



# Principles of Machine Learning Systems

## 2: Model Compression

# Roadmap for Today

---

## Architecture (*Re-*)Design

1. Early Evolution; Decomposition as a tool
2. Efficient Architectures; Example: MobileNet

## Architecture Compression

3. Parameter and Channel Pruning
4. Parameter Quantization
5. Knowledge Distillation
6. Compression pipelines



# Roadmap for Today

---

## **Architecture (Re-)Design**

1. Early Evolution; Decomposition as a tool
2. Efficient Architectures; Example: MobileNet

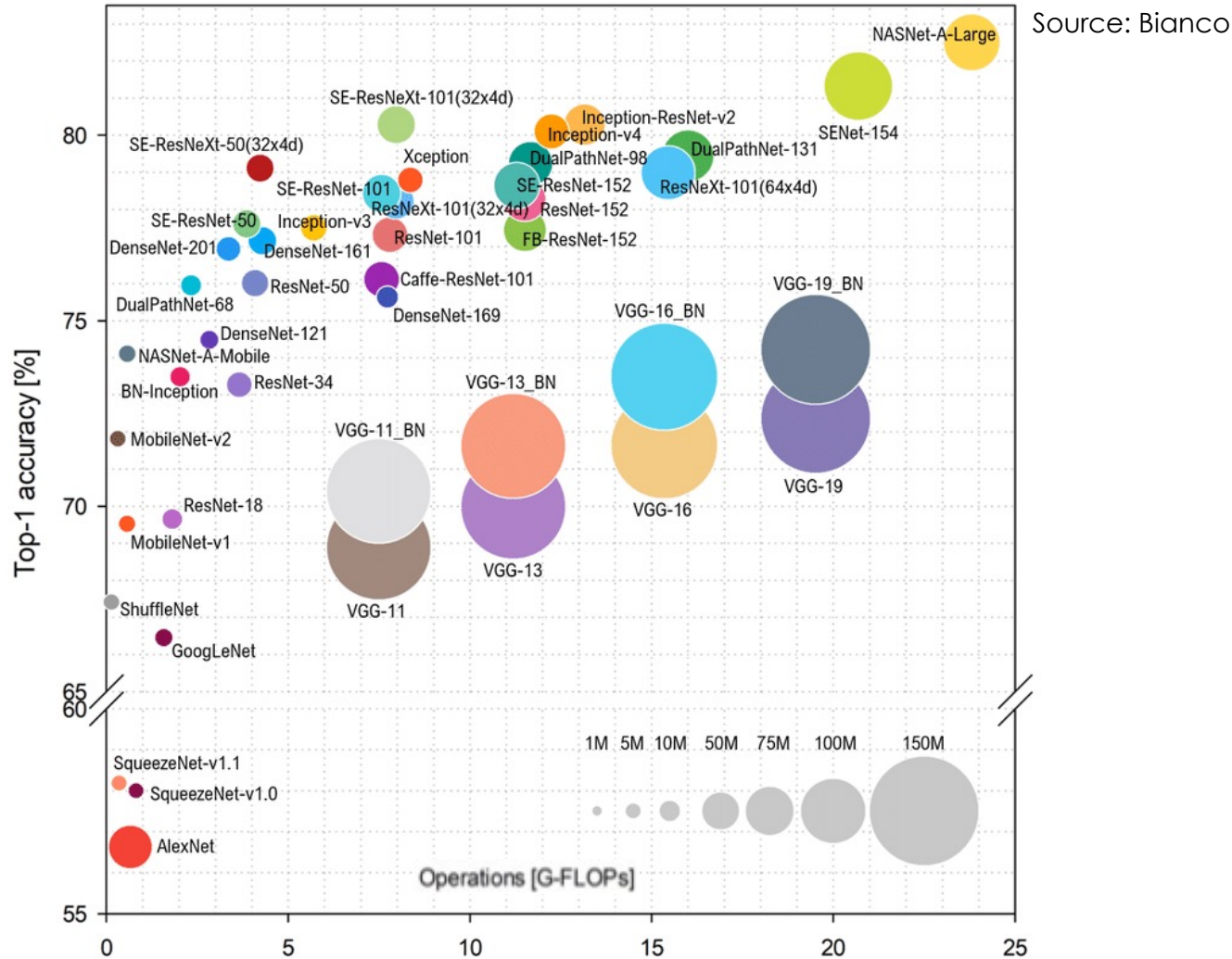
## **Architecture Compression**

3. Parameter and Channel Pruning
4. Parameter Quantization
5. Knowledge Distillation
6. Compression pipelines

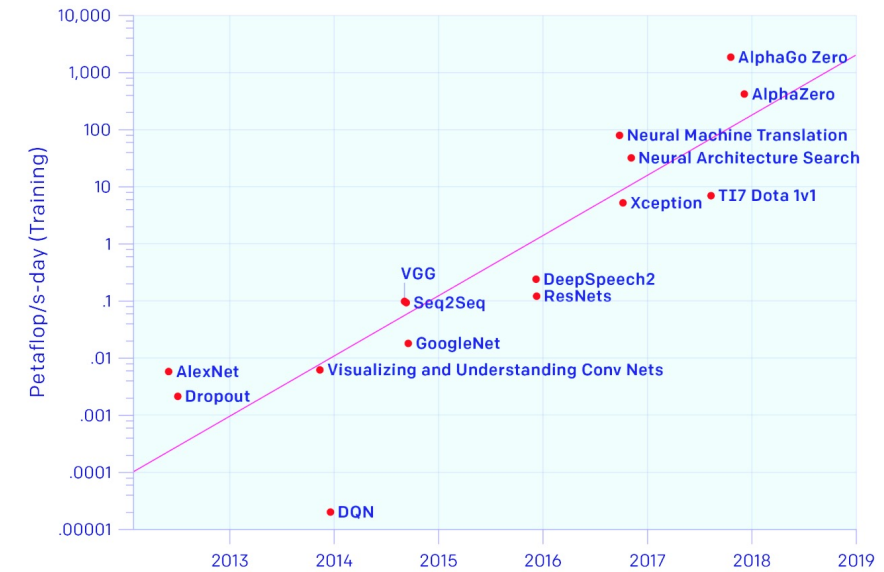




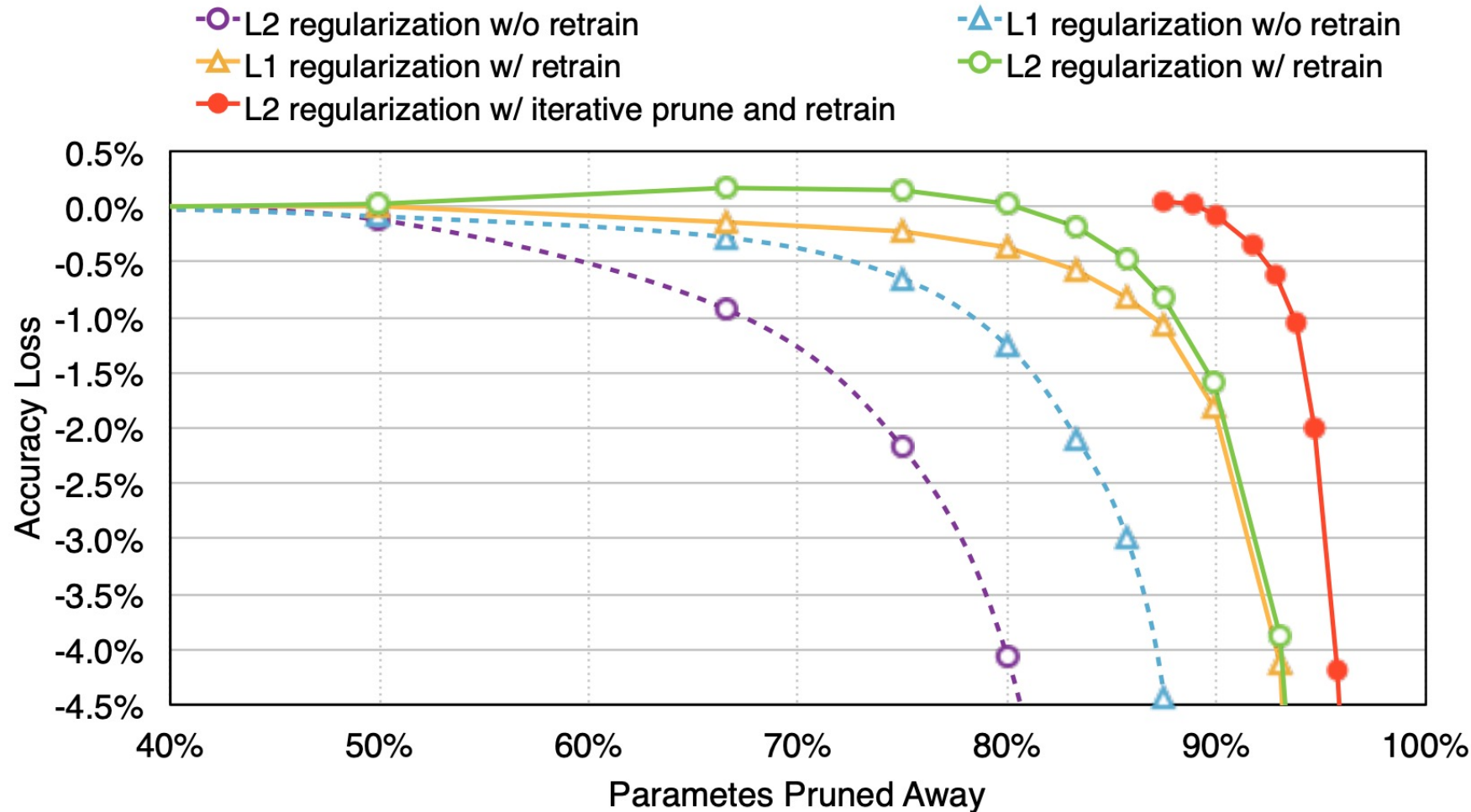
# Roadmap for Today



Source: OpenAI



# Roadmap for Today

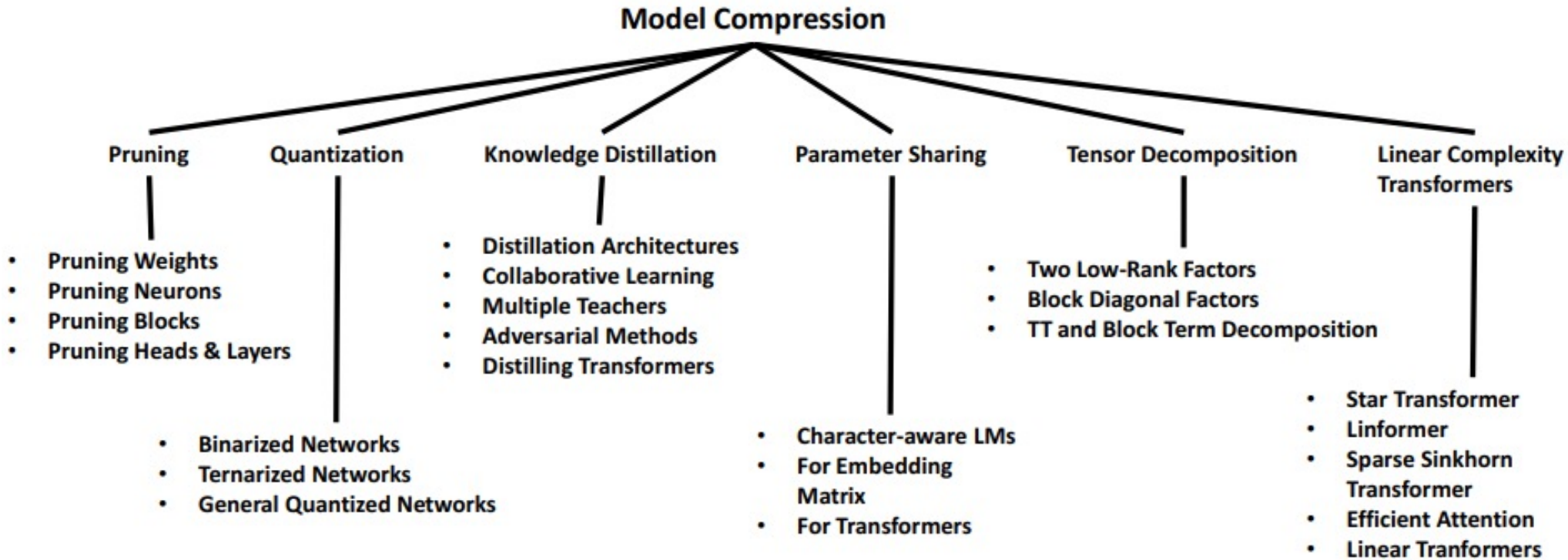


Source: Han



# Roadmap for Today

---



Source: Gupta



# Roadmap for Today

---

Architecture (*Re-*)Design

- 1. Early Evolution; Decomposition as a tool**
2. Efficient Architectures; Example: MobileNet

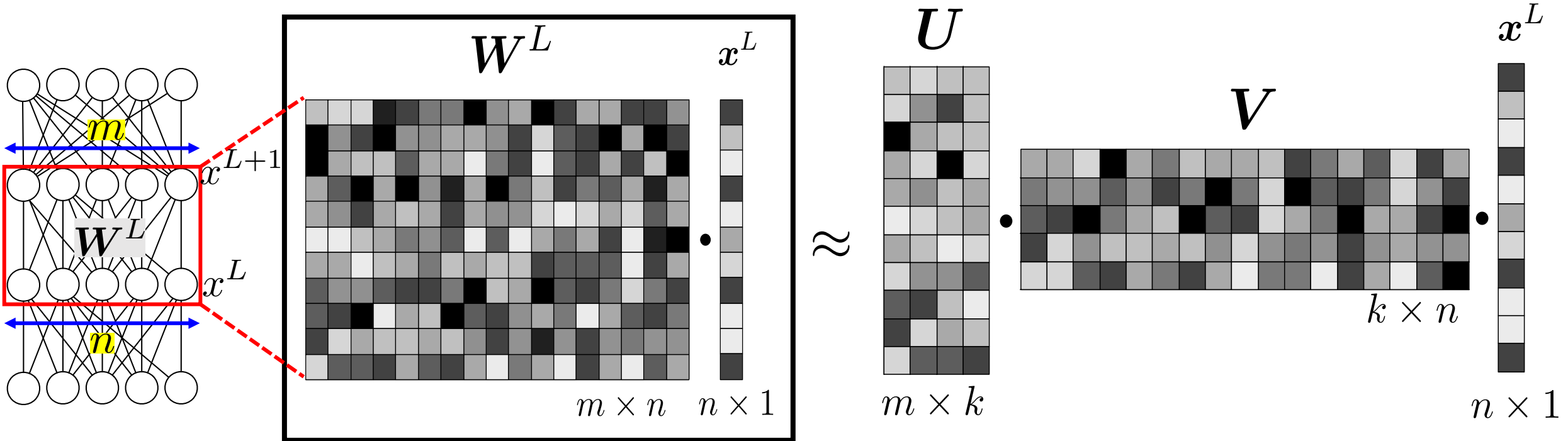
Architecture Compression

3. Parameter and Channel Pruning
4. Parameter Quantization
5. Knowledge Distillation
6. Compression pipelines



# Matrix Decomposition

Bottleneck: Repeated Matrix Multiplications



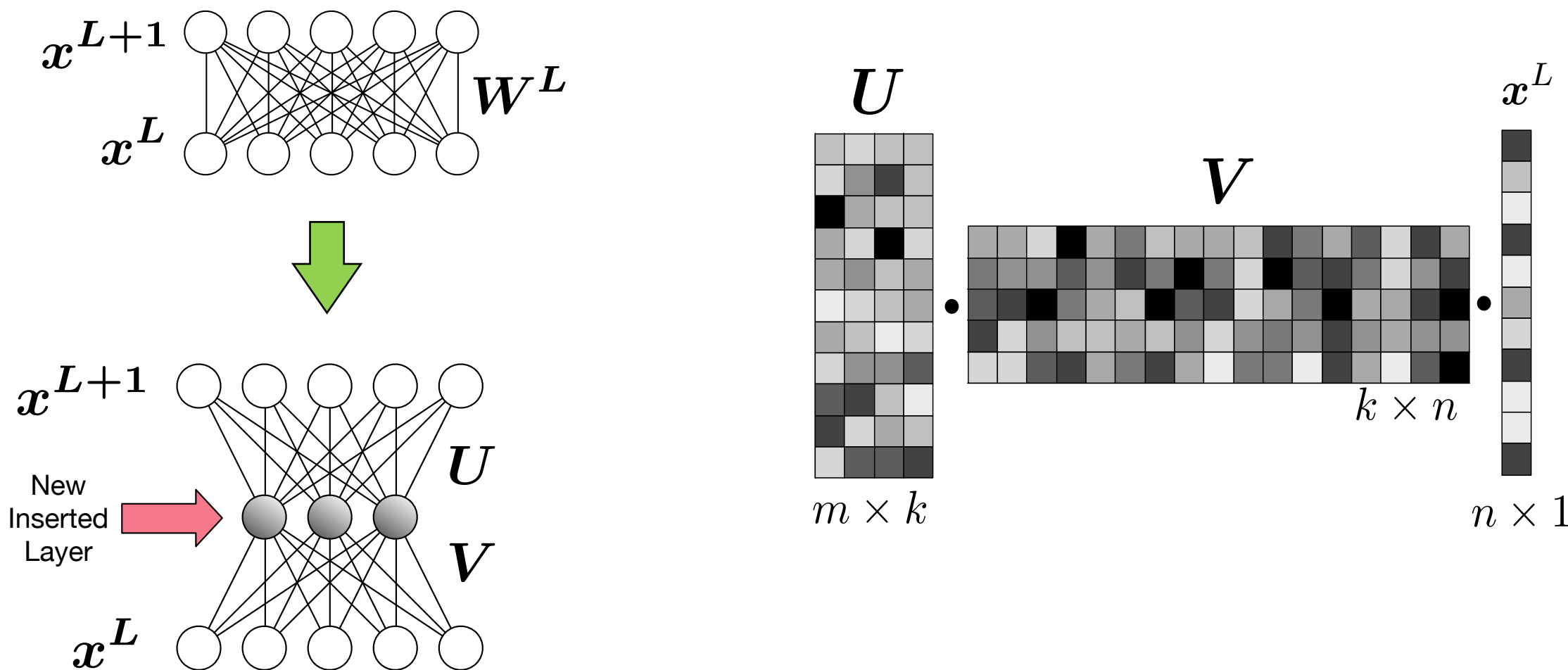
Total Operations:  $m \times n \times 1$

Total Operations:  $m \times k \times 1 + k \times n \times 1$

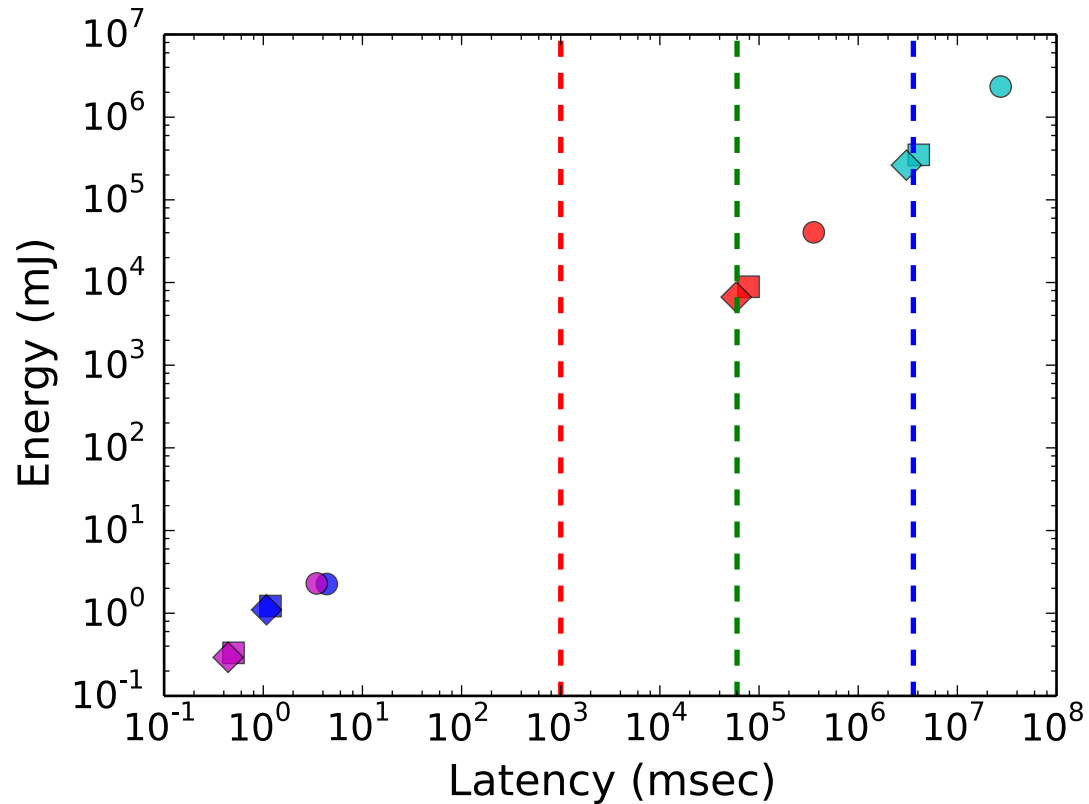
Gain in memory and computations when:  $k < \frac{m \times n}{m + n}$



# Matrix Decomposition

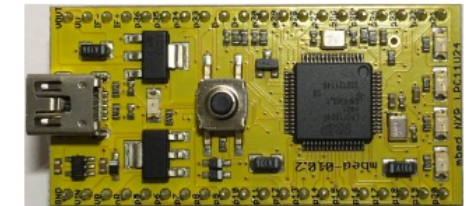
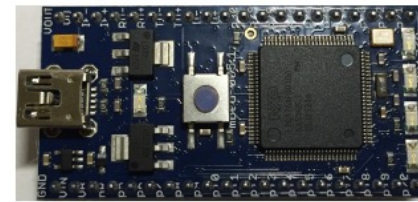


# Decomposition Benefits



32 KB

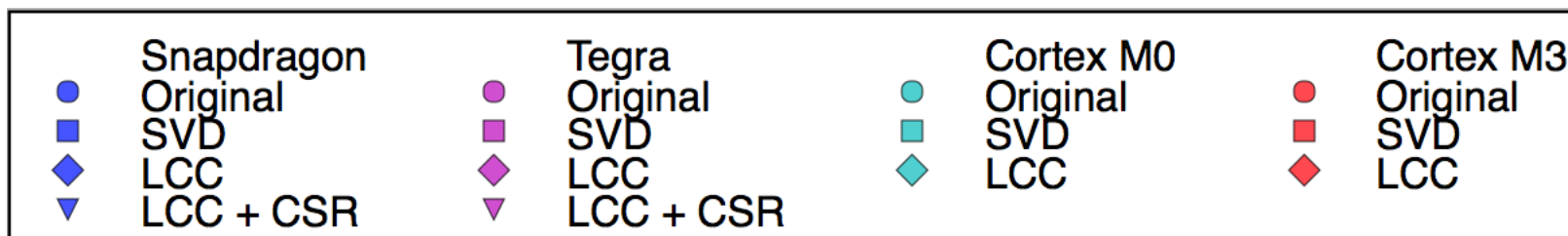
16 KB



ARM Cortex M3

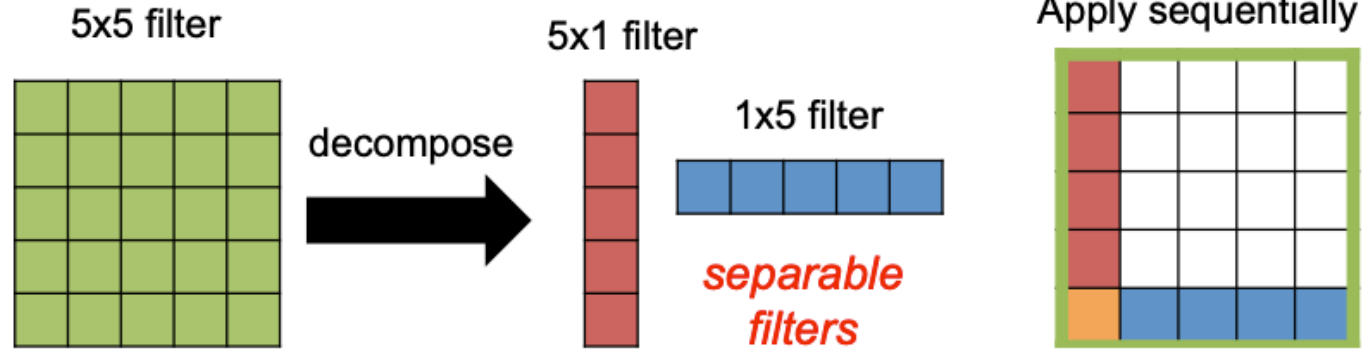
ARM Cortex M0

2-4% degradation in accuracy

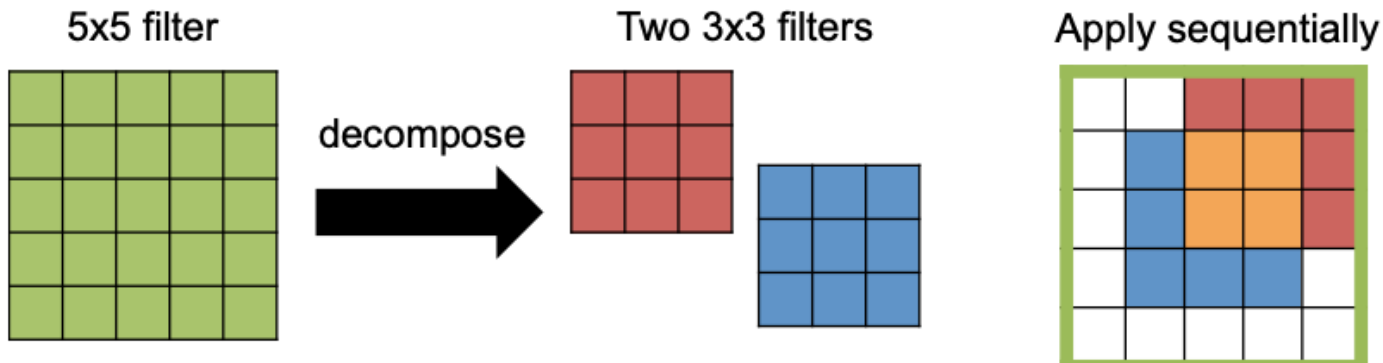


# CNN Kernel Decomposition

## Inception v3 (2014)



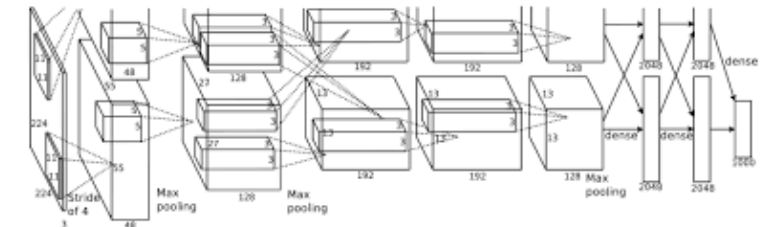
Replacing a large filter with a series of small filters



Source: V. Sze

## VGG (2014)

## AlexNet (2012)



# Roadmap for Today

---

## Architecture (*Re-*)Design

1. Early Evolution; Decomposition as a tool
- 2. Efficient Architectures; Example: MobileNet**

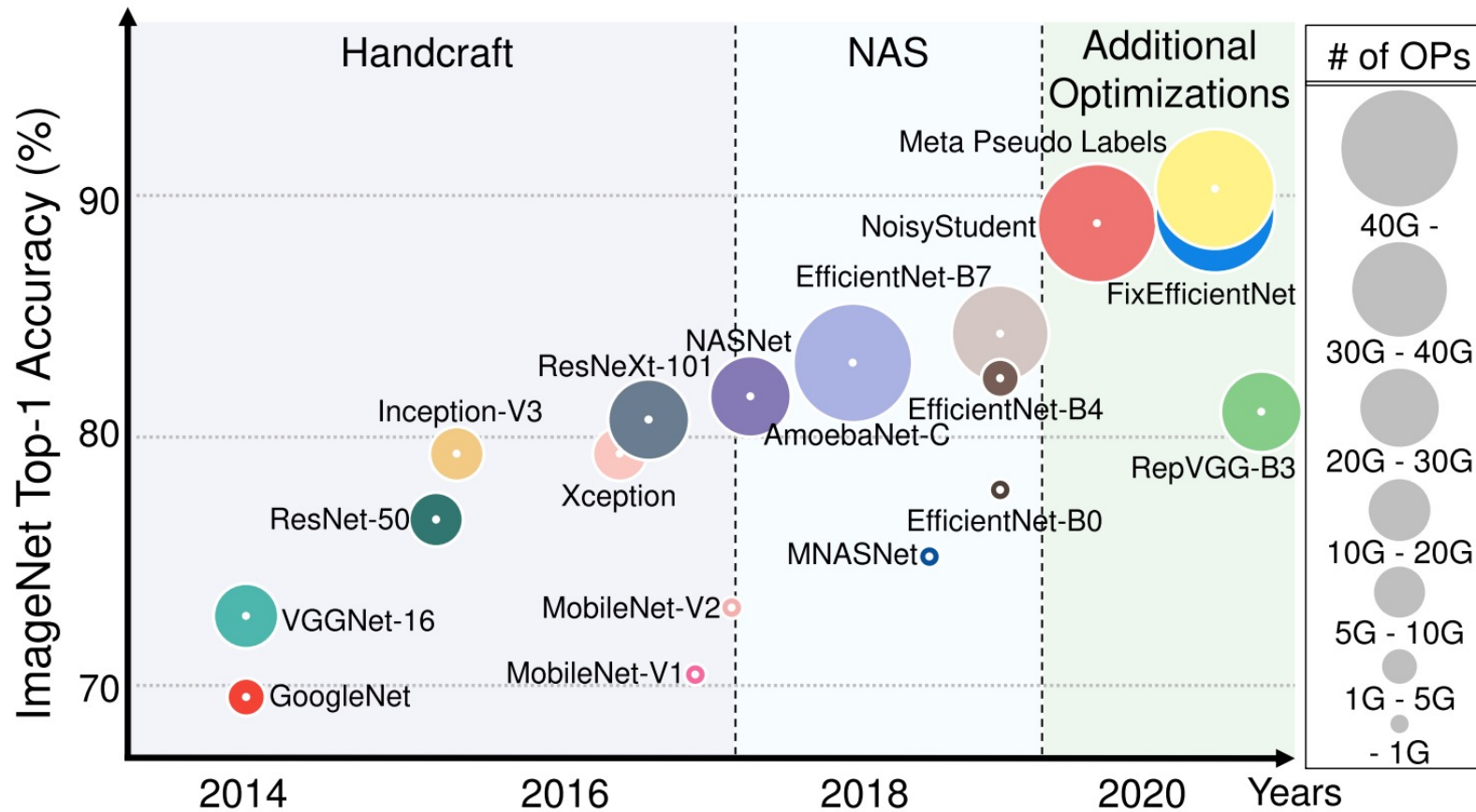
## Architecture Compression

3. Parameter and Channel Pruning
4. Parameter Quantization
5. Knowledge Distillation
6. Compression pipelines





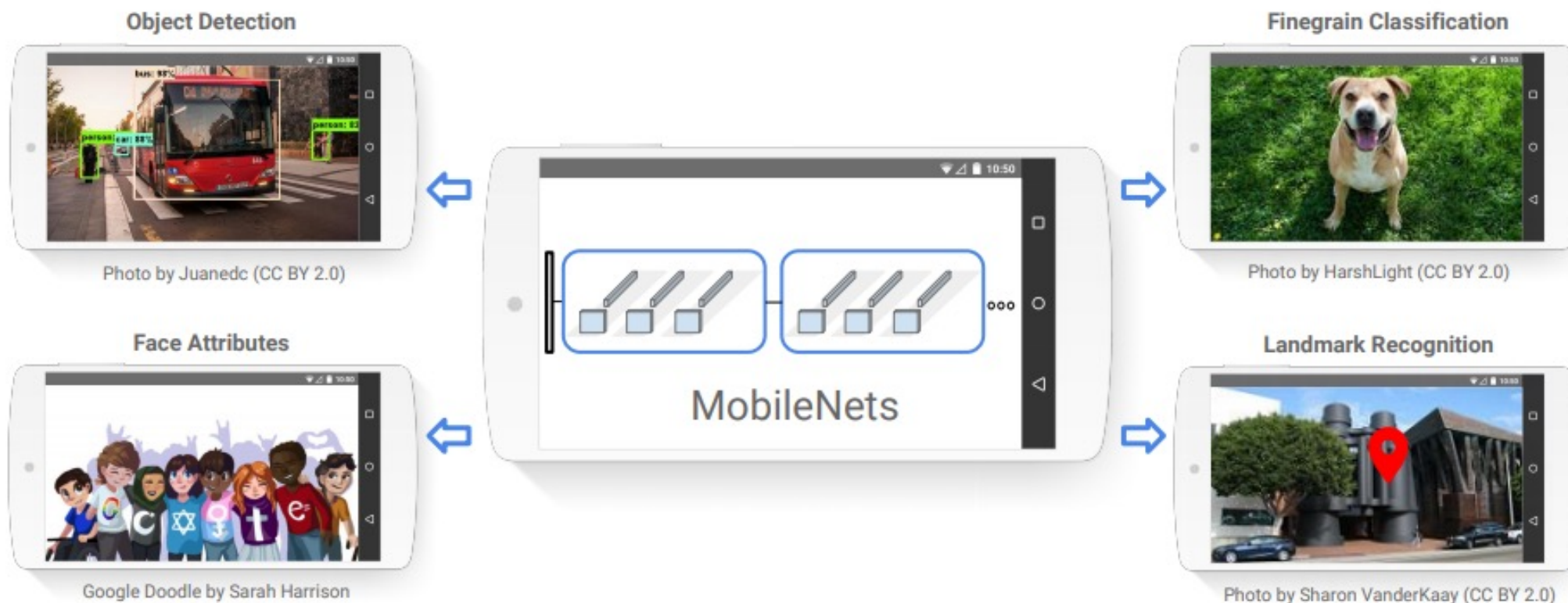
# Efficient Architecture Evolution



Source: Lee



# MobileNet Example

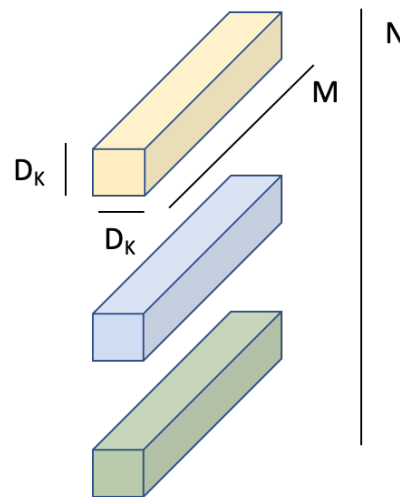


Source: Google

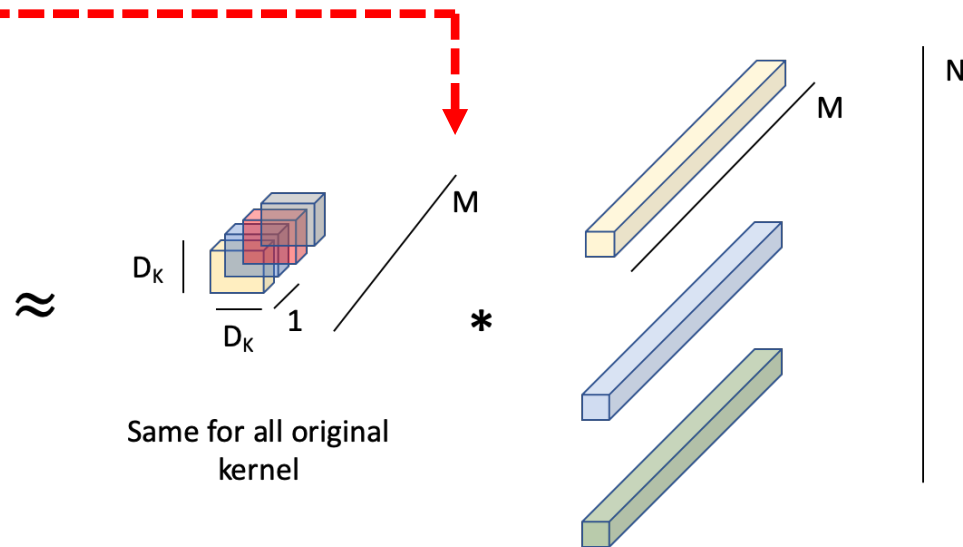


# MobileNet Block

Type / Stride	Filter Shape	Input Size
Conv / s2	$3 \times 3 \times 3 \times 32$	$224 \times 224 \times 3$
Conv dw / s1	$3 \times 3 \times 32$ dw	$112 \times 112 \times 32$
Conv / s1	$1 \times 1 \times 32 \times 64$	$112 \times 112 \times 32$
Conv dw / s2	$3 \times 3 \times 64$ dw	$112 \times 112 \times 64$
Conv / s1	$1 \times 1 \times 64 \times 128$	$56 \times 56 \times 64$
Conv dw / s1	$3 \times 3 \times 128$ dw	$56 \times 56 \times 128$
Conv / s1	$1 \times 1 \times 128 \times 128$	$56 \times 56 \times 128$
Conv dw / s2	$3 \times 3 \times 128$ dw	$56 \times 56 \times 128$
Conv / s1	$1 \times 1 \times 128 \times 256$	$28 \times 28 \times 128$
Conv dw / s1	$3 \times 3 \times 256$ dw	$28 \times 28 \times 256$
Conv / s1	$1 \times 1 \times 256 \times 256$	$28 \times 28 \times 256$
Conv dw / s2	$3 \times 3 \times 256$ dw	$28 \times 28 \times 256$
Conv / s1	$1 \times 1 \times 256 \times 512$	$14 \times 14 \times 256$
$5 \times$ Conv dw / s1	$3 \times 3 \times 512$ dw	$14 \times 14 \times 512$
Conv / s1	$1 \times 1 \times 512 \times 512$	$14 \times 14 \times 512$
Conv dw / s2	$3 \times 3 \times 512$ dw	$14 \times 14 \times 512$
Conv / s1	$1 \times 1 \times 512 \times 1024$	$7 \times 7 \times 512$
Conv dw / s2	$3 \times 3 \times 1024$ dw	$7 \times 7 \times 1024$
Conv / s1	$1 \times 1 \times 1024 \times 1024$	$7 \times 7 \times 1024$
Avg Pool / s1	Pool $7 \times 7$	$7 \times 7 \times 1024$
FC / s1	$1024 \times 1000$	$1 \times 1 \times 1024$
Softmax / s1	Classifier	$1 \times 1 \times 1000$



Original Independent Kernels



One for each kernel

Source: Howard



# MobileNet Benefits

Model	ImageNet Accuracy	Million Mult-Adds	Million Parameters
1.0 MobileNet-224	70.6%	569	4.2
GoogLeNet	69.8%	1550	6.8
VGG 16	71.5%	15300	138

Model	ImageNet Accuracy	Million Mult-Adds	Million Parameters
Conv MobileNet	71.7%	4866	29.3
MobileNet	70.6%	569	4.2

Source: Howard

## Parameter Reduction

$$\frac{\text{Depthwise Separable}}{\text{Regular Conv}} = \frac{M \cdot (N + DK^2)}{D_K^2 \cdot M \cdot N} = \frac{1}{D_K^2} + \frac{1}{N}$$

## MACs Reduction

$$\frac{\text{Depthwise Separable}}{\text{Regular Conv}} = \frac{D_F^2 \cdot M \cdot (N + DK^2)}{D_K^2 \cdot M \cdot D_F^2 \cdot N} = \frac{1}{D_K^2} + \frac{1}{N}$$





# Roadmap for Today

---

## Architecture (*Re-*)Design

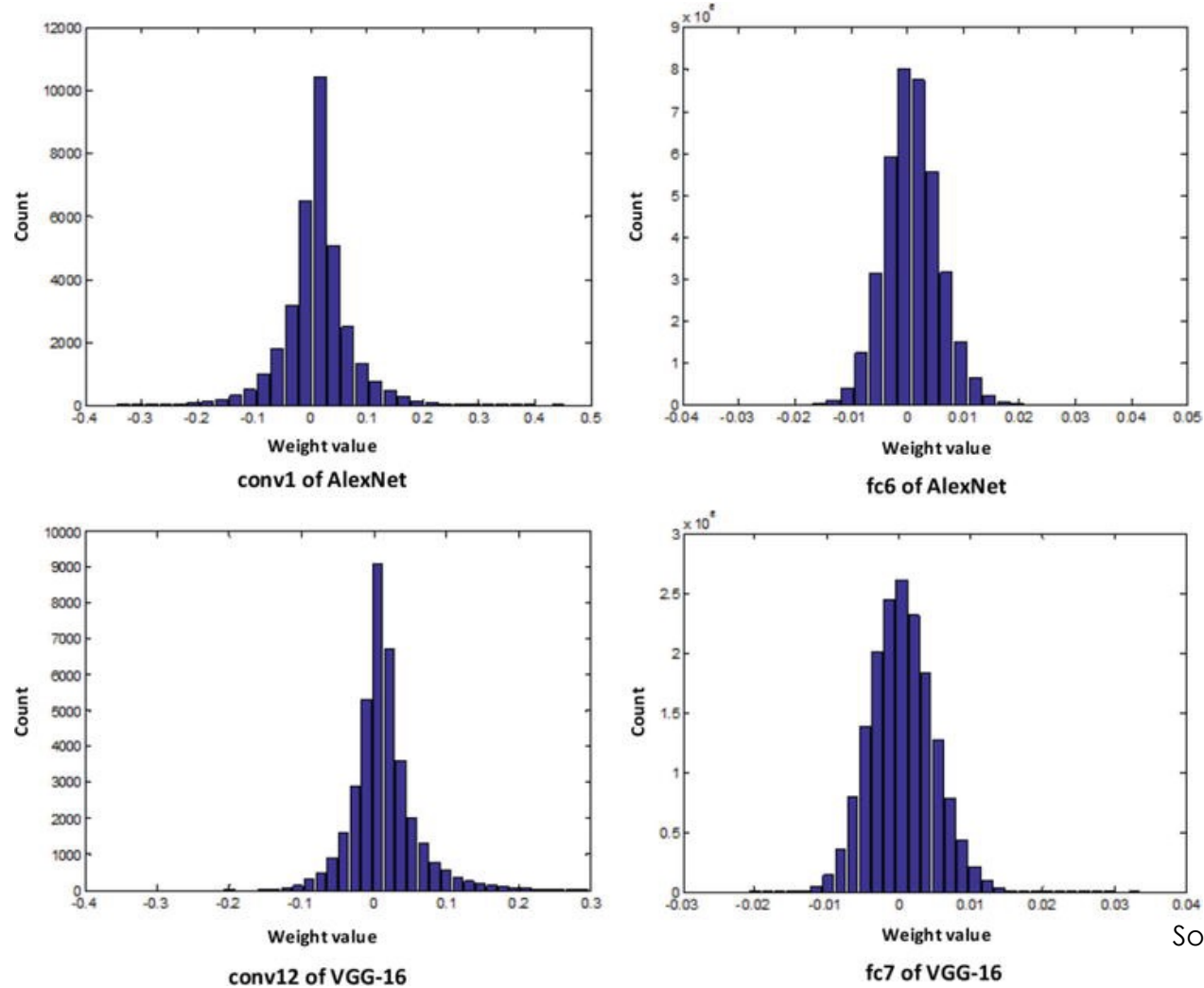
1. Early Evolution; Decomposition as a tool
2. Efficient Architectures; Example: MobileNet

## Architecture Compression

- 3. Parameter and Channel Pruning**
4. Parameter Quantization
5. Knowledge Distillation
6. Compression pipelines



# Pruning Justification

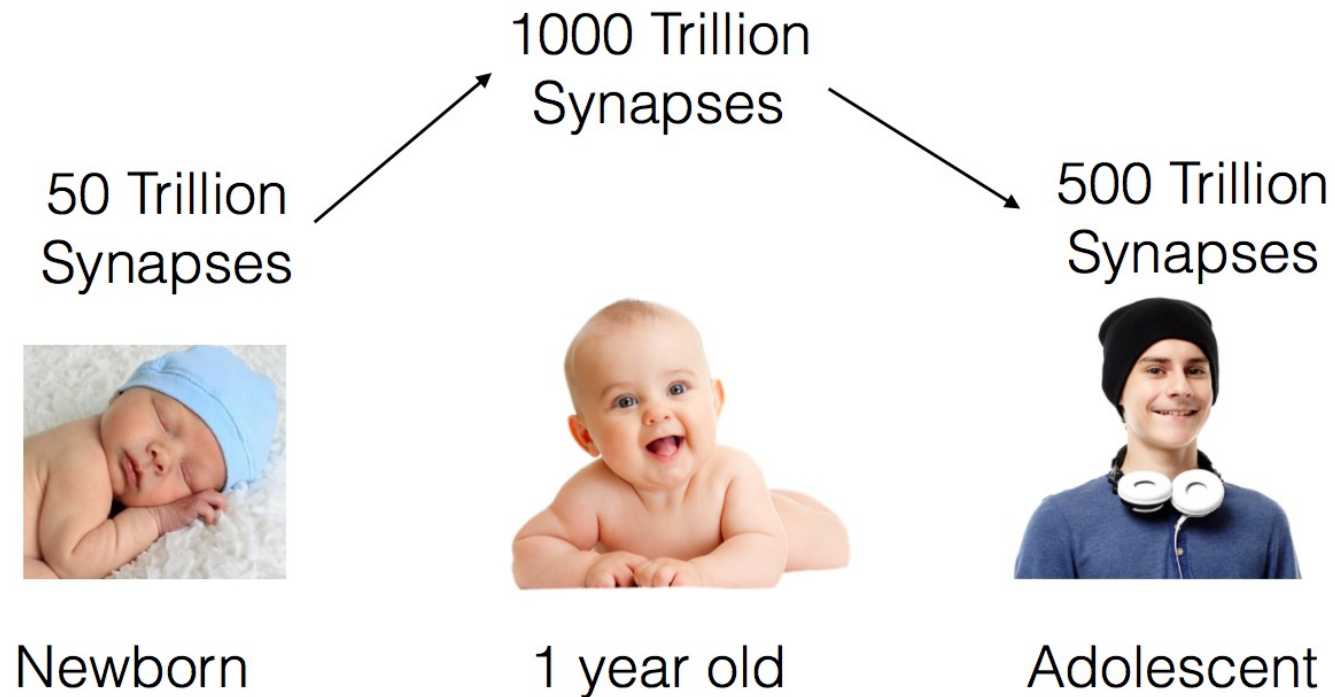


Source: Han



# Related Natural Phenomena

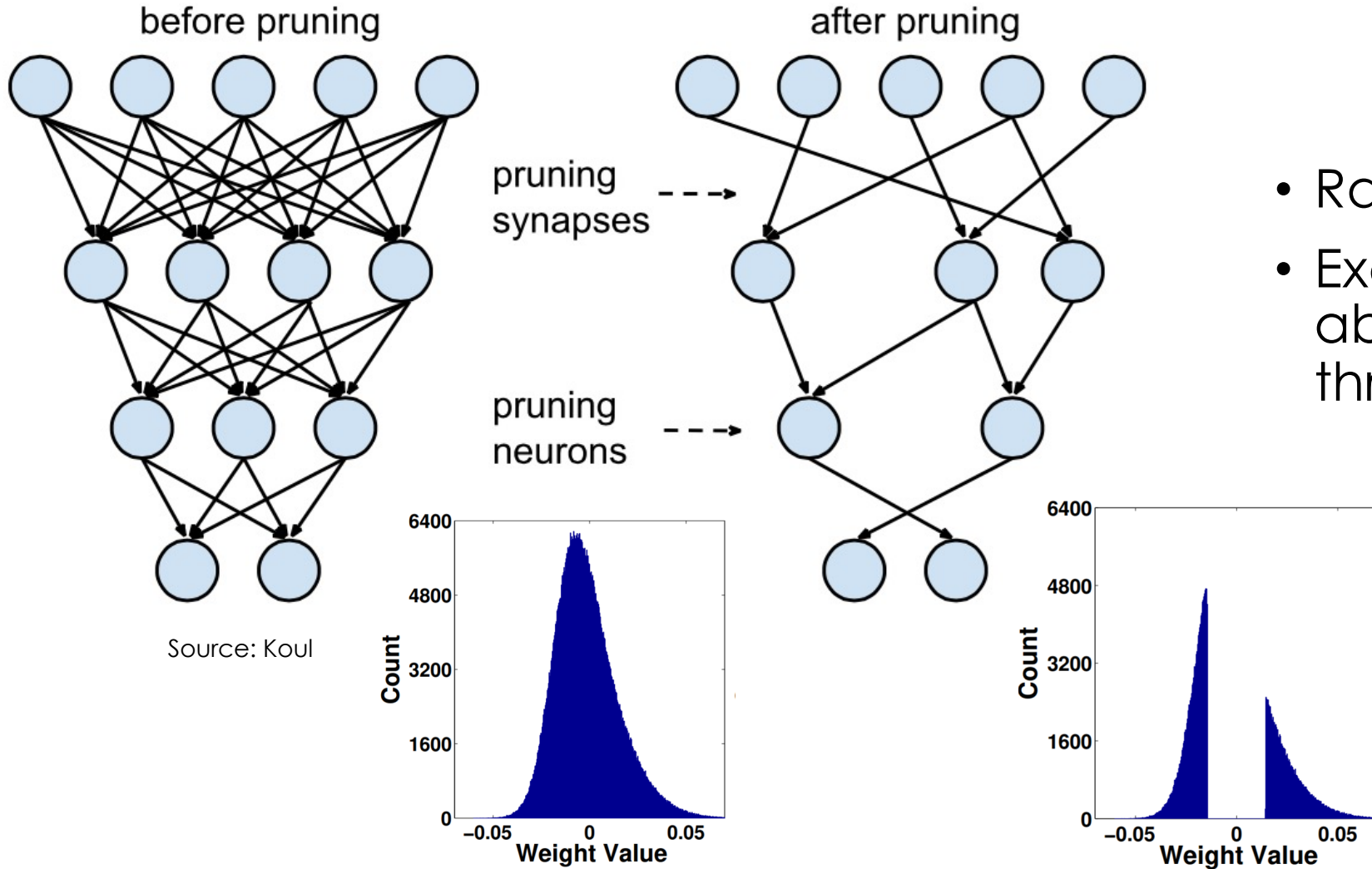
---



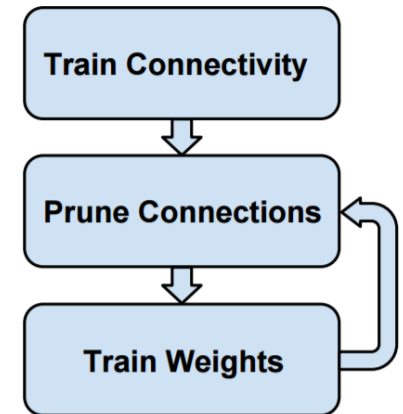
Source: Han



# Pruning Algorithm

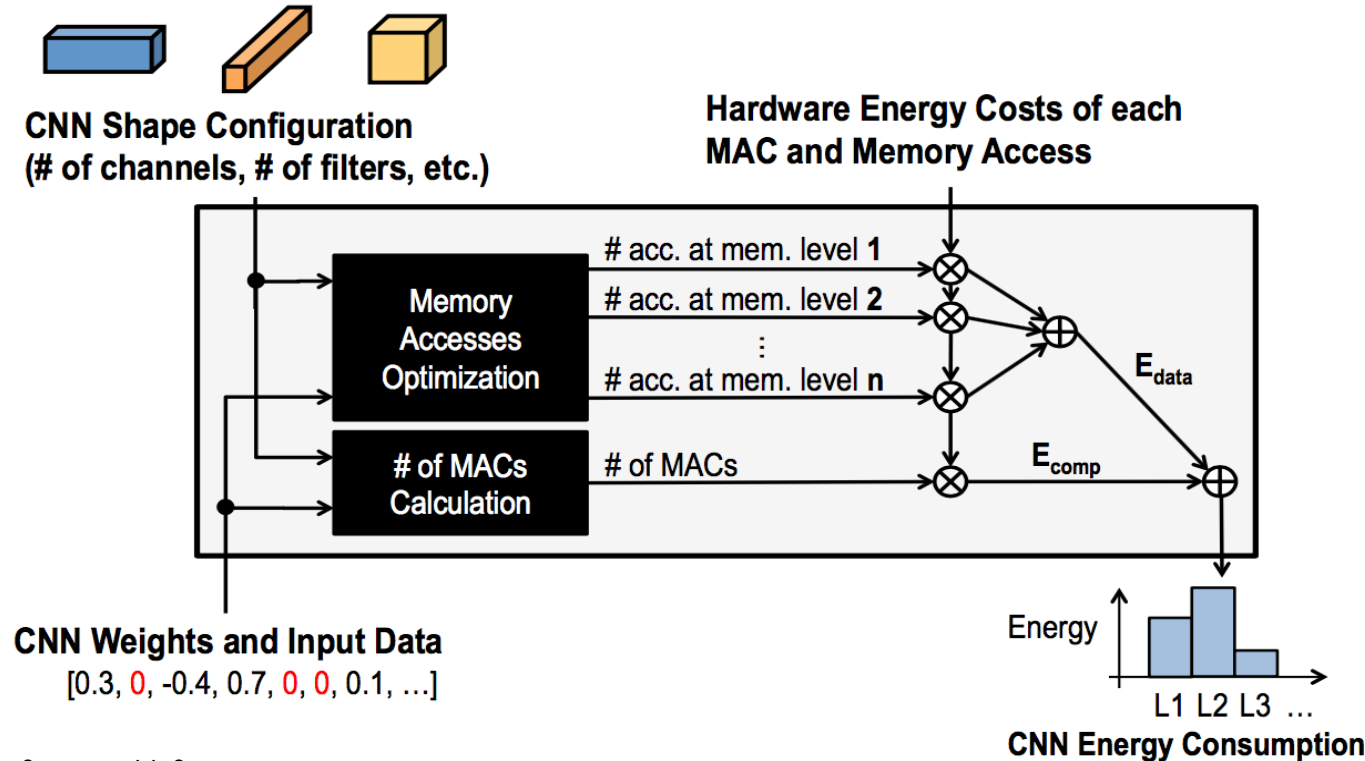


- Range of pruning criteria
- Example: Prune when absolute weight is  $<$  a threshold





# Alternative Pruning Criteria

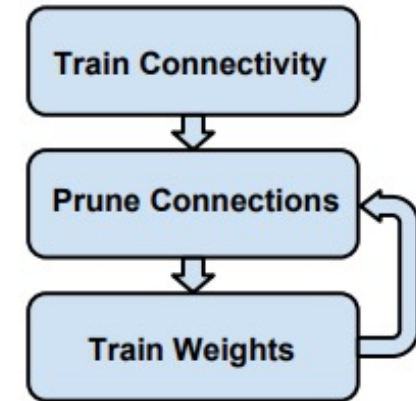
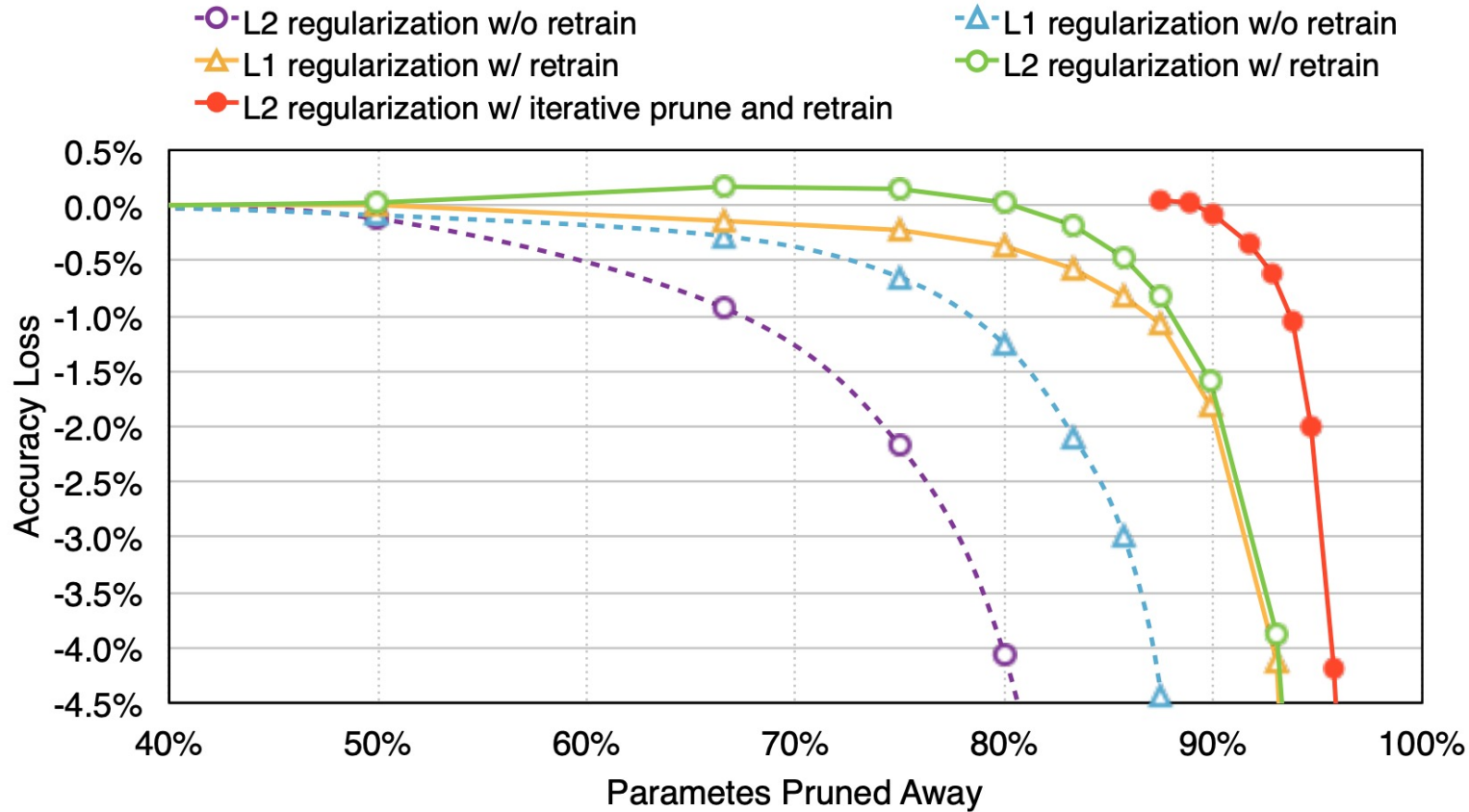


Source: V. Sze

- Magnitude
- Energy
  - Estimation
  - Measurement
- Random
- etc.



# Pruning Schedule

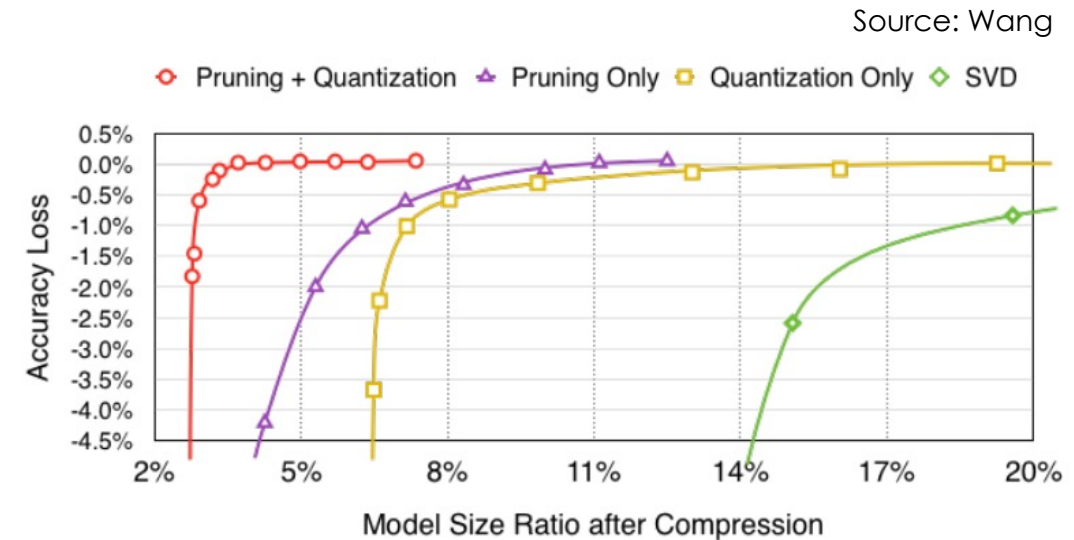


Source: Han



# Pruning Performance

Dataset	Model	Unpruned	Prune Ratio	Fine-tuned	Scratch-E	Scratch-B
CIFAR-10	VGG-19	93.50 ( $\pm 0.11$ )	30%	93.51 ( $\pm 0.05$ )	<b>93.71</b> ( $\pm 0.09$ )	93.31 ( $\pm 0.26$ )
			80%	93.52 ( $\pm 0.10$ )	<b>93.71</b> ( $\pm 0.08$ )	93.64 ( $\pm 0.09$ )
			95%	93.34 ( $\pm 0.13$ )	93.21 ( $\pm 0.17$ )	<b>93.63</b> ( $\pm 0.18$ )
	PreResNet-110	95.04 ( $\pm 0.15$ )	30%	95.06 ( $\pm 0.05$ )	94.84 ( $\pm 0.07$ )	<b>95.11</b> ( $\pm 0.09$ )
			80%	<b>94.55</b> ( $\pm 0.11$ )	93.76 ( $\pm 0.10$ )	94.52 ( $\pm 0.13$ )
			95%	<b>92.35</b> ( $\pm 0.20$ )	91.23 ( $\pm 0.11$ )	91.55 ( $\pm 0.34$ )
	DenseNet-BC-100	95.24 ( $\pm 0.17$ )	30%	95.21 ( $\pm 0.17$ )	95.22 ( $\pm 0.18$ )	<b>95.23</b> ( $\pm 0.14$ )
			80%	95.04 ( $\pm 0.15$ )	94.42 ( $\pm 0.12$ )	<b>95.12</b> ( $\pm 0.04$ )
			95%	<b>94.19</b> ( $\pm 0.15$ )	92.91 ( $\pm 0.22$ )	93.44 ( $\pm 0.19$ )
CIFAR-100	VGG-19	71.70 ( $\pm 0.31$ )	30%	71.96 ( $\pm 0.36$ )	72.81 ( $\pm 0.31$ )	<b>73.30</b> ( $\pm 0.25$ )
			50%	71.85 ( $\pm 0.30$ )	73.12 ( $\pm 0.36$ )	<b>73.77</b> ( $\pm 0.23$ )
			95%	70.22 ( $\pm 0.38$ )	70.88 ( $\pm 0.35$ )	<b>72.08</b> ( $\pm 0.15$ )
	PreResNet-110	76.96 ( $\pm 0.34$ )	30%	76.88 ( $\pm 0.31$ )	76.36 ( $\pm 0.26$ )	<b>76.96</b> ( $\pm 0.31$ )
			50%	<b>76.60</b> ( $\pm 0.36$ )	75.45 ( $\pm 0.23$ )	76.42 ( $\pm 0.39$ )
			95%	68.55 ( $\pm 0.51$ )	68.13 ( $\pm 0.64$ )	<b>68.99</b> ( $\pm 0.32$ )
	DenseNet-BC-100	77.59 ( $\pm 0.19$ )	30%	77.23 ( $\pm 0.05$ )	77.58 ( $\pm 0.25$ )	<b>77.97</b> ( $\pm 0.31$ )
			50%	77.41 ( $\pm 0.14$ )	77.65 ( $\pm 0.09$ )	<b>77.80</b> ( $\pm 0.23$ )
			95%	<b>73.67</b> ( $\pm 0.03$ )	71.47 ( $\pm 0.46$ )	72.57 ( $\pm 0.37$ )
ImageNet	VGG-16	73.37	30%	73.68	72.75	<b>74.02</b>
			60%	<b>73.63</b>	71.50	73.42
	ResNet-50	76.15	30%	<b>76.06</b>	74.77	75.70
			60%	<b>76.09</b>	73.69	74.91



# Roadmap for Today

---

## Architecture (*Re-*)Design

1. Early Evolution; Decomposition as a tool
2. Efficient Architectures; Example: MobileNet

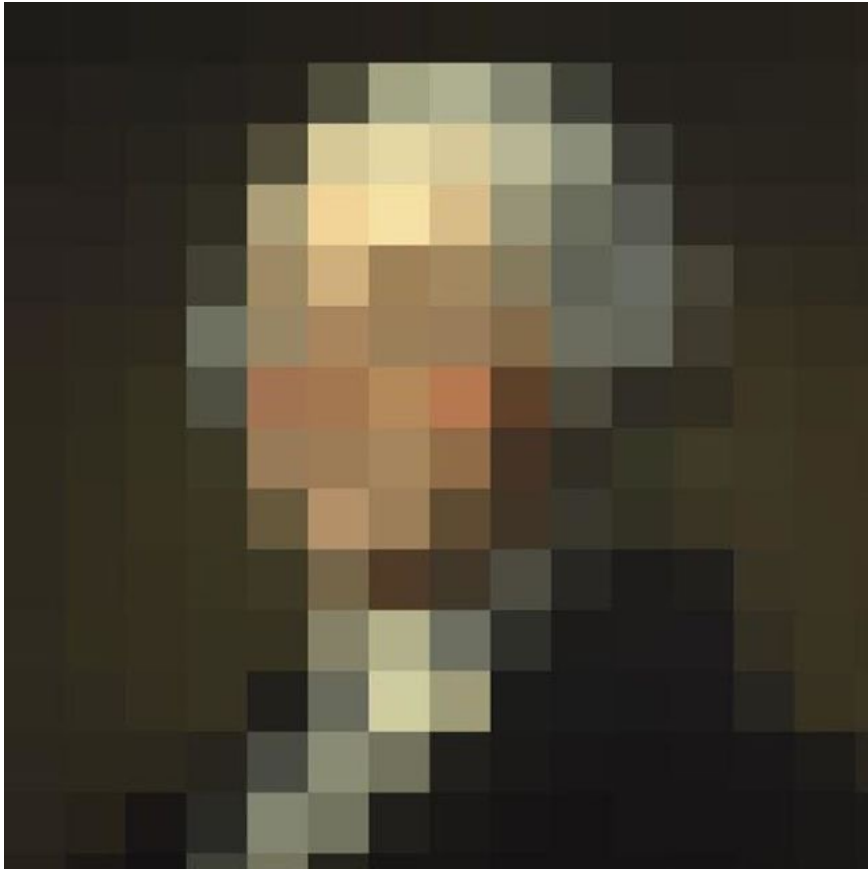
## Architecture Compression

3. Parameter and Channel Pruning
- 4. Parameter Quantization**
5. Knowledge Distillation
6. Compression pipelines



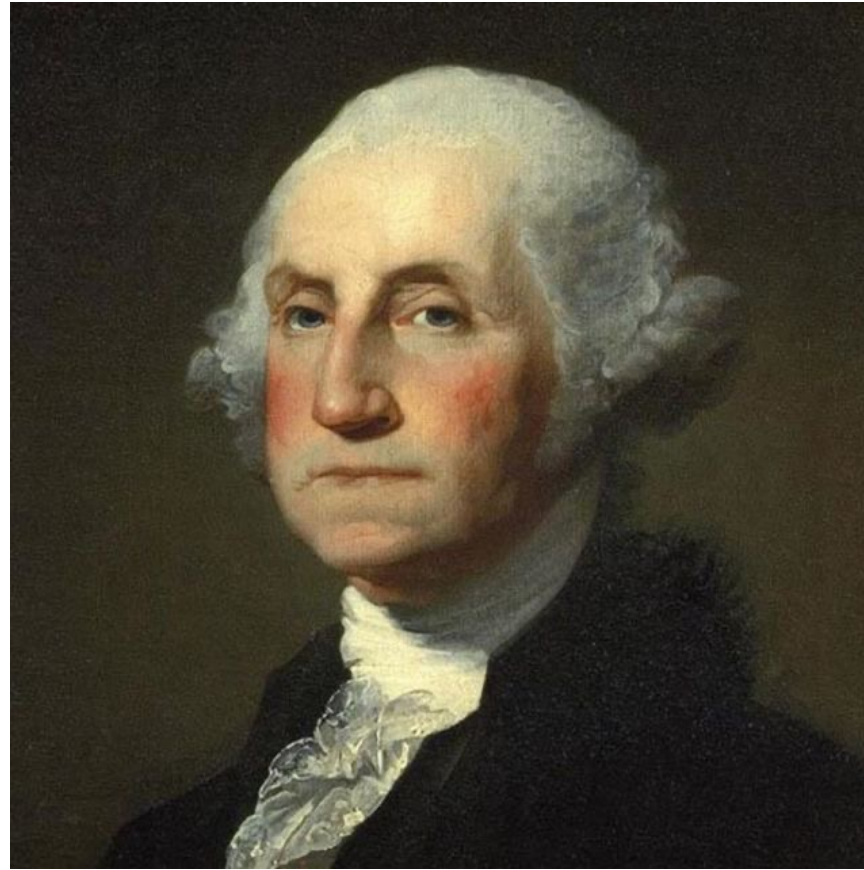
# Guess who

---



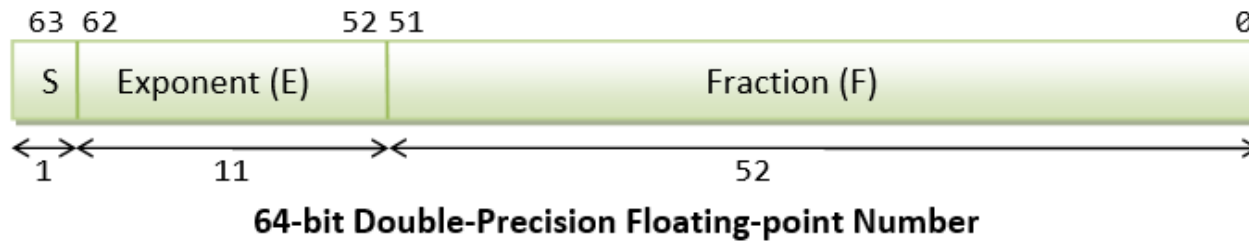
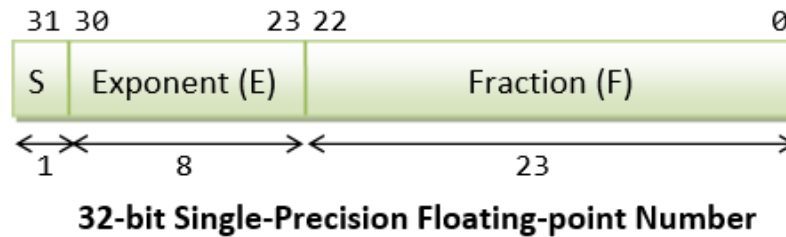
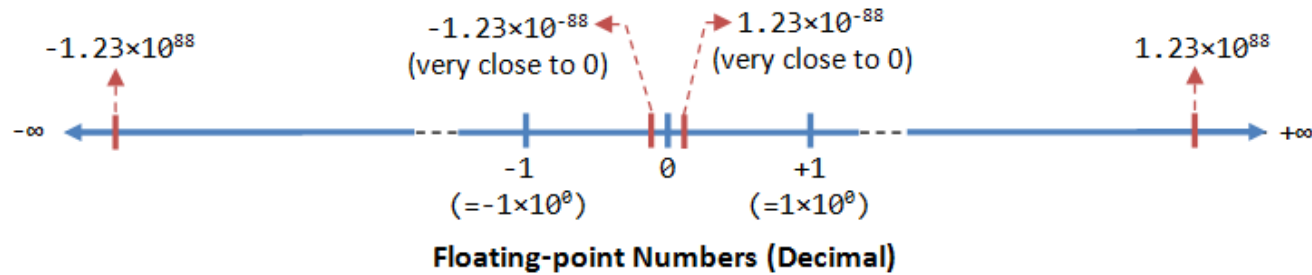
# Guess who

---

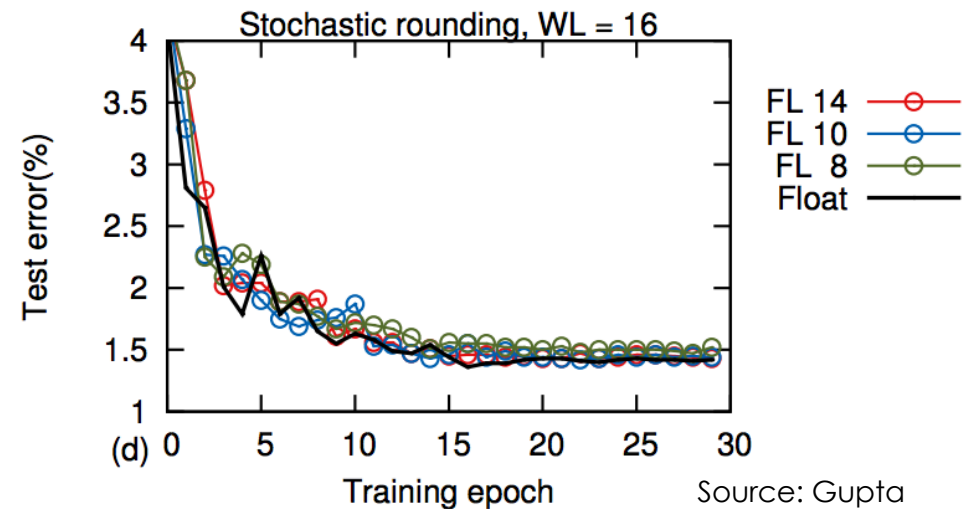
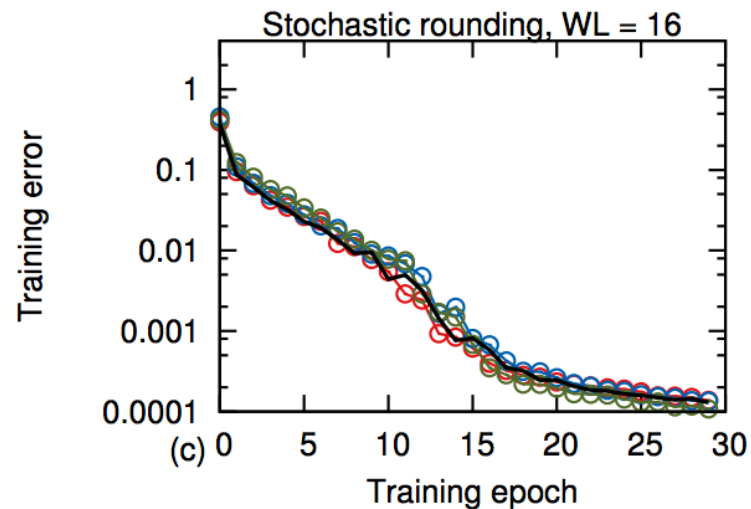
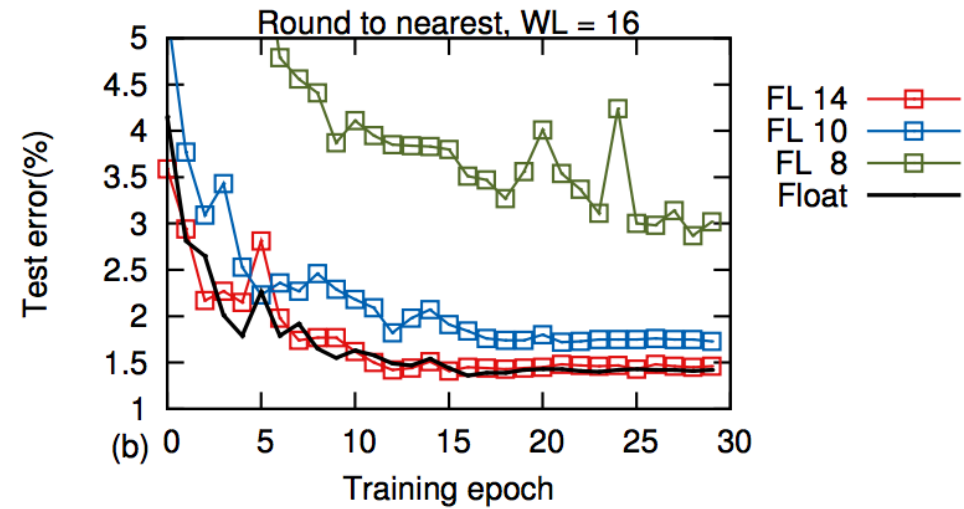
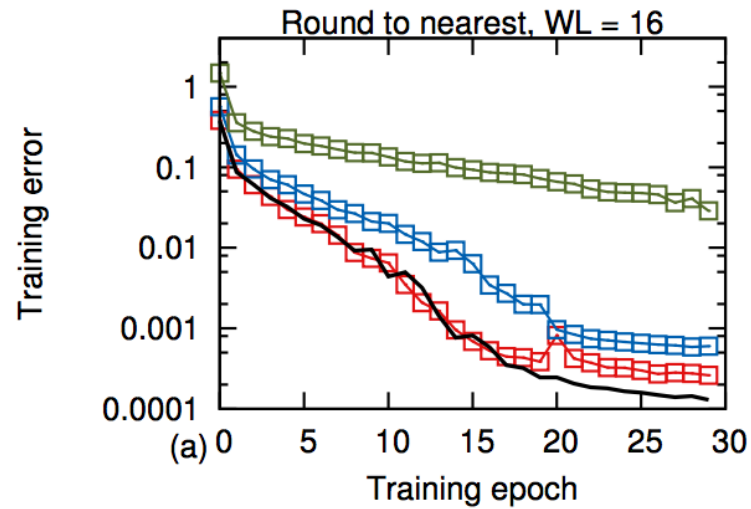




# Numerical Representation



# Low Precision Training

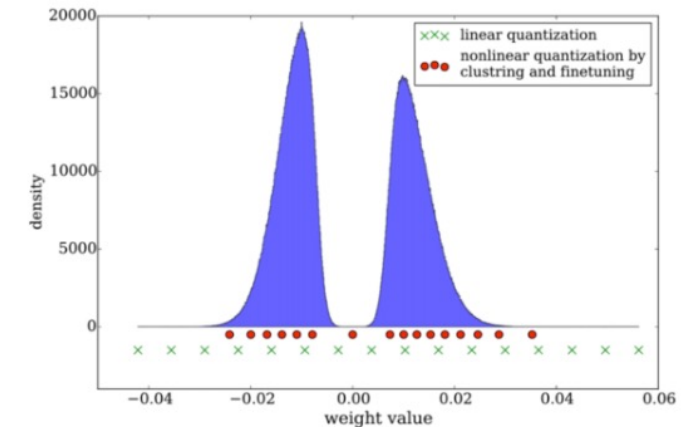


Source: Gupta



# Quantization Implementation

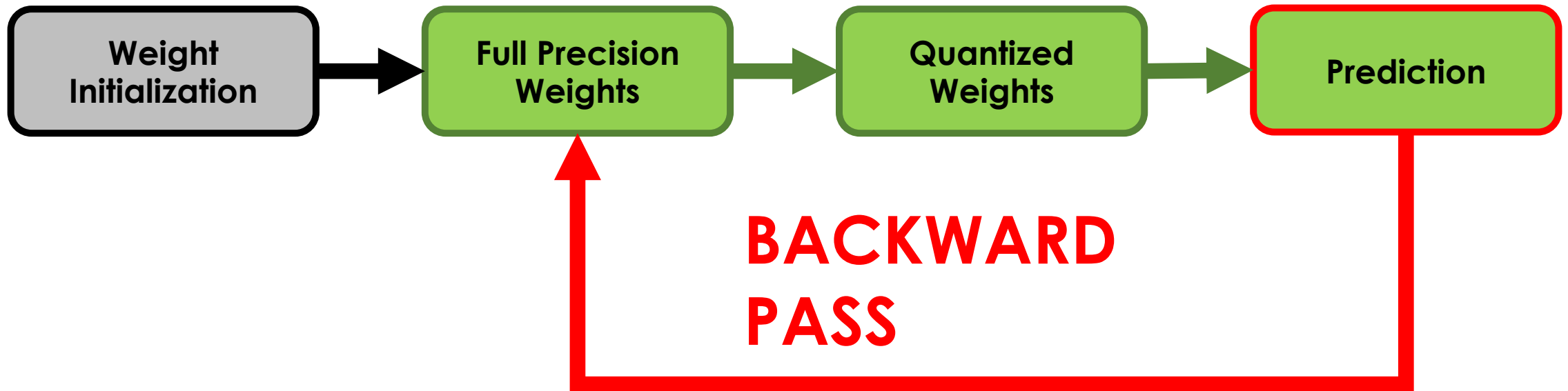
- Linear vs. Non-linear
  - Linear – representable space is divided equally “X”
  - Non-linear – representable space is divided un-equally “O”
- Round to nearest vs. Stochastic rounding
  - Round to nearest – each number is represented by the closest representable number
  - Stochastic rounding – each number is represented by the closest higher or lower value with equal probabilities



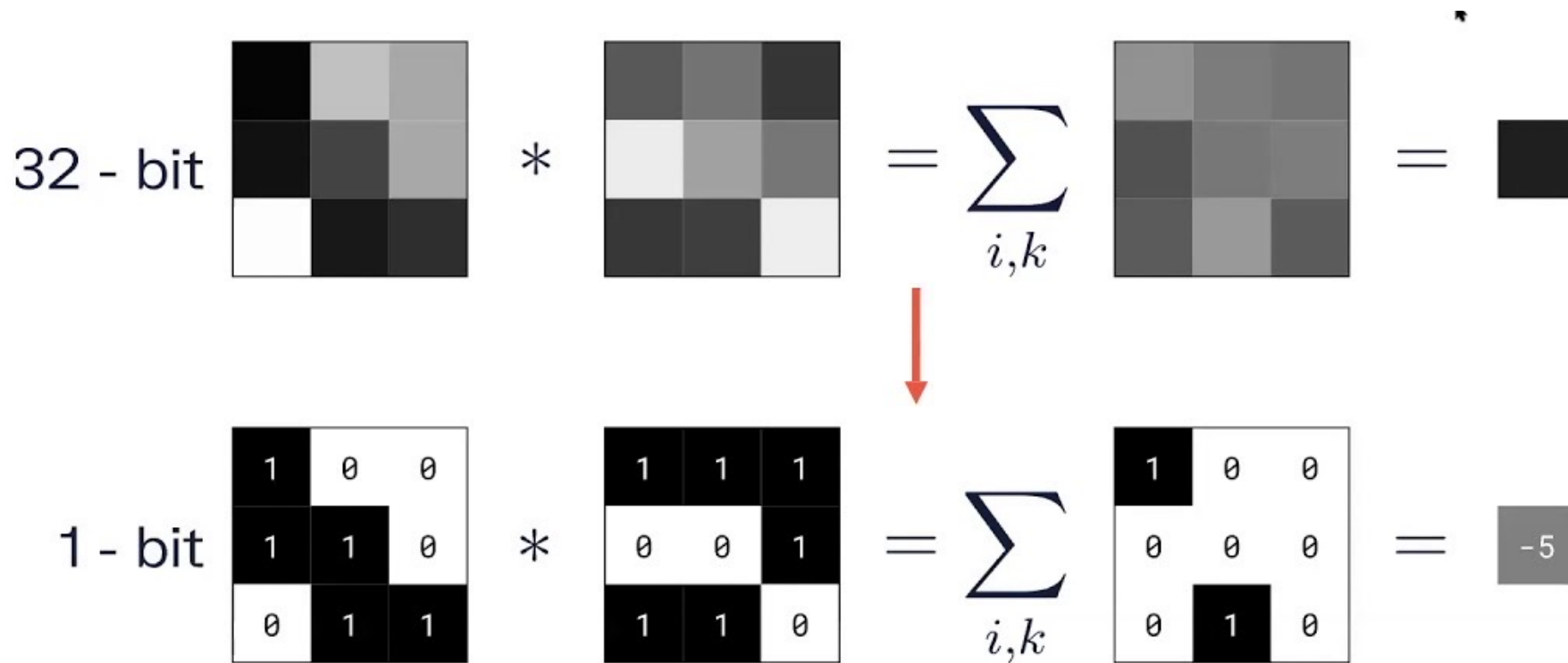
# Training Quantized Models

---

**FORWARD  
PASS**



# Binary: Extreme Quantization



Space  
Reduction  
32x

Compute  
Reduction  
No FP Ops!

Source: Geiger



# Roadmap for Today

---

## Architecture (*Re-*)Design

1. Early Evolution; Decomposition as a tool
2. Efficient Architectures; Example: MobileNet

## Architecture Compression

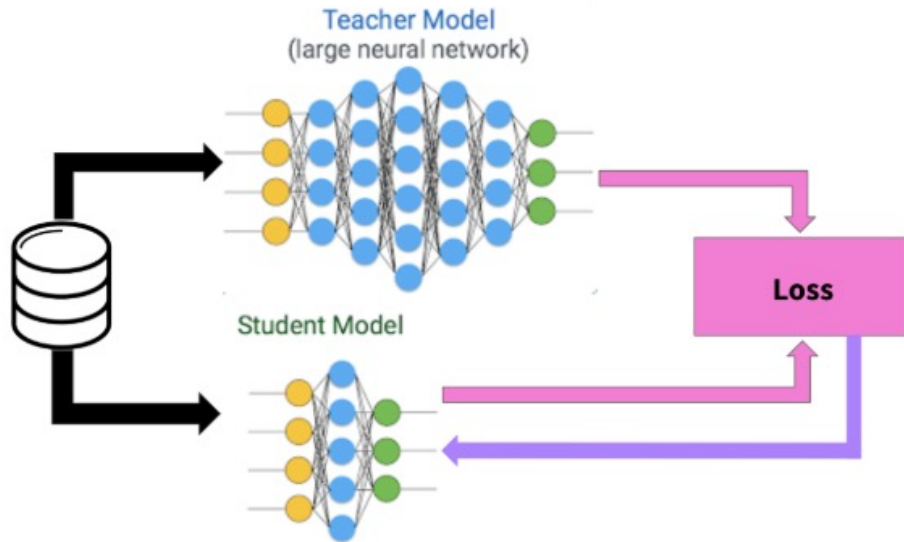
3. Parameter and Channel Pruning
4. Parameter Quantization
- 5. Knowledge Distillation**
6. Compression pipelines





# Knowledge Distillation

---

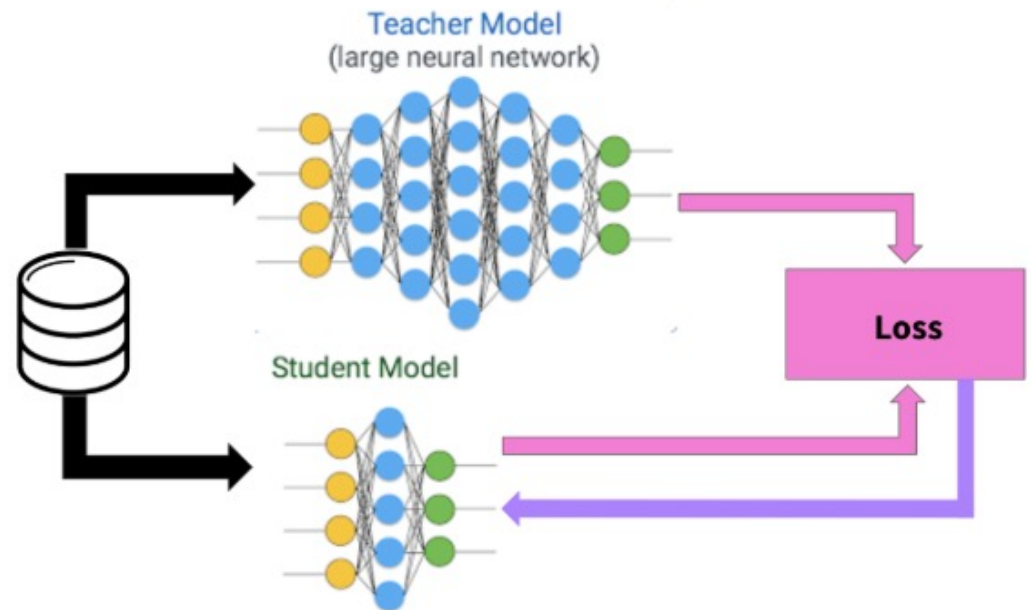


Distillation of a large model into a small one is performed in a teacher-student setup whereby the student network is trained to **replicate the teacher's raw output** (logits).



# Knowledge Distillation

- **Knowledge distillation does not require actual data for the student's training** since the goal is for it to match the teacher's output not the implied labels.
- Thus, students can be trained on random inputs just as well as on the original data! **The teacher provides the targets in real time.**



# Knowledge Distillation

#	Model	SST-2	QQP	MNLI-m	MNLI-mm
		Acc	F <sub>1</sub> /Acc	Acc	Acc
1	BERT <sub>LARGE</sub> (Devlin et al., 2018)	94.9	72.1/89.3	86.7	85.9
2	BERT <sub>BASE</sub> (Devlin et al., 2018)	93.5	71.2/89.2	84.6	83.4
3	OpenAI GPT (Radford et al., 2018)	91.3	70.3/88.5	82.1	81.4
4	BERT ELMo baseline (Devlin et al., 2018)	90.4	64.8/84.7	76.4	76.1
5	GLUE ELMo baseline (Wang et al., 2018)	90.4	63.1/84.3	74.1	74.5
6	Distilled BiLSTM <sub>SOFT</sub>	<b>90.7</b>	<b>68.2/88.1</b>	<b>73.0</b>	<b>72.6</b>
7	BiLSTM (our implementation)	86.7	63.7/86.2	68.7	68.3
8	BiLSTM (reported by GLUE)	85.9	61.4/81.7	70.3	70.8
9	BiLSTM (reported by other papers)	87.6 <sup>†</sup>	– /82.6 <sup>‡</sup>	66.9 <sup>*</sup>	66.9 <sup>*</sup>

Source: Tang

Achieve comparable results with ELMo, while using roughly 100 times fewer parameters and 15 times less inference time.



# Roadmap for Today

---

## Architecture (*Re-*)Design

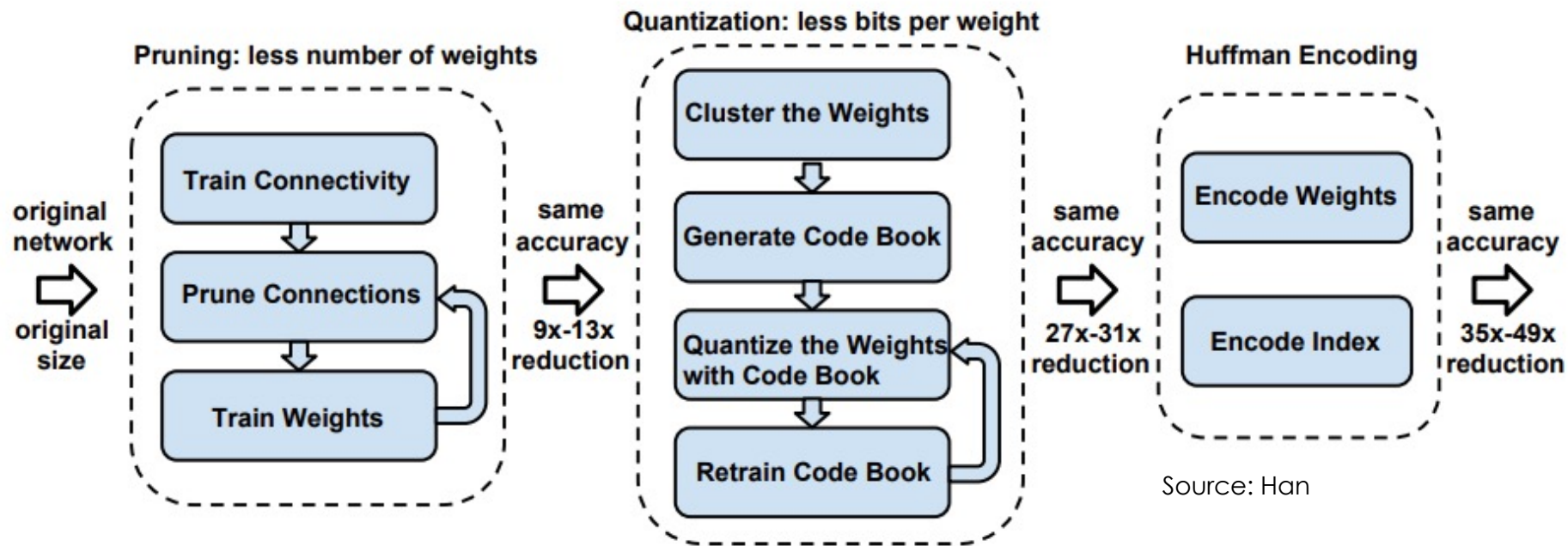
1. Early Evolution; Decomposition as a tool
2. Efficient Architectures; Example: MobileNet

## Architecture Compression

3. Parameter and Channel Pruning
4. Parameter Quantization
5. Knowledge Distillation
- 6. Compression pipelines**

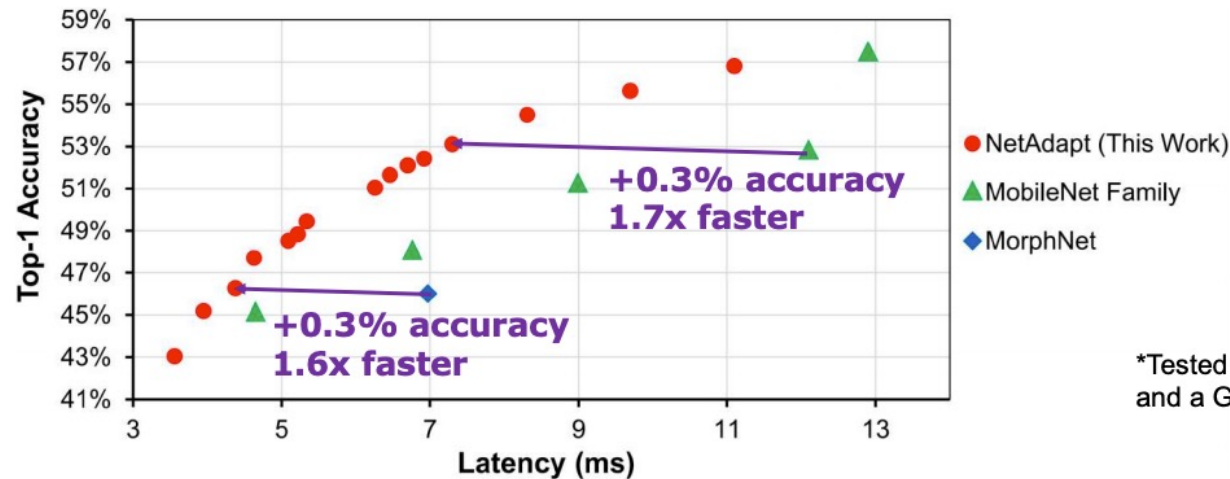
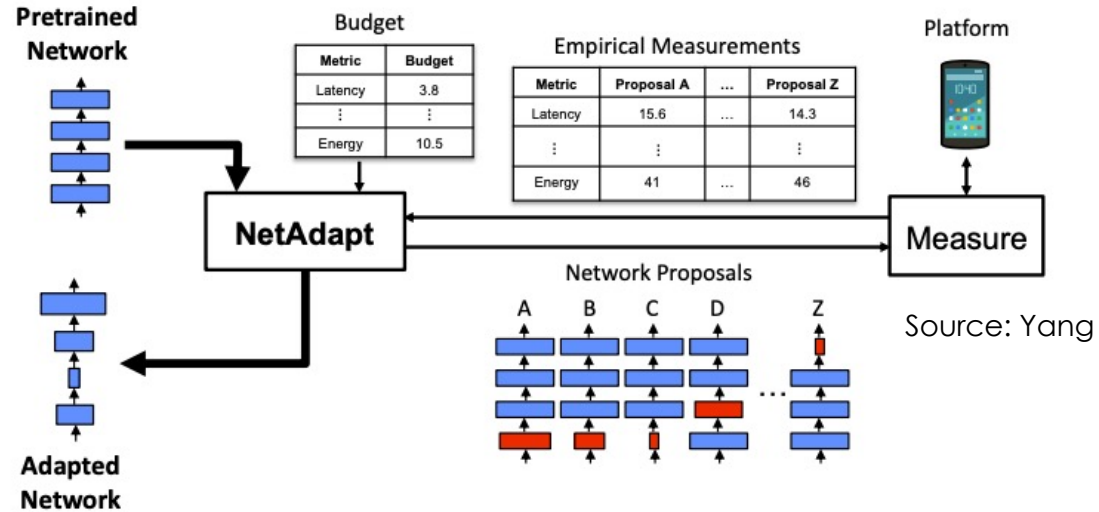


# Basic Compression Pipeline



# Resource-aware Pipeline

- **Automatically adapt DNN** to a mobile platform to reach a target latency or energy budget
- Use **empirical measurements** to guide optimization (avoid modeling of tool chain or platform architecture)
- Requires **very few hyperparameters** to tune



\*Tested on the ImageNet dataset and a Google Pixel 1 CPU





# Summary of the Day

---

## Architecture (*Re-*)Design

1. Early Evolution; Decomposition as a tool
2. Efficient Architectures; Example: MobileNet

## Architecture Compression

3. Parameter and Channel Pruning
4. Parameter Quantization
5. Knowledge Distillation
6. Compression pipelines

