# AI at the Edge

**Wei Gao**

# AI At the Edge



Advancement in AI Computing

juven

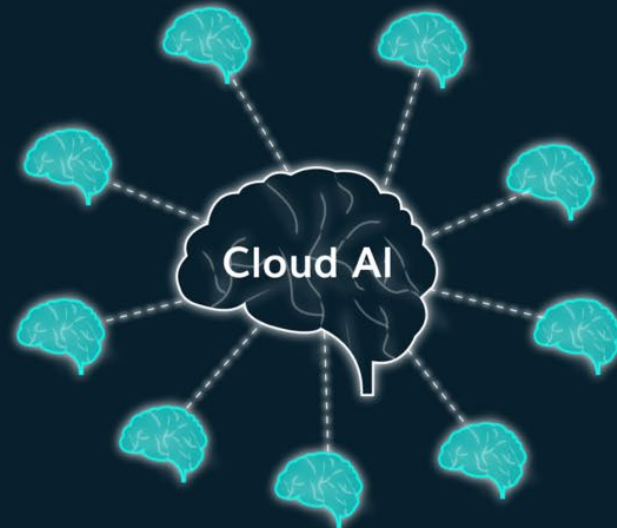**Today**

**Cloud Computing**

Cloud AI

Remote devices connect to an AI in the cloud which does the actual processing.

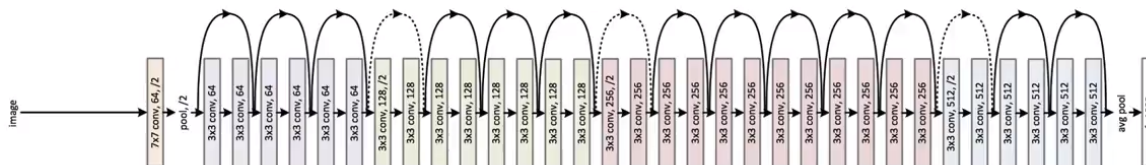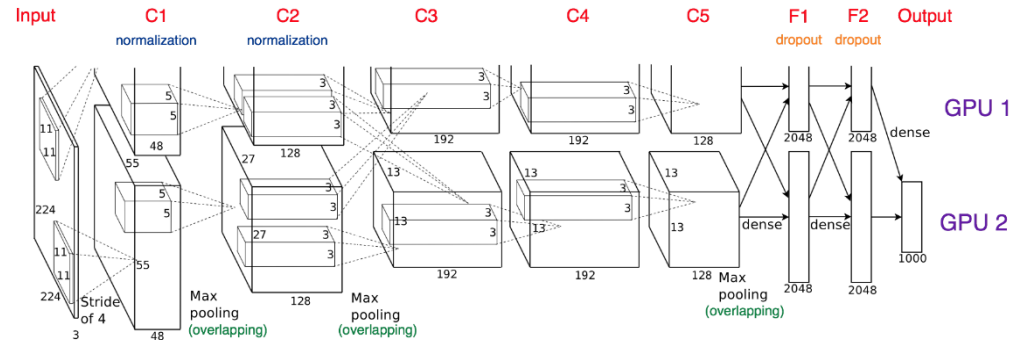**Tomorrow**

**Edge Computing**

Cloud AI

With the increasingly greater power and smaller size of the AI processor, devices are more self-reliant in data processing.

# Key Challenges

- Limited computing resources at the edge devices
  - Limited computing power
  - Limited memory space

- Slow speed of training and inference

# Huge NN Models

- Image Recognition

- AlexNet (2012) ILSVRC winner

  - 8 layers, 62Mparameters

  - 1.4 GFLOP inference

  - 16% error rate



- ResNet (2015) ILSVRC winner

  - 152 layers, 60Mparameters

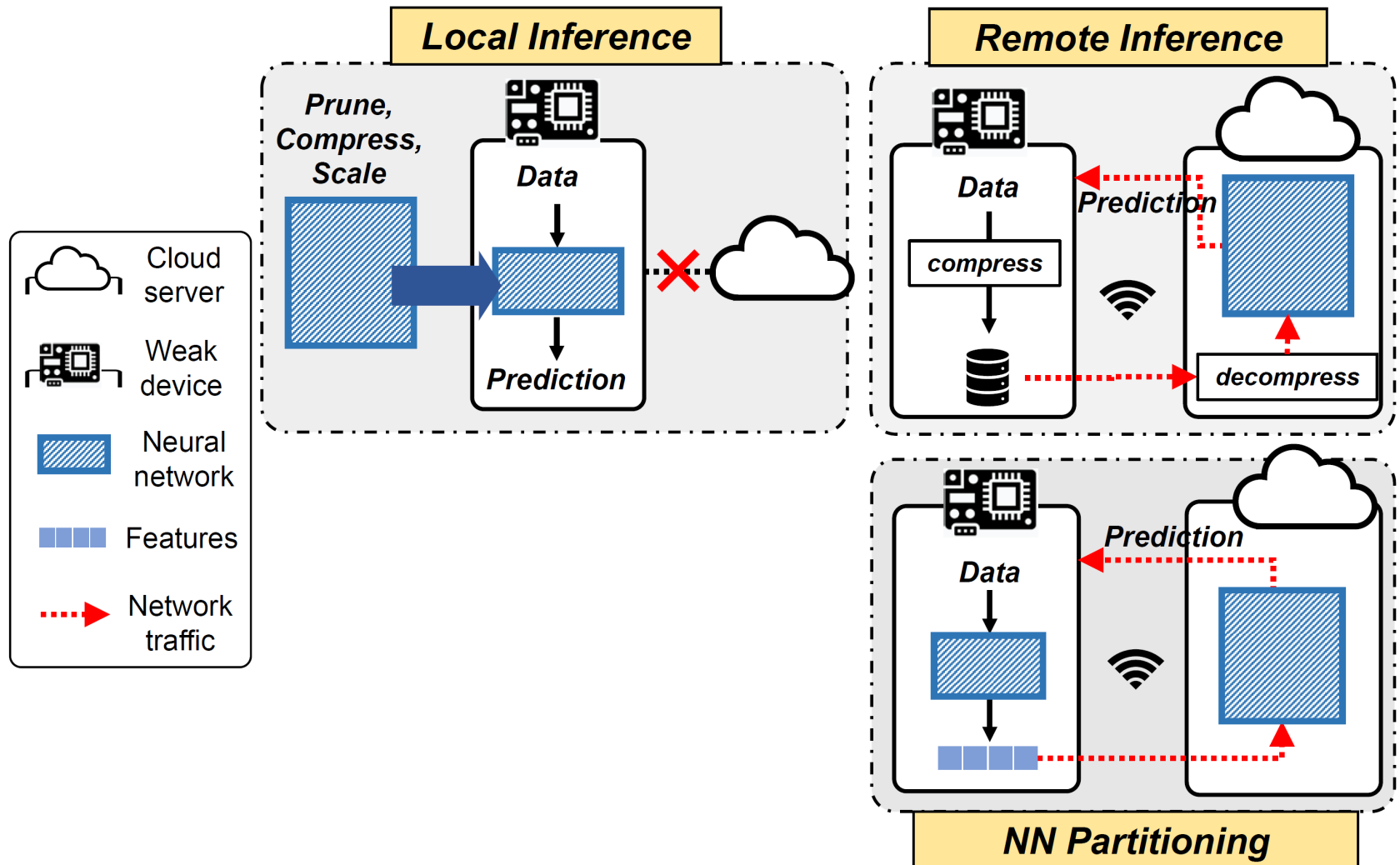  - 22.6 GFLOP inference

  - 6.16% error rate

# Huge NN Models

- Network design and training time have become a huge bottleneck

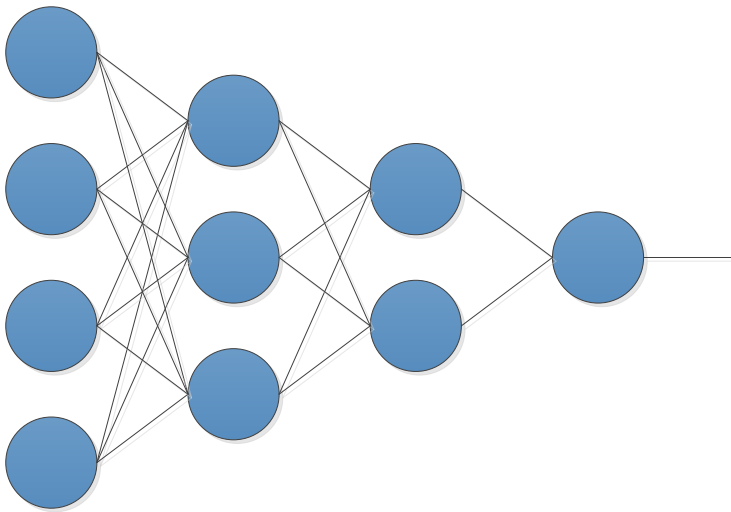|  | Error rate | Training time |
|---|---|---|
| ResNet 18: | 10.76% | 2.5 days |
| ResNet 50: | 7.02% | 5 days |
| ResNet 101: | 6.21% | 1 week |
| ResNet 152: | 6.16% | 1.5 weeks |

# Potential Solutions
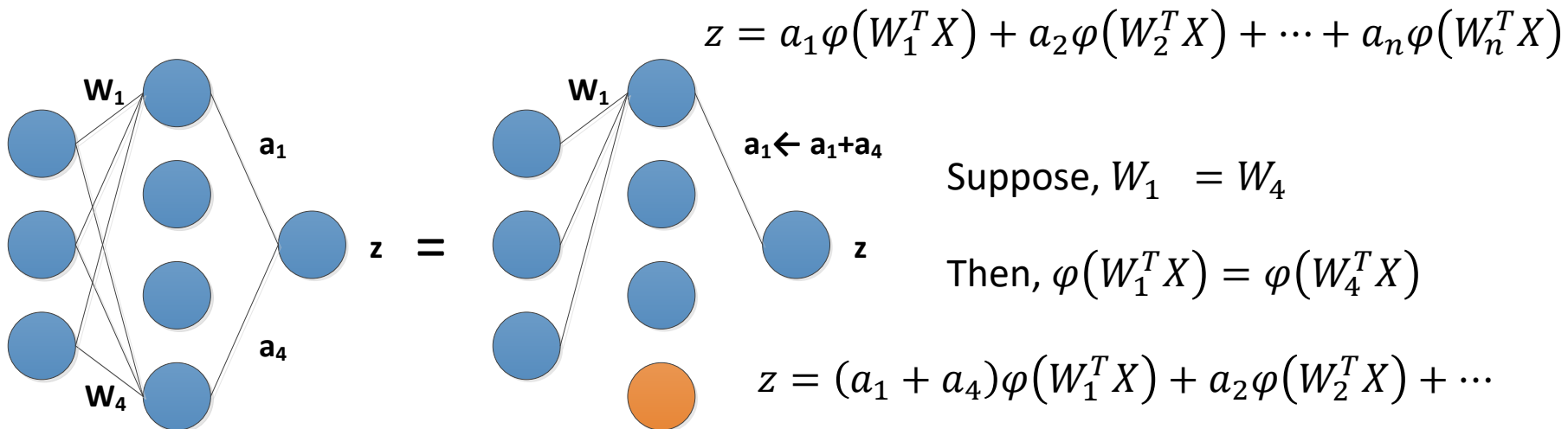
before pruning

after pruning

pruning synapses

pruning neurons

# Scenario 1: You only have a model

- Naïve pruning: Remove weights based on magnitude, weights close to zero are removed
  - No well-founded theory, error increases rapidly
- Data-Free parameter pruning based upon weight similarity

$$z = a_1\varphi(W_1^T X) + a_2\varphi(W_2^T X) + \cdots + a_n\varphi(W_n^T X)$$

$W_1$

$a_1$

$z$ =

$a_4$

$W_4$

$W_1$

$a_1 \leftarrow a_1 + a_4$

$z$

Suppose, $W_1 = W_4$

Then, $\varphi(W_1^T X) = \varphi(W_4^T X)$

$$z = (a_1 + a_4)\varphi(W_1^T X) + a_2\varphi(W_2^T X) + \cdots$$

# Data-Free pruning uses only the model sensitivity

- In practice neurons are different, $\|W_1 - W_2\| = \|\varepsilon_{1,2}\| \geq 0$

  - Compute errors for Weight replacement and naïve removal, so called **saliency** matrix **M**

    – Pick minimum entry in the list e.g. indices $(i', j')$, delete the $j'^{\text{th}}$ neuron and update $a_{i'} \leftarrow a_{i'} + a_{j'}$

    – Update **M** by removing $j'^{\text{th}}$ column and row, and update the $i'^{\text{th}}$ column for updated $a_{i'}$

- When to stop?

  - Saliency in line with test error

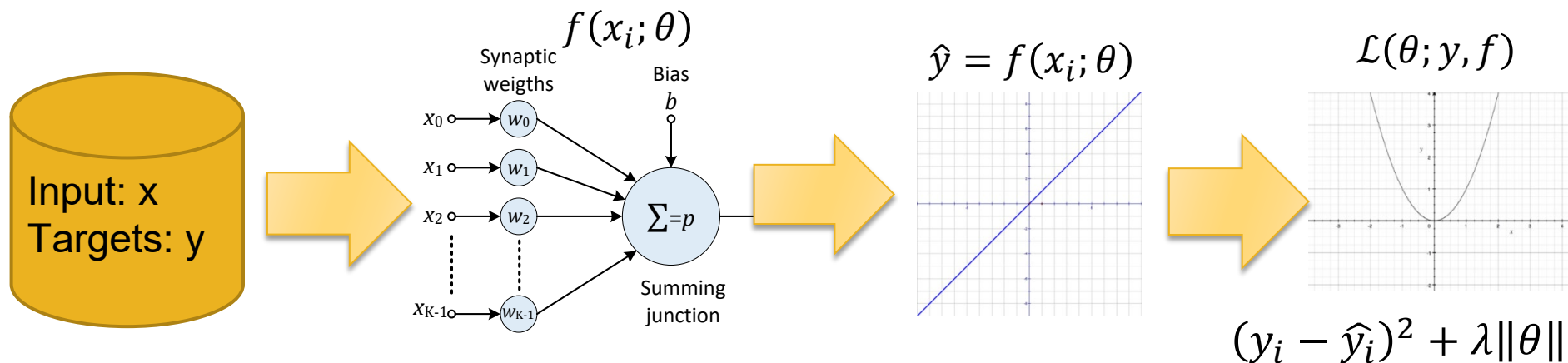  - Find the mode in the gauss like curve



9

# Scenario 2: You have data: how to prune aggressively?

- With access to training data, you can do a lot more
1. Train your network differently such that you have more zero weights
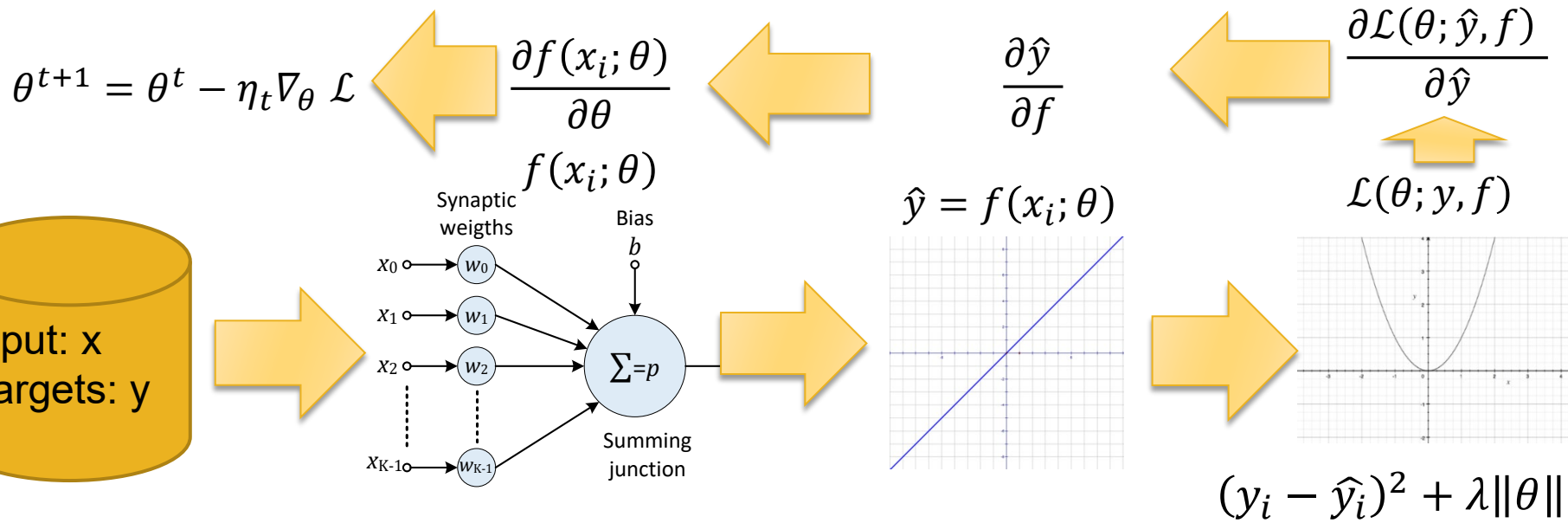2. Retrain you network after pruning to fix the errors

# Training a neural network: With a weight regularization

- The neural network is a function of inputs $x_i$ and weights $\theta$: $f(x_i; \theta)$

- Start with feed forward batch (i=1..64) through the network: $x_i \rightarrow \widehat{y_i}$

- Insert network results and desired target labels into a loss function: $\mathcal{L}(\theta; y, f)$

- Compute a score on how well the net performs, not only error also weight organization

$$f(x_i; \theta)$$

Synaptic weigths

Bias $b$

Input: x
Targets: y

$x_0$ → $w_0$

$x_1$ → $w_1$

$x_2$ → $w_2$

$x_{K-1}$ → $w_{K-1}$

$\Sigma = p$

Summing junction

$$\hat{y} = f(x_i; \theta)$$

$$\mathcal{L}(\theta; y, f)$$

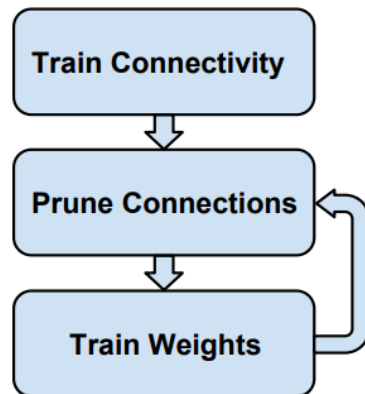$$(y_i - \widehat{y_i})^2 + \lambda \|\theta\|$$

# Tune the weights by gradient descent

- Compute the error gradients
- Update the coefficients to reduce error, also taking into account regularization
- Repeat

$$\theta^{t+1} = \theta^t - \eta_t \nabla_\theta \mathcal{L}$$

$$\frac{\partial f(x_i; \theta)}{\partial \theta}$$

$$\frac{\partial \hat{y}}{\partial f}$$

$$\frac{\partial \mathcal{L}(\theta; \hat{y}, f)}{\partial \hat{y}}$$

$$f(x_i; \theta)$$

Synaptic weigths

Bias
$b$

$\hat{y} = f(x_i; \theta)$

$\mathcal{L}(\theta; y, f)$

Input: x
Targets: y

$x_0$   $w_0$

$x_1$   $w_1$

$x_2$   $w_2$

$\Sigma = p$

$x_{K-1}$   $w_{K-1}$

Summing junction

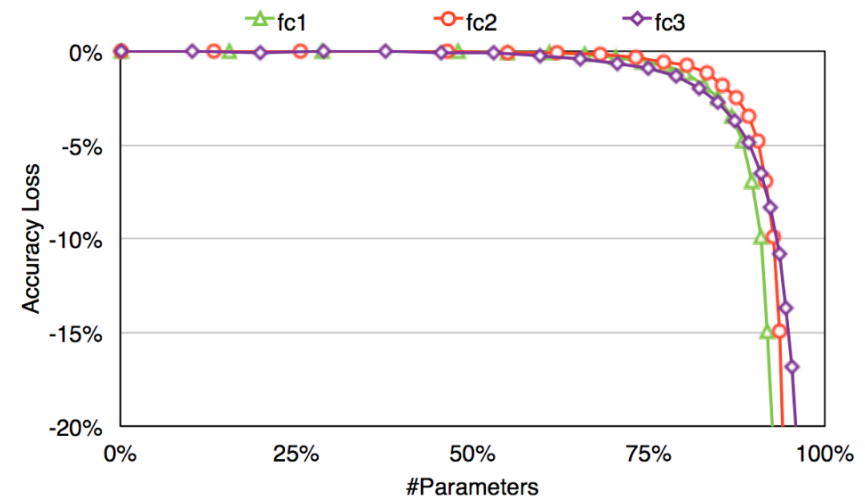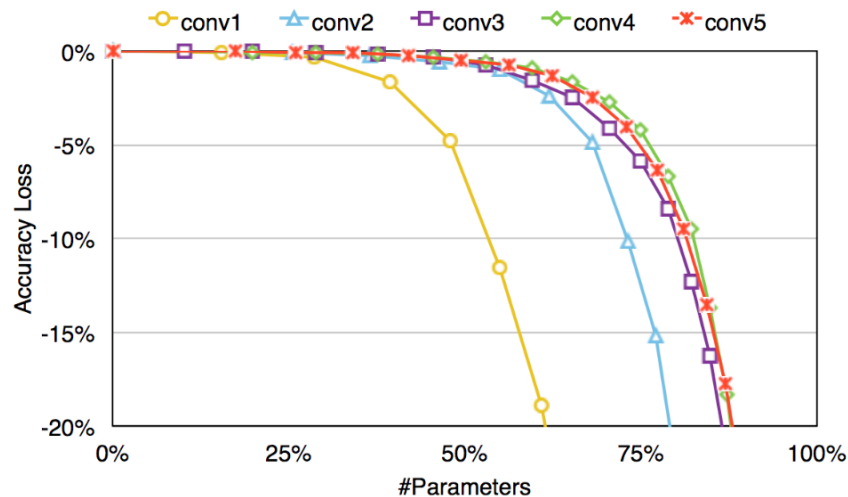$$(y_i - \hat{y}_i)^2 + \lambda \|\theta\|$$

12

# Iterative Pruning and Retraining

- Train a neural network until reasonable solution or download a pretrained net
  1. Prune the weights base on magnitudes that are less than a threshold
  2. Train the network until a reasonable solution is obtained
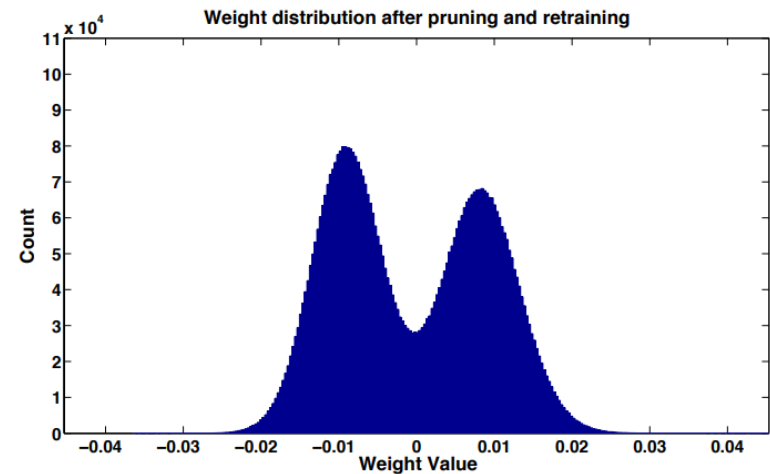  3. Iterate to step 1
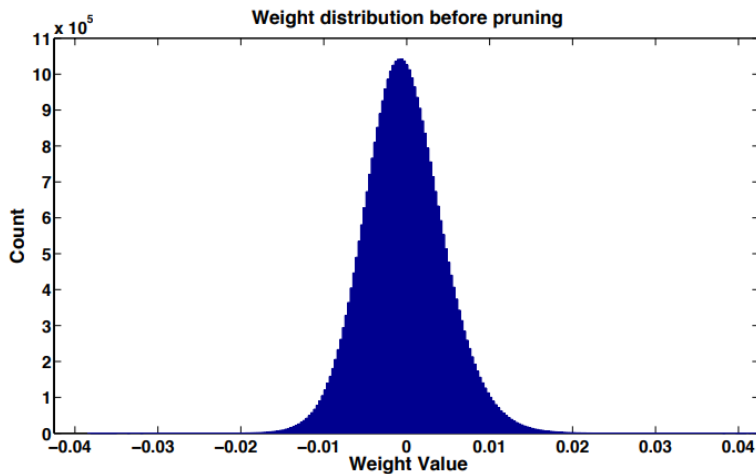
# Where does pruning help the most?

- Fully connected layers
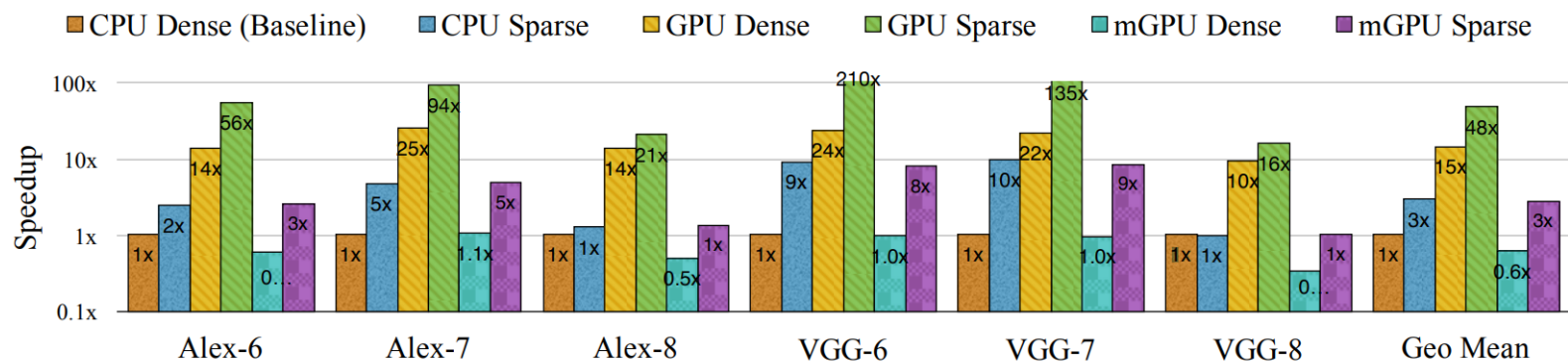


Pruning and Network Sparsity Improvements

# What happens to the weight distribution?

- Before: Most weights are close to zero; almost all between [-0.015, 0.015]
- After pruning: Bimodal distribution and more spread across x-axis, between [-0.025, 0.025]



Weight distribution before pruning
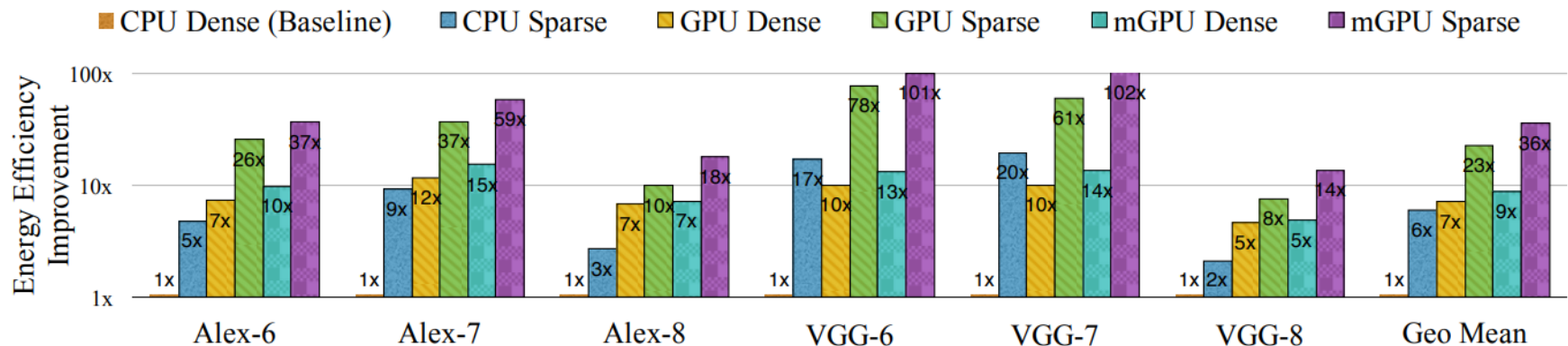
Weight distribution after pruning and retraining

# Using sparse Matrix Computations

- Use Intel Core i7 5930K, MKL CBLAS GEMV (full) vs MKS SPBLAS CSRMV (sparse)

- Use NVIDIA GTX Titan X, cuBLAS GEMV (full) vs cuSPARSE CSRMV (sparse)

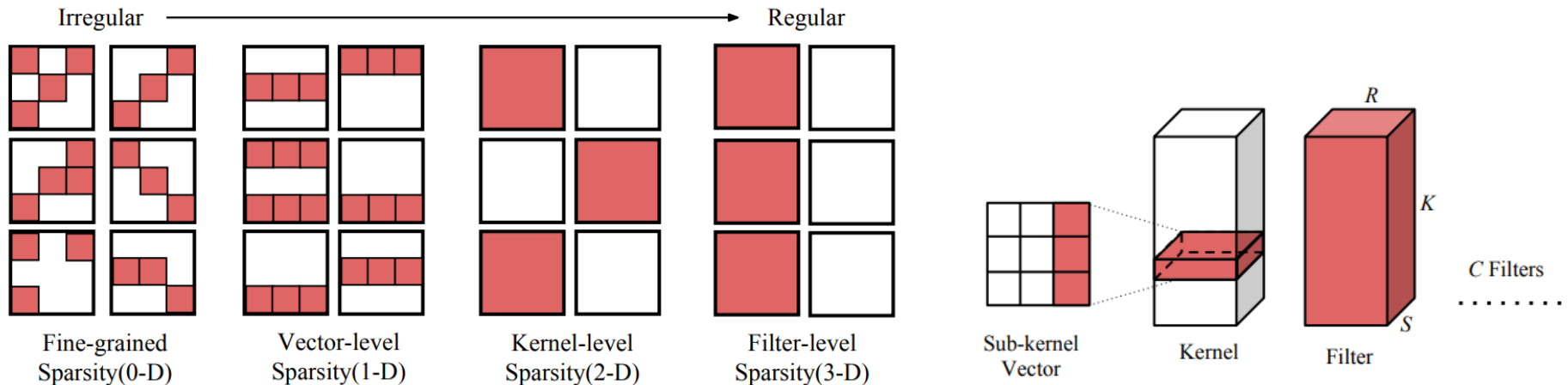- Use NVIDIA Tegra K1 as embedded GPU

# Energy Efficiency

- 6x improvement CPU    3.2x improvement GPU
       4x embedded GPU

- Difficult to exploit the large parameter reduction due to irregularity

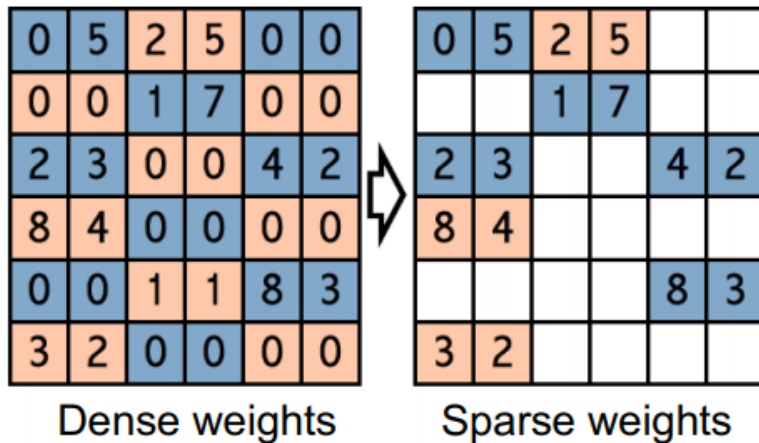  - Sparse matrices have also storage overhead; 16% for storing indices

# From Fine to Coarse-Grained Pruning

- Prune to match the underlying data-parallel hardware
  - E.g. prune by eliminating entire filter planes

# Structured Pruning
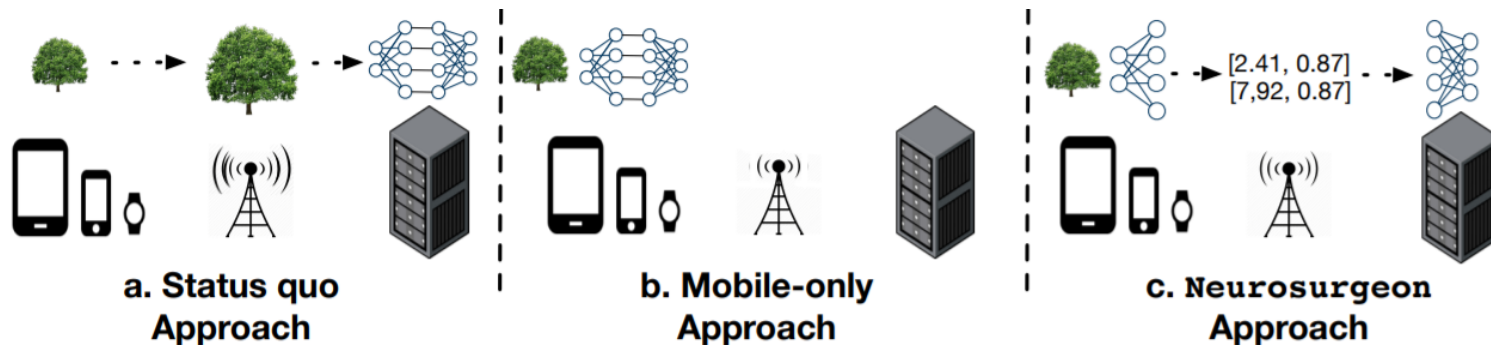
- Example 2-way SIMD
- Less storage overhead



Dense weights → Sparse weights

Fine-grained sparsity

| Weight | Index | Weight | Index | Weight | Index |

Coarse-grained sparsity

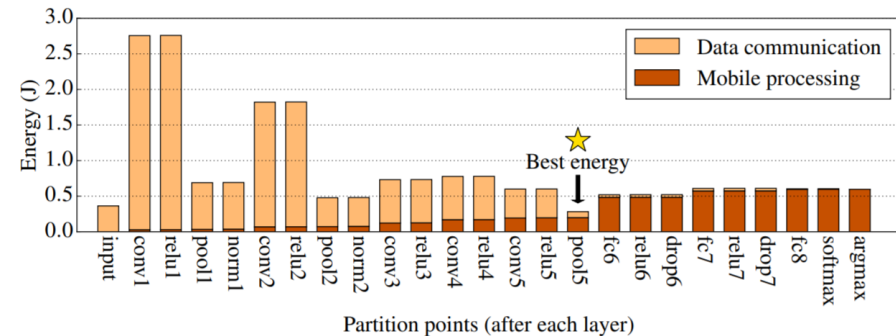| Weight | Weight | Weight | Index | Savings! |

# 2. NN Partitioning

- Why not offloading the work to the edge?

- Representative work: Neurosurgeon



a. Status quo Approach

b. Mobile-only Approach

c. Neurosurgeon Approach

- Partition the neural network in layers

# Key Question

- Where to partition?

# Practical Use



1) Generate prediction models

CONV   FC
POOL   ACT

CONV   FC
POOL   ACT

Prediction Model | Prediction Model | Prediction Model

**Deployment Phase**

1) Extract layer configurations

Target Application

2) Predict layer performance

Prediction Model

Prediction Model

3) Evaluate partition points

4) Partitioned Execution

**Runtime Phase**

# More Fine-Grained Partitioning?

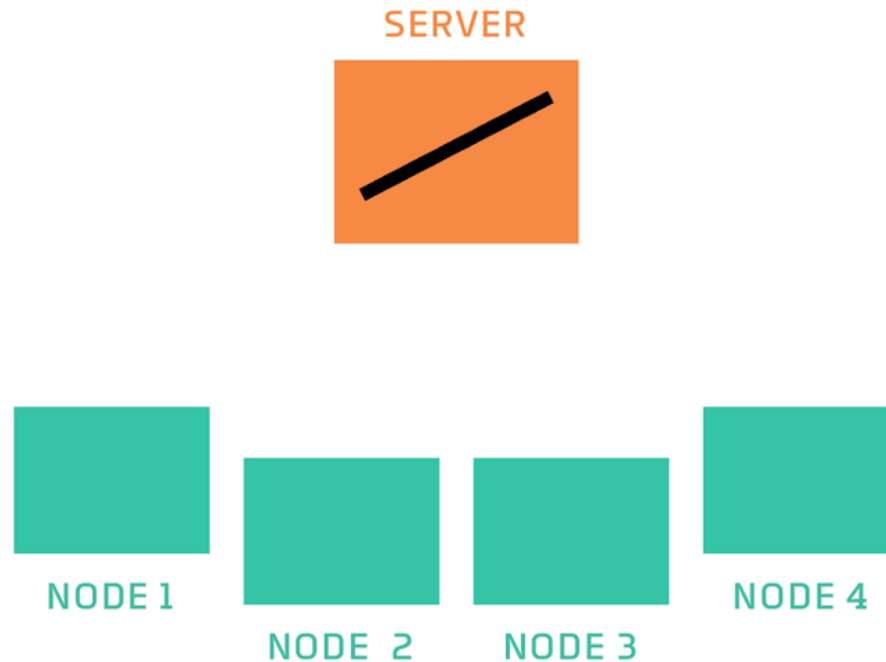- Vertical vs. Horizontal Partitioning
  - Partitioning the feature space

- Adaptive partitioning

# Federated and Distributed Learning



SERVER

NODE 1　　NODE 2　　NODE 3　　NODE 4

$$\min \left[ F(x) = \sum_{i=1}^{m} p_i F_i(x) \right]$$

| m | $n_i$ | n | $p_i$ | $F_i(x)$ | F(x) |
|---|---|---|---|---|---|
| Number of clients | Number of samples at client i | Total number of samples | $n_i / n$, relative sample size | Local objective function at client i | Global objective function |

# Challenges

**Expensive Communication:**

Communication in the network can be slower than local computation by many orders of magnitude.

Soultion: Smaller messages or sending less frequently

**Privacy Concerns:**

Sensitive information can still be revealed to third party or central server during the communication.

**Systems Heterogeneity**:

- Size of data
- Computational power
- Network stability
- Local solvers
- Learning rate

**Statistical Heterogeneity:**



$$\omega_g^{(t)} = \sum_{k=1}^{K} \frac{\alpha_k \cdot \omega_k^{(t-1)}}{\sum_{k=1}^{K} \alpha_k}$$

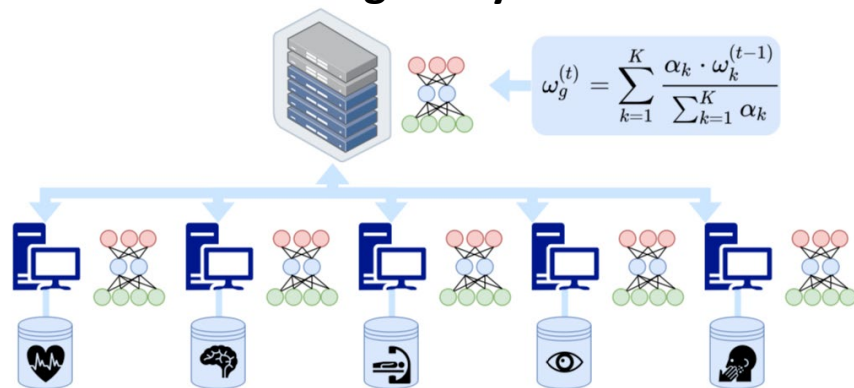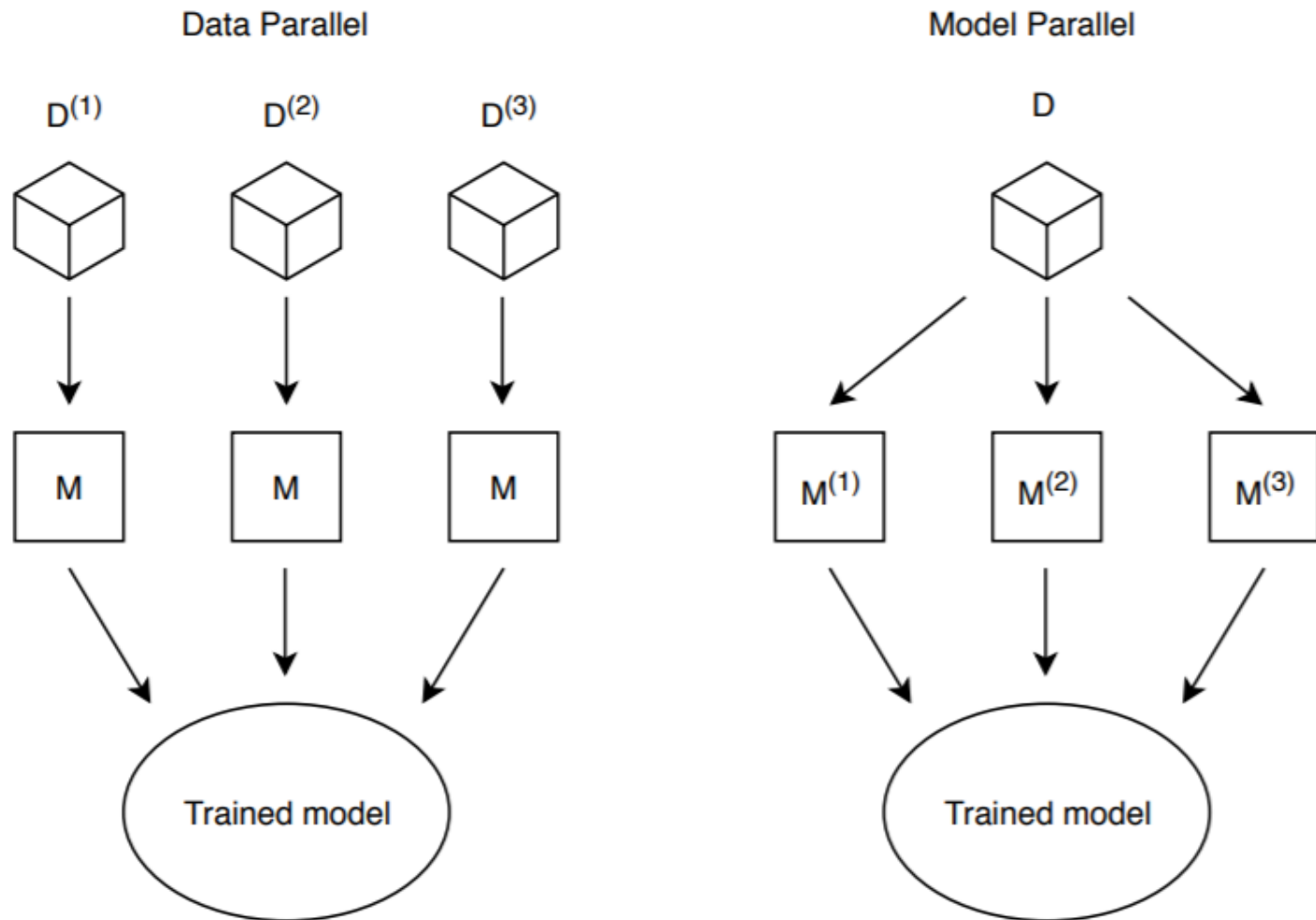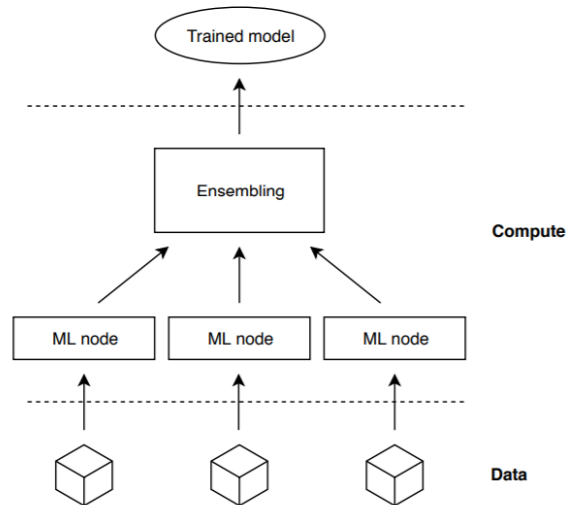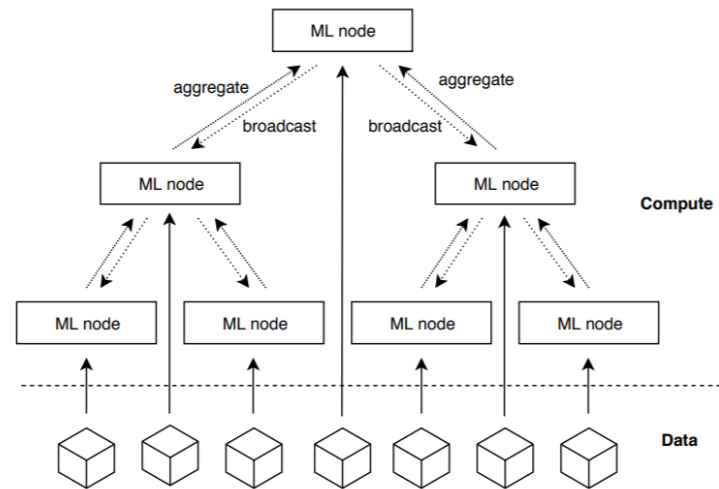Fig. 1: Federated learning with non-iid data - The data has different distributions among clients.

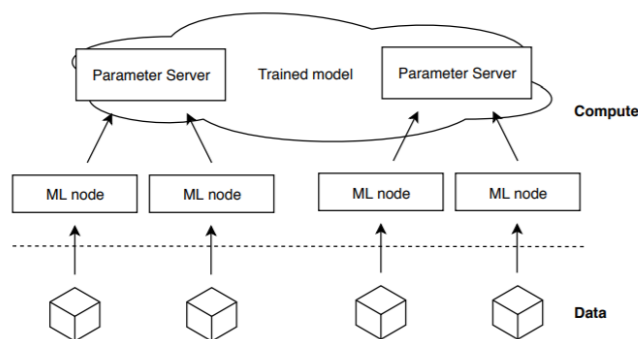# How to Achieve Parallelism

# Different Operational Modes



(a) Centralized (Ensembling)

(b) Decentralized (Tree)

(c) Decentralized (Parameter Server)

(d) Fully Distributed (Peer to Peer)

# The Ecosystem



**Distributed Machine Learning**

**General Purpose Distributed Computing Frameworks**
- Apache Hadoop
- Apache Spark
- Apache Flink
- etc.

Mahout →
MLlib →
Hadoop/Spark + AllReduce →

**Natively Distributed ML Systems**
- Caffe2
- CNTK
- DistBelief
- DIANNE
- Tensorflow
- MxNet
- etc.

← Keras
NVIDIA NCCL ←

**Single-Machine ML Systems and Libraries**
- Theano
- Caffe
- Scikit
- MLPack
- NVIDIA Libraries
- etc.

**Cloud Machine Learning**
- Google Cloud AI
- Microsoft Azure ML
- Amazon AWS ML
- IBM Watson Cloud
- etc.

# SGD Parallelization



Aggregation can be performed via:
- **Master node ("parameter server")**
- **MPI All-Reduce ("decentralized")**
  - **Shared-Memory**

$$\tilde{g}_n = \sum_i g^i / n$$

Stochastic gradient with **n times lower variance**

$\tilde{g}^1$    $\tilde{g}^2$    $\tilde{g}^3$     **...**     $\tilde{g}^n$

Dataset Partition **1**    Dataset Partition **2**    Dataset Partition **3**    Dataset Partition **n**

# Some more recent works

- TinyML / MCUNet
  - [https://mcunet.mit.edu/](https://mcunet.mit.edu/)


- Split learning
  - [http://splitlearning.mit.edu/](http://splitlearning.mit.edu/)