# Parameter-Efficient Finetuning

## CSE 493G/599G Recitation

Prepared by Scott Geng

# What is finetuning?



Take useful model that already knows a lot and update it slightly

Can build applications cheaper, better.



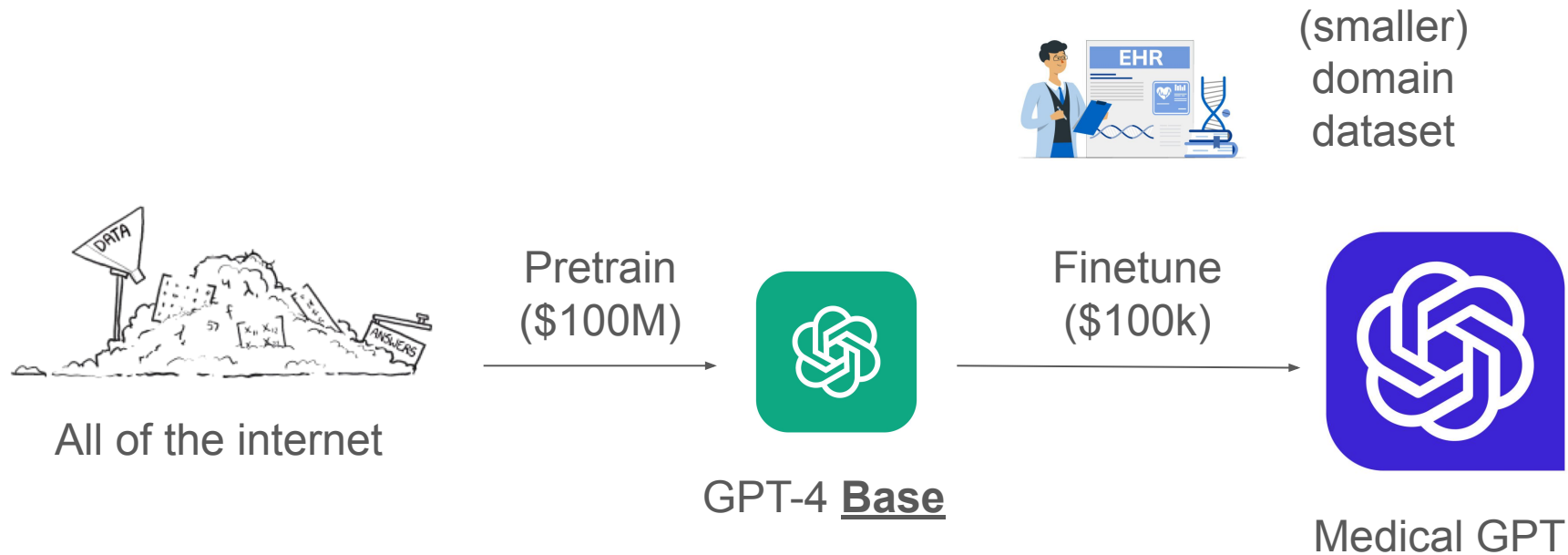Medical GPT

# Can build applications cheaper, better.

(smaller) domain dataset

All of the internet

Medical GPT

# Can build applications cheaper, better.



(smaller) domain dataset

All of the internet → Pretrain ($100M) → GPT-4 **Base** → Finetune ($100k) → Medical GPT

# Can build personalized applications.



Step 1

**Collect demonstration data and train a supervised policy.**

A prompt is sampled from our prompt dataset.
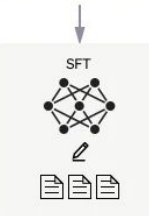
Explain reinforcement learning to a 6 year old.

A labeler demonstrates the desired output behavior.

We give treats and punishments to teach...

This data is used to fine-tune GPT-3.5 with supervised learning.

SFT

OpenAI.
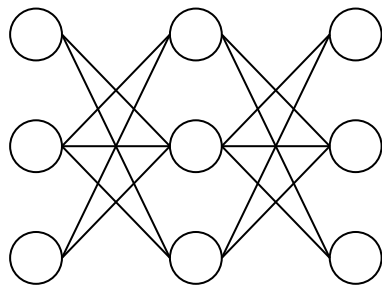
Input images

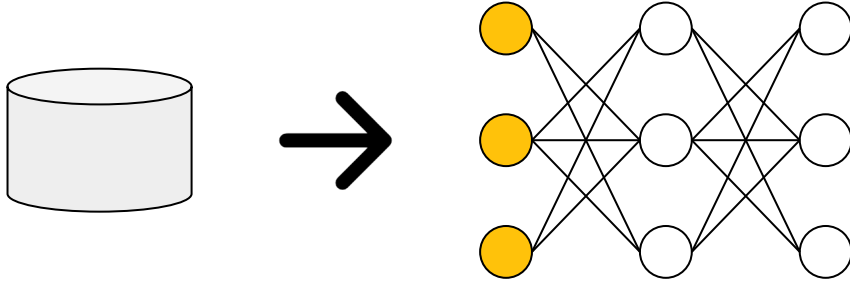w/o prior-preservation loss

Ours (full)

DreamBooth.

**Vision:** everyone should be able to easily adapt a very capable (very big) base model to whatever task they want
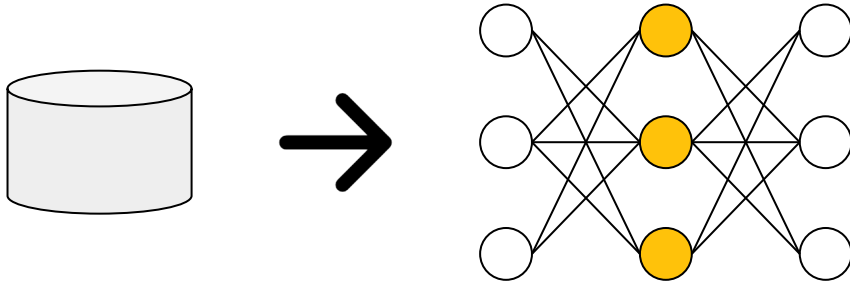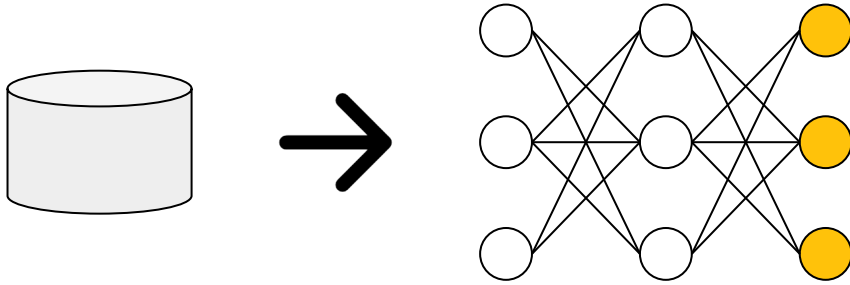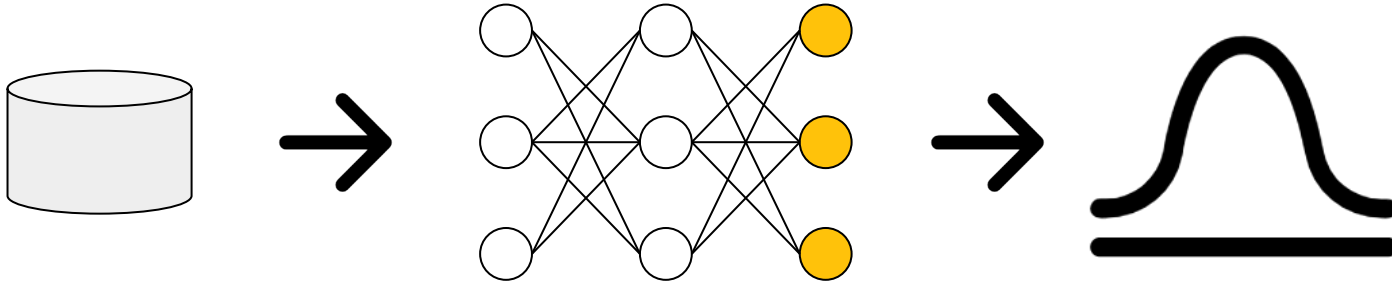
# How to finetune a model

Slide credit to Tim Dettmers

# How to finetune a model

Slide credit to Tim Dettmers

# How to finetune a model

Slide credit to Tim Dettmers

# How to finetune a model

Slide credit to Tim Dettmers

# How to finetune a model

Slide credit to Tim Dettmers

# How to finetune a model



**Error**

# How to finetune a model



**Error**

# How to finetune a model



$$\frac{\partial \mathbf{E}}{\partial \mathbf{W}_2}$$

Weight gradients

**Error**

Slide credit to Tim Dettmers

# How to finetune a model



$$\frac{\partial \mathbf{E}}{\partial \mathbf{W}_1} \qquad \frac{\partial \mathbf{E}}{\partial \mathbf{W}_2}$$

**Error**

Weight gradients

# Background: How to finetune a model

Update the weights



$$\frac{\partial \mathbf{E}}{\partial \mathbf{W}_1} \qquad \frac{\partial \mathbf{E}}{\partial \mathbf{W}_2}$$
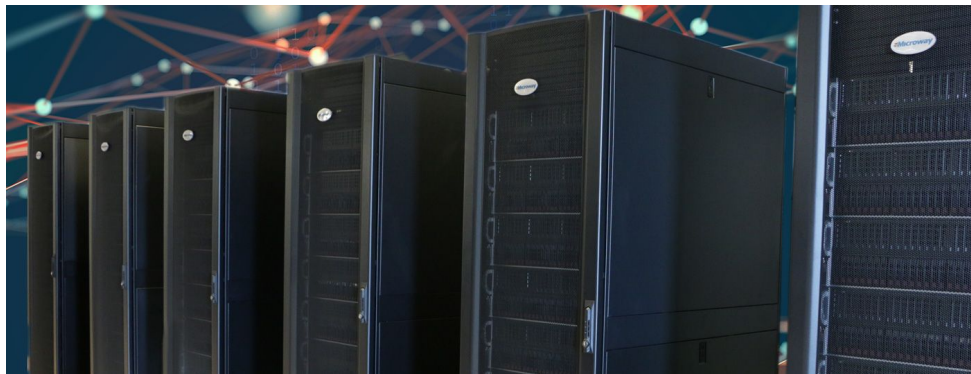
**Quality** ⬆

Slide credit to Tim Dettmers
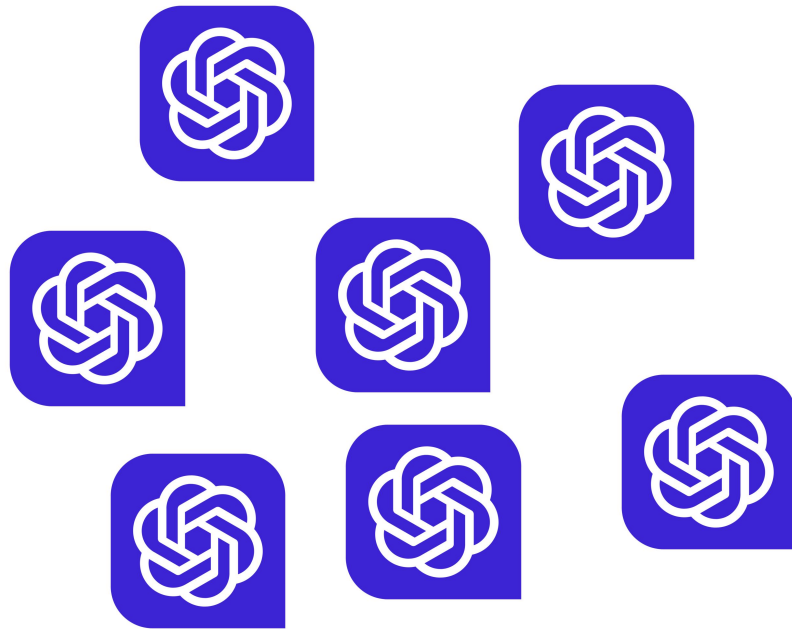
# Challenge: this is expensive compute wise.



100B parameters
base model
(~64-128 GPUs to train)

# Challenge: this is expensive storage wise.



100B parameters
base model
(~200GB)

Each finetuned copy is same size!

**Research problem:** how can we reduce the cost of (1) finetuning a model and (2) storing the updated copy?
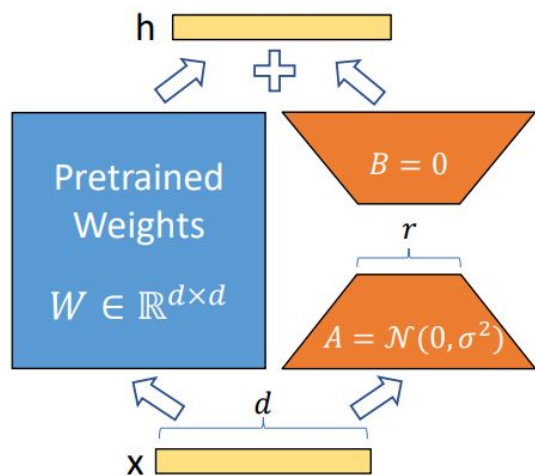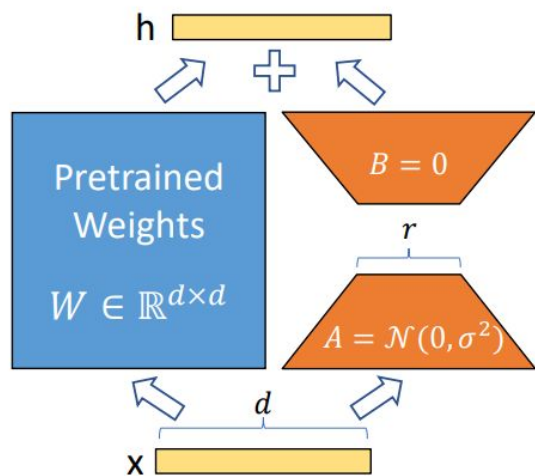
# **Lo**w **R**ank **A**daptation (LoRA)



Figure 1: Our reparametrization. We only train $A$ and $B$.

$$W_{\text{finetuned}} = W_{base} + \Delta W$$

$$h = W_{\text{finetuned}}(x) = W_{base}(x) + \Delta W(x)$$

# Low Rank Adaptation (LoRA)



Figure 1: Our reparametrization. We only train $A$ and $B$.
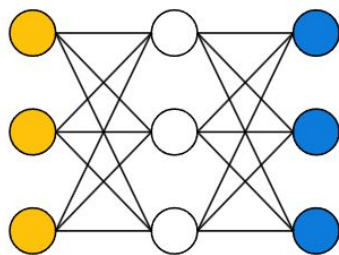
$$W_{\text{finetuned}} = W_{base} + \Delta W$$

$$h = W_{\text{finetuned}}(x) = W_{base}(x) + \Delta W(x)$$

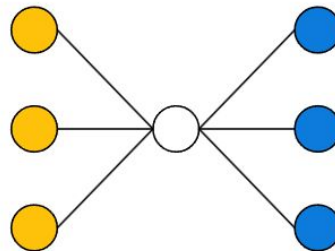**Key observation: deltaW has low rank, so that we can express it as a product of two simpler matrices**

$$\begin{bmatrix} 2 & 3 & -1 \\ 0 & 1 & 4 \\ 0 & 0 & 0 \end{bmatrix}$$ Find the Rank!

rank(A) = 2

Approximation through low-rank projection

$$M \approx L_k \times R_k^T$$

$m \times n$   $m \times k$   $k \times n$

Approximation
through
low-rank projection

# Low Rank Adaptation (LoRA)
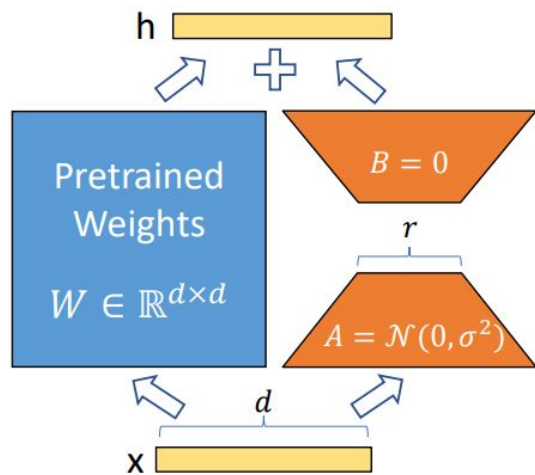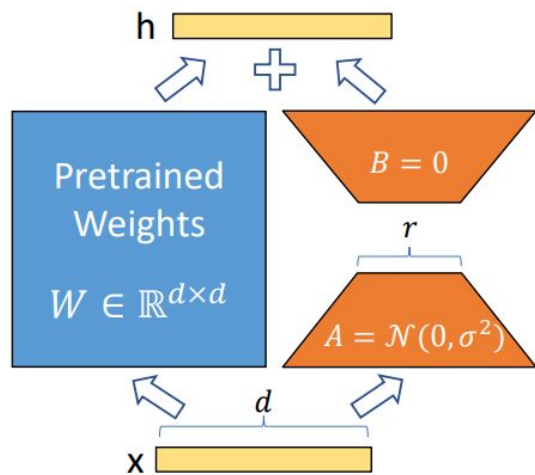


Figure 1: Our reparametrization. We only train $A$ and $B$.

$$W_{\text{finetuned}} = W_{base} + \Delta W$$

$$h = W_{\text{finetuned}}(x) = W_{base}(x) + \Delta W(x)$$

$$\Delta W = BA$$

$$h = W_{\text{finetuned}}(x) = W_{base}(x) + BAx$$

# Low Rank Adaptation (LoRA)



Figure 1: Our reparametrization. We only train $A$ and $B$.

$$W_{\text{finetuned}} = W_{base} + \Delta W$$

$$h = W_{\text{finetuned}}(x) = W_{base}(x) + \Delta W(x)$$

$$\Delta W = BA$$

$$h = W_{\text{finetuned}}(x) = W_{base}(x) + BAx$$
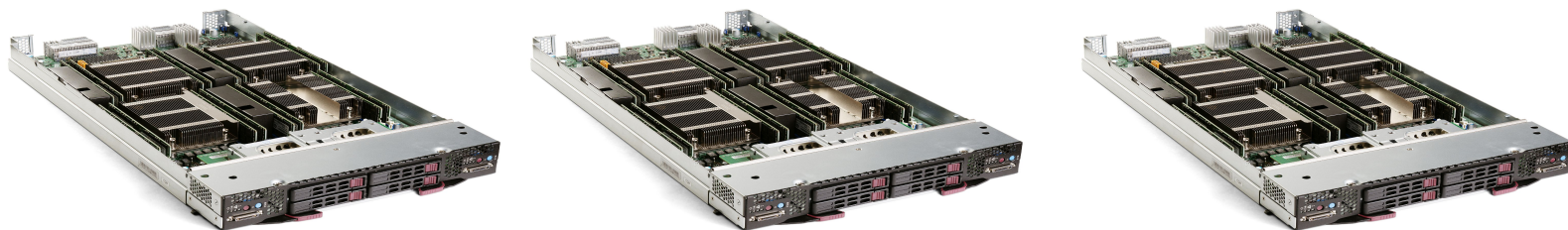
**Now, only need to train and store B, A**

Rank = 8

Rank = 128

# Finetuning a ~11B+ parameter model still requires multiple servers.

Dettmers et al., 2023

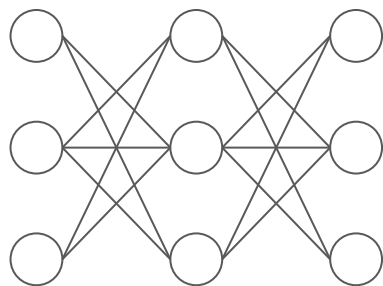# QLoRA: Finetuning large models on a single GPU.

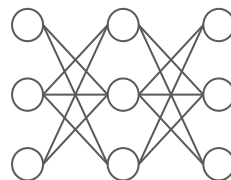

**QLoRA**

(4-bit finetuning)

Slide credit to Tim Dettmers

Dettmers et al., 2023
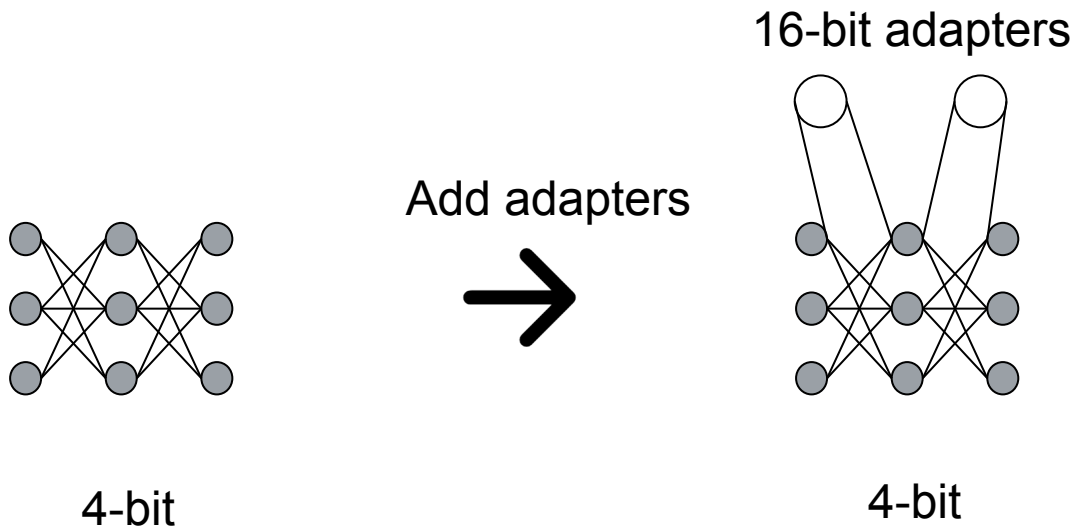
# Quantized Low-rank Adaptation (QLoRA)



16-bit $\rightarrow$ 4-bit

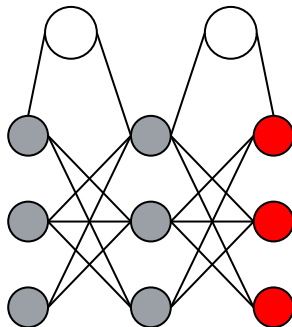# Quantized Low-rank Adaptation (QLoRA)

16-bit adapters

Add adapters

$\rightarrow$

4-bit

4-bit

Slide credit to Tim Dettmers

# Quantized Low-rank Adaptation (QLoRA)
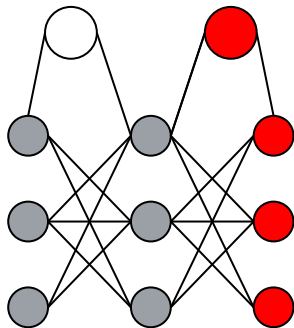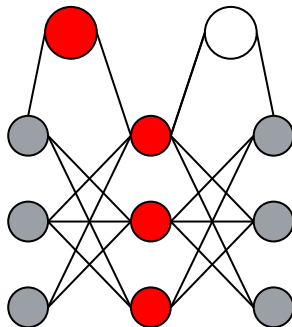
16-bit adapters

**4-bit Error**

4-bit model

# Quantized Low-rank Adaptation (QLoRA)

16-bit adapters



**4-bit Error**
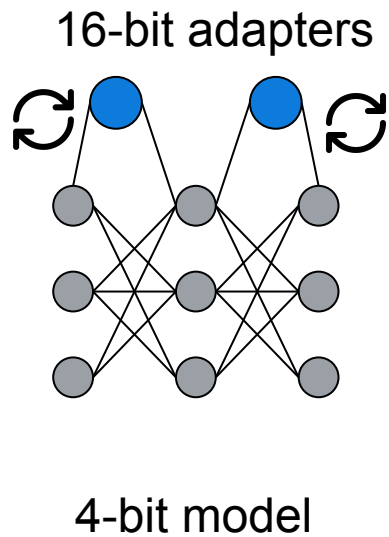
4-bit model

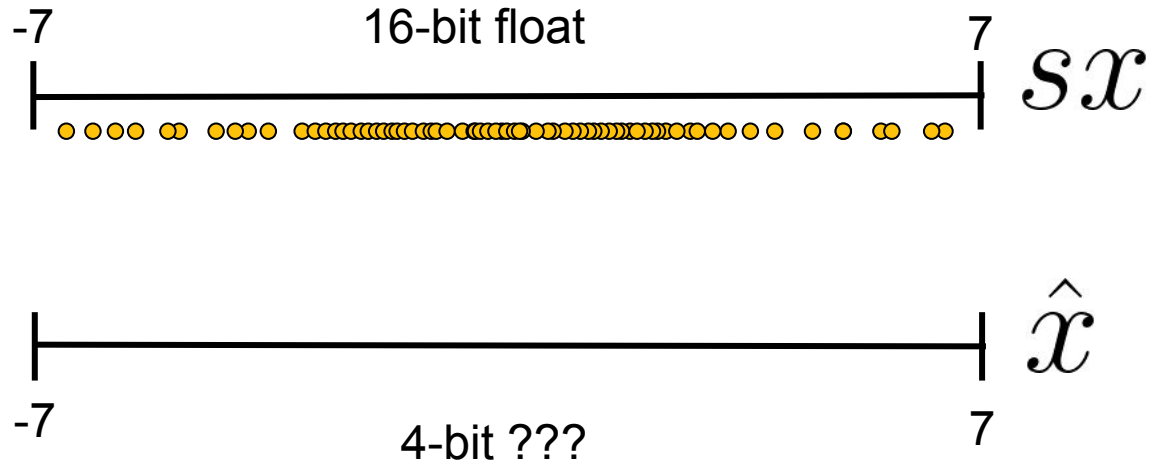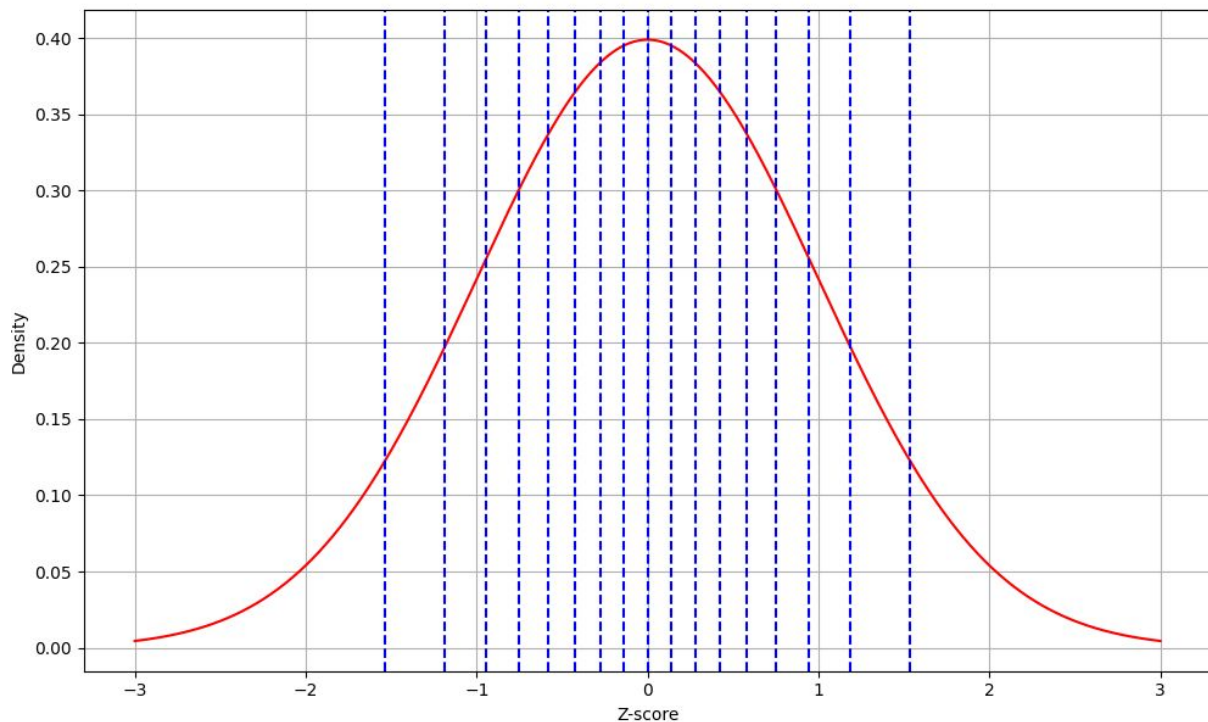# Quantized Low-rank Adaptation (QLoRA)

16-bit adapters



4-bit model

**4-bit Error**

# Quantized Low-rank Adaptation (QLoRA)

16-bit adapters

4-bit model

Slide credit to Tim Dettmers

# What 4-bit data type is information theoretically optimal?

-7              16-bit float             7 $sx$

$\hat{x}$

-7                  7

4-bit ???

# 4-bit NormalFloat (NF4) an information-theoretically optimal data type for normal distributions
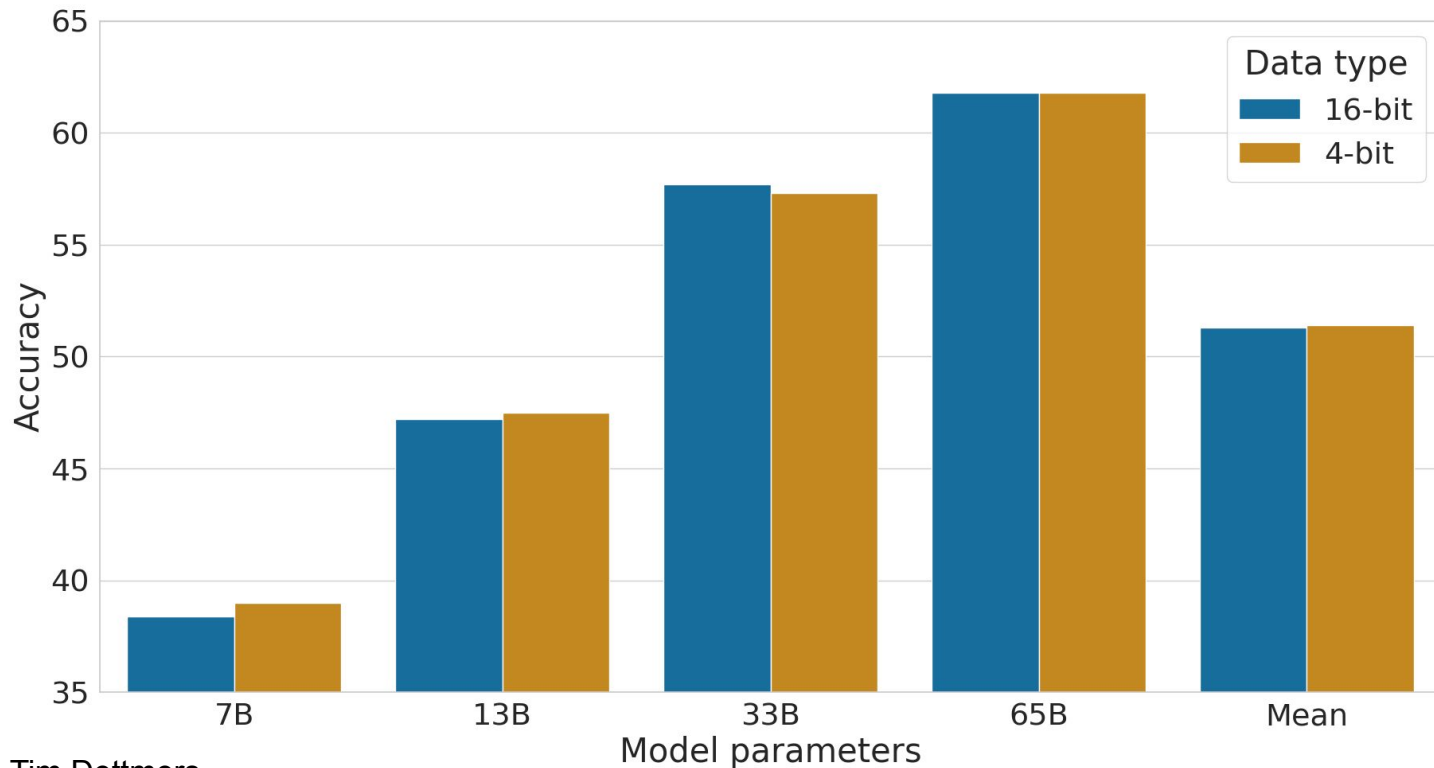
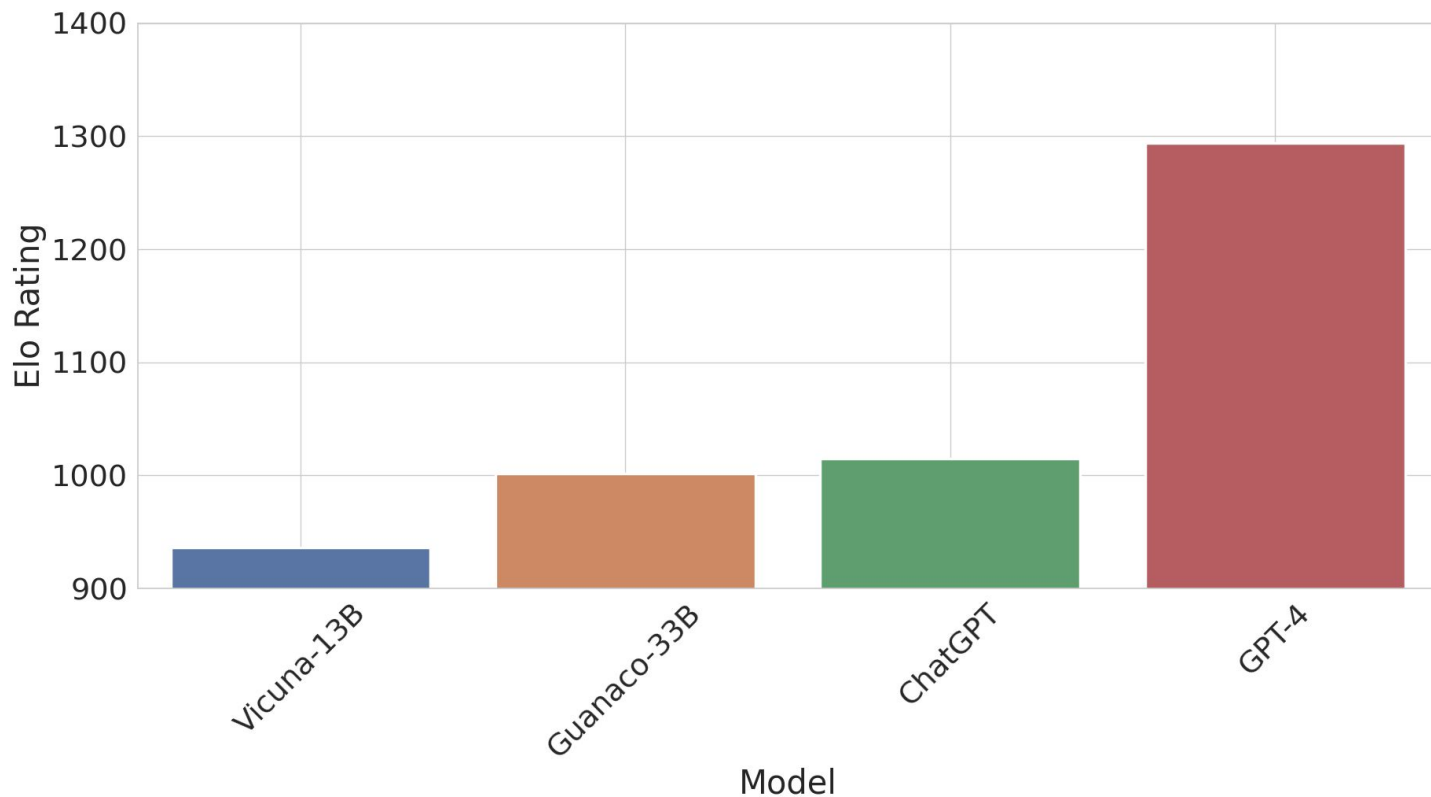# QLoRA systems contributions

- Double quantization
- GPU memory paging for optimizer

# Results

# QLoRA recovers lost performance through fine-tuning

# 4-bit Guanaco: A ChatGPT-quality 4-bit chatbot finetuned in 24h on a single GPU

# Take-away

4-bit finetuning is possible by passing gradients through a 4-bit neural network to 16-bit adapters.