

RLHF & “Small” LLMs

Instructor: Lara J. Martin (she/they)

TA: Omkar Kulkarni (he)

<https://laramartin.net/NLP-class/>

Slides modified from Yejin Choi, Bill Yuchen Lin, & Valentina Pyatkin

Learning Objectives

Describe what RLHF is and how it fits into alignment

Examining...

- Methods for shrinking pre-existing models
- Methods for mimicking pre-existing models with smaller models
- Methods for faster/smaller finetuning of pre-existing models
- Methods for training new models more efficiently

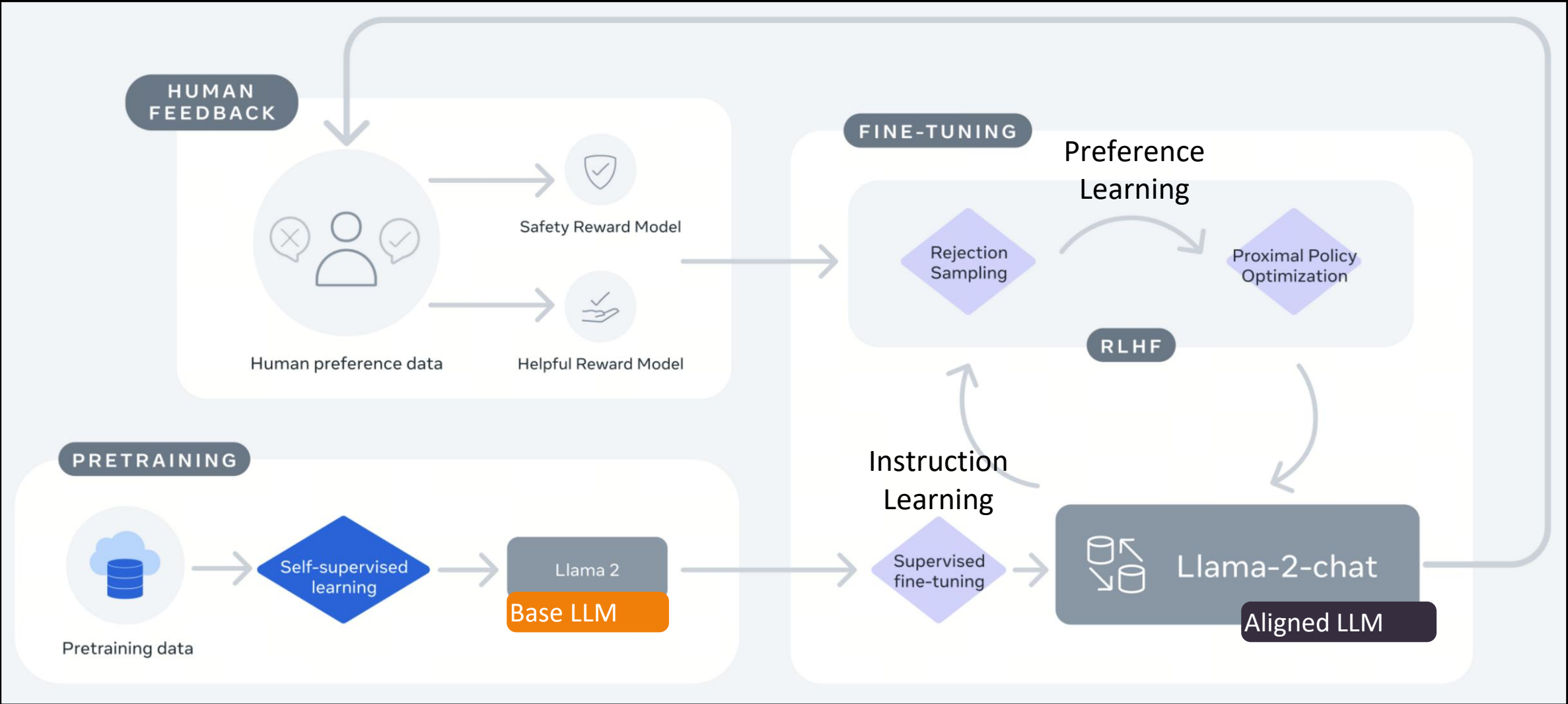
Finding where to implement these methods

Recognizing when to implement them

Review: What is Alignment?

AI concept for getting agents to match human values (value alignment), preferences, goals, etc.

Review: Llama-2's alignment



Review: Datasets for Instruction Learning

1. Synthetic Conversion
2. Human Annotation
3. Collected from ChatGPT/GPT-4
 - 3.1. Community Sharing
 - 3.2. Strategic Collecting

Review: Examples of Evaluations for Alignment

Benchmarking Datasets

Human Annotation

GPTs as Judges

Open LLM Evaluators

Safety Evaluation

Review: Helping out Instruction Tuning

Why do we need RLHF?

What makes one output better than the other? -> hard to define

What types of LM errors should be weighted more?

LM objective != human preferences

Review: “Learning to Summarize with Human Feedback”

1. Collect human feedback

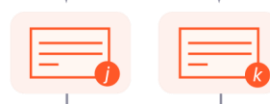
A Reddit post is sampled from the Reddit TL;DR dataset.



Various policies are used to sample N summaries.



Two summaries are selected for evaluation.



A human judges which is a better summary of the post.



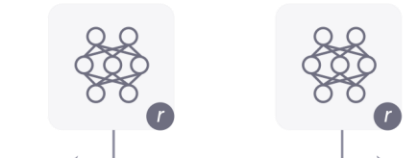
“j is better than k”

2. Train reward model

The post and summaries judged by the human are fed to the reward model.



The reward model calculates a reward r for each summary.



The loss is calculated based on the rewards and human label.



$$\text{loss} = \log(\sigma(r_j - r_k))$$

The loss is used to update the reward model.

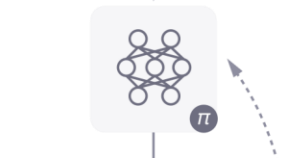
“j is better than k”

3. Train policy with PPO

A new post is sampled from the dataset.



The policy π generates a summary for the post.



The reward model calculates a reward for the summary.



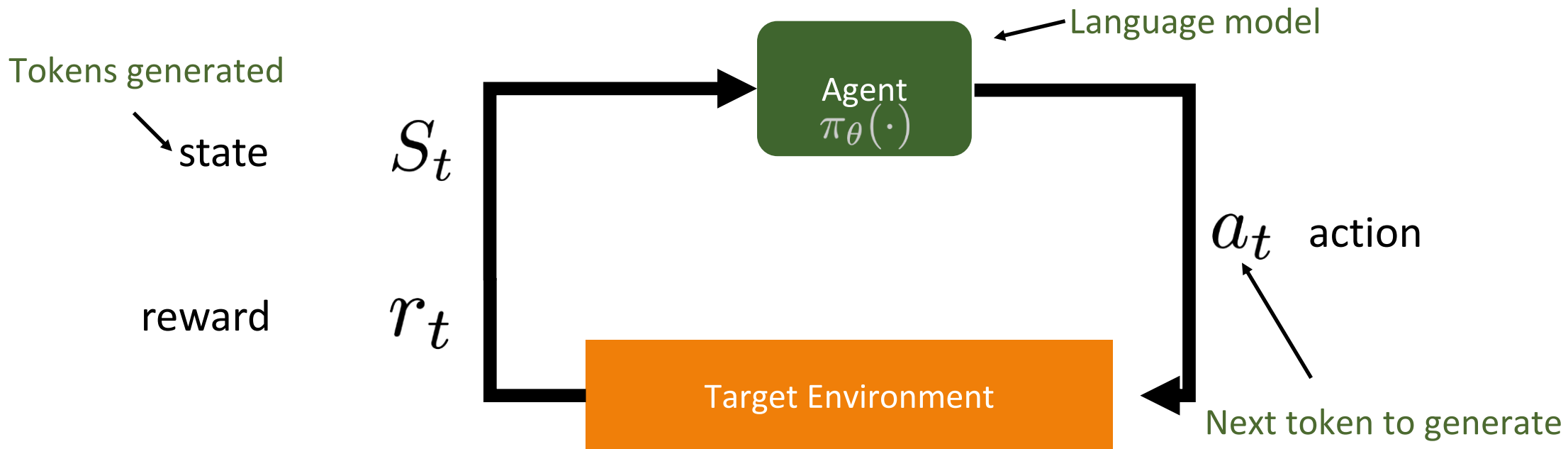
The reward is used to update the policy via PPO.



r_k

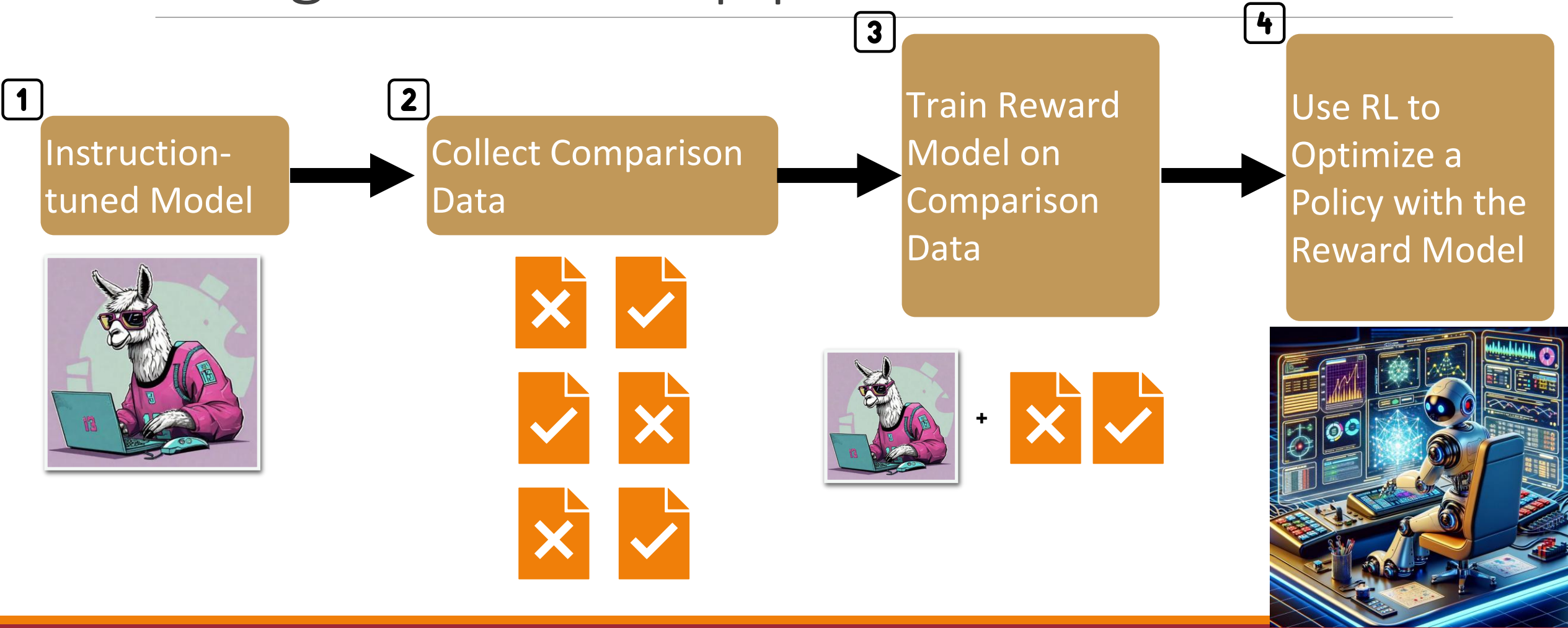
RL methods don't always assume “preference-based” (j is better than k) human feedback and reward model, but that's what's common with current “RLHF” approaches

Review: RL in the Context of Language Models...

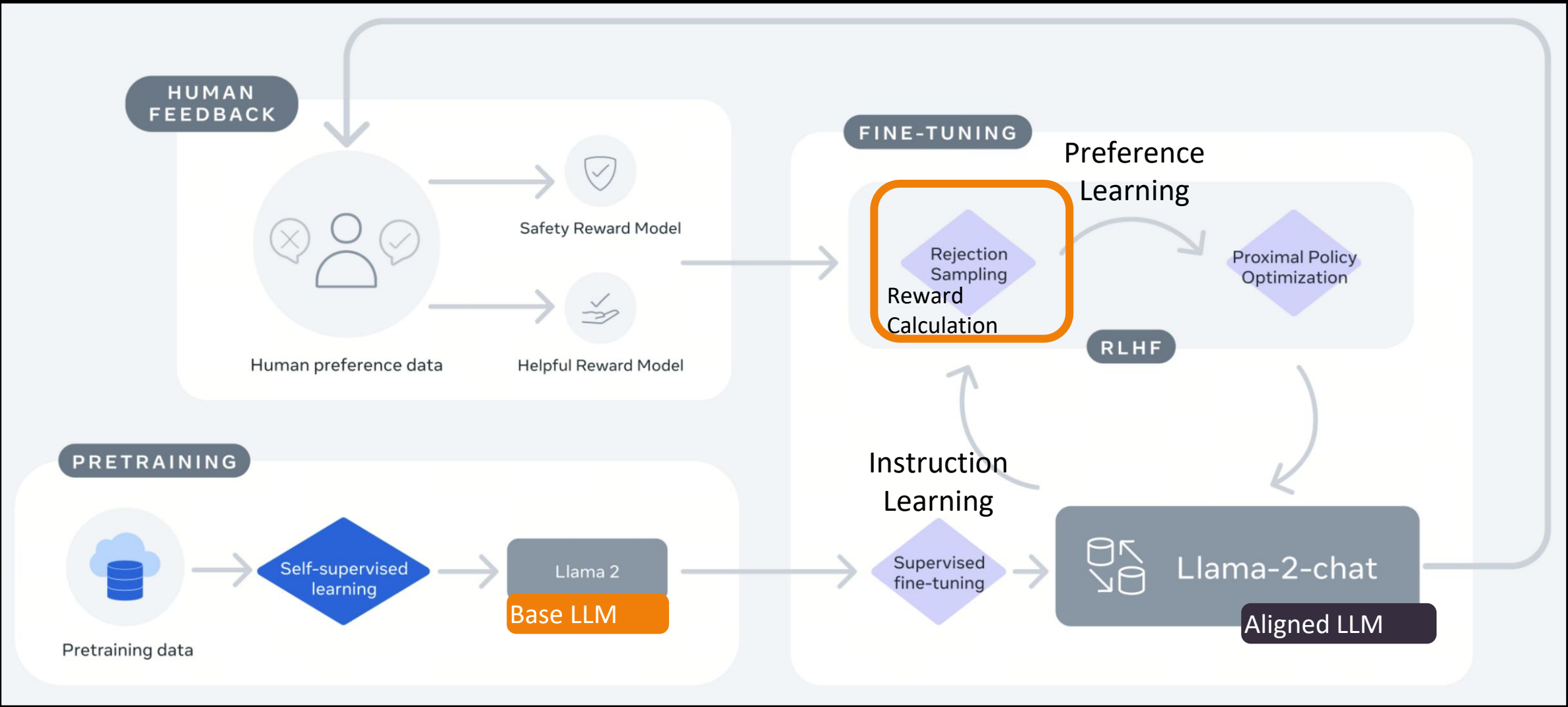


$$a_t \sim \pi_{\theta}(S_t) : \text{policy}$$

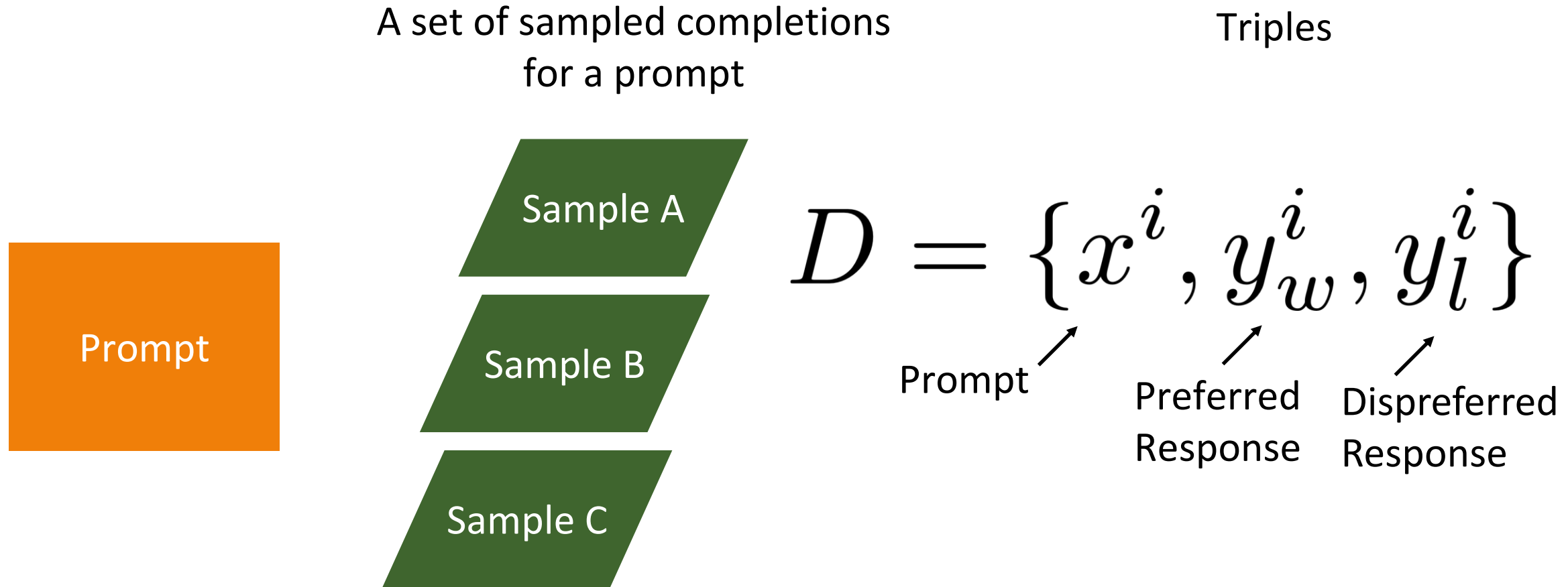
The general RLHF pipeline




Example: Llama-2's alignment




Where does the reward come from? Start with Human Preferences



 Can you help me write a resignation letter to my current employer, while leaving on good terms and expressing gratitude for the opportunities provided?

Example: Annotation

 Here are two responses from the chatbot. (Please scroll down on the content to see the entire response if it is too long)

A

Sure, here is a resignation letter template that you can use:

Dear [manager],

It is with a heavy heart that I submit this letter of resignation from my position as [position] at [company]. I have enjoyed my time here and am grateful for the opportunities that have been provided to me.

I have accepted a new opportunity that presents a greater challenge and is more aligned with my long-term career goals. I have learned so much during my time at

B

Of course! Here is a sample resignation letter that you can use as a guide:

Dear [Manager's Name],

I am writing to inform you that I have decided to resign from [Company Name]. This was not an easy decision, as I have enjoyed my time here and have learned a lot from my colleagues.

I want to express my sincere gratitude for the opportunities and support that [Company Name] has provided me over the years. I have grown both

Annotator needs to choose whether they prefer A or B.

$$D = \{x^i, y_w^i, y_l^i\}$$

Prompt ↙
Preferred ↙
Dispreferred ↙

Response
Response

Calculating reward

Bradley-Terry Model

Reward for preferred response

Reward for dispreferred response

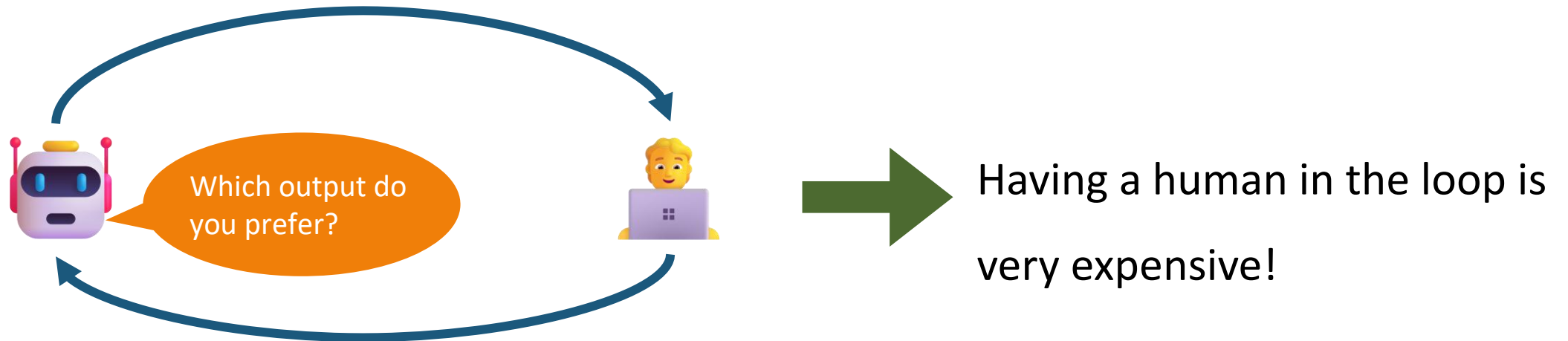
$$p(y_w > y_l | x) = \sigma(\underbrace{r(x, y_w)}_{\text{Reward for preferred response}} - \underbrace{r(x, y_l)}_{\text{Reward for dispreferred response}})$$

Logistic function; which is equivalent to using softmax:

$$\frac{1}{1 + e^{-x}}$$

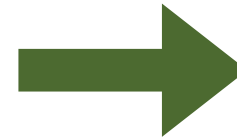
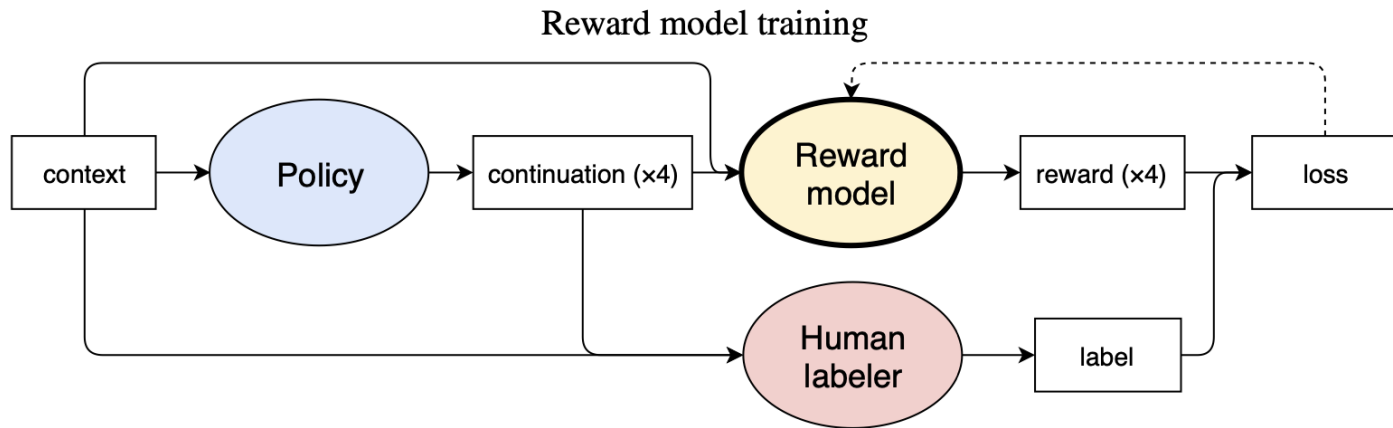
$$p(y_w > y_l | x) = \frac{\exp(r(x, y_w))}{\exp(r(x, y_w)) + \exp(r(x, y_l))}$$

Making things cheaper...



How can we more cheaply get the reward?

Reward Models



Instead: train a Reward Model (RM) on preference data to predict preferences!

Ziegler et al., 2019 "Fine-Tuning Language Models from Human Preferences"

Reward Modeling

Train on preference data.

Minimizing negative log likelihood.



Bradley-Terry Model



$$\mathcal{L}_R(\phi, D) = -\mathbb{E}_{(x, y_w, y_l) \sim D} [\log \sigma(r(x, y_w) - r(x, y_l))]$$


Train an LLM with an additional layer to minimize the neg. log likelihood

Fun Facts about Reward Models

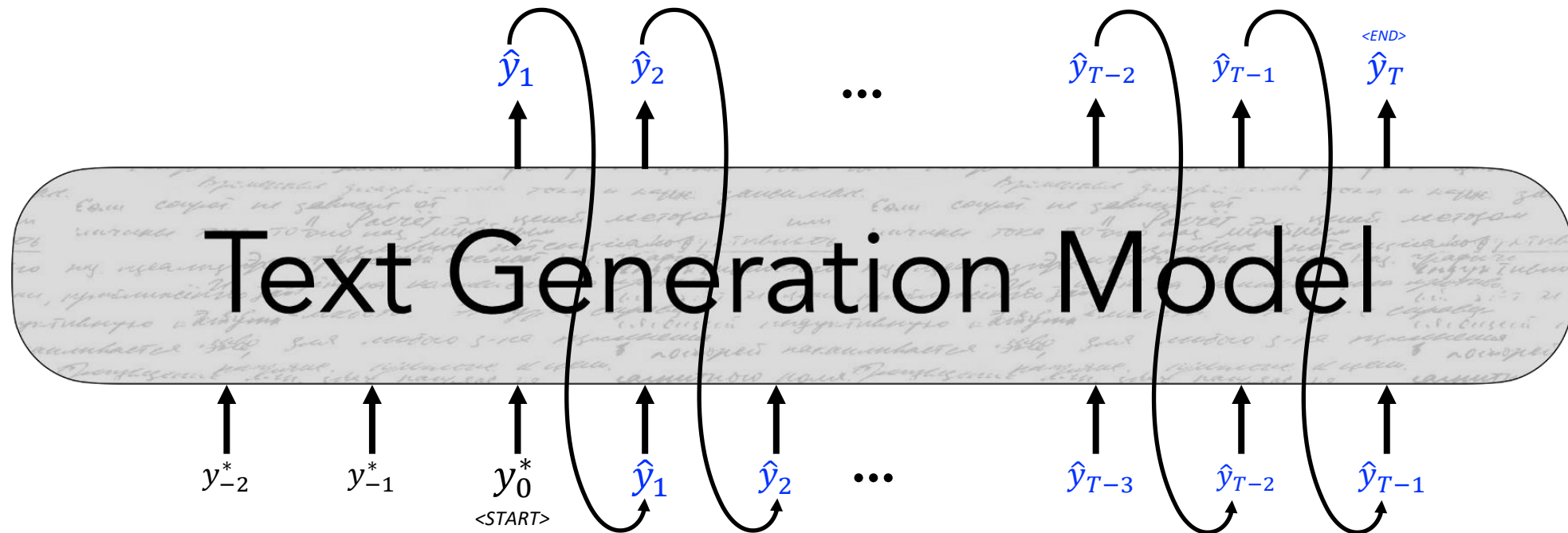
Trained for 1 epoch (to avoid overfitting)!

Evaluation often only has 65% - 75% agreement

Lambert et al., 2023

Review: REINFORCE

Sample a sequence from your model, score the sequence, and use the score to train the model.



REINFORCE

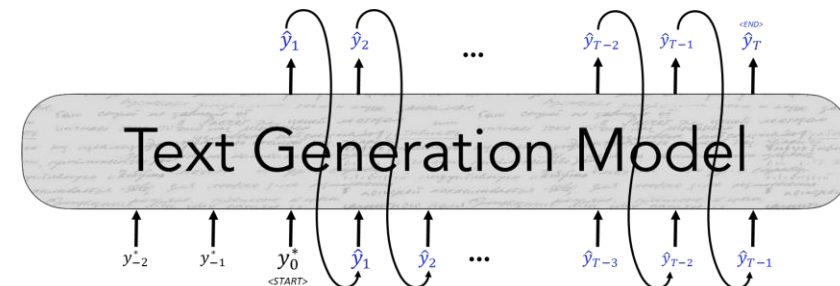
- Sample a sequence from your model, score the sequence, and use the score to train the model.

Next time, increase the probability of this sampled token in the same context.

$$L_{RL} = - \sum_{t=1}^T \underline{r(\hat{y}_t)} \log P(\hat{y}_t | \{y^*\}; \{\hat{y}\}_{<t})$$

... but increase it more if I get a higher reward from the reward function.

- $r(\cdot)$: Your reward model
- y^* : Input sequence given to the model
- \hat{y} : The sequence sampled from the model given y^*



Updating the weights using REINFORCE

$$\theta_{t+1} := \theta_t + \alpha \frac{1}{m} \sum_{i=1}^m R(S_i) \nabla_{\theta_t} \log p_{\theta_t}(S_i)$$

Simplified Intuition: good actions are reinforced and bad actions are discouraged.

Williams, 1992

Updating the weights using REINFORCE

$$\theta_{t+1} := \theta_t + \alpha \frac{1}{m} \sum_{i=1}^m R(S_i) \nabla_{\theta_t} \log p_{\theta_t}(S_i)$$

If: Reward is high/positive

Then: maximize this

Simplified Intuition: good actions are reinforced and bad actions are discouraged

Williams, 1992

Updating the weights using REINFORCE

$$\theta_{t+1} := \theta_t + \alpha \frac{1}{m} \sum_{i=1}^m R(S_i) \nabla_{\theta_t} \log p_{\theta_t}(S_i)$$

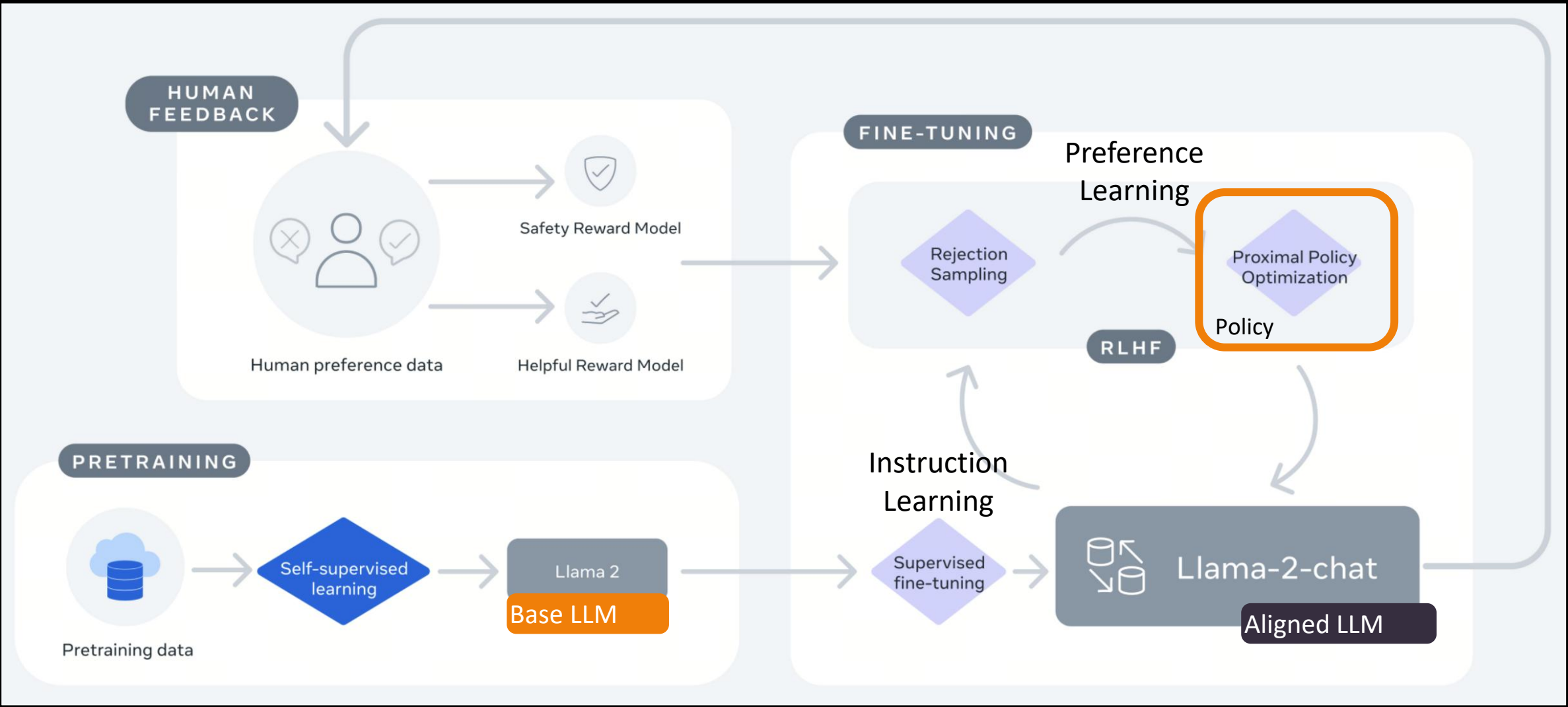
If: Reward is negative/low

Then: minimize this

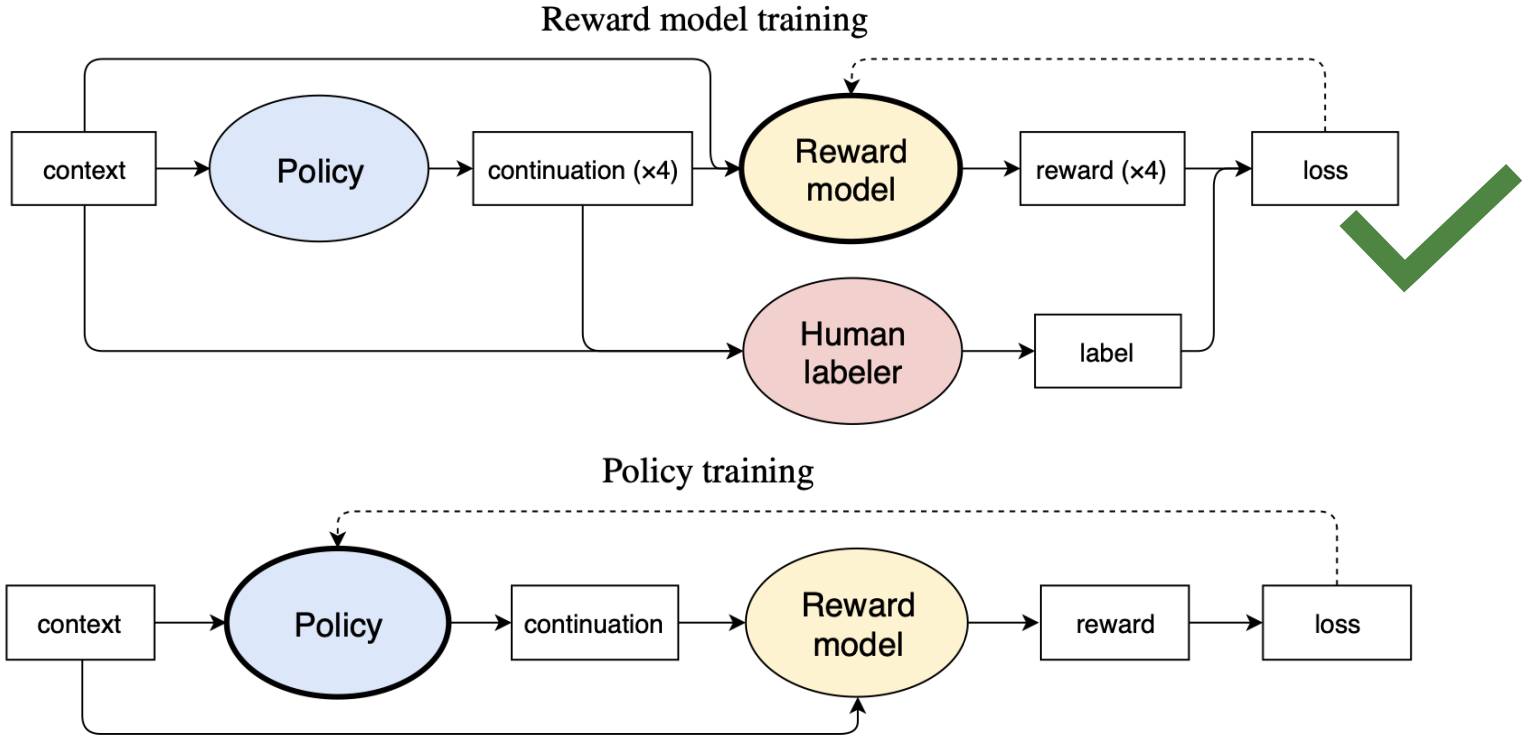
Simplified Intuition: good actions are reinforced and bad actions are discouraged

Williams, 1992

Example: Llama-2's alignment



Policy Training with Reward Model



Ziegler et al., 2019 "Fine-Tuning Language Models from Human Preferences"

Policy

We have: Reward Model

Next step: learn a **policy** to maximize the reward (minus KL regularization term) using the reward model

$$\max_{\pi_{\theta}} \mathbb{E}_{x \sim D, y \sim \pi_{\theta}(y|x)} [\underline{r_{\phi}(x, y)}] - \beta \underline{\mathbb{D}_{KL}[\pi_{\theta}(y|x) || \pi_{ref}(y|x)]}$$

Sampling from policy

Reward given prompt
and sampled generation



Should be high!

KL-divergence between original model's
generation and the sampled generation



Should be low!

PPO

Proximal Policy Optimization

Proximal Policy Optimization Algorithms

John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, Oleg Klimov
OpenAI
{joschu, filip, prafulla, alec, oleg}@openai.com

arxiv in July 2017

PPO: builds on Policy Gradient Methods

Gradient Estimator

$$\hat{g} = \hat{\mathbb{E}}_t[\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \hat{A}_t]$$

↓

Expectation: empirical average over a finite batch of samples

Objective / Loss:

$$L^{PG}(\theta) = \hat{\mathbb{E}}_t[\log \pi_{\theta}(a_t | s_t) \hat{A}_t]$$

➔ Often leads to (too) large policy updated

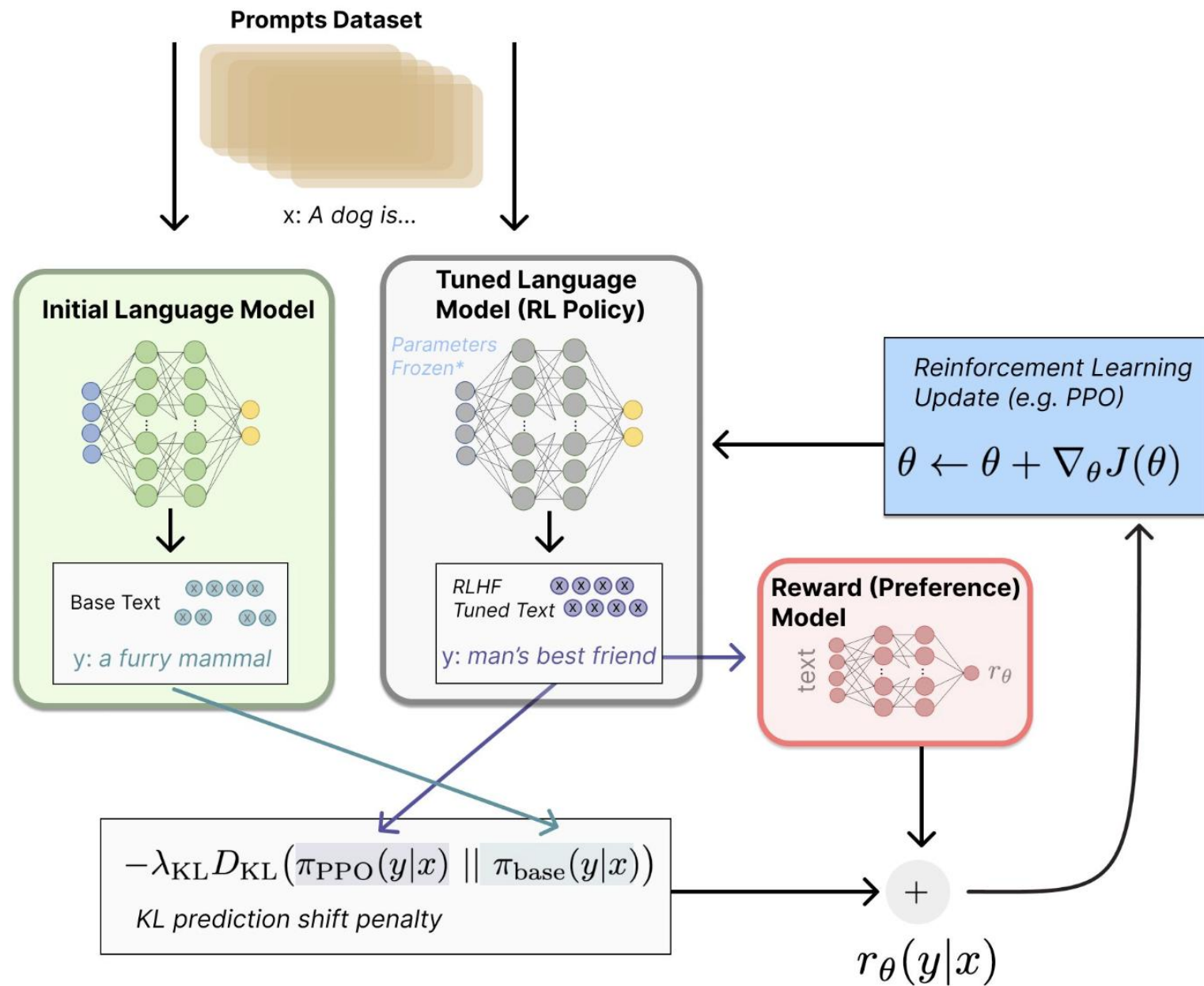
\hat{A}_t : estimator of the advantage function at timestep t (critic in actor-critic training)

π_{θ} : policy that we are trying to learn via PPO; this is initialized as a language model

$$\hat{A}_t = \hat{A}(s_t, a_t) = -V_{\phi}(t) + G_t = -V_{\phi}(t) + \sum_{t'=t}^T \gamma^{t'-t} r_{t'},$$

Schulman, 2017

PPO



Lambert, 2023

Evaluating the Learned Policy

Win Rate: How often does my policy's output win against a reference model's output, given the same instruction?

- Who compares the two outputs?
 - Humans
 - Simulated humans (and human variability!) using GPT-4 (**e.g.**, AlpacaFarm eval)

Dubois et al., 2023

Small LLMs

Efficient LLMs

Methods for shrinking pre-existing models

Methods for mimicking pre-existing models with smaller models

Methods for faster/smaller finetuning of pre-existing models

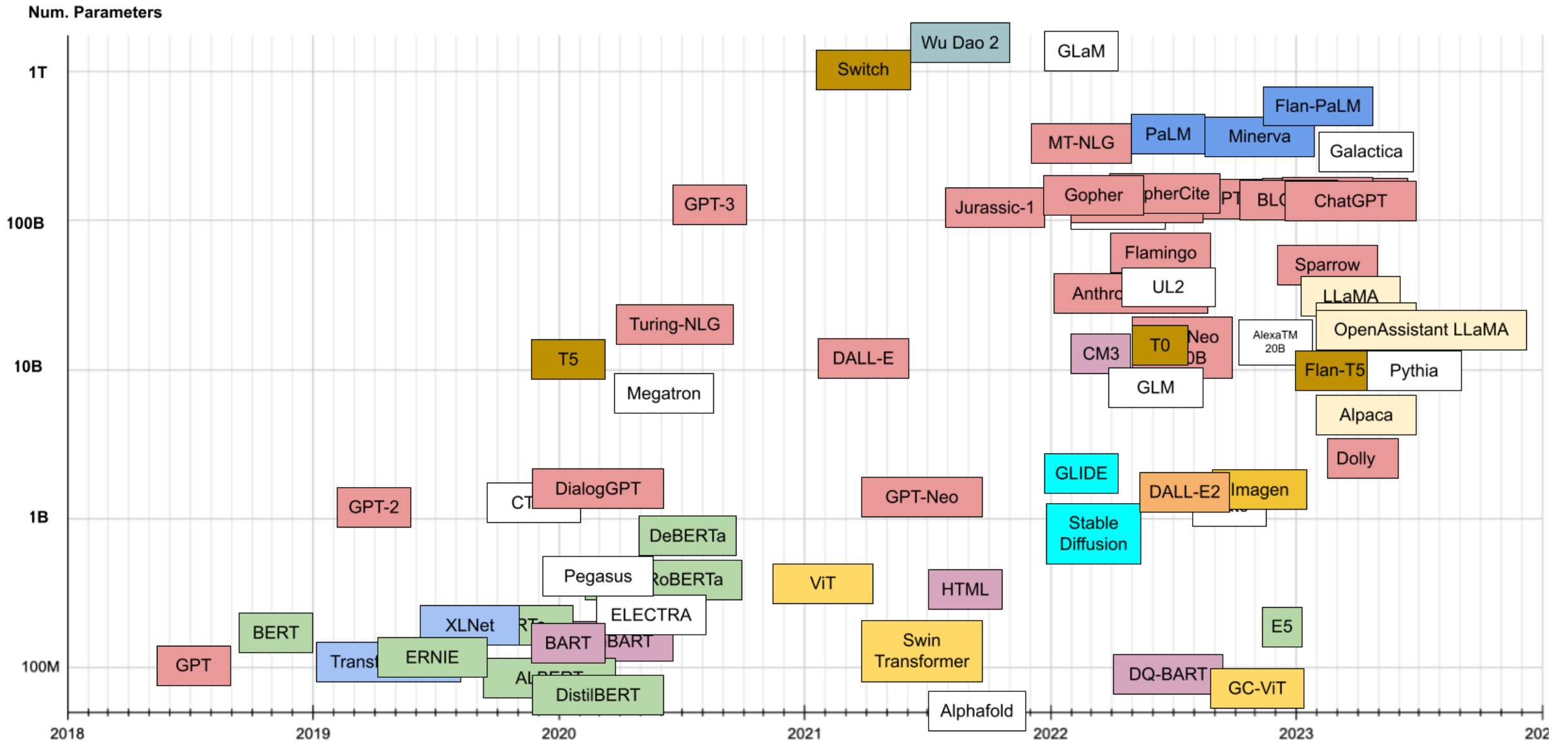
Methods for training new models more efficiently

Model Compression

Knowledge Distillation

PEFT

Efficient Training
e.g., Sparse Mixture of Experts



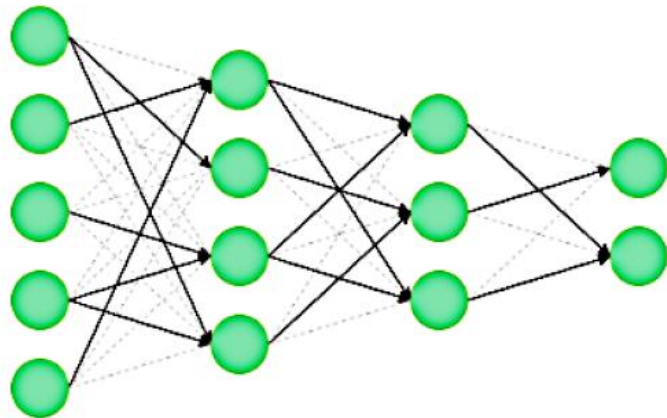
<https://amatriain.net/blog/transformer-models-an-introduction-and-catalog-2d1e9039f376/>

Model Compression

Pruning

Remove parts of the model

Unstructured Pruning



Structured Pruning

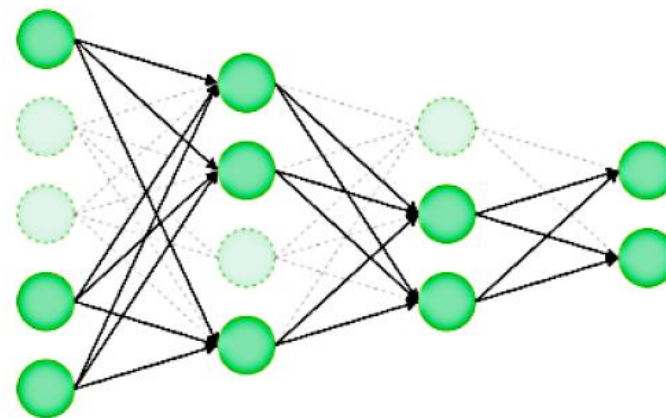
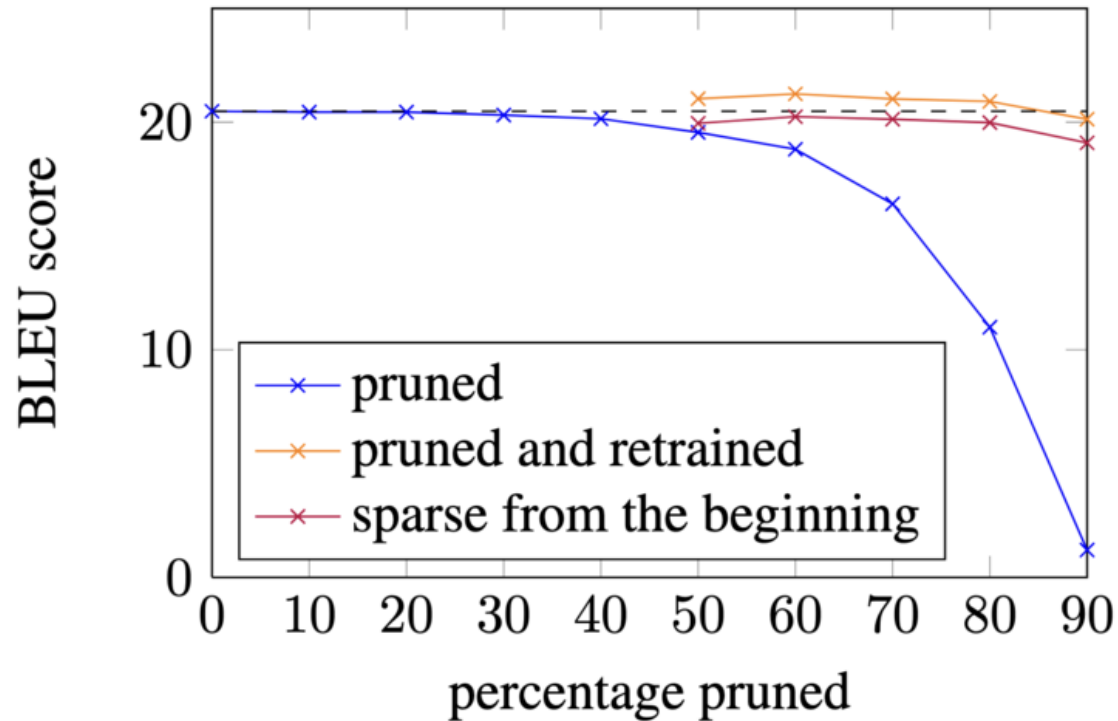


Image credits: neuralmagic.com

Implementation Tutorial: https://pytorch.org/tutorials/intermediate/pruning_tutorial.html

Magnitude Pruning



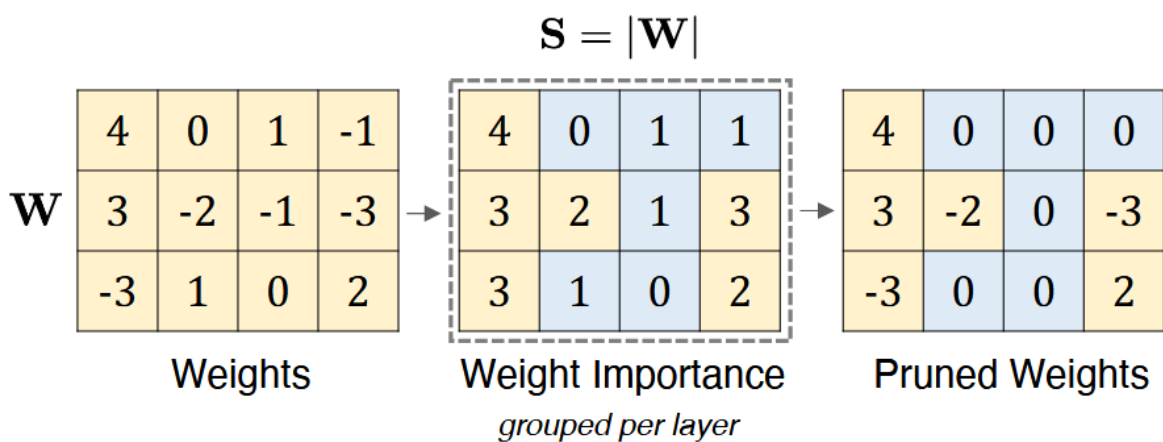
- prune weights with smallest absolute value
- prunes 40% of the weights with negligible performance loss
- by adding a retraining phase after pruning, we can prune 80% with no performance loss

Image credits: [See et al. 2016](#)

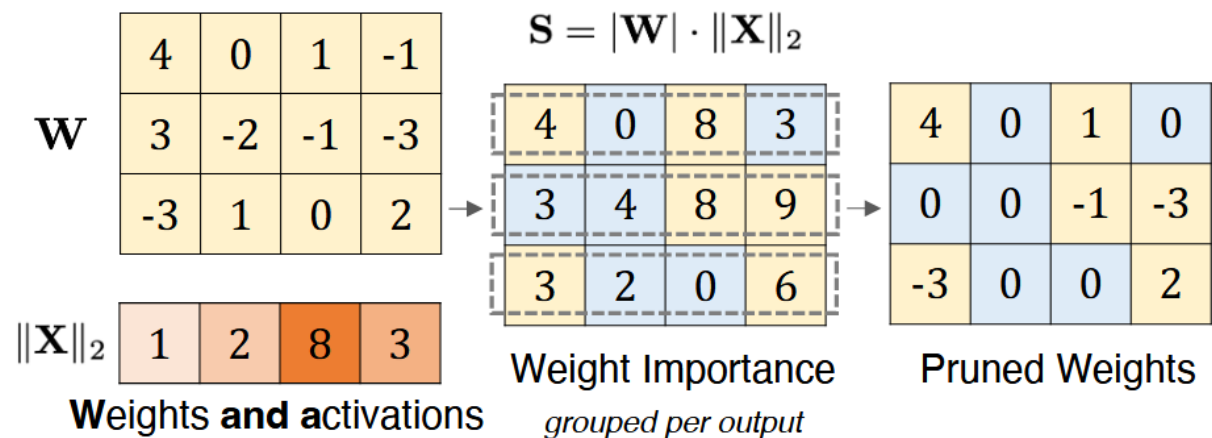
Slide by Dinesh Raghu

Wanda

Magnitude Pruning



Wanda



Quantizing Models

Compresses weights and activations from floating point numbers to integers (e.g., 4-bit, 8-bit)

Then use the scaling factor to get the “original” value

Implementation:

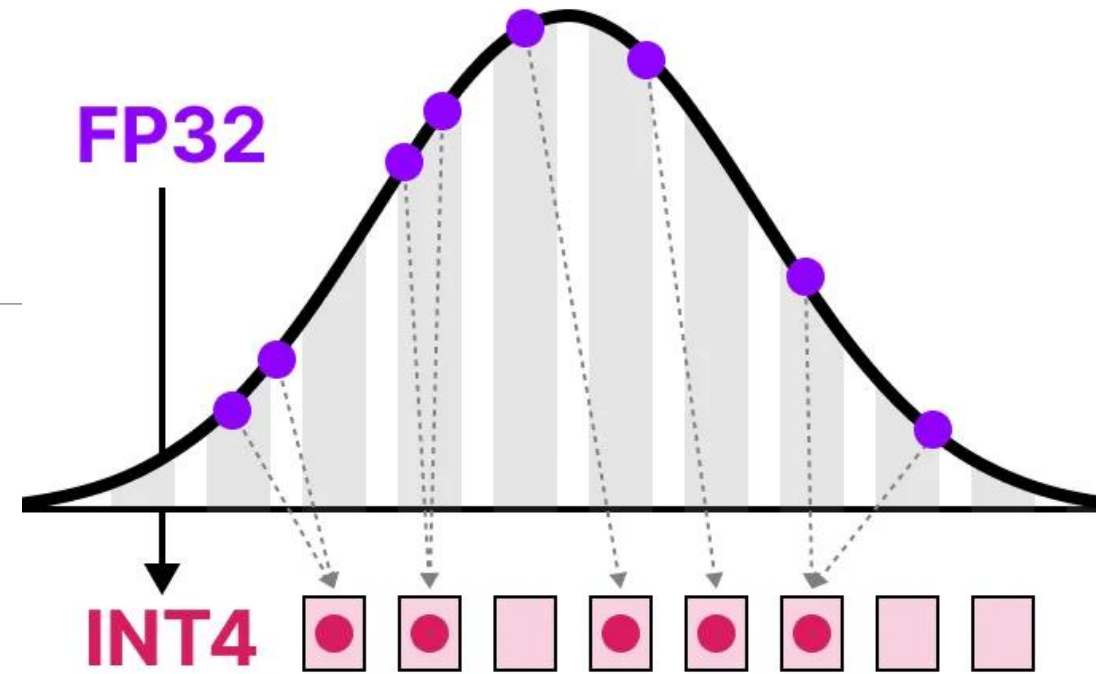
<https://pytorch.org/docs/stable/quantization.html>

<https://pypi.org/project/bitsandbytes/>

Learn more here:

<https://huggingface.co/blog/hf-bitsandbytes-integration>

Quantizing Models



Given a layer l with weight matrix W_l and layer input X_l , find quantized weight \hat{W}_l :

$$(\hat{W}_l^* = \operatorname{argmin}_{\hat{W}_l} \|W_l X - \hat{W}_l X\|_2^2)$$

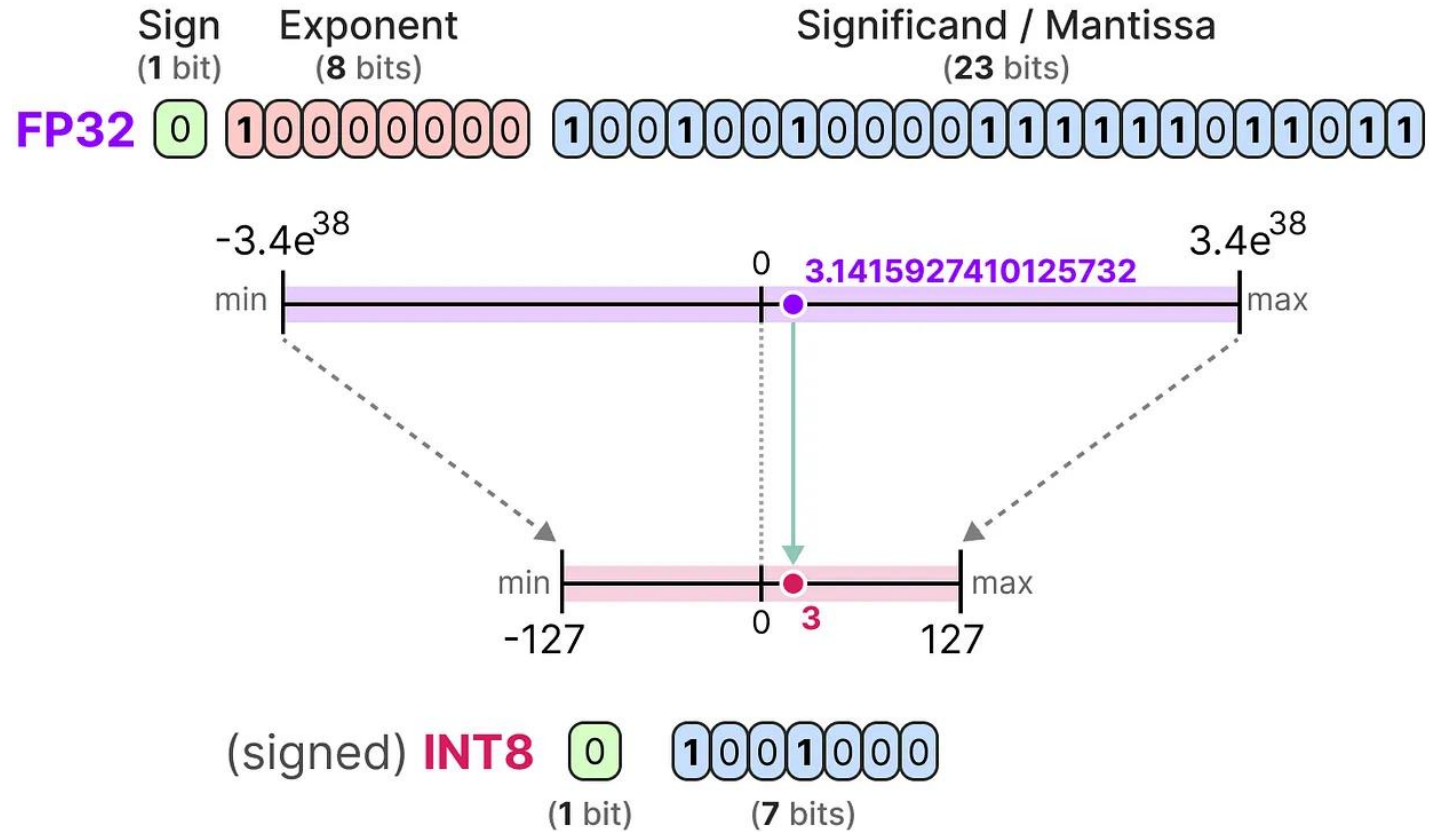
<https://huggingface.co/blog/merve/quantization>

<https://substack.com/home/post/p-145531349>

$$\mathbf{X}^{\text{Int8}} = \operatorname{round} \left(\frac{127}{\operatorname{absmax}(\mathbf{X}^{\text{FP32}})} \mathbf{X}^{\text{FP32}} \right) = \operatorname{round}(c^{\text{FP32}} \cdot \mathbf{X}^{\text{FP32}}), \text{ where } c \text{ is the quantization constant}$$

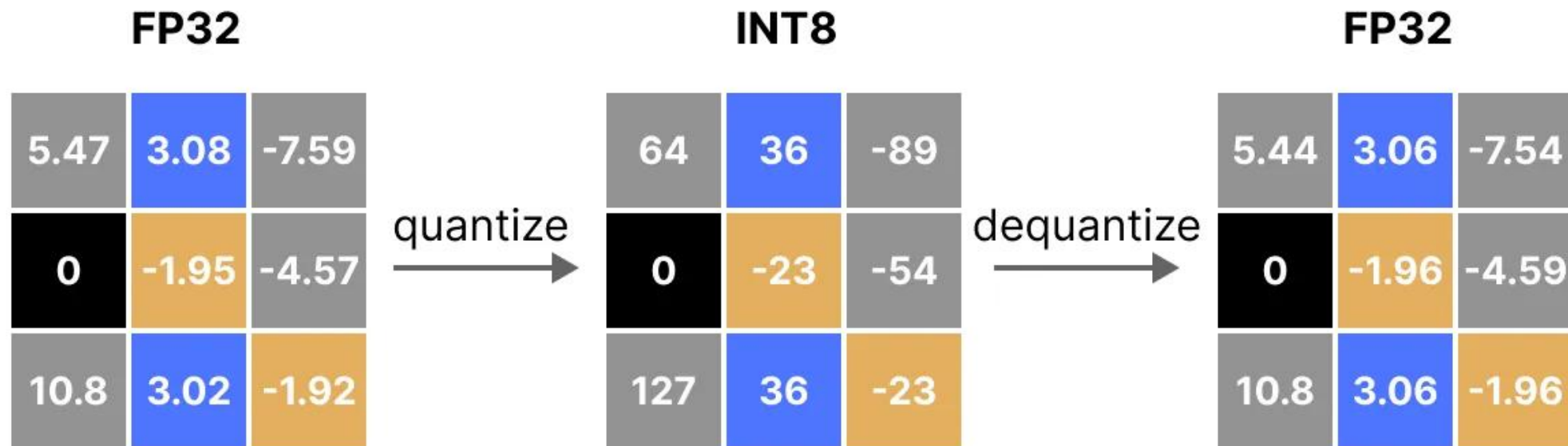
Dettmers, Tim, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. "QLoRA: Efficient Finetuning of Quantized LLMs." *Advances in Neural Information Processing Systems*

Mapping floating point to integer



<https://substack.com/home/post/p-145531349>

Quantizing → Dequantizing



Quantizing/Pruning (Model Compression)

PROS

More efficient

Less resources

Can add more parameters (if quantizing)

Model interpretability (if pruning)

CONS

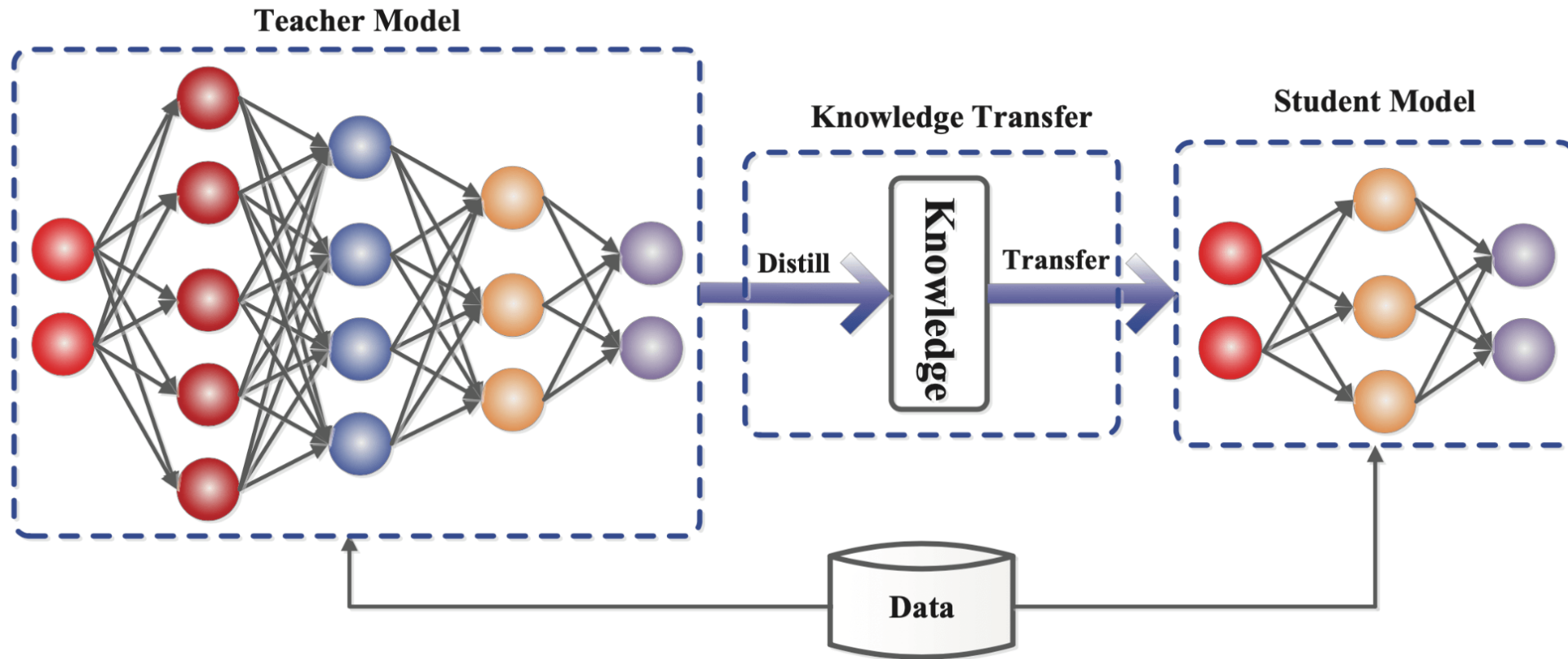
Losing information (lossy)

Possibly losing generalization

Adds a bit of time to dequantize?

Knowledge Distillation

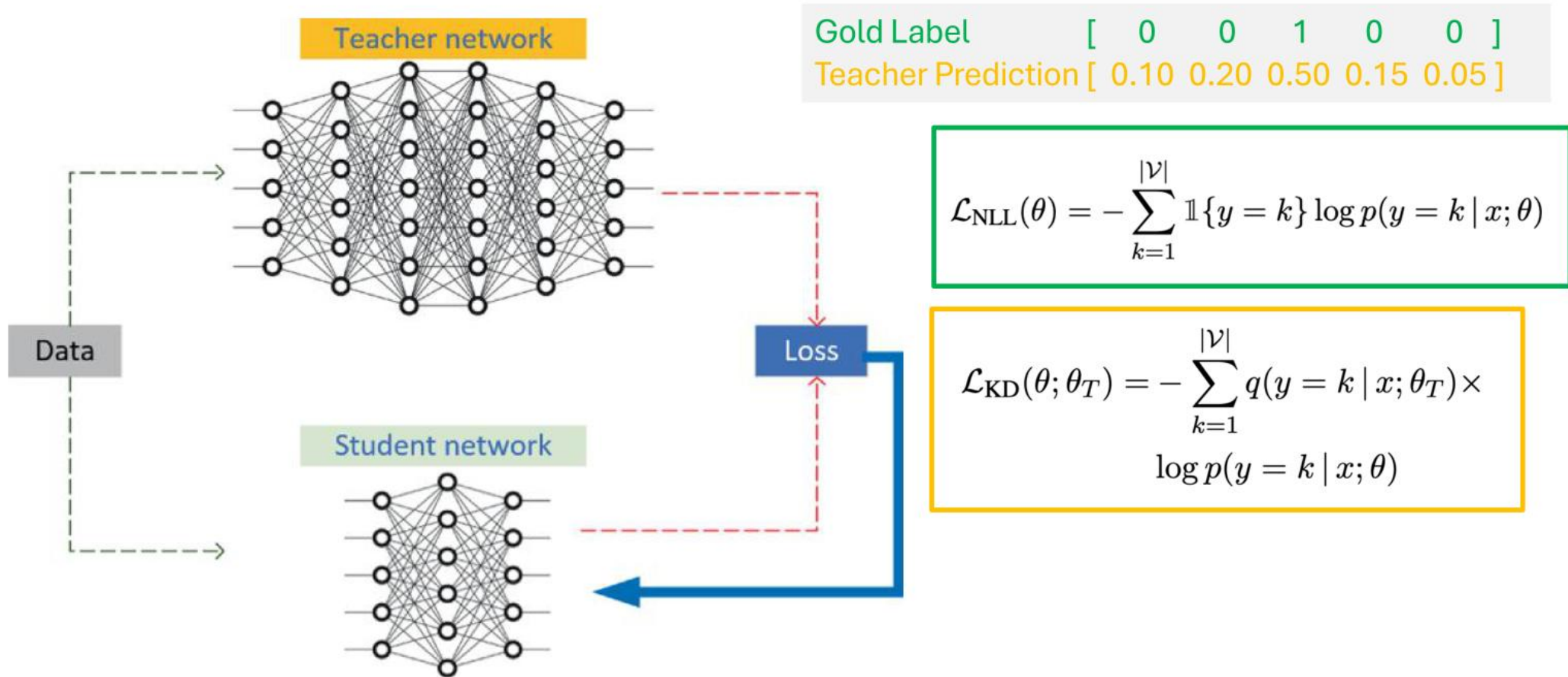
Knowledge Distillation



<https://neptune.ai/blog/knowledge-distillation>

Implementation Tutorial: https://pytorch.org/tutorials/beginner/knowledge_distillation_tutorial.html

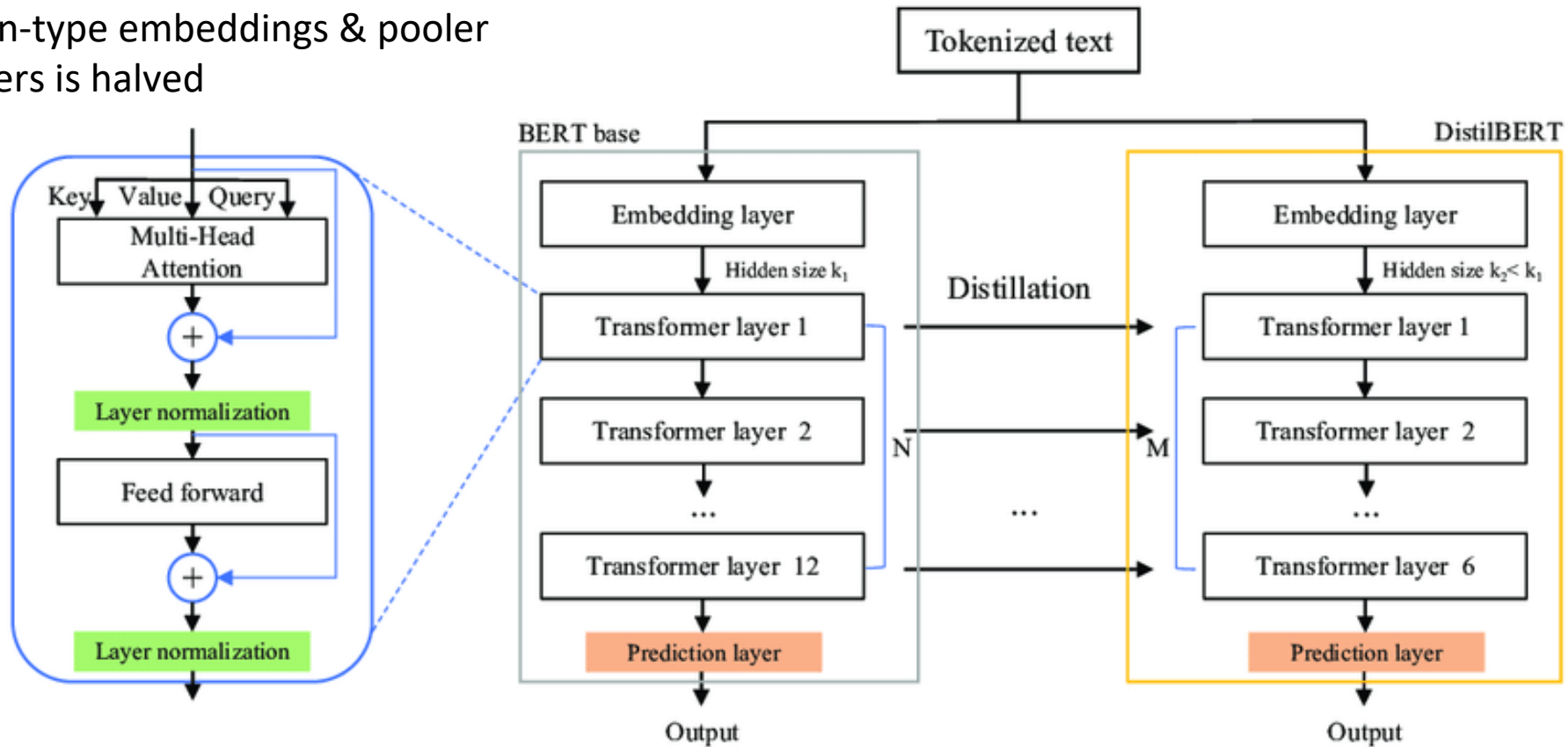
Training the Student Network



Slide by Dinesh Raghu

DistilBERT

Removed token-type embeddings & pooler
Number of layers is halved



https://www.researchgate.net/figure/The-DistilBERT-model-architecture-and-components_fig2_358239462

Knowledge Distillation

PROS

- Similar output but smaller architecture
- Simple to implement
- Flexibility in new model architecture

CONS

- Time consuming to train both models
- Time consuming to run through both models to train student
- Need to decide what the new architecture is
- Student picks up on teacher's biases
- Student isn't as generalizable
- Stealing?

PEFT

Parameter-efficient Fine-tuning (PEFT)

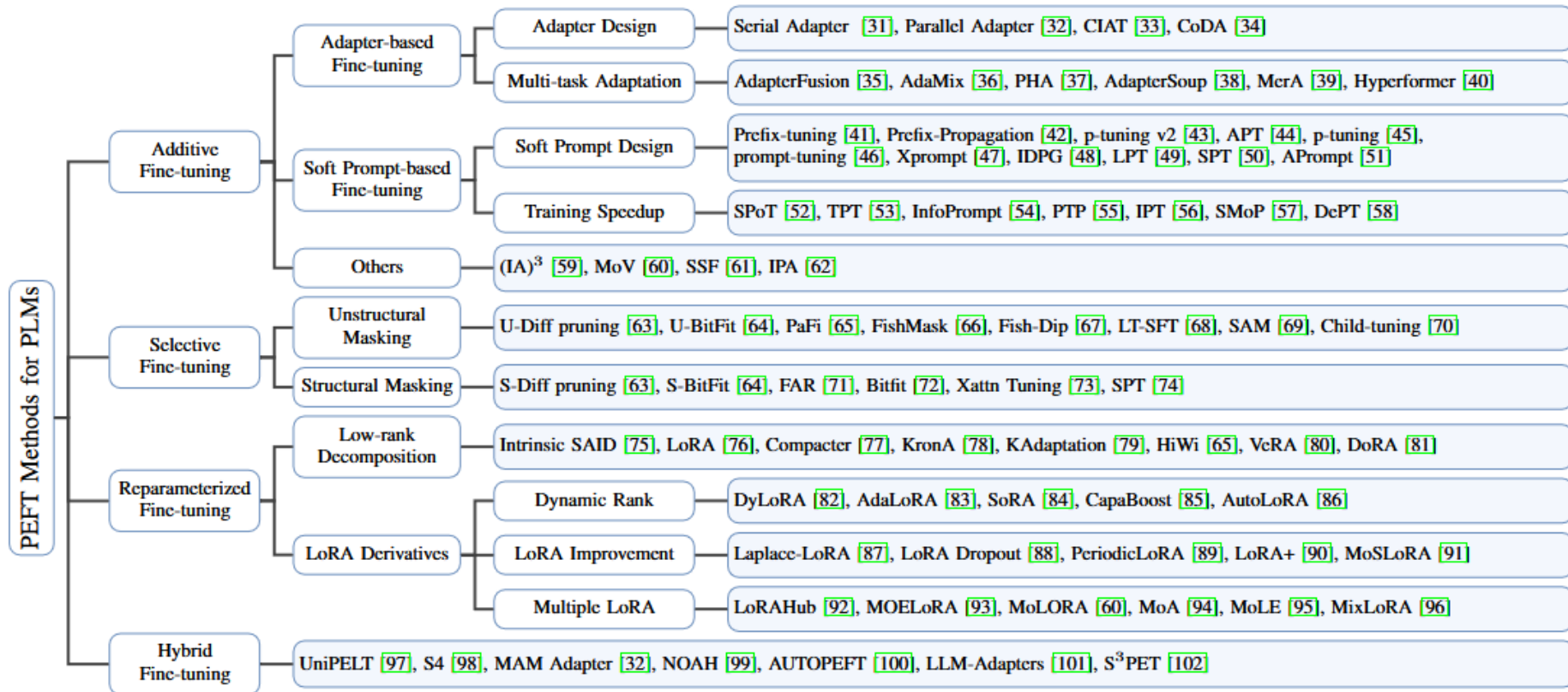


Fig. 3: Taxonomy of Parameter-Efficient Fine-Tuning Methods for Large Models.

LoRA (Low-Rank Adaptation)

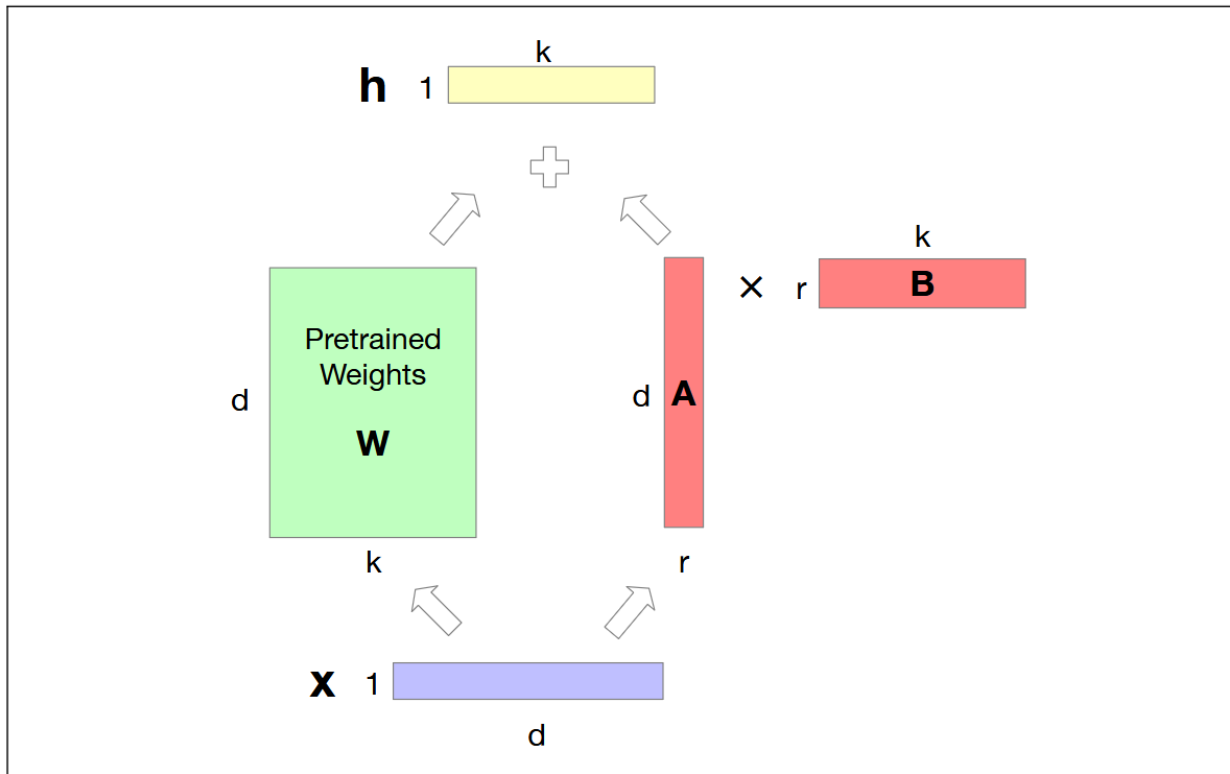


Figure 10.8 The intuition of LoRA. We freeze W to its pretrained values, and instead fine-tune by training a pair of matrices A and B , updating those instead of W , and just sum W and the updated AB .

Train a model using a pretrained LLM but give the new model fewer parameters \rightarrow a low-rank decomposition of the original weight matrix

$$h = xW + xAB$$
$$r \ll \min(d, N)$$

From SLP book Chapter 10

LoRA (Low-Rank Adaptation)

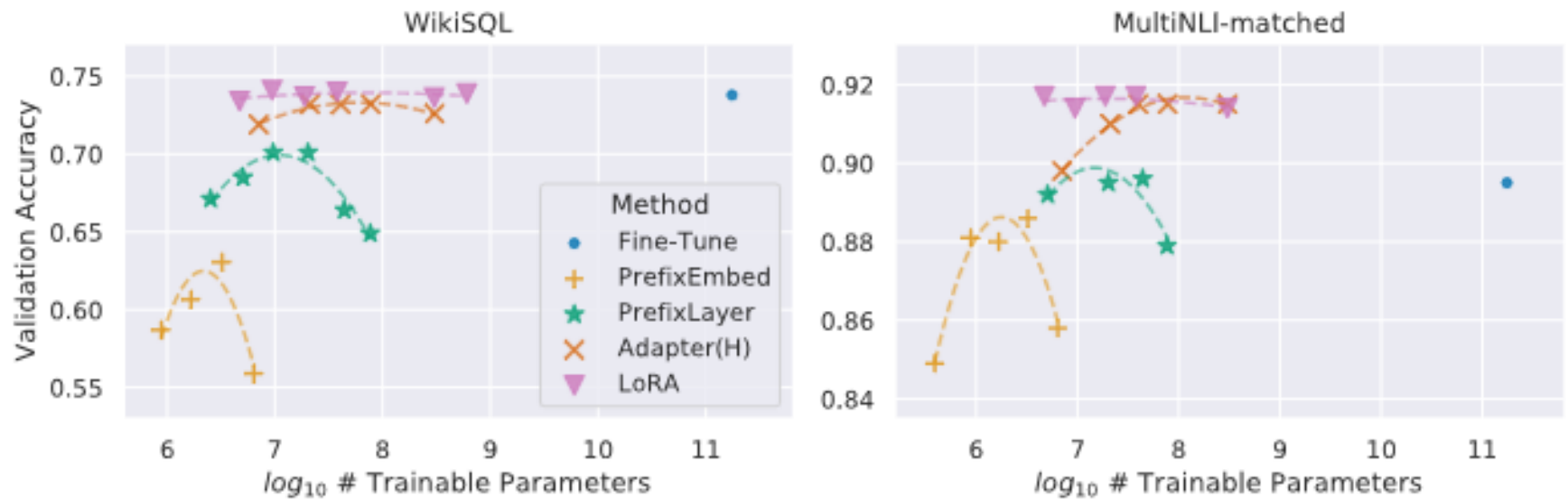
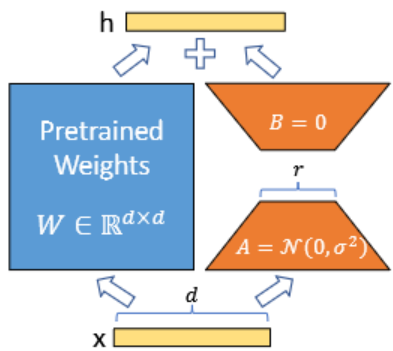


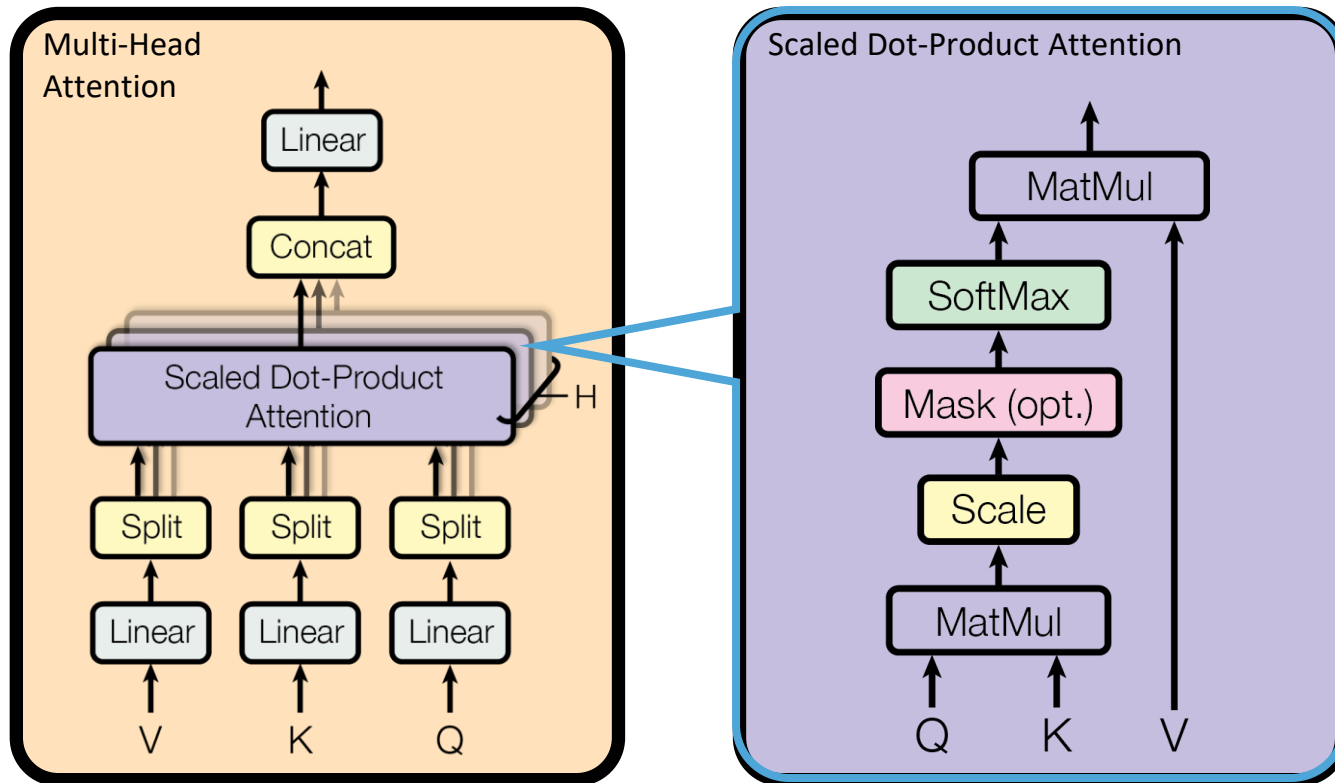
Figure 2: GPT-3 175B validation accuracy vs. number of trainable parameters of several adaptation methods on WikiSQL and MNLI-matched. LoRA exhibits better scalability and task performance. See [Section I.2](#) for more details on the plotted data points.

Implementation:

<https://github.com/microsoft/LoRA>

<https://huggingface.co/docs/diffusers/training/lora>

Review: Attention Mechanism



Original LoRA was just applied to the attention weights:
 W_Q , W_K , W_V , and W_O

Guanaco: QLoRA

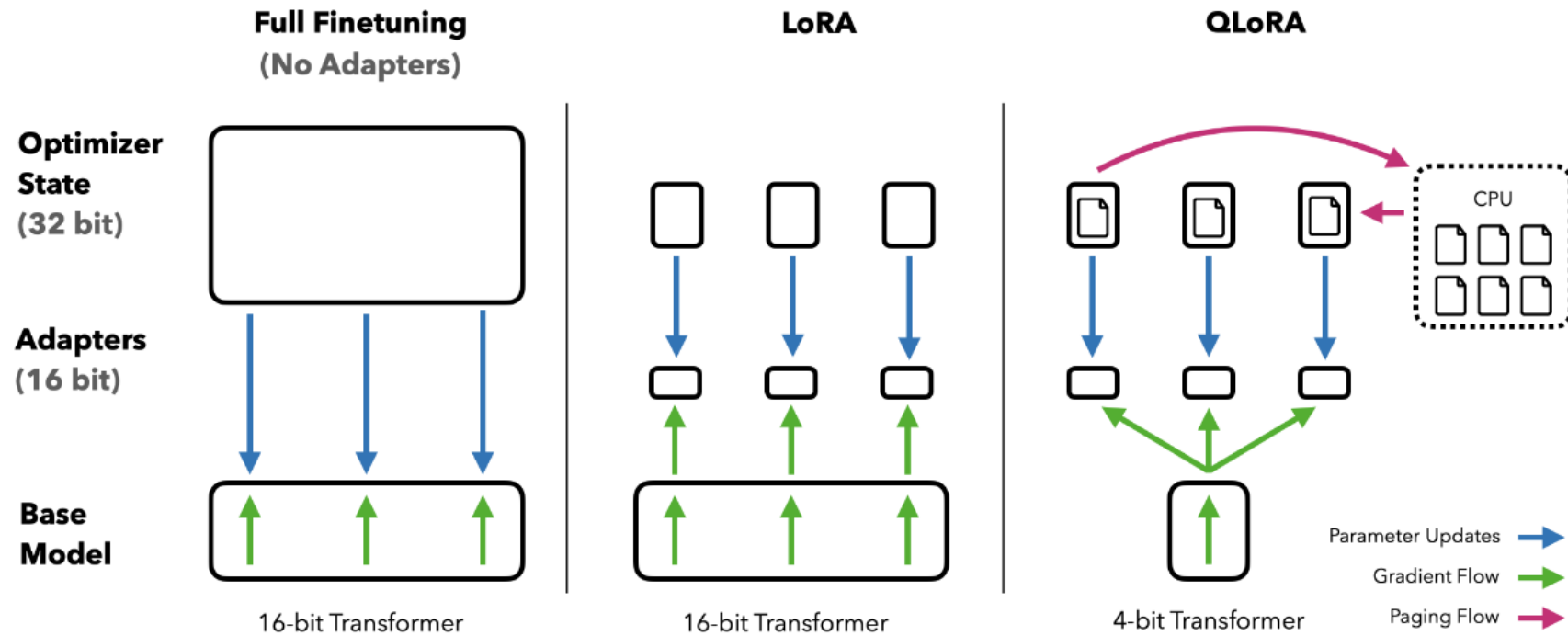


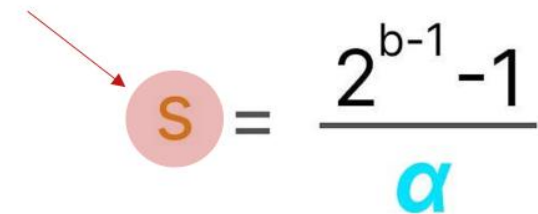
Figure 1: Different finetuning methods and their memory requirements. QLoRA improves over LoRA by quantizing the transformer model to 4-bit precision and using paged optimizers to handle memory spikes.

Guanaco: QLoRA

4-bit NormalFloat quantization

Uses *double quantization* (quantizing the quantization constants)

Stored as FP32


$$S = \frac{2^{b-1} - 1}{\alpha}$$

Dinesh Raghu

Guanaco: QLoRA

Table 3: Experiments comparing 16-bit BrainFloat (BF16), 8-bit Integer (Int8), 4-bit Float (FP4), and 4-bit NormalFloat (NF4) on GLUE and Super-NaturalInstructions. QLoRA replicates 16-bit LoRA and full-finetuning.

Dataset Model	GLUE (Acc.)	Super-NaturalInstructions (RougeL)				
	RoBERTa-large	T5-80M	T5-250M	T5-780M	T5-3B	T5-11B
BF16	88.6	40.1	42.1	48.0	54.3	62.0
BF16 replication	88.6	40.0	42.2	47.3	54.9	-
LoRA BF16	88.8	40.5	42.6	47.1	55.4	60.7
QLoRA Int8	88.8	40.4	42.9	45.4	56.5	60.7
QLoRA FP4	88.6	40.3	42.4	47.5	55.6	60.9
QLoRA NF4 + DQ	-	40.4	42.7	47.7	55.3	60.9

Table 8: Evaluation of biases on the CrowS dataset. A lower score indicates lower likelihood of generating biased sequences. Guanaco follows the biased pattern of the LLaMA base model.

	LLaMA-65B	GPT-3	OPT-175B	Guanaco-65B
Gender	70.6	62.6	65.7	47.5
Religion	79.0	73.3	68.6	38.7
Race/Color	57.0	64.7	68.6	45.3
Sexual orientation	81.0	76.2	78.6	59.1
Age	70.1	64.4	67.8	36.3
Nationality	64.2	61.6	62.9	32.4
Disability	66.7	76.7	76.7	33.9
Physical appearance	77.8	74.6	76.2	43.1
Socioeconomic status	71.5	73.8	76.2	55.3
Average	66.6	67.2	69.5	43.5

LoRA

PROS

Preserving a lot of the original model calculations

Does about the same performance-wise

Significantly smaller than the original

CONS

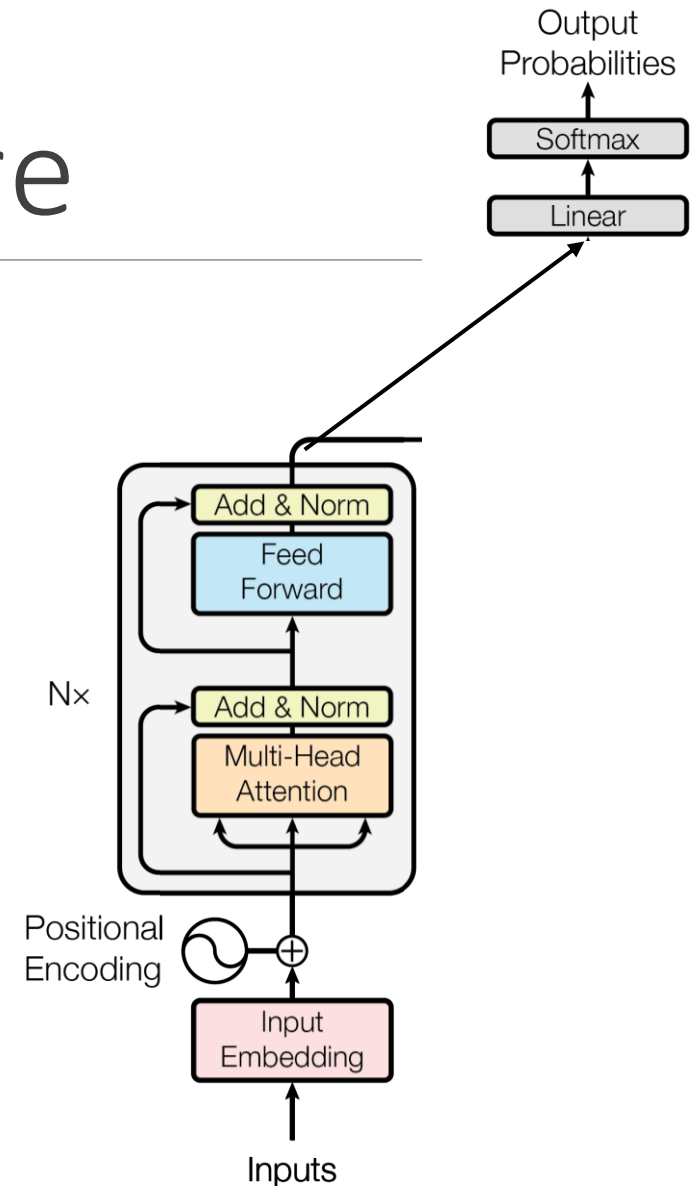
Potentially slightly worse for larger models

Efficient Training

Review: Transformer Architecture

Multi-Head Attention = $\text{MultiHeadAtt}(\mathbf{H}_i^{enc}, \mathbf{H}_i^{enc}, \mathbf{H}_i^{enc})$
Add & Norm = $\text{LayerNorm}(\text{Multi-Head Attention} + \mathbf{H}_i^{enc})$
Feed Forward = $\max(0, \text{Add & Norm } \mathbf{W}_1 + b_1) \mathbf{W}_2 + b_2$
Add & Norm (2) = $\text{LayerNorm}(\text{Feed Forward} + \mathbf{H}_i^{enc})$
 \mathbf{H}_{i+1}^{enc} = Add & Norm (2)

A decoder-only architecture is very similar to the encoder of the original transformer architecture



Sparse Mixture-of-Experts

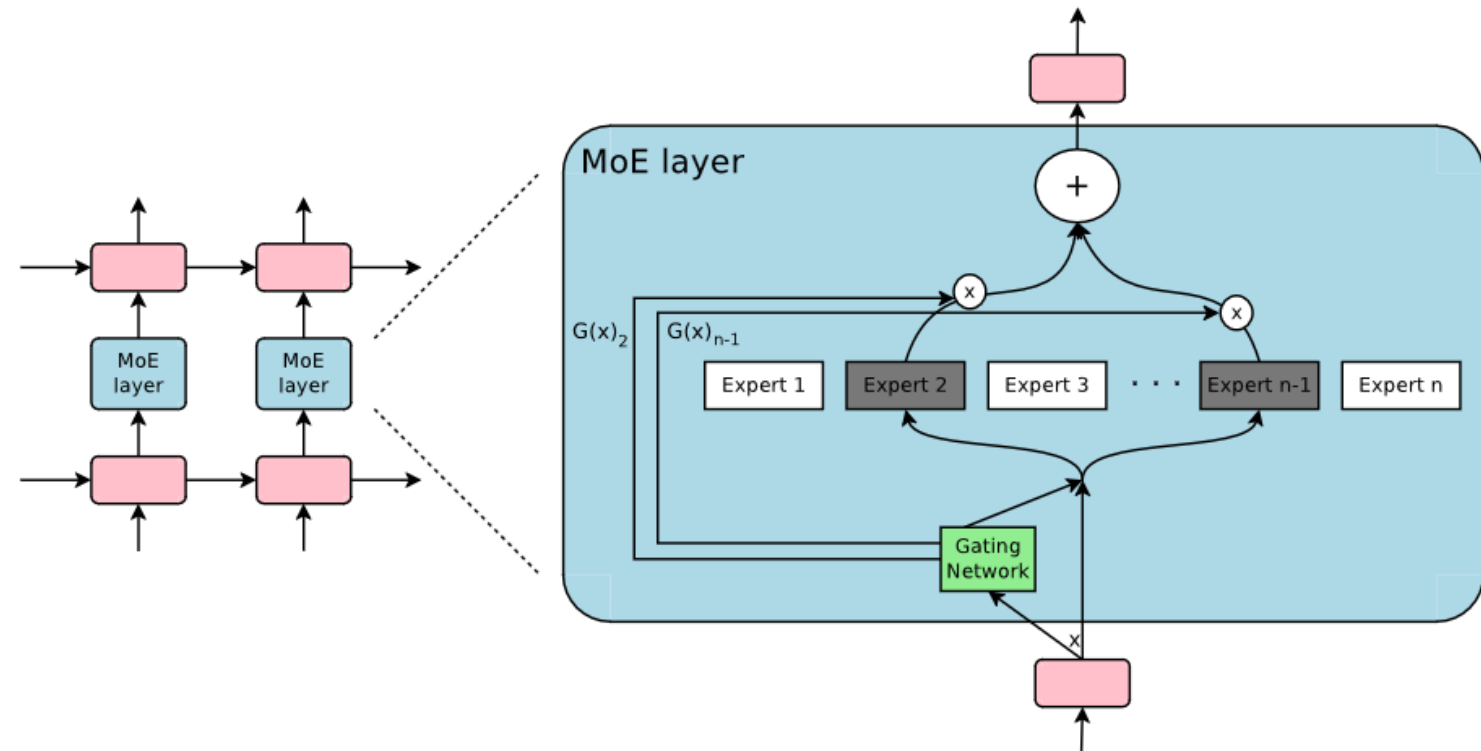


Figure 1: A Mixture of Experts (MoE) layer embedded within a recurrent language model. In this case, the sparse gating function selects two experts to perform computations. Their outputs are modulated by the outputs of the gating network.

Implementations

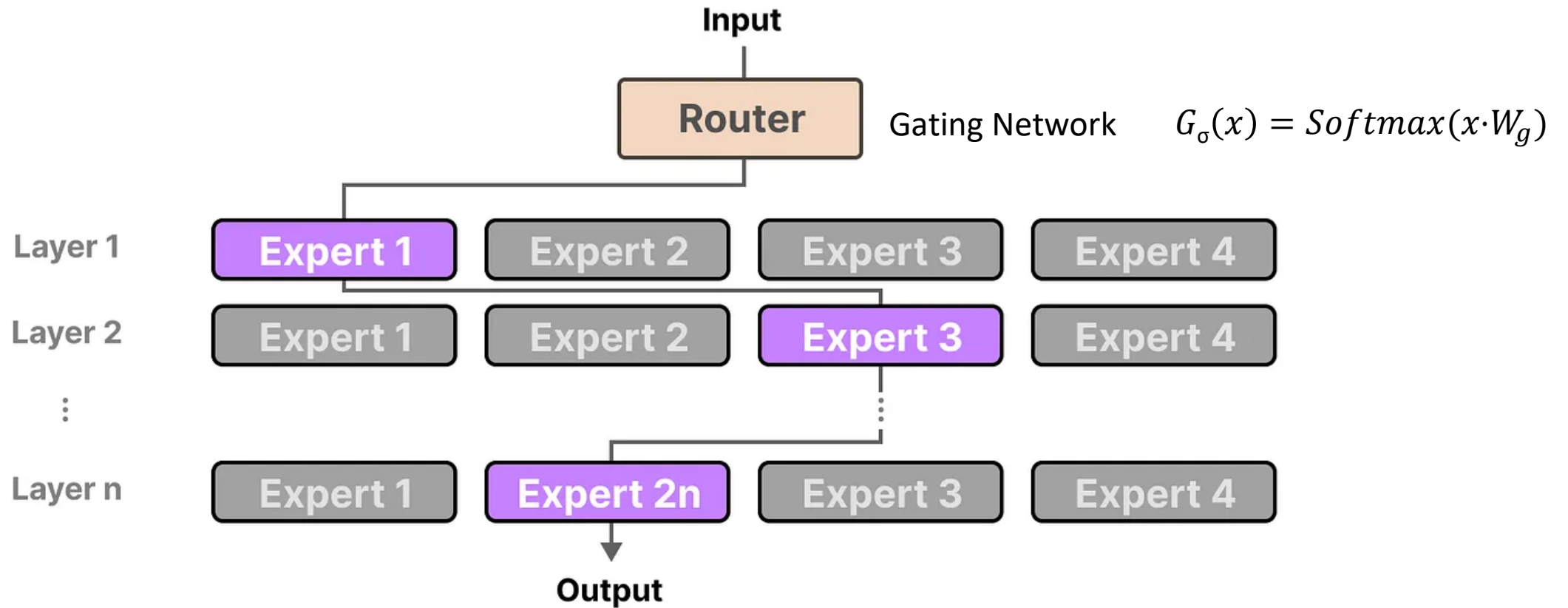
Megablocks: <https://github.com/stanford-futuredata/megablocks>

Fairseq: https://github.com/facebookresearch/fairseq/tree/main/examples/moe_lm

OpenMoE: <https://github.com/XueFuzhao/OpenMoE>

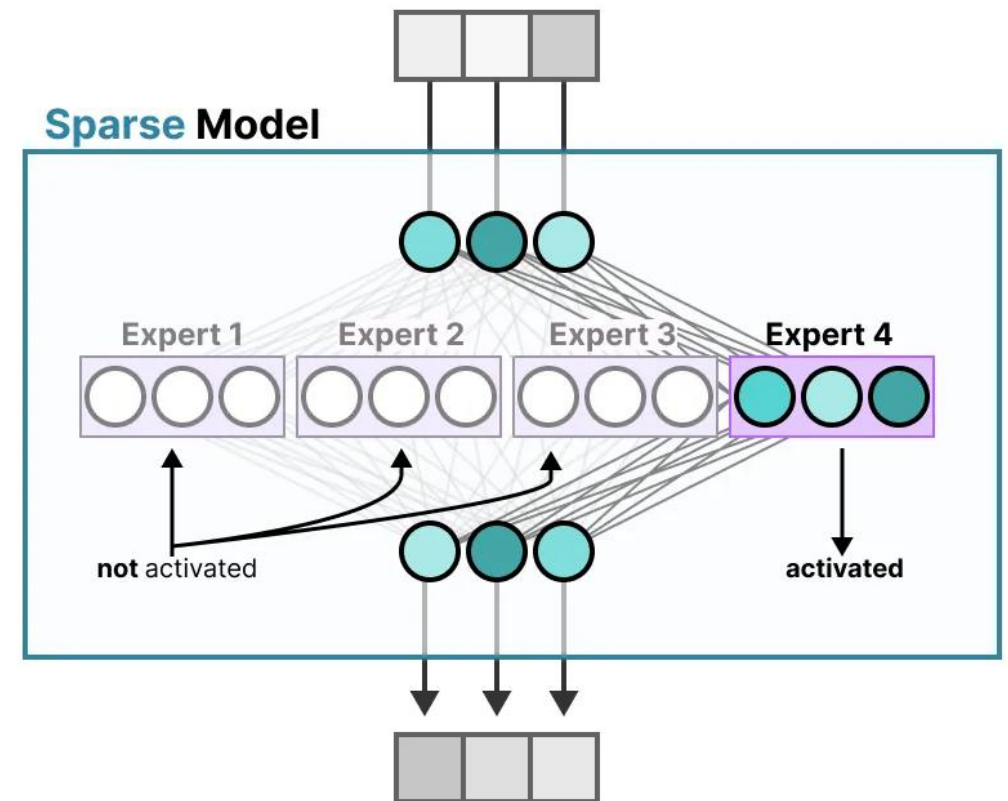
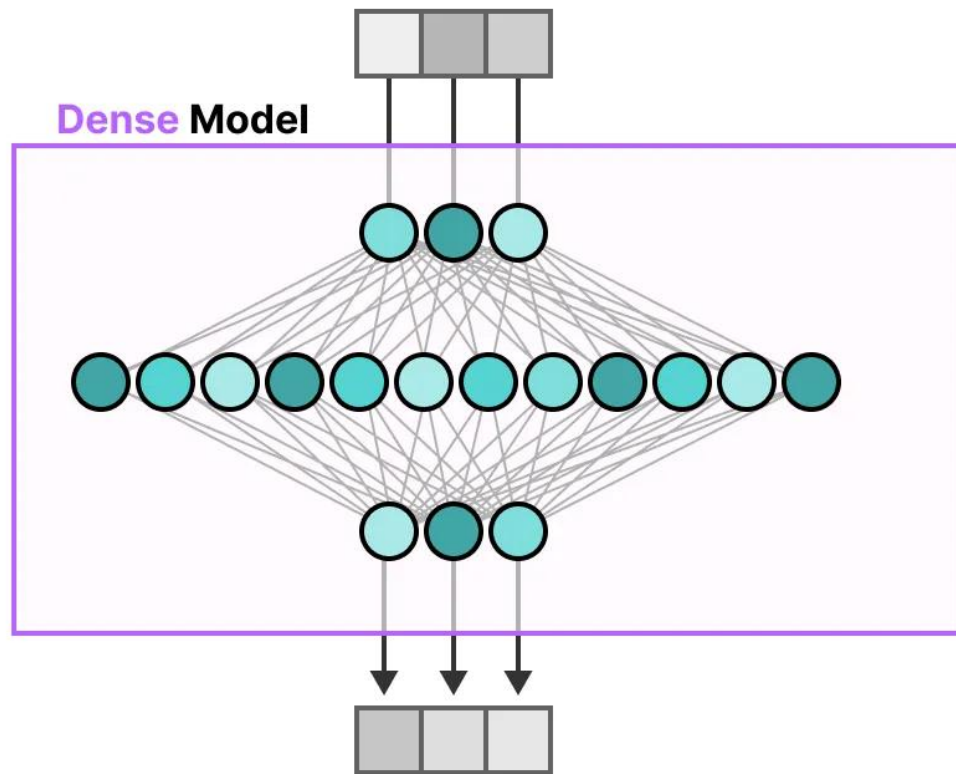
Mixture-of-Experts

$$y = \sum_{i=1}^n G(x)_i E_i(x)$$

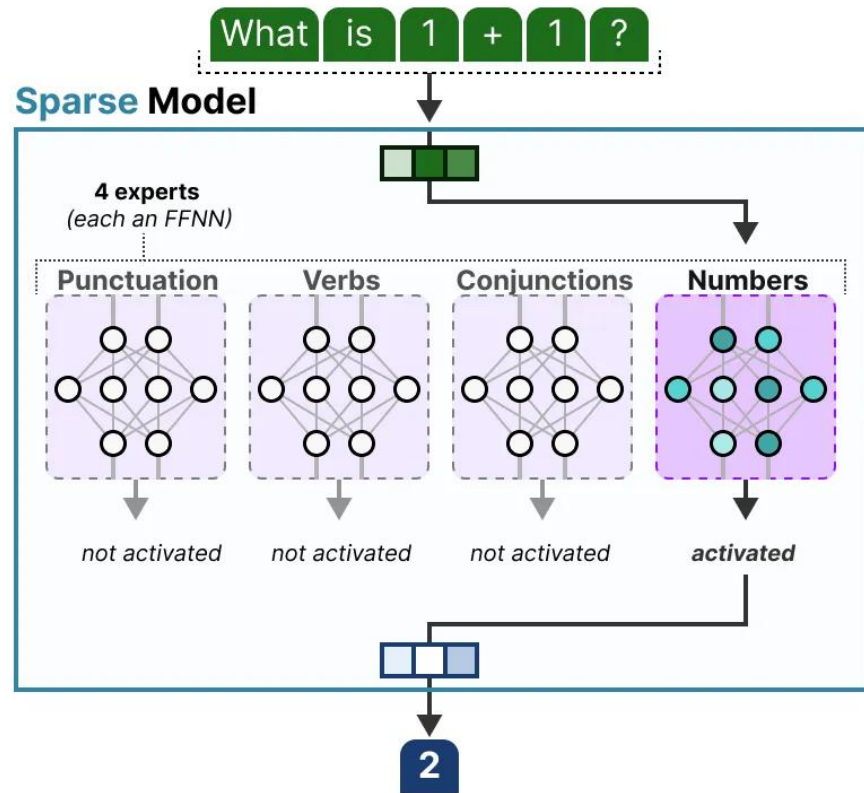


<https://substack.com/home/post/p-148217245>
Equations from <https://huggingface.co/blog/moe>

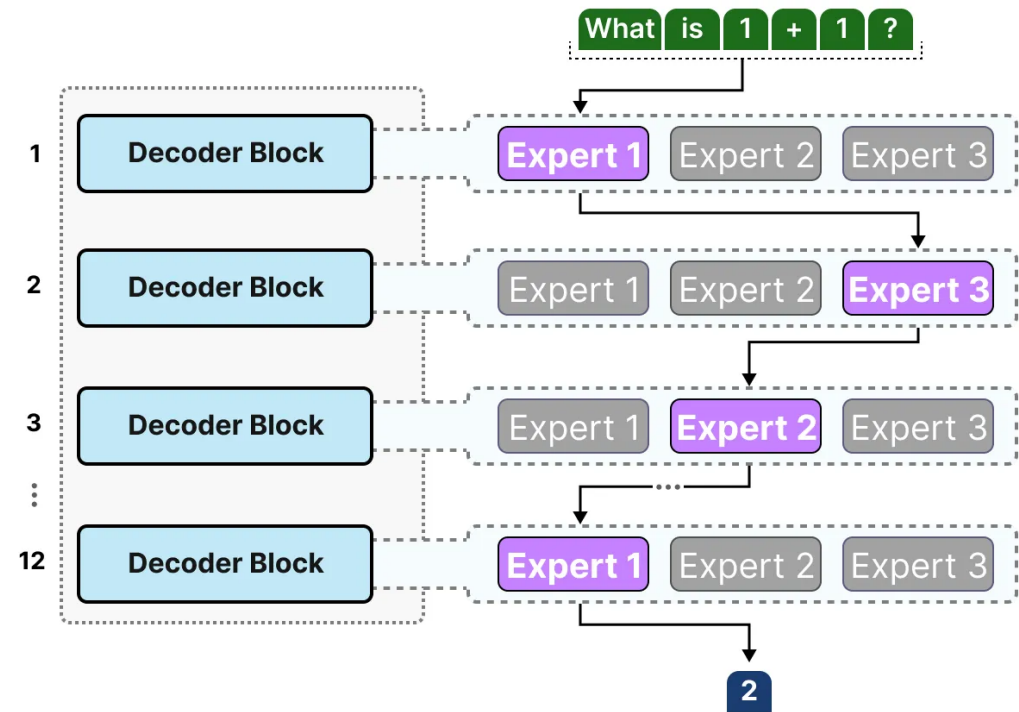
Sparse MoE

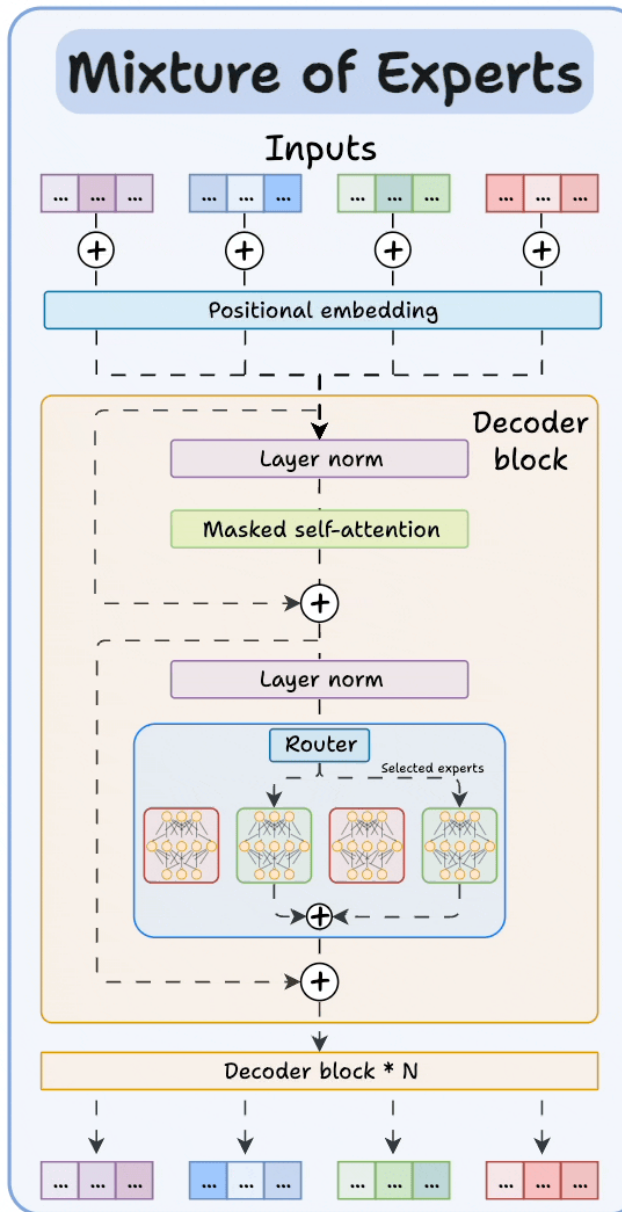
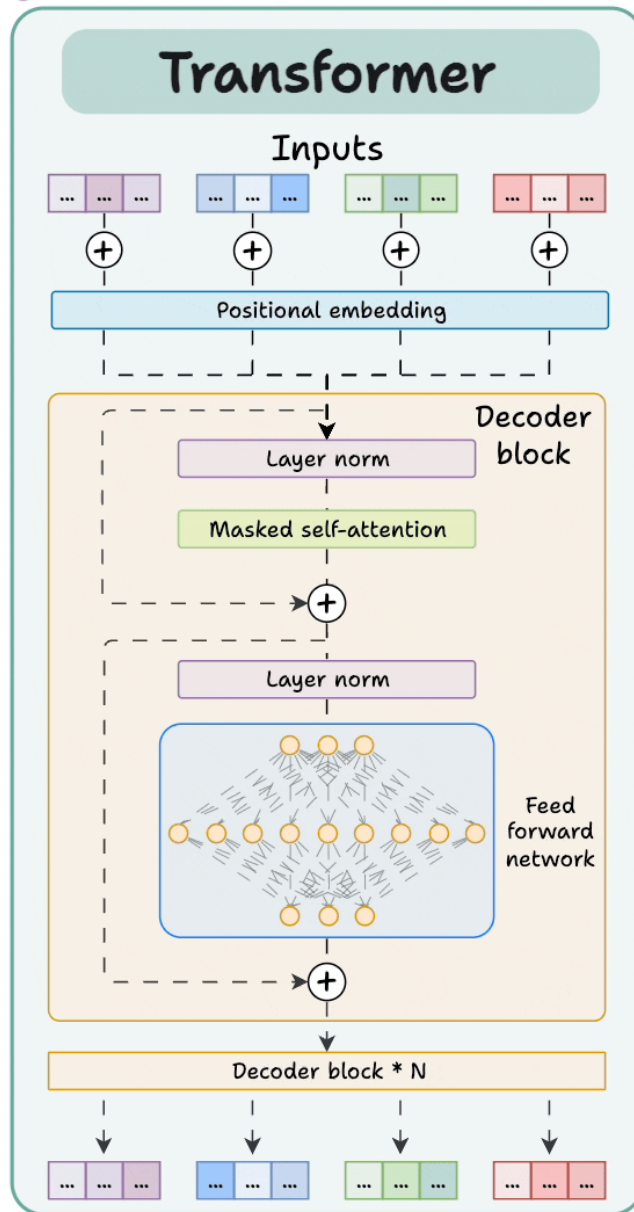


MoE Example



Most transformers have multiple decoder blocks





Mixtral 8x7B

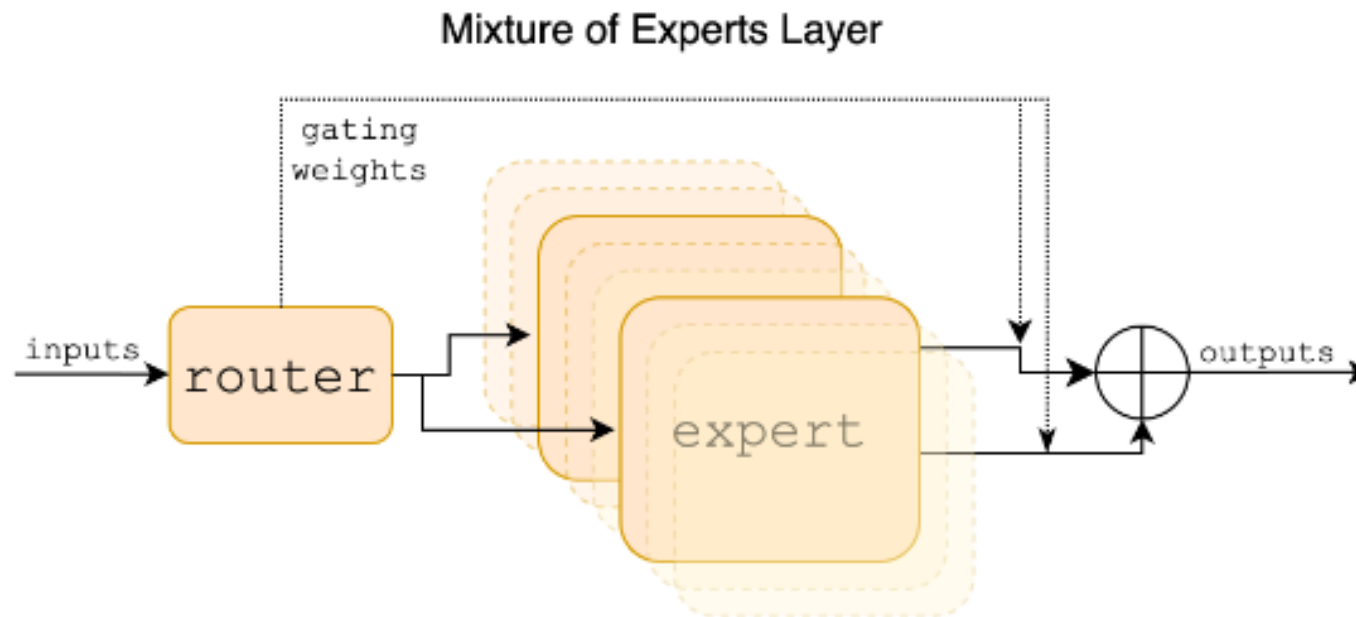


Figure 1: Mixture of Experts Layer. Each input vector is assigned to 2 of the 8 experts by a router. The layer’s output is the weighted sum of the outputs of the two selected experts. In Mixtral, an expert is a standard feedforward block as in a vanilla transformer architecture.

Mixtral 8x7B

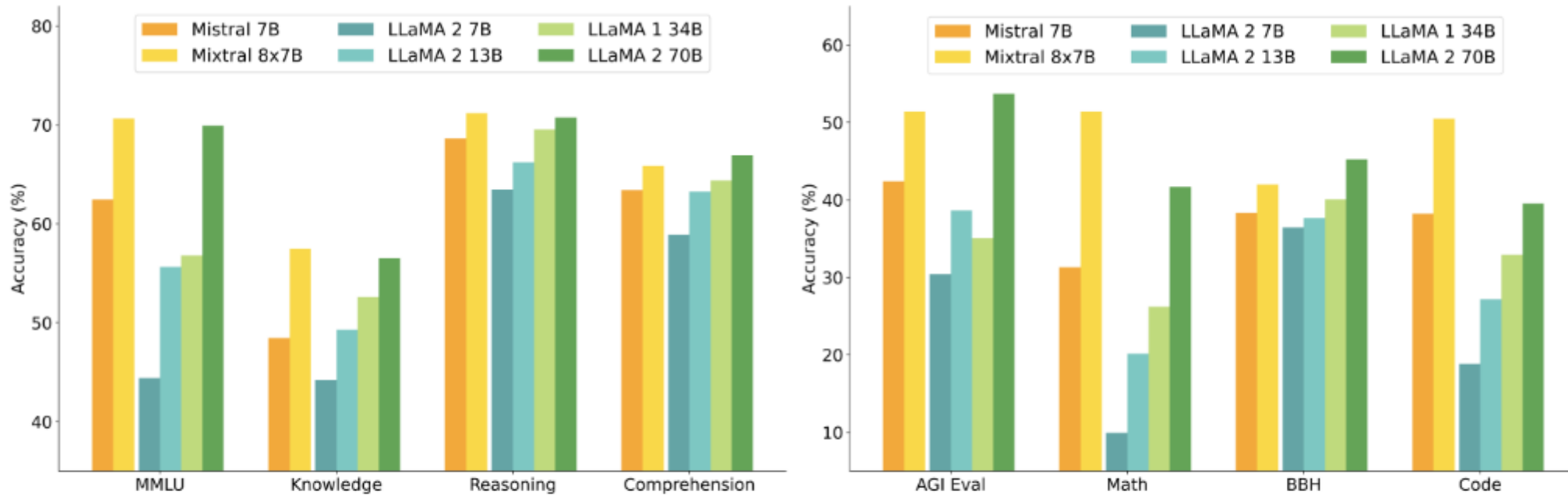
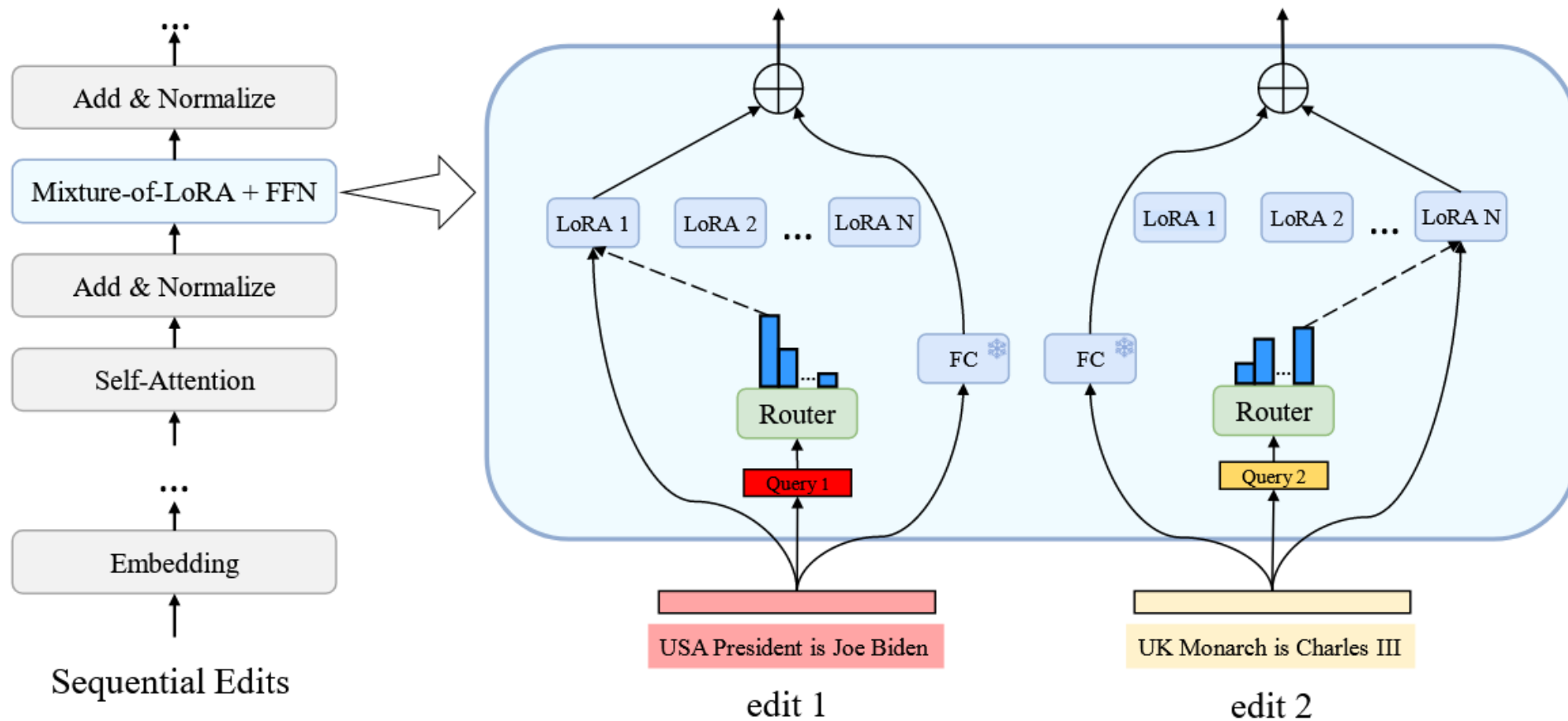


Figure 2: Performance of Mixtral and different Llama models on a wide range of benchmarks. All models were re-evaluated on all metrics with our evaluation pipeline for accurate comparison. Mixtral outperforms or matches Llama 2 70B on all benchmarks. In particular, it is vastly superior in mathematics and code generation.

ELDER: Mixture-of-LoRA



Mixture of Experts

PROS

Fewer parameters activated

Potentially more relevant answers

CONS

Same amount of parameters as original transformer

Potentially more errors if routed wrong

Might not work well if not clear categories to separate into experts