

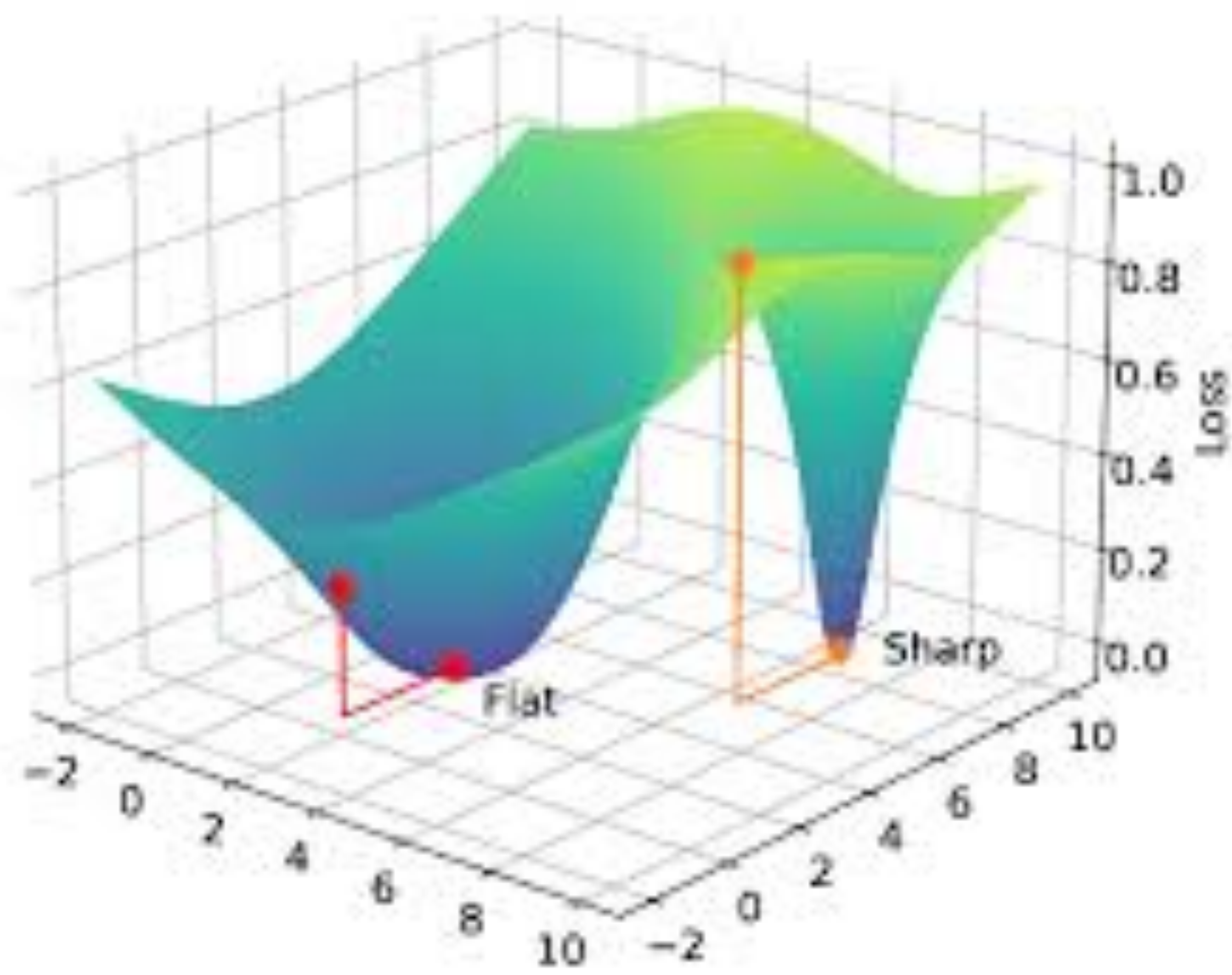
CSCI1470

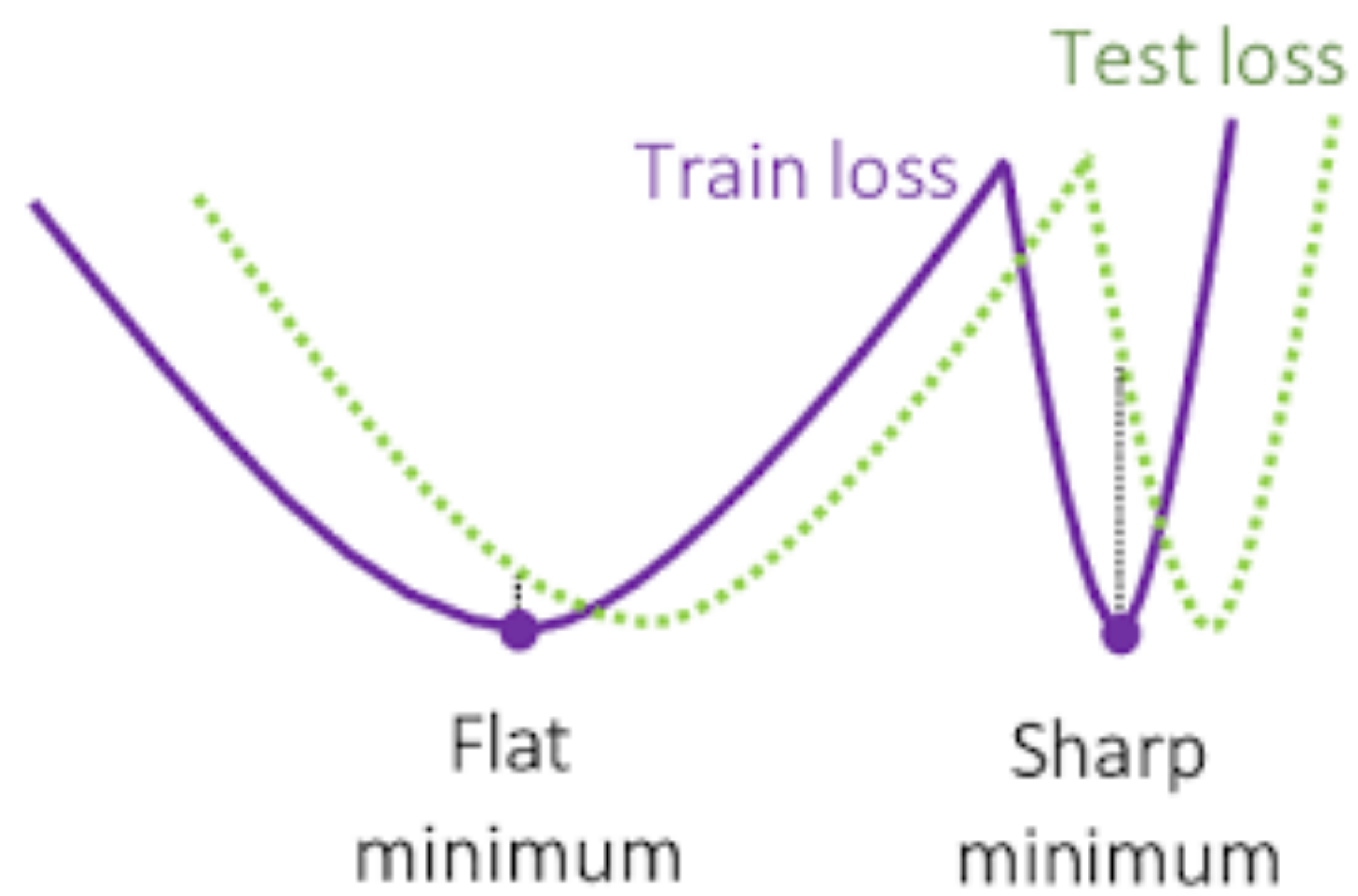
Deep Learning

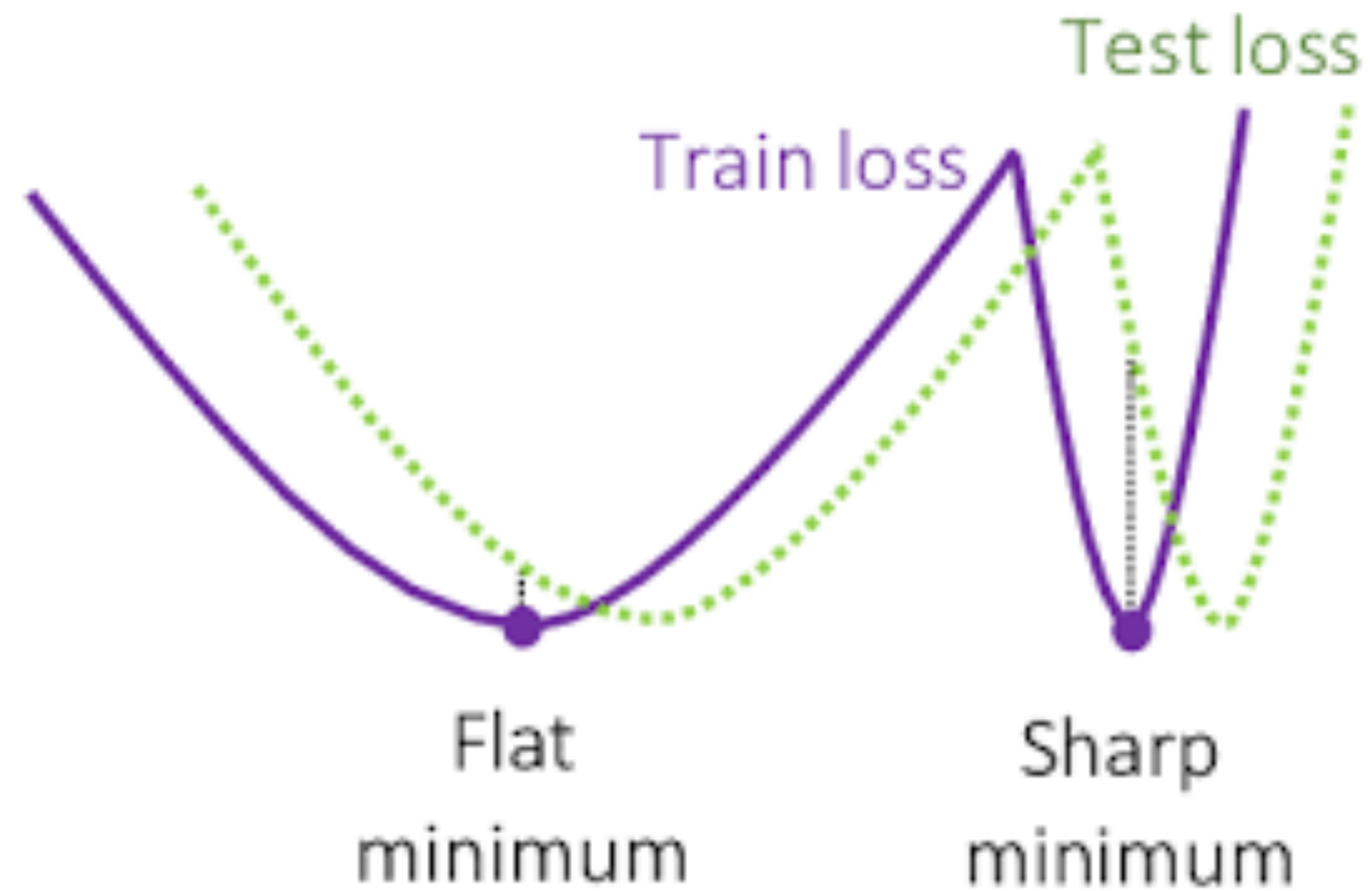
Randall Balestrieri

Recap

Optimization for Better Generalization?







Any idea?

The classic optimization framework of machine learning is empirical risk minimization

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N \ell(\mathbf{x}_i; \boldsymbol{\theta}) \quad (1)$$

where $\boldsymbol{\theta} \in \mathbb{R}^d$ is a vector of parameters, $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ is a training set of inputs $\mathbf{x}_n \in \mathbb{R}^D$, and $\ell(\mathbf{x}; \boldsymbol{\theta})$ is a loss function quantifying the performance of parameters $\boldsymbol{\theta}$ on \mathbf{x} . SGD samples a minibatch $\mathcal{S} \subset \{1, \dots, N\}$ of size $|\mathcal{S}| \ll N$ from the training set and updates the parameters through

$$\boldsymbol{\theta}_{t+1}^{\text{SGD}} = \boldsymbol{\theta}_t - \eta \mathbf{g}(\boldsymbol{\theta}_t), \text{ where } \mathbf{g}(\boldsymbol{\theta}) = \frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \nabla \ell(\boldsymbol{\theta}; \mathbf{x}_i), \quad (2)$$

for a length specified by η , the learning rate.

SHARPNESS-AWARE MINIMIZATION FOR EFFICIENTLY IMPROVING GENERALIZATION

Pierre Foret*

Google Research

pierre.pforet@gmail.com

Ariel Kleiner

Google Research

akleiner@gmail.com

Hossein Mobahi

Google Research

hmobahi@google.com

Behnam Neyshabur

Blueshift, Alphabet

neyshabur@google.com

Aware Minimization (SAM) problem:

$$\min_{\mathbf{w}} L_S^{SAM}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2 \quad \text{where} \quad L_S^{SAM}(\mathbf{w}) \triangleq \max_{\|\boldsymbol{\epsilon}\|_p \leq \rho} L_S(\mathbf{w} + \boldsymbol{\epsilon}), \quad (1)$$

SHARPNESS-AWARE MINIMIZATION FOR EFFICIENTLY IMPROVING GENERALIZATION

Pierre Foret *
Google Research
pierre.pforet@gmail.com

Ariel Kleiner
Google Research
akleiner@gmail.com

Hossein Mobahi
Google Research
hmobahi@google.com

Behnam Neyshabur
Blueshift, Alphabet
neyshabur@google.com

Aware Minimization (SAM) problem:

$$\min_{\mathbf{w}} L_S^{SAM}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2 \quad \text{where} \quad L_S^{SAM}(\mathbf{w}) \triangleq \max_{\|\boldsymbol{\epsilon}\|_p \leq \rho} L_S(\mathbf{w} + \boldsymbol{\epsilon}), \quad (1)$$

Remind you of anything?

Averaging Weights Leads to Wider Optima and Better Generalization

Pavel Izmailov*¹ **Dmitrii Podoprikin***^{2,3} **Timur Garipov***^{4,5} **Dmitry Vetrov**^{2,3} **Andrew Gordon Wilson**¹
¹Cornell University, ²Higher School of Economics, ³Samsung-HSE Laboratory,
⁴Samsung AI Center in Moscow, ⁵Lomonosov Moscow State University

Algorithm 1 Stochastic Weight Averaging [48]

Input: Loss function \mathcal{L} , training budget in number of iterations b , training dataset $\mathcal{D} := \cup_{i=1}^n \{\mathbf{x}_i\}$, mini-batch size $|\mathcal{B}|$, averaging start epoch E , averaging frequency ν , (scheduled) learning rate η , initial weights θ_0 .

for $k \leftarrow 1, \dots, b$ **do**

 Sample a mini-batch \mathcal{B} from \mathcal{D}

 Compute gradient $\mathbf{g} \leftarrow \nabla \mathcal{L}(\theta_t)$

 Update parameters $\theta_{t+1} \leftarrow \theta_t - \eta \mathbf{g}$

if $k \geq E$ and $\text{mod}(k, \nu) = 0$ **then**

$$\theta_{t+1}^{\text{SWA}} = (\theta_t^{\text{SWA}} \cdot l + \theta_{t+1}^{\text{SWA}}) / (l + 1)$$

end if

end for

return θ^{SWA}

Algorithm 2 Sharpness-Aware Minimization [22]

Input: Loss function \mathcal{L} , training budget in number of iterations b , training dataset $\mathcal{D} := \cup_{i=1}^n \{\mathbf{x}_i\}$, mini-batch size $|\mathcal{B}|$, neighborhood radius ρ , (scheduled) learning rate η , initial weights θ_0 .

for $k \leftarrow 1, \dots, b$ **do**

 Sample a mini-batch \mathcal{B} from \mathcal{D}

 Compute worst-case perturbation

$$\hat{\epsilon} \leftarrow \rho \frac{\nabla \mathcal{L}(\theta)}{\|\nabla \mathcal{L}(\theta)\|_2}$$

 Compute gradient $\mathbf{g} \leftarrow \nabla \mathcal{L}(\theta_t^{\text{SAM}} + \hat{\epsilon})$

 Update parameters $\theta_{t+1}^{\text{SAM}} \leftarrow \theta_t^{\text{SAM}} - \eta \mathbf{g}$

end for

return θ^{SAM}

Efficiently Seeking Flat Minima for Better Generalization in Fine-Tuning Large Language Models and Beyond

Jiaxin Deng^{1*}, Qingcheng Zhu^{2*}, Junbiao Pang^{1†}, Linlin Yang³, Zhongqian Fu⁴, Baochang Zhang^{5‡}

¹School of Information Science and Technology, Beijing University of Technology, Beijing, China

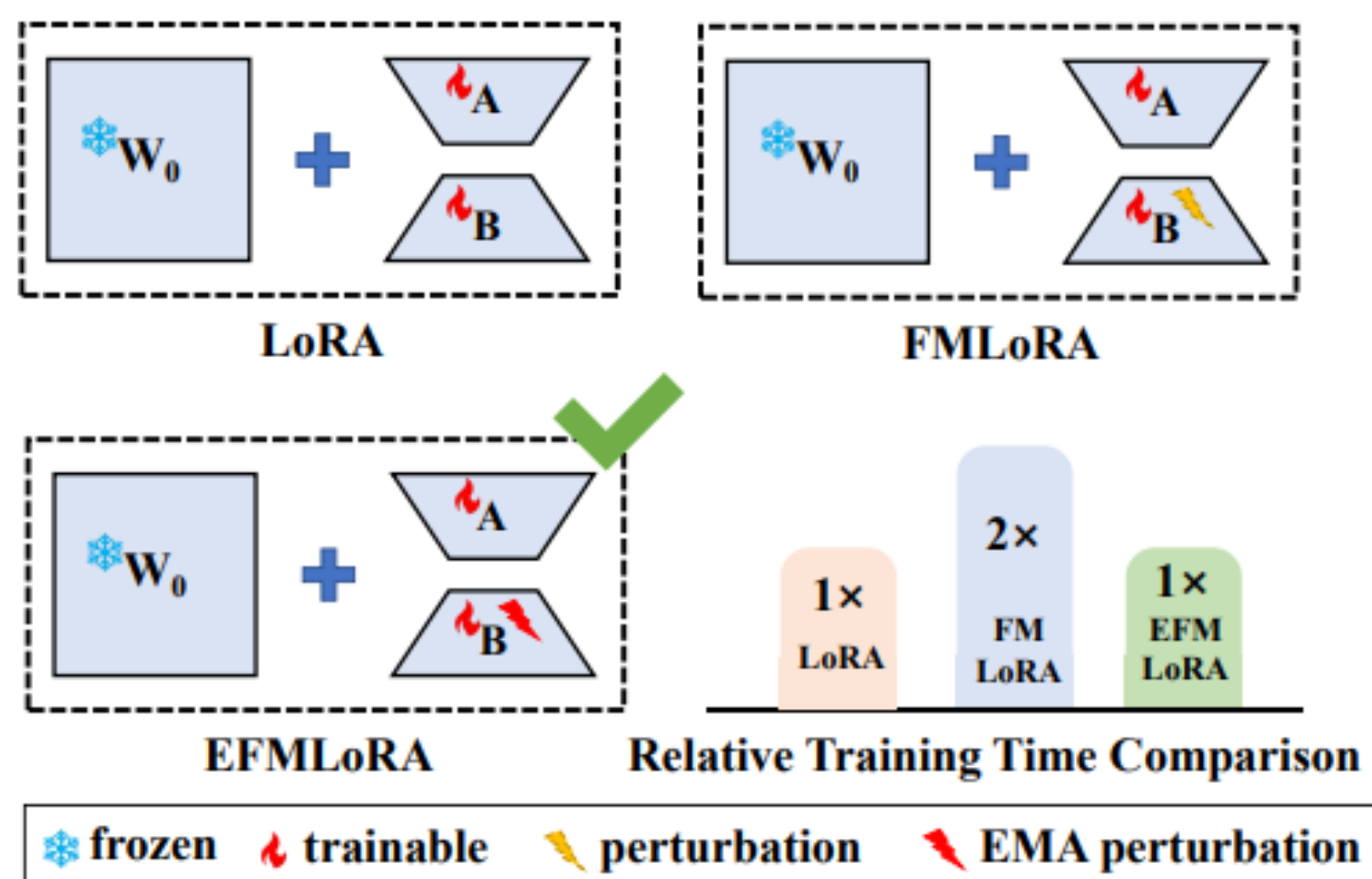
²School of Electronic Information Engineering, Beihang University, Beijing, China

³State Key Laboratory of Media Convergence and Communication, Communication University of China, Beijing, China

⁴Huawei Noah's Ark Lab, China

⁵Hangzhou Research Institute, School of Artificial Intelligence, Beihang University, China

dengjiaxin@emails.bjut.edu.cn, zhuqc@buaa.edu.cn, junbiao_pang@bjut.edu.cn, lyang@cuc.edu.cn, fuzhongqian@huawei.com, bczhang@buaa.edu.cn



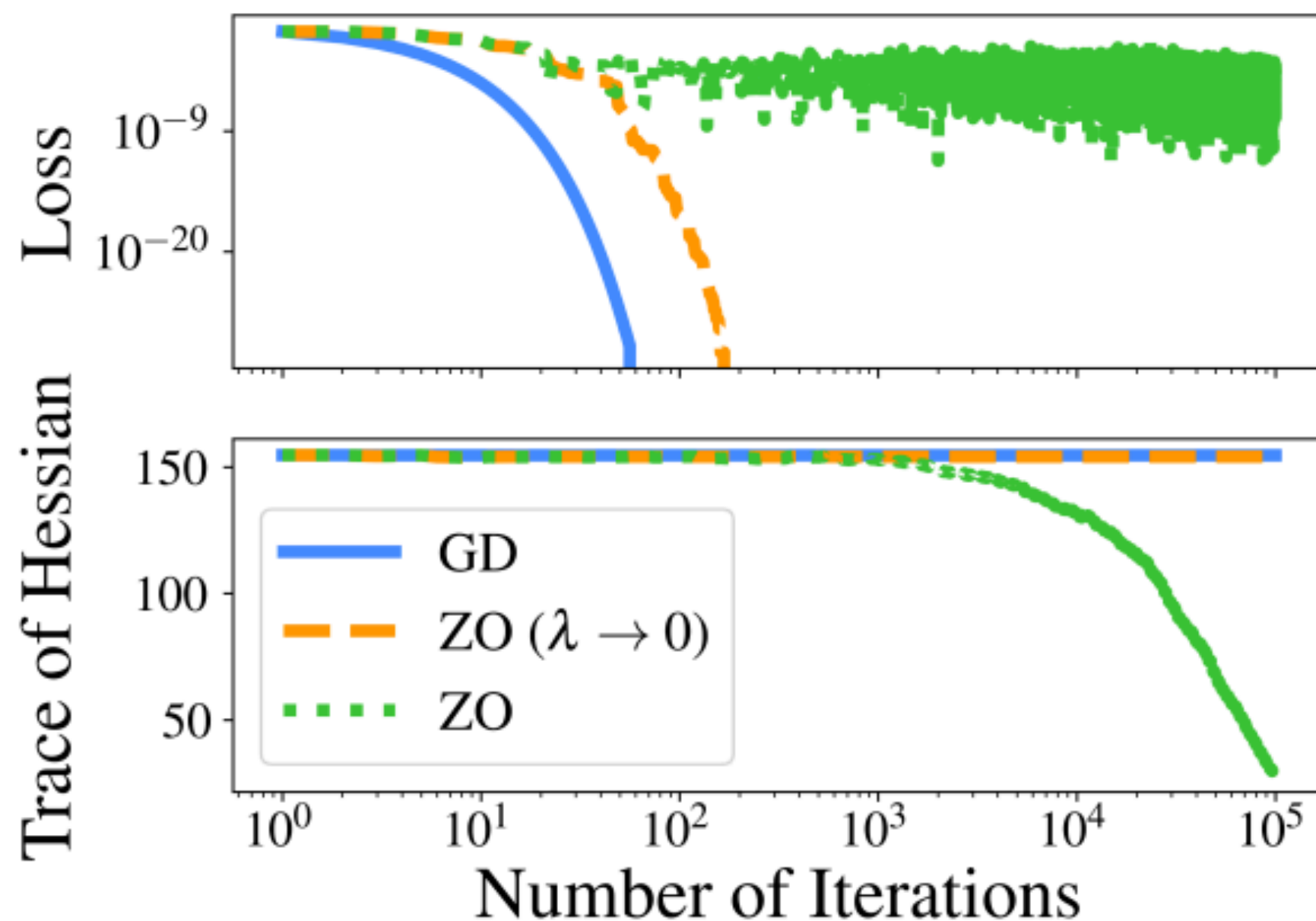
Zeroth-Order Optimization Finds Flat Minima

Liang Zhang^{1,2}, Bingcong Li¹, Kiran Koshy Thekumparampil³,
Sewoong Oh⁴, Michael Muehlebach², and Niao He¹

To be more specific, for the optimization problem $\min_{x \in \mathbb{R}^d} f(x)$ with parameters $x \in \mathbb{R}^d$ and a loss function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, zeroth-order optimization with the standard two-point gradient estimator [74] (Algorithm 1) iteratively updates x by substituting the computationally intractable gradient with

$$g_\lambda(x, u) := \frac{f(x + \lambda u) - f(x - \lambda u)}{2\lambda} u, \quad (1)$$

where $u \sim \mathcal{N}(0, I_d)$ is a standard Gaussian random vector and $\lambda \in \mathbb{R}$ is a smoothing parameter. It



(a) Comparison of GD and ZO.

Switching to deployment!

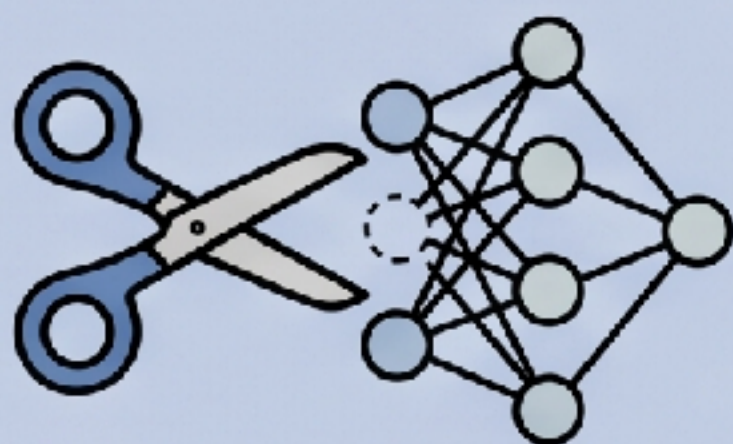
**What do we want to optimize
for?**

What do we want to optimize for?

- Latency
- Memory
- Hardware

OPTIMIZING A TRAINED MODEL FOR FAST & EFFICIENT DEPLOYMENT

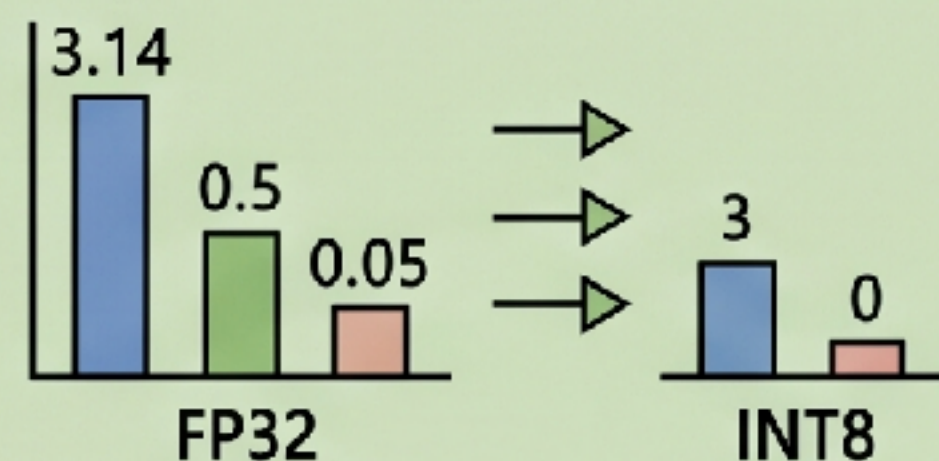
1. PRUNING



Reduce model size and computation by removing redundant weights and neurons.

- Remove connections with low contribution
- Creates a sparse model
- Achieves smaller memory footprint

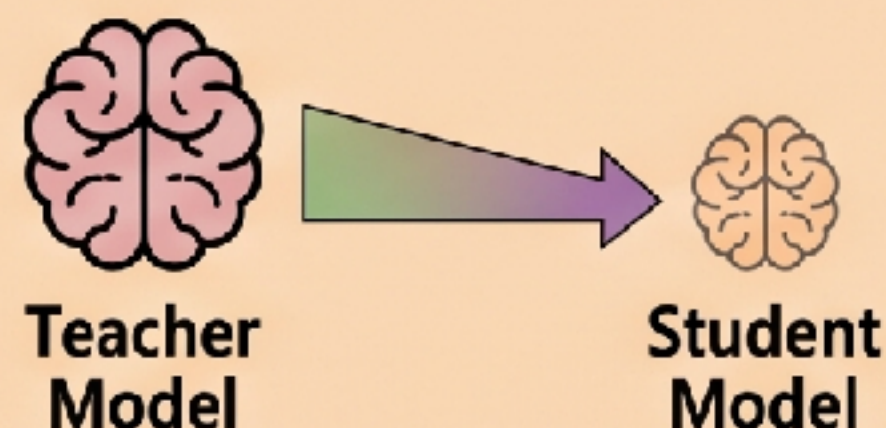
2. QUANTIZATION



Reduce precision of weights and activations with minimal accuracy loss.

- Convert from FP32 to smaller types (INT8, FP16)
- Significantly increases inference speed on compatible hardware
- Smaller memory usage

3. KNOWLEDGE DISTILLATION



Train a compact 'Student' model to mimic the performance of a large, complex 'Teacher' model.

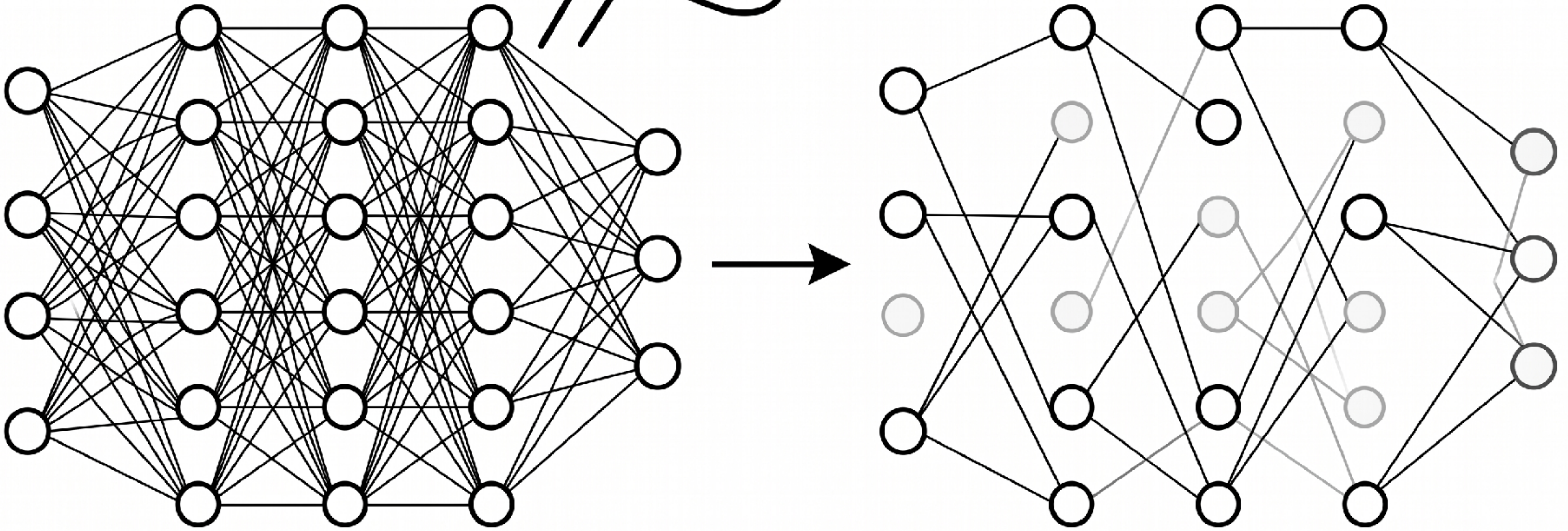
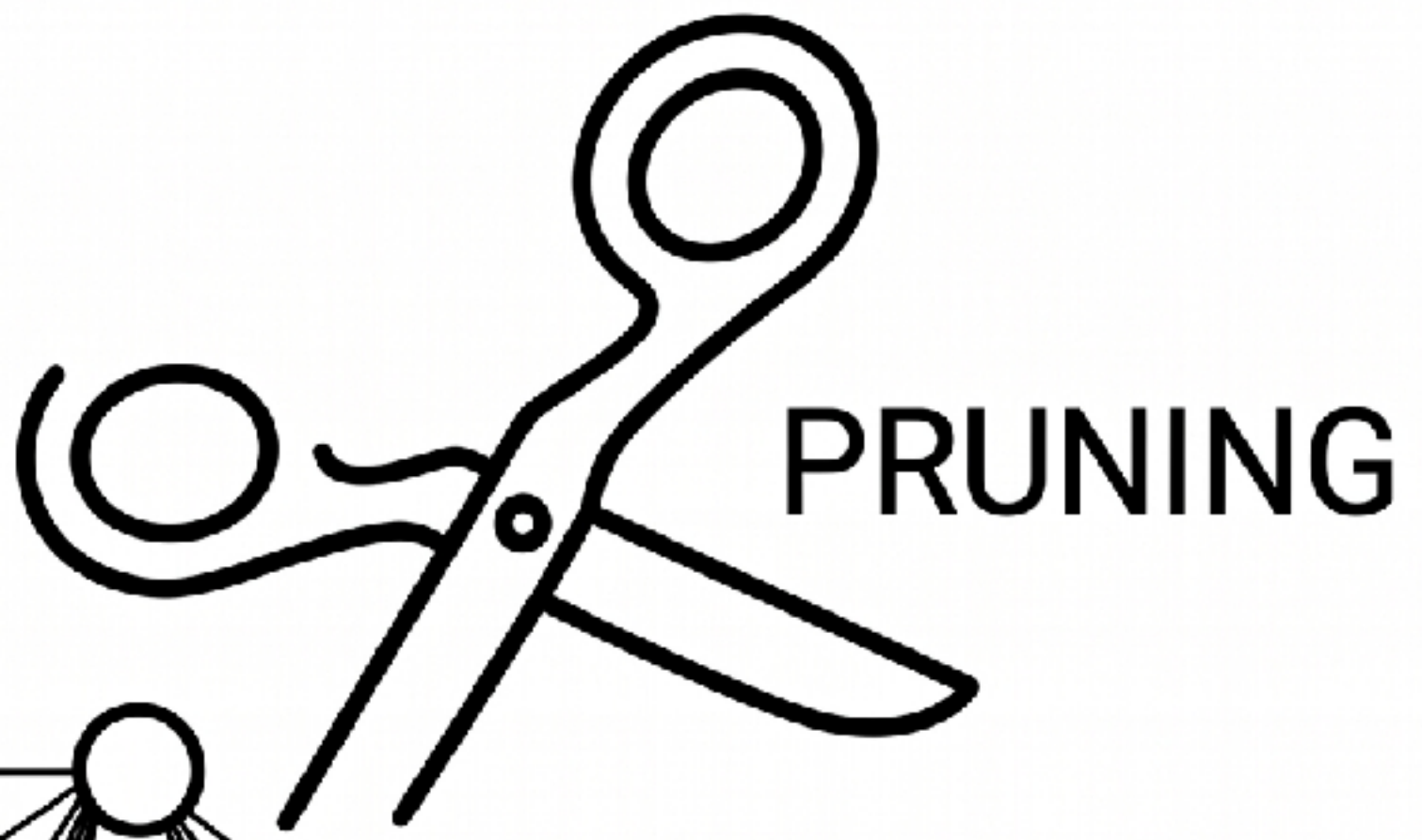
- Student model learns output distributions and representations
- Achieves comparable accuracy with much lower latency
- Ideal for edge devices

4. COMPILATION (e.g., TENSORRT)



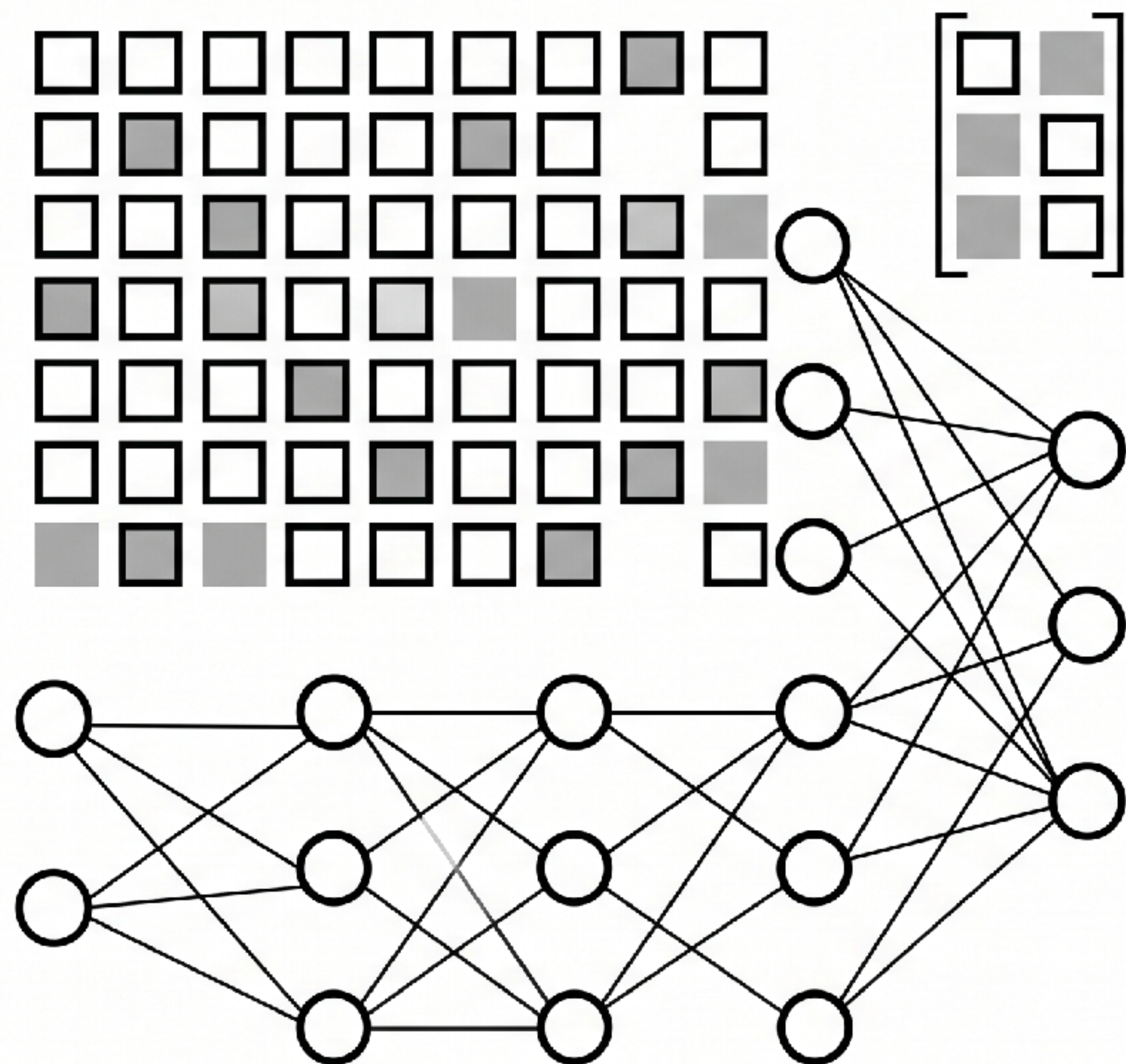
Apply a hardware-specific optimization compiler for maximum performance.

- Layer fusion, constant folding, and precision selection
- Tailored execution for target GPUs, CPUs, or NPUs
- Directly optimizes for target deployment environment



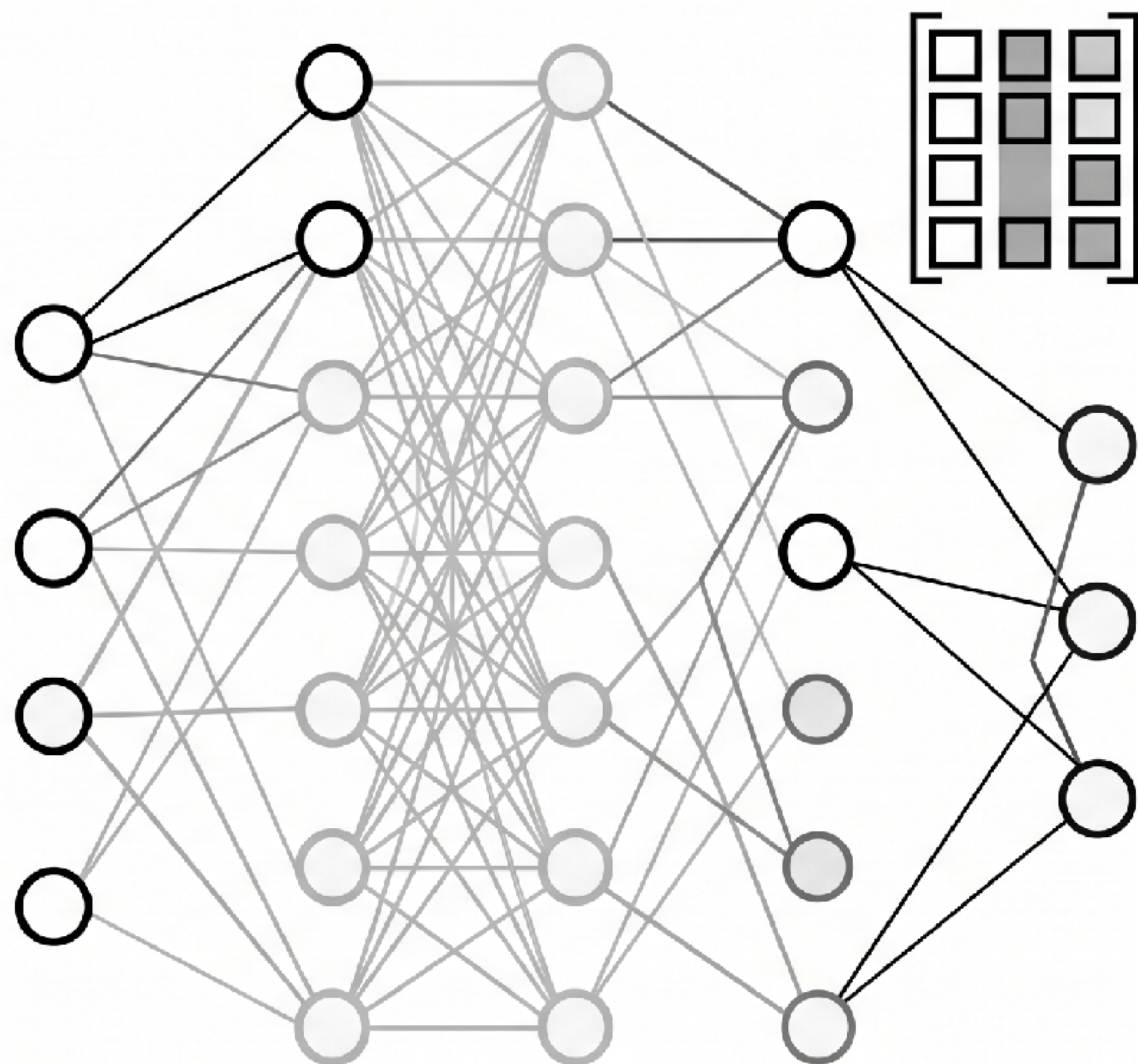
What is Model Pruning? (Simplifying the Model)

A Unstructured Pruning



Remove individual weights anywhere
(matrix becomes sparse)

B Structured Pruning



Remove entire neurons, channels, or layers
(retains structured connectivity)

Thank you!
See you Wednesday!