# Llama3 & DeepSeek v3

## Architecture | Fine-tuning | Inference



Dr. Karthik Mohan, Feb 26 2025

## 1. Types of Training

### 3. DeepSeekV3





## 4. Notebook Walkthrough









# Types of LLM training

**Foundation Model/Pre-Trained Models** 

**Post-Training** 





**Reinforcement Learning with Human** Feedback (RLHF)/Direct Preference **Optimization (DPO)** 

# Types of LLM training

**Foundation Model/Pre-Trained Models** 

**Supervised Fine-tuning** 

**Instruction Fine-tuning** 



**Foundation Model/Pre-Trained** Models

**Instruction Fine-tuning** 

**Reinforcement Learning with Human** Feedback (RLHF)/Direct Preference **Optimization (DPO)** 

# Types of LLM training

### **Supervised Fine-tuning**







## Why not just do a single training instead of multiple trainings for LLMs?

1.

2. 3. 4.



Single training won't work Multiple trainings also different sources of high quality vs medium quality data Single Training has no foundation to build on **Multiple Training is faster** 

# **Pre-Trained vs Instruct Mode**

## Pre-Trained Models are trained as a Masked Language Model and for Next Token Prediction or Multiple Token Prediction

## Instruct Fine-tuning is fine-tuning a pretrained model to follow instructions

# **Pre-Trained vs Instruct Mode**

Pre-Trained Models are trained as a Masked Language Model and for Next Token Prediction or Multiple Token Prediction

Instruct Fine-tuning is fine-tuning a pretrained model to follow instructions on a wide variety of tasks

# **Pre-Trained vs Instruct Fine-tuned Model**

**Pre-Trained** Input: "The red fox \_\_\_\_" Output: "chased" Input: "The red fox chased the \_\_\_\_\_" Output: "blue" Input: "The red fox chased the blue \_\_\_\_" Output: "bird"

## **Instruct Fine-tuned**

Input: "You are to complete the following sentence. Sentence: 'The red fox ' " Output: "The red fox chased the blue bird. And the bird flew away in the nick of time!"

# **Instruct Fine-tuning vs Supervised Fine-Tuning**

**Instruct Fine-tuning** is fine-tuning a pre-trained model to follow instructions on a wide variety of tasks

Supervised Fine-tuning is fine-tuning the pretrained LLM on specific tasks: Sentiment analysis, text summarization, question answering, intent detection, etc



# **Instruct Fine-tuning vs Supervised Fine-Tuning**

**Instruct Fine-tuning** is fine-tuning a pre-trained model to follow instructions on a wide variety of tasks

**Supervised Fine-tuning** is fine-tuning the pretrained LLM on specific tasks: Sentiment analysis, text summarization, question answering, intent detection, etc



# **SFT vs Instruct Fine-tuned Model**

## **Supervised Fine-tuning**

Input: "I am not feeling good today" Output: "Unhappy" Input: "I would love to go to New York and spend time on Times Square" Output: "New York, Times Square" Input: "What is the tallest mountain in the world?" Output: "Mount Everest"

## **Instruct Fine-tuning**

Input: "You are to complete the following sentence. Sentence: 'The red fox ' " Output: "The red fox chased the blue bird. And the bird flew away in the nick of time!"

# **Instruct Fine-tuning vs Supervised Fine-Tuning**

# **SFT** trains a pre-trained LLM to do well on specific NLP tasks.

Instruction Fine-tuning looks at the ability of an LLM to follow instructions on variety of tasks. There is no clear indication of task + makes the LLM behave more like an AI agent





# What is the loss function for SFT?

## What is the loss function for Instruct models?

1. 2.

1. 2.



**Cross-entropy Quadratic loss** 

**Cross-entropy Quadratic loss** 

# Lama 3 Instruct Performance

### Meta Llama 3 Instruct model performance

	Meta Llama 3 8B	<b>Gemma</b> 7B - It Measured	Mistral 7B Instruct Measured
<b>MMLU</b> 5-shot	68.4	53.3	58.4
<b>GPQA</b> 0-shot	34.2	21.4	26.3
<b>HumanEval</b> 0-shot	62.2	30.5	36.6
<b>GSM-8K</b> 8-shot, CoT	79.6	30.6	39.9
<b>MATH</b> 4-shot, CoT	30.0	12.2	11.0

Reference: <u>https://ai.meta.com/blog/meta-llama-3/</u>

	Meta	<b>Gemini</b>	Claude 3
	Llama 3	Pro 1.5	Sonnet
	70B	Published	Published
<b>MMLU</b> 5-shot	82.0	81.9	79.0
<b>GPQA</b>	39.5	<b>41.5</b>	<b>38.5</b>
0-shot		сот	Сот
<b>HumanEval</b> 0-shot	81.7	71.9	73.0
<b>GSM-8K</b>	93.0	<b>91.7</b>	<b>92.3</b>
8-shot, CoT		11-shot	0-shot
<b>MATH</b> 4-shot, CoT	50.4	<b>58.5</b> Minerva prompt	40.5

# Lama 3 Instruct Performance



Reference: <u>https://ai.meta.com/blog/meta-llama-3/</u>



## If you had to fine-tune a LLM model that can detect an emotion from a review - which would you pick?

Fine-tune an instruct model Fine-tune the pre-trained model Fine-tune the SFT model

1.

2. 3.

# Lama2vs Lama3

7 times larger pre-train data set. 15 Trillion **Tokens** of data ~ 150 million books High-quality filters to filter out bad data in training - Use Llama2 Better "data mix" - Trivia, STEM, coding, historical knowledge

Larger model means better performance (8B vs **70B**) But more data = better performance (also avoids over-fitting). Log-linear improvement from 200B to 15T tokens



### **Data Parallelization**

### **Model Parallelization**

Latent Attention





### **Data Parallelization**

### **Model Parallelization**

Latent Attention



### **Data Parallelization**

### Model Parallelization

Latent Attention



### **Data Parallelization**

### Model Parallelization

**Latent Attention** 

In the model parallelism regime - Assume a new **DeepFetch** model got released on the market with 1 trillion parameters. Assume that for pretraining, you are using H100 GPUs with 80 GB ram. How many GPUs would you need to have to hold the model in memory?

1. 2. 3. 4.

## **CE #3**

25 50 75

100