

# Efficient training & inference

CS 4804: Introduction to AI

*Fall 2025*

<https://tuvllms.github.io/ai-fall-2025/>

Tu Vu



# Logistics

- HW 2 released due 11/18
- Final presentations: 12/4 & 12/9

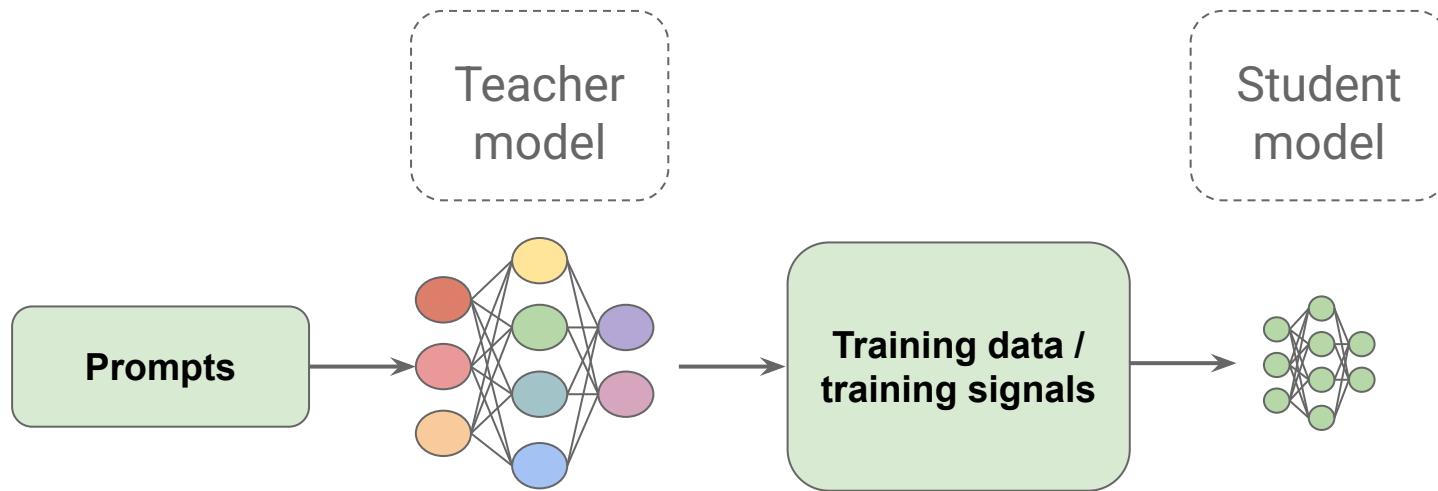
# Today's lecture

- Knowledge distillation
- Model merging / Model editing
- Multi-query attention / Grouped-query attention
- KV caching
- Steering vectors
- Pruning
- Quantization

# Future lectures

- Agents & Compound AI systems (2 lectures)
  - Retrieval-Augmented Generation (RAG)
  - Tool-use models
  - Modular AI software with DSPy
  - Mixture of agents
  - Agentic memory
  - Model Context Protocol (MCP) & Agentic AI workflows
- Diffusion models
- Ethics and safety

# Knowledge distillation



---

# **Distilling the Knowledge in a Neural Network**

---

**Geoffrey Hinton<sup>\*†</sup>**

Google Inc.

Mountain View

geoffhinton@google.com

**Oriol Vinyals<sup>†</sup>**

Google Inc.

Mountain View

vinyals@google.com

**Jeff Dean**

Google Inc.

Mountain View

jeff@google.com

---

# **DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter**

---

**Victor SANH, Lysandre DEBUT, Julien CHAUMOND, Thomas WOLF**  
Hugging Face  
{victor,lysandre,julien,thomas}@huggingface.co

$$\text{Loss} = \lambda_{\text{ce}} \cdot \mathcal{L}_{\text{ce}} + \lambda_{\text{kd}} \cdot \mathcal{L}_{\text{kd}}$$

$$\text{Loss} = \lambda_{\text{ce}} \cdot \left( - \sum_{i=1}^N y_i \log(p_i) \right) + \lambda_{\text{kd}} \cdot D_{\text{KL}}(q_{\text{teacher}}(x) \| q_{\text{student}}(x))$$

Where:

- $y_i$  is the true label for token  $i$ ,
- $p_i$  is the predicted probability for the correct token for token  $i$ ,
- $N$  is the number of tokens,
- $D_{\text{KL}}(q_{\text{teacher}}(x) \| q_{\text{student}}(x))$  is the Kullback-Leibler divergence between the teacher and student models' probability distributions,
- $q_{\text{teacher}}(x)$  and  $q_{\text{student}}(x)$  are the output probability distributions from the teacher and student models, respectively,
- $\lambda_{\text{ce}}$  and  $\lambda_{\text{kd}}$  are the weighting hyperparameters for the cross-entropy and knowledge distillation losses, respectively.

Assume two different distributions for predicting the next word:

- $P$  (from Model 1):
  - $mat \rightarrow 0.7$
  - $floor \rightarrow 0.2$
  - $chair \rightarrow 0.1$
- $Q$  (from Model 2):
  - $mat \rightarrow 0.5$
  - $floor \rightarrow 0.3$
  - $chair \rightarrow 0.2$

## Kullback–Leibler (KL) Divergence Calculation

KL divergence measures how much  $P$  diverges from  $Q$ :

$$D_{KL}(P||Q) = \sum_i P(i) \log \frac{P(i)}{Q(i)}$$

Substituting the values:

$$D_{KL}(P||Q) = 0.7 \log \frac{0.7}{0.5} + 0.2 \log \frac{0.2}{0.3} + 0.1 \log \frac{0.1}{0.2}$$

**Training loss** The student is trained with a distillation loss over the soft target probabilities of the teacher:  $L_{ce} = \sum_i t_i * \log(s_i)$  where  $t_i$  (resp.  $s_i$ ) is a probability estimated by the teacher (resp. the student). This objective results in a rich training signal by leveraging the full teacher distribution. Following Hinton et al. [2015] we used a *softmax-temperature*:  $p_i = \frac{\exp(z_i/T)}{\sum_j \exp(z_j/T)}$  where  $T$  controls the smoothness of the output distribution and  $z_i$  is the model score for the class  $i$ . The same temperature  $T$  is applied to the student and the teacher at training time, while at inference,  $T$  is set to 1 to recover a standard *softmax*.

The final training objective is a linear combination of the distillation loss  $L_{ce}$  with the supervised training loss, in our case the *masked language modeling* loss  $L_{mlm}$  [Devlin et al., 2018]. We found it beneficial to add a *cosine embedding* loss ( $L_{cos}$ ) which will tend to align the directions of the student and teacher hidden states vectors.

# DistilBERT reduces BERT's size by 40%, while retaining 97% of its performance and being 60% faster

Model	Score	CoLA	MNLI	MRPC	QNLI	QQP	RTE	SST-2	STS-B	WNLI
ELMo	68.7	44.1	68.6	76.6	71.1	86.2	53.4	91.5	70.4	56.3
BERT-base	79.5	56.3	86.7	88.6	91.8	89.6	69.3	92.7	89.0	53.5
DistilBERT	77.0	51.3	82.2	87.5	89.2	88.5	59.9	91.3	86.9	56.3

Model	IMDb (acc.)	SQuAD (EM/F1)
BERT-base	93.46	81.2/88.5
DistilBERT	92.82	77.7/85.8
DistilBERT (D)	-	79.1/86.9

Model	# param. (Millions)	Inf. time (seconds)
ELMo	180	895
BERT-base	110	668
DistilBERT	66	410

# **Distilling Step-by-Step! Outperforming Larger Language Models with Less Training Data and Smaller Model Sizes**

**Cheng-Yu Hsieh<sup>1\*</sup>, Chun-Liang Li<sup>2</sup>, Chih-Kuan Yeh<sup>3</sup>, Hootan Nakhost<sup>2</sup>,  
Yasuhiba Fujii<sup>3</sup>, Alexander Ratner<sup>1</sup>, Ranjay Krishna<sup>1</sup>, Chen-Yu Lee<sup>2</sup>, Tomas Pfister<sup>2</sup>**

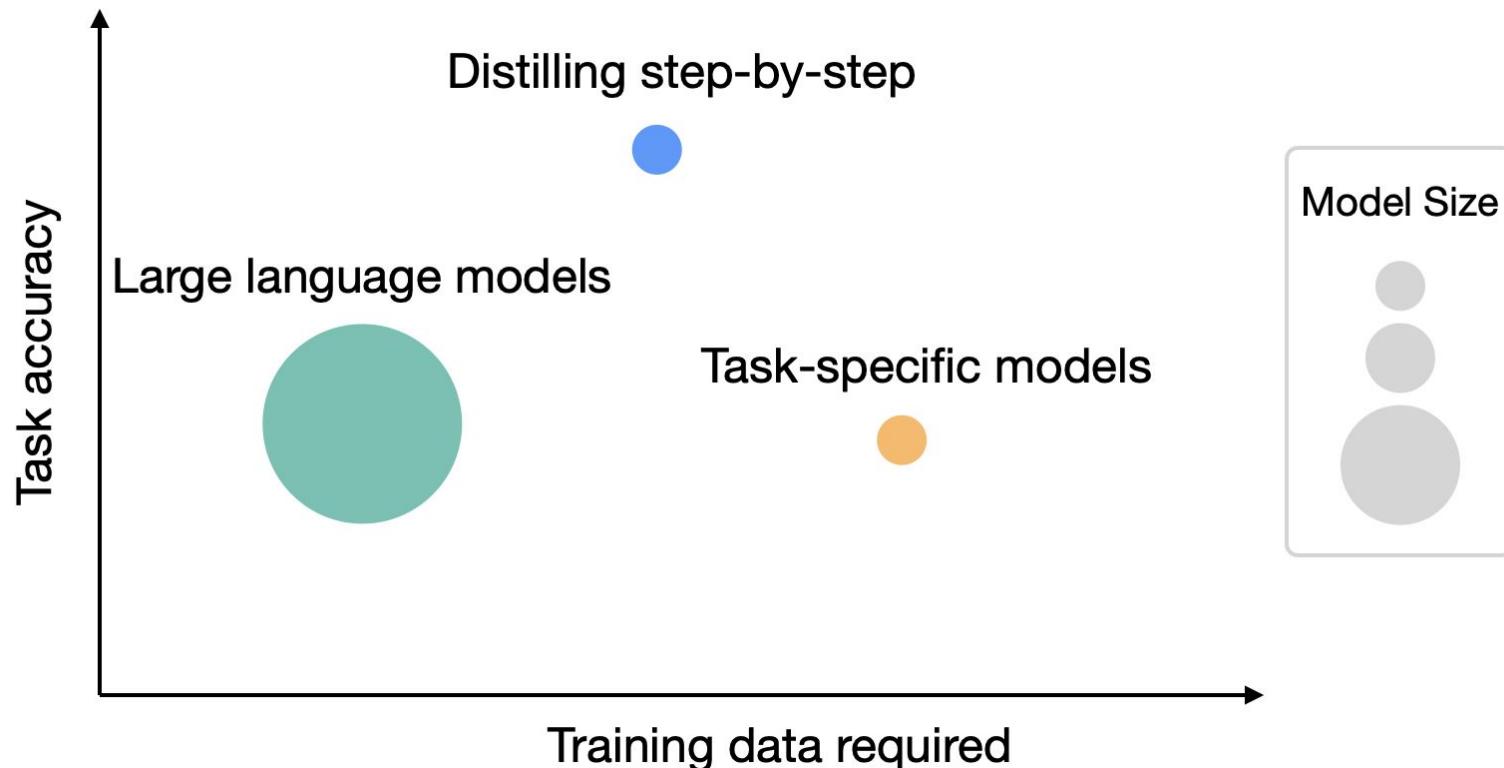
<sup>1</sup>University of Washington, <sup>2</sup>Google Cloud AI Research, <sup>3</sup>Google Research  
cydhsieh@cs.washington.edu

# **Distilling Step-by-Step! Outperforming Larger Language Models with Less Training Data and Smaller Model Sizes**

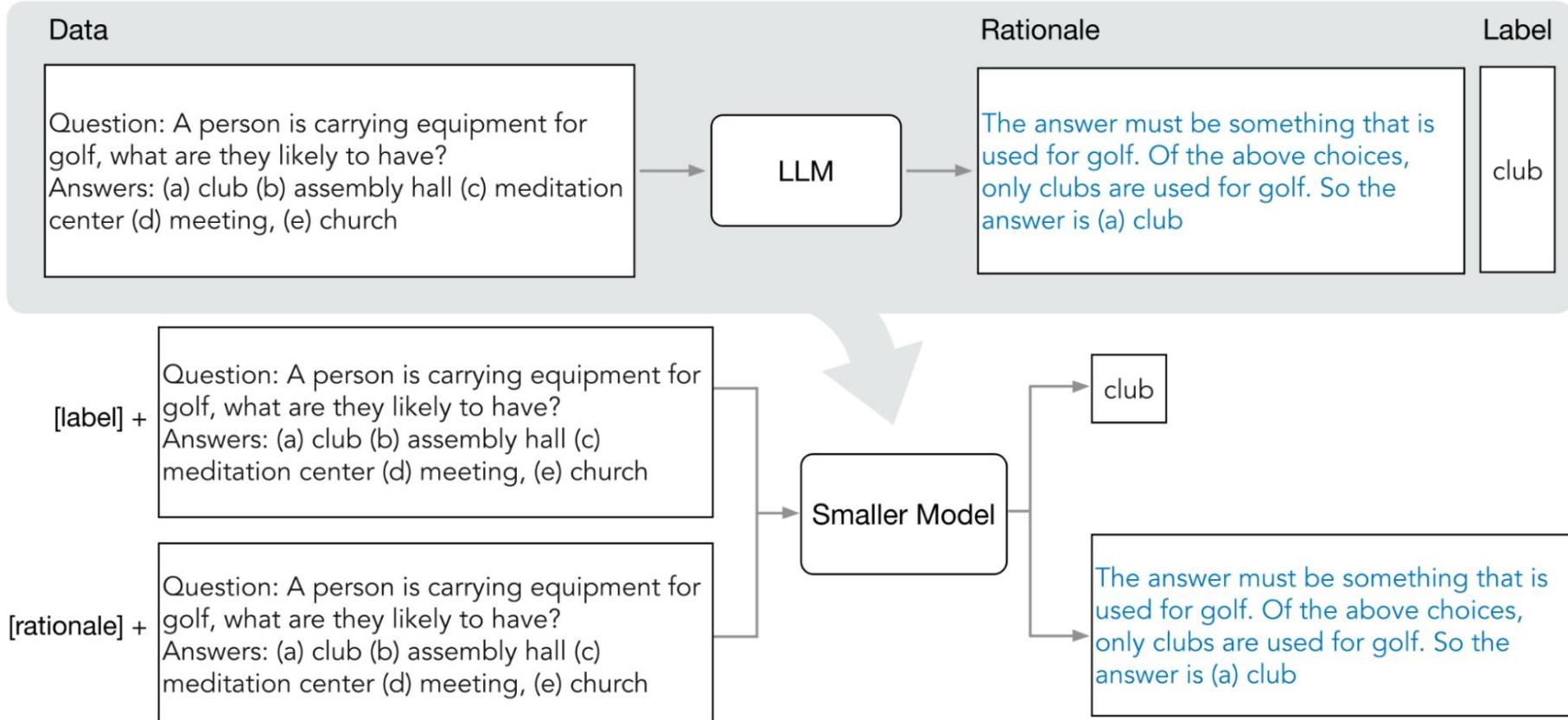
**Cheng-Yu Hsieh<sup>1\*</sup>, Chun-Liang Li<sup>2</sup>, Chih-Kuan Yeh<sup>3</sup>, Hootan Nakhost<sup>2</sup>,  
Yasuhiba Fujii<sup>3</sup>, Alexander Ratner<sup>1</sup>, Ranjay Krishna<sup>1</sup>, Chen-Yu Lee<sup>2</sup>, Tomas Pfister<sup>2</sup>**

<sup>1</sup>University of Washington, <sup>2</sup>Google Cloud AI Research, <sup>3</sup>Google Research  
cydhsieh@cs.washington.edu

# Enabling a 770M parameter T5 model to outperform the few-shot prompted 540B PaLM model



# Distilling step-by-step



# Leveraging few-shot CoT prompting to extract rationales from LLMs

Few-shot CoT

Question: Sammy wanted to go to where the people are. Where might he go?  
Answer Choices: (a) populated areas, (b) race track, (c) desert, (d) apartment, (e) roadblock  
Answer: The answer must be a place with a lot of people. Of the above choices, only populated areas have a lot of people. So the answer is (a) populated areas.

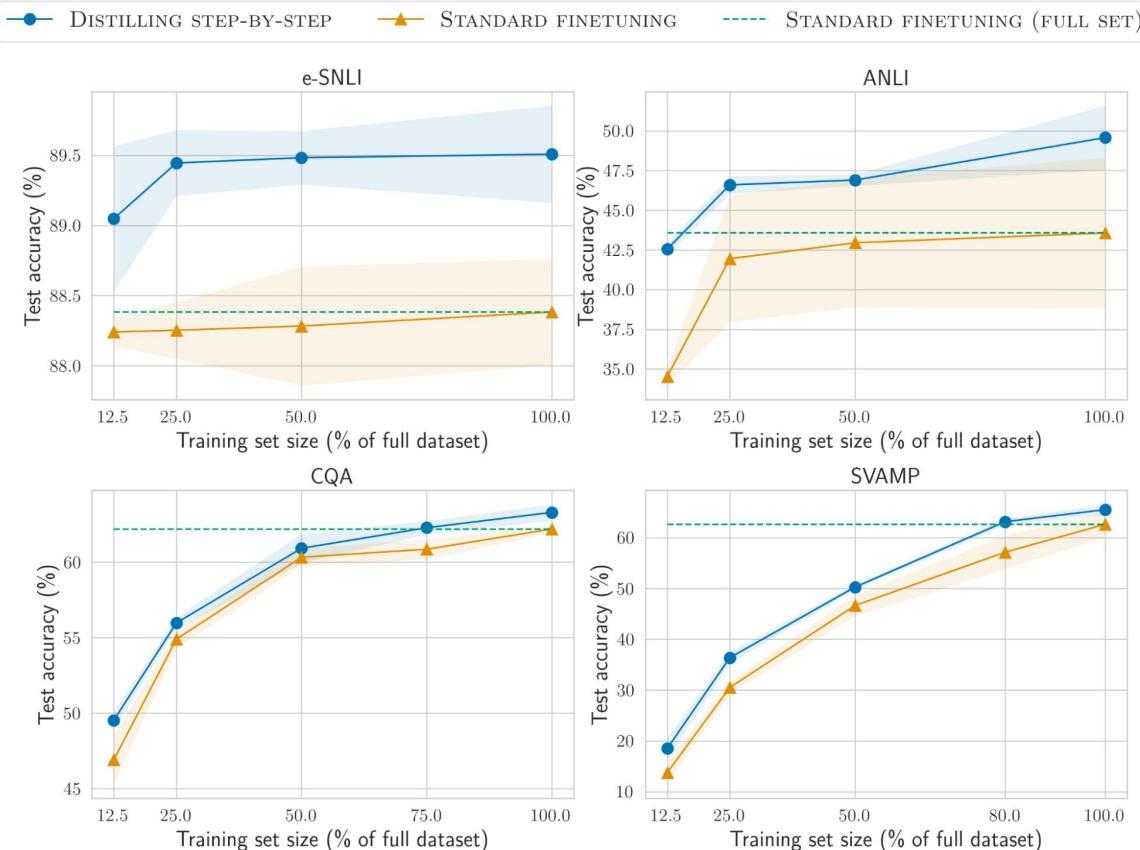
Input

Question: A person is carrying equipment for golf. What are they likely to have?  
Answer Choices: (a) club, (b) assembly hall, (c) meditation center, (d) meeting, (e) church  
Answer:

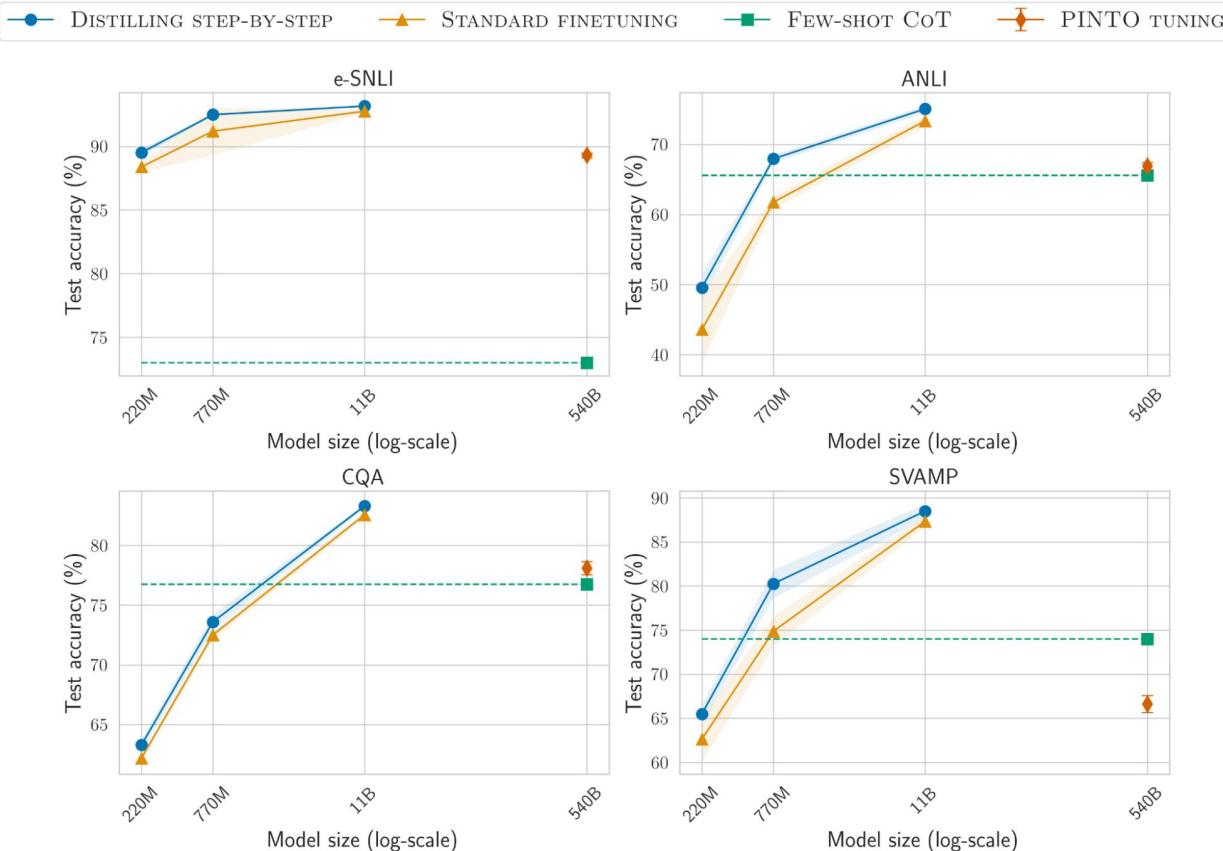
Output

The answer must be something that is used for golf. Of the above choices, only clubs are used for golf. So the answer is (a) club.

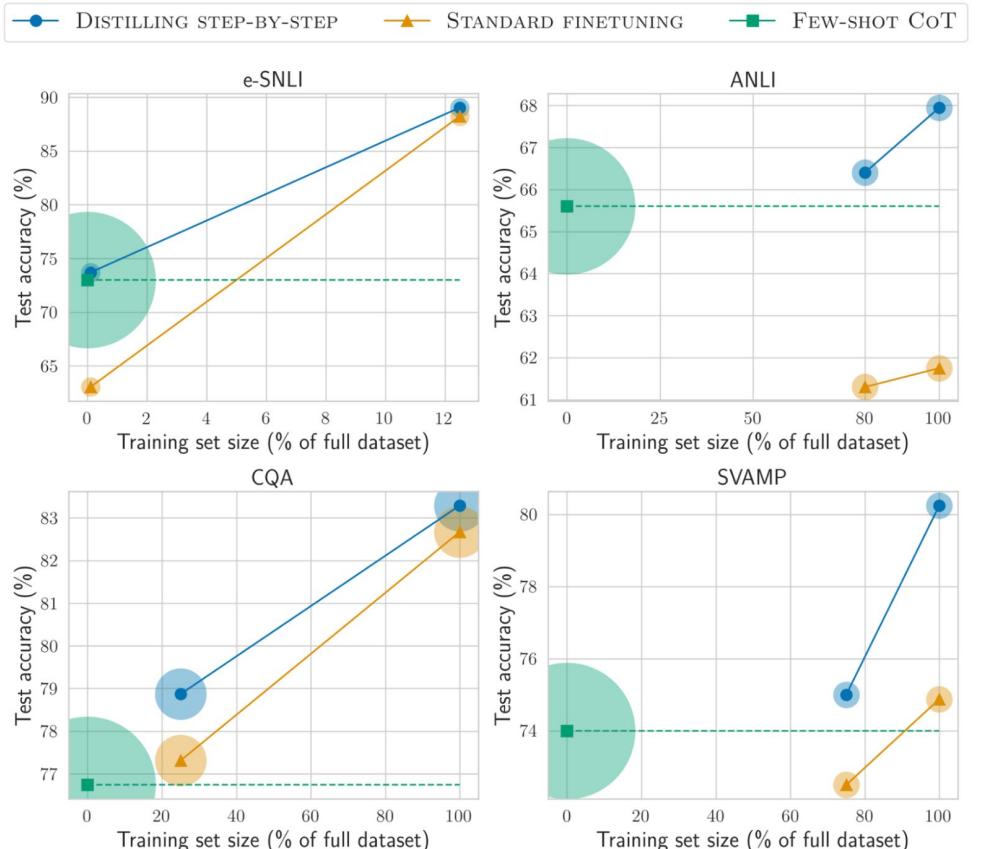
# Less training data



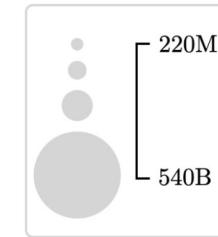
# Smaller deployed model size



# Distilling step-by-step outperforms few-shot LLMs with smaller models using less data



MODEL SIZE





---

# **DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning**

DeepSeek-AI

[research@deepseek.com](mailto:research@deepseek.com)

Model	AIME 2024		MATH-500	GPQA Diamond	LiveCode Bench	CodeForces
	pass@1	cons@64	pass@1	pass@1	pass@1	rating
<b>GPT-4o-0513</b>	9.3	13.4	74.6	49.9	32.9	759
<b>Claude-3.5-Sonnet-1022</b>	16.0	26.7	78.3	65.0	38.9	717
<b>OpenAI-o1-mini</b>	63.6	80.0	90.0	60.0	53.8	<b>1820</b>
<b>QwQ-32B-Preview</b>	50.0	60.0	90.6	54.5	41.9	1316
<b>DeepSeek-R1-Distill-Qwen-1.5B</b>	28.9	52.7	83.9	33.8	16.9	954
<b>DeepSeek-R1-Distill-Qwen-7B</b>	55.5	83.3	92.8	49.1	37.6	1189
<b>DeepSeek-R1-Distill-Qwen-14B</b>	69.7	80.0	93.9	59.1	53.1	1481
<b>DeepSeek-R1-Distill-Qwen-32B</b>	<b>72.6</b>	83.3	94.3	62.1	57.2	1691
<b>DeepSeek-R1-Distill-Llama-8B</b>	50.4	80.0	89.1	49.0	39.6	1205
<b>DeepSeek-R1-Distill-Llama-70B</b>	70.0	<b>86.7</b>	<b>94.5</b>	<b>65.2</b>	<b>57.5</b>	1633

Table 5 | Comparison of DeepSeek-R1 distilled models and other comparable models on reasoning-related benchmarks.

Published as a conference paper at ICLR 2023

---

# EDITING MODELS WITH TASK ARITHMETIC

**Gabriel Ilharco<sup>\*1</sup> Marco Tulio Ribeiro<sup>2</sup> Mitchell Wortsman<sup>1</sup> Suchin Gururangan<sup>1</sup>**

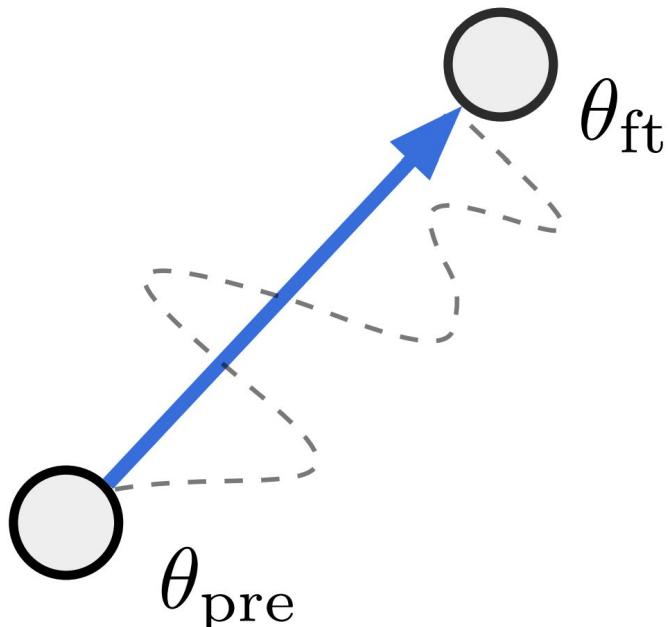
**Ludwig Schmidt<sup>1,3</sup> Hannaneh Hajishirzi<sup>1,3</sup> Ali Farhadi<sup>1</sup>**

<sup>1</sup>University of Washington <sup>2</sup>Microsoft Research <sup>3</sup>Allen Institute for AI

# Why do we want to edit LLMs?

- improve performance on downstream tasks
- mitigate biases or unwanted behavior
- align models with human preferences
- update models with new information

# The notion of task vectors



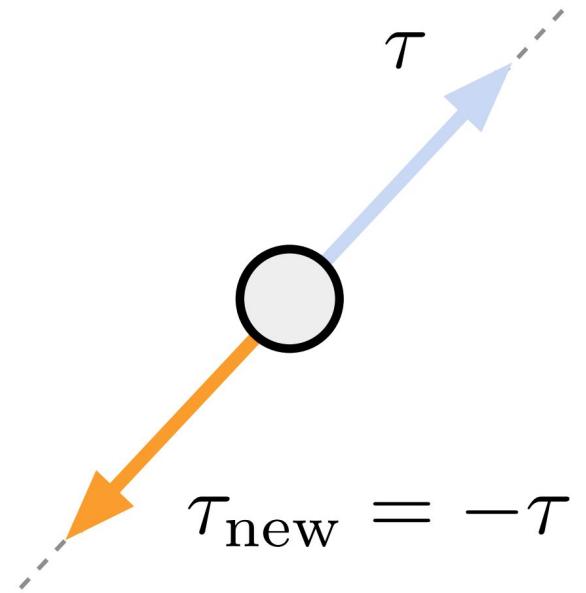
$$\tau = \theta_{\text{ft}} - \theta_{\text{pre}}$$

$$\theta_{\text{new}} = \theta + \tau$$

*In practice, we have an optional scaling term  $\lambda$*

$$\theta_{\text{new}} = \theta + \lambda\tau$$

# Forgetting via negation



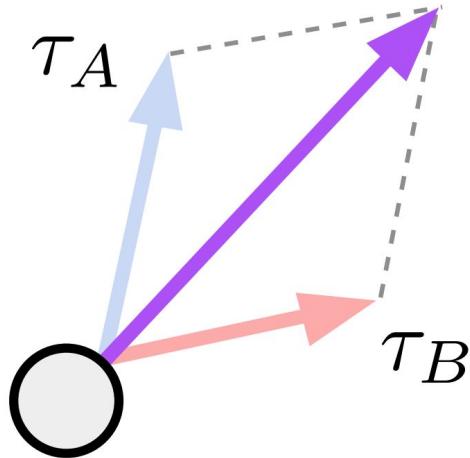
Example: making a language model produce less toxic content

$$\theta_{\text{new}} = \theta - \tau = \theta - (\theta_{ft} - \theta)$$

*In practice, we have an optional scaling term  $\lambda$*

# Learning via addition

$$\tau_{\text{new}} = \tau_A + \tau_B$$



$$\theta_{\text{new}} = \theta + \tau = \theta + (\tau_A + \tau_B)$$

$$= \theta + (\theta_A - \theta) + (\theta_B - \theta)$$

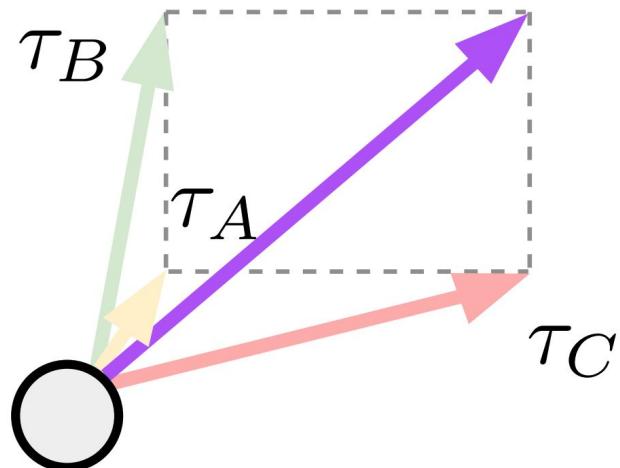
Example: building a multi-task model

*In practice, we have optional scaling terms  $\lambda_A, \lambda_B$*

# Task analogies

*"A is to B as C is to D"*

$$\tau_{\text{new}} = \tau_C + (\tau_B - \tau_A)$$



$$\tau_B - \tau_A = \tau_D - \tau_C$$

$$\tau_D = \tau_C + (\tau_B - \tau_A)$$

$$\theta_{\text{new}} = \theta + \tau_C + (\tau_B - \tau_A)$$

$$= \theta + (\theta_C - \theta) + (\theta_B - \theta) - (\theta_A - \theta)$$

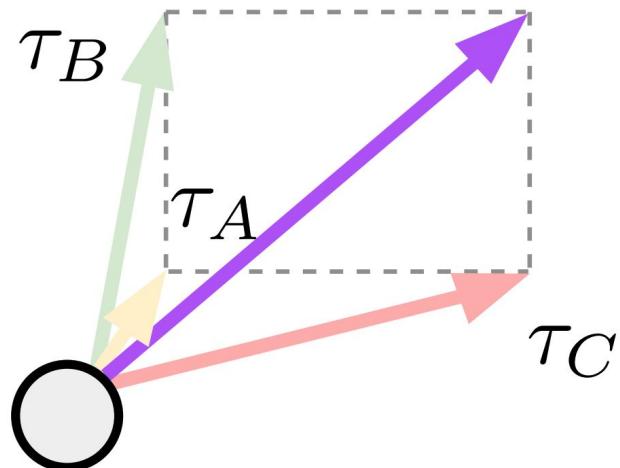
Example: improving  
domain generalization

*In practice, we have optional  
scaling terms  $\lambda_A, \lambda_B, \lambda_c$*

# Task analogies

*"A is to B as C is to D"*

$$\tau_{\text{new}} = \tau_C + (\tau_B - \tau_A)$$



$$\tau_B - \tau_A = \tau_D - \tau_C$$

$$\tau_D = \tau_C + (\tau_B - \tau_A)$$

$$\theta_{\text{new}} = \theta + \tau_C + (\tau_B - \tau_A)$$

$$= \theta + (\theta_C - \theta) + (\theta_B - \theta) - (\theta_A - \theta)$$

Example: improving  
domain generalization

*In practice, we have optional  
scaling terms  $\lambda_A, \lambda_B, \lambda_c$*

# Making language models less toxic with negative task vectors

Method	% toxic generations (↓)	Avg. toxicity score (↓)	WikiText-103 perplexity (↓)
Pre-trained	4.8	0.06	16.4
Fine-tuned	57	0.56	16.6
Gradient ascent	0.0	0.45	>10 <sup>10</sup>
Fine-tuned on non-toxic	1.8	0.03	17.2
Random vector	4.8	0.06	16.4
Negative task vector	0.8	0.01	16.9

# What Matters for Model Merging at Scale?

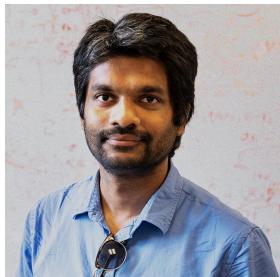


Prateek  
Yadav

Tu Vu

Jonathan Lai

Alexandra  
Chronopoulou

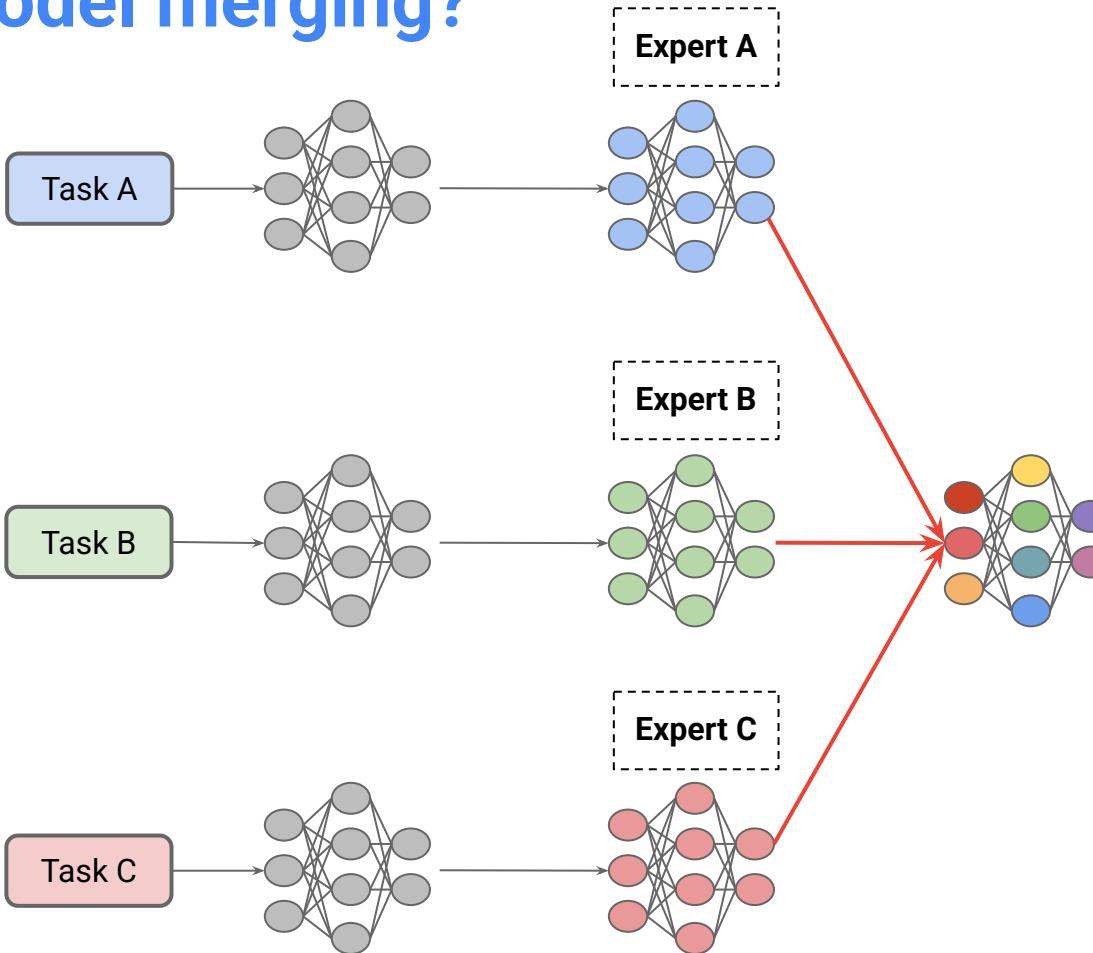


Manaal Faruqui

Mohit Bansal

Tsendsuren  
Munkhdalai

# What is model merging?



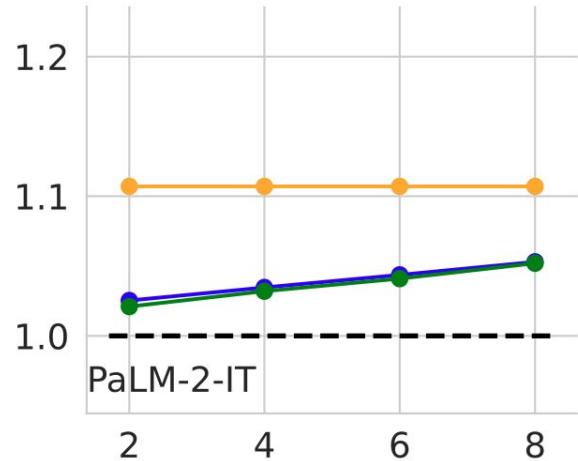
# Why model merging?

- dramatically reduces storage and serving costs by reusing a single model across tasks
- enables compositional combination of capabilities from expert models, which can improve generalization to novel tasks
- supports decentralized and modular model development by allowing multiple contributors to independently build models and later combine them together

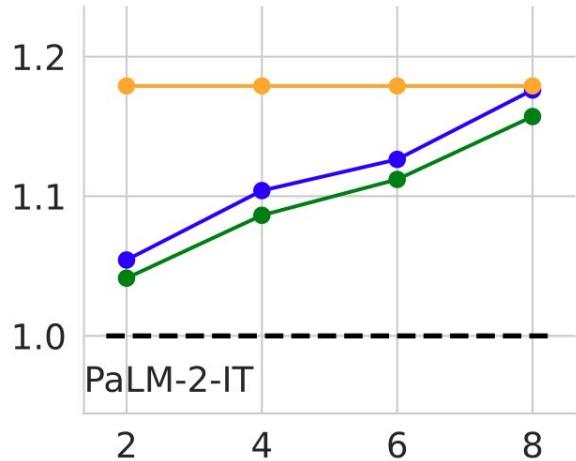
# Merging boosts zero-shot generalization

Held-Out Performance

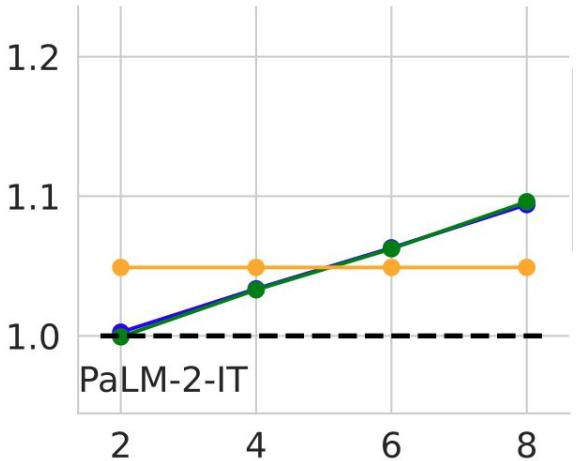
**PaLM-2-IT, 1B**



**PaLM-2-IT, 24B**



**PaLM-2-IT, 64B**



# of Experts

- Merging Method**
- Task Arithmetic
  - TIES
  - Multitask

# **Efficient Model Development through Fine-tuning Transfer**

**Pin-Jie Lin<sup>1</sup>**

*pinjie@vt.edu*

**Rishab Balasubramanian<sup>1</sup>**

*rishbb@vt.edu*

**Fengyuan Liu<sup>2</sup>**

*fy.liu@mail.utoronto.ca*

**Nikhil Kandpal<sup>2</sup>**

*nkandpa2@cs.toronto.edu*

**Tu Vu<sup>1</sup>**

*tuvu@vt.edu*

<sup>1</sup> *Virginia Tech*    <sup>2</sup> *University of Toronto & Vector Institute*

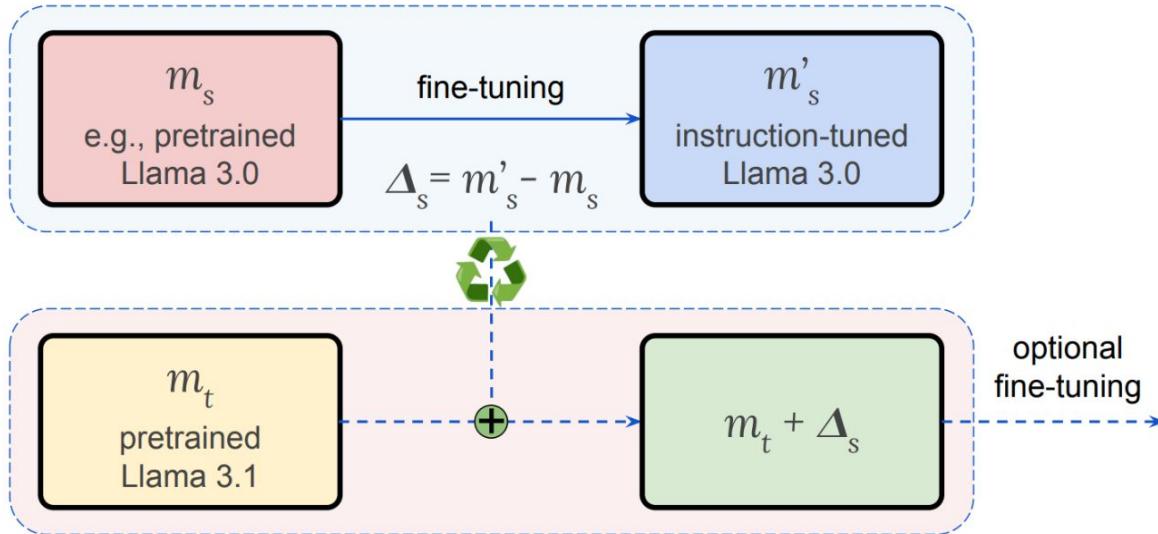


Figure 1: To transfer fine-tuning (e.g., instruction tuning) from a *source* model version  $s$  (e.g., Llama 3.0) to a *target* version  $t$  (Llama 3.1), we first compute the diff vector  $\Delta_s = m'_s - m_s$  from version  $s$ , where  $m'_s$  is the fine-tuned model (instruction-tuned Llama 3.0) and  $m_s$  is the base model (pretrained Llama 3.0). Then, we add  $\Delta_s$  to the target base model (pretrained Llama 3.1) to approximate the fine-tuned model in version  $t$  (instruction-tuned Llama 3.1). We explore two scenarios: (1) *recycling*—transferring from an older model version to a newer one to reduce retraining, and (2) *backporting*—transferring from a newer version to an older one to take advantage of the newer fine-tuning while maintaining optimization for specific use cases.

# Transferring fine-tuning updates

Model	GSM8K	MATH	ARCC	GPQA	MMLU	IFEval
Llama 3.0 8B Instruct	81.1	28.8	82.4	<b>31.5</b>	64.9	<b>76.6</b>
Llama 3.0 8B + $\Delta_{3.1}$	55.6 <b>82.8</b>	17.3 <b>44.7</b>	79.7 <b>83.0</b>	22.3 25.9	66.7 <b>70.0</b>	34.5 <b>76.6</b>
Llama 3.1 8B Instruct	<b>86.5</b>	<b>50.3</b>	<b>83.8</b>	31.3	<b>72.9</b>	80.5
Llama 3.1 8B + $\Delta_{3.0}$	56.6 79.8	19.3 29.9	79.2 82.9	21.9 <b>32.6</b>	66.8 65.1	36.4 <b>83.3</b>

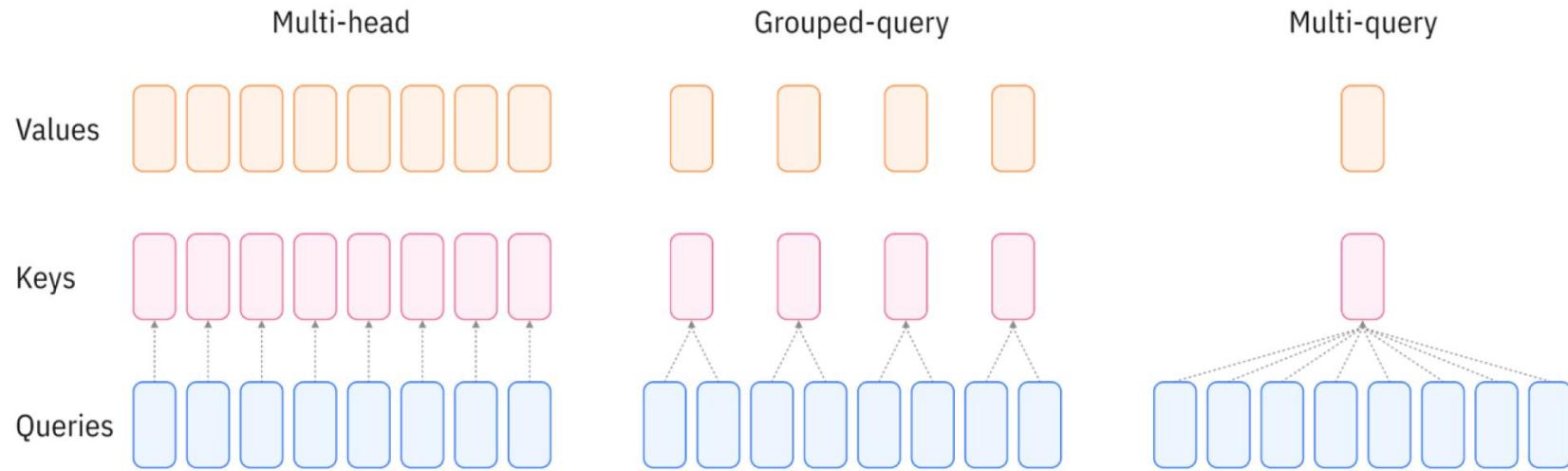
Table 1: Fine-tuning transfer significantly improves the performance of the target base model across various tasks, achieving results comparable to its fine-tuned counterpart in many cases. Here,  $\Delta_{3.0}$  and  $\Delta_{3.1}$  represent the diff vectors between Llama Instruct and Llama for versions 3.0 and 3.1, respectively. Notably, adding the diff vector  $\Delta_s$  from a different model version can effectively transform a non-instruction-tuned model (e.g., Llama 3.0 or Llama 3.1) into one that follows instructions well (Llama 3.0 +  $\Delta_{3.1}$  or Llama 3.1 +  $\Delta_{3.0}$ ) without further training. Additional results for OLMo and Tülu can be found in Appendix A, where we additionally find that advanced LLM capabilities, attained through alignment tuning stages such as Supervised Fine-Tuning (SFT), Direct Preference Optimization (DPO), or Group Relative Policy Optimization (GRPO), can be successfully transferred across different model versions.

# Multilingual model development

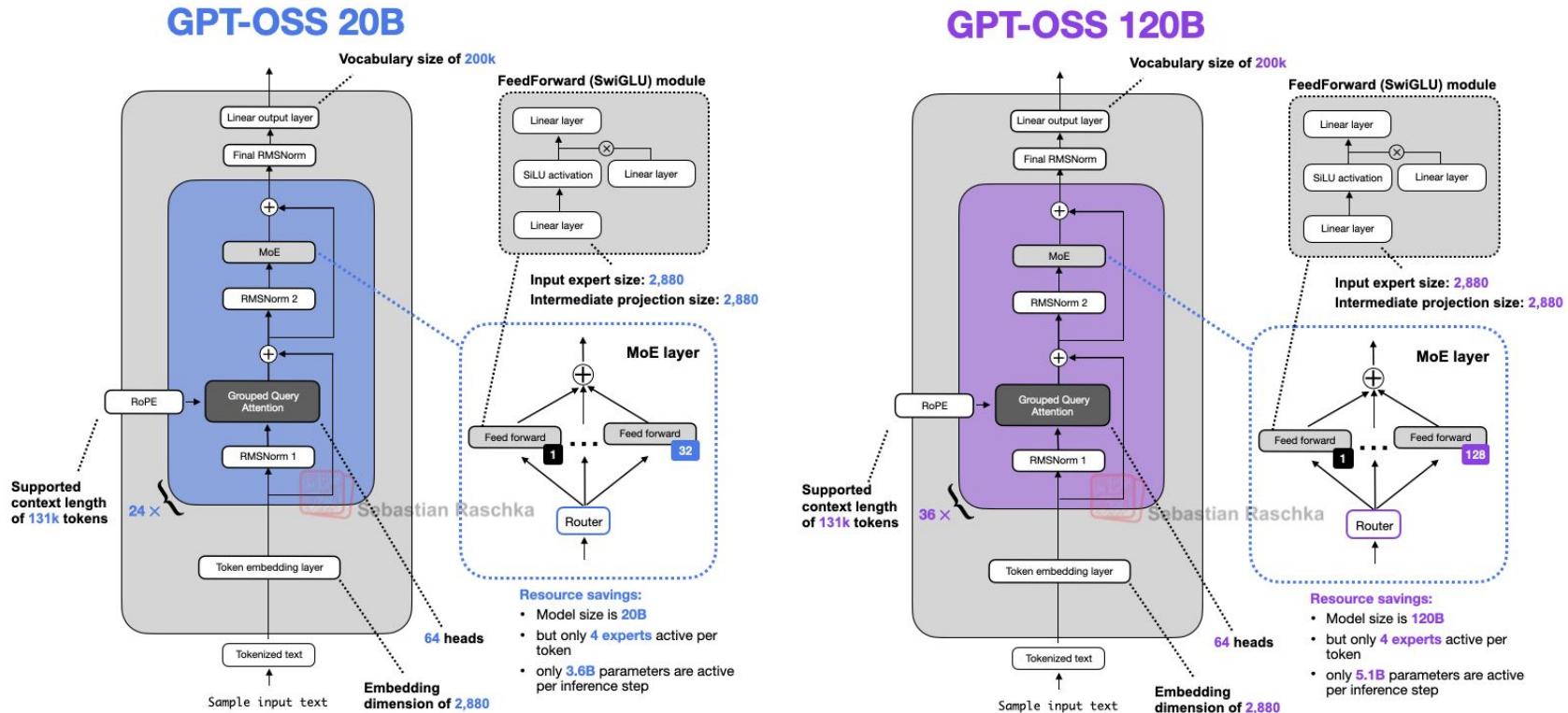
Model	Malagasy	Sinhala	Turkish
Llama 3.0 8B Instruct	23.1	23.3	30.8
+ FT	30.8	29.0	43.2
Llama 3.1 8B Instruct	27.6	<b>33.0</b>	27.7
+ $\Delta_{3.0}$	<b>32.3</b>	32.3	<b>43.2</b>

Table 2: Recycling fine-tuning updates improves multilingual performance on Global MMLU without re-training, yielding a 4.7% and 15.5% absolute improvement for Malagasy and Turkish, respectively, compared to Llama 3.1 8B Instruct.  $\Delta_{3.0}$  represents the diff vector between Llama 3.0 Instruct and its monolingual fine-tuned (FT) version.

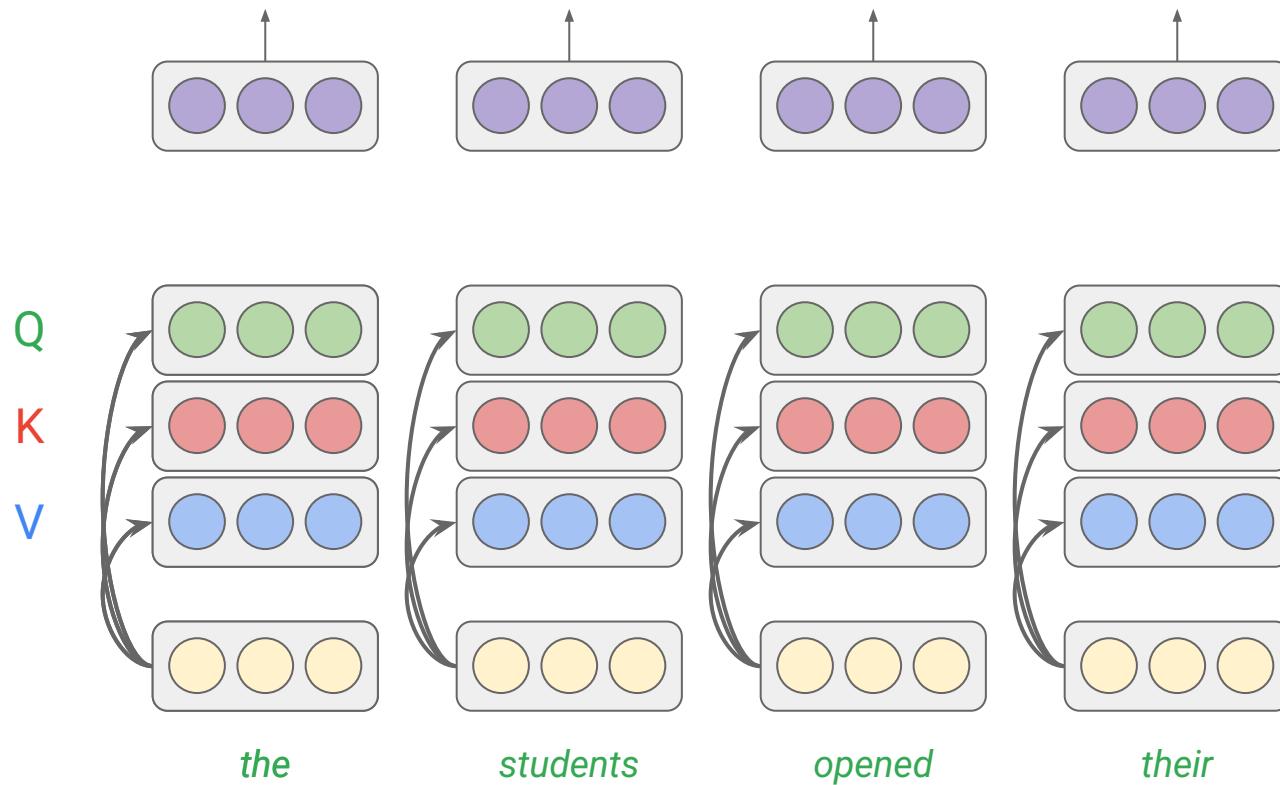
# Multi-query attention / grouped-query attention



# Multi-query attention / grouped-query attention



# KV caching



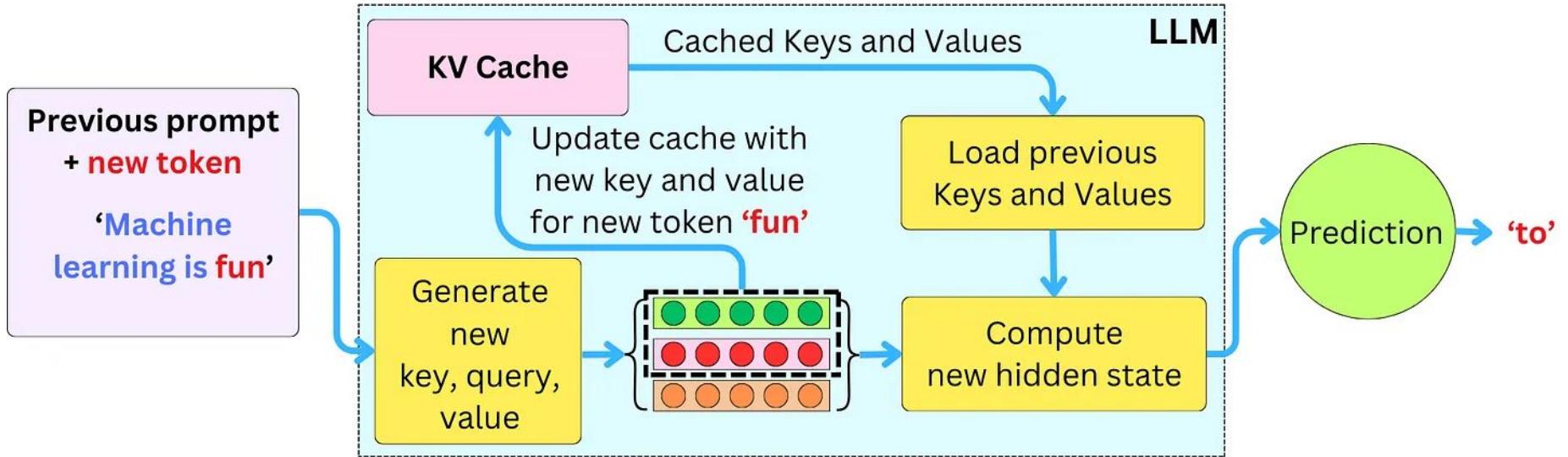
$$Q = X \cdot W_Q$$

$$K = X \cdot W_K$$

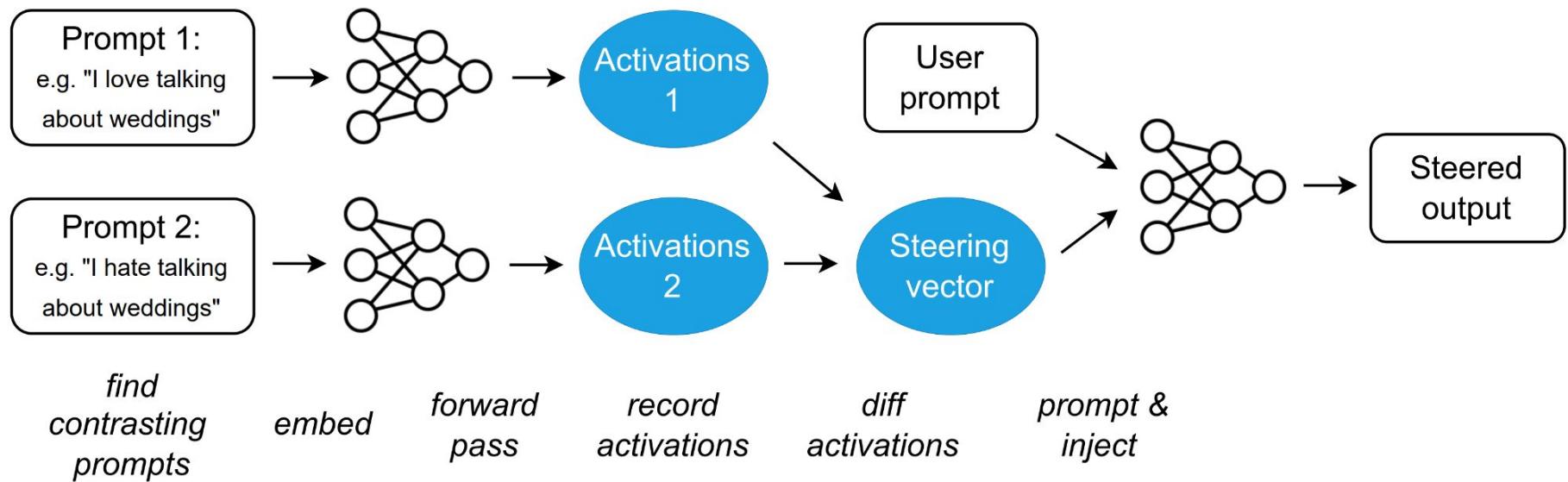
$$V = X \cdot W_V$$

linear  
projections

# KV caching (cont'd)



# Steering vectors



# Steering vectors (cont'd)

---

Prompt	+	steering	=	completion
I hate you because...	[None]			...you are the most disgusting thing I have ever seen.
	ActAdd (love)			...you are so beautiful and I want to be with you forever.
I went up to my friend and said...	[None]			...“I’m sorry, I can’t help you.” “No,” he said. “You’re not.”
	ActAdd (weddings)			...“I’m going to talk about the wedding in this episode of Wedding Season. I think it’s a really good episode. It’s about how you’re supposed to talk about weddings.”

---

# Pruning

- Remove parameters from the model after training

Published as a conference paper at ICLR 2024

---

# A SIMPLE AND EFFECTIVE PRUNING APPROACH FOR LARGE LANGUAGE MODELS

**Mingjie Sun**<sup>1\*</sup>

**Zhuang Liu**<sup>2\*</sup>

**Anna Bair**<sup>1</sup>

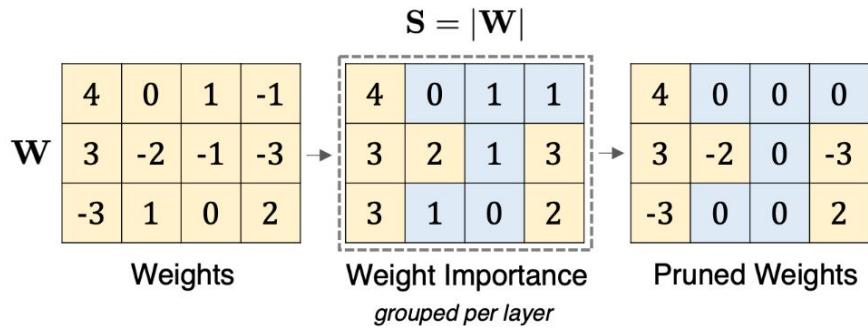
**J. Zico Kolter**<sup>1,3</sup>

<sup>1</sup>Carnegie Mellon University

<sup>2</sup>Meta AI Research

<sup>3</sup>Bosch Center for AI

## Magnitude Pruning



## Wanda

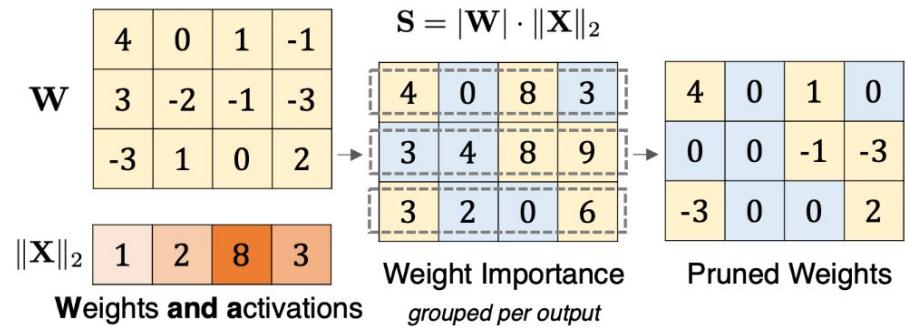


Figure 1: Illustration of our proposed method Wanda (Pruning by **Weights and activations**), compared with the magnitude pruning approach. Given a weight matrix  $\mathbf{W}$  and input feature activations  $\mathbf{X}$ , we compute the weight importance as the elementwise product between the weight magnitude and the norm of input activations ( $|\mathbf{W}| \cdot \|\mathbf{X}\|_2$ ). Weight importance scores are compared on a *per-output* basis (within each row in  $\mathbf{W}$ ), rather than globally across the entire matrix.

---

# Are Sixteen Heads Really Better than One?

---

**Paul Michel**

Language Technologies Institute  
Carnegie Mellon University  
Pittsburgh, PA  
pmichel1@cs.cmu.edu

**Omer Levy**

Facebook Artificial Intelligence Research  
Seattle, WA  
omerlevy@fb.com

**Graham Neubig**

Language Technologies Institute  
Carnegie Mellon University  
Pittsburgh, PA  
gneubig@cs.cmu.edu

Published as a conference paper at ICLR 2019

---

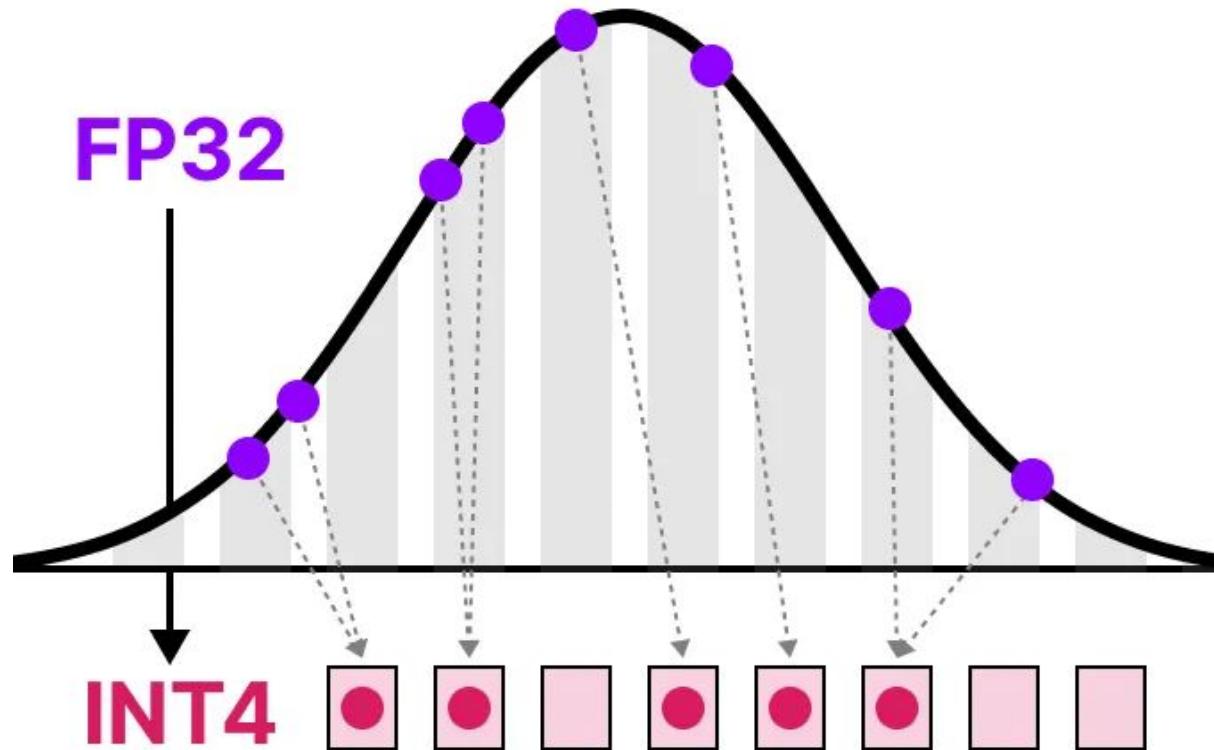
# THE LOTTERY TICKET HYPOTHESIS: FINDING SPARSE, TRAINABLE NEURAL NETWORKS

**Jonathan Frankle**  
MIT CSAIL  
[jfrankle@csail.mit.edu](mailto:jfrankle@csail.mit.edu)

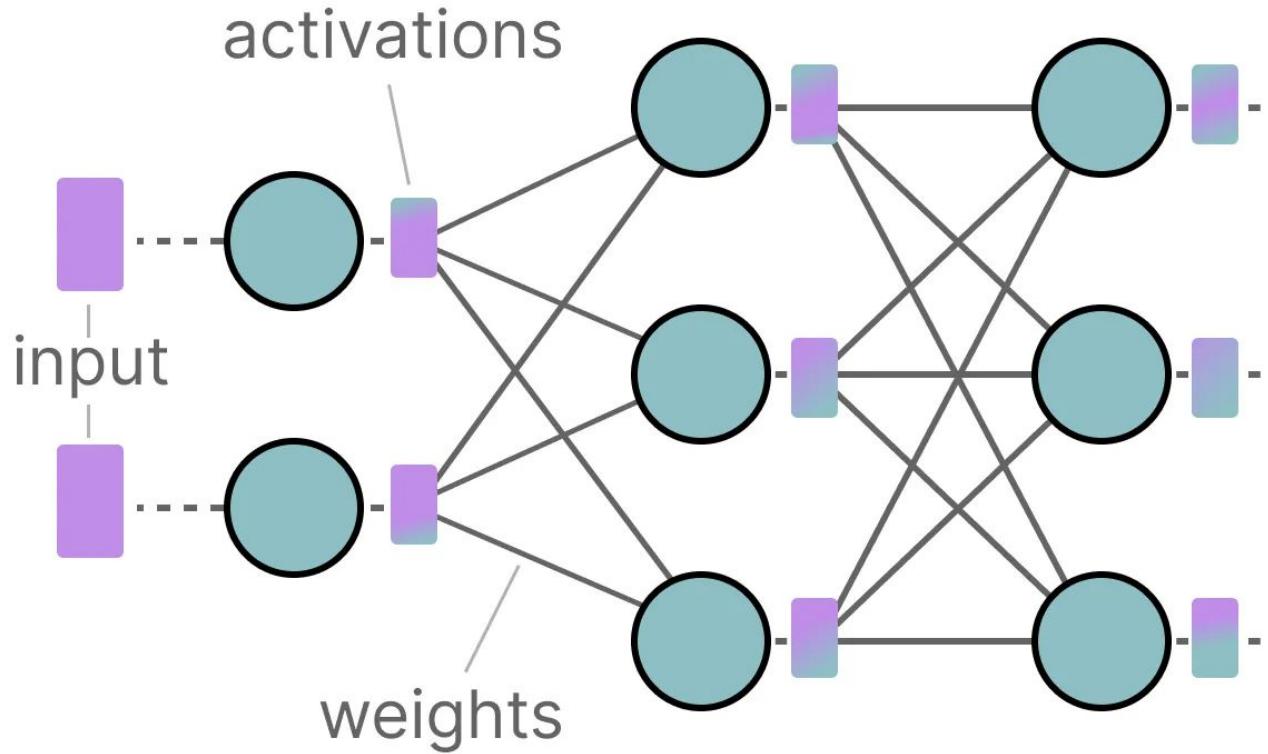
**Michael Carbin**  
MIT CSAIL  
[mcarbin@csail.mit.edu](mailto:mcarbin@csail.mit.edu)

*Training a pruned randomly-initialized networks can be better than training the full randomly-initialized network*

# Quantization

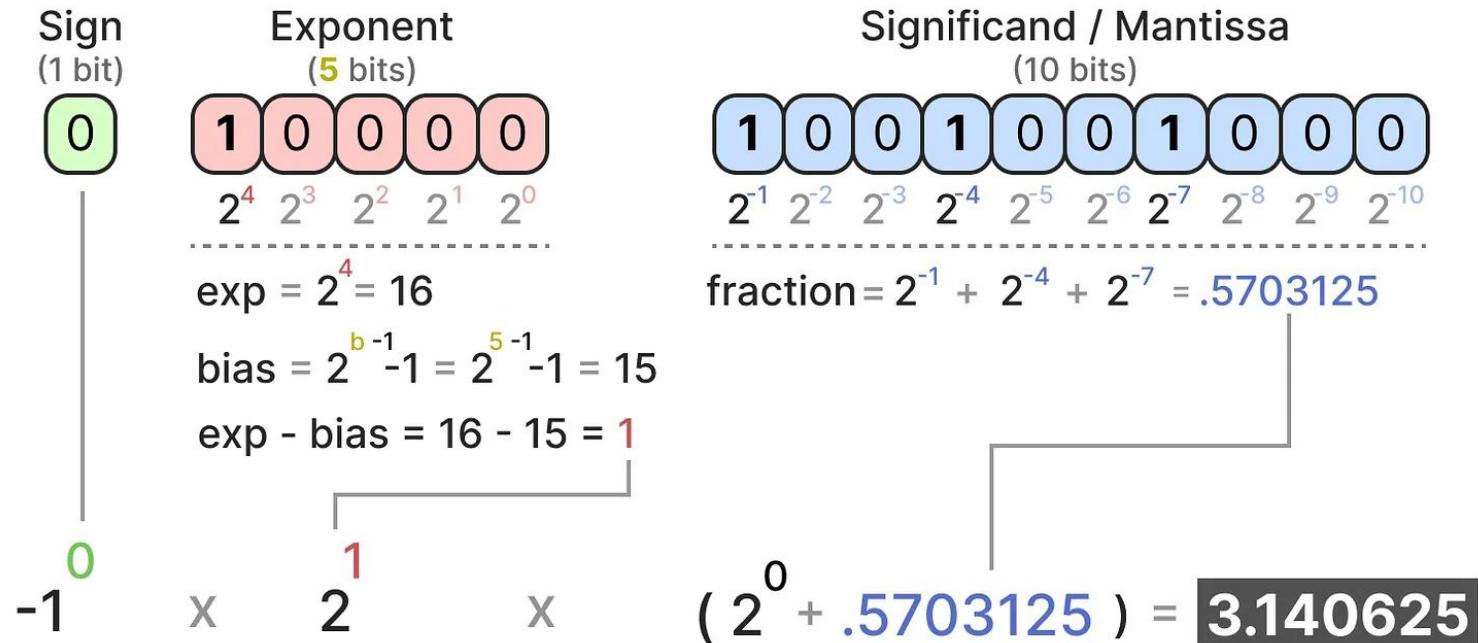


# Quantizing both the weights and activations



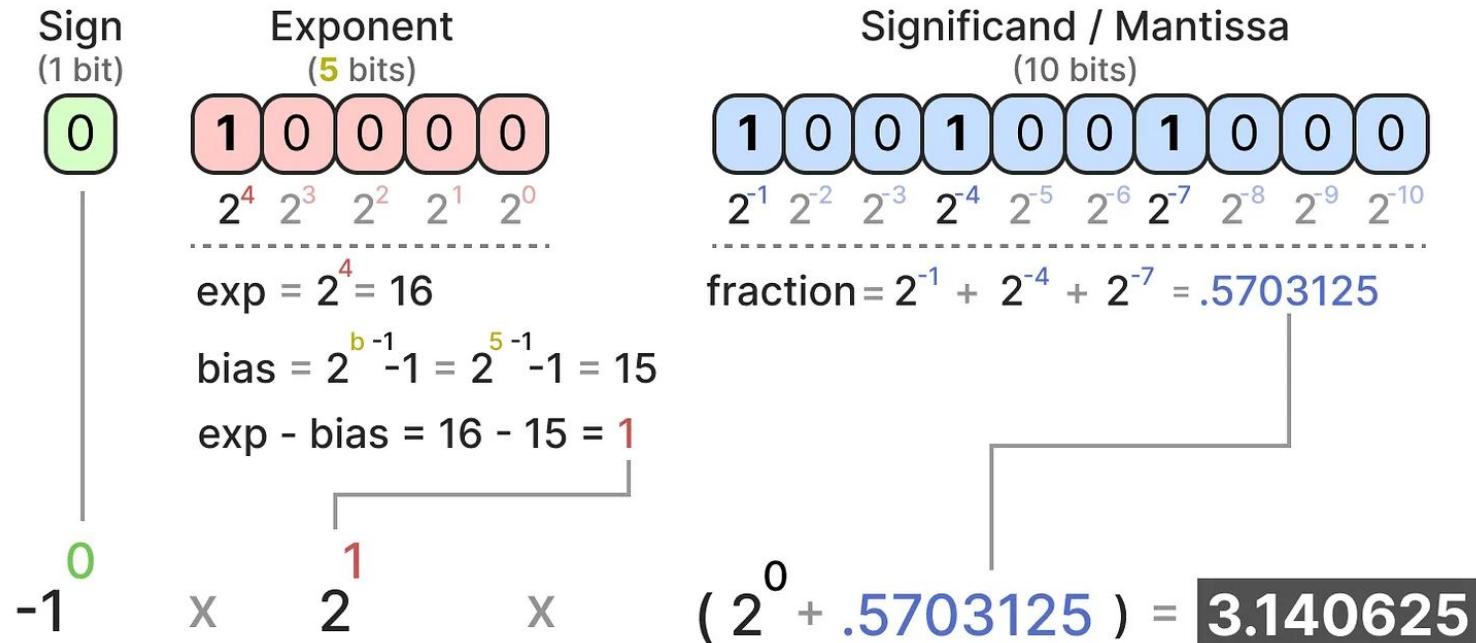
# How to represent numerical values

## Float 16-bit (FP16)



# How to represent numerical values (cont'd)

## Float 16-bit (FP16)



# How to represent numerical values (cont'd)

## Float 32-bit (FP32)

0 10000000 1001001000011111011011

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

higher precision

## Float 16-bit (FP16)

0 10000 1001001000

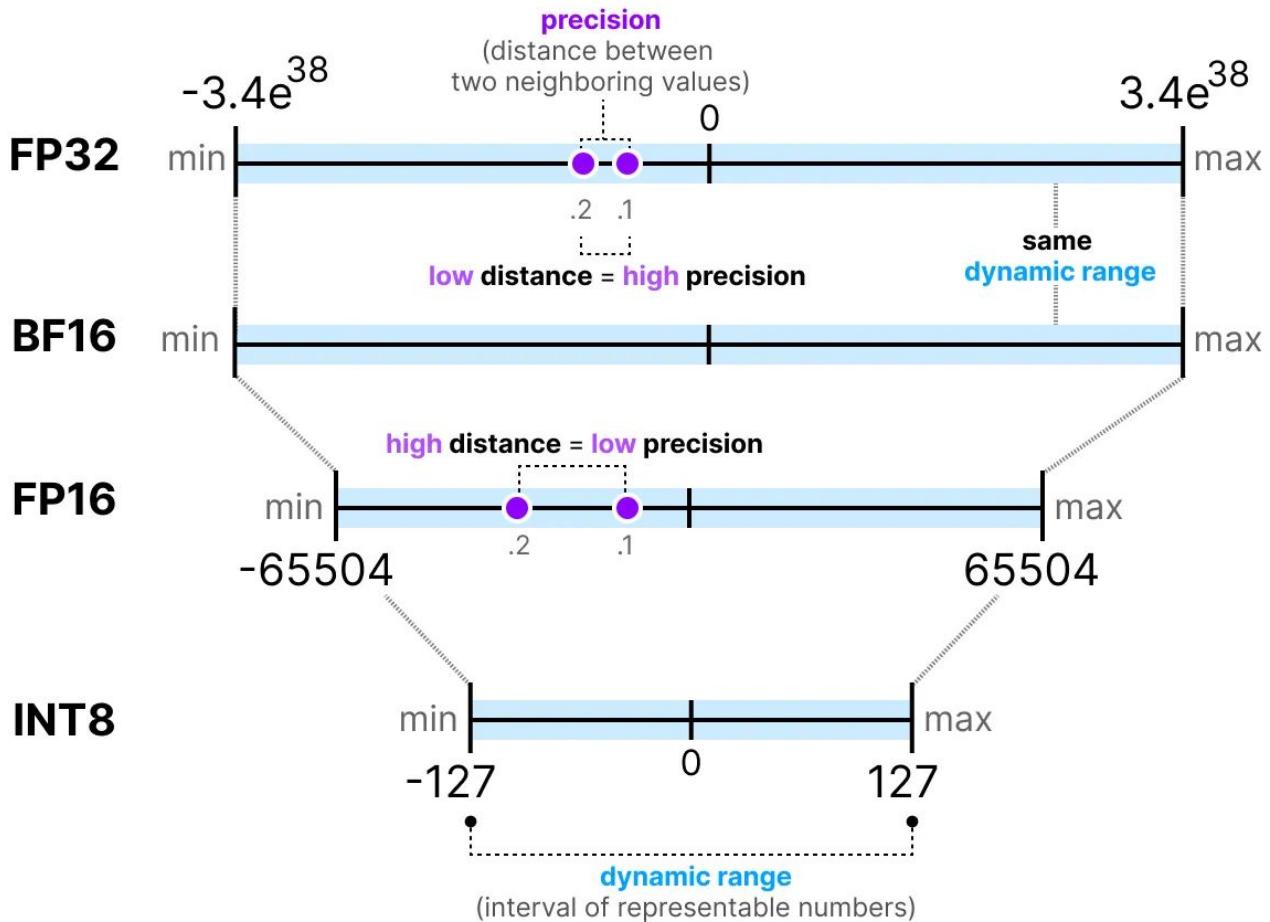
$$(-1)^0 \times 2^1 \times 1.5703125 = 3.140625$$

lower precision

original value

**3.1415927**

# Memory constraints



## Memory constraints (cont'd)

$$\text{memory} = \frac{\text{nr\_bits}}{8} \times \text{nr\_params}$$

---

$$\mathbf{64\text{-bits}} = \frac{64}{8} \times 70\text{B} \approx \mathbf{560 \text{ GB}}$$

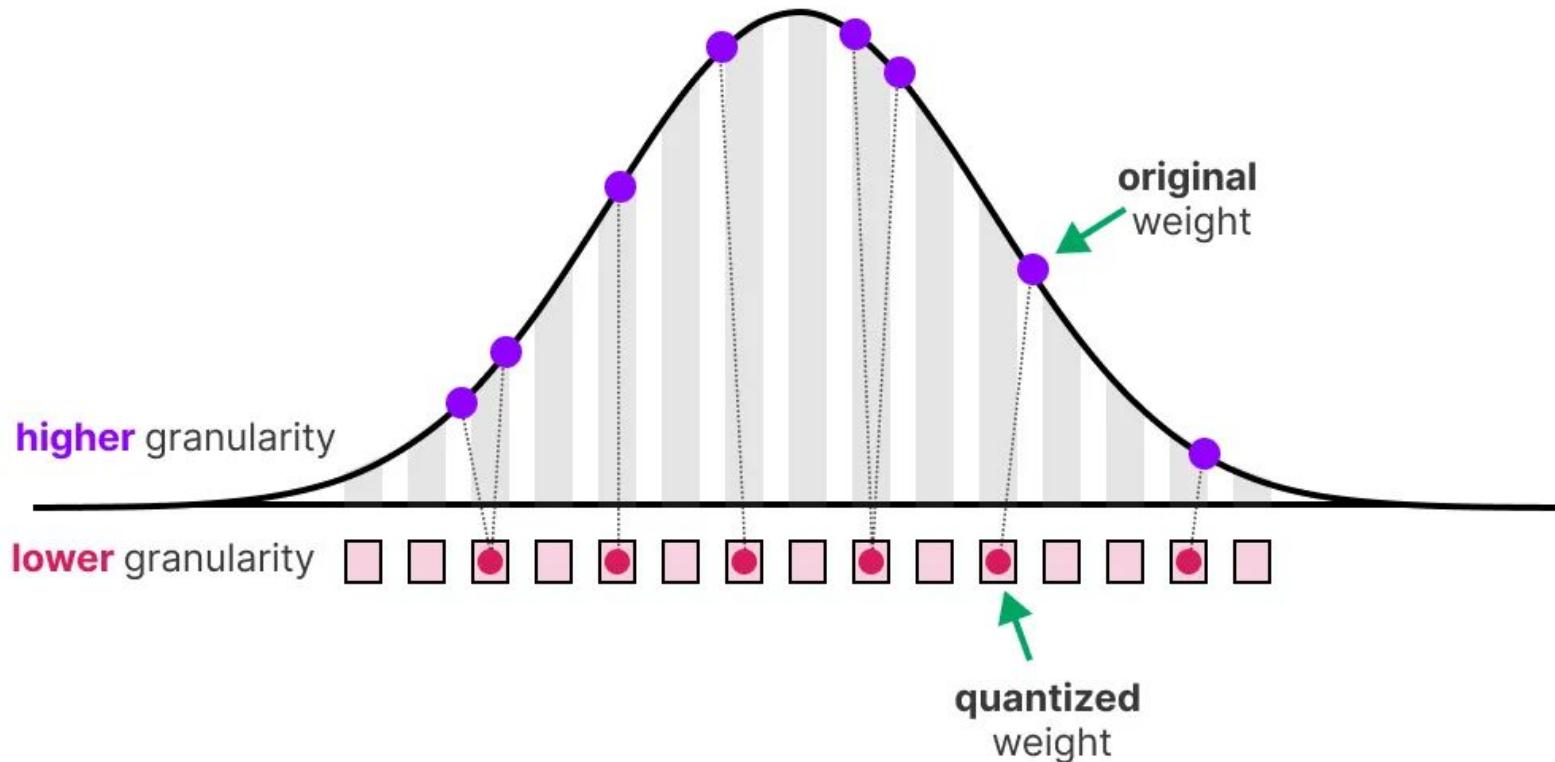
---

$$\mathbf{32\text{-bits}} = \frac{32}{8} \times 70\text{B} \approx \mathbf{280 \text{ GB}}$$

---

$$\mathbf{16\text{-bits}} = \frac{16}{8} \times 70\text{B} \approx \mathbf{140 \text{ GB}}$$

# Quantization



# Quantization

Original Image

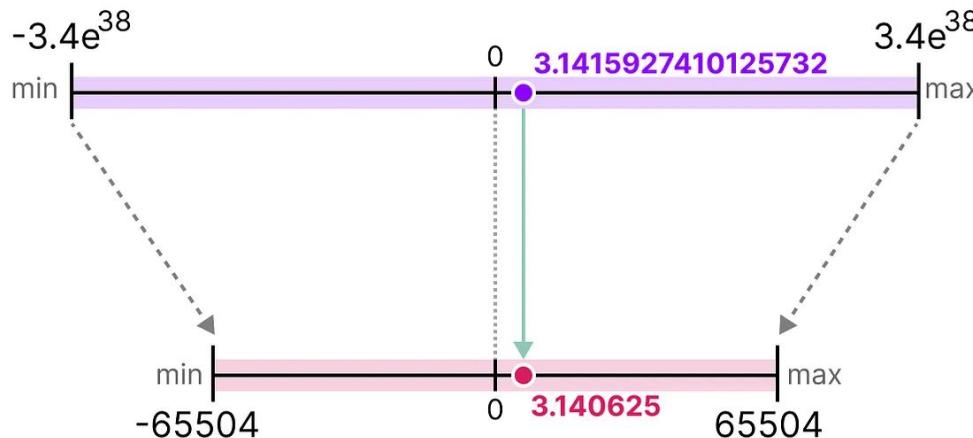


“Quantized” Image



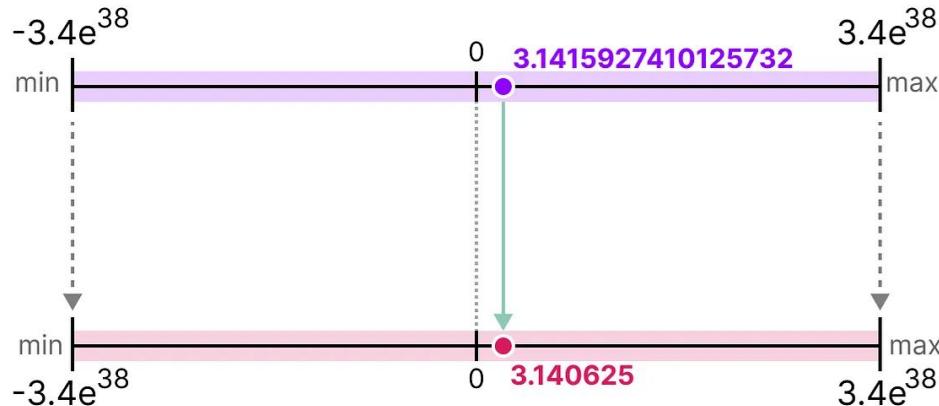
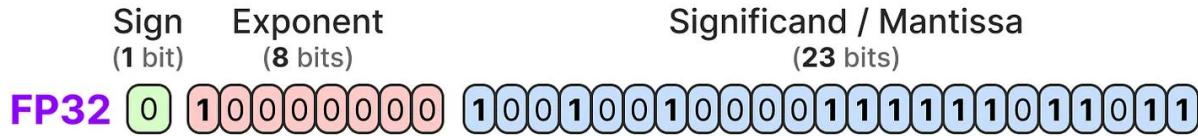
# Common data types: FP16 (half precision)

Sign (1 bit)	Exponent (8 bits)	Significand / Mantissa (23 bits)
FP32	0 10000000	1001001000011111011011



FP16    0 100000 1001001000  
(1 bit)    (5 bits)    (10 bits)

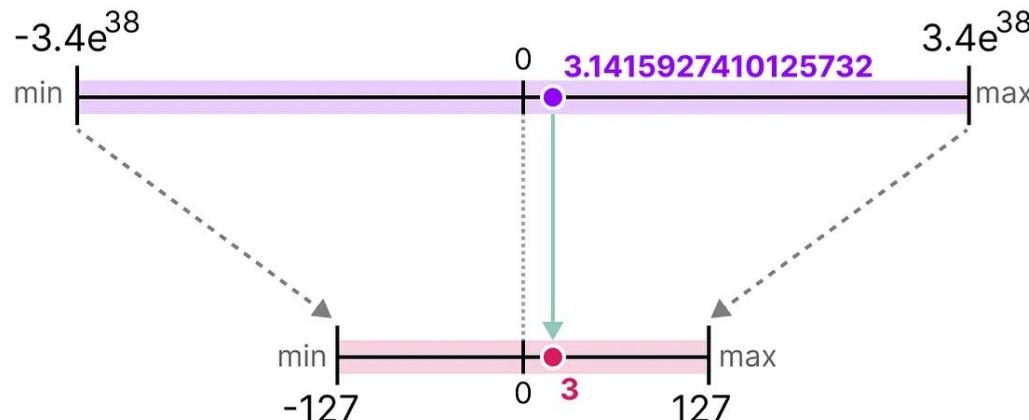
# Common data types: BF16



# Common data types: INT8

FP32      Sign (1 bit)      Exponent (8 bits)      Significand / Mantissa (23 bits)

0 10000000 1001001000011111011011

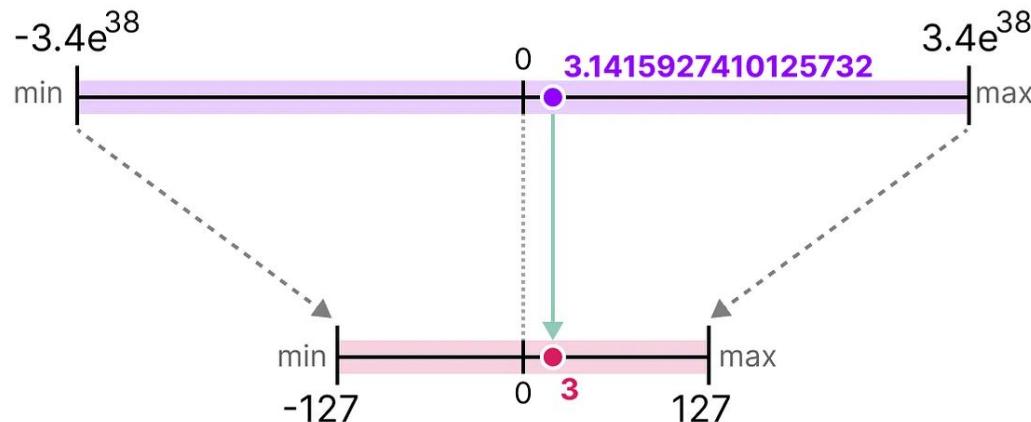


(signed) INT8      0 1001000  
(1 bit)      (7 bits)

# Common data types: INT8

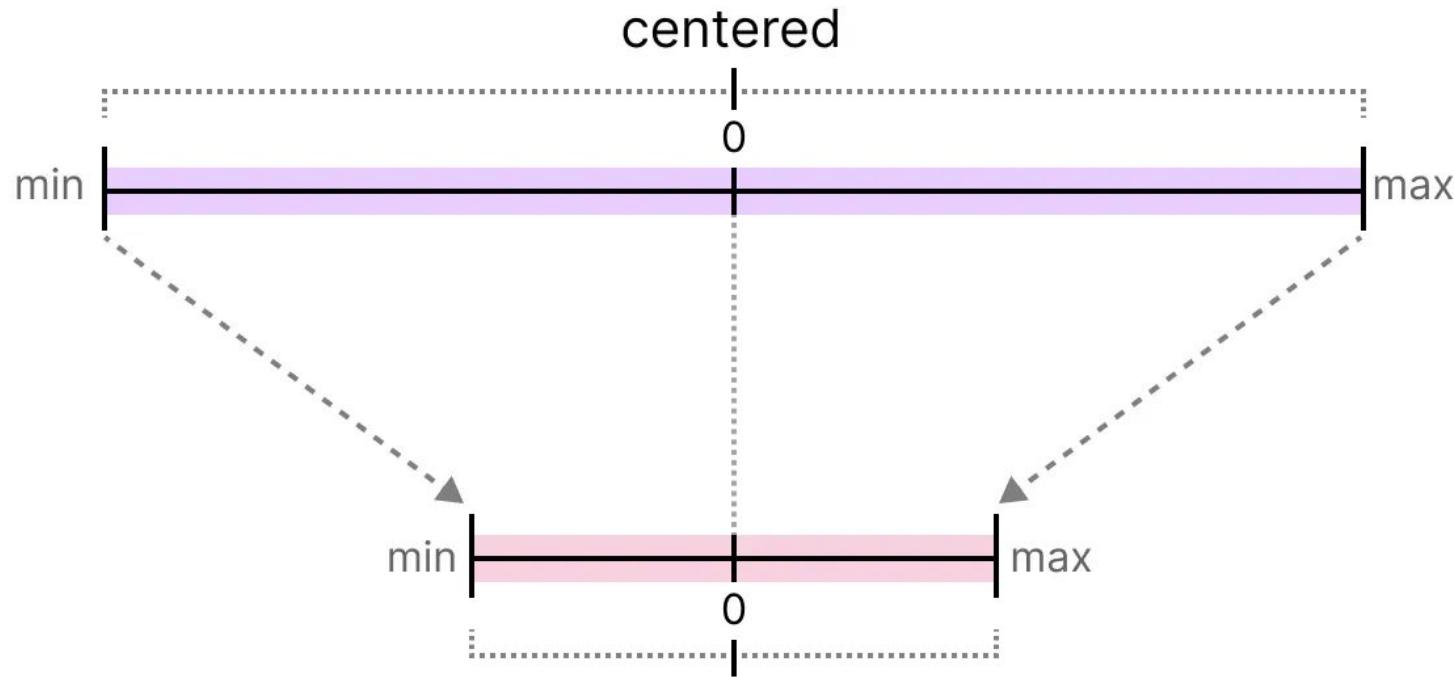
FP32      Sign (1 bit)      Exponent (8 bits)      Significand / Mantissa (23 bits)

0 10000000 1001001000011111011011



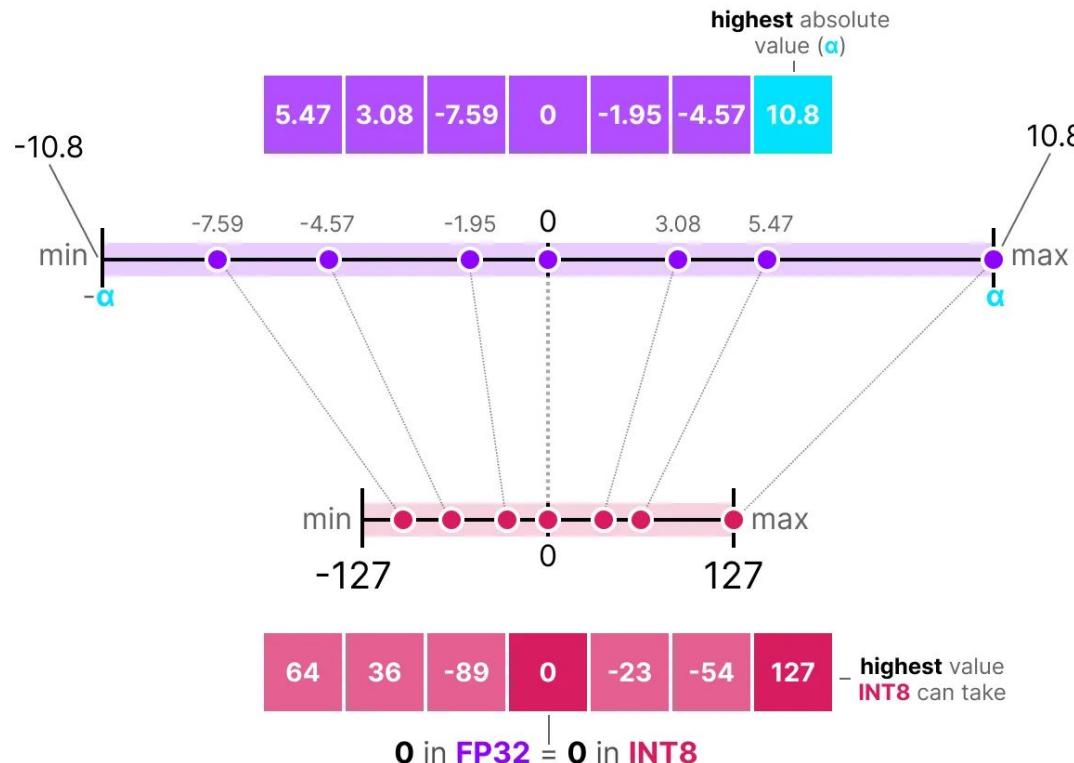
(signed) INT8      0 1001000  
(1 bit)      (7 bits)

# Symmetric quantization



**0 in FP32 = 0 in INT8**

# Absolute maximum (absmax) quantization



# Absolute maximum (absmax) quantization

We first calculate a scale factor ( $s$ ) using:

- $b$  is the number of bytes that we want to quantize to (8),
- $\alpha$  is the *highest* absolute value,

Then, we use the  $s$  to quantize the input  $x$ :

$$s = \frac{2^{b-1} - 1}{\alpha} \quad (\text{scale factor})$$

$$x_{\text{quantized}} = \text{round}(s \cdot x) \quad (\text{quantization})$$

Filling in the values would then give us the following:

$$s = \frac{127}{10.8} = 11.76 \quad (\text{scale factor})$$

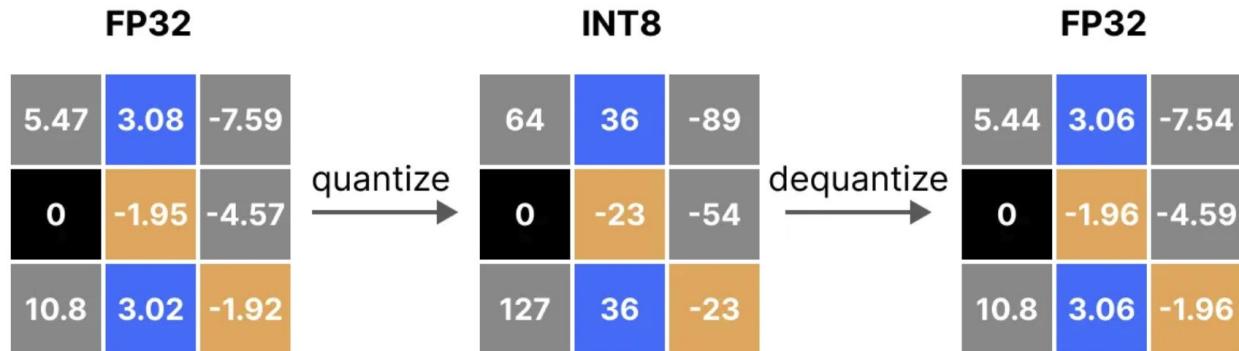
$$x_{\text{quantized}} = \text{round}(11.76 \cdot \text{█████}) \quad (\text{quantization})$$

# Dequantization

To retrieve the original FP32 values, we can use the previously calculated *scaling factor* ( $s$ ) to *dequantize* the quantized values.

$$X_{\text{dequantized}} = \frac{\text{█}}{s} \quad (\text{dequantize})$$

Applying the quantization and then dequantization process to retrieve the original looks as follows:



# Dequantization

FP32 (original)			FP32 (dequantized)			Quantization error		
5.47	3.08	-7.59	5.44	3.06	-7.54	.03	.02	.05
0	-1.95	-4.57	0	-1.96	-4.59	0	-.01	-.02
10.8	3.02	-1.92	10.8	3.06	-1.96	0	-.04	-.04

# Dequantization

FP32 (original)		
5.47	3.08	-7.59
0	-1.95	-4.57
10.8	3.02	-1.92

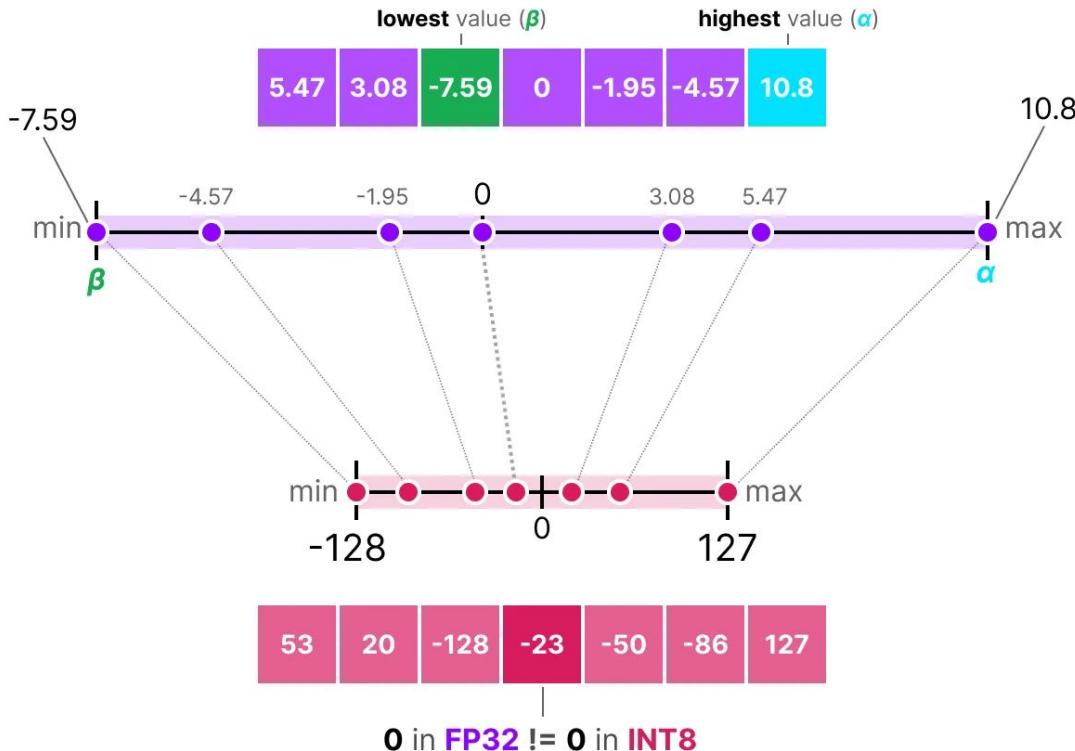
FP32 (dequantized)		
5.44	3.06	-7.54
0	-1.96	-4.59
10.8	3.06	-1.96

Quantization error		
.03	.02	.05
0	-.01	-.02
0	-.04	-.04

-

=

# Asymmetric quantization



## Asymmetric quantization (cont'd)

$$S = \frac{128 - -127}{\alpha - \beta} \quad (\text{scale factor})$$

---

$$Z = \text{round}(-S \cdot \beta) - 2^{b-1} \quad (\text{zeropoint})$$

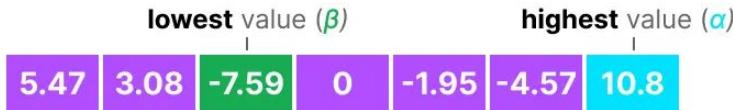
---

$$X_{\text{quantized}} = \text{round}(S \cdot X + Z) \quad (\text{quantization})$$

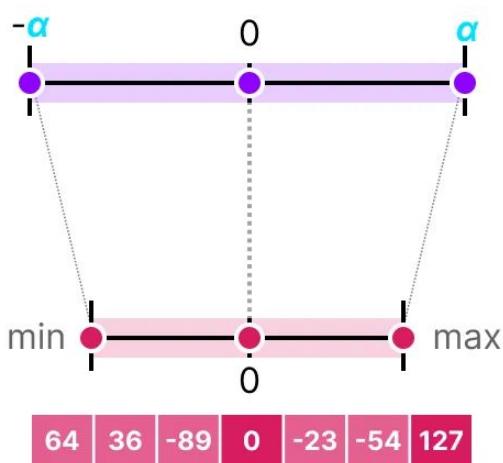
# Asymmetric quantization (cont'd)

$$x_{\text{dequantized}} = \frac{\text{[pink bar]} - z}{s} \quad (\text{dequantize})$$

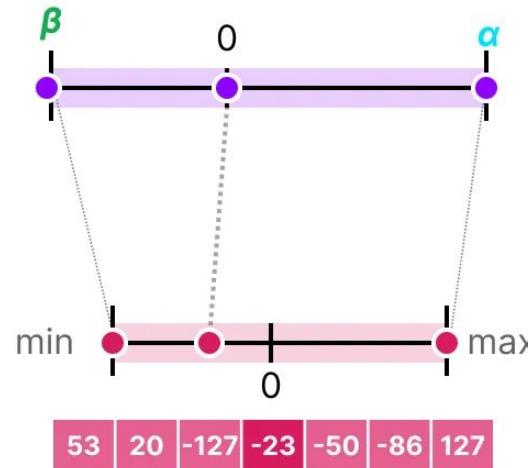
# Symmetric vs. Asymmetric quantization



**Symmetric**  
[-10.8, 10.8]



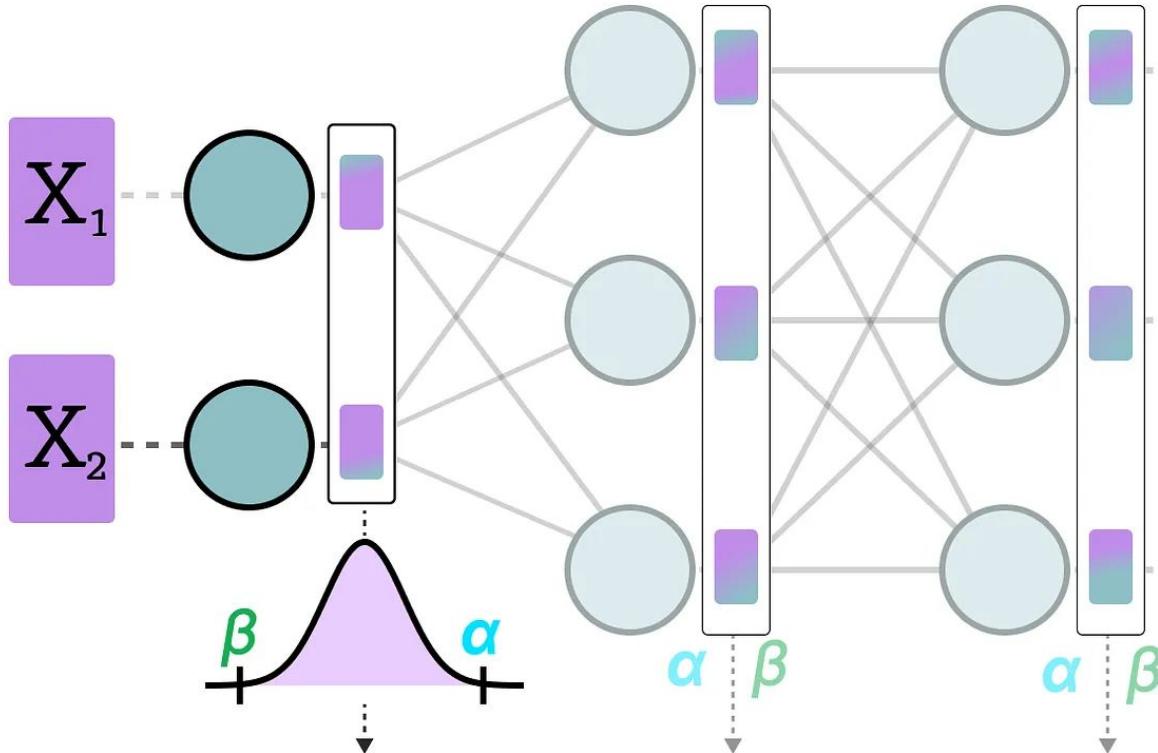
**Asymmetric**  
[-7.59, 10.8]



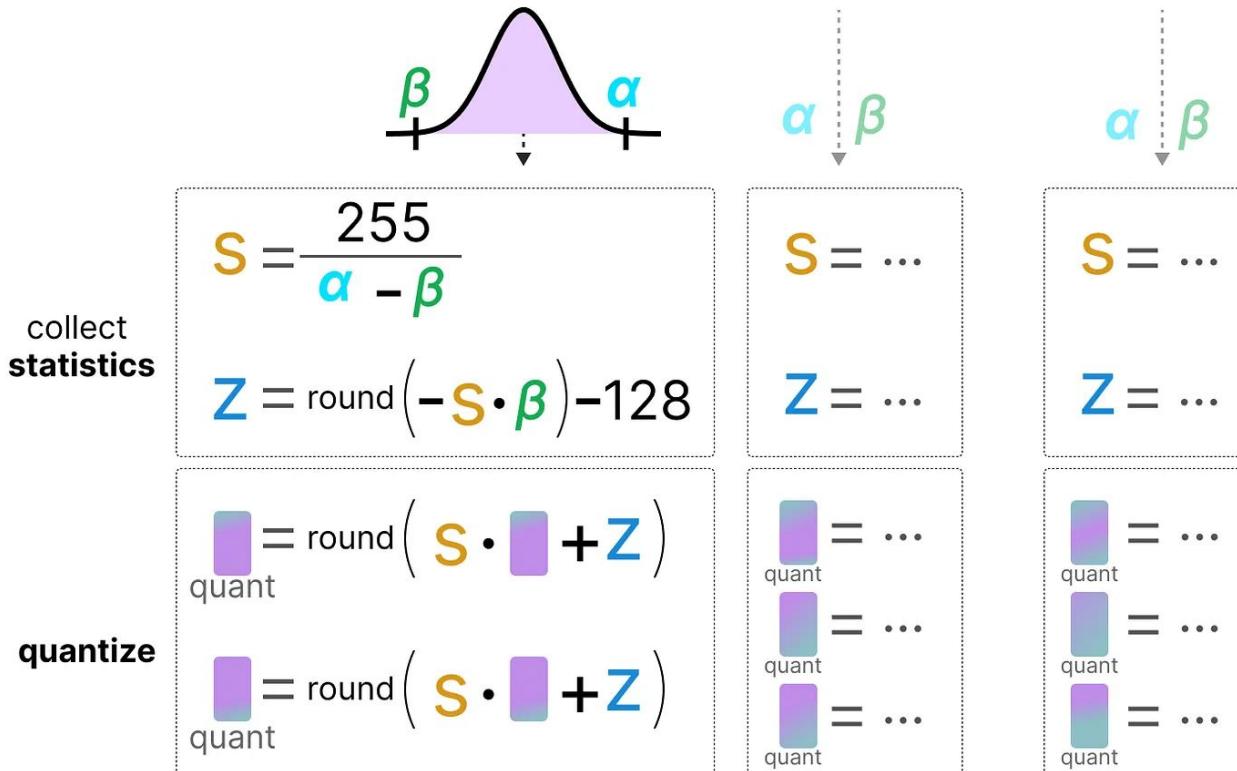
# Post-training quantization

- Dynamic Quantization
- Static Quantization

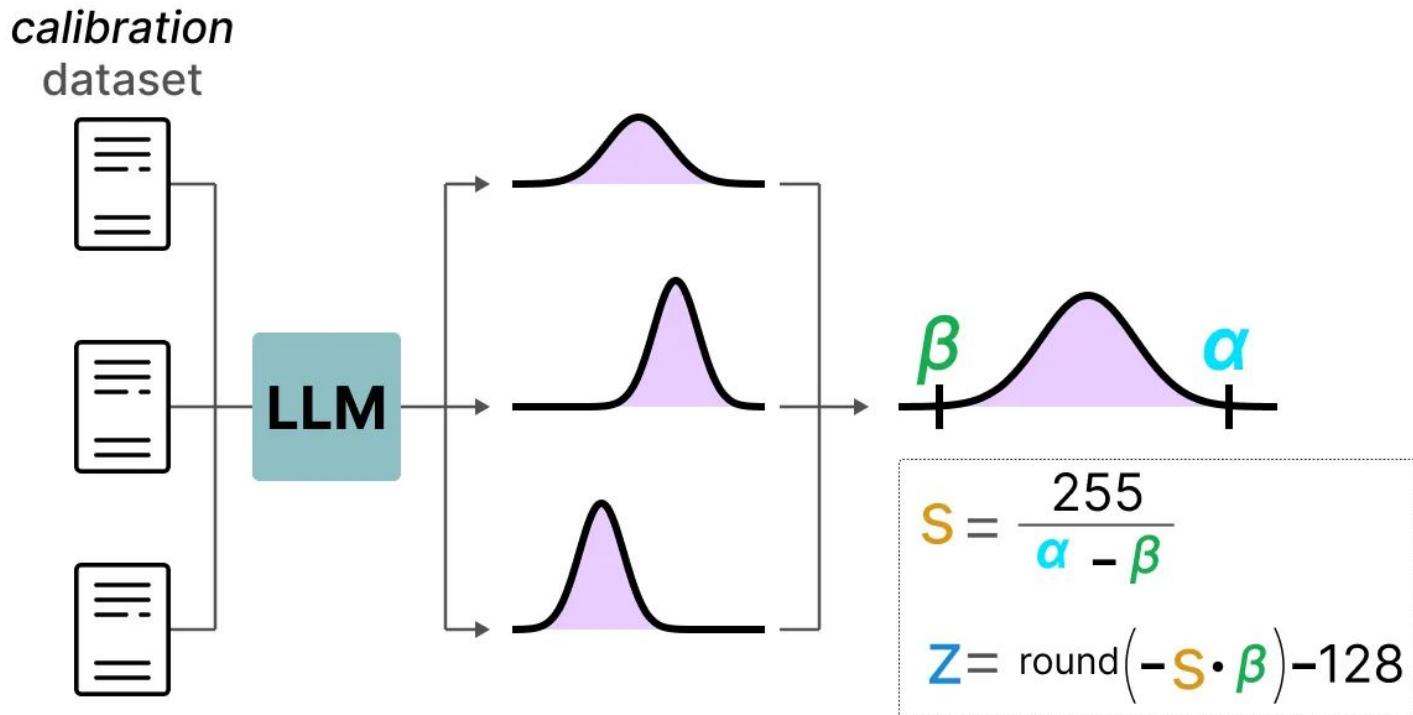
# Dynamic quantization



# Dynamic quantization (cont'd)



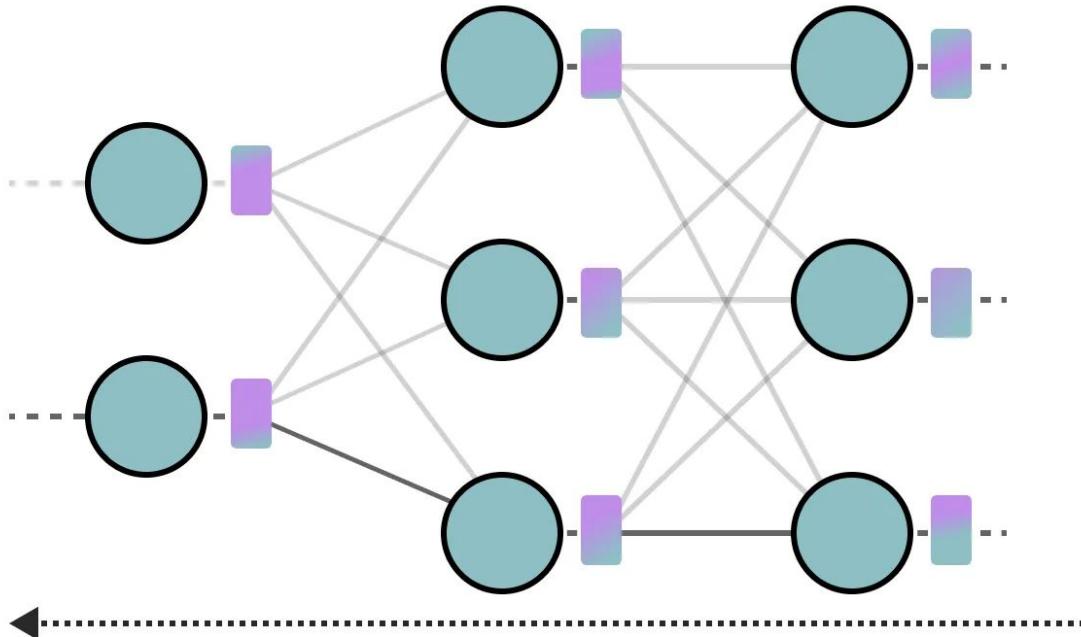
# Static quantization



# The realm of 4-bit quantization

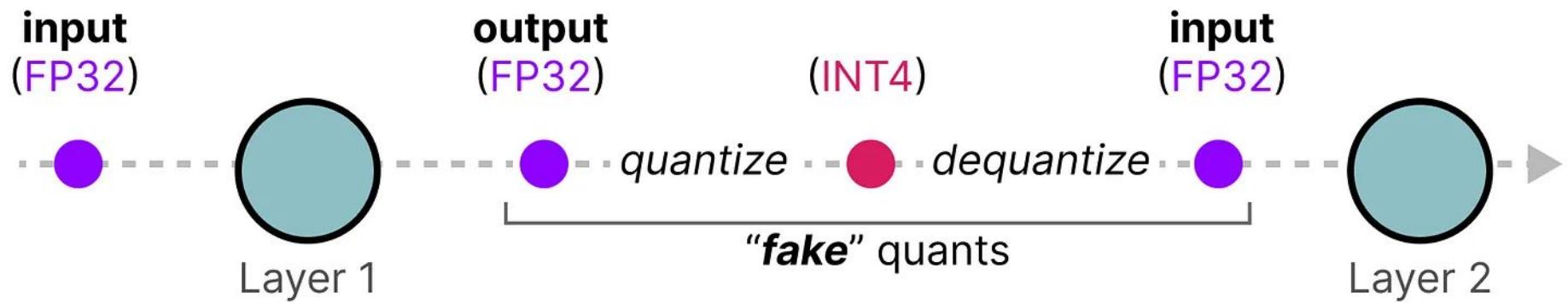
- GPTQ (full model on GPU)
- GGUF (potentially offload layers on the CPU)

# Quantization aware training



Learn **quantization parameters** ( $s$ ,  $\alpha$ ,  $\beta$ ,  $z$ )  
during **backward pass**

# Quantization aware training (cont'd)



---

# QLoRA: Efficient Finetuning of Quantized LLMs

---

**Tim Dettmers\***

**Artidoro Pagnoni\***

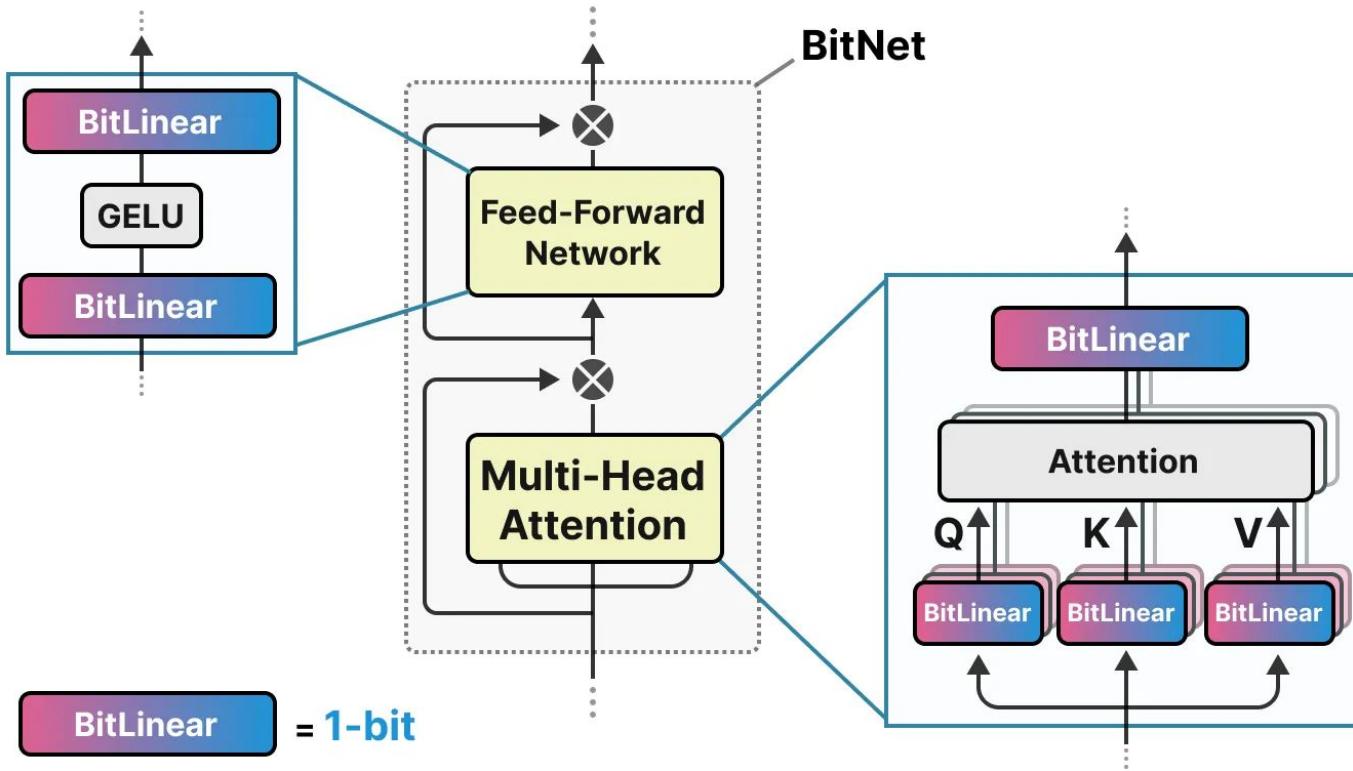
**Ari Holtzman**

**Luke Zettlemoyer**

University of Washington

{dettmers, artidoro, ahai, lsz}@cs.washington.edu

# The era of 1-bit LLMs: BitNet



**Thank you!**