# DS 4440 "Post" Training
## Or, "Aligning" LLMs

**Large Language Models (LLMs)**
are stacks of Transformer layers

2018    BERT (10/23 lecture)
340M params

2019    GPT-2
1.5b params

2020    GPT-3
~175b params

More capable but problematic
+ LM objective only not super
useful for most tasks.

Pre-Training provides the base.

Post-Training is about making models useful, or — if you're fancy — alignment.

Instruction Tuning (SFT) is a simple approach which aims to teach models to follow instructions.

Compile standard supervised tasks and pre-pend INSTRUCTIONS

$X_1$ Today the Trump Campaign ...

$Y_1$ Trump ...
(Summerization)

$X_2$ This movie rocked ...

$Y_2$ Positive
(classification)

Summarize the following document:
Today the Trump campaign

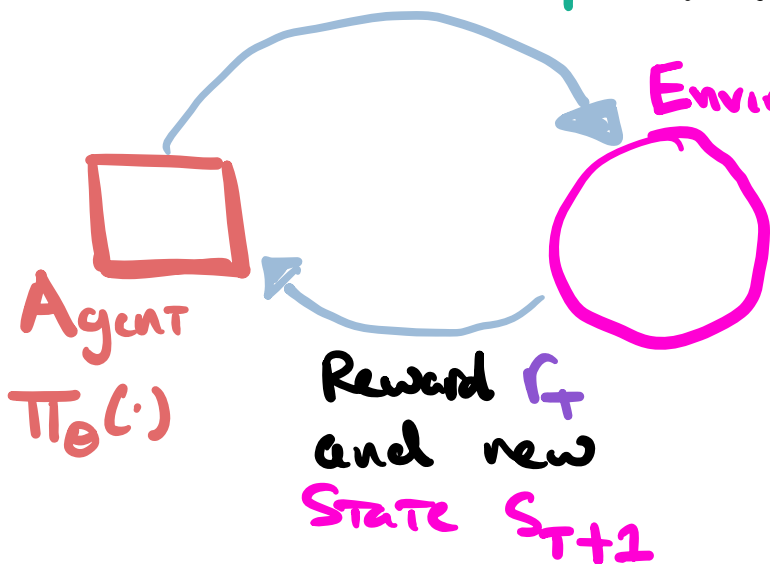Classify this review as positive or

negative. This movie rocked.

(See Slides)

Reinforcement Learning from Human feedback (RLHF)

Why RL? Not clear how to supervise learn our way to avoid TOXIC outputs or produce funnier jokes.
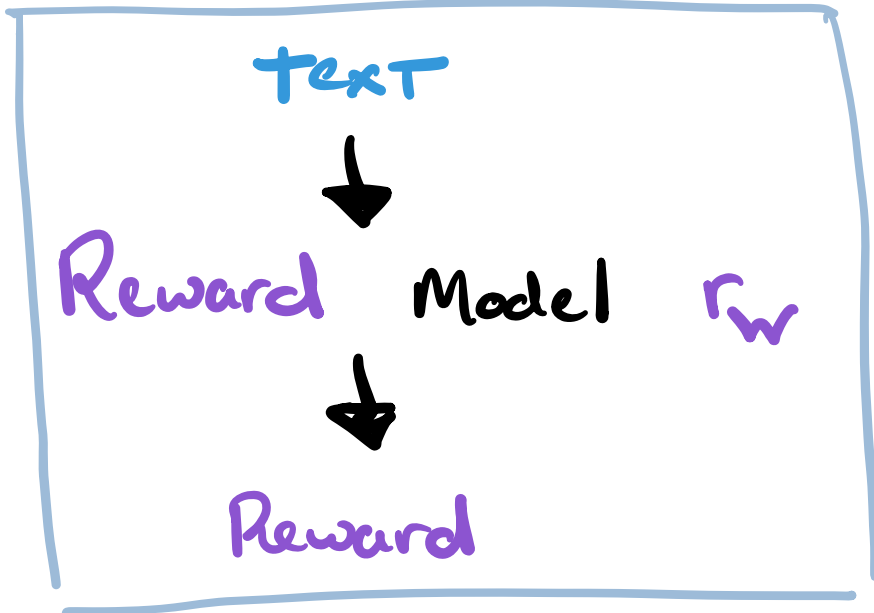
## RL

Take action $a_T$ from State $S_T$

Environment

Agent $\pi_\theta(\cdot)$

Reward $r_T$ and new State $S_{T+1}$

$$R(T) = \sum_{T=1}^{T} \gamma^T r_T$$

Total discounted future reward

$a_T \sim \pi_\theta (s_t)$    Agent policy picks actions to maximize $R(T)$.

Learn reward Model based on human Preference feedack $y_j \succ y_k$

TEXT

↓

Reward Model $r_W$

↓

Reward

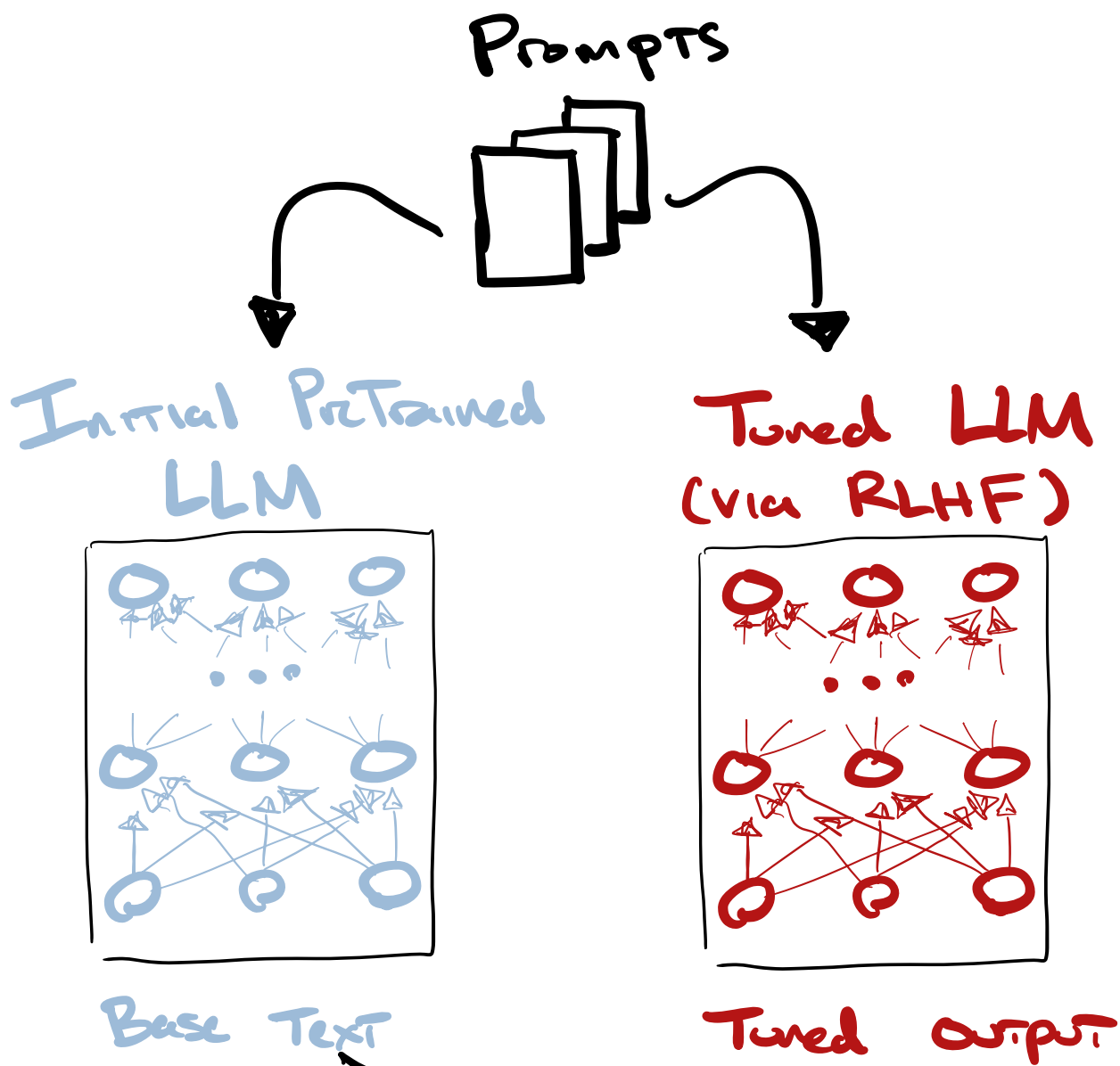Train $r_W$ — Separate from LLM!

$\square \rightarrow$ 🕸️ $\rightarrow \hat{r}$

So $r_W$ provides rewards from environment. States & actions both language.

Once we have $r_w$, update LLM to yield outputs w/ bigger rewards.

**P**roblem The LLM might go off the rails — optimize $r$ but produce jibberish.

Prompts

Initial PreTrained LLM

Tuned LLM (via RLHF)

Base Text

Tuned output

$-\lambda \, KL(\text{Tuned}, \text{Base})$

Penalty for diverging from
Base pretrained LLM

$\nabla_\theta J$ (policy gradient)

Adjust $\theta$
to $\uparrow$ reward

$r_W(\text{Tuned})$

$$\nabla_\theta J \overset{\cdot\cdot}{=}$$

$$\mathbb{E}_{\tau \sim \pi_\theta} \left[ \sum_{\tau=0}^{T} \nabla_\theta \log \pi_\theta(a_\tau | s_\tau) \, r_W(\tau) \right]$$

(Likelihood
under LLM)

Nothing new

full output
Trajectory

Push token probs up in proportion
to reward score

# Direct Preference Optimization (DPO)

But who wants to deal w/
RL? Rafailov et al. (2023)
present DPO to learn directly
from prefence feedback — no
reward model or RL required.

(See Slides)