# CME 295: Transformers &
# Large Language Models
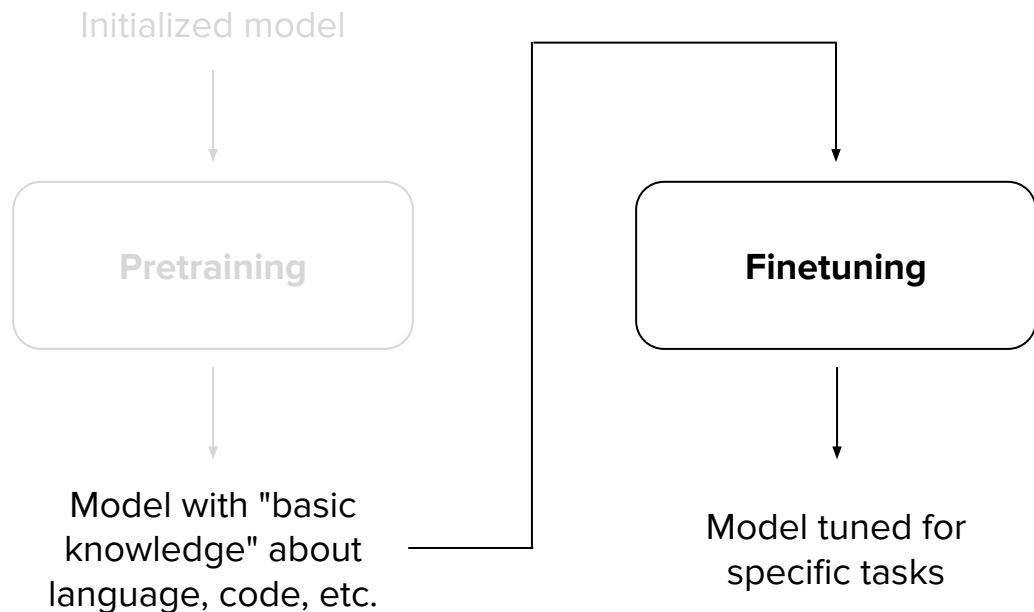
**Afshine Amidi** & **Shervine Amidi**
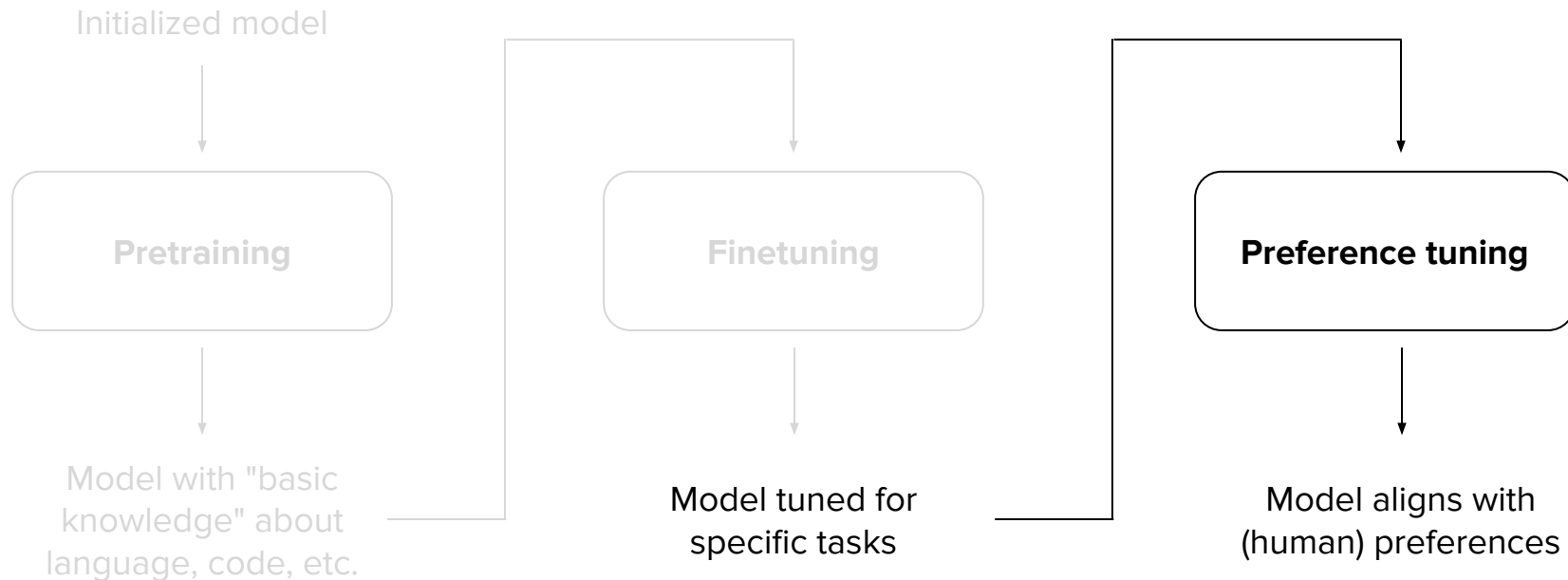
Initialized model

↓

**Pretraining**

↓

Model with "basic knowledge" about language, code, etc.

ICME

# Recap of last episodes...



Initialized model

Pretraining

Finetuning

Model with "basic knowledge" about language, code, etc.

Model tuned for specific tasks

ICME

# Recap of last episodes...



Initialized model

**Pretraining**

Model with "basic knowledge" about language, code, etc.

**Finetuning**

Model tuned for specific tasks

**Preference tuning**

Model aligns with (human) preferences

Stanford University

ICME

# Today's focus



Initialized model

**Pretraining**

Model with "basic knowledge" about language, code, etc.

**Finetuning**

**Model tuned for specific tasks**

**Preference tuning**

**Model aligns with (human) preferences**

ICME

# Transformers & Large Language Models

**Preference tuning**

Data collection

RLHF

DPO

ICME

# Preference tuning

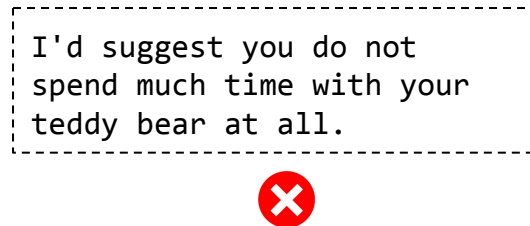**Context**. Model may misbehave. Need to inject negative signals.

```
Suggest a new activity I
could do with my teddy bear.
```
→ LLM$_{SFT}$ →
```
I'd suggest you do not
spend much time with your
teddy bear at all.
```
❌

Stanford University

ICME

# Preference tuning

**Context**. Model may misbehave. Need to inject negative signals.

```
Suggest a new activity I
could do with my teddy bear.
```
→ LLM$_{SFT}$ →
```
I'd suggest you do not
spend much time with your
teddy bear at all.
```
❌

**Idea**. Collect preference pairs and train the model on it.

```
I'd suggest you do not
spend much time with your
teddy bear at all.
```
❌

# Preference tuning

**Context**. Model may misbehave. Need to inject negative signals.

```
Suggest a new activity I
could do with my teddy bear.
```
→ LLM~SFT~ →
```
I'd suggest you do not
spend much time with your
teddy bear at all.
```
❌

**Idea**. Collect preference pairs and train the model on it.

```
Of course! Teddy bears not only make awesome
companions for a delightful sleep, but can
also be great buddies for fun activities. How
about you both watch a movie together?
```
✅

```
I'd suggest you do not
spend much time with your
teddy bear at all.
```
❌

- **Easier** to **compare** (e.g. A better than B) than **generate** (e.g. generate A from scratch)

ICME

# Why preference tuning?

- **Easier** to **compare** (e.g. A better than B) than **generate** (e.g. generate A from scratch)

- **Distribution** is very important for SFT: easy to "mess up"

ICME

- **Easier** to **compare** (e.g. A better than B) than **generate** (e.g. generate A from scratch)

- **Distribution** is very important for SFT: easy to "mess up"

- **Not scalable**: data quality is very important and hard to get

# Why preference tuning?

- **Easier** to **compare** (e.g. A better than B) than **generate** (e.g. generate A from scratch)

- **Distribution** is very important for SFT: easy to "mess up"

- **Not scalable**: data quality is very important and hard to get

However, "model misbehaving" can also be a good wake-up call to check **SFT data quality**

ICME

# Transformers & Large Language Models

Preference tuning

**Data collection**

RLHF

DPO

ICME

Observation = (prompt $x$, response $\widehat{y}$)

# Preference data

Observation = (prompt $x$, response $\widehat{y}$)

**Pointwise**
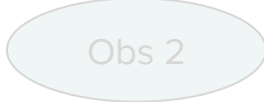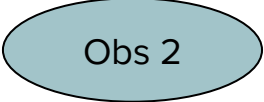
Obs 1     0.4

Obs 2     0.9

Obs 3     0.1

Obs 4     0.2

# Preference data

Observation = (prompt $x$, response $\widehat{y}$)

| Pointwise | | Pairwise | | |
|---|---|---|---|---|
| Obs 1 | 0.4 | Obs 1 | < | Obs 2 |
| Obs 2 | 0.9 | Obs 1 | > | Obs 3 |
| Obs 3 | 0.1 | Obs 1 | > | Obs 4 |
| Obs 4 | 0.2 | Obs 2 | > | Obs 3 |

ICME

# Preference data

Observation = (prompt $x$, response $\widehat{y}$)

| Pointwise | | Pairwise | | | Listwise | |
|---|---|---|---|---|---|---|
| Obs 1 | 0.4 | Obs 1 | < | Obs 2 | Obs 2 | 1 |
| Obs 2 | 0.9 | Obs 1 | > | Obs 3 | Obs 1 | 2 |
| Obs 3 | 0.1 | Obs 1 | > | Obs 4 | Obs 4 | 3 |
| Obs 4 | 0.2 | Obs 2 | > | Obs 3 | Obs 3 | 4 |

# Preference data

Observation = (prompt $x$, response $\widehat{y}$)

1. **Generate** pair of responses $(\widehat{y}_1, \widehat{y}_2)$ for the same prompt $x$

   - Input $x$ via logs / reference distribution
   - Output $\widehat{y}$ via SFT model with $T > 0$ / synthetic / rewrites

ICME

1. **Generate** pair of responses $(\widehat{y}_1, \widehat{y}_2)$ for the same prompt $x$

   - Input $x$ via logs / reference distribution
   - Output $\widehat{y}$ via SFT model with $T > 0$ / synthetic / rewrites

2. **Label** $(x, \widehat{y}_1)$ and $(x, \widehat{y}_2)$

   - Human rating
   - Proxies (e.g. LLM-as-a-judge, BLEU, ROUGE, etc.)
   - Variants: binary scale (better or worse) vs "nuanced" scale

# Transformers & Large Language Models

Preference tuning

Data collection

**RLHF**

DPO

ICME

Agent

Environment

Reward $r_t$

State $s_t$

Action $a_t$

Policy $\pi_\theta(a_t | s_t)$

ICME

LLM

State $s_t$

Action $a_t$

Reward $r_t$

Environment

Policy $\pi_\theta(a_t|s_t)$

Stanford University

ICME

LLM

**Input so far** $s_t$

Reward $r_t$

Action $a_t$

Environment

Policy $\pi_\theta(a_t | s_t)$

LLM

Input so far $s_t$

**Next token** $a_t$

Reward $r_t$

Environment

Policy $\pi_\theta(a_t|s_t)$

LLM

**Tokens**

Reward $r_t$

Input so far $s_t$

Next token $a_t$

Policy $\pi_\theta(a_t|s_t)$

LLM

Tokens

Input so far $s_t$

Next token $a_t$

Reward $r_t$

**Probability of next token** $\pi_\theta(a_t|s_t)$

LLM

**Human preference** $r_t$

Input so far $s_t$

Next token $a_t$

Tokens

Probability of next token $\pi_\theta(a_t|s_t)$

**Idea**. Learn $\theta$ so that $\pi_\theta$ aligns with human preferences



LLM

Human
preference $r_t$

Input so far $s_t$

Next token $a_t$

Tokens

Probability of next token $\pi_\theta(a_t|s_t)$

**RLHF** = **R**einforcement **L**earning from **H**uman **F**eedback

*"Training language models to follow instructions with human feedback", Ouyang et al., 2022.*

# **RLHF** = **R**einforcement **L**earning from **H**uman **F**eedback

**Step 1 – Reward modeling**: Distinguish good from bad!

- Input: (prompt $x$, response $\widehat{y}$)

- Output: quantitative score $r(x, \widehat{y})$

*"Training language models to follow instructions with human feedback", Ouyang et al., 2022.*

ICME

**RLHF** = **R**einforcement **L**earning from **H**uman **F**eedback

**Step 1 – Reward modeling**: Distinguish good from bad!

- Input: (prompt $x$, response $\widehat{y}$)
- Output: quantitative score $r(x, \widehat{y})$

**Step 2 – Reinforcement learning**: Align the model!

- Input: prompt $x$
- Output: response $\widehat{y}$

# Step 1: Reward modeling

**Idea**. Know which answers are bad and which are good via **Reward Model**.

*"Training language models to follow instructions with human feedback", Ouyang et al., 2022.*

**Idea**. Know which answers are bad and which are good via **Reward Model**.

```
Suggest a new activity I could do with my teddy bear.
```
prompt

ICME

**Idea**. Know which answers are bad and which are good via **Reward Model**.

```
Suggest a new activity I could do with my teddy bear.
```
prompt

Of course! Teddy bears not only make awesome companions for a delightful sleep, but can also be great buddies for fun activities. How about you both watch a movie together?

→ RM → 👍

*"Training language models to follow instructions with human feedback", Ouyang et al., 2022.*

**Idea**. Know which answers are bad and which are good via **Reward Model**.

```
Suggest a new activity I could do with my teddy bear.
```
prompt

Of course! Teddy bears not only make awesome companions for a delightful sleep, but can also be great buddies for fun activities. How about you both watch a movie together?

RM 👍

I'd suggest you do not spend much time with your teddy bear at all.

RM 👎

*"Training language models to follow instructions with human feedback", Ouyang et al., 2022.*

ICME

**Bradley-Terry formulation**. Probability that $y_i$ better than $y_j$ is **defined** as:

$$p(y_i > y_j) = \frac{e^{r_i}}{e^{r_i} + e^{r_j}} = \sigma(r_i - r_j)$$

where $r_i, r_j$ rewards of $y_i, y_j$ respectively

*"The rank analysis of incomplete block designs: I. the method of paired comparisons"*, Bradley et al., 1952.

ICME

**Bradley-Terry formulation**. Probability that $y_i$ better than $y_j$ is **defined** as:

$$p(y_i > y_j) = \frac{e^{r_i}}{e^{r_i} + e^{r_j}} = \sigma(r_i - r_j)$$

where $r_i, r_j$ rewards of $y_i, y_j$ respectively



$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

*"The rank analysis of incomplete block designs: I. the method of paired comparisons", Bradley et al., 1952.*   ICME

**Training**. Learn $r$ based on **pairwise** preference data

$$(x, \hat{y}_w) \longrightarrow \boxed{\text{RM}} \longrightarrow r(x, \hat{y}_w)$$
$$\updownarrow \mathscr{L}$$
$$(x, \hat{y}_l) \longrightarrow \phantom{\boxed{\text{RM}}} \longrightarrow r(x, \hat{y}_l)$$

**Training**. Learn $r$ based on **pairwise** preference data

$$(x, \hat{y}_w) \longrightarrow \boxed{\text{RM}} \longrightarrow r(x, \hat{y}_w)$$

$$\updownarrow \mathscr{L}$$

$$(x, \hat{y}_l) \longrightarrow \boxed{\text{RM}} \longrightarrow r(x, \hat{y}_l)$$

$$\boxed{\mathscr{L}(\theta) = -\mathbb{E}\big[\log(\sigma(r(x, \widehat{y}_w) - r(x, \widehat{y}_l)))\big]}$$

**Data**.

- O(10,000) observations
- label = human rating (which is where the "HF" from RLHF comes from)

**Data**.

- O(10,000) observations
- label = human rating (which is where the "HF" from RLHF comes from)

**Model**.

- Pretrained LLM with classification head (instead of next token prediction)
- Encoder-only: BERT and the like via [CLS] projection

Stanford University     *"RewardBench: Evaluating Reward Models for Language Modeling", Lambert et al., 2024.*     ICME

# Step 2: Reinforcement learning

**Idea**. Change weights of LLM to penalize bad answers and promote good answers via **Reinforcement Learning** using the **Reward Model**.

*"Training language models to follow instructions with human feedback", Ouyang et al., 2022.*

ICME

**Idea**. Change weights of LLM to penalize bad answers and promote good answers via **Reinforcement Learning** using the **Reward Model**.

```
┌ ─ ─ ─ ─ ─ ─ ─ ┐
  Suggest a new
  activity I           ┌──────────┐
  could do with  ──→   │   LLM    │
  my teddy bear.       └──────────┘
└ ─ ─ ─ ─ ─ ─ ─ ┘
```

**Idea**. Change weights of LLM to penalize bad answers and promote good answers via **Reinforcement Learning** using the **Reward Model**.

```
Suggest a new
activity I          →    LLM    →    I'd suggest you
could do with                        do not spend
my teddy bear.                       much time with
                                     your teddy bear
                                     at all.
```

*"Training language models to follow instructions with human feedback", Ouyang et al., 2022.*

ICME

**Idea**. Change weights of LLM to penalize bad answers and promote good answers via **Reinforcement Learning** using the **Reward Model**.



*"Training language models to follow instructions with human feedback", Ouyang et al., 2022.*

**Idea**. Change weights of LLM to penalize bad answers and promote good answers via **Reinforcement Learning** using the **Reward Model**.

**Idea**. Change weights of LLM to penalize bad answers and promote good answers via **Reinforcement Learning** using the **Reward Model**.



via RL

*"Training language models to follow instructions with human feedback", Ouyang et al., 2022.*

**Idea**. Change weights of LLM to penalize bad answers and promote good answers via **Reinforcement Learning** using the **Reward Model**.



🔥 **Trained**

```
Suggest a new
activity I
could do with
my teddy bear.
```

LLM

```
I'd suggest you
do not spend
much time with
your teddy bear
at all.
```

❄️ **Frozen**

RM

via RL

*"Training language models to follow instructions with human feedback", Ouyang et al., 2022.*

# Step 2: Reinforcement learning

**Idea**. Change weights of LLM to penalize bad answers and promote good answers via **Reinforcement Learning** using the **Reward Model**.



Suggest a new activity I could do with my teddy bear. → LLM → I'd suggest you do not spend much time with your teddy bear at all. → RM → 👎

via RL

Objective function optimizes for **higher rewards without going too far** from the **base model**

*"Training language models to follow instructions with human feedback", Ouyang et al., 2022.*

ICME

**Data**.

- O(100,000) observations
- label = score given by reward model

*"Training language models to follow instructions with human feedback", Ouyang et al., 2022.*

**Data**.

- O(100,000) observations

- label = score given by reward model

**Model**. Initialized at SFT model

*"Training language models to follow instructions with human feedback", Ouyang et al., 2022.*

**Data**.

- O(100,000) observations
- label = score given by reward model

**Model**. Initialized at SFT model

**Training**. Change weights of policy (LLM) via objective function:

$$\mathscr{L}(\theta) = \boxed{\text{Maximize rewards}} + \boxed{\begin{array}{c}\text{Don't deviate too much}\\\text{from base model}\end{array}}$$

*"Training language models to follow instructions with human feedback", Ouyang et al., 2022.*

ICME

**Data**.

- O(100,000) observations
- label = score given by reward model

**Model**. Initialized at SFT model

**Training**. Change weights of policy (LLM) via objective function:

Avoid **"reward hacking"** + **training instability**

$$\mathscr{L}(\theta) = \boxed{\text{Maximize rewards}} + \boxed{\begin{array}{c}\text{Don't deviate too much}\\\text{from base model}\end{array}}$$

*"Training language models to follow instructions with human feedback", Ouyang et al., 2022.*

ICME

**PPO** = **P**roximal **P**olicy **O**ptimization

$$\mathscr{L}(\theta) = \boxed{\text{Maximize rewards}} + \boxed{\begin{array}{c}\text{Don't deviate too much}\\\text{from base model}\end{array}}$$

*"Proximal Policy Optimization Algorithms", Schulman et al., 2017.*

ICME

**PPO** = **P**roximal **P**olicy **O**ptimization

$$\mathscr{L}(\theta) = \boxed{\text{Maximize rewards}} + \boxed{\text{Don't deviate too much from base model}}$$

*"Proximal Policy Optimization Algorithms", Schulman et al., 2017.*

**PPO** = **P**roximal **P**olicy **O**ptimization

$$\mathscr{L}(\theta) = -\left[r(x,\widehat{y}) - \lambda\mathrm{KL}(\pi_\theta(\widehat{y}|x)||\pi_{\mathrm{ref}}(\widehat{y}|x))\right]$$

**PPO** = **P**roximal **P**olicy **O**ptimization

$$\mathscr{L}(\theta) = -\left[ \underline{r(x,\widehat{y})} - \underline{\lambda\mathrm{KL}(\pi_\theta(\widehat{y}|x)||\pi_{\mathrm{ref}}(\widehat{y}|x))} \right]$$

Maximize rewards

Don't deviate too much from base model

*"Training language models to follow instructions with human feedback", Ouyang et al., 2022.*

**ICME**

**PPO** = **P**roximal **P**olicy **O**ptimization

$$\mathscr{L}(\theta) = -\left[ r(x, \widehat{y}) - \lambda \mathrm{KL}(\pi_\theta(\widehat{y}|x) || \pi_{\mathrm{ref}}(\widehat{y}|x)) \right]$$



$$\mathrm{KL}(P||Q) = \sum_{i=1}^{n} p_i \log\left(\frac{p_i}{q_i}\right)$$

*Figure from "Super Study Guide: Transformers and Large Language Models", Amidi et al., 2024.*

$$\mathscr{L}(\theta) = \boxed{\text{Maximize rewards}} + \boxed{\begin{array}{c}\text{Don't deviate too much}\\\text{from base model}\end{array}}$$

$$\mathscr{L}(\theta) = \boxed{\text{Maximize } \textbf{advantages}} + \boxed{\text{Don't deviate too much from base model}}$$

**Advantage** ~ Reward - Baseline

# PPO actually computes advantages (and not just rewards)

$\mathscr{L}(\theta) =$ ⟦ Maximize **advantages** ⟧ $+$ ⟦ Don't deviate too much from base model ⟧

**Advantage** ~ Reward - Baseline ← Value function

$$\mathscr{L}(\theta) = \boxed{\text{Maximize \textbf{advantages}}} + \boxed{\text{Don't deviate too much from base model}}$$

$$\boxed{\textbf{Advantage} \sim \text{Reward - Baseline}} \longleftarrow \quad \text{Value function}$$

**Value function**.

- Token-level

- What would be the reward if follow the policy

- Trained jointly with policy

- Label = reward

$$\mathcal{L}(\theta) = \boxed{\text{Maximize } \textbf{advantages}} + \boxed{\begin{array}{c}\text{Don't deviate too much}\\\text{from base model}\end{array}}$$

$$\boxed{\textbf{Advantage} \sim \text{Reward - Baseline}} \longleftarrow \quad \text{Value function}$$

**Value function.**

- Token-level

- What would be the reward if follow the policy

- Trained jointly with policy

- Label = reward

GAE method

ICME

**Idea**. Clip ratio between new and old policy to prevent large updates

$$L^{CLIP}(\theta) = \hat{\mathbb{E}}_t \left[ \min(r_t(\theta)\hat{A}_t, \operatorname{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)\hat{A}_t) \right]$$

$$\text{with} \quad r_t(\theta) = \frac{\pi_\theta(a_t \mid s_t)}{\pi_{\theta_{\text{old}}}(a_t \mid s_t)}$$

*"Proximal Policy Optimization Algorithms"*, *Schulman et al., 2017.*

ICME

**Idea**. Clip ratio between new and old policy to prevent large updates

$$L^{CLIP}(\theta) = \hat{\mathbb{E}}_t \Big[ \min(r_t(\theta)\hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)\hat{A}_t) \Big]$$

$$\text{with} \quad r_t(\theta) = \frac{\pi_\theta(a_t \mid s_t)}{\pi_{\theta_{\text{old}}}(a_t \mid s_t)}$$

**Terminology**. Confusing since it is an objective function ("**maximize**") and NOT a loss.

*"Proximal Policy Optimization Algorithms"*, *Schulman et al., 2017.*

ICME

**Idea**. Clip ratio between new and old policy to prevent large updates

$$L^{CLIP}(\theta) = \hat{\mathbb{E}}_t \left[ \min(r_t(\theta)\hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)\hat{A}_t) \right]$$

with $\quad r_t(\theta) = \dfrac{\pi_\theta(a_t \mid s_t)}{\pi_{\theta_{\text{old}}}(a_t \mid s_t)}$

**Terminology**. Confusing since rewards noted "r".
Here we are talking about the ratio.

**Idea**. Clip ratio between new and old policy to prevent large updates

$$L^{CLIP}(\theta) = \hat{\mathbb{E}}_t \left[ \min(r_t(\theta)\hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)\hat{A}_t) \right]$$

with $\quad r_t(\theta) = \frac{\pi_\theta(a_t \mid s_t)}{\pi_{\theta_{\text{old}}}(a_t \mid s_t)}$

$A > 0$ $\qquad\qquad$ $A < 0$



*"Proximal Policy Optimization Algorithms"*, Schulman et al., 2017.

ICME

**Idea**. Clip ratio between new and old policy to prevent large updates

$$L^{CLIP}(\theta) = \hat{\mathbb{E}}_t\Big[\min(r_t(\theta)\hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)\hat{A}_t)\Big]$$

with $\quad r_t(\theta) = \dfrac{\pi_\theta(a_t \mid s_t)}{\pi_{\theta_{\text{old}}}(a_t \mid s_t)}$

**Idea**. Penalize difference in policy distributions

$$L^{KLPEN}(\theta) = \hat{\mathbb{E}}_t \left[ \frac{\pi_\theta(a_t \mid s_t)}{\pi_{\theta_{\text{old}}}(a_t \mid s_t)} \hat{A}_t - \beta \, \text{KL}[\pi_{\theta_{\text{old}}}(\cdot \mid s_t), \pi_\theta(\cdot \mid s_t)] \right]$$

**Idea**. Penalize difference in policy distributions

$$L^{KLPEN}(\theta) = \hat{\mathbb{E}}_t\left[\frac{\pi_\theta(a_t \mid s_t)}{\pi_{\theta_{\text{old}}}(a_t \mid s_t)}\hat{A}_t - \beta\,\mathrm{KL}[\pi_{\theta_{\text{old}}}(\cdot \mid s_t), \pi_\theta(\cdot \mid s_t)]\right]$$

**Idea**. Penalize difference in policy distributions

$$L^{KLPEN}(\theta) = \hat{\mathbb{E}}_t\left[\frac{\pi_\theta(a_t \mid s_t)}{\pi_{\theta_{\text{old}}}(a_t \mid s_t)}\hat{A}_t - \beta\,\mathrm{KL}[\pi_{\theta_{\text{old}}}(\cdot \mid s_t), \pi_\theta(\cdot \mid s_t)]\right]$$

**Terminology**.

- old = model from previous RL iteration
- ref = base model

*"Proximal Policy Optimization Algorithms", Schulman et al., 2017.*

ICME

**Idea**. Penalize difference in policy distributions

$$L^{KLPEN}(\theta) = \hat{\mathbb{E}}_t \left[ \frac{\pi_\theta(a_t \mid s_t)}{\pi_{\theta_{\text{old}}}(a_t \mid s_t)} \hat{A}_t - \beta \, \mathrm{KL}[\pi_{\theta_{\text{ref}}}(\cdot \mid s_t), \pi_\theta(\cdot \mid s_t)] \right]$$

Nowadays, KL divergence is with respect to **ref** (base model)

*"Proximal Policy Optimization Algorithms", Schulman et al., 2017.*

**Limitations of PPO**.

- Need 4 models (policy, value, reward model, base)

- Is it worth it?

ICME

**Limitations of PPO**.

- Need 4 models (policy, value, reward model, base)

- Is it worth it?

**Variants**.

- REINFORCE

- GRPO

- ... and many more!

- Requires training a reward model (**2-stage process**)

# Challenges with RL-based approach

- Requires training a reward model (**2-stage process**)

- Many hyperparameters to tune

ICME

- Requires training a reward model (**2-stage process**)

- Many hyperparameters to tune

- Training instability

ICME

# Challenges with RL-based approach

- Requires training a reward model (**2-stage process**)

- Many hyperparameters to tune

- Training instability

- Metric to monitor training

Stanford University

ICME

- Requires training a reward model (**2-stage process**)

- Many hyperparameters to tune

- Training instability

- Metric to monitor training

- Need diversity in completions!

Stanford University

ICME

# Challenges with RL-based approach

- Requires training a reward model (**2-stage process**)

- Many hyperparameters to tune

- Training instability

- Metric to monitor training

- Need diversity in completions!

- Not abundantly clear why preference tuning absolutely needs RL

ICME

**BoN** = **B**est **o**f **N**

**BoN** = **B**est **o**f **N**

**Idea**. Skip the RL step and leverage the reward model scores

ICME

**BoN** = **B**est **o**f **N**

**Idea**. Skip the RL step and leverage the reward model scores

**Strategy**.

- Given a prompt, generate several outputs with SFT model

- Rank output with score given by reward model

- Take the best one

Suggest a new
activity I
could do with
my teddy
bear.

# BoN in action

Suggest a new activity I could do with my teddy bear. → LLM$_{SFT}$

# BoN in action

Suggest a new activity I could do with my teddy bear.

→ LLM_SFT →

Of course! Teddy bears not only make awesome companions for a delightful sleep, but can also be great buddies for fun activities. How about you both watch a movie together?

I'd suggest you do not spend much time with your teddy bear at all.

Take your teddy bear on a picnic in your backyard.

# BoN in action



Suggest a new activity I could do with my teddy bear.

LLM$_{SFT}$

Of course! Teddy bears not only make awesome companions for a delightful sleep, but can also be great buddies for fun activities. How about you both watch a movie together?

RM

I'd suggest you do not spend much time with your teddy bear at all.

RM

Take your teddy bear on a picnic in your backyard.

RM

ICME

# BoN in action

# BoN in action



Suggest a new activity I could do with my teddy bear.

$LLM_{SFT}$

Of course! Teddy bears not only make awesome companions for a delightful sleep, but can also be great buddies for fun activities. How about you both watch a movie together?

RM → **0.8** **1**

I'd suggest you do not spend much time with your teddy bear at all.

RM → **-2** **3**

Take your teddy bear on a picnic in your backyard.

RM → **0.2** **2**

ICME

# BoN in action

Suggest a new activity I could do with my teddy bear.

→ LLM$_{SFT}$ →

Of course! Teddy bears not only make awesome companions for a delightful sleep, but can also be great buddies for fun activities. How about you both watch a movie together?

→ RM → **0.8** — ①

I'd suggest you do not spend much time with your teddy bear at all.

→ RM → **-2** — ③

Take your teddy bear on a picnic in your backyard.

→ RM → **0.2** — ②

ICME

# Transformers & Large Language Models

Preference tuning

Data collection

RLHF

**DPO**

ICME

- Limitations using RL

$$L^{KLPEN}(\theta) = \hat{\mathbb{E}}_t \left[ \frac{\pi_\theta(a_t \mid s_t)}{\pi_{\theta_{\text{old}}}(a_t \mid s_t)} \hat{A}_t - \beta \, \mathrm{KL}[\pi_{\theta_{\text{old}}}(\cdot \mid s_t), \pi_\theta(\cdot \mid s_t)] \right]$$

- Limitations using RL

- Best-of-N is costly at inference time

- Limitations using RL

- Best-of-N is costly at inference time

- **Why don't we train in a supervised fashion?**

**DPO** = **D**irect **P**reference **O**ptimization

Rewrite the **loss function** in a supervised way:

$$\mathcal{L}_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[ \log \sigma \left( \beta \log \frac{\pi_\theta(y_w \mid x)}{\pi_{\text{ref}}(y_w \mid x)} - \beta \log \frac{\pi_\theta(y_l \mid x)}{\pi_{\text{ref}}(y_l \mid x)} \right) \right]$$

*"Direct Preference Optimization: Your Language Model is Secretly a Reward Model"*, *Rafailov et al., 2023.*

ICME

**DPO** = **D**irect **P**reference **O**ptimization

Rewrite the **loss function** in a supervised way:

$$\mathcal{L}_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[ \log \sigma \left( \beta \log \frac{\pi_\theta(y_w \mid x)}{\pi_{\text{ref}}(y_w \mid x)} - \beta \log \frac{\pi_\theta(y_l \mid x)}{\pi_{\text{ref}}(y_l \mid x)} \right) \right]$$

No r(x, y)!

- No need to train a separate reward model

*"Direct Preference Optimization: Your Language Model is Secretly a Reward Model", Rafailov et al., 2023.*

ICME

## DPO = Direct Preference Optimization

Rewrite the **loss function** in a supervised way:

$$\mathcal{L}_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[ \log \sigma \left( \beta \log \frac{\pi_\theta(y_w \mid x)}{\pi_{\text{ref}}(y_w \mid x)} - \beta \log \frac{\pi_\theta(y_l \mid x)}{\pi_{\text{ref}}(y_l \mid x)} \right) \right]$$

- No need to train a separate reward model

- Operates directly on preference data

*"Direct Preference Optimization: Your Language Model is Secretly a Reward Model", Rafailov et al., 2023.*

ICME

## **DPO** = **D**irect **P**reference **O**ptimization

Rewrite the **loss function** in a supervised way:

$$\mathcal{L}_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}) = -\mathbb{E}_{(x,y_w,y_l)\sim\mathcal{D}}\left[\log\sigma\left(\beta\log\frac{\pi_\theta(y_w\mid x)}{\pi_{\text{ref}}(y_w\mid x)} - \beta\log\frac{\pi_\theta(y_l\mid x)}{\pi_{\text{ref}}(y_l\mid x)}\right)\right]$$

- No need to train a separate reward model    $r_\theta(x,y) = \beta\frac{\pi_\theta(y|x)}{\pi_{\text{ref}}(y|x)}$

- Operates directly on preference data

- Similar to the Bradley-Terry formulation with a special kind of reward!

*"Direct Preference Optimization: Your Language Model is Secretly a Reward Model", Rafailov et al., 2023.*  ICME

**DPO** = **D**irect **P**reference **O**ptimization

Rewrite the **loss function** in a supervised way:

$$\mathcal{L}_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[ \log \sigma \left( \beta \log \frac{\pi_\theta(y_w \mid x)}{\pi_{\text{ref}}(y_w \mid x)} - \beta \log \frac{\pi_\theta(y_l \mid x)}{\pi_{\text{ref}}(y_l \mid x)} \right) \right]$$

$$\mathcal{L}_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[ \log \sigma \left( r_\theta(x, y_w) - r_\theta(x, y_l) \right) \right]$$

*"Direct Preference Optimization: Your Language Model is Secretly a Reward Model", Rafailov et al., 2023.*

ICME

① **Start from PPO objective**

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

*"Direct Preference Optimization: Your Language Model is Secretly a Reward Model", Rafailov et al., 2023.*

① Start from PPO objective

② **Derive optimal policy**

$$\pi^*(y \mid x) = \frac{1}{Z(x)} \pi_{\text{ref}}(y \mid x) \exp\left(\frac{1}{\beta} r^*(x, y)\right)$$

*"Direct Preference Optimization: Your Language Model is Secretly a Reward Model", Rafailov et al., 2023.*

ICME

①  Start from PPO objective

②  Derive optimal policy

③  **Identify a "reward" term**

$$r^*(x, y) = \beta \log \frac{\pi^*(y \mid x)}{\pi_{\text{ref}}(y \mid x)} + \beta \log Z(x)$$

*"Direct Preference Optimization: Your Language Model is Secretly a Reward Model", Rafailov et al., 2023.*

ICME

# Where does the DPO formulation come from?

① Start from PPO objective

② Derive optimal policy

③ Identify a "reward" term

④ **Write Bradley-Terry formulation for this "reward"**

$$p^*(y_w \succ y_\ell \mid x) = \frac{1}{1 + \exp\left(\beta \log \frac{\pi^*(y_\ell \mid x)}{\pi_{\text{ref}}(y_\ell \mid x)} - \beta \log \frac{\pi^*(y_w \mid x)}{\pi_{\text{ref}}(y_w \mid x)}\right)}$$

*"Direct Preference Optimization: Your Language Model is Secretly a Reward Model", Rafailov et al., 2023.*   ICME

①  Start from PPO objective

②  Derive optimal policy

③  Identify a "reward" term

④  Write Bradley-Terry formulation for this "reward"

⑤  **"Infer" DPO loss function**

$$\mathcal{L}_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[ \log \sigma \left( \beta \log \frac{\pi_\theta(y_w \mid x)}{\pi_{\text{ref}}(y_w \mid x)} - \beta \log \frac{\pi_\theta(y_l \mid x)}{\pi_{\text{ref}}(y_l \mid x)} \right) \right]$$

*"Direct Preference Optimization: Your Language Model is Secretly a Reward Model", Rafailov et al., 2023.*

ICME

# Use PPO-based RLHF or DPO?

**Ease of implementation**.

| RLHF | DPO |
|------|-----|
| <ul><li>Multi-stage training</li><li>Needs extra models: reward model, value model, base model</li></ul> | <ul><li>Supervised learning</li><li>Base model is the only extra model needed</li></ul> |

**Performance**. No common absolute consensus. Varies from task to task and sensitive to implementation.

*"Is DPO Superior to PPO for LLM Alignment? A Comprehensive Study", Xu et al., 2024.*     ICME

Can I put my teddy bear in the washer?

$\rightarrow$

**Pretrained** + **instruction tuned LLM**
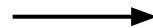
$\rightarrow$

No, it might get damaged. Try hand washing instead.

ICME

Can I put my teddy bear in the washer?

→

**Pretrained** + **instruction tuned** + **preference tuned LLM**

→

**?**

ICME

Can I put my teddy bear in the washer?

→

**Pretrained** + **instruction tuned** + **preference tuned LLM**

→

It's better not to. Your teddy could get hurt! A gentle hand wash is safer.

Thank you for your attention!