The University of Texas at Austin
Electrical and Computer Engineering
Cockrell School of Engineering

**Fall 2021**

# ADVANCED TOPICS IN COMPUTER VISION

**Atlas Wang**

Assistant Professor, The University of Texas at Austin

**Visual Informatics Group@UT Austin**
https://vita-group.github.io/

# Deep Learning on the Edge

- Deploying CNNs on resource-constrained platforms/at the edge
- Two Scenarios: **Inference** (pre-trained model), and **Training** (online adaptation)



## Real-Time Machine Learning (RTML)

**PROGRAM SOLICITATION**
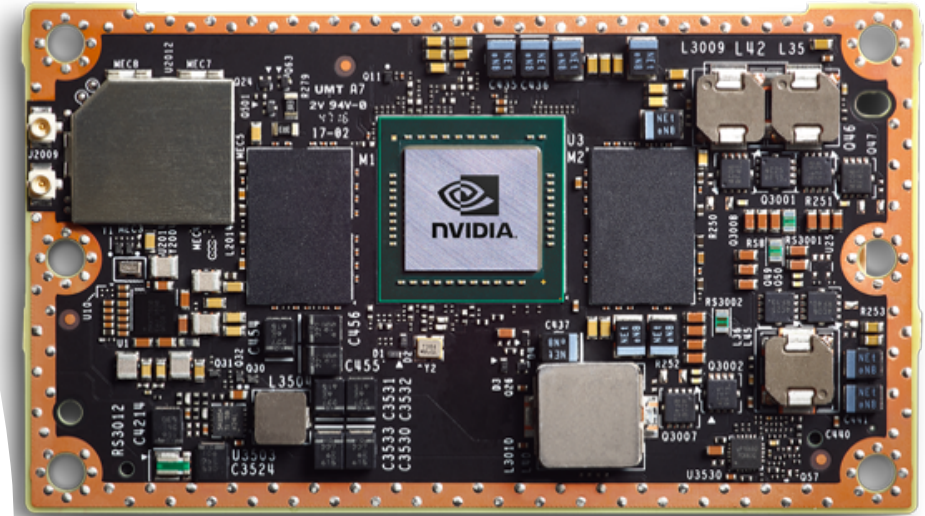NSF 19-566

**National Science Foundation**

Directorate for Computer and Information Science and Engineering
Division of Computing and Communication Foundations

Directorate for Engineering
Division of Electrical, Communications and Cyber Systems

**RTML Program goal:** *"for next-generation **co-design** of RTML algorithms and hardware, with the principal focus on developing novel hardware architectures and learning algorithms in which **all stages of training** (including incremental training, hyperparameter estimation, and deployment) can be performed in real time."*

# Deep Learning on the Edge



- Three Top Concerns:
  - **Storage and Memory**
  - **Speed or Latency**
  - **Energy Efficiency**

- The three goals all pursue "light weight"
- … but they are often **not aligned**\*
- … so need to **consider all** in implementation
- … and for both **Inference** and **Training**

- Broad economic viability requires energy efficient AI
- Energy efficiency of a brain is **100x better** than current SOTA hardware!



**Common carbon footprint benchmarks**

in lbs of CO2 equivalent

| | |
|---|---|
| Roundtrip flight b/w NY and SF (1 passenger) | 1,984 |
| Human life (avg. 1 year) | 11,023 |
| American life (avg. 1 year) | 36,156 |
| US car including fuel (avg. 1 lifetime) | 126,000 |
| Transformer (213M parameters) w/ neural architecture search | 626,155 |

Chart: MIT Technology Review • Source: Strubell et al. • Created with Datawrapper

# Model Compression

- Training Phase:
  - The easiest way to extract a lot of knowledge from the training data is to learn many different models in parallel.
  - 3B: Big Data, Big Model, Big Ensemble
  - Imagenet: 1.2 million pictures in 1,000 categories.
  - AlexNet: ~ 240Mb, VGG16: ~550Mb
- Testing Phase:
  - Want small and specialist models.
  - Minimize the amount of computation and the memory footprint.
  - Real time prediction
  - Even able to run on mobile devices.

# Two Main Streams

- **"Transfer":** How to transfer knowledge from big general model (teacher) to small specialist models (student)?
  - Example: "Distilling the Knowledge in a Neural Network", G. Hinton et. al., 2015

- **"Compress":** How to reduce the size of the same model, during or after training, without losing much accuracy.
  - Example: "Deep Compression: Compressing Deep Neural Networks with Pruning, Trained Quantization and Huffman Coding", S. Han et. al., 2016

- **Comparison:** Knowledge Transfer provides a way to train a <u>new small model</u> inheriting from big general models, while Deep Compression Directly does the surgery on big models, using a pipeline: pruning, quantization & Huffman coding.

# Knowledge Transfer/"Distillation": Main Idea

- Introduce "Soft targets" as one way to transfer the knowledge from big models.
  - Classifiers built from a softmax function have a great deal more information contained in them than just a classifier;
  - The correlations in the softmax outputs are very informative.

- Hard Target: the ground truth label (one-hot vector)
- Soft Target:  $q_i = \dfrac{exp(z_i/T)}{\sum_j exp(z_j/T)}$     T is "temperature", z is logit

- More information in soft targets

| cow | dog | cat | car | |
|-----|-----|-----|-----|--|
| 0 | 1 | 0 | 0 | original hard targets |

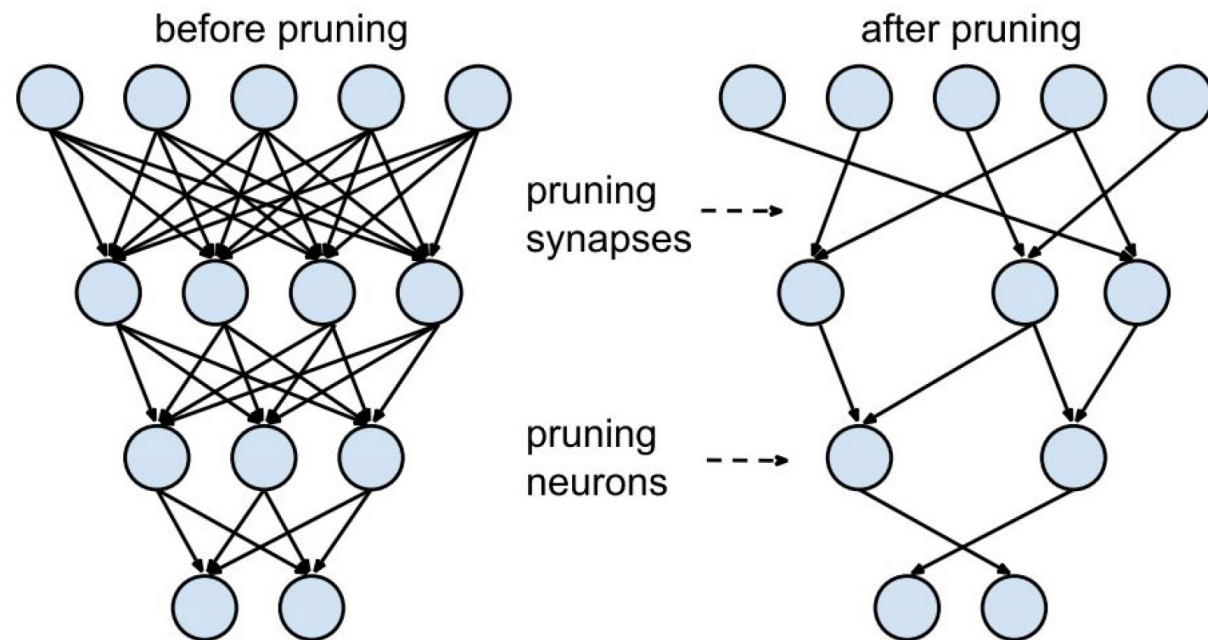| cow | dog | cat | car | |
|-----|-----|-----|-----|--|
| .05 | .3 | .2 | .005 | softened output of ensemble |

**Hinton's Observation:** If we can extract the knowledge from the data using very big models or ensembles of models, it is quite easy to distill most of it into a much smaller model for deployment.

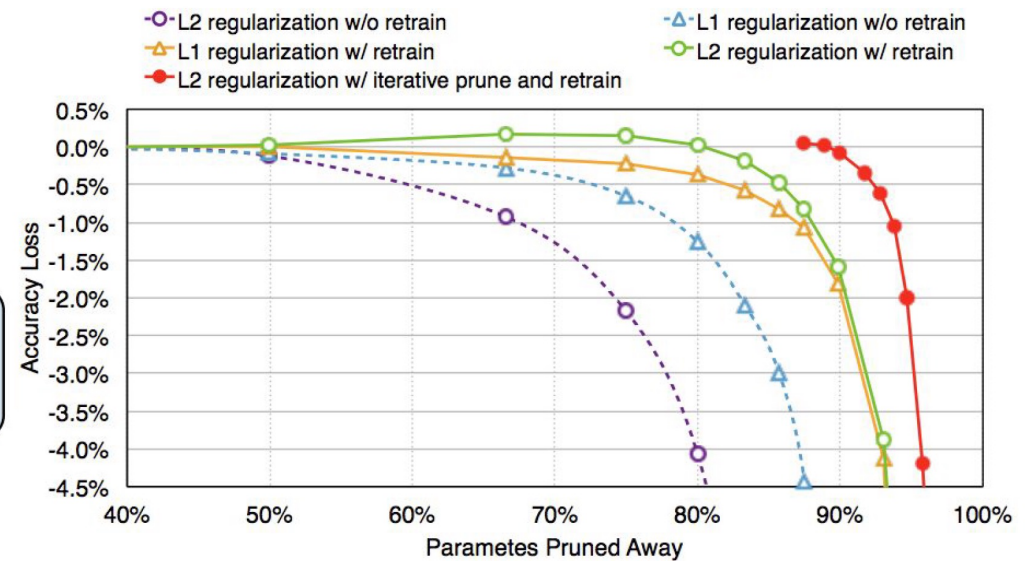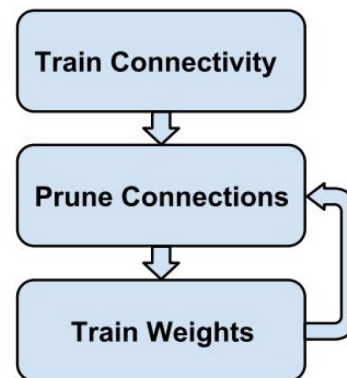**More follow-up observations:** teachers can be weak, or even the same as student …

Deep Compression: Main Idea (i)

# Pruning



before pruning  →  pruning synapses  →  after pruning

pruning neurons

# Deep Compression: Main Idea (ii)

## Retrain to Recover Accuracy



Network pruning can save 9x to 13x parameters without drop in accuracy

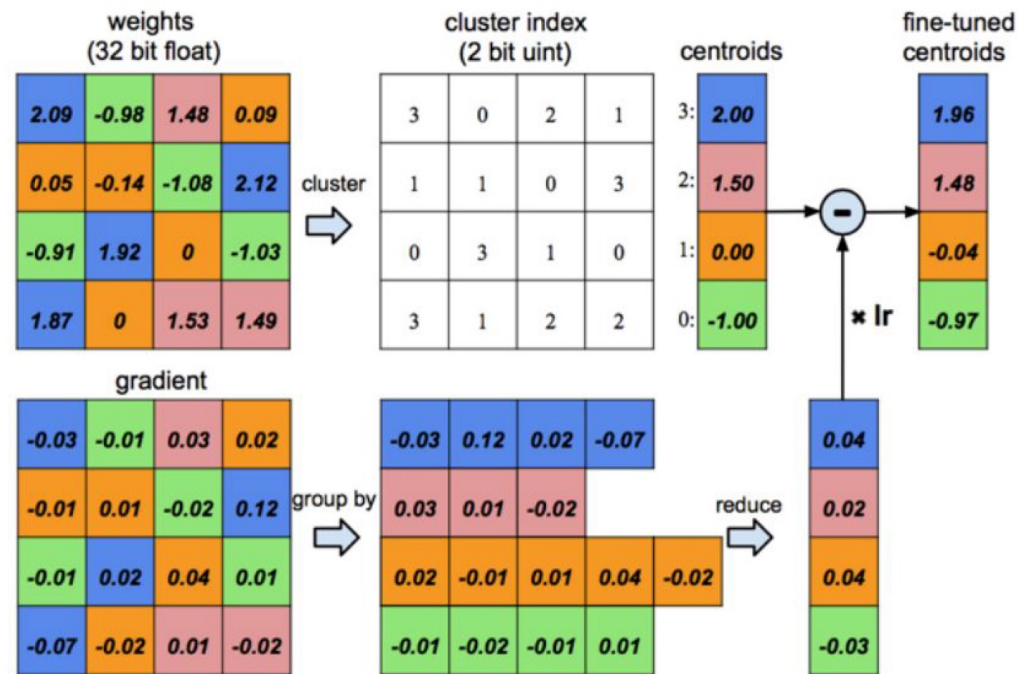# Deep Compression: Main Idea (iii)
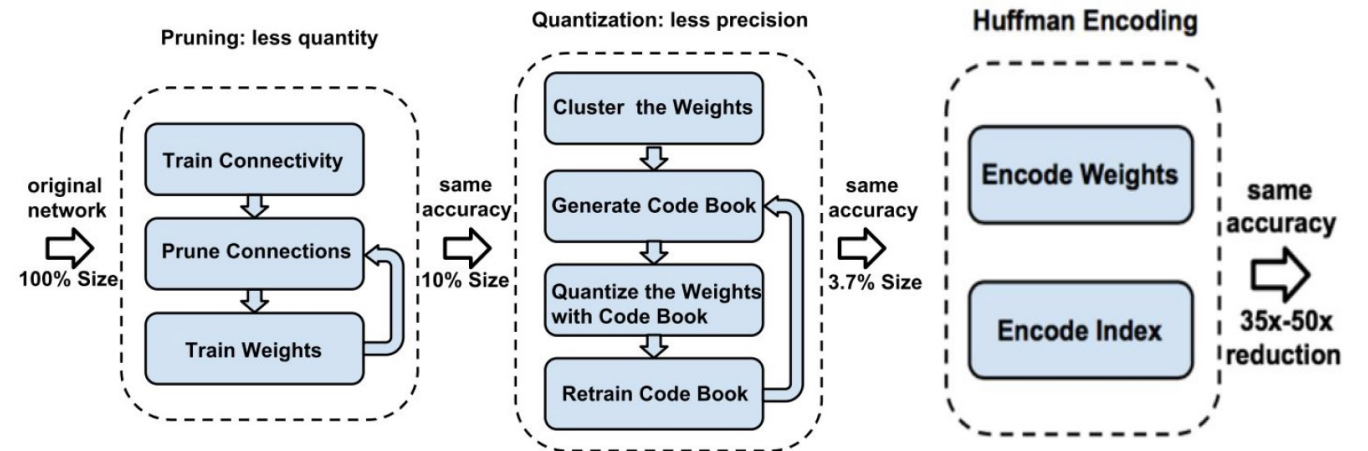
## Weight Sharing (Trained Quantization)



Figure 3: Weight sharing by scalar quantization (top) and centroids fine-tuning (bottom)

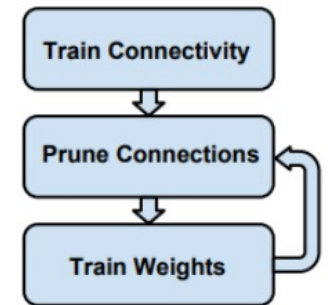# Deep Compression: Main Idea (iv)
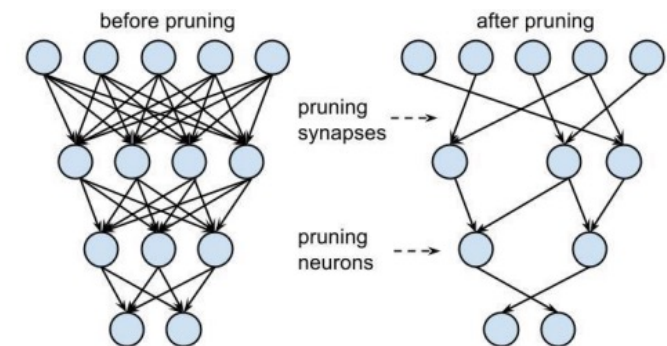
## Huffman Coding

## More About Pruning

- Determining **low-saliency parameters**, given a pre-trained network
- Follows the framework proposed by LeCun et al. (1990):

1. **Train** a deep model until convergence
2. **Delete** "unimportant" connections w.r.t. a certain criteria
3. **Re-train** the network
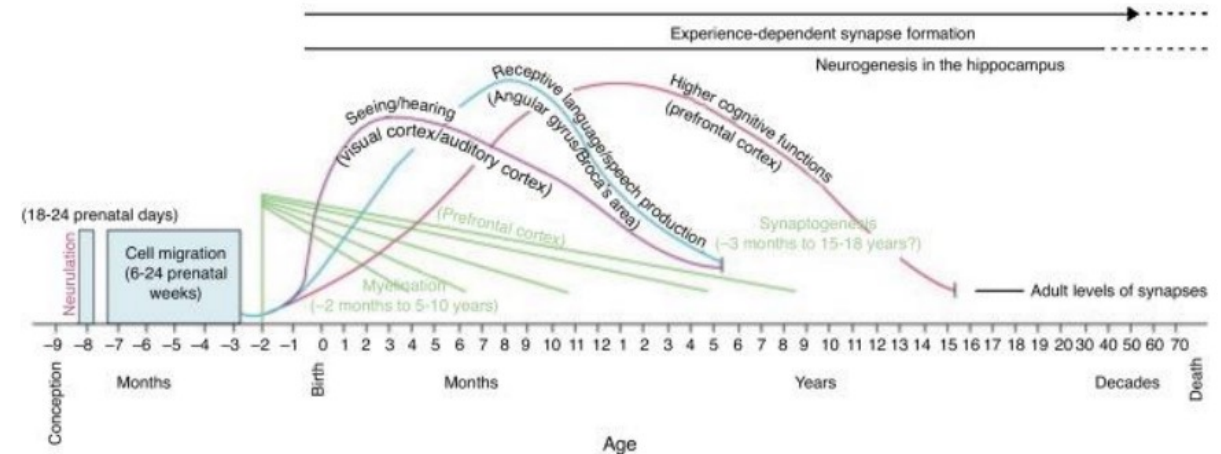4. **Iterate** to step 2, or stop



- Defining **which connection is unimportant** can vary
  - Weight magnitudes ($L^2$, $L^1$, ...)
  - Mean activation [Molchanov et al., 2016]
  - Avg. % of Zeros (APoZ) [Hu et al., 2016]
  - Low entropy activation [Luo et al., 2017]
  - ...

# Human Brain Prunes too!

- Human brains are also using pruning schemes as well
- **Synaptic pruning** removes redundant synapses in the brain during lifetime

## Optimal Brain Damage (OBD)

- Network pruning **perturbs weights W** by **zeroing** some of them
- How the loss $L$ would be changed when **W** is perturbed?

- **OBD** approximates $L$ by the 2nd order Taylor series:

$$\delta L \simeq \underbrace{\sum_i \frac{\partial L}{\partial w_i} \delta w_i}_{\text{1st order}} + \underbrace{\frac{1}{2} \sum_i \frac{\partial^2 L}{\partial w_i{}^2} \delta w_i^2 + \frac{1}{2} \sum_{i,j} \frac{\partial^2 L}{\partial w_i \partial w_j} \delta w_i \delta w_j}_{\text{2nd order}} + O(\|\delta \mathbf{W}\|^3)$$

- **Problem:** Computing $H = \left( \frac{\partial L}{\partial w_i \partial w_j} \right)_{i,j}$ is usually intractable
  - Requires $O(n^2)$ on **# weights**
  - Neural networks usually have enormous number of weights
    - e.g. AlexNet: **60M** parameters $\Rightarrow H$ consists $\approx \mathbf{3.6 \times 10^{15}}$ elements

## Optimal Brain Damage (OBD)

- **Problem:** Computing $H = \left( \frac{\partial L}{\partial w_i \partial w_j} \right)_{i,j}$ is usually intractable

- Two additional assumptions for tractability

  1. **Diagonal** approximation: $H = \frac{\partial^2 L}{\partial w_i \partial w_j} = 0$ if $i \neq j$

  2. **Extremal** assumption: $\frac{\partial L}{\partial w_i} = 0$ $\forall i$
     - **W** would be in a local minima if it's pre-trained

- Now we get: $\delta L \simeq \frac{1}{2} \sum_i \frac{\partial^2 L}{\partial w_i^2} \delta w_i^2 + O(\|\delta \mathbf{W}\|^3)$
  - It only needs $\mathrm{diag}(H) := \left( \frac{\partial^2 L}{\partial w_i^2} \right)_i$

- **diag($H$)** can be computed in $O(n)$, allowing a backprop-like algorithm
  - For details, see [LeCun et al., 1987]
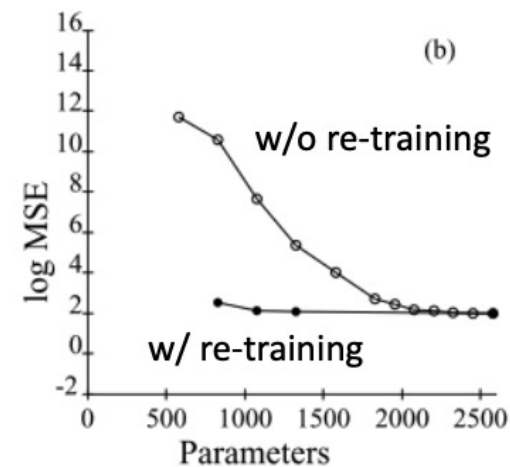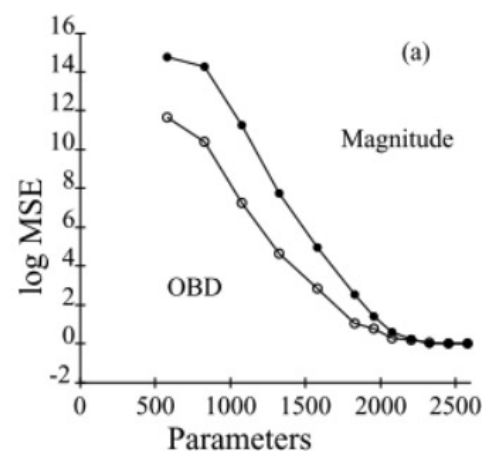
## Optimal Brain Damage (OBD)

- How the loss $L$ would be changed when $\mathbf{W}$ is perturbed?

$$L(\delta \mathbf{W}) \simeq \frac{1}{2} \sum_i \frac{\partial^2 L}{\partial w_i^2} \delta w_i^2 =: \sum_i \frac{1}{2} h_{ii} \delta w_i^2$$

- The **saliency** for each weight $\Rightarrow$ $s_i := \frac{1}{2} h_{ii} |w_i|^2$
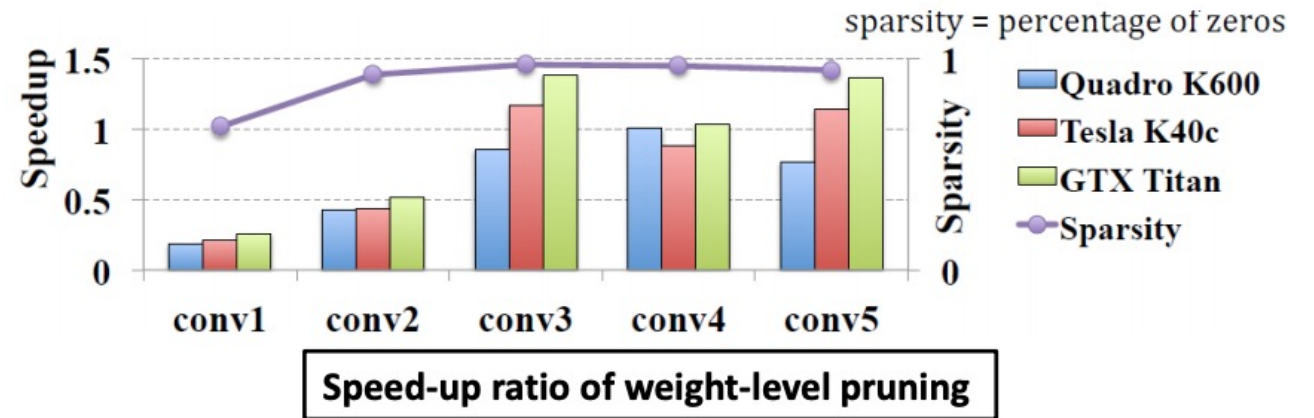
$$s_i := |w_i|$$

- OBD shows robustness on pruning compared to magnitude-based deletion

- After re-training, the original test accuracy is recovered

# Structured Sparsity

- "Un-structured" **weight-level pruning** may not engage a **practical speed-up**
  - Despite of extremely high sparsity, actual speed-ups in GPU is limited



sparsity = percentage of zeros

Speed-up ratio of weight-level pruning

Non-structured sparsity (poor data pattern)

Structured sparsity (regular data pattern)

5× speedup after concatenation of nonzero rows and columns

# Structured sparsity

- **Structured sparsity** can be induced by adding group-lasso regularization

$$\min_{\mathbf{W}} \mathcal{L}(\mathbf{W}) + \lambda \sum_{l=1}^{L} R_g(\mathbf{W}^{(l)}), \ R_g(\mathbf{w}) = \sum_{g=1}^{G} \|\mathbf{w}^{(g)}\|_2$$

- **Filter-wise** and **channel-wise:**   # filters   # channels

$$R_g(\mathbf{W}^{(l)}) = \sum_{n_l=1}^{N_l} \|\mathbf{W}_{n_l,:,:,:}^{(l)}\|_2 + \sum_{c_l=1}^{C_l} \|\mathbf{W}_{:,c_l,:,:}^{(l)}\|_2$$

Table 1: Results after penalizing unimportant filters and channels in *LeNet*

| LeNet # | Error | Filter # [§] | Channel # [§] | FLOP [§] | Speedup [§] |
|---|---|---|---|---|---|
| 1 (*baseline*) | 0.9% | 20—50 | 1—20 | 100%—100% | 1.00×—1.00× |
| 2 | 0.8% | 5—19 | 1—4 | 25%—7.6% | 1.64×—5.23× |
| 3 | 1.0% | 3—12 | 1—3 | 15%—3.6% | 1.99×—7.44× |

[§]In the order of *conv1—conv2*
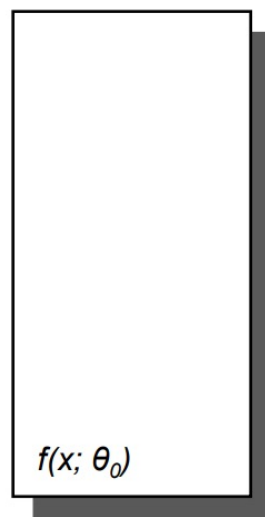
LeNet 1

LeNet 2

LeNet 3

conv1 filters

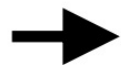Fewer but smoother feature extractors

# Lottery Ticket Hypothesis

**The Lottery Ticket Hypothesis.** *A randomly-initialized, dense neural network contains a subnetwork that is initialized such that—when trained in isolation—it can match the test accuracy of the original network after training for at most the same number of iterations.*
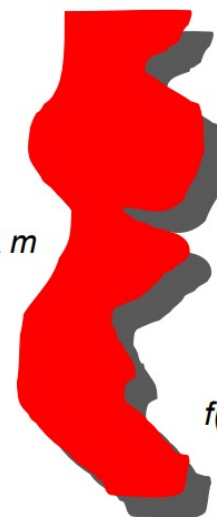
Original network

Winning Ticket

Prune $p\%$

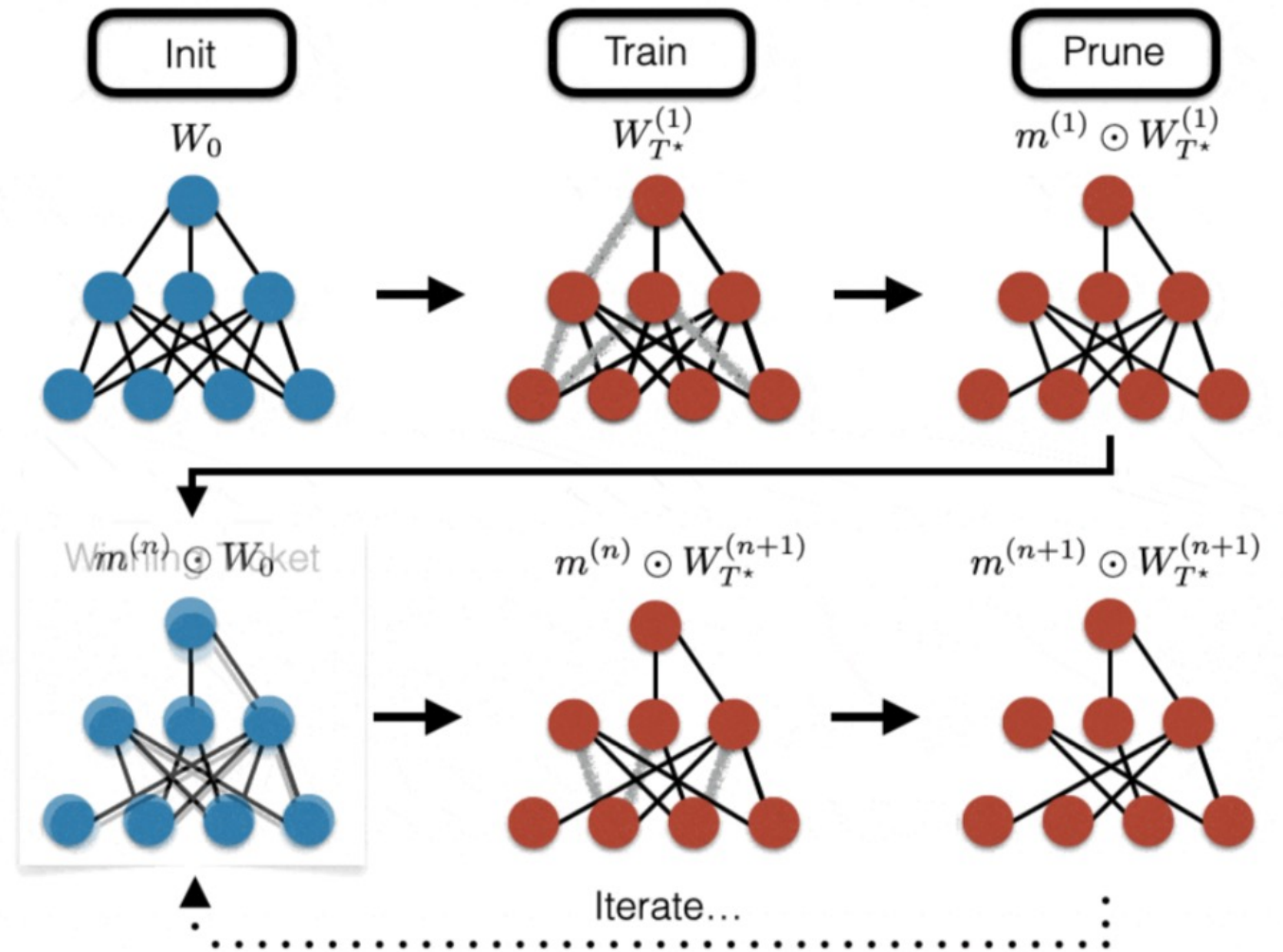Mask $m$

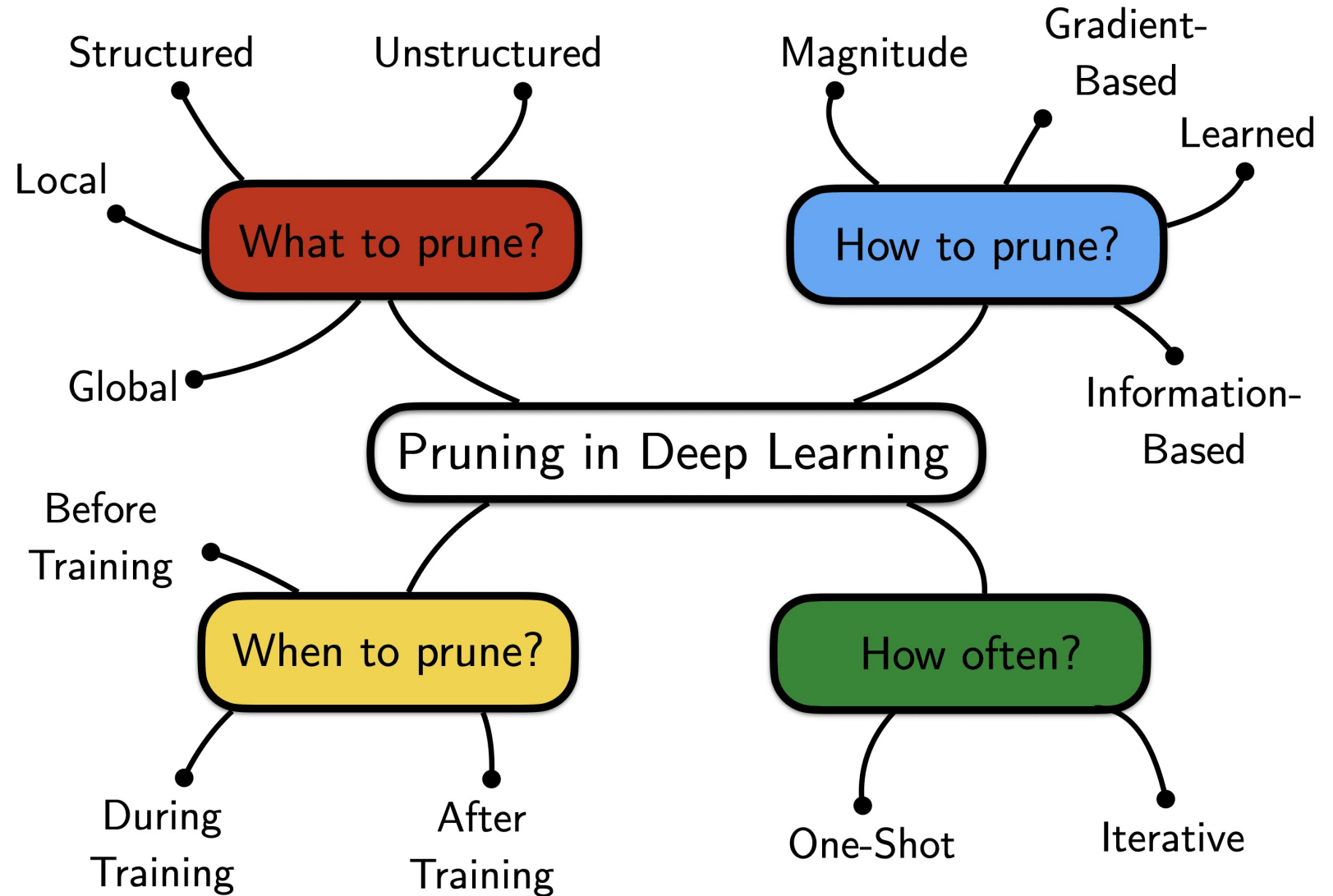$f(x; \theta_0)$

$f(x; m \odot \theta_0)$

- Winning Ticket gives
  - Better or same results
  - Shorter or same training time
  - Notably fewer parameters
  - Is trainable from the beginning

Lottery Ticket Hypothesis

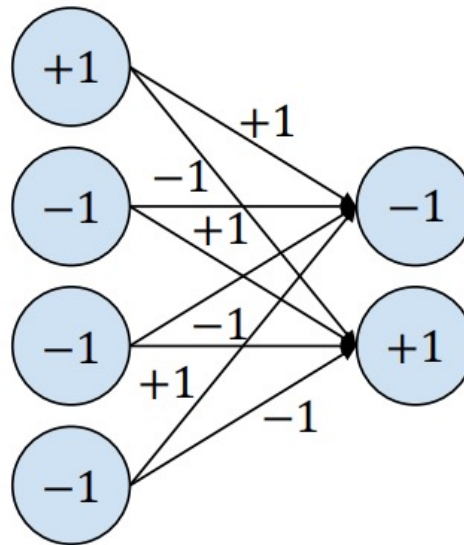Searching for Tickets: Iterative Magnitude Pruning

Init — $W_0$
Train — $W_{T^\star}^{(1)}$
Prune — $m^{(1)} \odot W_{T^\star}^{(1)}$

Winning Ticket — $m^{(n)} \odot W_0$
$m^{(n)} \odot W_{T^\star}^{(n+1)}$
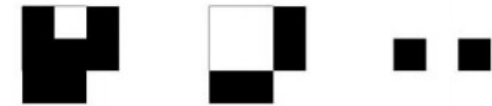$m^{(n+1)} \odot W_{T^\star}^{(n+1)}$

Iterate…

# More About Quantization

- Neural networks can be even **binarized** (**+1** or **-1**)
  - DNNs trained to use binary weights and binary activations

- Expensive **32-bit MAC** (**M**ultiply-**AC**cumulate) ⇒ Cheap **1-bit XNOR-Count**
  - "MAC == XNOR-Count": when the weights and activations are $\pm 1$

# 1s in bits



**Binarized weights**

**Binarized feature maps**

## Binary Neural Networks

- **Idea**: Training real-valued nets ($W_r$) treating binarization ($W_b$) as noise
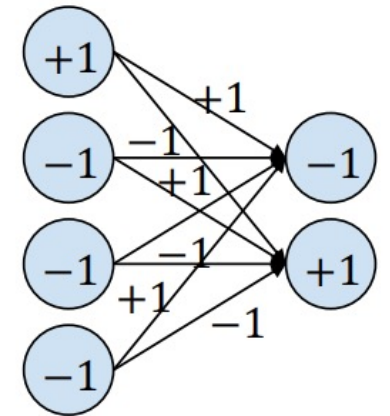  - Training $W_r$ is done by **stochastic gradient descent**

- **Binarization** ($W_r \rightarrow W_b$) occurs for each forward propagation
  - On each of **weights**: $W_b = \mathrm{sign}(W_r)$
  - ... also on each **activation**: $a_b = \mathrm{sign}(a_r)$

- Gradients for $W_r$ is estimated from $\frac{\partial L}{\partial W_b}$ [Bengio et al., 2013]
  - "Straight-through estimator": Ignore the binarization during backward!

$$\frac{\partial L}{\partial W_r} = \frac{\partial L}{\partial W_b} \mathbf{1}_{|W_r| \leq 1}$$
$$\frac{\partial L}{\partial a_r} = \frac{\partial L}{\partial a_b} \mathbf{1}_{|a_r| \leq 1}$$

  - Cancelling gradients for better performance
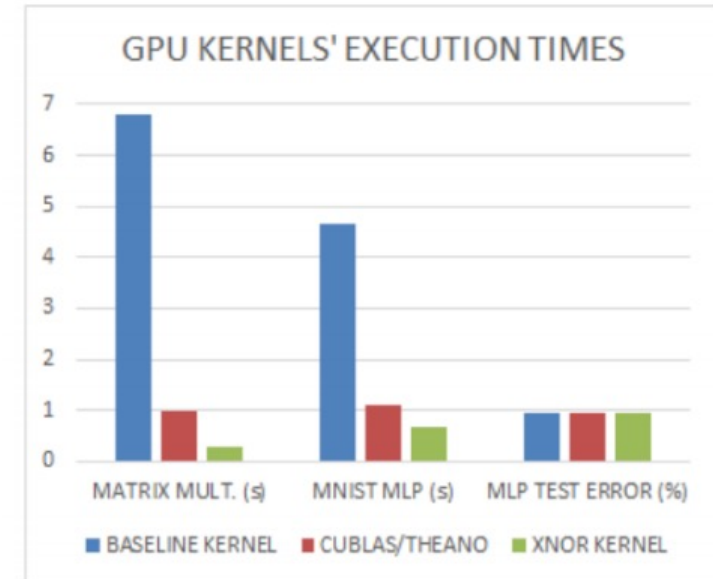    - When the value is too large

## Binary Neural Networks

- BNN yields **32x less memory** compared to the baseline 32-bit DNNs
  - … also expected to reduce energy consumption drastically

- **23x faster** on kernel execution times
  - BNN allows us to use XNOR kernels
  - **3.4x** faster than cuBLAS

| Operation | MUL | ADD |
|---|---|---|
| 8bit Integer | 0.2pJ | 0.03pJ |
| 32bit Integer | 3.1pJ | 0.1pJ |
| 16bit Floating Point | 1.1pJ | 0.4pJ |
| 32tbit Floating Point | 3.7pJ | 0.9pJ |



GPU KERNELS' EXECUTION TIMES
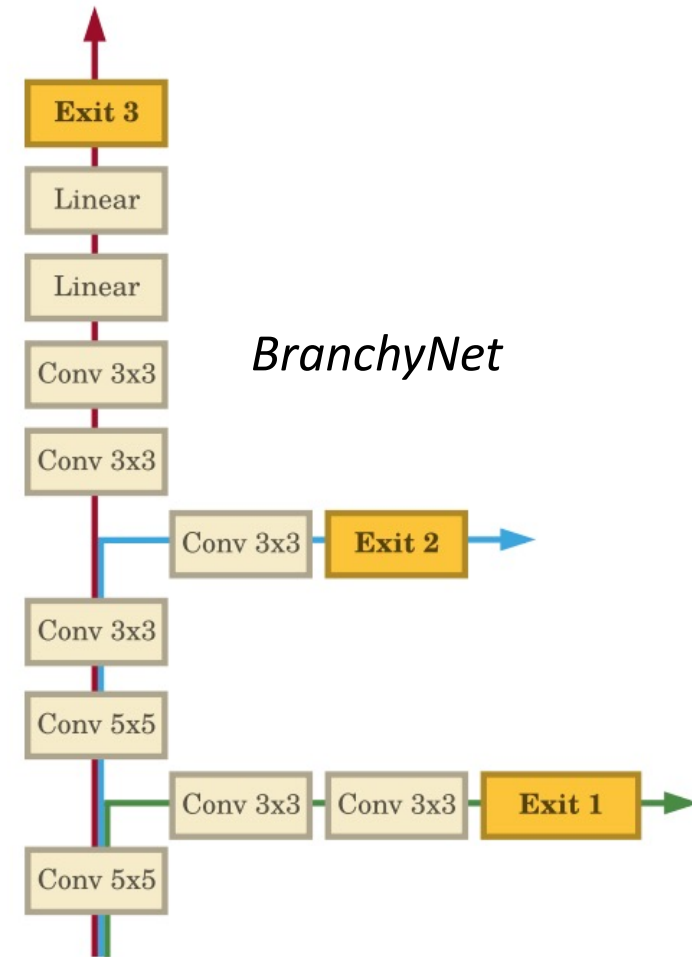
BASELINE KERNEL   CUBLAS/THEANO   XNOR KERNEL

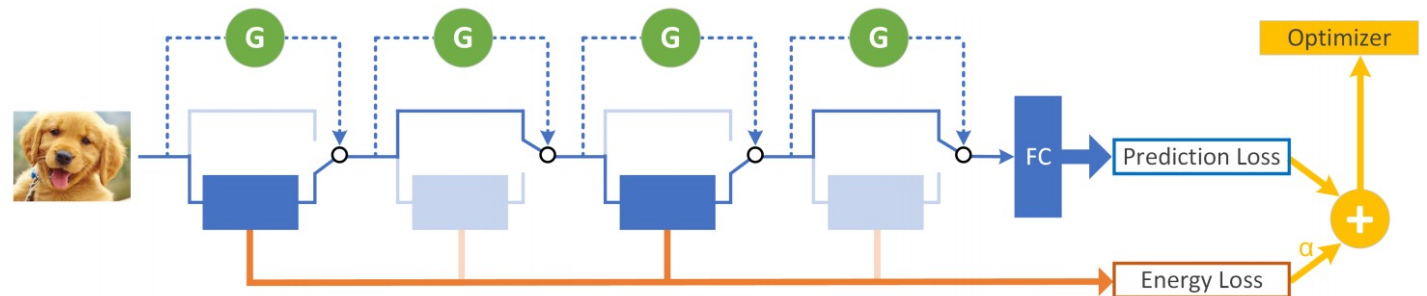- **BNN** achieves comparable error rates over existing DNNs

# Dynamic Inference

- Only execute a fraction of the network per needed

- Can enable both "input-dependent" and "resource-dependent" forms

*BranchyNet*



*SkipNet*

# Real-World Efficient ML: Way to Go

- Jointly utilizing several compression means
  - Also, can choose efficient "by-design" models (MobileNets, or even non-deep models, etc.)
  - Channel pruning is in fact very similar to NAS
- **Data processing** is often a key concern, maybe more important
- **Hardware co-design** is another key concern
- Resource constraints & user demands often **change over time**
- From single task to multi-task and lifelong learning ...

# Demo: Energy-Efficient UAV-Based Text Spotting System

- **Task:** accurate detecting signs and recognizing texts in the video, captured by an unmanned aerial vehicle (UAV), with minimal energy cost as possible (**Hardware:** Raspberry Pi 3B+)

- Our solution won 2nd prize in the high-visibility IEEE CVPR **2020 Low-Power Computer Vision (LPCV) Challenge,** among 11 university & company teams that submitted 84 independent solutions.

① **Only homogeneous regions**

**Dropped**

② **Unlikely with texts**

**Dropped**

③ **Poor-quality texts**

**Dropped**

④ **High-quality, likely with texts**
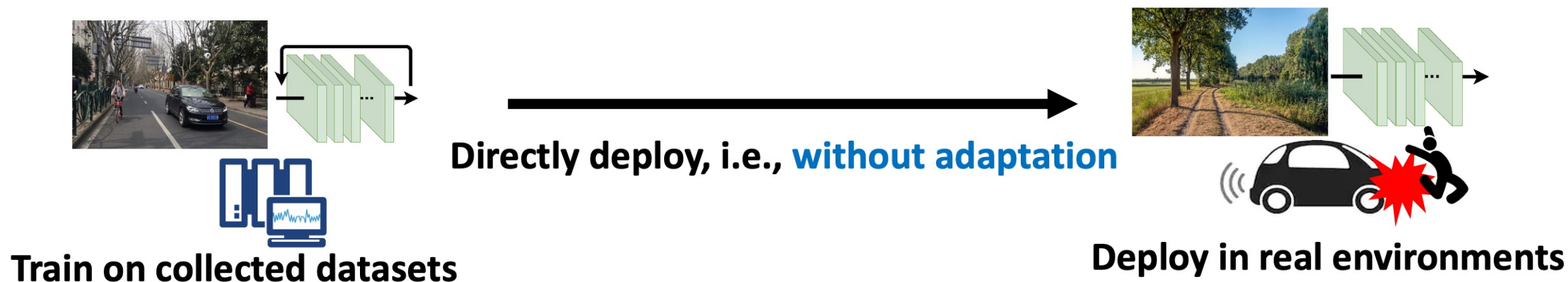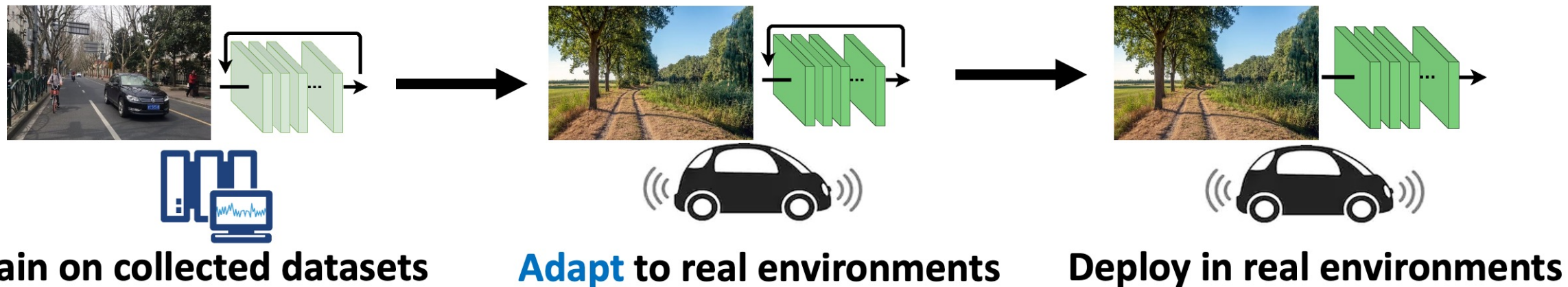
**Processed by OCR**

# Energy-Efficient Training: Prevailing Demands





- Shifting model training from the cloud to the edge
  - Facilitating personalization; saving bandwith/communication energy; protecting privacy

- Deep learning has a terrible carbon footprint
  - "**Training a single AI model can emit as much carbon as five cars in their lifetimes**", MIT Tech Review

# On-Device Training (Adaptation) is on Growing Demand



**Train on collected datasets** → **Adapt to real environments** → **Deploy in real environments**

**Train on collected datasets** → Directly deploy, i.e., **without adaptation** → **Deploy in real environments**
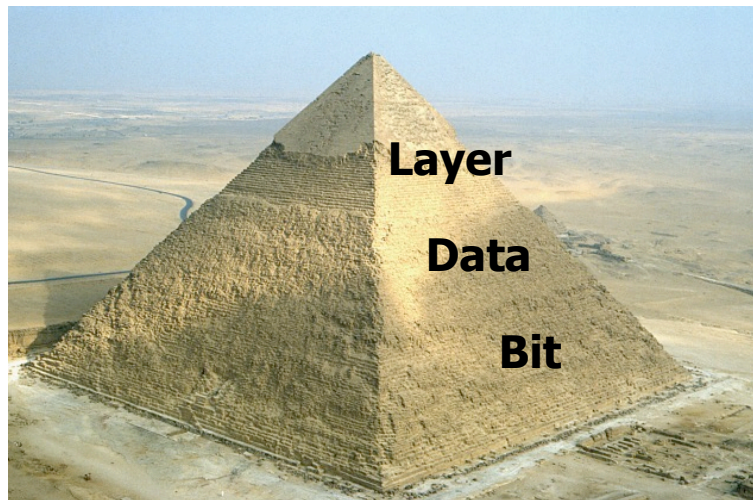
# Problem Setting

- We consider the most basic CNN training, assuming both the model structure and the dataset to be pre-given, training from scratch
  - Trim down the total <span style="color:red">energy cost</span> for <span style="color:red">in-situ, resource-constrained</span> training.
  - not usually the realistic IoT case, but address it as a starting point

- Many existing works are on accelerated CNN training
  - … they mostly focus on reducing the total <span style="color:blue">training time</span> in <span style="color:blue">resource-rich</span> settings, such as by <span style="color:blue">distributed</span> training in **large-scale** GPU clusters

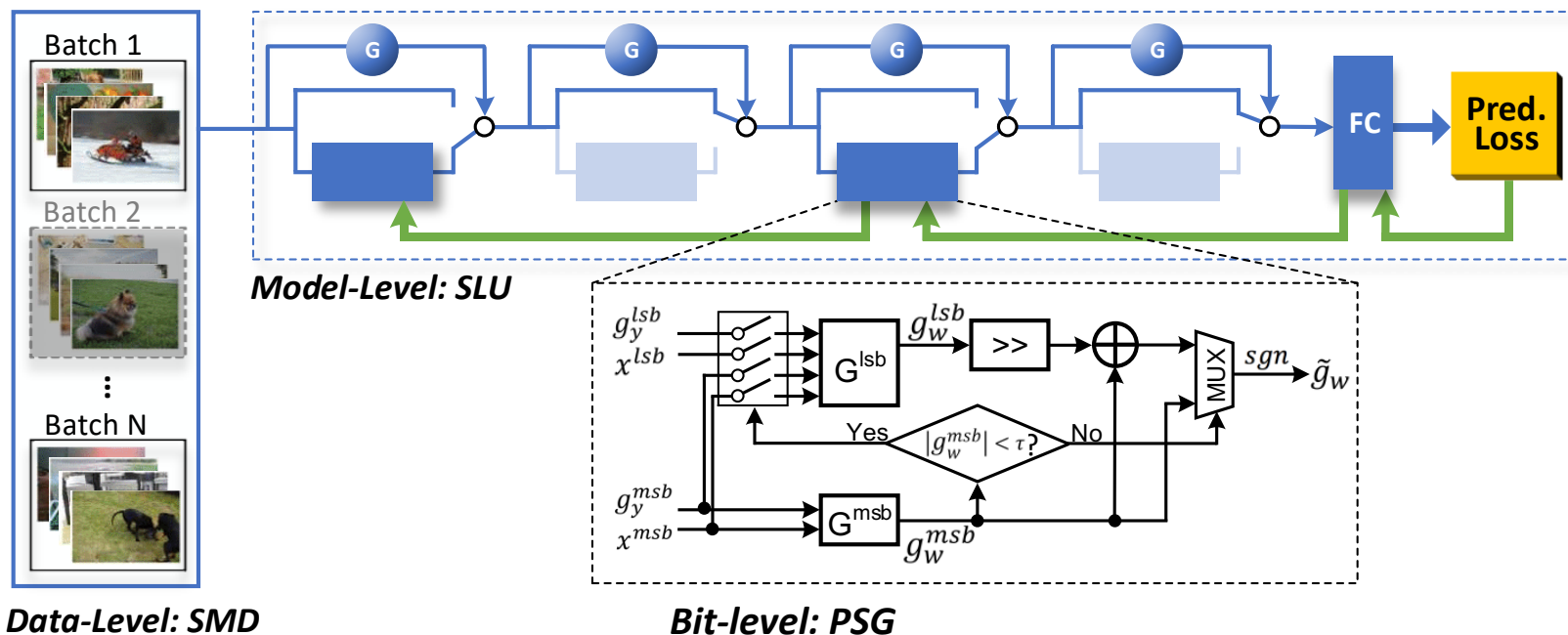# From Inference to Training: Lessons and Challenges

- **Training v.s. Inference:** one-pass feedforward v.s. iterative forward + backward

- **Lessons that we learned from Inference:**
  - Model parameters are not born equally, and many redundancies do exist
  - _Know your specific goal:_ saving memory, latency and energy are often not aligned
  - To achieve energy goal, realistic energy models and/or hardware measurements are very helpful
  - Consider a more "end-to-end" effort beyond just the model itself (data, hardware, architecture…)

- **New Challenges posed for Training:**
  - Saving per-sample (mini-batch) complexity (both feed-forward and backward)
  - The empirical convergence (how many iterations needed) matters more than per-MB complexity
  - Data access/movement bottlenecks are (even more) crucial

# E2-Train: Energy-Efficient CNN Training [NeurIPS'2019]

**Motivation:**



Layer

Data

Bit



*Model-Level: SLU*

*Data-Level: SMD*

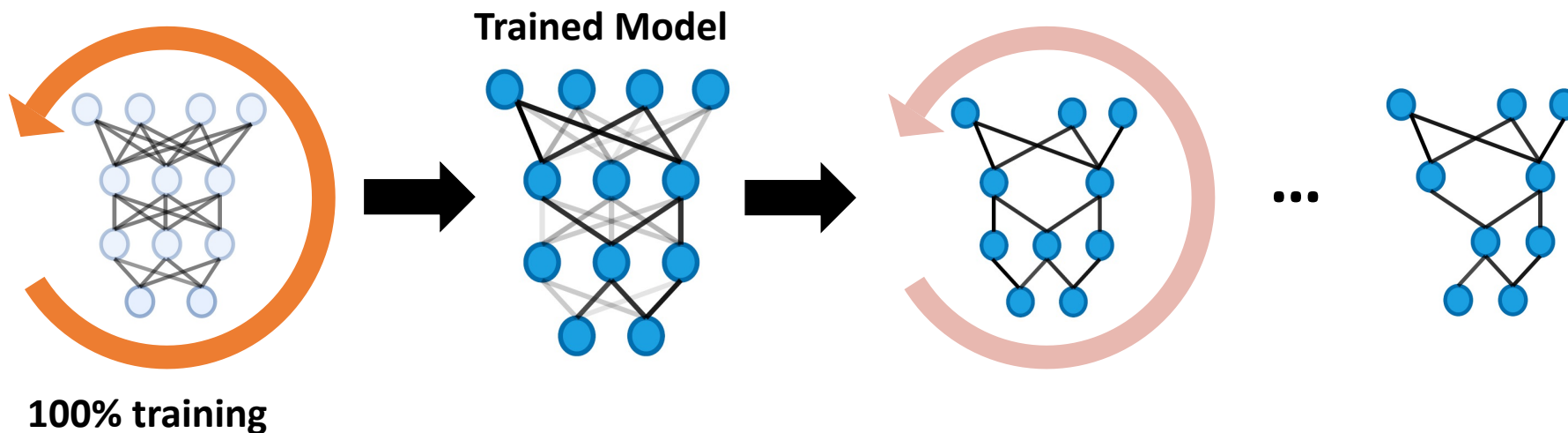*Bit-level: PSG*

**"Three-Pronged" Approach:**
- **Data-Level:** stochastic mini-batch dropping
- **Layer-Level:** selective layer update
- **Bit-Level:** predictive sign gradient descent

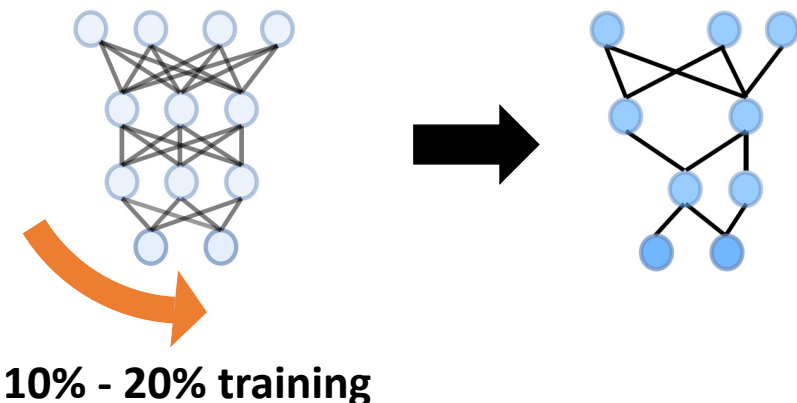| Datasets | Models | Accuracy (vs. Original One) | Energy Savings |
|----------|--------|------------------------------|----------------|
| CIFAR-10 | MobileNetV2 | 92.06% (vs. 92.47%) | **88%** |
|          | ResNet-110 | 93.01% (vs. 93.57%) | **83%** |
| CIFAR-100 | MobileNetV2 | 71.61% (vs. 71.91%) | **88%** |
|           | ResNet-110 | 71.63% (vs. 71.60%) | **84%** |

Energy savings is quantified based on **FPGA implementation**

- **Progressive Pruning and Training** (e.g., [J. Frankle, ICLR 2019])

**Trained Model**



**100% training**

- **Early-Bird Train (Proposed)**



**10% - 20% training**

**For the first time:**

1. We **discover the existence** of Early-Bird (EB) Tickets

2. We **propose a detector** of low cost to detect EB Tickets

3. We leverage the existence of EB Tickets to **develop an efficient training scheme**

➢ **5.8× - 10.7× reduced training energy** with a comparable or even better accuracy over the most competitive baseline

32

The University of Texas at Austin

# Electrical and Computer Engineering

*Cockrell School of Engineering*