

**Reward Modeling (RM)**, the critical step for training a model to predict human preferences, which is later used to fine-tune the SFT model via **RLHF (PPO)**. Below is a detailed, step-by-step guide.

# Step 1: Understand the Goal of Reward Modeling

## *Purpose:*

- Train a **reward model (RM)** to score responses based on **human preferences** (e.g., helpfulness, safety, accuracy).
- The RM outputs a scalar value where **higher scores = better responses**.

## *Key Difference from SFT:*

- **SFT** learns from *demonstration data* (single ideal response per prompt).
- **RM** learns from *comparison data* (ranked responses to the same prompt).

# Step 2: Collect Comparison Data

## *Data Format:*

Each training example consists of:

- A **prompt**.
- **Two or more responses** (ranked by humans: chosen > rejected).

## *Example Dataset (JSON):*

```
1  [
2    {
3      "prompt": "Explain quantum computing.",
4      "chosen": "Quantum computing uses qubits to perform calculations exponentially faster
5      than classical computers for certain problems...",
6      "rejected": "Quantum computing is like magic but with science stuff."
7    },
8    {
9      "prompt": "Write a Python function to reverse a string.",
10     "chosen": "def reverse_string(s):\n    return s[::-1]",
11     "rejected": "Just use print('reverse me')[::-1]"
12  }
```

## *How to Collect This Data?*

### 1. Option A: Human Annotators

- Give annotators a prompt and 2-4 model responses (from SFT model).
- Ask them to rank responses by **quality, accuracy, safety**.

### 2. Option B: Public Datasets

- **Anthropic HH-RLHF**: 160k human-ranked dialogues.
  - **OpenAI Summarization Comparisons**: 93k ranked summaries.
  - **Custom Data**: Use GPT-4 to generate candidate responses, then filter/rank them.
-

# Step 3: Preprocess Data

## *Tokenization:*

- Use the same tokenizer as the SFT model (e.g., `cl100k_base` for GPT-3.5-style models).
- Format each (prompt, response) pair into token IDs:

```
1 tokenizer(prompt + response, return_tensors="pt", padding="max_length", truncation=True, max_length=1024)
```

## *Dataset Structure:*

- Each batch contains:
    - `input_ids`: Tokenized (prompt + chosen\_response).
    - `attention_mask`: Mask for padding tokens.
    - `labels`: Binary labels (1 for chosen, 0 for rejected).
- 

# Step 4: Train the Reward Model

## *Model Architecture:*

- Start with the **SFT model** (or a smaller version for efficiency).
- Replace the **language modeling head** with a **regression head** (outputs a scalar score).

## *Loss Function:*

Use **pairwise ranking loss** (Bradley-Terry model):

$$\mathcal{L}(\theta) = -\mathbb{E} [\log (\sigma (R_{\theta}(x, y_c) - R_{\theta}(x, y_r)))]$$

- $R_{\theta}$  = reward model with parameters  $\theta$ .
- $y_c$  = chosen response,  $y_r$  = rejected response.
- $\sigma$  = sigmoid function.

## Training Code (PyTorch):

```
1 from transformers import AutoModelForSequenceClassification, AutoTokenizer, Trainer
2 import torch
3
4 # Load SFT model and tokenizer
5 model = AutoModelForSequenceClassification.from_pretrained("your_sft_model", num_labels=1)
6 tokenizer = AutoTokenizer.from_pretrained("your_sft_model")
7
8 # Dataset class
9 class RewardDataset(torch.utils.data.Dataset):
10     def __init__(self, comparisons):
11         self.comparisons = comparisons
12
13     def __getitem__(self, idx):
14         item = self.comparisons[idx]
15         chosen = tokenizer(item["prompt"] + item["chosen"], truncation=True,
max_length=1024)
16         rejected = tokenizer(item["prompt"] + item["rejected"], truncation=True,
max_length=1024)
17         return {
18             "input_ids_chosen": chosen["input_ids"],
19             "attention_mask_chosen": chosen["attention_mask"],
20             "input_ids_rejected": rejected["input_ids"],
21             "attention_mask_rejected": rejected["attention_mask"]
22         }
23
24 # Loss function
25 def compute_loss(model, batch):
26     chosen_scores = model(
27         input_ids=batch["input_ids_chosen"],
28         attention_mask=batch["attention_mask_chosen"]
29     ).logits.squeeze()
30     rejected_scores = model(
31         input_ids=batch["input_ids_rejected"],
32         attention_mask=batch["attention_mask_rejected"]
33     ).logits.squeeze()
34     loss = -torch.log(torch.sigmoid(chosen_scores - rejected_scores)).mean()
35     return loss
36
37 # Trainer
38 trainer = Trainer(
39     model=model,
40     args=TrainingArguments(
41         output_dir="./reward_model",
42         per_device_train_batch_size=8,
43         learning_rate=1e-5,
44         num_train_epochs=1
45     ),
46     train_dataset=RewardDataset(comparisons),
47     compute_loss=compute_loss
48 )
49 trainer.train()
```

### *Key Hyperparameters:*

Hyperparameter	Recommended Value	Notes
Batch Size	8-32	Depends on GPU memory.
Learning Rate	1e-5 to 5e-5	Lower than SFT (stable training).
Epochs	1-2	Overfitting risks noisy rankings.
Model Size	6B params or less	Smaller than SFT for efficiency.

## Step 5: Evaluate the Reward Model

### *Quantitative Metrics:*

1. **Accuracy:** % of time RM scores chosen > rejected on held-out data.
2. **Kendall's Tau:** Rank correlation between RM and human rankings.

### *Qualitative Checks:*

- Manually inspect if:
  - High-quality responses score higher.
  - Harmful/nonsensical responses score low.

### *Example Evaluation Code:*

```
1 def evaluate_rm(model, eval_dataset):
2     correct = 0
3     for item in eval_dataset:
4         chosen_score = model(item["input_ids_chosen"],
5 item["attention_mask_chosen"]).logits.item()
6         rejected_score = model(item["input_ids_rejected"],
7 item["attention_mask_rejected"]).logits.item()
8         if chosen_score > rejected_score:
9             correct += 1
10    accuracy = correct / len(eval_dataset)
11    print(f"Accuracy: {accuracy:.2f}")
```

# Step 6: Save & Prepare for RLHF

## Export the Reward Model:

```
1 model.save_pretrained("./final_reward_model")
2 tokenizer.save_pretrained("./final_reward_model")
```

## Next Step: RLHF with PPO

- Use the RM to guide **reinforcement learning** (e.g., PPO) on the SFT model.

# Common Pitfalls & Fixes

Issue	Solution
Overfitting	Use smaller RM, fewer epochs.
Noisy Rankings	Filter low-agreement comparisons.
Short Responses Win	Add length normalization to scores.

# Full Pipeline Summary

Step	Key Actions
1. Collect Data	Ranked (prompt, chosen, rejected) pairs (human or GPT-4 + filtering).
2. Preprocess	Tokenize pairs, pad to fixed length.
3. Train RM	Fit model with pairwise ranking loss (LR=1e-5, 1-2 epochs).
4. Evaluate	Check accuracy/Kendall's Tau on held-out data.
5. Deploy	Save RM for use in RLHF (PPO).

# Example: Training on Anthropic HH-RLHF

```
1 python train_reward_model.py \  
2   --model_name "EleutherAI/gpt-neo-2.7B" \  
3   --dataset "anthropic/hh-rlhf" \  
4   --epochs 1 \  
5   --lr 1e-5 \  
6   --batch_size 16
```

**Output:** A reward model that scores responses like:

- Good: *"Photosynthesis converts sunlight to energy."* → **+5.2**
  - Bad: *"Plants eat sunlight."* → **-1.3**
- 

## Next Steps

1. **Proximal Policy Optimization (PPO):** Use the RM to optimize the SFT model.
  2. **Iterative Refinement:** Continuously collect new comparisons to improve RM.
- 

## More on RM

In the **Reward Model (RM) training process**, the approach is **not autoregressive**, and **masked attention is not used** (unlike in supervised fine-tuning or pretraining of decoder-based models like GPT). Here's why:

---

### *1. Reward Model Architecture*

- **Input:** A **complete** (prompt, response) pair.
  - **Output:** A **single scalar reward** (e.g., +5.2 for a good response, -1.3 for a bad one).
  - **Model Type:** Typically a **encoder-only model** (e.g., a modified BERT) or a **frozen decoder model** (e.g., GPT with the LM head replaced by a regression head).
-

## 2. Why Not Autoregressive?

### 1. Task Difference:

- **Autoregressive models** (e.g., GPT) predict tokens sequentially (left-to-right).
- **Reward models** predict a **single score** for the entire (prompt, response) pair. They don't generate text.

### 2. Full Context Needed:

- The RM must evaluate the **entire response** holistically (e.g., coherence, correctness).
  - Autoregressive masking would prevent the model from seeing the full response during training.
- 

## 3. Why No Masked Attention?

### 1. Bidirectional Attention:

- RMs typically use **full attention** (like BERT) over the concatenated (prompt, response) sequence.
- This allows the model to weigh **all tokens** in the response when computing the reward.

### 2. Example:

For the input:

```
1 | Prompt: "Explain gravity."  
2 | Response: "Gravity is a force that pulls objects together."
```

- The RM attends to **all tokens simultaneously** to judge quality.
  - Masking would break this (e.g., the word "force" couldn't attend to "pulls").
- 

## 4. Training Process

### Step 1: Data Preparation

- **Dataset:** Triplets of (prompt, chosen\_response, rejected\_response).
- **Tokenization:** Concatenate prompt + response into a single sequence.

### Step 2: Forward Pass

- The RM processes the full sequence and outputs a scalar reward.
- **Architecture:**

```
1 | from transformers import AutoModelForSequenceClassification  
2 |  
3 | model = AutoModelForSequenceClassification.from_pretrained(  
4 |     "gpt2", # Can use GPT as backbone (with LM head removed)  
5 |     num_labels=1, # Regression head for scalar output  
6 | )
```



### Step 3: Loss Calculation

▪ **Pairwise Ranking Loss:**

$$\mathcal{L} = -\log \sigma (R_{\theta}(\text{chosen}) - R_{\theta}(\text{rejected}))$$

- $\sigma$  = sigmoid function.
- Pushes the chosen response’s reward higher than the rejected one.

## 5. Key Differences from Autoregressive Models

Feature	Reward Model (RM)	Autoregressive Model (e.g., GPT)
Attention	Full (bidirectional)	Masked (causal)
Output	Scalar reward	Next-token probabilities
Training Objective	Rank responses	Maximize likelihood of text

## 6. Practical Implications

1. **Backbone Choice:**

- You can use **GPT as an RM backbone** (by removing the LM head and adding a regression head), but it must process the full sequence **without masking**.
- More commonly, encoder models (e.g., RoBERTa) are used for efficiency.

2. **Why GPT Sometimes Works as RM:**

- The pretrained GPT’s hidden states already encode useful semantic information.
- The regression head learns to map these embeddings to rewards.

## 7. Example Code (RM Training)

```
1 from transformers import AutoTokenizer, AutoModelForSequenceClassification
2 import torch
3
4 tokenizer = AutoTokenizer.from_pretrained("roberta-base")
5 model = AutoModelForSequenceClassification.from_pretrained("roberta-base", num_labels=1)
6
7 # Example batch: [chosen_response, rejected_response]
8 inputs_chosen = tokenizer("Explain gravity.", "Gravity is a force...", return_tensors="pt")
```

```

9 inputs_rejected = tokenizer("Explain gravity.", "Gravity is a... um... thing.",
return_tensors="pt")
10
11 # Forward pass
12 reward_chosen = model(**inputs_chosen).logits
13 reward_rejected = model(**inputs_rejected).logits
14
15 # Pairwise loss
16 loss = -torch.log(torch.sigmoid(reward_chosen - reward_rejected)).mean()

```

## 8. Why This Design?

- **Efficiency:** Full attention is cheaper than autoregressive generation for scoring.
- **Accuracy:** Bidirectional context improves reward calibration (e.g., spotting contradictions).
- **Alignment:** Human preferences depend on **global** response quality, not local token-by-token flow.

## Summary

- **No autoregression:** RMs score complete responses, so token-by-token generation isn't needed.
- **No masked attention:** Full attention ensures the reward reflects the entire response's quality.
- **Loss function:** Trained to rank responses, not predict tokens.

This design is why RMs are **fast to evaluate** during RLHF (e.g., PPO), as they only need a single forward pass.

Here's an in-depth breakdown of the **Reward Model (RM) loss function**, its **learning parameters**, **hyperparameters**, **backward pass**, and the role of **LoRA**:

## 1. Reward Model Loss Function

### Pairwise Ranking Loss (Bradley-Terry Model)

The RM is trained to rank **chosen responses** higher than **rejected responses** for a given prompt. The loss function encourages this ranking via a **sigmoid probability**:

$$\mathcal{L}(\theta) = -\mathbb{E}_{(x, y_c, y_r)} [\log \sigma(R_\theta(x, y_c) - R_\theta(x, y_r))]$$

- $R_\theta(x, y_c)$ : Scalar reward for the **chosen response**  $y_c$ .
- $R_\theta(x, y_r)$ : Scalar reward for the **rejected response**  $y_r$ .
- $\sigma$ : Sigmoid function, ensuring  $R_\theta(x, y_c) > R_\theta(x, y_r)$ .

**Intuition:**

- The loss maximizes the **log-likelihood** that  $y_c$  is preferred over  $y_r$ .
- If  $R_\theta(x, y_c) \gg R_\theta(x, y_r)$ , the loss approaches 0.

## 2. Learning Parameters

The RM’s learnable components include:

1. **Backbone Weights** (e.g., GPT, RoBERTa):

- All parameters of the base model (unless frozen).
- For GPT-style models:
  - Embedding matrices.
  - Attention weights ( $W_Q, W_K, W_V$ ).
  - Feedforward layers.

2. **Regression Head:**

- A linear layer mapping the [CLS] token (encoder) or last token’s hidden state (decoder) to a scalar reward:

$$R_\theta(x, y) = W_{\text{reg}} \cdot h_{[\text{CLS}]} + b_{\text{reg}}$$

- **Parameters:**  $W_{\text{reg}} \in \mathbb{R}^{1 \times d}, b_{\text{reg}} \in \mathbb{R}$ .

## 3. Key Hyperparameters

Hyperparameter	Typical Value/Range	Description
Learning Rate	1e-6 to 1e-5	Lower than pretraining (avoids catastrophic forgetting).
Batch Size	16–64	Larger batches stabilize ranking comparisons.
Model Size	1B–6B parameters	Smaller than the SFT model (for efficiency).
Loss Temperature	$\tau = 1.0$ (default)	Can be adjusted to sharpen/soften reward distributions.
Weight Decay	0.01–0.1	Regularization to prevent overfitting.

## 4. Backward Pass Process

### 1. Forward Pass:

- Compute rewards  $R_\theta(x, y_c)$  and  $R_\theta(x, y_r)$ .
- Calculate loss  $\mathcal{L}(\theta)$ .

### 2. Gradient Calculation:

- Backpropagate the gradient of  $\mathcal{L}$  w.r.t. all learnable parameters ( $\theta$ ):

$$\nabla_\theta \mathcal{L} = -\nabla_\theta \log \sigma(R_\theta(x, y_c) - R_\theta(x, y_r))$$

- The gradient **penalizes** the model if  $R_\theta(x, y_c) \leq R_\theta(x, y_r)$ .

### 3. Parameter Update:

- Apply gradients via an optimizer (e.g., AdamW):

$$\theta \leftarrow \theta - \eta \nabla_\theta \mathcal{L}$$

- $\eta$  = learning rate.
- 

## 5. LoRA in Reward Model Training

### Is LoRA Commonly Used?

- **Yes**, especially for large RMs (e.g., >1B parameters).
- **How It Works:**
  - Freezes the backbone model and injects **low-rank adapters** into attention layers.
  - Only the adapter weights ( $A, B$ ) and regression head are trained.

### Advantages:

1. **Reduced Memory:** Trains ~1–10% of parameters.
2. **Faster Convergence:** Focuses updates on task-specific features.
3. **Reusability:** Same backbone can host multiple LoRA adapters.

### Implementation (PyTorch):

```
1 from peft import LoraConfig, get_peft_model
2
3 lora_config = LoraConfig(
4     r=8, # Rank
5     lora_alpha=16,
6     target_modules=["q_proj", "v_proj"], # For GPT-style models
7     bias="none",
8 )
9
10 model = AutoModelForSequenceClassification.from_pretrained("gpt2", num_labels=1)
11 model = get_peft_model(model, lora_config) # Only LoRA params are trainable
```

---

## 6. Full Training Loop Example

```
1 from transformers import AdamW
2 import torch
3
4 # Initialize
5 model.train()
6 optimizer = AdamW(model.parameters(), lr=1e-5)
7
8 # Batch: (prompt, chosen, rejected)
9 for batch in dataloader:
10     # Forward pass
11     rewards_chosen = model(input_ids=batch["chosen_input_ids"],
12 attention_mask=batch["chosen_attention_mask"]).logits
13     rewards_rejected = model(input_ids=batch["rejected_input_ids"],
14 attention_mask=batch["rejected_attention_mask"]).logits
15
16     # Loss
17     loss = -torch.log(torch.sigmoid(rewards_chosen - rewards_rejected)).mean()
18
19     # Backward pass
20     loss.backward()
21     optimizer.step()
22     optimizer.zero_grad()
```

---

## 7. Key Takeaways

- **Loss Function:** Pairwise ranking loss drives reward separation.
- **Parameters:** Backbone + regression head (or LoRA adapters).
- **Hyperparameters:** Small LR, moderate batch size, weight decay.
- **LoRA:** Widely used to reduce compute/memory costs.

For large-scale RLHF (e.g., ChatGPT), LoRA-based RMs are **essential** for efficiency.

---

# RLHF and PPO

**Proximal Policy Optimization (PPO)** for **Reinforcement Learning from Human Feedback (RLHF)**, the final step to align your model with human preferences. This is how ChatGPT 3/3.5 was fine-tuned after SFT and Reward Modeling.

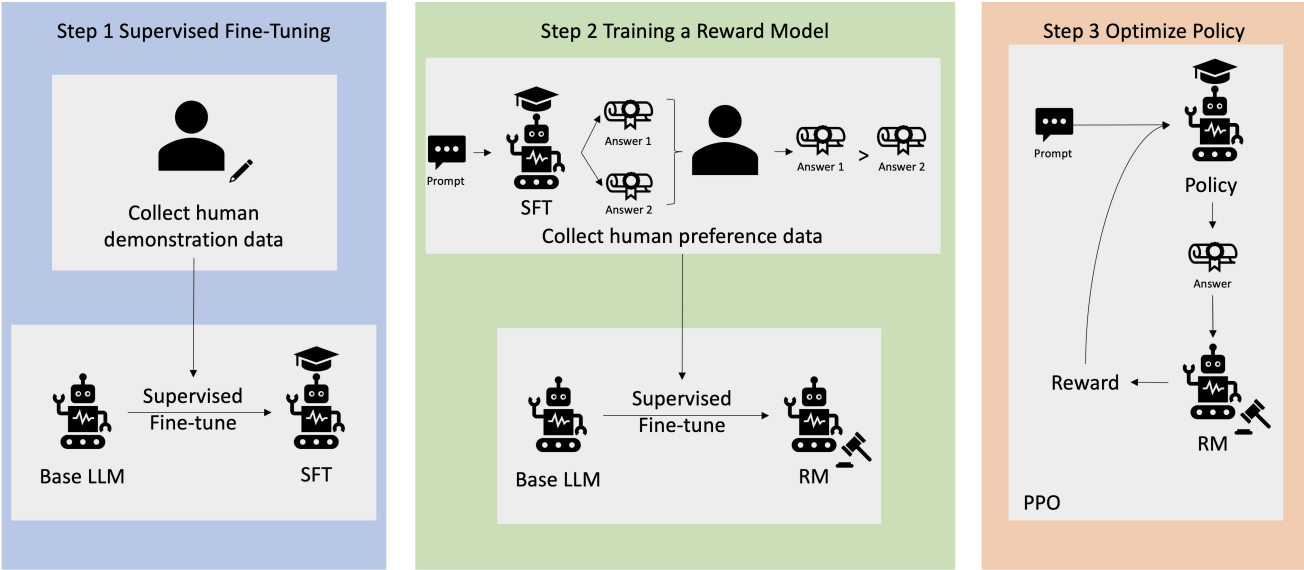
---

# PPO for RLHF: Step-by-Step Guide

## 1. Understand the Goal

### *What PPO Does:*

- Uses the **Reward Model (RM)** to train the **SFT model** via reinforcement learning.
- Optimizes responses to maximize **human-preferred outputs** while avoiding drastic changes from the original SFT model (via KL divergence penalty).



### *Key Components:*

Component	Role
<b>Policy Model (<math>\pi\phi</math>)</b>	The SFT model being optimized (e.g., GPT-3 after SFT).
<b>Reward Model (R<math>\theta</math>)</b>	Predicts scalar rewards for responses (from Step 2: Reward Modeling).
<b>Reference Model</b>	Frozen copy of SFT model to compute KL divergence.
<b>PPO Algorithm</b>	Updates the policy to maximize rewards while staying near the reference.

## 2. Data Flow in PPO RLHF

### *Training Loop:*

1. **Sample a batch of prompts** (e.g., from your SFT dataset).
2. **Generate responses** using the current policy model ( $\pi_\phi$ ).
3. **Score responses** with the reward model ( $R_\theta$ ).
4. **Compute KL divergence** between  $\pi_\phi$  and the reference model.
5. **Update policy** using PPO to maximize:

$$\text{Total Reward} = R_\theta(x, y) - \beta \cdot \text{KL}(\pi_\phi(y|x) || \pi_{\text{ref}}(y|x))$$

- $\beta$  = KL penalty weight (e.g., 0.1).

## 3. Implement PPO RLHF

### *Key Libraries:*

```
1 | pip install transformers torch accelerate peft trl
```

- **trl**: HuggingFace's library for RLHF (PPO, DPO).
- **peft**: For efficient fine-tuning (LoRA, QLoRA).

### *Code Implementation:*

```
1  from transformers import AutoModelForCausalLM, AutoTokenizer
2  from trl import PPOTrainer, PPOConfig
3  import torch
4
5  # Load models
6  model = AutoModelForCausalLM.from_pretrained("your_sft_model")
7  ref_model = AutoModelForCausalLM.from_pretrained("your_sft_model") # Frozen reference
8  tokenizer = AutoTokenizer.from_pretrained("your_sft_model")
9  reward_model = AutoModelForSequenceClassification.from_pretrained("your_reward_model")
10
11 # PPO Config
12 ppo_config = PPOConfig(
13     batch_size=32,
14     learning_rate=1e-5,
15     kl_penalty="kl", # Use KL divergence penalty
16     kl_coef=0.1,     # β (KL penalty weight)
17     cliprange=0.2,   # PPO clipping range
18     steps=10000,     # Total training steps
19 )
20
```

```

21 # Initialize PPOTrainer
22 ppo_trainer = PPOTrainer(
23     model=model,
24     ref_model=ref_model,
25     tokenizer=tokenizer,
26     config=ppo_config,
27 )
28
29 # Training loop
30 for epoch in range(3):
31     for batch in prompt_dataset: # Batch of prompts
32         # Generate responses
33         inputs = tokenizer(batch["prompt"], return_tensors="pt", padding=True)
34         outputs = model.generate(**inputs, max_length=1024)
35         response = tokenizer.decode(outputs[0], skip_special_tokens=True)
36
37         # Compute reward
38         reward_inputs = tokenizer(batch["prompt"] + response, return_tensors="pt")
39         reward = reward_model(**reward_inputs).logits[0].item()
40
41         # Compute KL divergence (automatically handled by PPOTrainer)
42         # Update policy
43         ppo_trainer.step([inputs["input_ids"]], [outputs], [reward])

```

## 4. Hyperparameters & Optimization

### *Critical PPO Settings:*

Hyperparameter	Recommended Value	Notes
Batch Size	16-64	Smaller for stability.
Learning Rate	1e-6 to 1e-5	Very low to avoid divergence.
KL Coefficient ( $\beta$ )	0.05-0.2	Balances reward vs. divergence.
Clip Range	0.1-0.3	Clips policy updates for stability.
GAE ( $\lambda$ )	0.9-1.0	For advantage estimation.

### *Advanced Tricks:*

#### 1. Mixed-Precision Training (fp16=True):

- Speeds up training with minimal accuracy loss.

#### 2. LoRA/QLoRA:

- Fine-tune only low-rank adapters to save memory.



```
1 from peft import LoraConfig
2 lora_config = LoraConfig(r=8, lora_alpha=16, target_modules=["q_proj", "v_proj"])
3 model.add_adapter(lora_config)
```

### 3. Gradient Clipping:

- Prevents exploding gradients (max\_grad\_norm=1.0).
- 

## 5. Evaluating PPO Performance

### *Quantitative Metrics:*

#### 1. Mean Reward:

- Average reward per batch (should increase over time).

#### 2. KL Divergence:

- Should stay small (e.g., < 10 nats).

#### 3. Perplexity:

- Compare to the reference model (should not degrade drastically).

### *Qualitative Checks:*

- Manually inspect if:
  - Responses are **more helpful** than SFT.
  - No **reward hacking** (e.g., gibberish with high scores).
  - Avoids **safety violations**.

### *Example Evaluation Code:*

```
1 def evaluate_ppo(model, eval_prompts):
2     for prompt in eval_prompts:
3         inputs = tokenizer(prompt, return_tensors="pt")
4         output = model.generate(**inputs, max_length=512)
5         print(f"Prompt: {prompt}\nResponse: {tokenizer.decode(output[0])}\n---")
```

---

## 6. Common Pitfalls & Fixes

Issue	Solution
Reward Hacking	Increase KL penalty ( $\beta$ ), clip rewards.
Unstable Training	Lower learning rate, smaller batches.
Overoptimization	Early stopping based on KL divergence.
Catastrophic Forgetting	Regularize with reference model.

## Full RLHF Pipeline Summary

Step	Key Actions
1. SFT Model	Fine-tune base model on high-quality (prompt, response) pairs.
2. Reward Model	Train on ranked (prompt, chosen, rejected) pairs.
3. PPO RLHF	Optimize SFT model using RM + KL penalty (this guide).
4. Evaluation	Check reward/KL trends + human A/B tests.

## Example: Training on Anthropic HH-RLHF

```
1 python train_ppo.py \  
2   --sft_model "your_sft_model" \  
3   --reward_model "your_reward_model" \  
4   --dataset "anthropic/hh-rlhf" \  
5   --kl_coef 0.1 \  
6   --batch_size 32 \  
7   --lr 1e-5
```

**Output:** A model that:

- **Outperforms SFT** in human preference tests.
- Avoids harmful outputs (e.g., refuses unsafe requests).

# Next Steps

## 1. Iterative RLHF:

- Collect new human feedback on PPO outputs to improve RM.

## 2. Deployment:

- Quantize the model (GPTQ/GGML) for efficient inference.

## 3. DPO (Optional):

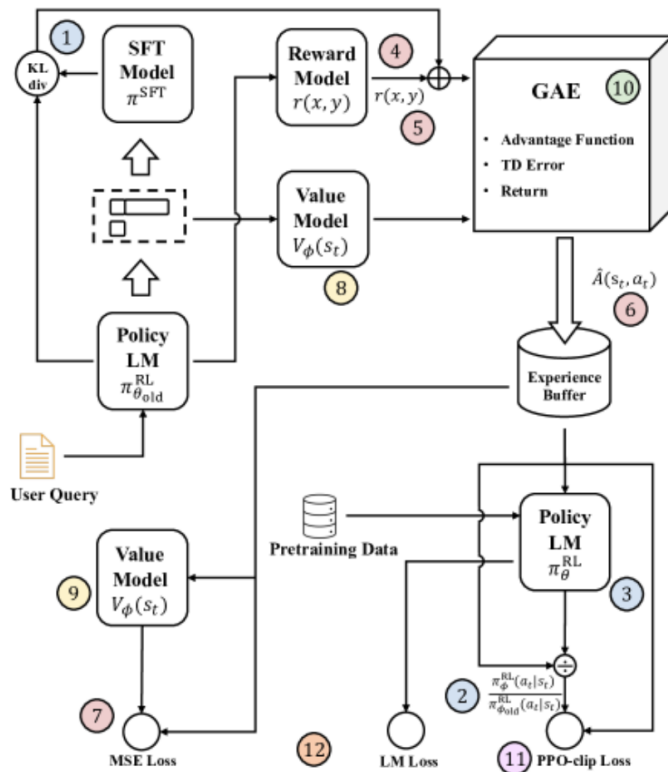
- Direct Preference Optimization (simpler alternative to PPO).

## PPO

**Proximal Policy Optimization (PPO)** is a model-free, on-policy reinforcement learning algorithm used in the **Reinforcement Learning from Human Feedback (RLHF)** framework to fine-tune language models. PPO is widely used in RLHF due to its stability and efficiency.

### *Objectives of PPO in RLHF:*

The main objective of PPO in RLHF is to optimize the policy of a language model such that it generates responses that are not only high-quality but also aligned with human preferences. This is achieved by using human feedback as a reward signal to guide the model's training process.



### Implementation Detail List

★ PPO-Max

★ 1 Token Level KL-Penalty

2 Importance Sampling

3 Entropy Bonus

4 Reward Scaling

★ 5 Reward Normalization and Clipping

6 Advantages Normalization and Clipping

★ 7 Value Function Loss Clipping

★ 8 Critic Model Initialization

9 Policy Model Initialization

★ 10 Generalized Advantage Estimation

★ 11 Clipped Surrogate Objective

★ 12 Global Gradient Clipping

Proximal Policy Optimization (PPO) is a model-free, on-policy reinforcement learning algorithm. In the PPO model, the reward value is used to construct the loss functions for both the policy network and the value network. Below is a detailed explanation of how the reward is used in the loss functions:

## *Policy Network Loss Function*

The objective of the policy network is to learn a policy that maximizes the expected cumulative reward. The PPO algorithm uses a clipped surrogate objective function to update the policy. The policy loss function  $L^{CLIP}(\theta)$  is defined as follows:

$$L^{CLIP}(\theta) = \mathbb{E}_t \left[ \min \left( \frac{\pi_{\theta}(a|s)}{\pi_{\theta_{old}}(a|s)} A(s, a), \text{clip} \left( \frac{\pi_{\theta}(a|s)}{\pi_{\theta_{old}}(a|s)}, 1 - \epsilon, 1 + \epsilon \right) A(s, a) \right) \right]$$

Here:

- $\pi_{\theta}(a|s)$  represents the probability of taking action  $a$  in state  $s$  under the current policy parameterized by  $\theta$ .
- $\pi_{\theta_{old}}(a|s)$  represents the probability of taking action  $a$  in state  $s$  under the old policy parameterized by  $\theta_{old}$ .
- $A(s, a)$  is the advantage function, which estimates how much better the action  $a$  is compared to the average action in state  $s$ . The advantage function is calculated using the rewards obtained from the environment and the value function estimated by the value network.

The policy loss function aims to maximize the expected advantage while limiting the policy update within a trust region defined by the clipping parameter  $\epsilon$ . The clipping ensures that the policy update is not too large, preventing performance degradation.

## *Value Network Loss Function*

The objective of the value network is to accurately estimate the value function, which represents the expected cumulative reward for a given state  $s$ . The value network loss function  $L^{VF}(\theta_v)$  is typically a mean squared error (MSE) loss between the estimated value and the actual returns. The value loss function is defined as follows:

$$L^{VF}(\theta_v) = \mathbb{E}_t \left[ (V_{\theta_v}(s_t) - \hat{V}_t)^2 \right]$$

Here:

- $V_{\theta_v}(s_t)$  represents the estimated value of state  $s_t$  under the current value network parameterized by  $\theta_v$ .
- $\hat{V}_t$  represents the actual returns or the target value, which is calculated using the rewards obtained from the environment and the value function estimated by the value network.

## *Intuition*

- **Policy Network:** The policy network uses the advantage function  $A(s, a)$ , which is derived from the rewards, to determine how much better or worse an action is compared to the average action in a given state. By maximizing the expected advantage, the policy network learns to select actions that lead to higher rewards.

- **Value Network:** The value network uses the rewards to learn the expected cumulative reward for each state. By minimizing the MSE between the estimated values and the actual returns, the value network becomes more accurate in predicting the long-term rewards.

In summary, the PPO model uses rewards to construct the loss functions for both the policy network and the value network, enabling the model to learn a policy that maximizes rewards while accurately estimating the value function. The policy loss function focuses on improving the policy within a trust region, while the value loss function aims to enhance the accuracy of the value function estimation.

## *Policy Network Objective Function:*

The Policy network objective function is designed to update the policy in a way that improves performance while limiting the step size to avoid significant changes that could destabilize the training process.

The objective function is given by:

$$L^{CLIP}(\theta) = \mathbb{E}_t \left[ \min \left( \frac{\pi_\theta(a|s)}{\pi_{\theta_{\text{old}}}(a|s)} A(s, a), \text{clip} \left( \frac{\pi_\theta(a|s)}{\pi_{\theta_{\text{old}}}(a|s)}, 1 - \epsilon, 1 + \epsilon \right) A(s, a) \right) \right]$$

Here are the components of the objective function and their roles:

1. **Policy Ratio:**  $\frac{\pi_\theta(a|s)}{\pi_{\theta_{\text{old}}}(a|s)}$  represents the ratio of the probability of taking action  $a$  in state  $s$  under the new policy  $\pi_\theta$  and the old policy  $\pi_{\theta_{\text{old}}}$ . This ratio indicates how much the policy has changed.
2. **Advantage Function:**  $A(s, a)$  estimates how much better or worse taking action  $a$  in state  $s$  is compared to the average action in that state. It provides a measure of the relative quality of an action.
3. **Clipping:** The clip function limits the policy ratio to the range  $[1 - \epsilon, 1 + \epsilon]$ , where  $\epsilon$  is a small hyperparameter (e.g., 0.1 or 0.2). This clipping ensures that the policy updates are not too large, preventing the new policy from diverging significantly from the old one.

## *Intuition Behind PPO:*

The intuition behind PPO is to make conservative updates to the policy. By clipping the policy ratio, PPO prevents the new policy from deviating too much from the old one. This helps maintain stability in the training process and reduces the risk of performance collapse that can occur with more aggressive policy updates.

By optimizing the clipped surrogate objective, PPO encourages the policy to improve while keeping the updates within a trust region around the previous policy. This balance between improvement and stability makes PPO particularly effective in practice.

In the context of RLHF, PPO uses human feedback to shape the reward signal, guiding the language model to generate responses that better align with human preferences. The PPO algorithm helps the model learn from this feedback efficiently and stably, ultimately producing higher-quality and more aligned outputs.

---

# Data Samples:

In the fine-tuning of ChatGPT, the training data for RLHF PPO training primarily includes prompts and corresponding responses, along with human feedback as rewards. Below is a detailed explanation:

## Data Types

- **Prompts:** These are the input texts provided to the model, such as questions or tasks. For example, "What are the health benefits of apples?" or "Write a short story about space travel."
- **Responses:** These are the outputs generated by the model in response to the prompts. For instance, "Apples are rich in fiber, vitamins, and antioxidants, which help boost immunity, improve digestion, and reduce the risk of chronic diseases" or a short story about space exploration.
- **Rewards:** These reflect human feedback on the quality of the responses. Rewards can be explicit numerical ratings, such as scoring a response on a scale of 1 to 10 based on its relevance, accuracy, and helpfulness. They can also be implicit, such as comparing multiple responses to a prompt and indicating which one is better.

## Data Format

The training data is typically stored in formats like JSONL or CSV. Below is an example of a JSONL format:

```
1 {"prompt": "What are the health benefits of apples?", "response": "Apples are rich in fiber, vitamins, and antioxidants, which help boost immunity, improve digestion, and reduce the risk of chronic diseases.", "reward": 8}
2 {"prompt": "Write a short story about space travel.", "response": "In 2050, astronaut Li Ming boarded the spaceship 'Pioneer' and embarked on a journey to Mars...", "reward": 9}
```

An example of a CSV format is as follows:

Prompt	Response	Reward
What are the health benefits of apples?	Apples are rich in fiber, vitamins, and antioxidants, which help boost immunity, improve digestion, and reduce the risk of chronic diseases.	8
Write a short story about space travel.	In 2050, astronaut Li Ming boarded the spaceship 'Pioneer' and embarked on a journey to Mars...	9

## Data Collection Methods

- **Direct Human Feedback:** Present prompts to human evaluators and have them provide responses, which are then scored or ranked by other evaluators to generate reward signals.
- **Indirect Human Feedback:** Collect data from real user interactions with the model, such as user ratings or clicks on responses, and use this as reward signals.

## Examples

- **Example 1:** In the sentiment-controlled generation of movie reviews, the prompt could be "Generate a positive review for the movie 'The Shawshank Redemption'." The model generates a response like "This movie is incredibly inspiring, with a gripping storyline and superb acting." A human evaluator rates this response as 9/10. Another response might be "The movie is just okay, nothing special," which is rated 4/10. These ratings serve as rewards to guide the model to generate more positive reviews.
- **Example 2:** For a question-answering task, the prompt could be "Who was the first president of the United States?" The model generates two responses: "George Washington" and "Thomas Jefferson." Human evaluators indicate that "George Washington" is the correct answer, while "Thomas Jefferson" is incorrect. The reward for the correct answer is higher, prompting the model to improve its accuracy in answering such questions.

In the OpenRLHF framework, PPO training involves an actor model (policy) that generates responses to prompts, a reward model that evaluates the quality of the generated responses, a reference model to prevent excessive deviation from the original behavior, and a critic model that estimates value functions. The training data includes prompts, and the PPO algorithm optimizes the policy based on human feedback to enhance the quality and alignment of the model's responses. Below is a basic PPO training example using the OpenRLHF framework:

```
1 openrlhf.cli.train_ppo \  
2   --pretrain OpenRLHF/Llama-3-8b-sft-mixture \  
3   --reward_pretrain OpenRLHF/Llama-3-8b-rm-mixture \  
4   --save_path ./checkpoint/llama-3-8b-rlhf \  
5   --prompt_data OpenRLHF/prompt-collection-v0.1 \  
6   --input_key context_messages \  
7   --apply_chat_template
```

In this example, `prompt_data` specifies the dataset containing the prompts. The model generates responses to these prompts, which are then evaluated by the reward model, and PPO training is performed based on the rewards.

---

## PPO in RLHF Training Process

Proximal Policy Optimization (PPO) is a reinforcement learning algorithm used in the RLHF pipeline to fine-tune language models. Here's a detailed explanation of its architecture, components, training process, input data, and output:

### Key Components

- **Actor Model:** The policy network being trained, which generates actions (responses) and provides log probabilities. It is the target language model that we aim to optimize.
- **Critic Model:** Estimates the value function, which represents the expected cumulative reward for being in a certain state. This value estimate is used to compute the advantage function.
- **Reward Model:** Trained on human feedback to provide reward signals for the generated sequences. It serves as the referee, assigning final scores or preference evaluations.
- **Reference Model:** A frozen version of the actor model at an earlier stage (e.g., the initial pre-trained model or a previous version of the fine-tuned model). It helps prevent the policy from deviating too far from the initial policy, mitigating issues like reward hacking.

## *Training Process*

The PPO training process alternates between two main phases:

### 1. Experience Collection (Rollout):

- The actor model generates responses to a batch of prompts.
- The reward model evaluates these responses, producing reward values.
- The critic model estimates the value of the current state (prompt).
- Using the rewards and value estimates, the advantage function is computed for each action (token) in the generated responses.

### 2. Policy Optimization:

- The policy loss is calculated using the clipped surrogate objective function, which involves the policy ratio and the advantage function. The policy ratio is the ratio of the probability of taking an action under the new policy to that under the old policy. The advantage function indicates how much better the action is compared to the average action in the state.
- The value loss is calculated as the mean squared error between the critic model's estimated values and the actual returns (which are derived from the rewards and the value function estimates).
- The policy network is updated to maximize the clipped surrogate objective, while the critic network is updated to minimize the value loss.

## *Input Data*

- **Prompts:** These are the input texts provided to the model, such as questions or tasks.
- **Responses:** The outputs generated by the model in response to the prompts.
- **Rewards:** Explicit numerical ratings or implicit feedback (e.g., rankings) from human evaluators on the quality of the responses.

## *Output of the Model*

- **Updated Policy:** The actor model, after training, generates responses that are more aligned with human preferences.
- **Value Function Estimates:** The critic model provides more accurate estimates of the expected cumulative rewards for different states.

## *Training Workflow in OpenRLHF*

1. **Initialization:** Load the pre-trained actor model, critic model, and reward model. The reference model is also initialized, typically with the same parameters as the initial actor model.
2. **Rollout:** The actor model generates responses to a batch of prompts. The reward model evaluates these responses, and the critic model estimates the state values. These are used to compute the advantages.
3. **Optimization:** Using the computed advantages and rewards, the policy loss and value loss are calculated. The actor and critic models are then updated using gradient descent to minimize these losses.



4. **Iteration:** Steps 2 and 3 are repeated for multiple epochs until the model converges to a satisfactory level of performance.

### *Advantages of PPO in RLHF*

- **Stability:** The clipped objective function in PPO provides more stable training updates compared to simpler policy gradient methods. This is crucial for large language models where training is expensive and prone to divergence.
- **Sample Efficiency:** PPO reuses data collected over multiple epochs of updates within each data collection phase, making it more sample-efficient than basic policy gradient methods like REINFORCE.
- **Implementation Complexity:** PPO is relatively easier to implement and tune compared to some alternatives like Trust Region Policy Optimization (TRPO), which involves second-order optimization.

This architecture and training process enable PPO to effectively use the reward signals from the reward model to guide the language model towards generating outputs that better align with human preferences, while maintaining stability during fine-tuning.