

$$W \rightarrow \text{huge numbers}$$

entire weights

frozen (freeze)

$\leftarrow A$

$\leftarrow B$

$W' = \underline{W} + \Delta W$

$\Delta W = \underline{(A * B)}$

$$A(\gamma \times d)$$

$$B(d \times \gamma)$$

$\hookrightarrow A$

$d \Rightarrow$  dimensions  
 $\gamma \Rightarrow$  low rank factor

Interview question

LORA

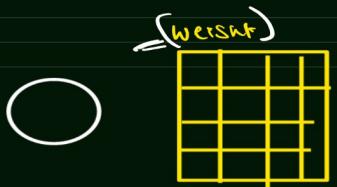
32 GB RAM

LORA =  
12-24 GB RAM

Transformer  $\Rightarrow$  Self Attention + NN

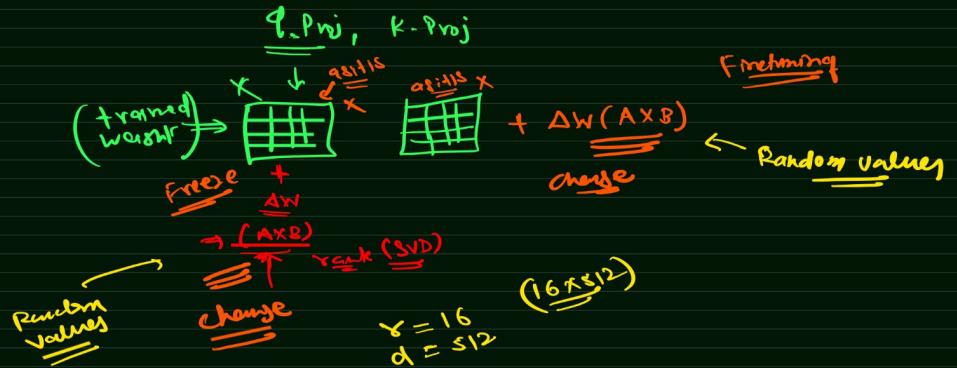
huge number weights

LORA subset



$(A \times B)$

Rank  $\Rightarrow$  (SVD)



$\Rightarrow [1 \text{ BIT} \Rightarrow 0 \text{ or } 1]$   $\xrightarrow{\text{System}} \underline{\underline{0 \text{ or } 1}}$

$1 \text{ BYTE} \Rightarrow 8 \text{ BIT}$

$1 \text{ KB} \Rightarrow 1024 \text{ BYTE}$

$1 \text{ MB} \Rightarrow 1024 \text{ KB}$

$1 \text{ GB} \Rightarrow 1024 \text{ MB}$

Data type  $\Rightarrow$

$\text{float 32} = 32 \text{ bit} = 4 \text{ byte}$   
 $\text{float 16} = 16 \text{ bit} = 2 \text{ byte}$   
 $\text{float 8} = 8 \text{ bit} = 1 \text{ byte}$   
 $\text{Int 8} = 8 \text{ bit} = 1 \text{ byte}$

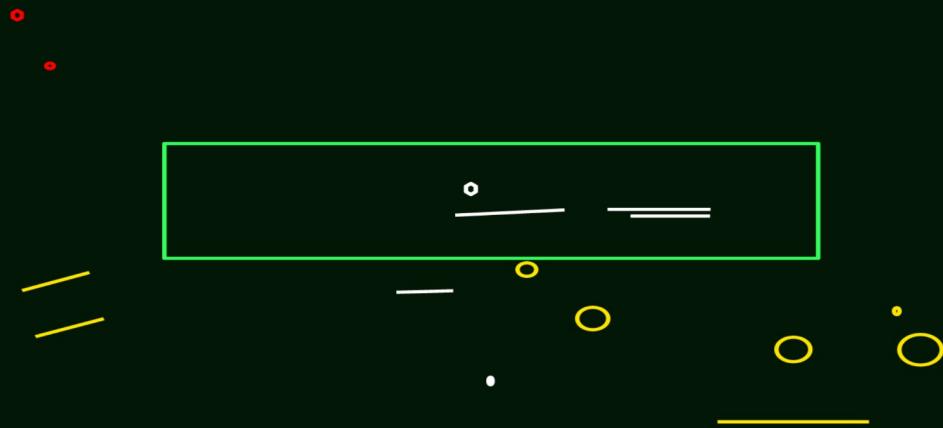
Data type

$\text{weight} \rightarrow 2.87 \rightarrow \underline{\underline{\text{float 32}}} \rightarrow 4 \text{ byte}$

$\xrightarrow{\text{GPT 3.5}} \underline{\underline{125 \text{B}}} \text{ weight} = 125 \text{B} \times 4 \text{ byte}$

Quantization  $\rightarrow$  reduce precision

$\underline{\underline{\text{FP32}}} \rightarrow \underline{\underline{\text{Int 8}}}$   
 $\underline{\underline{1 \text{ Byte}}} \rightarrow \underline{\underline{1 \text{ Byte}}}$



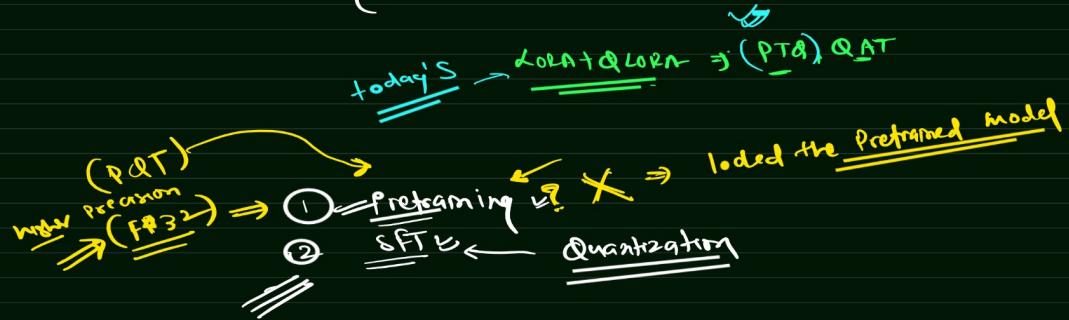
(PTQ)

## 1 Post-training quantization

1 First trained in full precision (FP32)

2 After training the w & b are converted to lower precision

{ FP16 }  
INT 8



(QAT)

## 2 Quantization-aware training

1 While training use quantization effects

2 Instead of training on full precision the model will train with quantization effects

PTQ

Post Quantization technique

1 unsupervised  $\leftarrow$  FP32  $\rightarrow$  full precision

2 SFT  $\leftarrow$  quantization

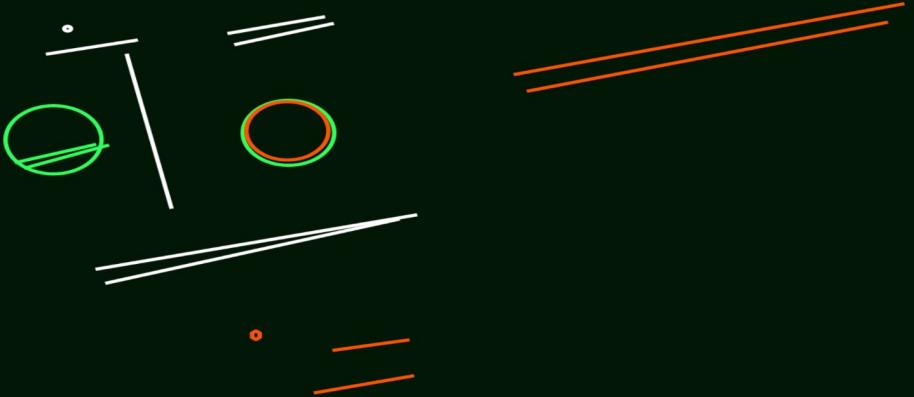
`bnb_config=BitsAndBytesConfig(`

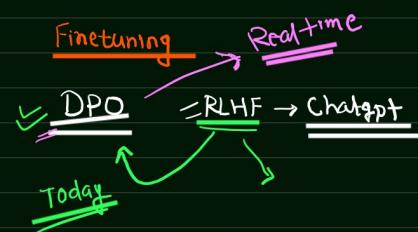
`load_in_4bit=True,`

`bnb_4bit_use_double_quant=True,`

`bnb_4bit_quant_type="nfp4",`

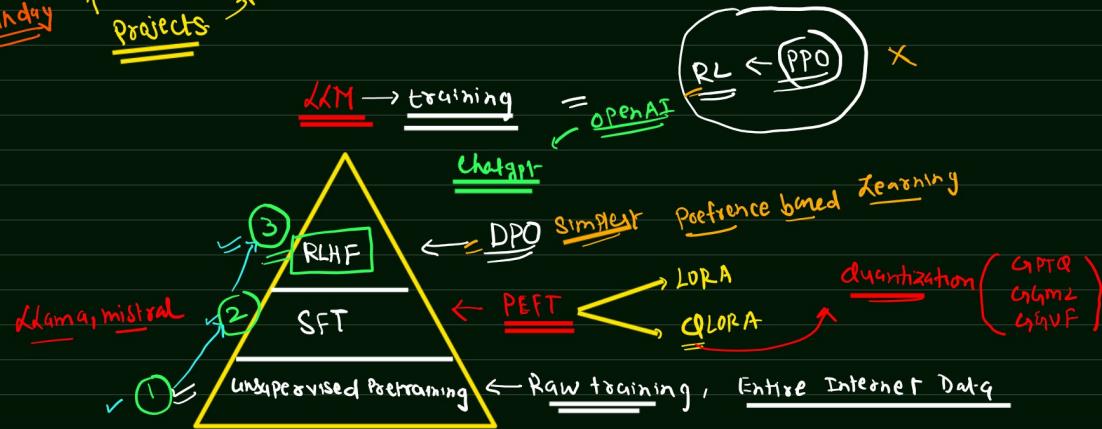
normal float 4 bit (Nfp4bit)





= RAH (Realtime) = Agents

Sunday ↑ Projects ↑



RLHF

RL → X

DPO → ✓

Roadmap

Bigger picture

RL

Agent ← model ← software

Agent State

Action Reward

Environment



- 1950  
 - Bellman eq.  
 - Dynamic prog  
 $\begin{array}{c} \text{Agent}, \text{Action, Reward} \\ \hline \text{Environment} \end{array}$

1990 - Monte Carlo method

1990-2000 - Q-learning

(2015) Google Deep Mind

ANN  $\rightarrow$  DL + RL  $\Rightarrow$  Deep Reinforcement Learning

1995  
 \* Policy based method  
 & Policy gradient method

PPO  $\Rightarrow$  OpenAI  $\leftarrow$  (ChatSPT, RLHF)

Proximal Policy Optimization

State of Art

RL

Data (Supervised, unsupervised data)

DL  $\rightarrow$  FP + Loss + BP

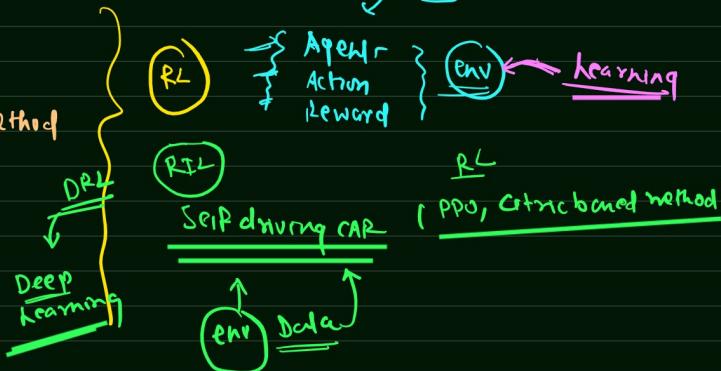
X Data X Supervised / Unsupervised

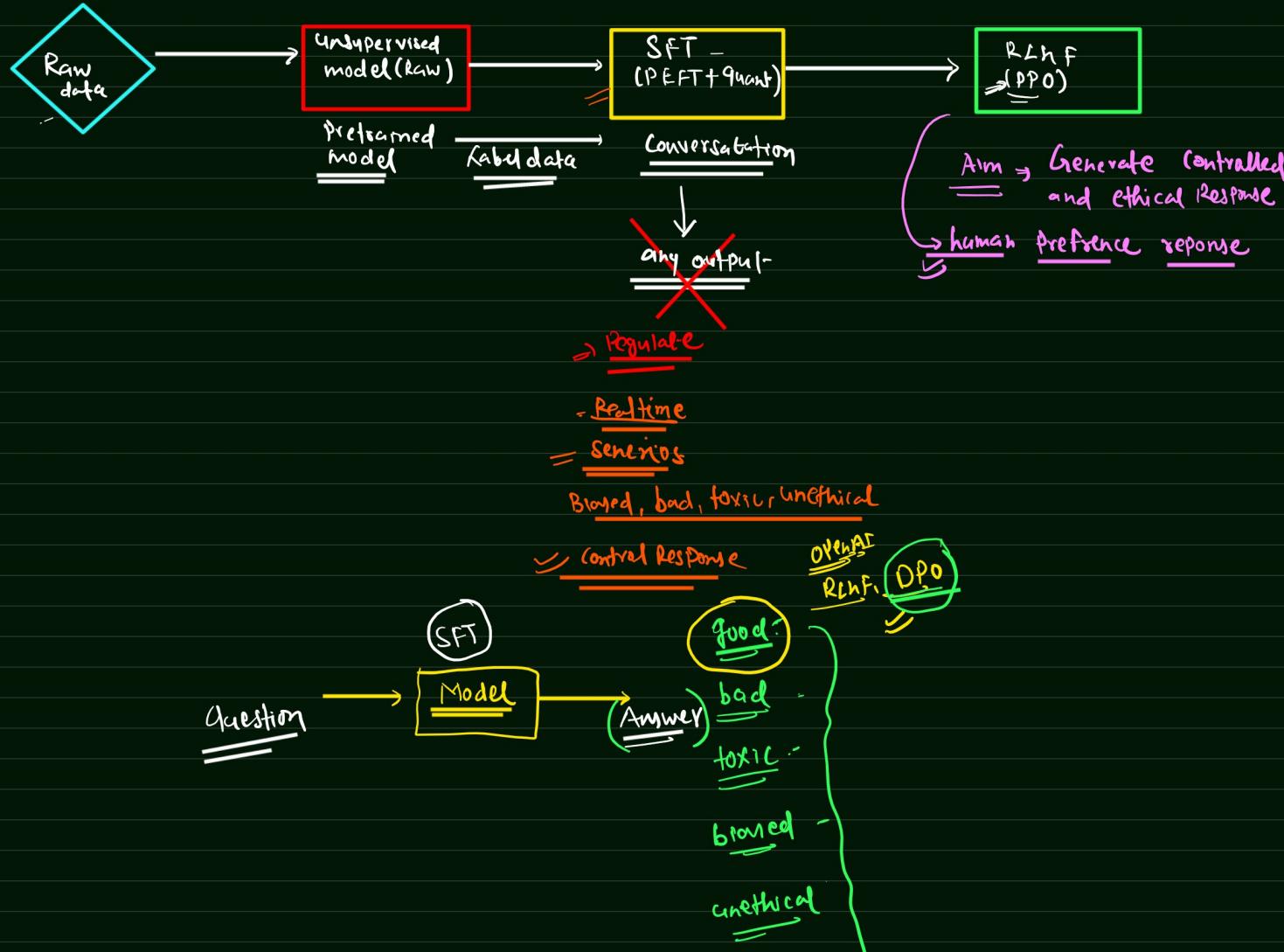
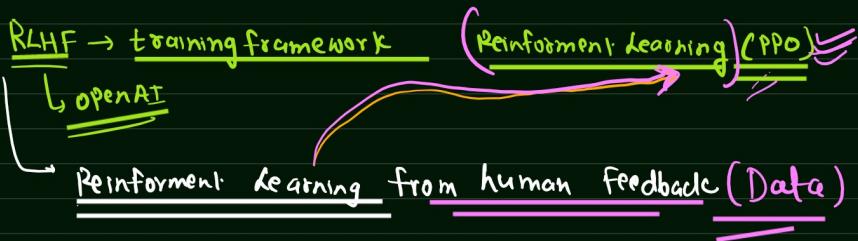
## Reinforcement Learning

- (1) Q Learning
- (2) DRL
- (3) Policy based method :- Policy gradient method
- (4) PPO
- (5) Actor-Critic method

OPENAI

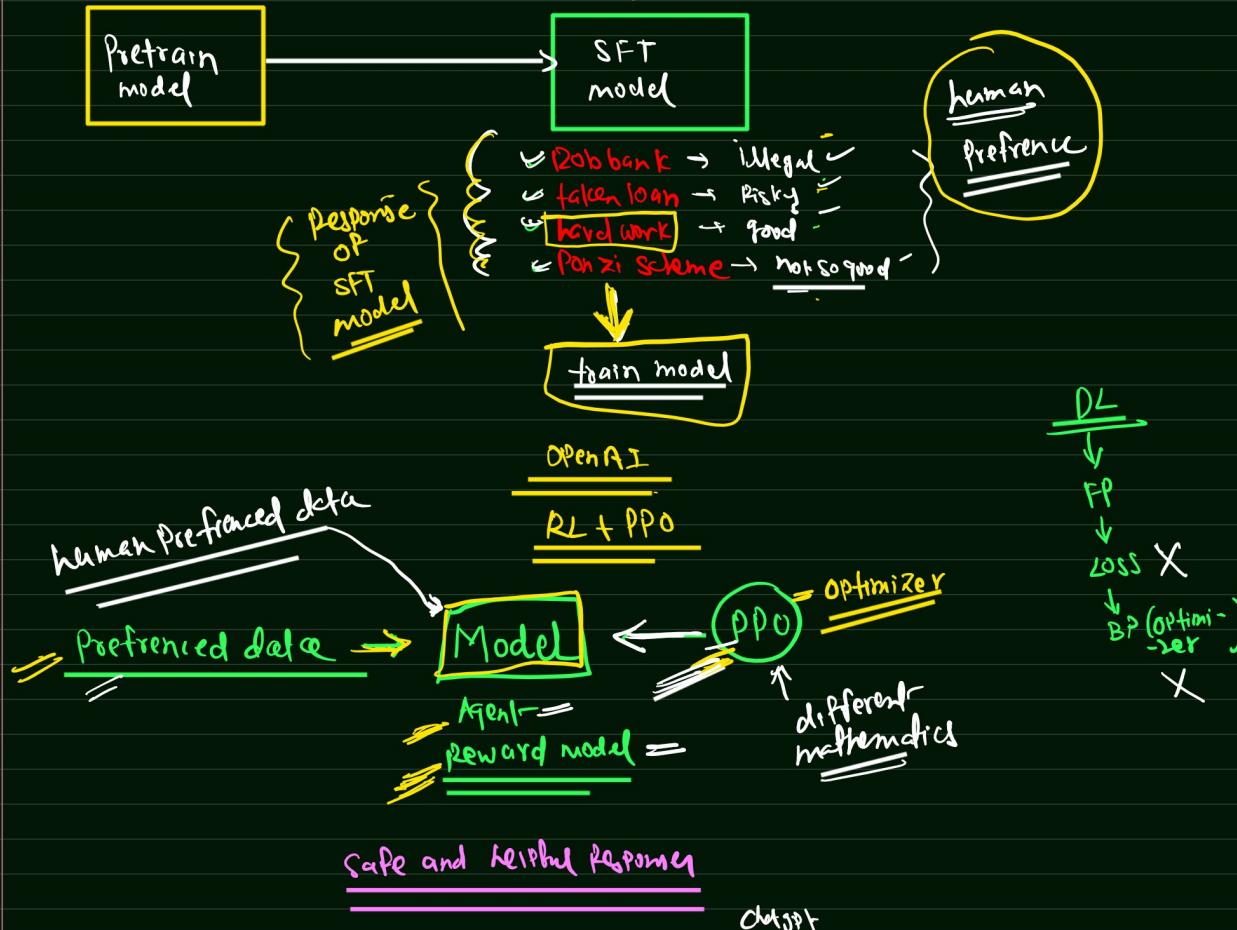
(6) KTO    KLHF  
(Kahneman - Tversky optimization)





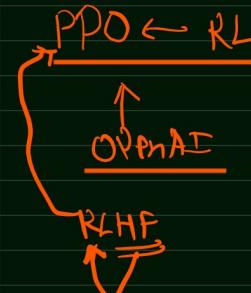
RLHF ← OpenAI ← ChatGPT

how to make money?



#### • RLHF Vs Normal PPO

Feature	Normal PPO (Reinforcement Learning)	RLHF (Human Feedback-Based RL)
Goal	Maximize numeric rewards (Game score, robotics control)	Maximize human preference-based rewards
Reward Source	Direct environment rewards	Human-labeled reward model
Application	Games, robotics, simulations	AI assistants, ChatGPT, LLaMA models
Output Control	Less control	More ethical & human-friendly control



DPO

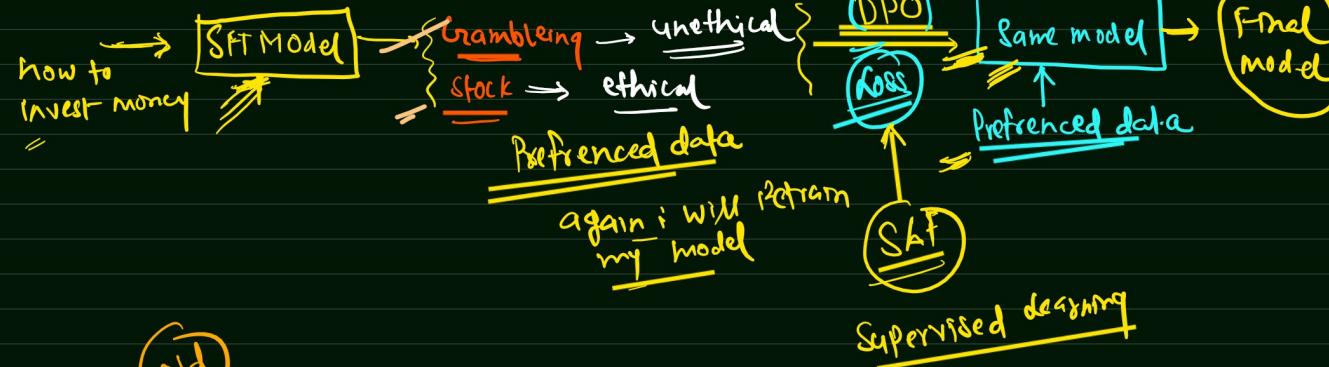
Direct preference optimization



Chatspl

Supervised loss function

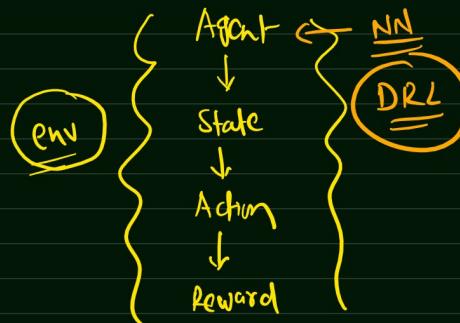
DPO directly work for preference based learning



Old

OpenAI 2022-23

Proximal Policy optimization

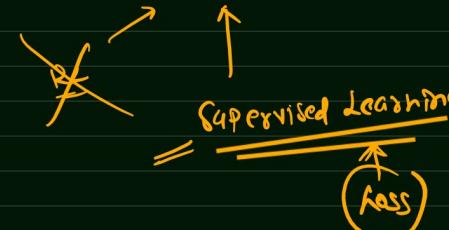


& learning, KTO, PPO, critic method

Latest

DPO

direct Preference optimization



BP ⇒ optimize = update weight

QD, SGD, MinGD, ADAGRAD, ADAM

Feature	Reinforcement Learning (RL)	DPO (Direct Preference Optimization)
Learning Type	Reward-based learning	Direct supervised learning
Algorithm Used	PPO, Q-Learning, Actor-Critic	Simple loss function
Training Complexity	Zyada tuning aur instability	Simple aur stable
Compute Requirement	High (Zyada tuning chahiye)	Low (Efficient & fast)

More