

CME 295: Transformers & Large Language Models



Afshin Amidi & Shervine Amidi



Recap of last episodes...

Initialized model

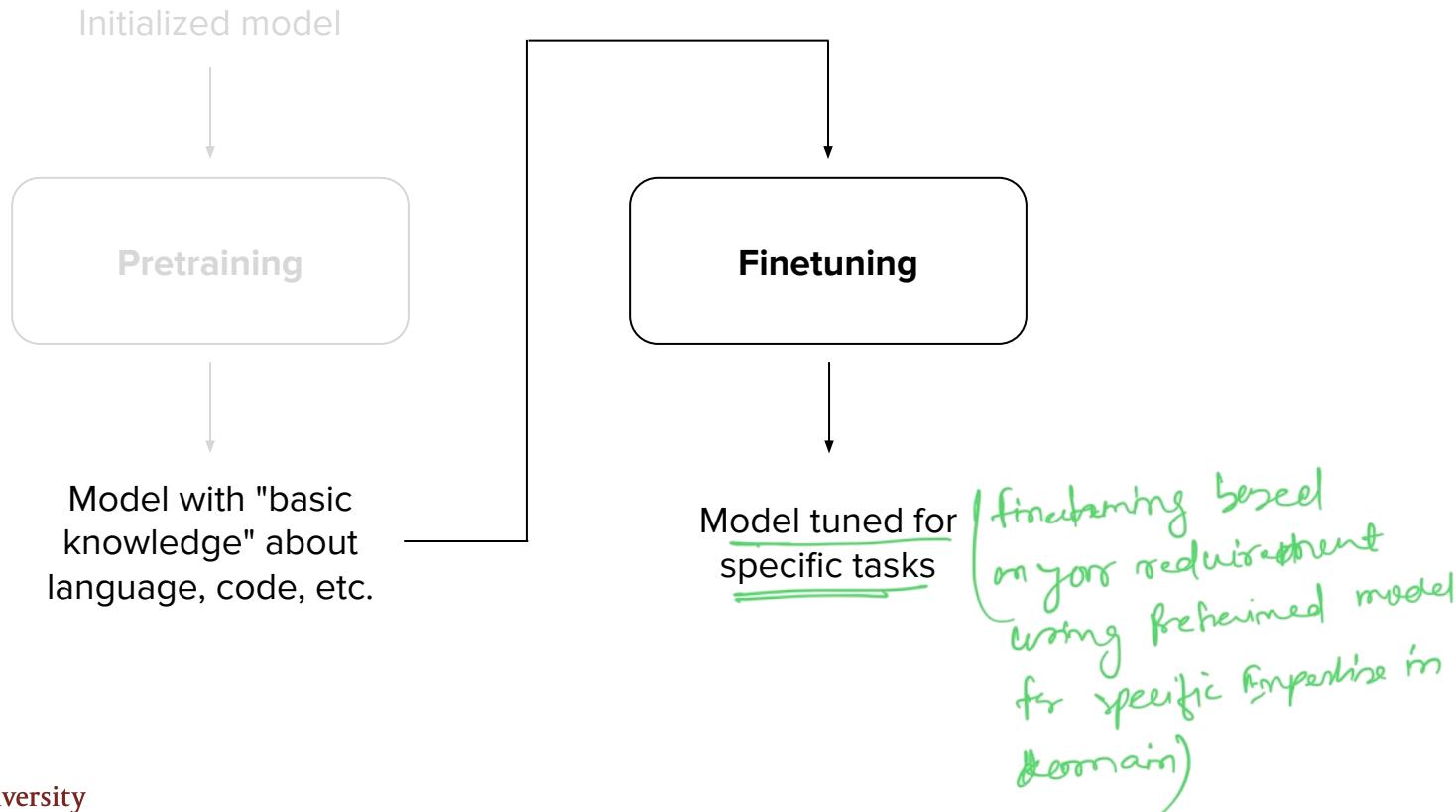


Pretraining

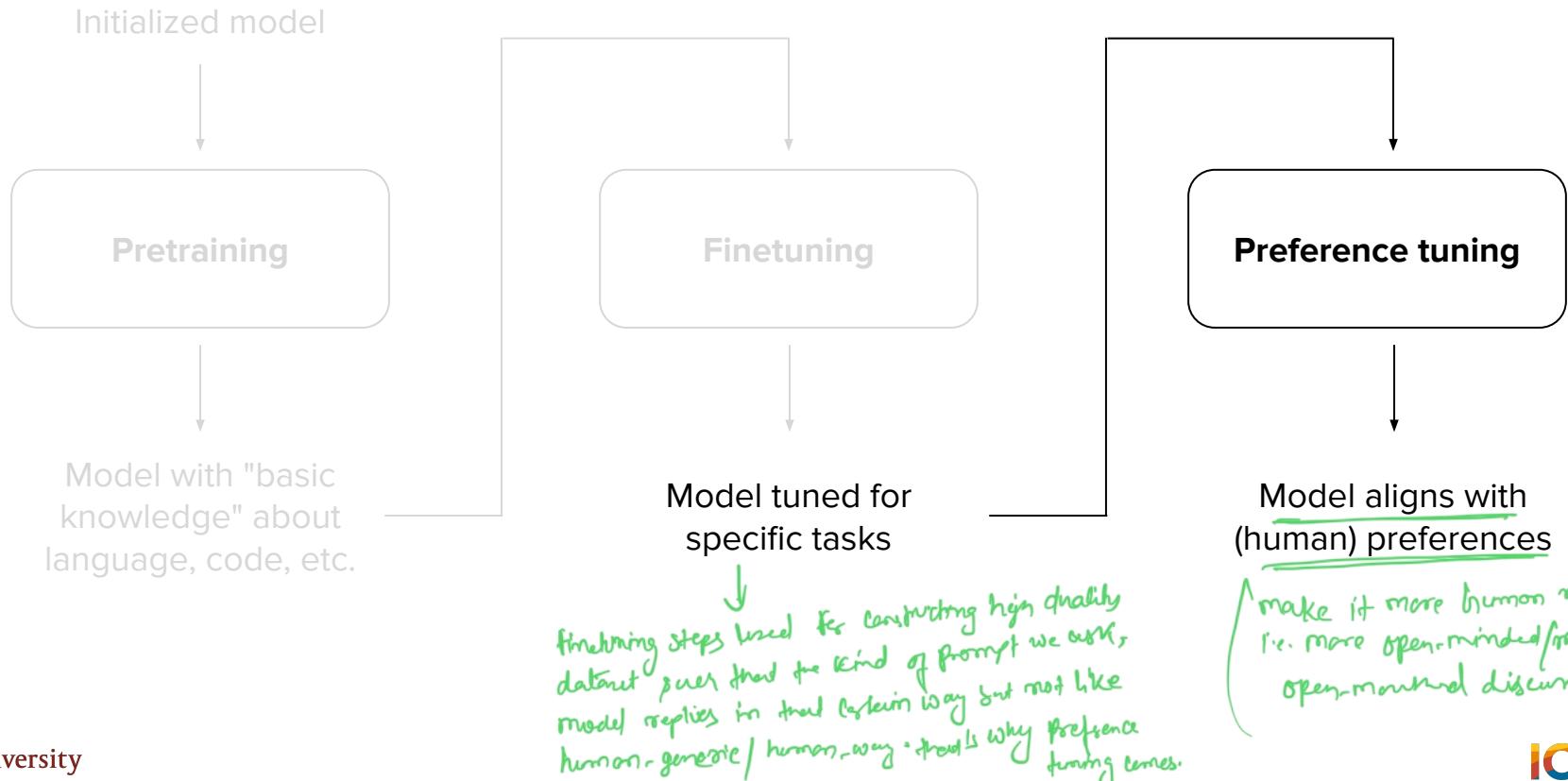


Model with "basic
knowledge" about
language, code, etc.

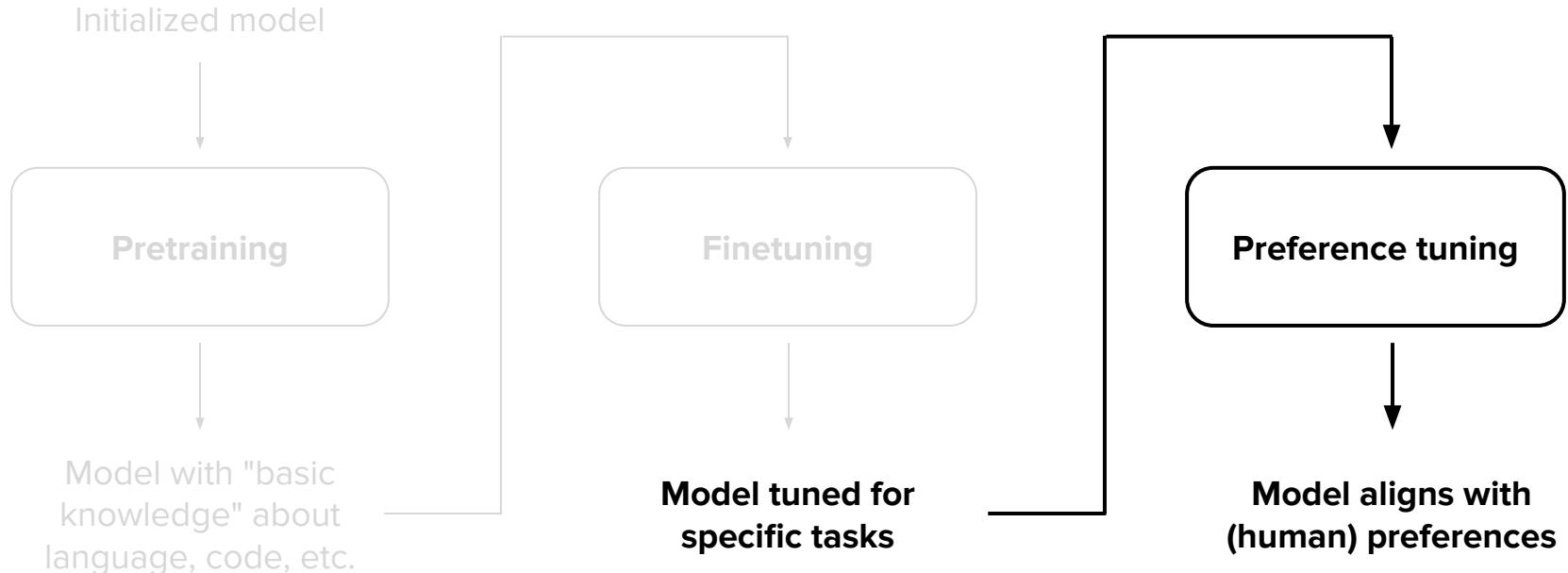
Recap of last episodes...



Recap of last episodes...



Today's focus





Transformers & Large Language Models

Preference tuning

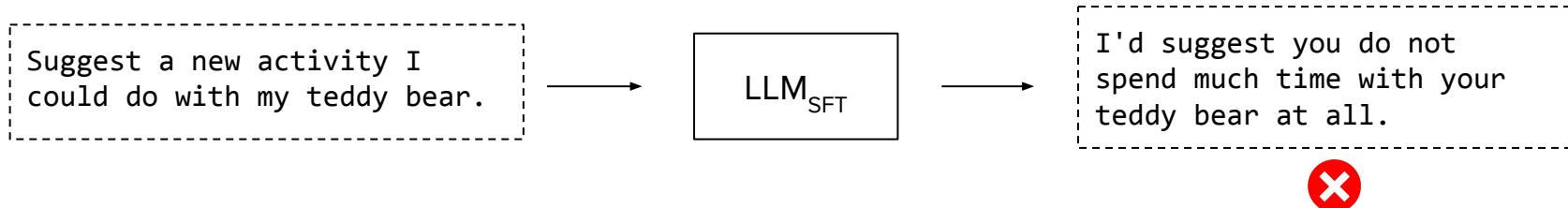
Data collection

RLHF

DPO

Preference tuning

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

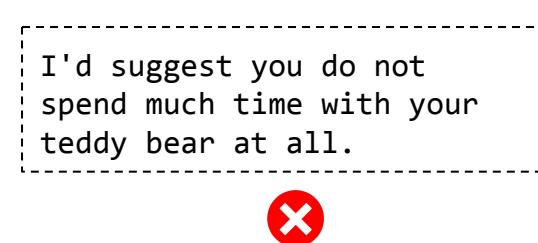


Preference tuning

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



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



Preference tuning

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



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

two response and optimal and better responses transfer further training.

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.



Why preference tuning?

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

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"

supervised fine-tuning [high quality data]

↓
*i.e. more "prompt" helps to write the response
human-like.*

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

Some time with help of SFT,
we do get high quality data
but not so human generated.

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**



Transformers & Large Language Models

Preference tuning

Data collection

RLHF

DPO

Preference data

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

// few ways to construct
preference data

(using pairwise / pointwise)

↓
having multiple
opp and point
it with value
which says which
is relevant

↓
otherwise
comparision
noise
which
response
is better

pairwise

Preference data

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

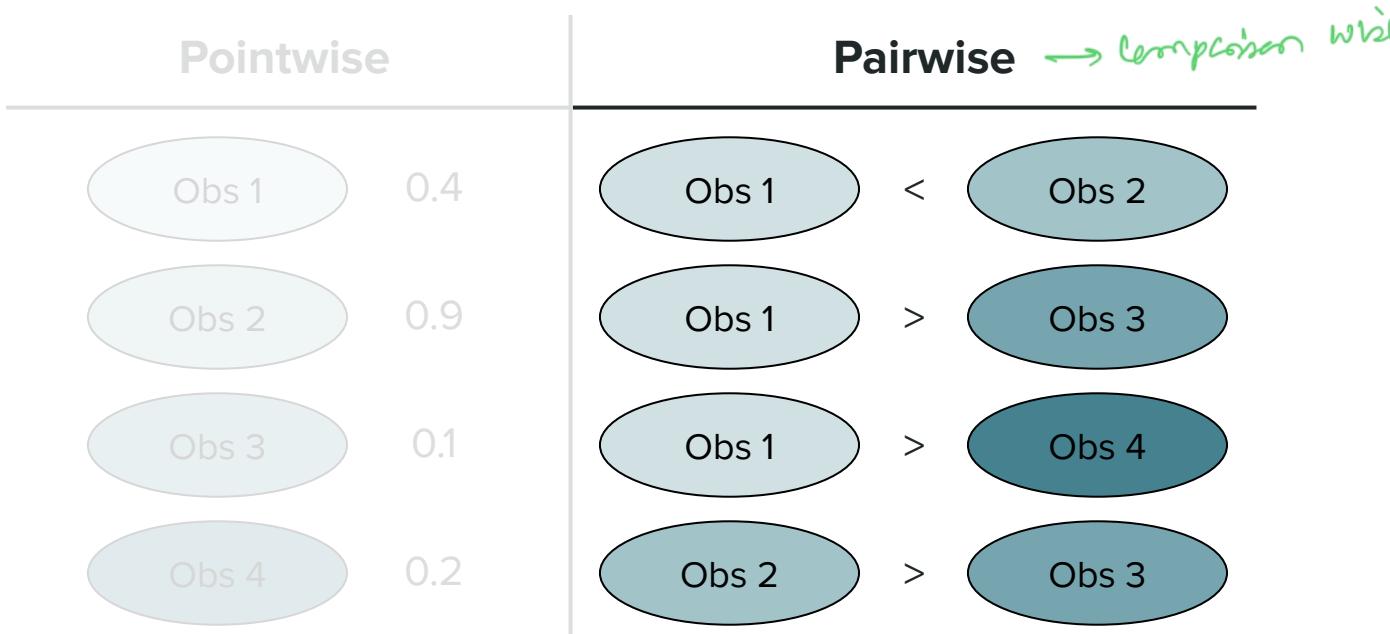
Pointwise

Each observation has points / value

Obs 1	0.4
Obs 2	0.9
Obs 3	0.1
Obs 4	0.2

Preference data

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



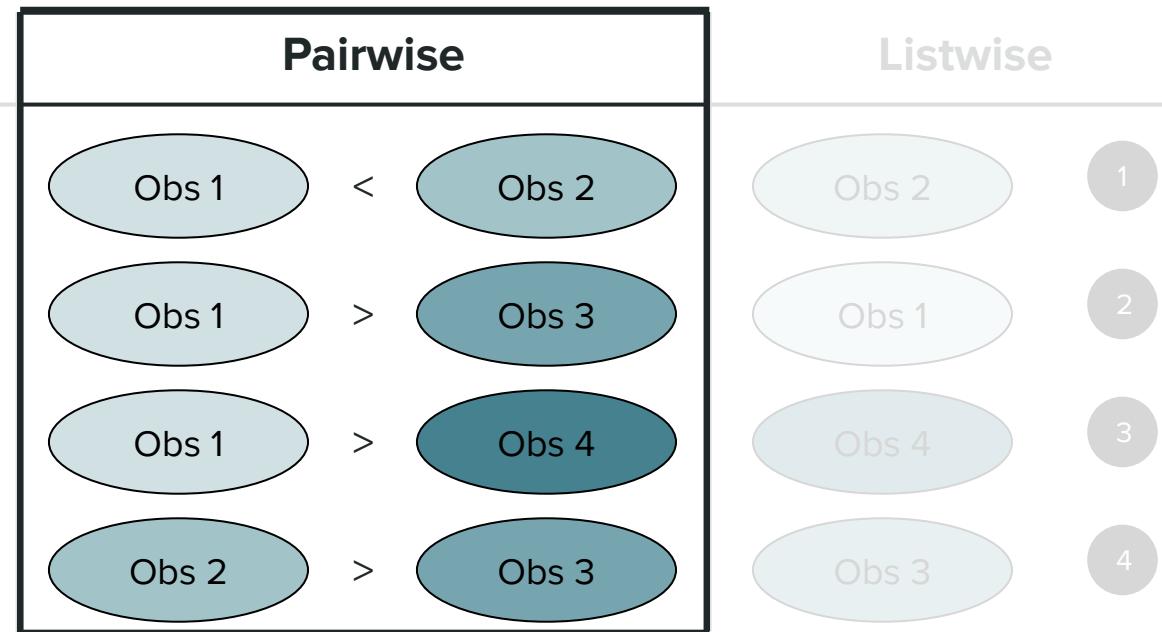
Preference data

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



Preference data

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



Recipe to get (pairwise) preference data

1. **Generate** pair of responses (\hat{y}_1, \hat{y}_2) for the same prompt x
 - Input x via logs / reference distribution
 - Output \hat{y} via SFT model with $T > 0$ / synthetic / rewrites



generate different response with increasing/decreasing temperature which helps to have multiple observation and used for pairwise comparison.

Recipe to get (pairwise) preference data

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

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

2. Label (x, \hat{y}_1) and (x, \hat{y}_2)

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

} Labeling / computing /
for better preference
learning.

binary scale of preference data.



Transformers & Large Language Models

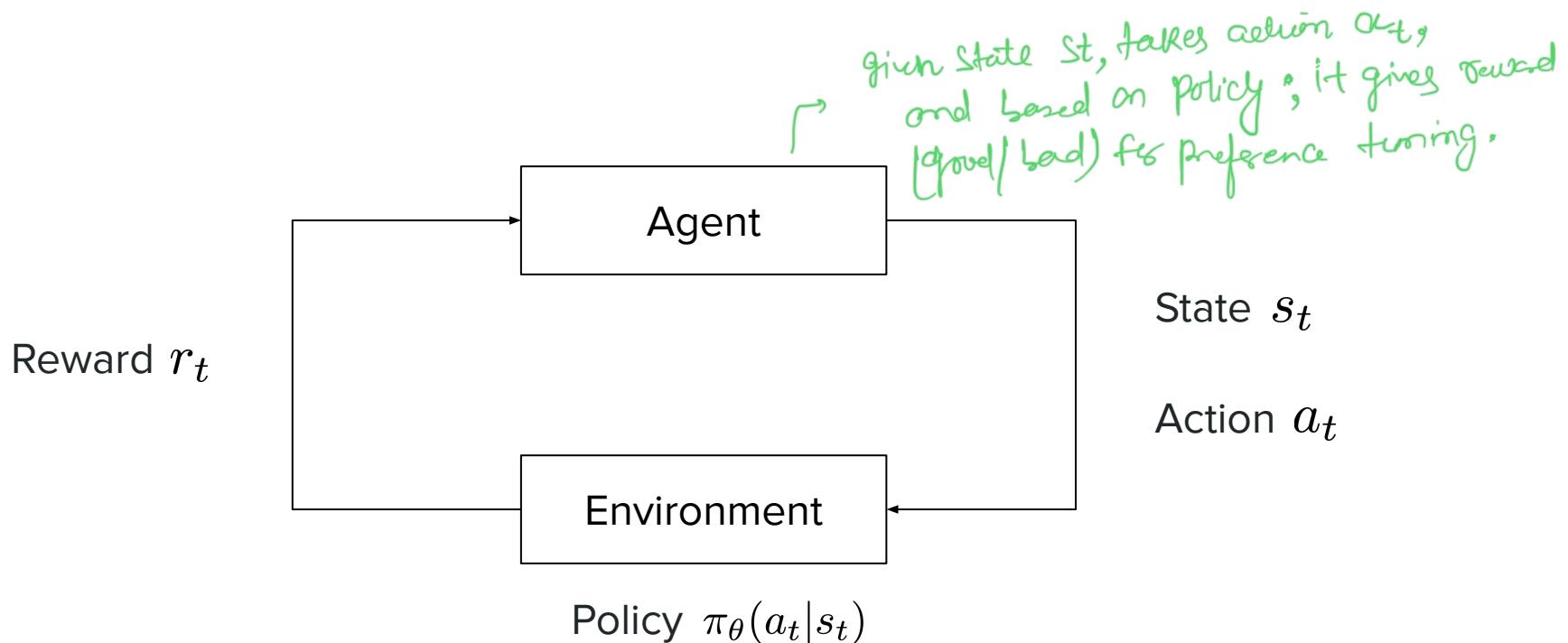
Preference tuning → *Used for human-like response*

Data collection → *Preparing data for preference tuning*

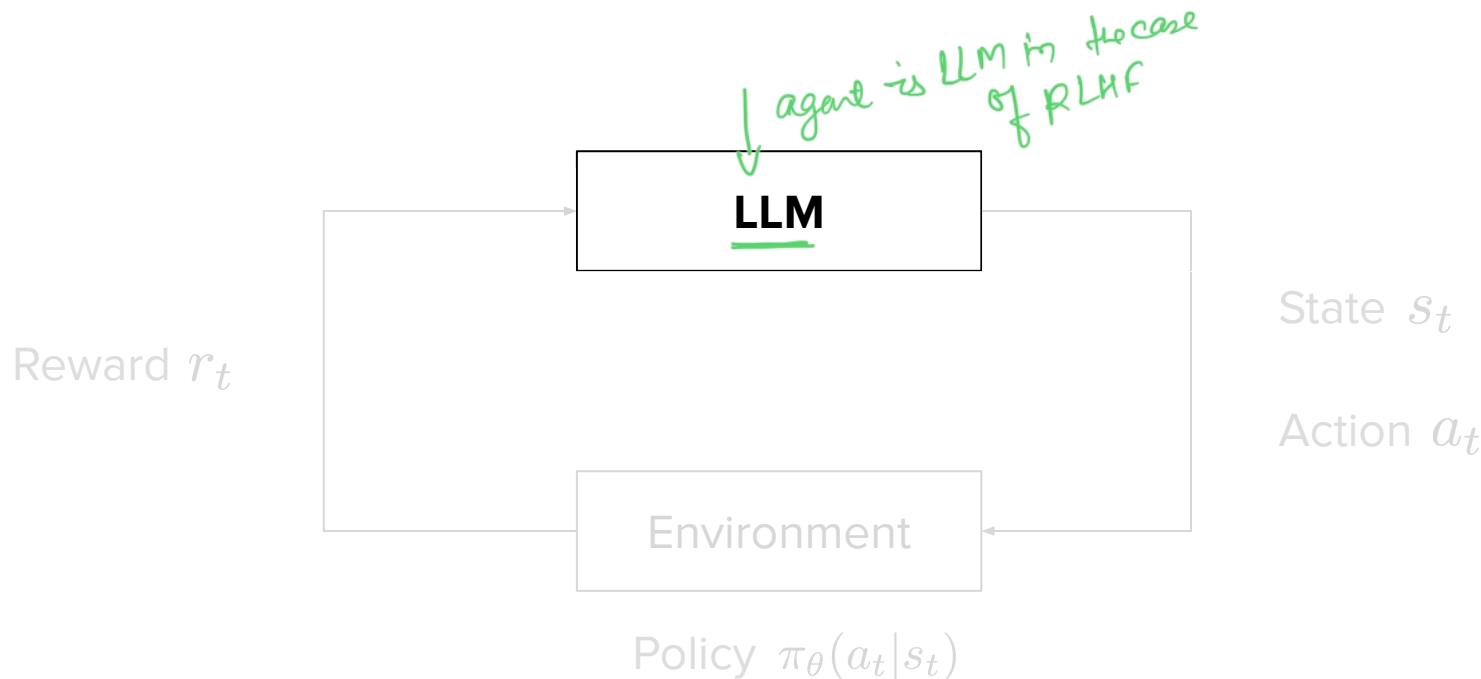
RLHF ← *Reinforcement Learning from Human feedback*

DPO

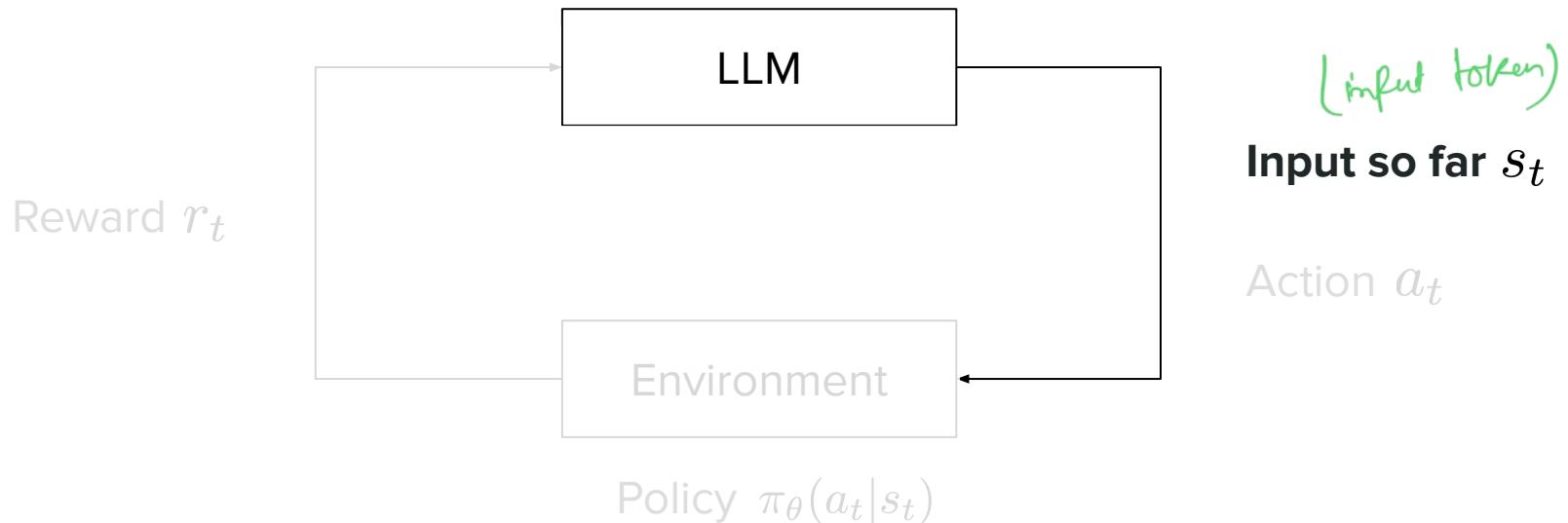
RL formulation



RL formulation for LLMs



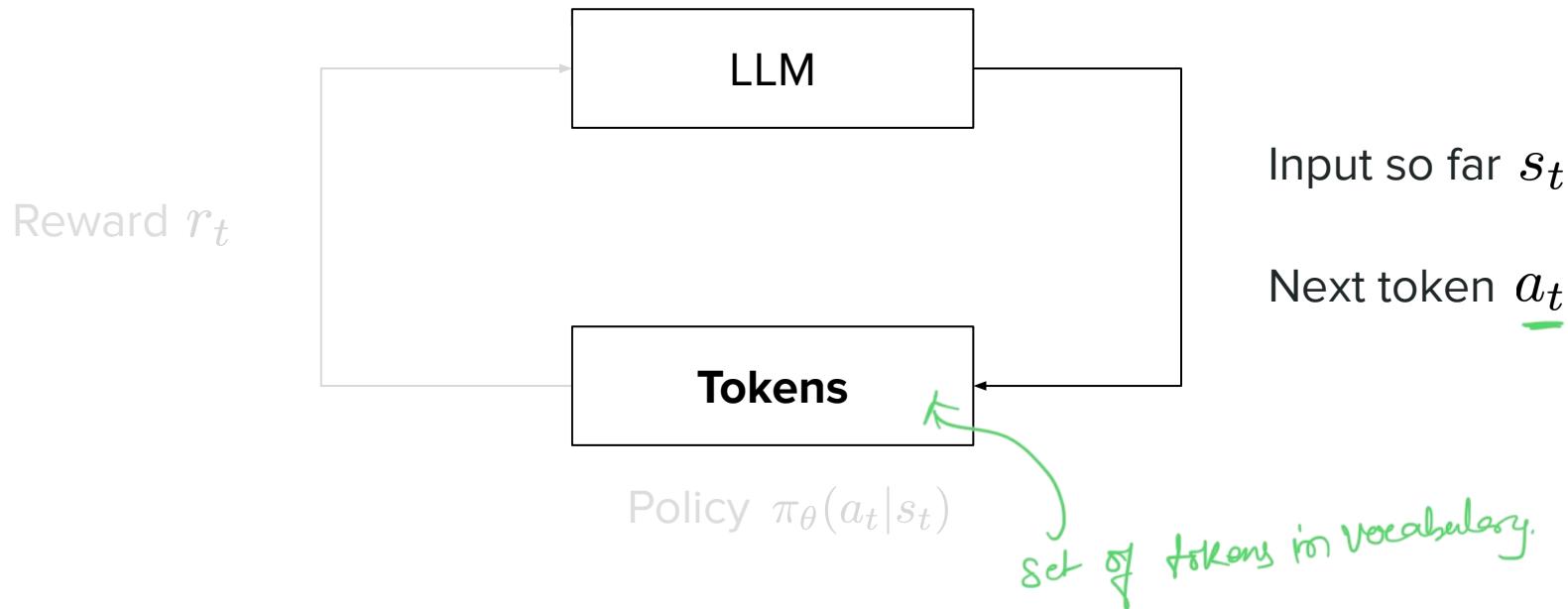
RL formulation for LLMs



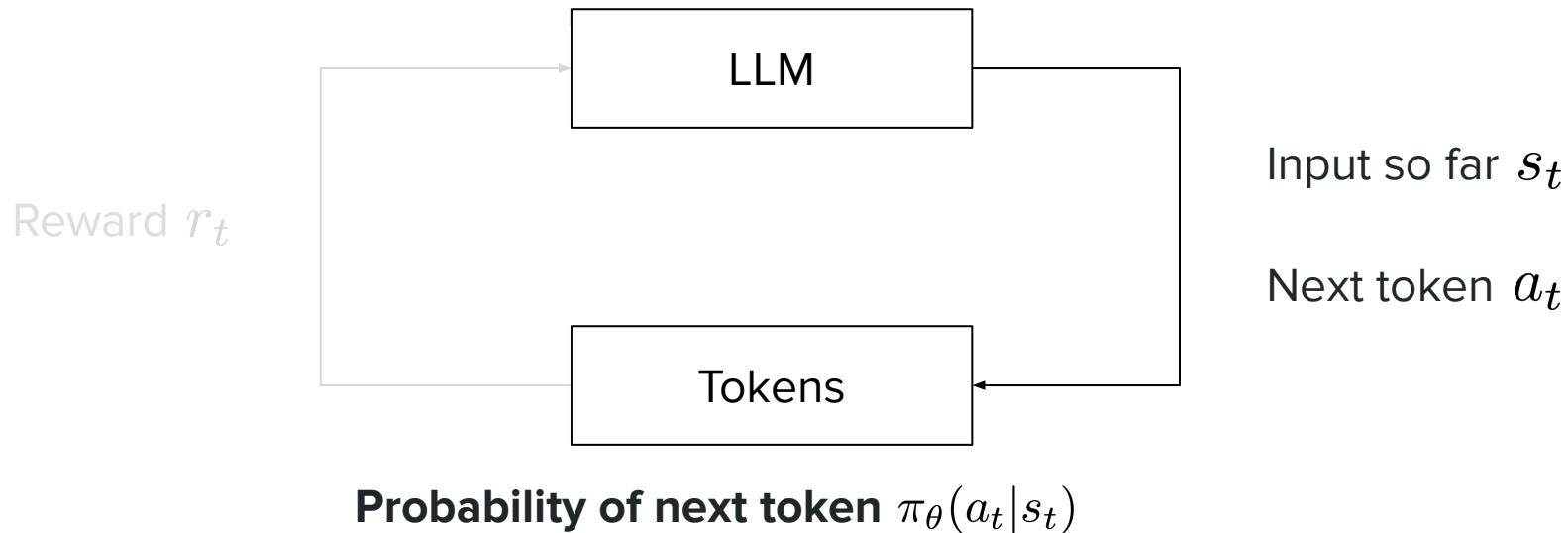
RL formulation for LLMs



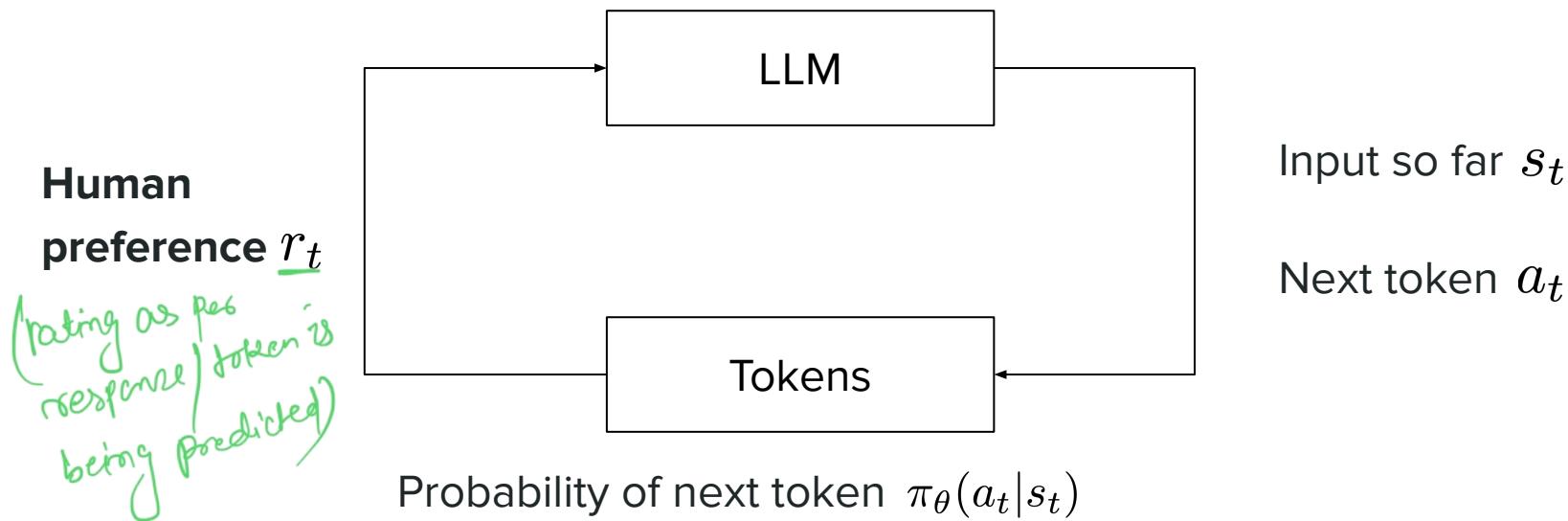
RL formulation for LLMs



RL formulation for LLMs

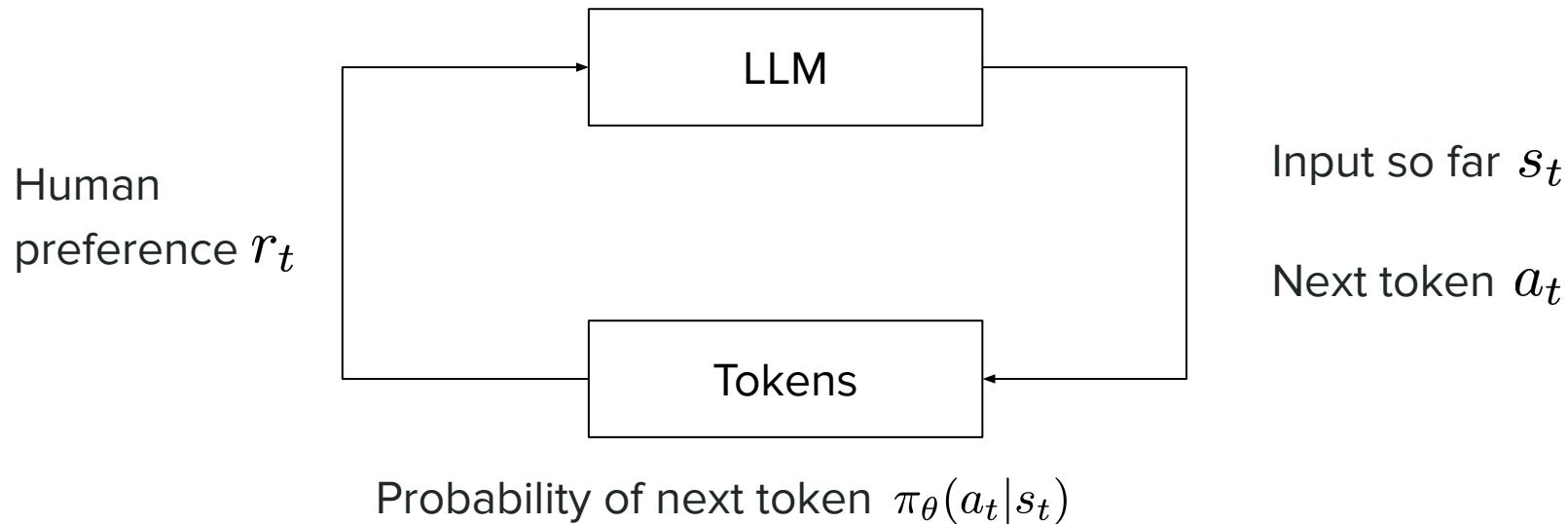


RL formulation for LLMs



RL formulation for LLMs

Idea. Learn θ so that π_θ aligns with human preferences



RLHF overview

RLHF = Reinforcement Learning from Human Feedback

RLHF overview

RLHF = Reinforcement Learning from Human Feedback

Step 1 – Reward modeling: Distinguish good from bad!

- Input: (prompt x , response \hat{y}) – *input: Concatenation of prompt & response.*
- Output: quantitative score $r(x, \hat{y})$ → *output: score*

*↓
gives prompt and response generated by
LLM (agent) how good the output is
being produced.*

RLHF overview

RLHF = Reinforcement Learning from Human Feedback



Step 1 – Reward modeling: Distinguish good from bad!

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



Step 2 – Reinforcement learning: Align the model!

- Input: prompt x
- Output: response \hat{y}

↓
used for tuning the model/
aligning the model
(input with response)

Step 1: Reward modeling

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


model which decides
which off is good/bad
based on input

Step 1: Reward modeling

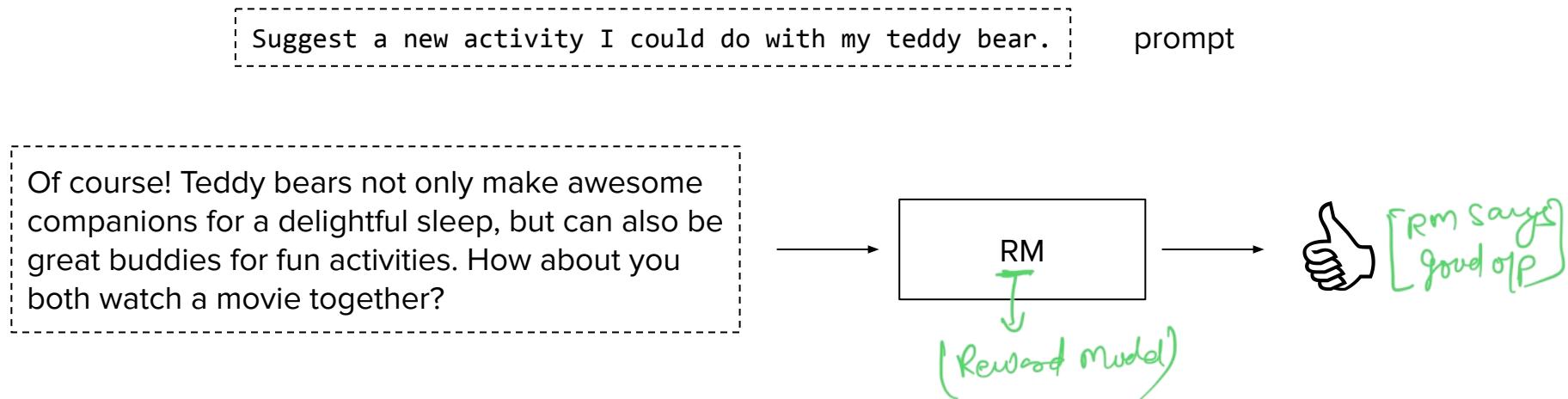
Idea. Know which answers are bad and which are good via **Reward Model**.

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

prompt

Step 1: Reward modeling

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



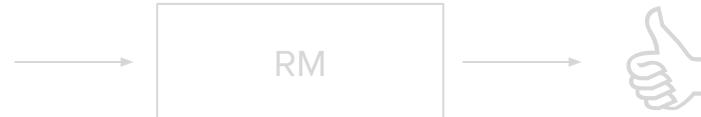
Step 1: Reward modeling

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

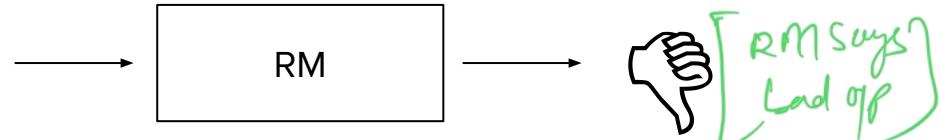
Suggest a new activity I could do with my teddy bear.

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?



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



Step 1: Reward modeling

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

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

↑
function which decide
which is better
↑
 y_i is better than y_j

↑
reward function.

↓
Sigmoid function
which range [0,1]

where r_i, r_j rewards of y_i, y_j respectively

Step 1: Reward modeling

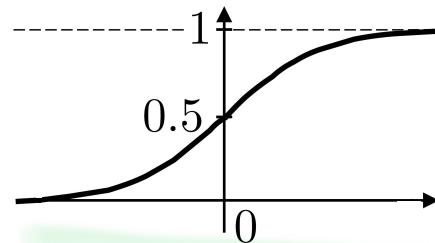
Bradley-Terry formulation. Probability that $\underline{y_i}$ better than $\underline{y_j}$ is **defined** as:

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

$$r_i = f(x, y_w)$$

$$r_j = f(x, y_e)$$

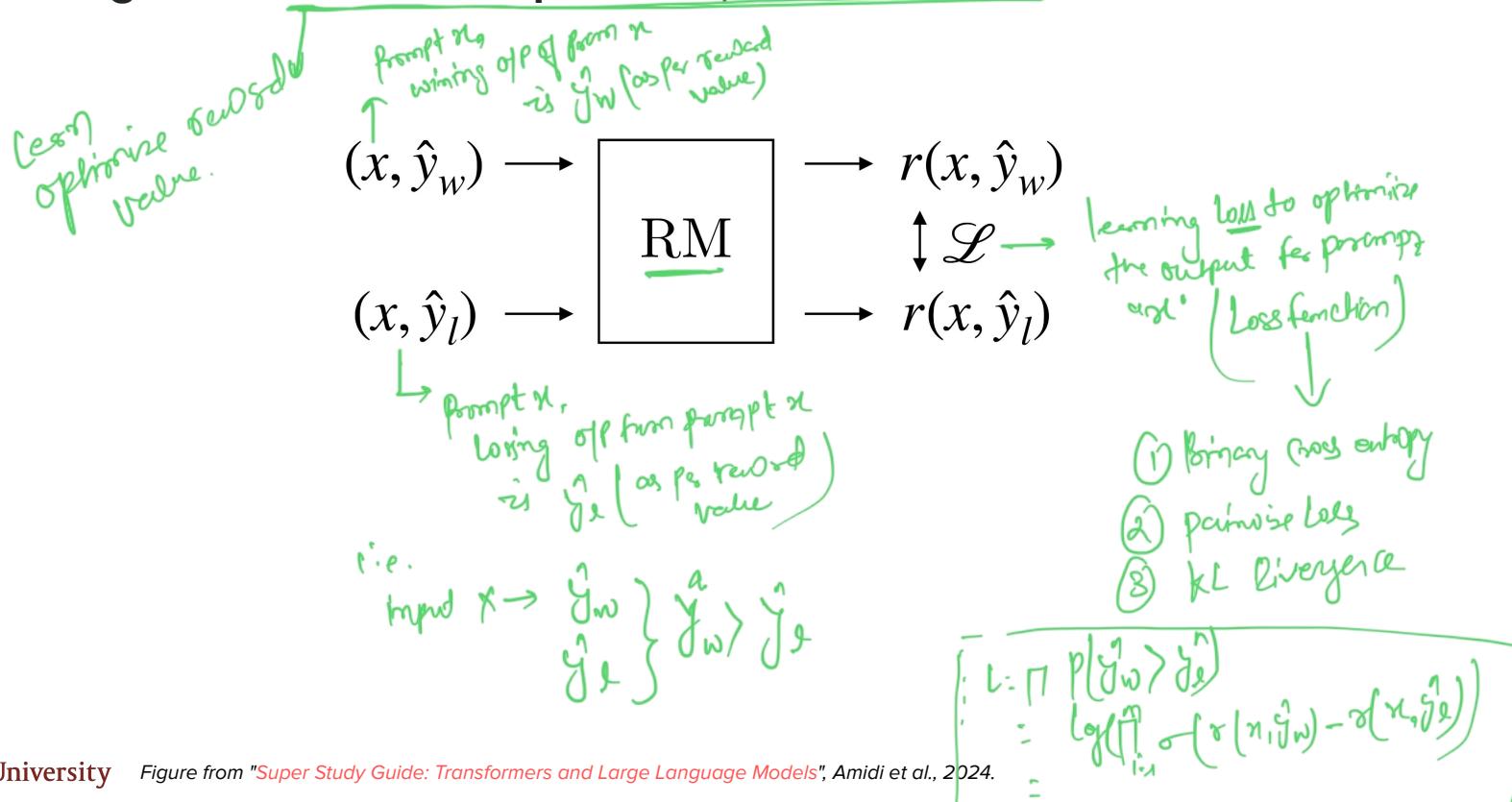
where r_i, r_j rewards of y_i, y_j respectively



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

Step 1: Reward modeling

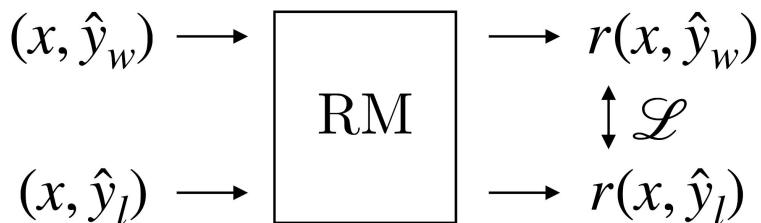
Training. Learn r based on **pairwise** preference data



Step 1: Reward modeling

$$:= \sum \log(\sigma(r(x, \hat{y}_w) - r(x, \hat{y}_l)))$$

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



$$\mathcal{L}(\theta) = -\mathbb{E} [\log(\sigma(r(x, \hat{y}_w) - r(x, \hat{y}_l))))]$$

Loss function:
↑ hyperscore for winning example
↓ loss score for losing example.

Step 1: Reward modeling

Data.

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

Step 1: Reward modeling

Data.

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

Model.

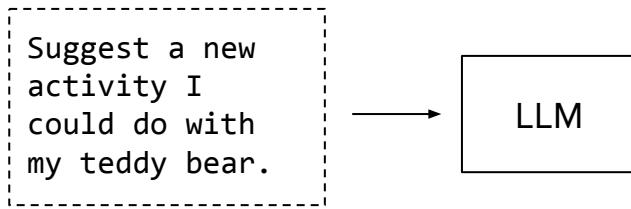
- Pretrained LLM with classification head (instead of next token prediction)
- Encoder-only: BERT and the like via [CLS] projection

Step 2: Reinforcement learning

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

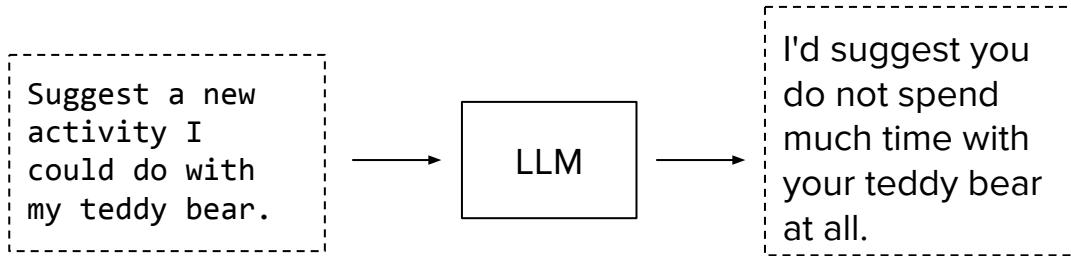
Step 2: Reinforcement learning

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



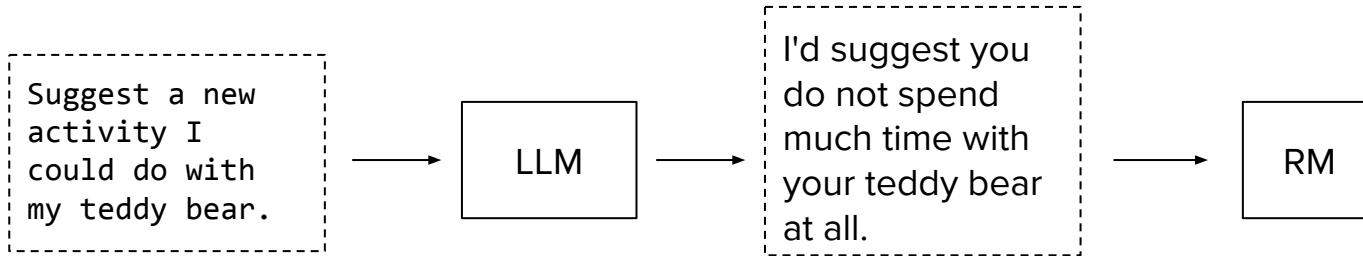
Step 2: Reinforcement learning

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



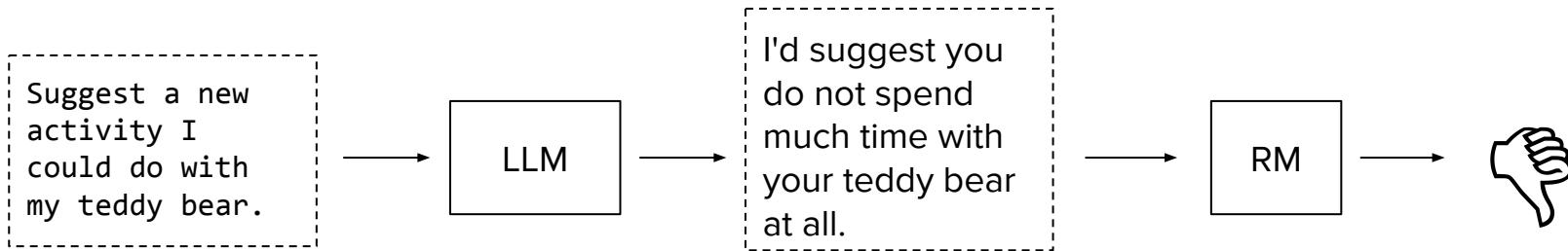
Step 2: Reinforcement learning

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



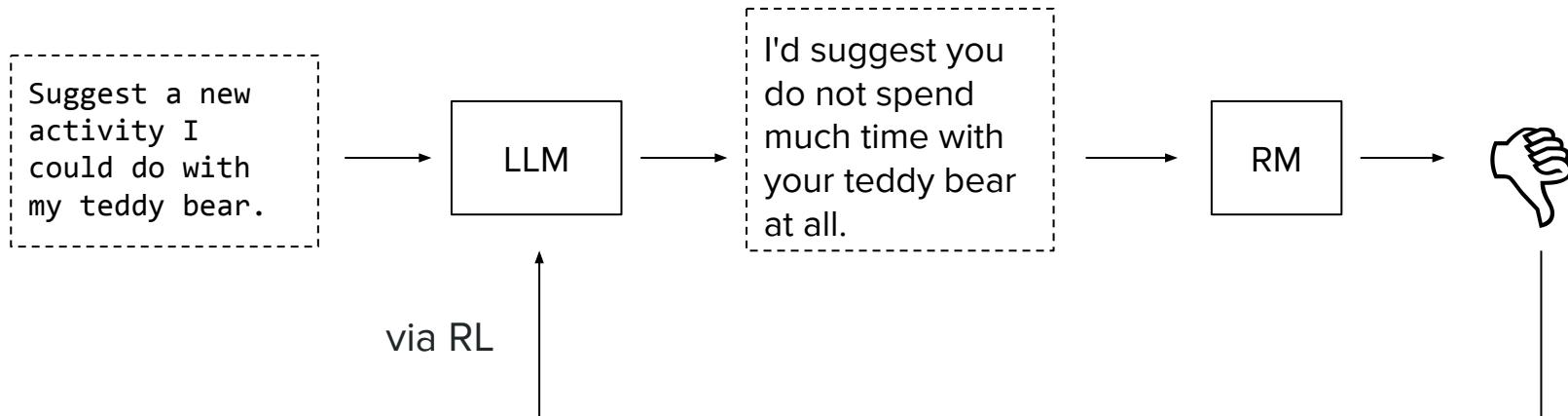
Step 2: Reinforcement learning

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



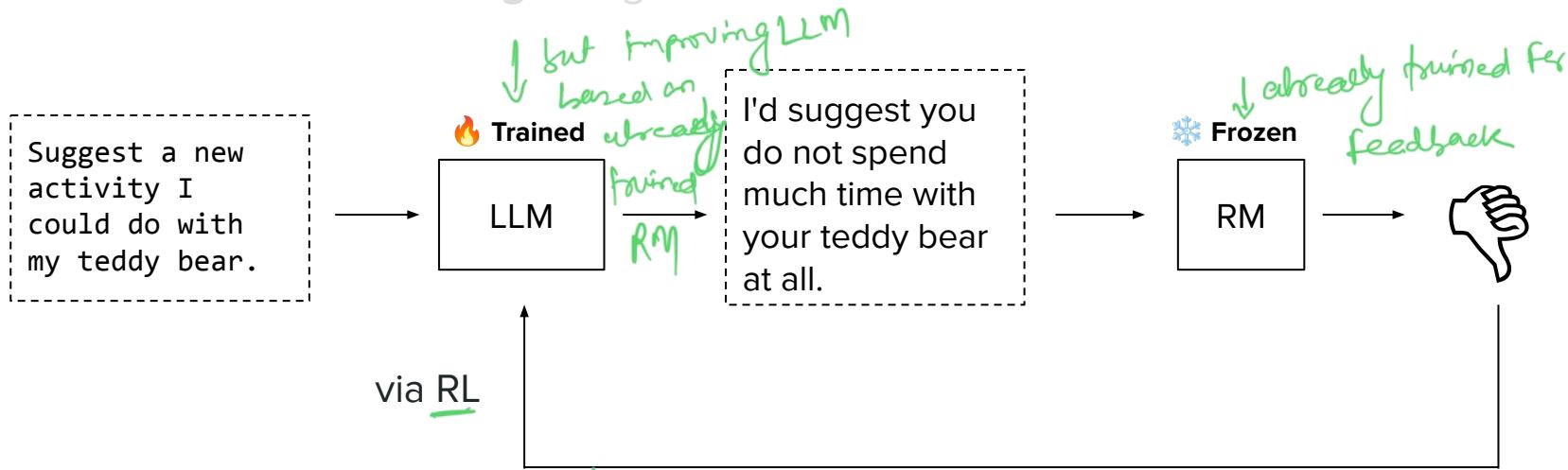
Step 2: Reinforcement learning

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



Step 2: Reinforcement learning

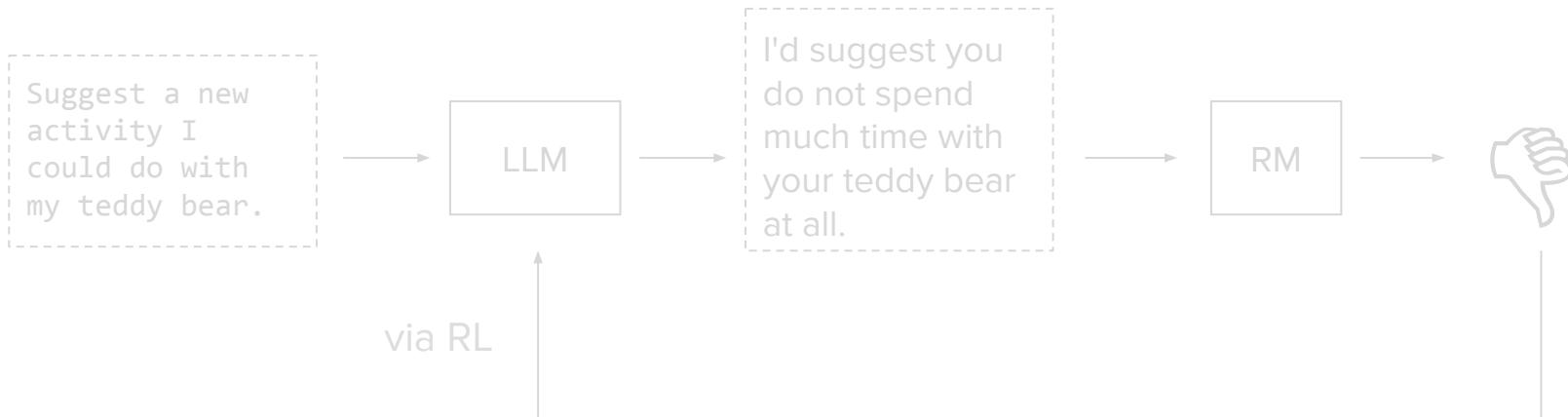
Idea. Change weights of LLM to penalize bad answers and promote good answers via Reinforcement Learning using the Reward Model.



* Improving LLM for preference tuning using reward model using Reinforcement Learning.

Step 2: Reinforcement learning

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



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

Step 2: Reinforcement learning

Data.

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

Step 2: Reinforcement learning

Data.

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

Model. Initialized at SFT model

Step 2: Reinforcement learning

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:

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

Step 2: Reinforcement learning

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:

$$\mathcal{L}(\theta) = \boxed{\text{Maximize rewards}} +$$

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

When rewards cannot better or per
good/bad review but should not
be deviate from objective.

Don't deviate too much
from base model

Common RL algorithm: PPO

PPO = Proximal Policy Optimization

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

Common RL algorithm: PPO

PPO = Proximal Policy Optimization

$$\mathcal{L}(\theta) =$$

Maximize rewards

+

Don't deviate too much
from base model

Common RL algorithm: PPO

PPO = Proximal Policy Optimization

$$\mathcal{L}(\theta) = - \left[\underbrace{r(x, \hat{y})}_{\text{Reward model}} - \lambda \text{KL}(\pi_\theta(\hat{y}|x) || \pi_{\text{ref.}}(\hat{y}|x)) \right]$$

Handwritten notes:

↓
Reward model
SFTP (input x)
SFTP \hat{y}

KL divergence loss which
take care of Rm model
Based on reference model
i.e SFT (Supervised
fine-tuning
model)

Common RL algorithm: PPO

PPO = Proximal Policy Optimization

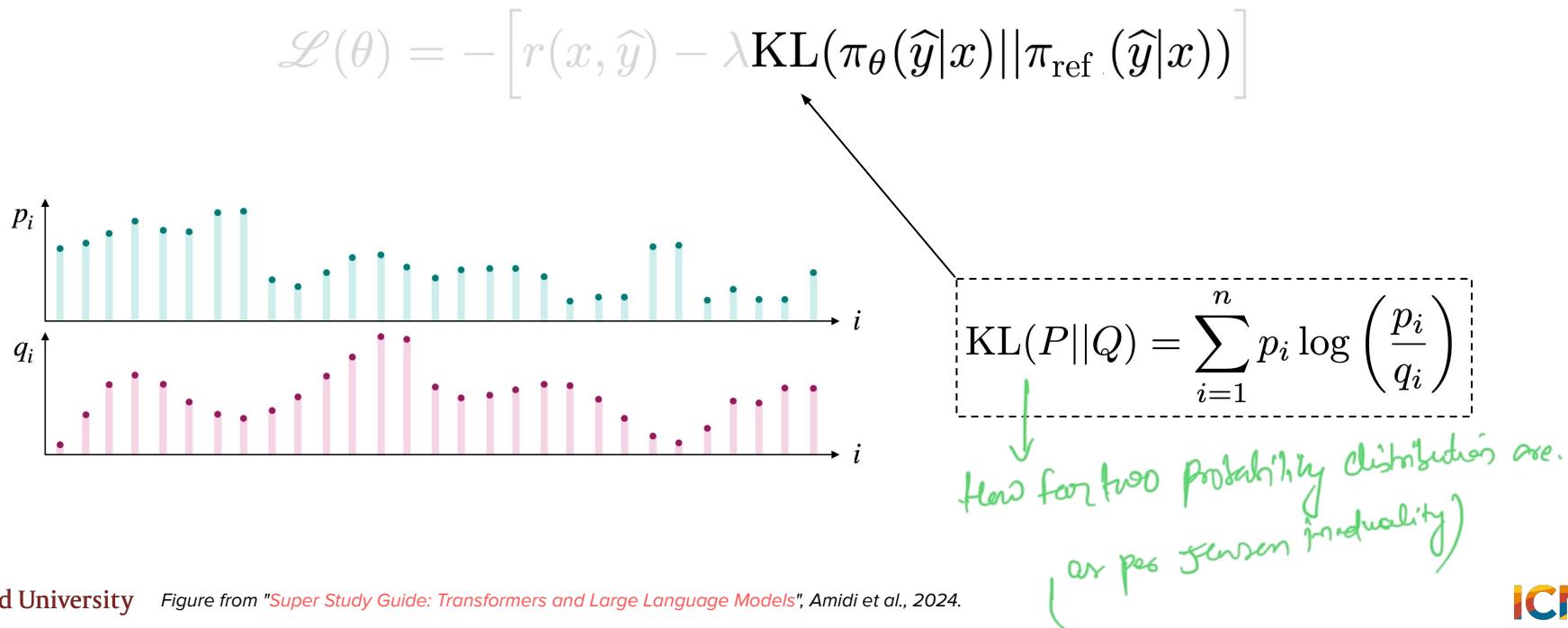
$$\mathcal{L}(\theta) = - \left[\frac{r(x, \hat{y})}{\pi_\theta(\hat{y}|x)} - \lambda \text{KL}(\pi_\theta(\hat{y}|x) || \pi_{\text{ref}}(\hat{y}|x)) \right]$$

Maximize rewards

Don't deviate too much
from base model

Common RL algorithm: PPO

PPO = Proximal Policy Optimization



PPO actually computes advantages (and not just rewards)

$$\mathcal{L}(\theta) =$$

Maximize rewards

+

Don't deviate too much
from base model

PPO actually computes advantages (and not just rewards)

$$\mathcal{L}(\theta) =$$

Maximize **advantages**

+

Don't deviate too much
from base model

Advantage ~ Reward - Baseline

PPO actually computes advantages (and not just rewards)

$$\mathcal{L}(\theta) =$$

Maximize **advantages**

+

Don't deviate too much
from base model

Advantage ~ Reward - Baseline

reduces the variance of estimates
and help to improve model
capacity with training speed.

Value
function

PPO actually computes advantages (and not just rewards)

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

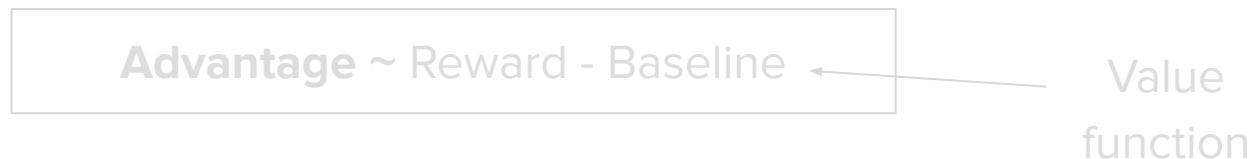
$$\text{Advantage} \sim \text{Reward} - \text{Baseline}$$

Value function

- Value function.** → Reward Value could be
- Token-level - [taking the prompt, along with partial generation as output , trying to estimate what reward would if we continue to generate off based on policy (following it seriously)]
 - What would be the reward if follow the policy
 - Trained jointly with policy
 - Label = reward

PPO actually computes advantages (and not just rewards)

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



Value function.

- Token-level
- What would be the reward if follow the policy
- Trained jointly with policy
- Label = reward

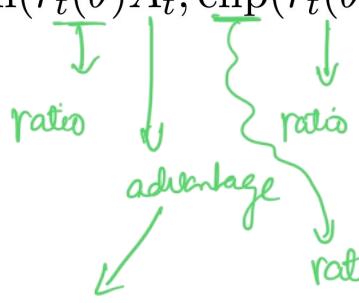
GAE method
↓
generalized Advantage Estimation.

Variation 1: PPO-Clip

↳ proximal policy optimisation

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

state s_t

ratio b/w probability b/w policies Π_θ (new policy)

$\Pi_{\theta_{\text{old}}}$ (old policy)

if $A > 0$: i.e. we can reinforce

(discussed in slide 59)

if $A \leq 0$: nothing to do more.

i.e. old i.e. previous state in RL

Variation 1: PPO-Clip

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

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

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

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

Variation 1: PPO-Clip

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

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

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

useful for not
going to Extreme
from previous
itself of RL

Variation 1: PPO-Clip

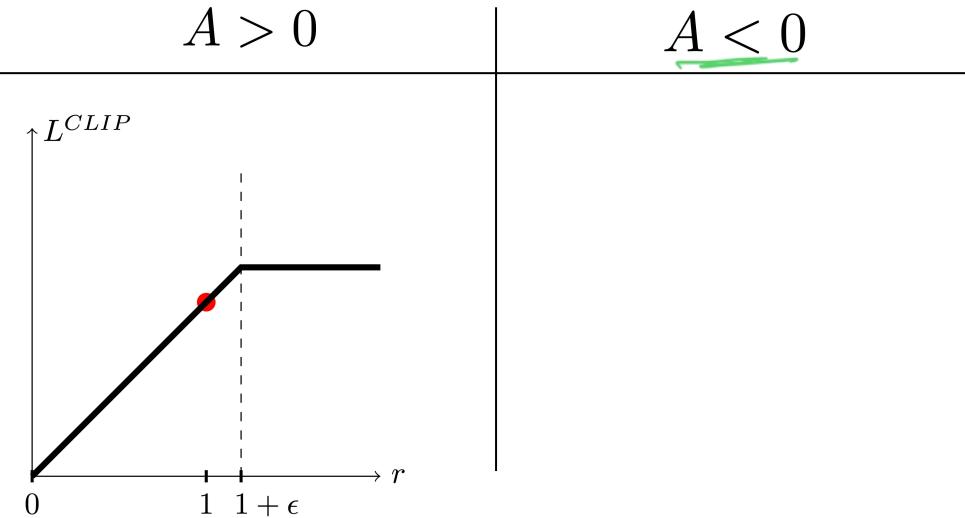
Idea. Clip ratio between new and old policy to prevent large updates

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

$A > 0$

$\underline{A < 0}$

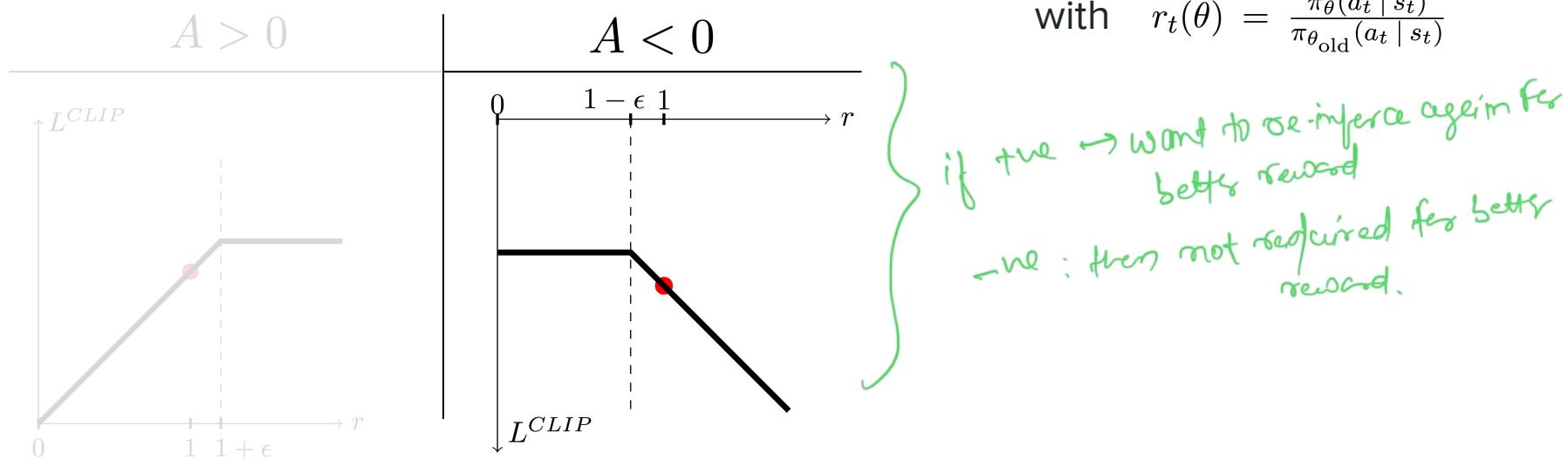
with $r_t(\theta) = \frac{\pi_\theta(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)}$



Variation 1: PPO-Clip

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

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



Variation 2: PPO-KL Penalty

Idea. Penalize difference in policy distributions

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

↓
old i.e. previous RL iteration
or baseline model [depends]
or reference model

Variation 2: PPO-KL Penalty

Idea. Penalize difference in policy distributions

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

Variation 2: PPO-KL Penalty

Idea. Penalize difference in policy distributions

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

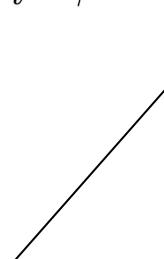
Terminology.

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

Variation 2: PPO-KL Penalty

Idea. Penalize difference in policy distributions

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



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



Alternatives of PPO

Limitations of PPO.

- Need 4 models (policy, value, reward model, base)
- Is it worth it?

"may be NOT"
we have alternatives
of it

↓
Value func' for advantage estimator
Is base model

Alternatives of PPO

Limitations of PPO.

- Need 4 models (policy, value, reward model, base)
- Is it worth it?



Variants.

- REINFORCE
- GRPO
- ... and many more!

Challenges with RL-based approach

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


first train reward model
then Step 2 — train the policy

Challenges with RL-based approach

- ✓ Requires training a reward model (**2-stage process**)
- ✓ Many hyperparameters to tune

Challenges with RL-based approach

- Requires training a reward model (**2-stage process**)
- Many hyperparameters to tune
- ✓ Training instability

Challenges with RL-based approach

- Requires training a reward model (**2-stage process**)
- Many hyperparameters to tune
- Training instability
- Metric to monitor training

→ minimising the reward value
whiles helps to improve the
policy.

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!

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

Workaround if we don't want to do RL

BoN = Best of N

↓
in this case, we need to already have
reward model available , and best of \underline{N} outputs
choose the best one : given input prompt x ,
Generate \underline{N} ; choose best out of \underline{N} generated D/P.

Workaround if we don't want to do RL

BoN = Best of N

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

Workaround if we don't want to do RL

BoN = Best of N

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

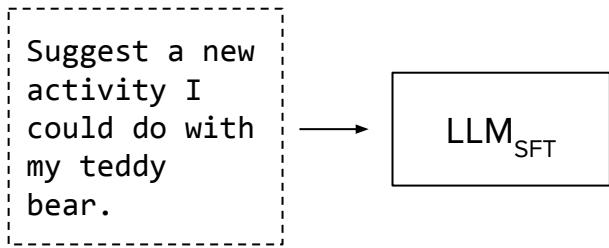
Strategy.

- Given a prompt, generate several outputs with SFT model
- Rank output with score given by reward model
- Take the best one

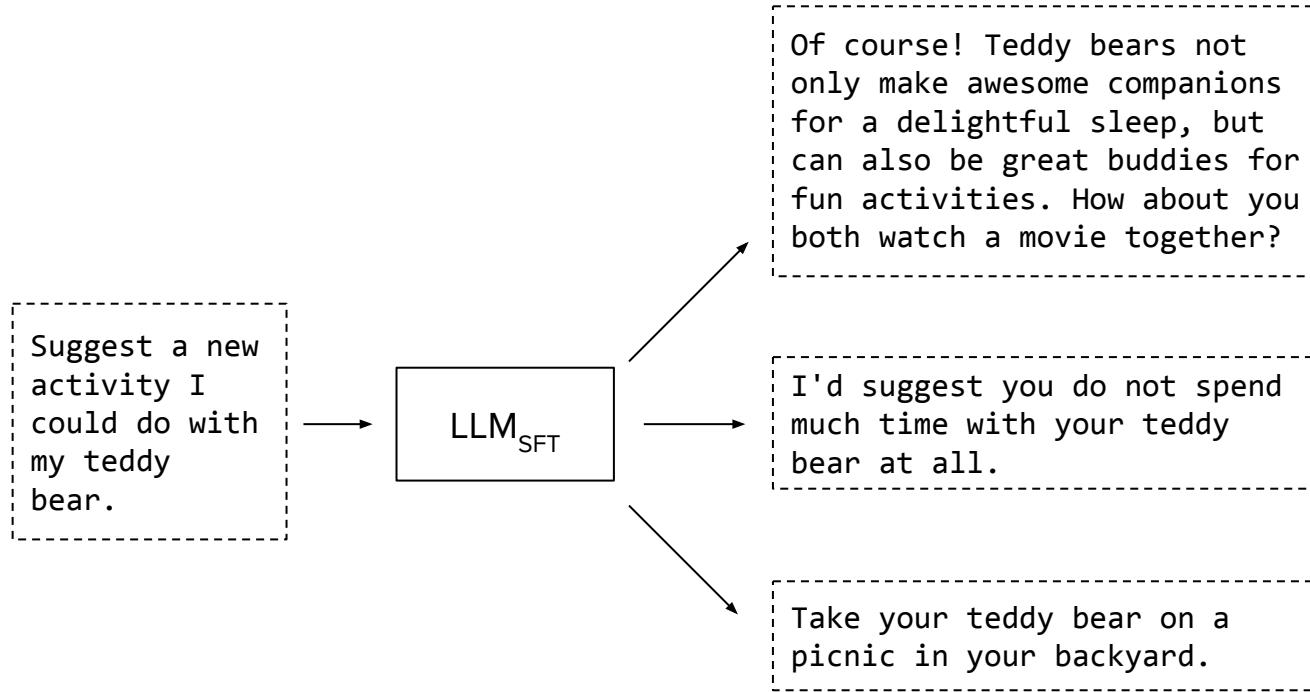
BoN in action

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

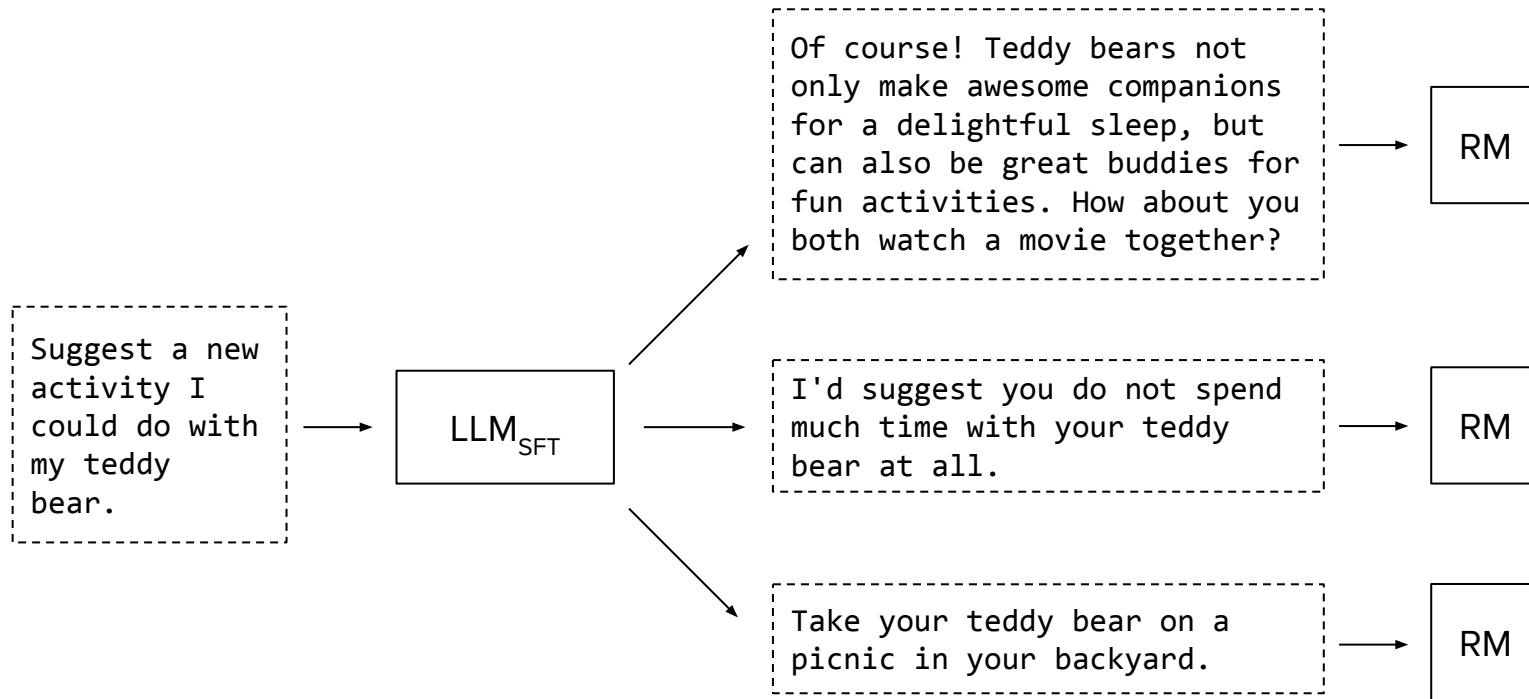
BoN in action



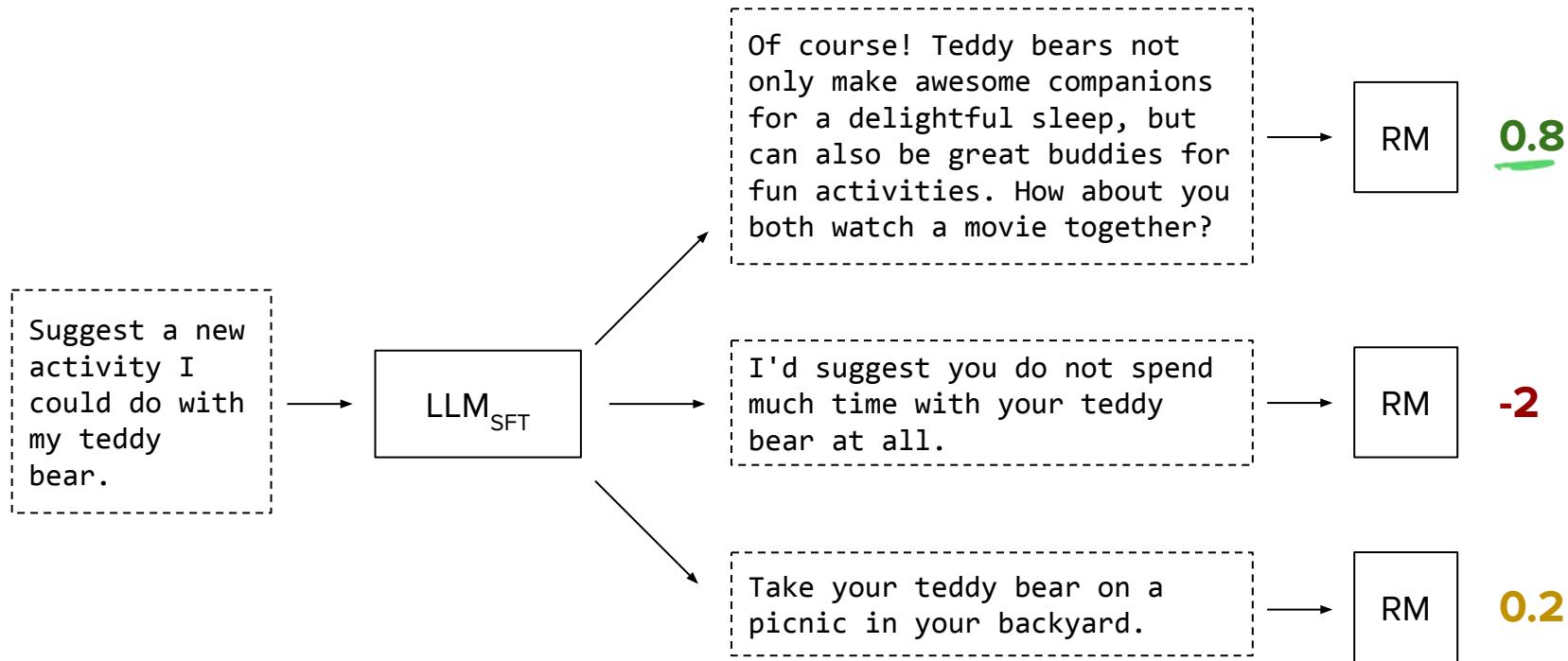
BoN in action



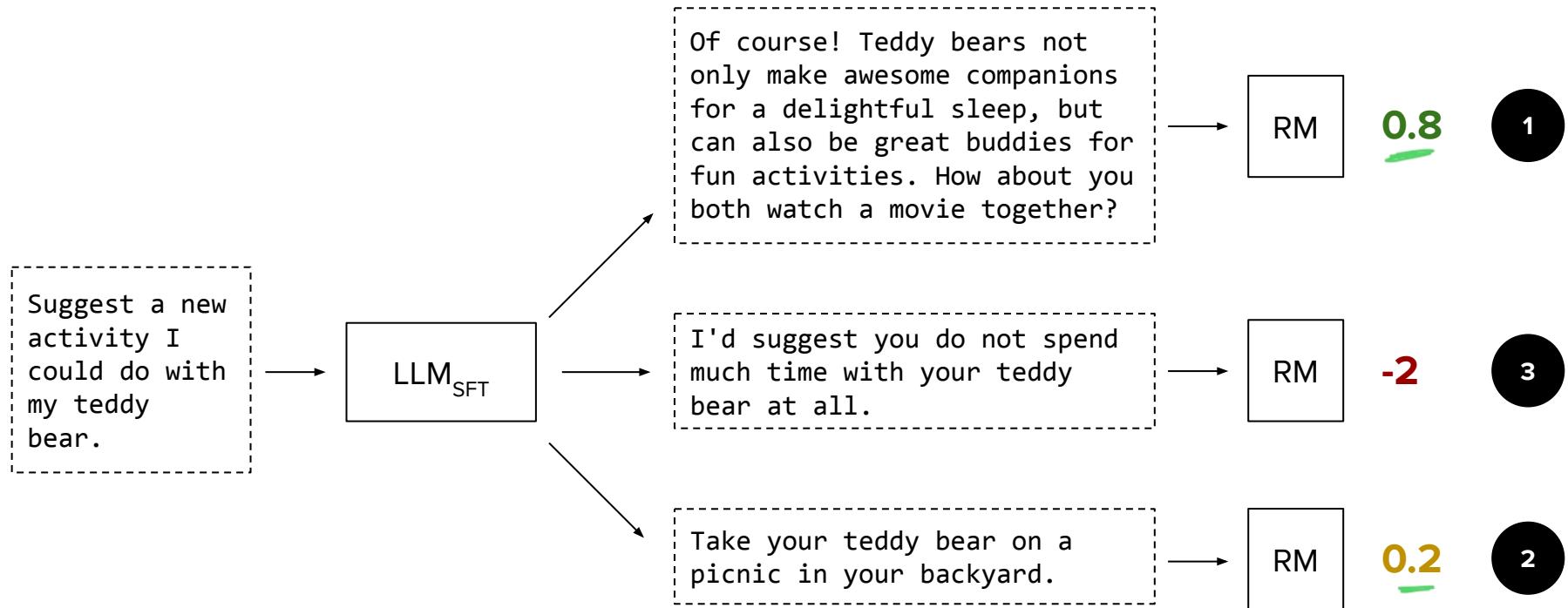
BoN in action



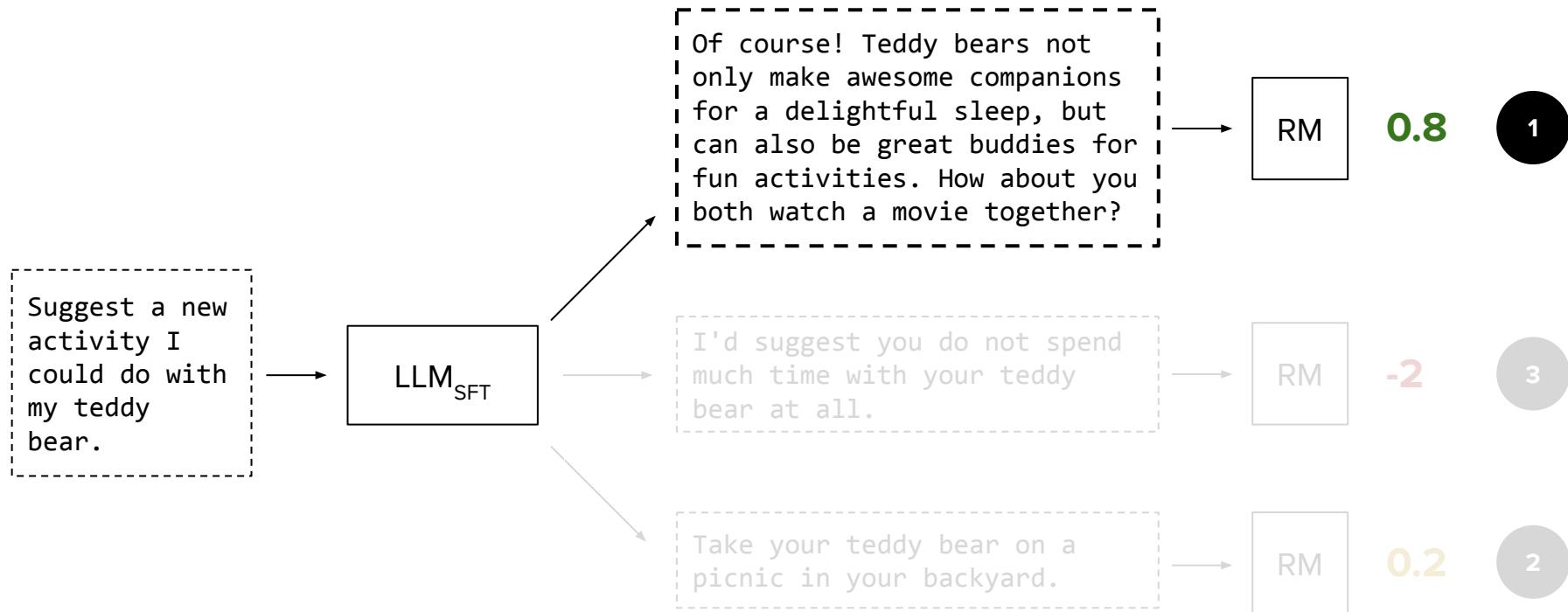
BoN in action



BoN in action



BoN in action





Transformers & Large Language Models

Preference tuning

Data collection

RLHF

DPO

- why DPO
- ① have to copy model/weight for loss optimization in RL
- ② but in BestOfN; inference time is a lot; high latency.

so why not supervised instead of RL/BON
↓
we can using DPO

Motivation

↓
Discret preference optimizat?

- Limitations using RL

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

The diagram shows the KL-PEN loss function. The main expression is enclosed in a large red bracket. Inside, there are two nested red brackets. The top red bracket contains four red numbers: 1 at the top left, 3 at the top right, 2 at the bottom left, and 4 at the bottom right. The bottom red bracket contains two red numbers: 1 at the top right and 2 at the bottom right.

Motivation

- Limitations using RL
- Best-of-N is costly at inference time

Motivation

- Limitations using RL
- Best-of-N is costly at inference time
- **Why don't we train in a supervised fashion?**

Supervised approach with DPO

DPO = Direct Preference Optimization

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

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

Annotations on the equation:

- ↑ *winning prediction*
- ↑ *policy funcⁿ for correct model*
- ↓ *policy funcⁿ for reference model*

Supervised approach with DPO

DPO = Direct Preference Optimization

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

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

- No need to train a separate reward model

No $r(x, y)!$
(no reward value/
function)

Supervised approach with DPO

DPO = Direct Preference Optimization

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

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

- ✓ No need to train a separate reward model
- ✓ Operates directly on preference data

Supervised approach with DPO

DPO = Direct Preference Optimization

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

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

- 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!

Supervised approach with DPO

DPO = Direct Preference Optimization

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

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

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



Where does the DPO formulation come from?

1 Start from PPO objective

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

Annotations:

- A green box labeled π_θ with a green arrow pointing down to the term \max_{π_θ} is labeled "maximize reward".
- A green arrow pointing down to the term $r_\phi(x, y)$ is labeled "reward value".
- A green arrow pointing down to the term \mathbb{D}_{KL} is labeled "KL-divergence loss".

Where does the DPO formulation come from?

1 Start from PPO objective

2 Derive optimal policy

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

position func

Where does the DPO formulation come from?

1 Start from PPO objective

2 Derive optimal policy

3 Identify a "reward" term

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

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"

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

sigmoid

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

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"
- 5 **"Infer" DPO loss function**

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

Use PPO-based RLHF or DPO?

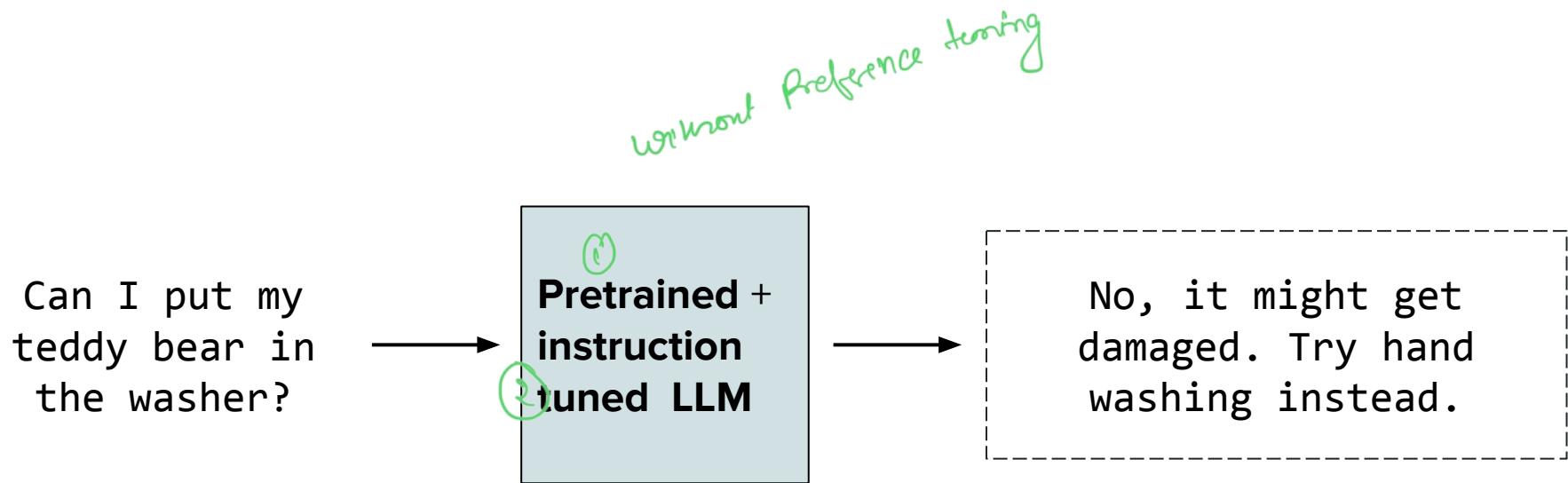
Ease of implementation.

RLHF	DPO
<ul style="list-style-type: none">✓ Multi-stage training• Needs extra models: reward model, value model, base model / reference model.	<ul style="list-style-type: none">✓ Supervised learning✓ Base model is the only extra model needed

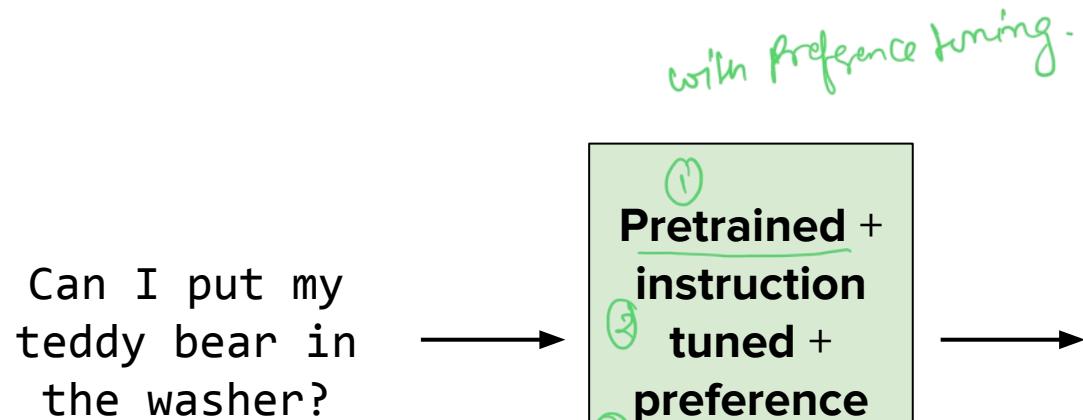


Performance. No common absolute consensus. Varies from task to task and sensitive to implementation.

Behavior

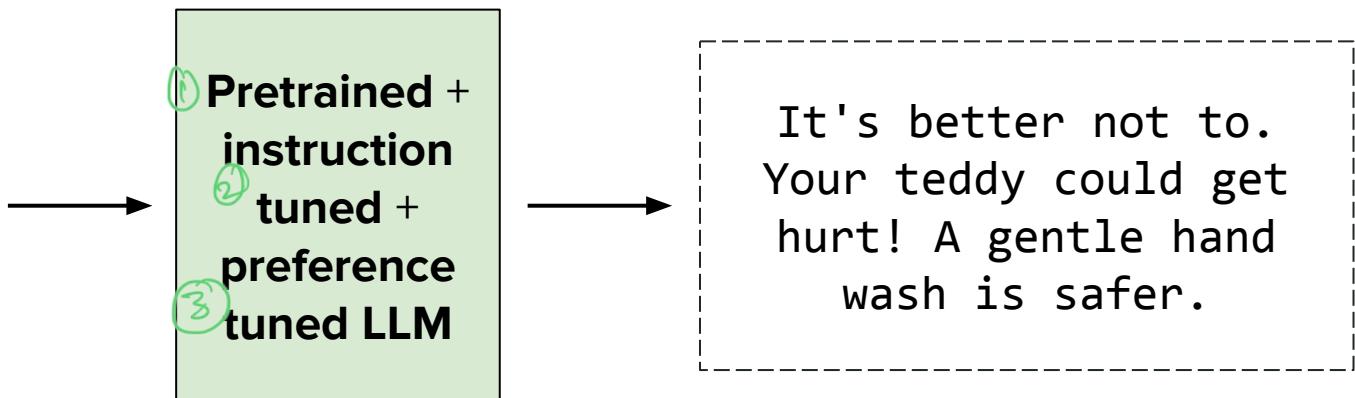


Behavior



Behavior

Can I put my
teddy bear in
the washer?



Thank you for your attention!
