# PPO

Let's get confused together :)

# RLHF

**❶ Collect human feedback**

A Reddit post is sampled from the Reddit TL;DR dataset.

Various policies are used to sample a set of summaries.

Two summaries are selected for evaluation.

A human judges which is a better summary of the post.

*"j is better than k"*

**❷ Train reward model**

One post with two summaries judged by a human are fed to the reward model.

The reward model calculates a reward $r$ for each summary.

$r_j$

$r_k$

The loss is calculated based on the rewards and human label, and is used to update the reward model.

$$loss = log(\sigma(r_j - r_k))$$

*"j is better than k"*

**❸ Train policy with PPO**

A new post is sampled from the dataset.

The policy $\pi$ generates a summary for the post.

The reward model calculates a reward for the summary.

The reward is used to update the policy via PPO.

$r$

Figure 2: Diagram of our human feedback, reward model training, and policy training procedure.

2020

**Step 1**
**Collect demonstration data, and train a supervised policy.**

A prompt is sampled from our prompt dataset.

A labeler demonstrates the desired output behavior.

This data is used to fine-tune GPT-3 with supervised learning.

**Step 2**
**Collect comparison data, and train a reward model.**

A prompt and several model outputs are sampled.

A labeler ranks the outputs from best to worst.

This data is used to train our reward model.

**Step 3**
**Optimize a policy against the reward model using reinforcement learning.**

A new prompt is sampled from the dataset.

The policy generates an output.

The reward model calculates a reward for the output.

The reward is used to update the policy using PPO.

2022

Reward model training

context → Policy → continuation (×4) → Reward model → reward (×4) → loss

Human labeler → label

Policy training

context → Policy → continuation → Reward model → reward → loss
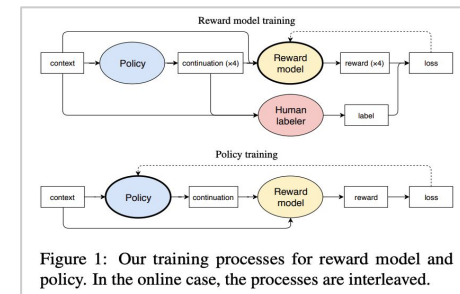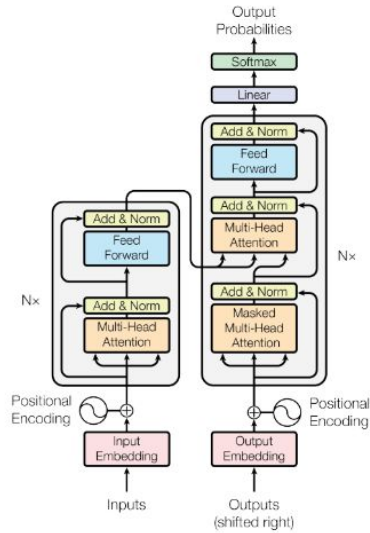
Figure 1: Our training processes for reward model and policy. In the online case, the processes are interleaved.

2019

# Problems

## Able to store massive amounts of data



## Efficient knowledge extraction

# AGI with RL?

# Overview

# Reward

Bradley–Terry

$$\mathcal{L}(\theta) = -\mathbb{E}_{(x, y_w, y_l)}\Big[\log \sigma\big(r_\theta(x, y_w) - r_\theta(x, y_l)\big)\Big] = -\mathbb{E}_{(x, y_w, y_l)} \log\Big[\frac{e^{r_\theta(x, y_w)}}{e^{r_\theta(x, y_w)} + e^{r_\theta(x, y_l)}}\Big]$$

.     .     .     .     4.2                          .     .     .     .     3.1

| Reward |          | Reward |

Hello    my    name    is    Bes                    Hello    I    call    myself    Bes

$$s_t = (x, y_{1:t-1}) \tag{1}$$

$$a_t = y_t \tag{2}$$

$$r_t = -\beta \, D_{\mathrm{KL}}\big(\pi_\theta(\cdot \mid s_t) \,\|\, \pi_{\mathrm{ref}}(\cdot \mid s_t)\big) + 1\!\!1_{\{t=T\}} \, r_{\mathrm{score}}(x, y). \tag{3}$$

$$D_{\mathrm{KL}}^{(t)} \approx \log \pi_\theta(a_t \mid s_t) - \log \pi_{\mathrm{ref}}(a_t \mid s_t). \tag{4}$$

$$\delta_t = r_t + \gamma \, V_\phi(s_{t+1}) - V_\phi(s_t). \tag{5}$$

$$\hat{A}_t = \sum_{l=0}^{T-t-1} (\gamma\lambda)^l \, \delta_{t+l}. \tag{6}$$

$$\hat{R}_t = \hat{A}_t + V_\phi(s_t). \tag{7}$$

$$r_t(\theta) = \frac{\pi_\theta(a_t \mid s_t)}{\pi_{\mathrm{old}}(a_t \mid s_t)}. \tag{8}$$

$$r_t(\theta) = \frac{\pi_\theta(a_t \mid s_t)}{\pi_{\mathrm{old}}(a_t \mid s_t)}. \tag{8}$$

$$J_{\mathrm{clip}}(\theta) = \mathbb{E}_t\Big[\min\Big(r_t(\theta)\,\hat{A}_t, \; \mathrm{clip}(r_t(\theta), 1-\epsilon, 1+\epsilon)\,\hat{A}_t\Big)\Big]. \tag{9}$$

$$L_{\mathrm{policy}}(\theta) = -J_{\mathrm{clip}}(\theta). \tag{10}$$

$$L_{\mathrm{value}}(\phi) = \mathbb{E}_t\Big[\big(V_\phi(s_t) - \hat{R}_t\big)^2\Big]. \tag{11}$$

$$L_{\mathrm{ent}}(\theta) = -\mathbb{E}_t\big[\mathcal{H}(\pi_\theta(\cdot \mid s_t))\big]. \tag{12}$$

$$L_{\mathrm{ptx}}(\theta) = -\mathbb{E}_{(x,y)\sim\mathcal{D}_{\mathrm{ptx}}} \sum_{t=1}^{T} \log \pi_\theta(y_t \mid x, y_{1:t-1}). \tag{13}$$

$$L_{\mathrm{total}}(\theta, \phi) = L_{\mathrm{policy}}(\theta) + c_v \, L_{\mathrm{value}}(\phi) + c_{\mathrm{ent}} \, L_{\mathrm{ent}}(\theta) + c_{\mathrm{ptx}} \, L_{\mathrm{ptx}}(\theta). \tag{14}$$

# Policy & Friends



yes, everything is calculated at the token level, it is a form of reward shaping

only the last token gets the extra reward score, previous tokens get only the approximate KL divergence

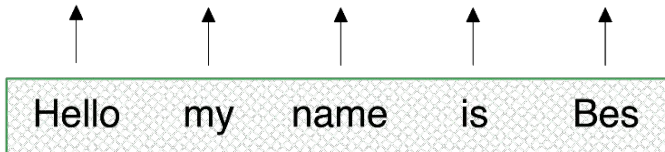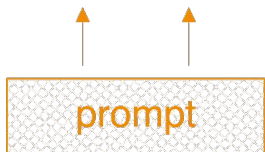this is calculated for every single token, correct?

what is this?

Autoregressive Run

Policy

prompt

Hello my name is Bes

$$r_t = -\beta \, D_{\mathrm{KL}}\big(\pi_\theta(\cdot \mid s_t) \,\|\, \pi_{\mathrm{ref}}(\cdot \mid s_t)\big) + \mathbb{1}_{\{t=T\}} \, r_{\mathrm{score}}(x, y). \qquad (3)$$
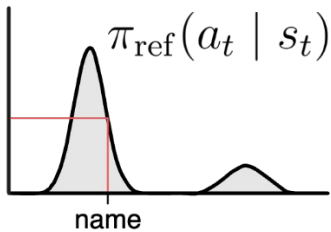
$$D_{\mathrm{KL}}^{(t)} \approx \log \pi_\theta(a_t \mid s_t) - \log \pi_{\mathrm{ref}}(a_t \mid s_t). \qquad (4)$$
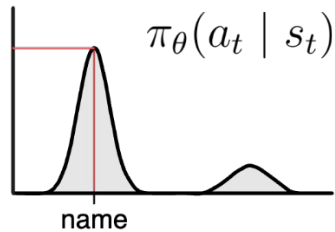
why an approximate?

approximated via sampling, because the true KL divergence would require the all distribution or all the possible actions

# KL

$$D_{\mathrm{KL}}^{(t)} \approx \log \pi_\theta(a_t \mid s_t) - \log \pi_{\mathrm{ref}}(a_t \mid s_t)$$

$$r_t(\theta) = \frac{\pi_\theta(a_t \mid s_t)}{\pi_{\mathrm{old}}(a_t \mid s_t)}$$

$\pi_{\mathrm{ref}}(a_t \mid s_t)$

name

$\pi_\theta(a_t \mid s_t)$

name

$\pi_\theta(a_t \mid s_t)$

name

$\pi_{\mathrm{old}}(a_t \mid s_t)$

name

Ref

Policy

Policy

Old

# clip

$$r_t(\theta) = \frac{\pi_\theta(a_t \mid s_t)}{\pi_{\text{old}}(a_t \mid s_t)}. \tag{8}$$

$$J_{\text{clip}}(\theta) = \mathbb{E}_t\left[\min\left(r_t(\theta)\,\hat{A}_t,\ \text{clip}(r_t(\theta),\,1-\epsilon,\,1+\epsilon)\,\hat{A}_t\right)\right]. \tag{9}$$

$$L_{\text{policy}}(\theta) = -J_{\text{clip}}(\theta). \tag{10}$$

$$L_{\text{value}}(\phi) = \mathbb{E}_t\left[\left(V_\phi(s_t) - \hat{R}_t\right)^2\right]. \tag{11}$$

$$L_{\text{ent}}(\theta) = -\mathbb{E}_t\left[\mathcal{H}(\pi_\theta(\cdot \mid s_t))\right]. \tag{12}$$

$$L_{\text{ptx}}(\theta) = -\mathbb{E}_{(x,y)\sim\mathcal{D}_{\text{ptx}}}\sum_{t=1}^{T}\log\pi_\theta(y_t \mid x, y_{1:t-1}). \tag{13}$$

$$L_{\text{total}}(\theta,\phi) = L_{\text{policy}}(\theta) + c_v\,L_{\text{value}}(\phi) + c_{\text{ent}}\,L_{\text{ent}}(\theta) + c_{\text{ptx}}\,L_{\text{ptx}}(\theta). \tag{14}$$

$$\text{clip}(x, a, b) = \begin{cases} a, & x < a, \\ x, & a \leq x \leq b, \\ b, & x > b. \end{cases}$$