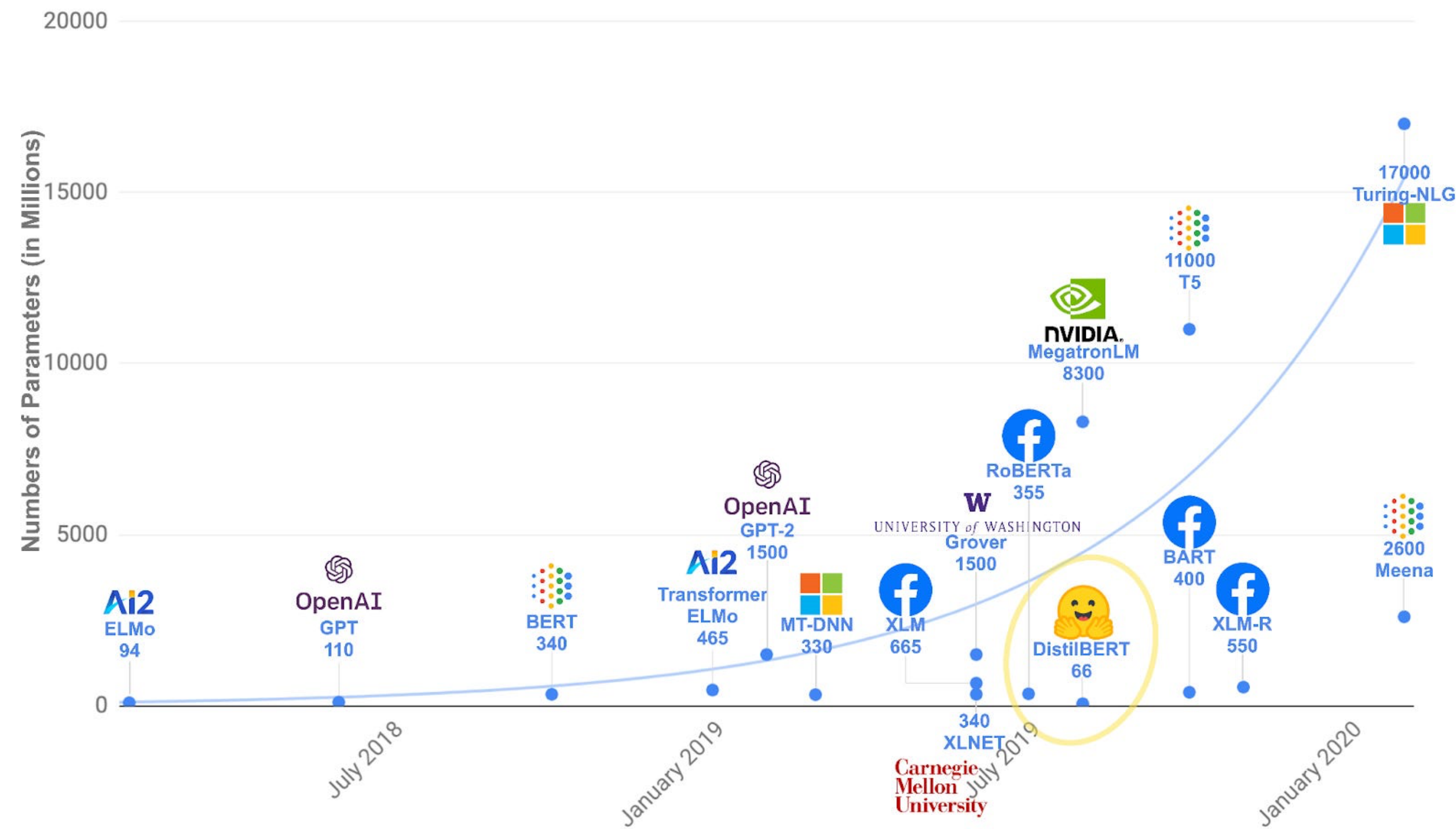# Efficient Inference

**Divyam Madaan**

October 17, 2024

# Overview of efficient inference

- Efficiency challenge

- Quantization

  - What is quantization and how to quantize?

- Pruning

  - Pruning before, during, after training

- Knowledge distillation

  - Distillation on outputs, weights and features.

- SCP tutorial

# Efficiency challenge

Size of models makes _____

# How does quantization look like in real-world?

Goal is to **reduce** the **number of bits (colors)** while **preserving the precision**
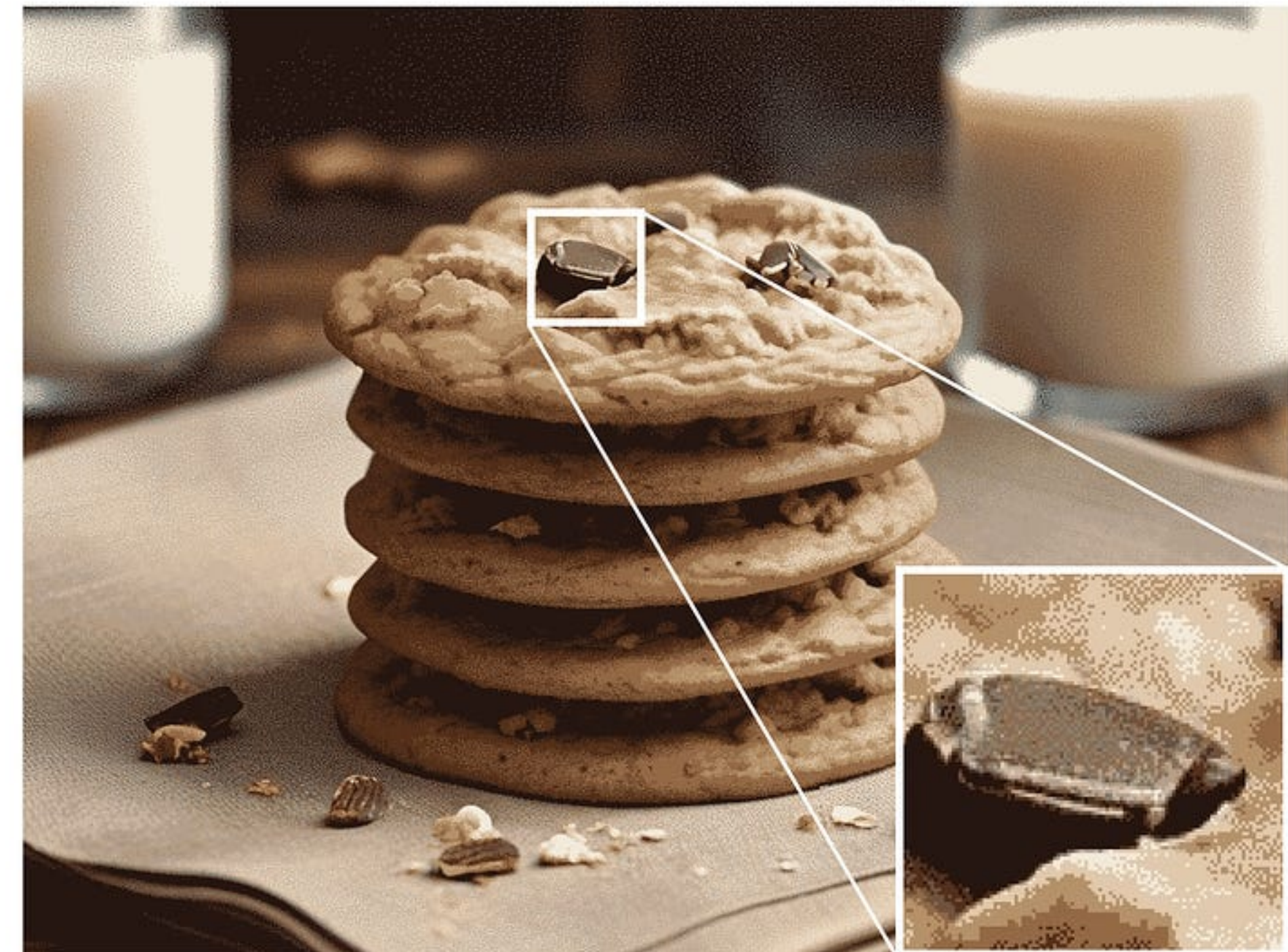

Original Image

# How does quantization look like in real-world?

Reducing the image to use just eight colors leads in a **loss of detail and precision**

# Common data types

Bits use _____ to represent a value



**Float 32-bit** (FP32)
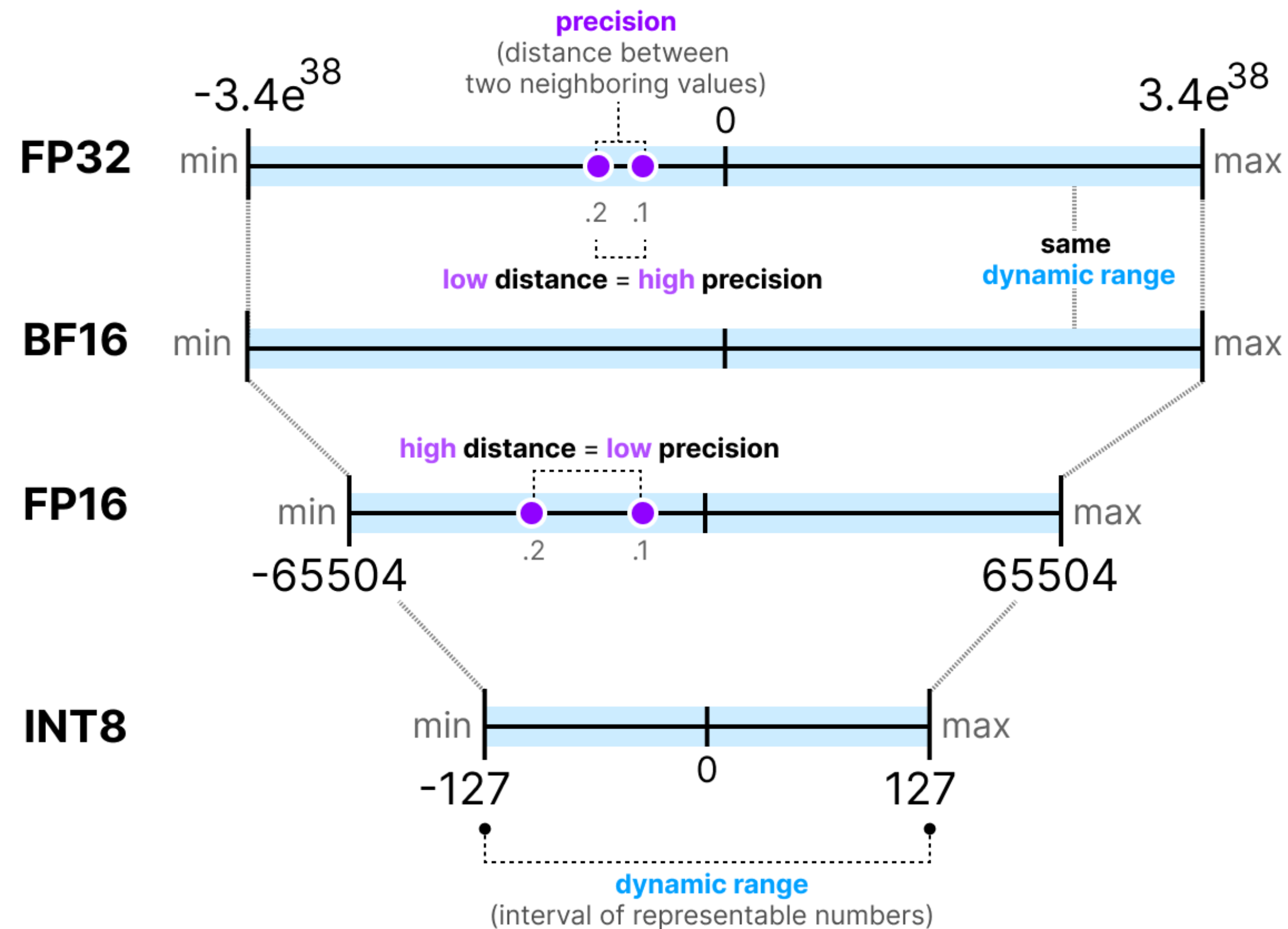
0  10000000  10010010000011111101011011

$(-1)^0 \times 2^1 \times 1.5707964 = $ **3.1415927410125732**

**higher** precision

$$(-1)^{\text{sign}} \times \text{base}^{\text{exponent}} \times \text{fraction}$$
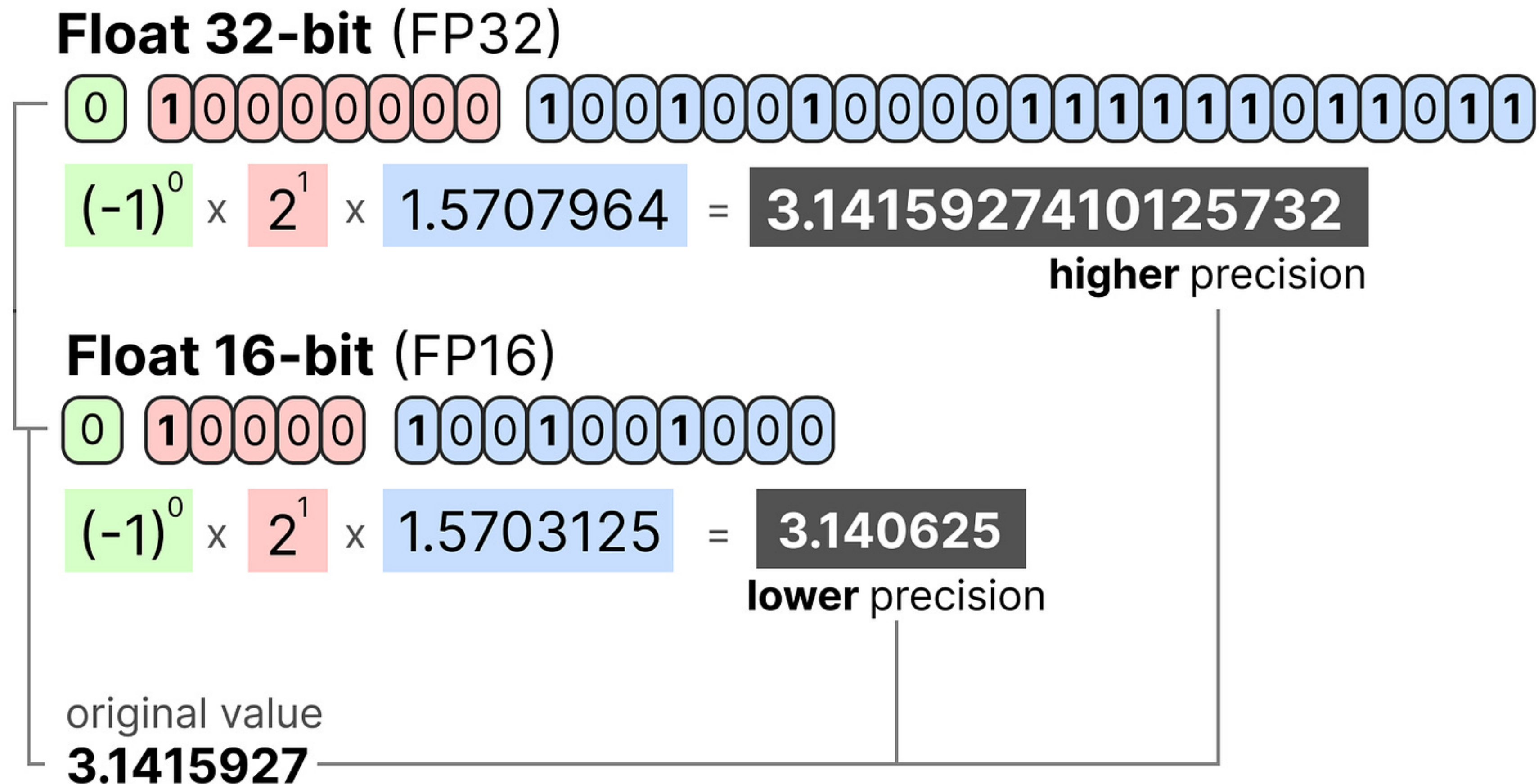
# What is precision?

**Precision** is a measure of how precisely a number can be represented



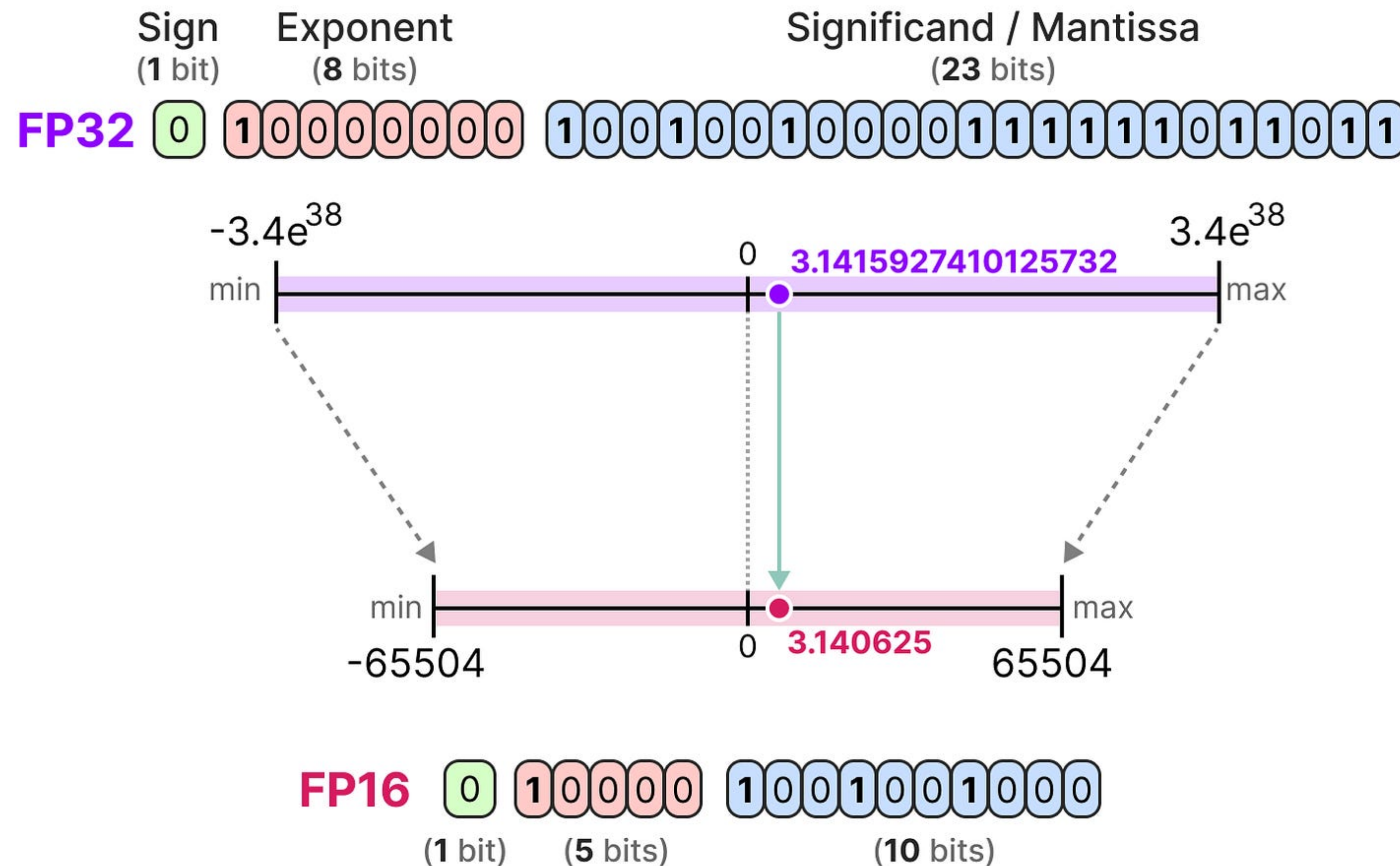What is the difference between BF16 and FP16?

# Common data types

Using more bits **increases the precision** of a value

**Float 32-bit** (FP32)

$0$ $10000000$ $10010010000111111011011$

$(-1)^0 \times 2^1 \times 1.5707964 = 3.1415927410125732$

**higher** precision

**Float 16-bit** (FP16)

$0$ $10000$ $1001001000$

$(-1)^0 \times 2^1 \times 1.5703125 = 3.140625$

**lower** precision

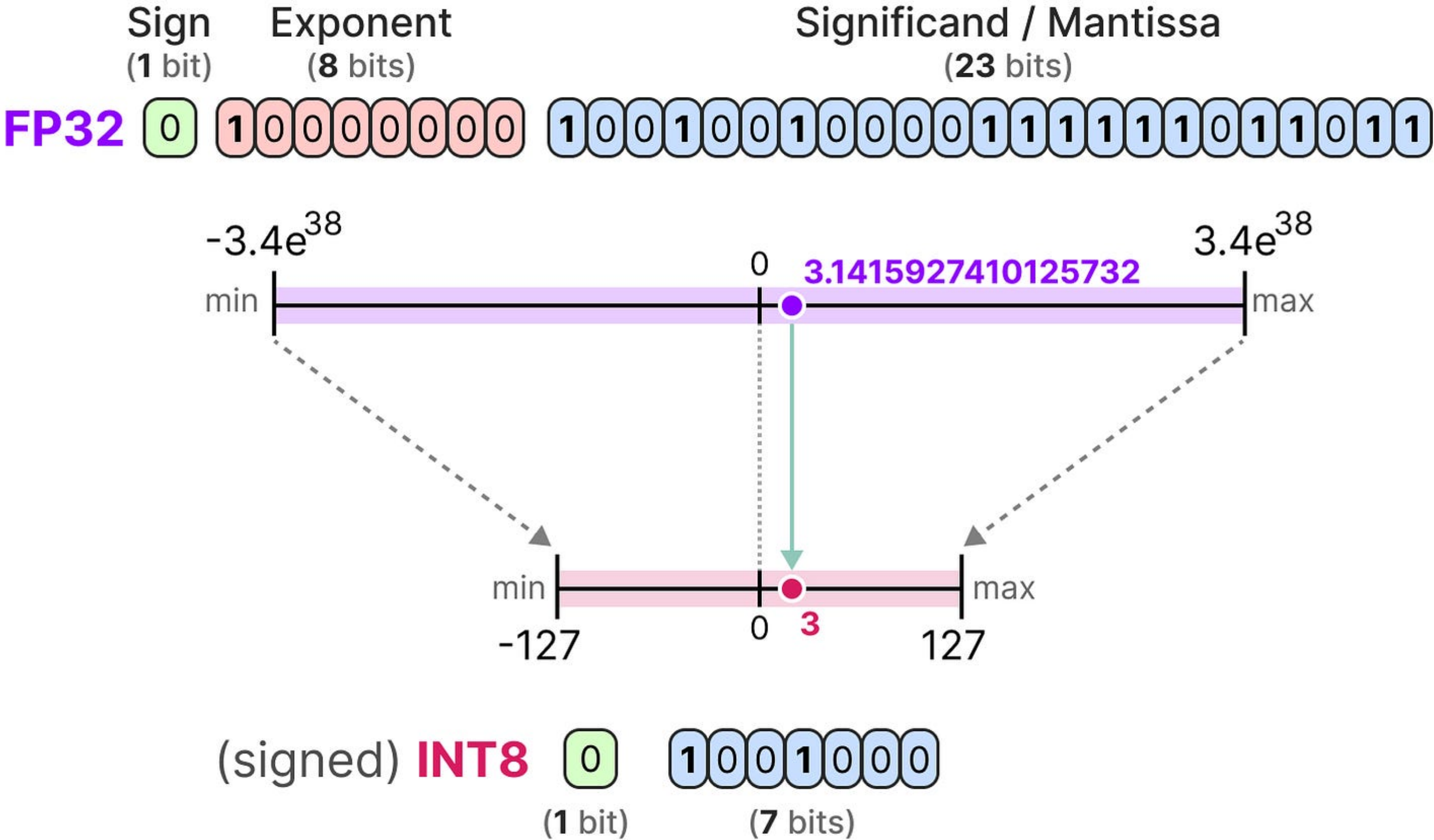original value
**3.1415927**

# Common data types

The range of values **reduce with quantization**

# Common data types

We can convert FP32 to **8-bit integer**-based representations to **save memory**
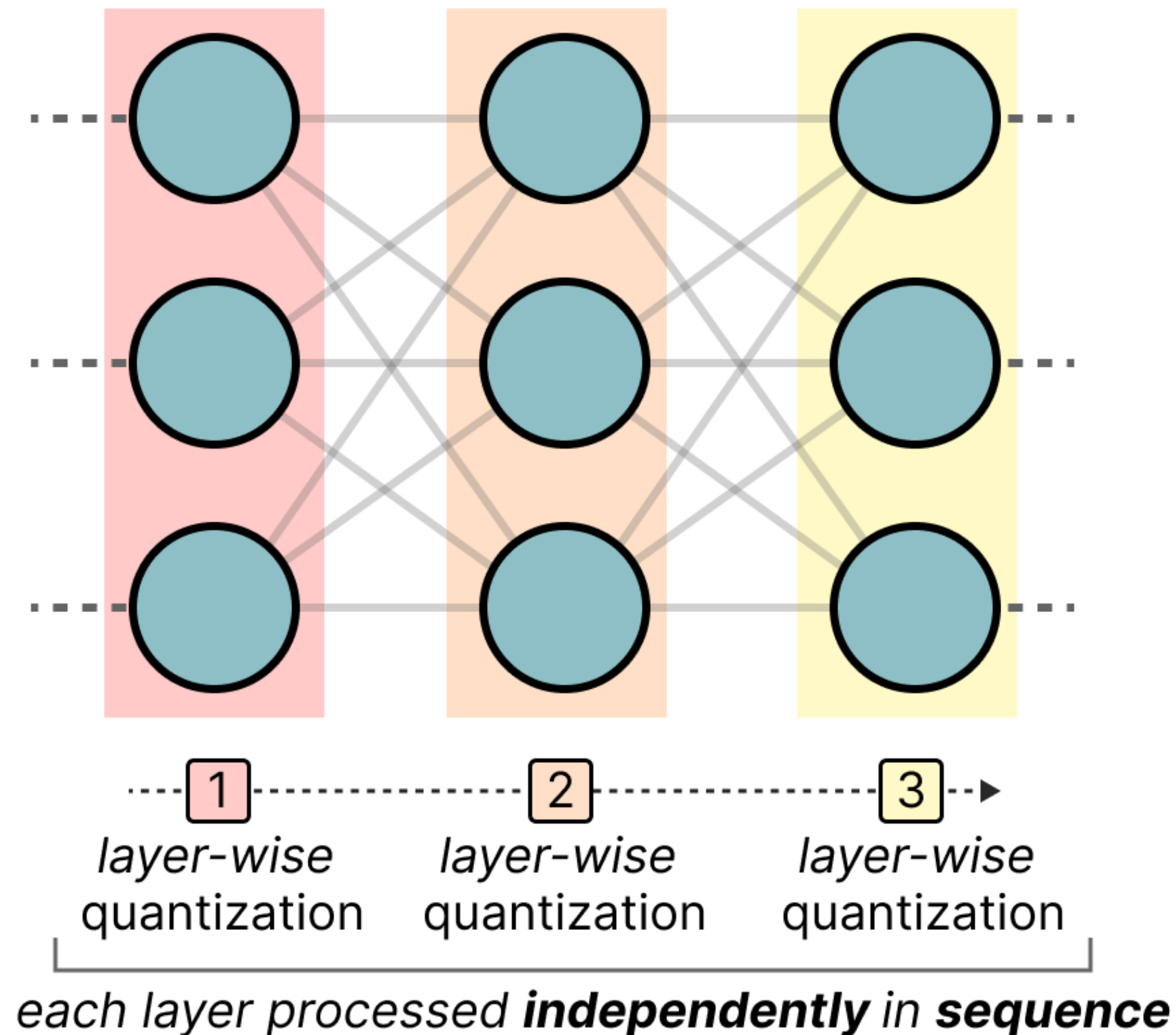
# Absolute maximum (absmax) quantization

Original number is _____ to scale it into the range [-127, 127].

$$X_{quant} = round\left( \quad\quad\quad\quad \right)$$

$$X_{dequant} = \frac{\max |X|}{127} \cdot X_{quant}$$
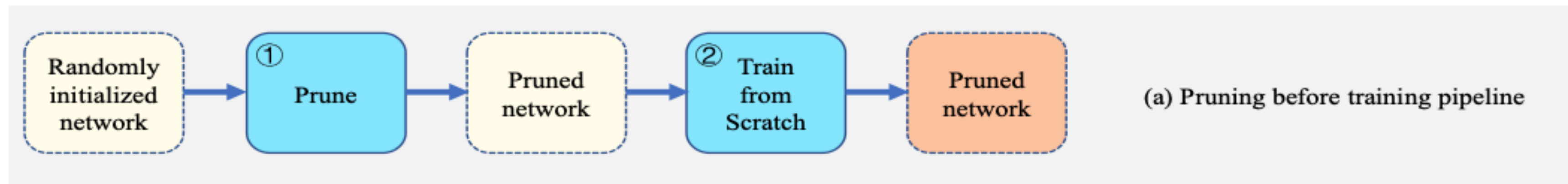
# GPTQ

Each layer is **quantized independently**. Given a layer $l$ with a weight matrix $W_l$, find quantized weights $\widehat{W_l}$



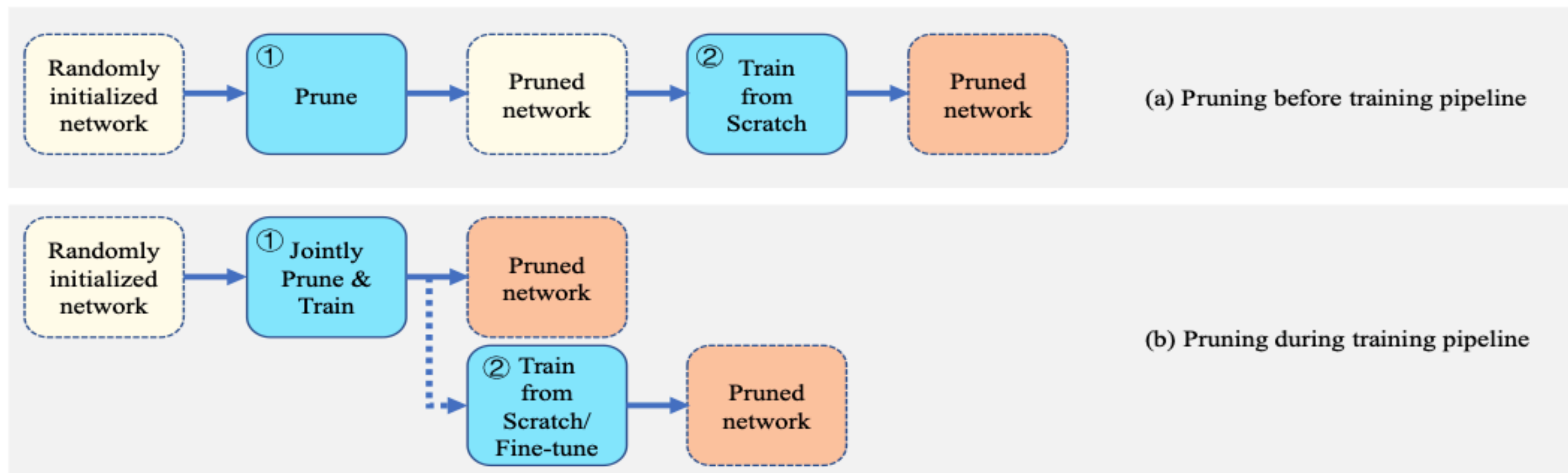$$\widehat{W_l^*} = \arg\min_{W_l} \| W_l X - \widehat{W_l} X \|^2$$



1    2    3

*layer-wise* quantization    *layer-wise* quantization    *layer-wise* quantization

*each layer processed **independently** in **sequence***

# Pruning before training

Procedure: _____



(a) Pruning before training pipeline

$$f(x_t, W_0 \odot M') \rightarrow f(x_t, W_t \odot M')$$
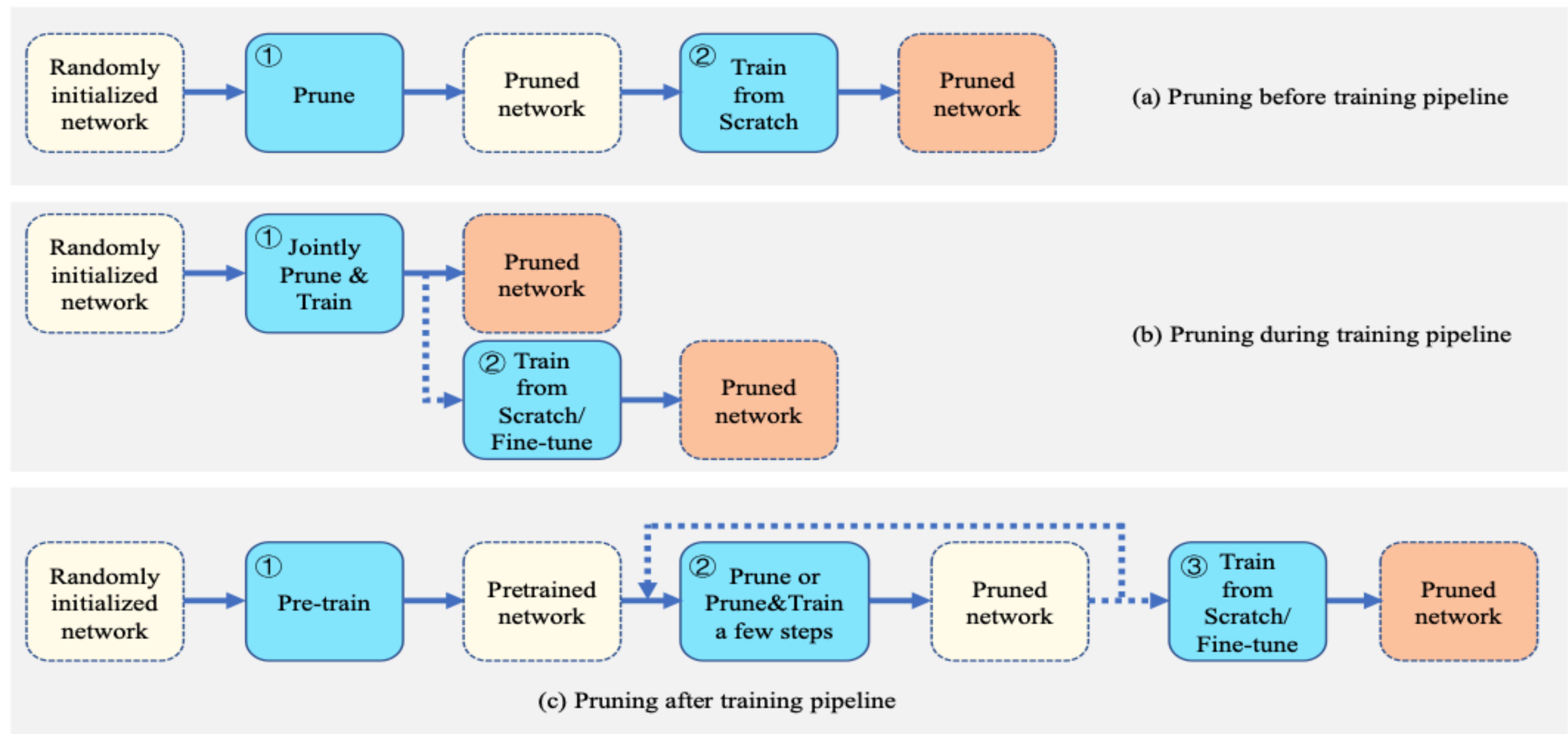
# Pruning during training

Procedure: _____



(a) Pruning before training pipeline

$$f(x_t, W_0 \odot M') \rightarrow f(x_t, W_t \odot M')$$

(b) Pruning during training pipeline

$$f(x_t, W_0) \rightarrow f(x_t, W_t \odot M_t)$$

# Pruning after training

Procedure: _____



$$f(x_t, W_0 \odot M') \rightarrow f(x_t, W_t \odot M')$$

(a) Pruning before training pipeline

$$f(x_t, W_0) \rightarrow f(x_t, W_t \odot M_t)$$

(b) Pruning during training pipeline

$$f(x_t, W_0) \rightarrow f(x_t, W_t)$$
$$\rightarrow f(x_t, W_t' \odot M')$$

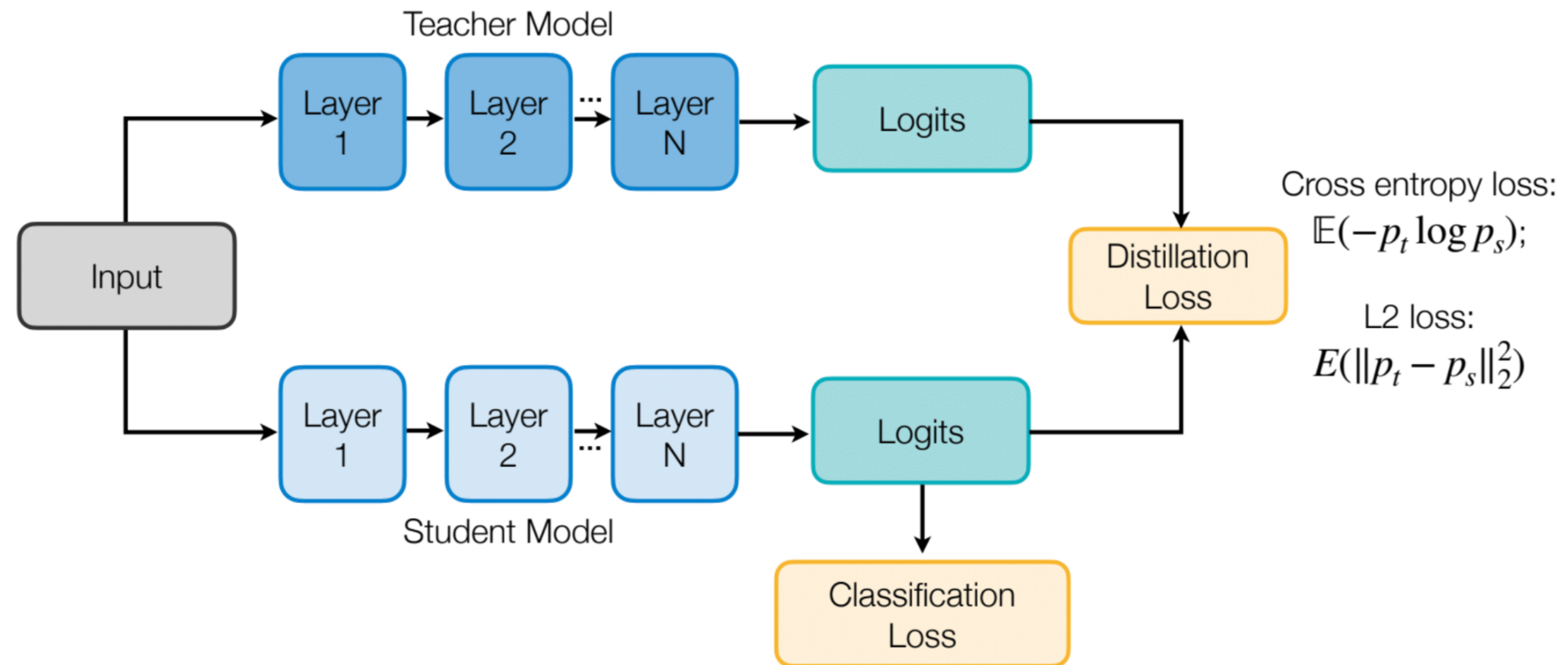(c) Pruning after training pipeline

# Knowledge Distillation

Goal is to **transfer knowledge** from a **larger model** to a **student model**



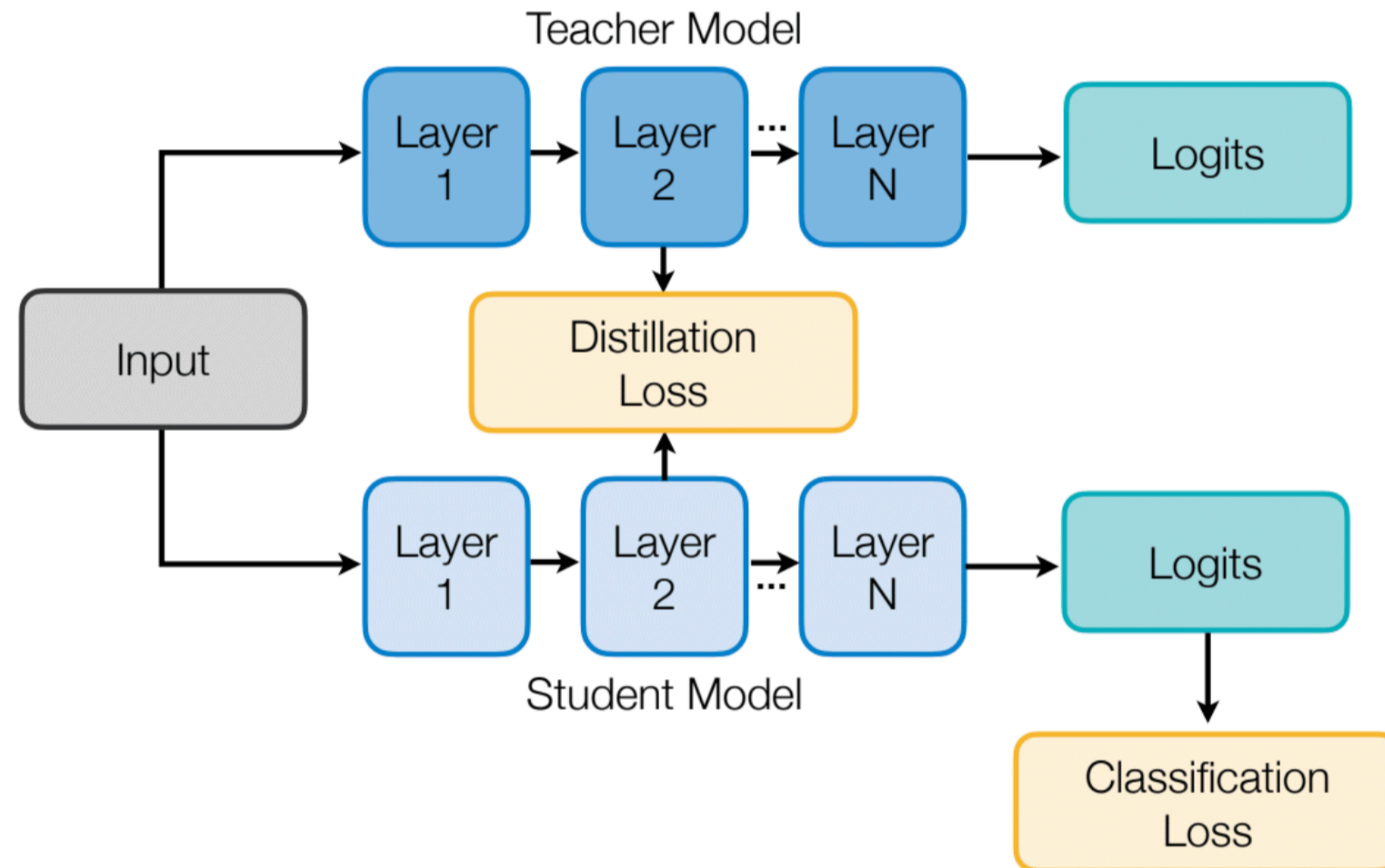Why not train the **student model from scratch**?

# Knowledge Distillation

Simplest way is to **match the outputs** using a distance metric



Cross entropy loss:
$$\mathbb{E}(-p_t \log p_s);$$
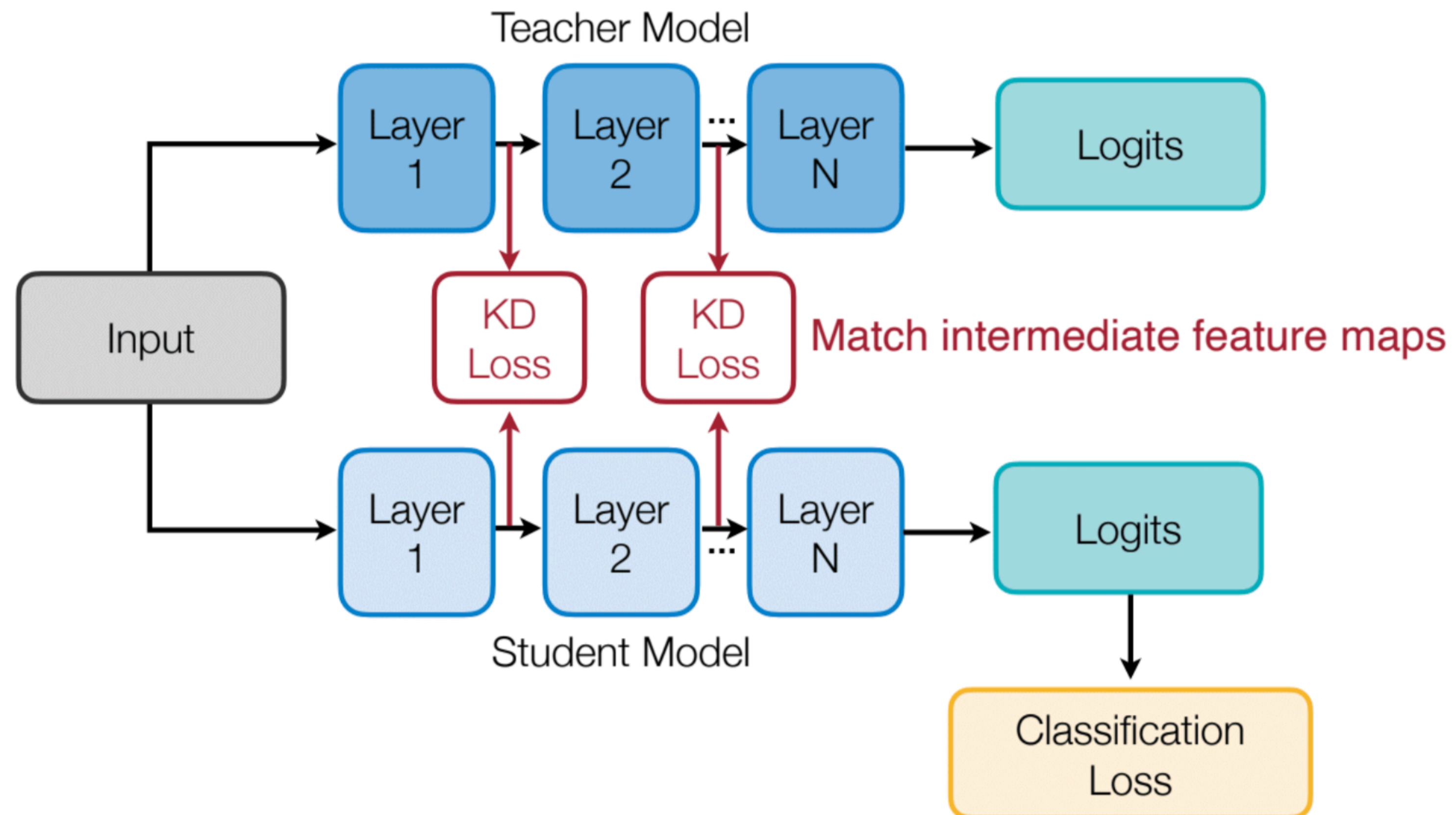
L2 loss:
$$E(\|p_t - p_s\|_2^2)$$

# Knowledge Distillation

Match the **weights in intermediate layers** using a distance metric

# Knowledge Distillation

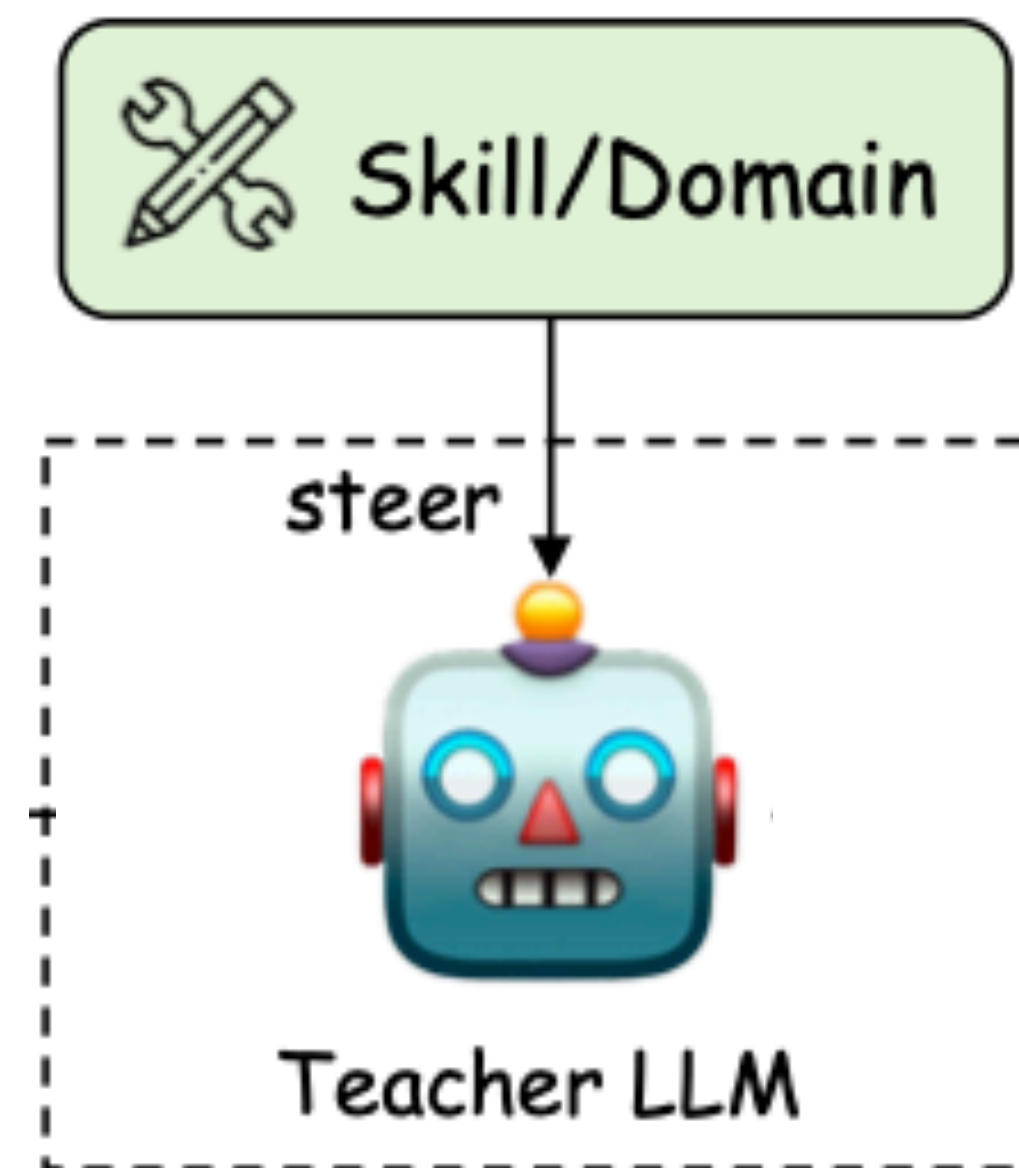Match the **intermediate feature maps** using a distance metric

# Knowledge Distillation in LLM Era

Due to **inaccessible parameters**, we want to **transfer knowledge** from LLMs
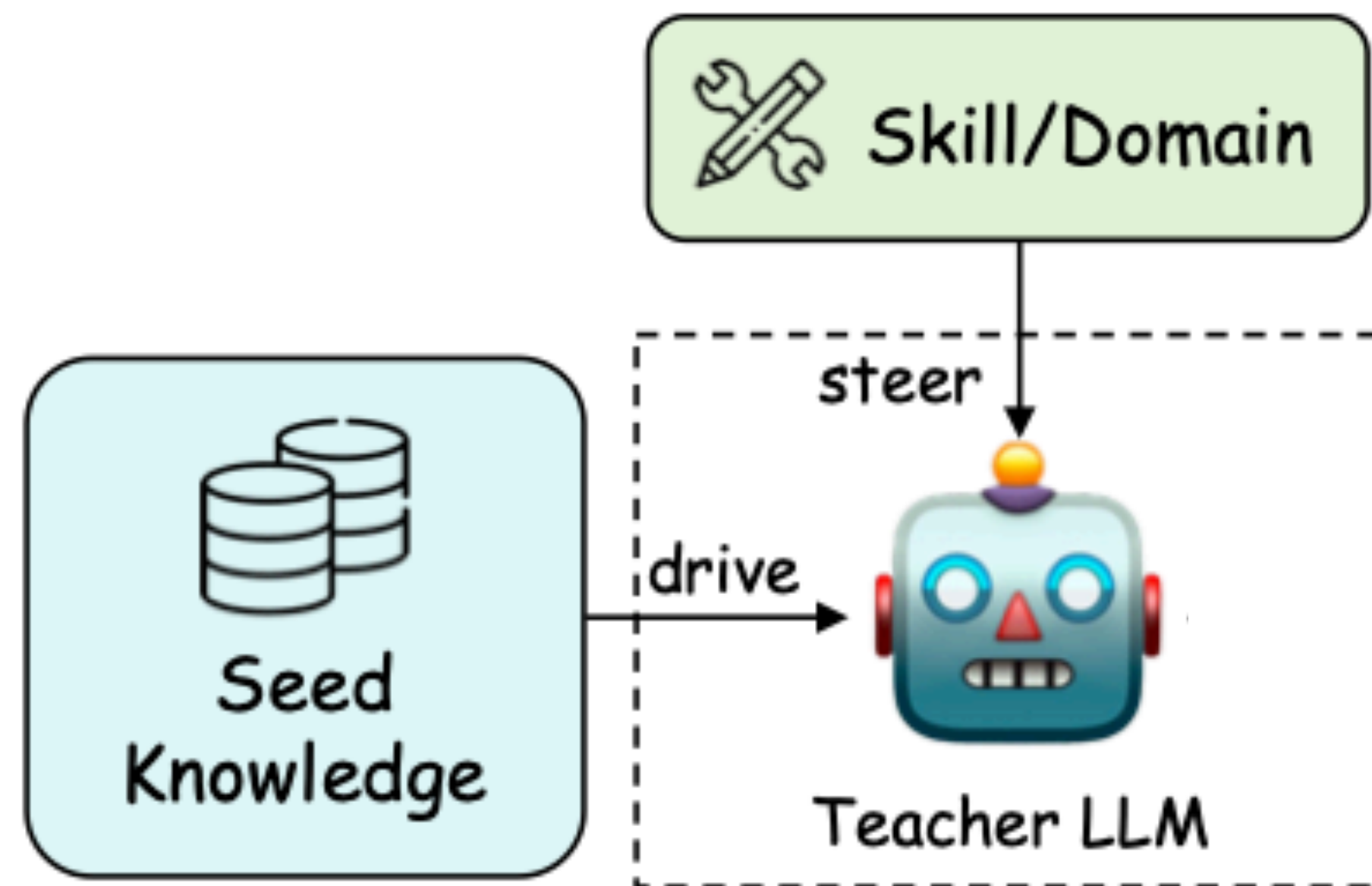

Teacher LLM

# Knowledge Distillation in LLM Era

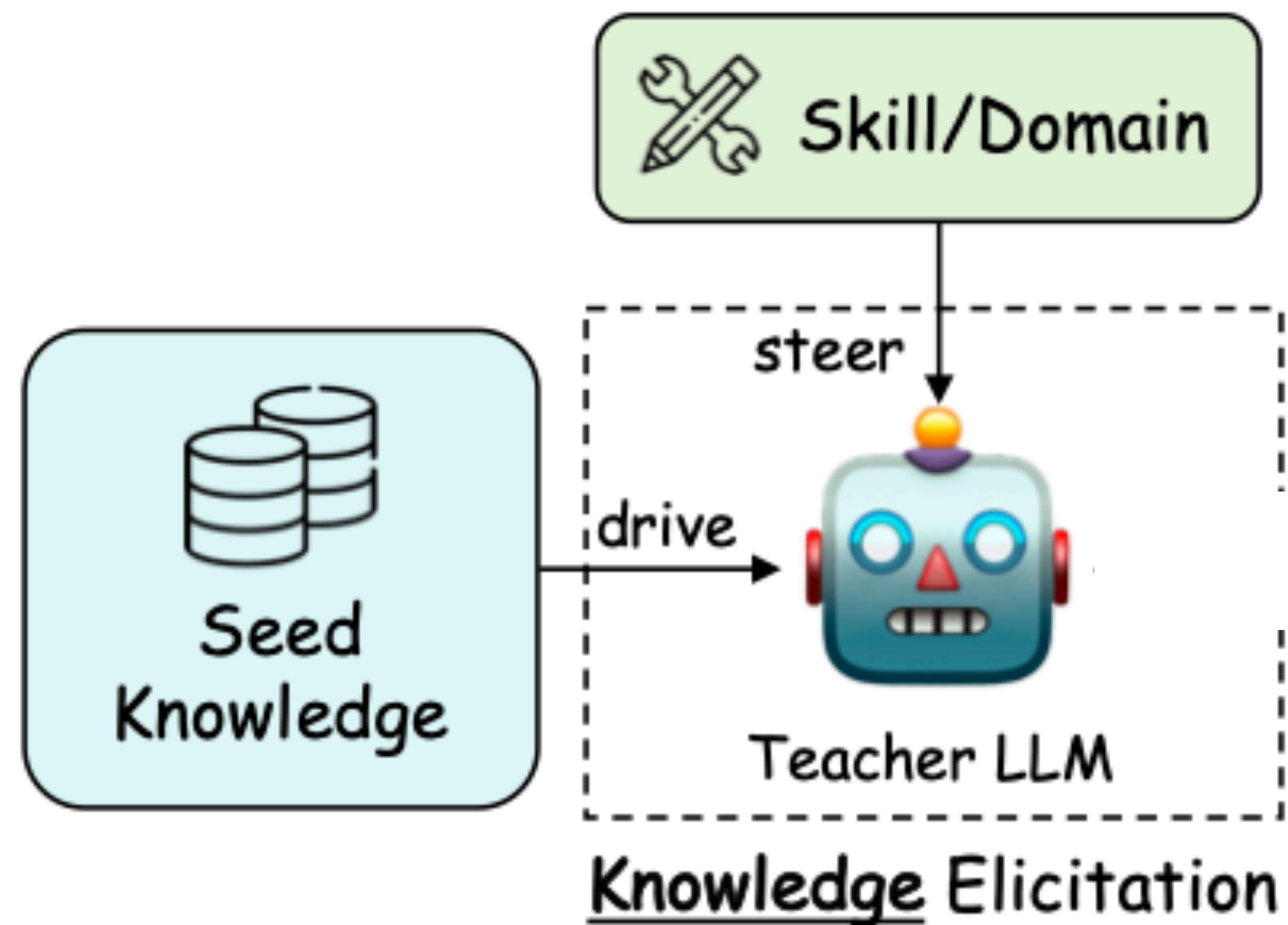**Steer the LLM** for a target skill or domain

# Knowledge Distillation in LLM Era

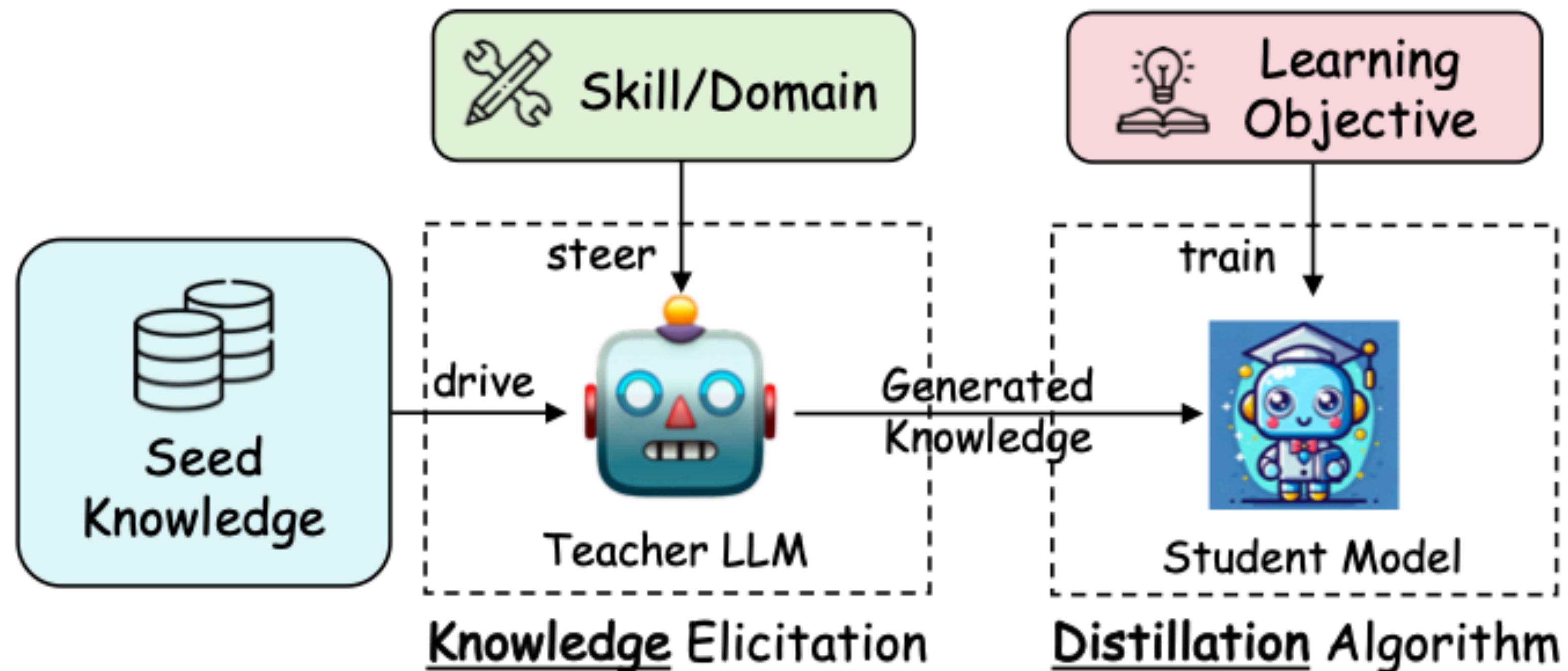Feed the LLM with a small data as **seed knowledge**

# Knowledge Distillation in LLM Era

Feed the LLM with a small data as **seed knowledge**

# Knowledge Distillation in LLM Era

**Generate knowledge** from the teacher LLM and **emulate teacher skills**

# Takeaways