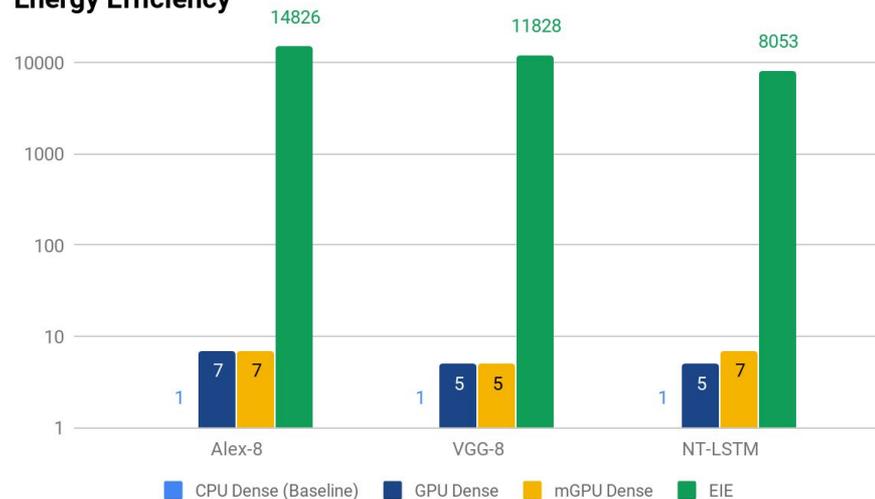
3x less computation

EIE: Speedup and Energy Efficiency

Speedup



Energy Efficiency



Break

