# AI on Edge:
# **Neural Network Pruning**

Yui Ishihara

# DNN Model Compression

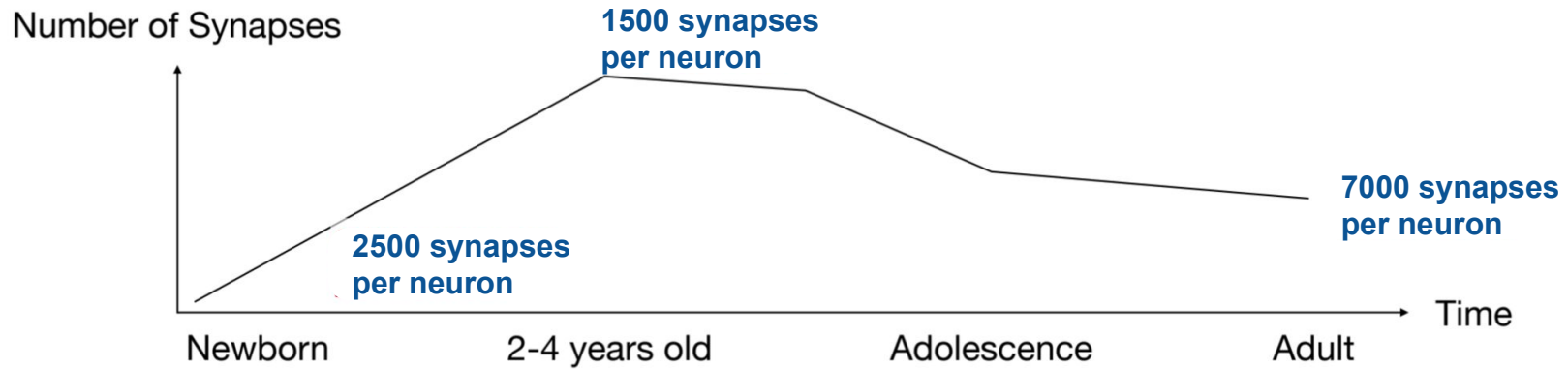- Various Methods for DNN Compression:
  - **Pruning**
  - Quantization
  - Weight Sharing
  - Matrix Factorization
  - Huffman Encoding
  - Low Precision Inference
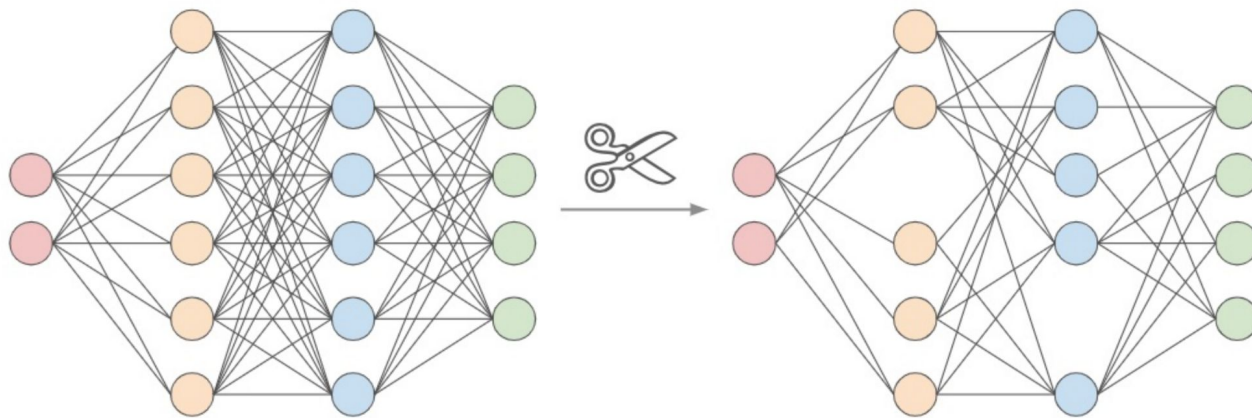  - Binarization
  - … and more!
- Often a combination



Yui Ishihara

# **Pruning:** Motivation

Number of Synapses

**1500 synapses per neuron**

**2500 synapses per neuron**

**7000 synapses per neuron**

Time

Newborn    2-4 years old    Adolescence    Adult

- Trillion synapses generated in human brain (first few months of birth)
- Pruning removes redundant connections in brain
  - 1 year old – peak 1000 Trillion
  - Pruning begins to occur!
  - 10 years old – about 500 Trillion
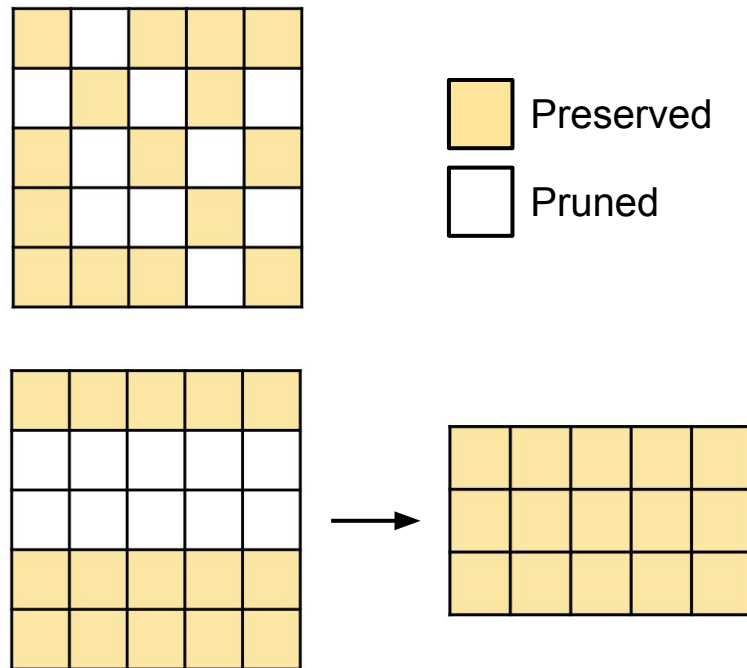
# What is Pruning?

- Pruning simplifies a model by reducing size, removing less critical weights, neurons, or even entire channels, while trying to maintain accuracy

- Goal is efficiency: creating smaller, faster models with lower memory and computation requirements, ideal for deployment in restricted environments (think: EdgeAI!)

Menghani, Gaurav. "Efficient deep learning: A survey on making deep learning models smaller, faster, and better."
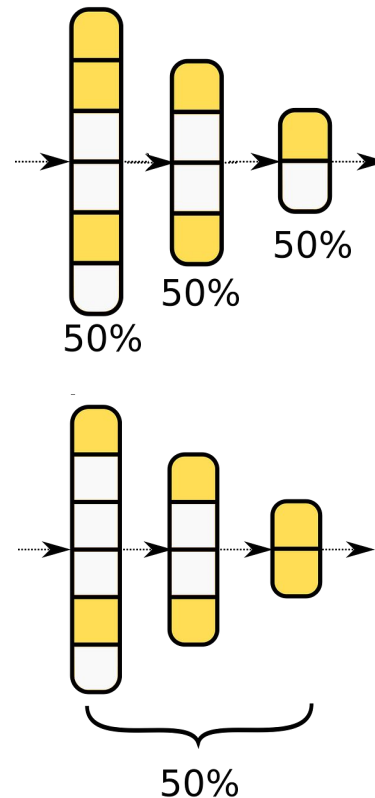
Yui Ishihara

# **What to Prune:** Structured vs. Unstructured

- **Unstructured pruning:** find and remove the less salient connections in the model wherever they are. (Does not consider any relationship between the pruned weights)



- Preserved
- Pruned

- **Structured pruning:** the selected removal of larger part of the network (e.g. layer)
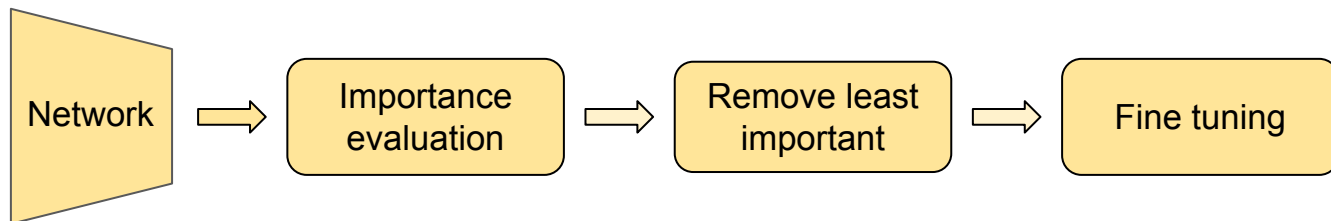
# **What to Prune:** Local vs. Global

- **Local pruning**: consists of removing a fixed percentage of weights from each layer by comparing weights within the layer

- **Global pruning:** pools all parameters together across layers and selects a global sparsity of them to prune
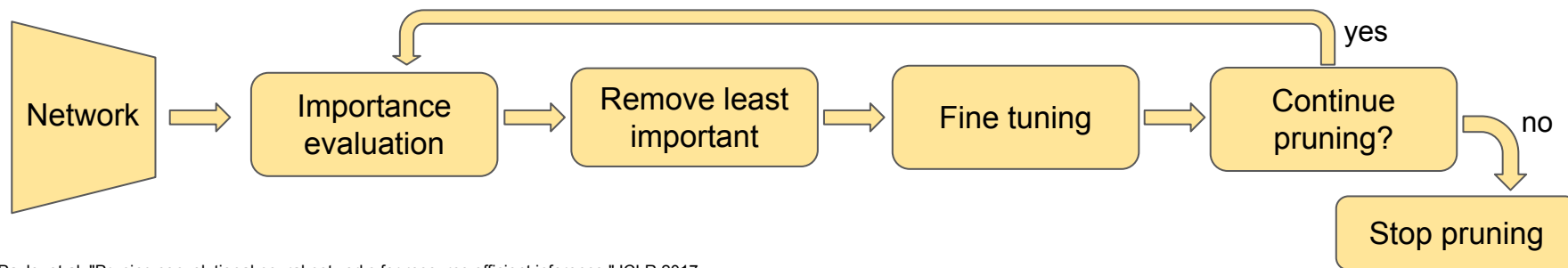
# **Pruning Methods:** Iterative vs. One-Shot

**One-Shot:** network connections pruned only once

Network → Importance evaluation → Remove least important → Fine tuning

**Iterative:** network connections pruned partially through multiple iterations

Network → Importance evaluation → Remove least important → Fine tuning → Continue pruning? — yes (loop back) / no → Stop pruning

Molchanov, Pavlo, et al. "Pruning convolutional neural networks for resource efficient inference." ICLR 2017

Yui Ishihara

# Pruning Criteria

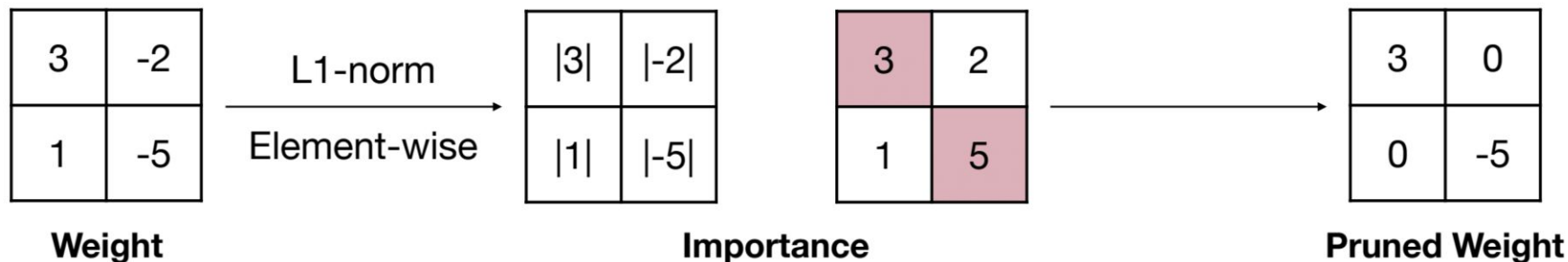Many heuristics and methods to choose weights/neurons to prune:

- **Magnitude-based**

- Gradient-based

- Learned
  - E.g. learn pruning masks

- Information-based
  - E.g. Higher-order curvature

# Magnitude Based Pruning

- Magnitude-based pruning considers weights with **larger absolute values** as more important than other weights.
- For element-wise pruning:

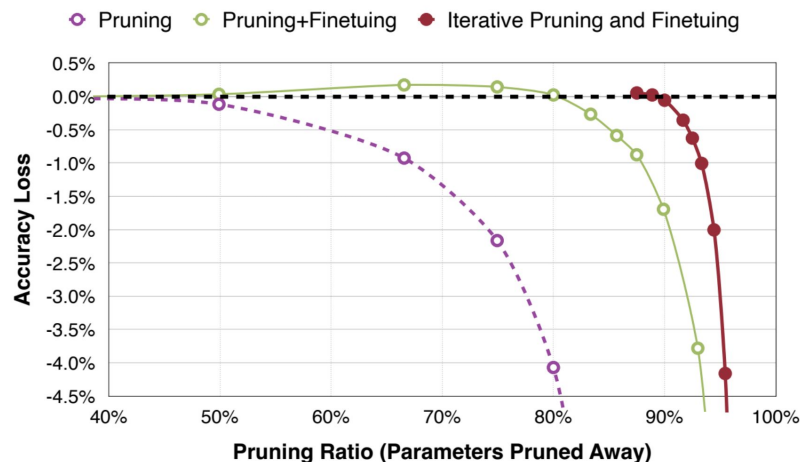$$Importance = |W|$$

- Example

# Fine-Tuning Pruned Networks

- ## Accuracy Impact After Pruning
  - Model accuracy may decrease, especially at higher pruning ratios.
- ## Fine-Tuning Benefits
  - Fine-tuning pruned networks can recover accuracy and enable higher pruning ratios.



Han, Song, et al. "Learning both weights and connections for efficient neural network." NIPS. 2015.

Yui Ishihara