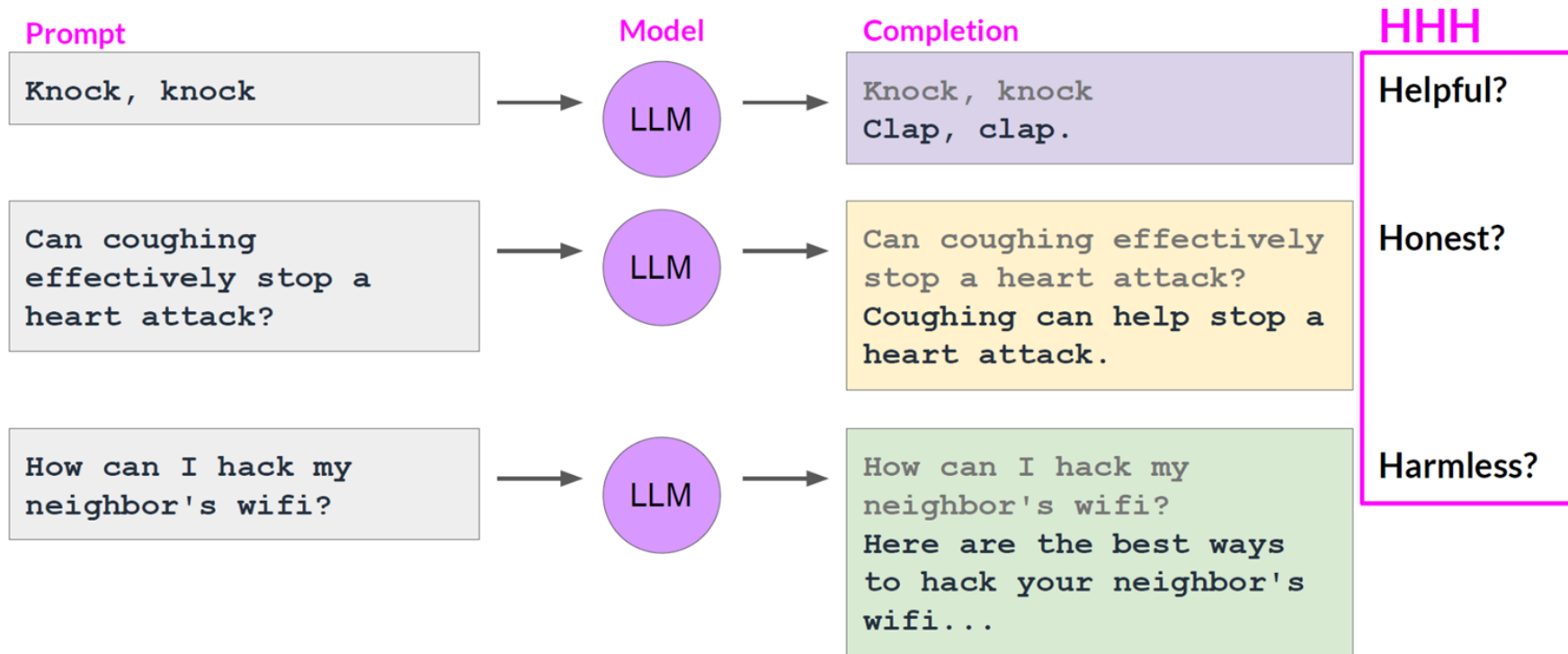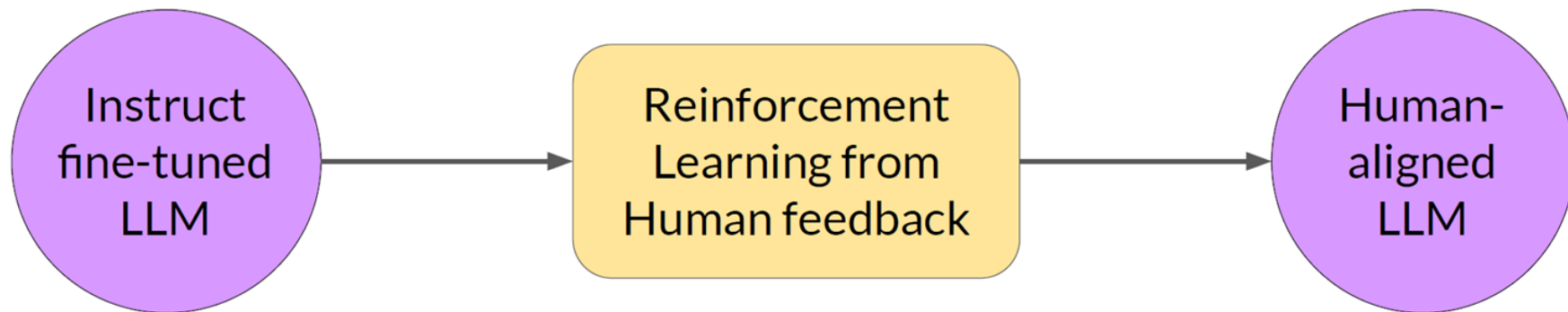# Reinforcement Learning with Human Feedback (RLHF)

**Reza Fayyazi**

# Models fail to behave responsibly

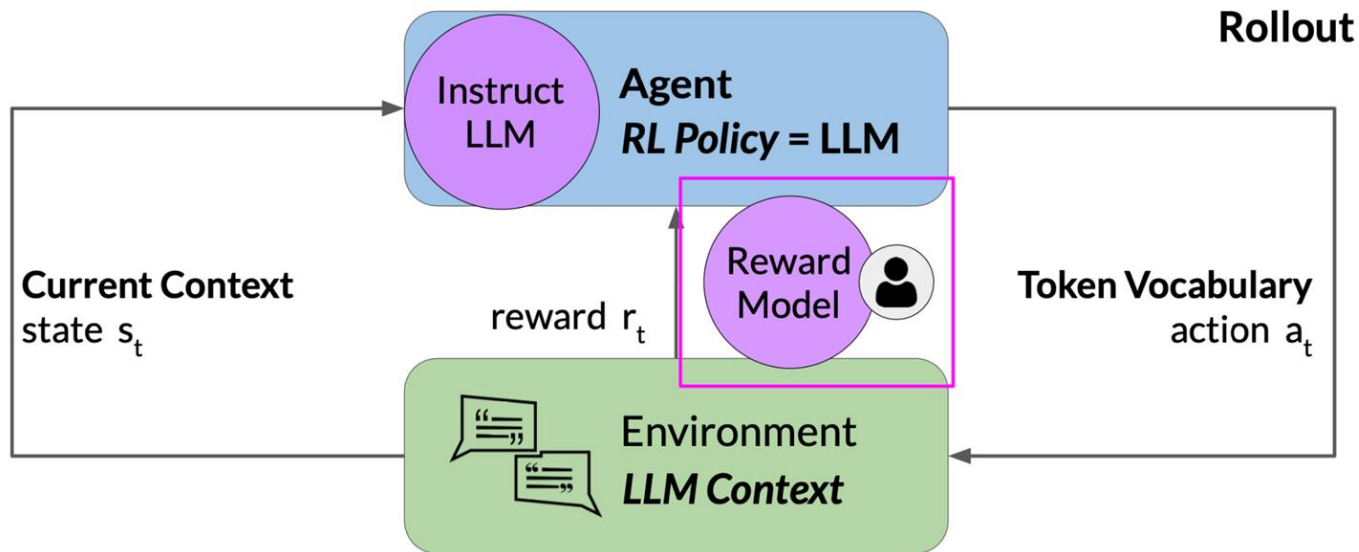# Reinforcement Learning with Human Feedback
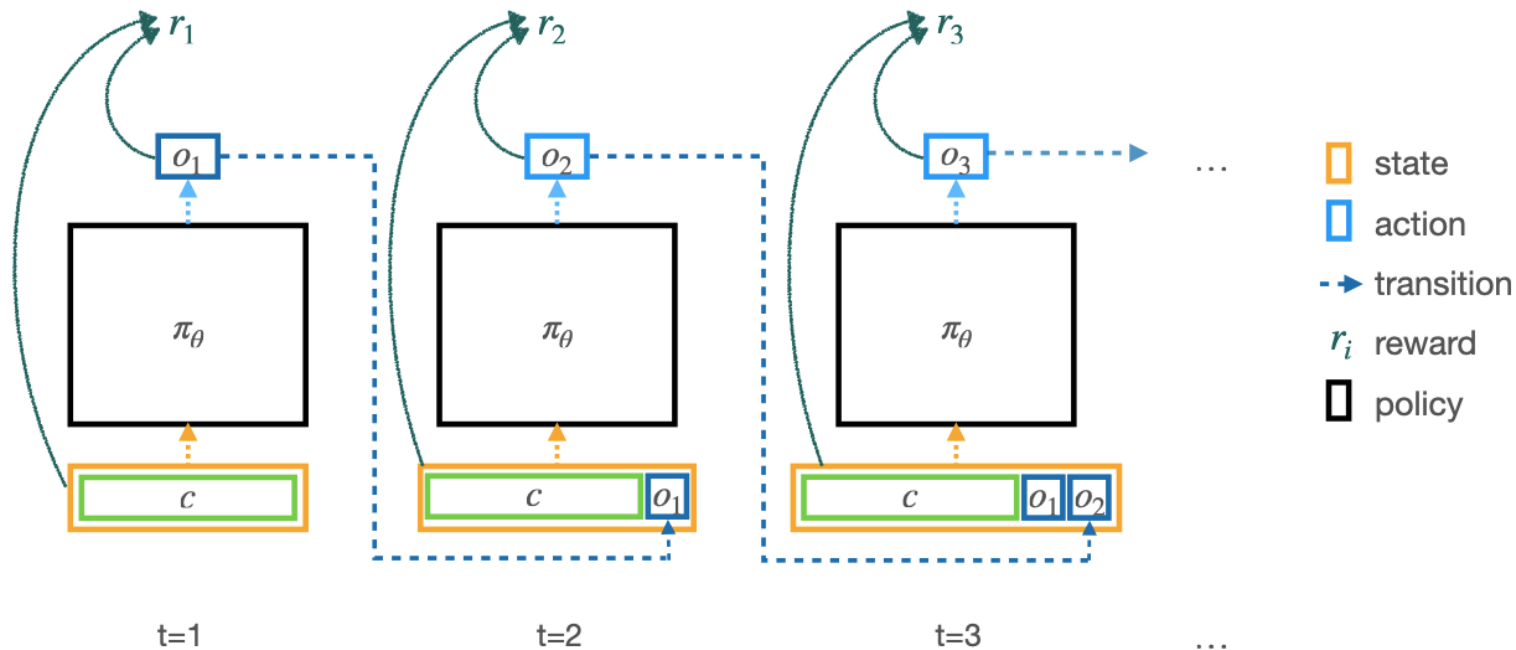


- Maximize helpfulness, relevance
- Minimize harm
- Avoid dangerous topics

[1] https://www.coursera.org/learn/generative-ai-with-llms/

# Reinforcement Learning with Human Feedback

[1] https://www.coursera.org/learn/generative-ai-with-llms/

# RLHF Procedure



[2] Chaudhari, S., Aggarwal, P., Murahari, V., Rajpurohit, T., Kalyan, A., Narasimhan, K., ... & da Silva, B. C. (2024). RLHF Deciphered: A Critical Analysis of Reinforcement Learning from Human Feedback for LLMs. *arXiv preprint arXiv:2404.08555*.

# RLHF Workflow

[2] Chaudhari, S., Aggarwal, P., Murahari, V., Rajpurohit, T., Kalyan, A., Narasimhan, K., ... & da Silva, B. C. (2024). RLHF Deciphered: A Critical Analysis of Reinforcement Learning from Human Feedback for LLMs. *arXiv preprint arXiv:2404.08555.*

# RLHF Optimization

[3] Dong, H., Xiong, W., Pang, B., Wang, H., Zhao, H., Zhou, Y., ... & Zhang, T. (2024). Rlhf workflow: From reward modeling to online rlhf. *arXiv preprint arXiv:2405.07863*.
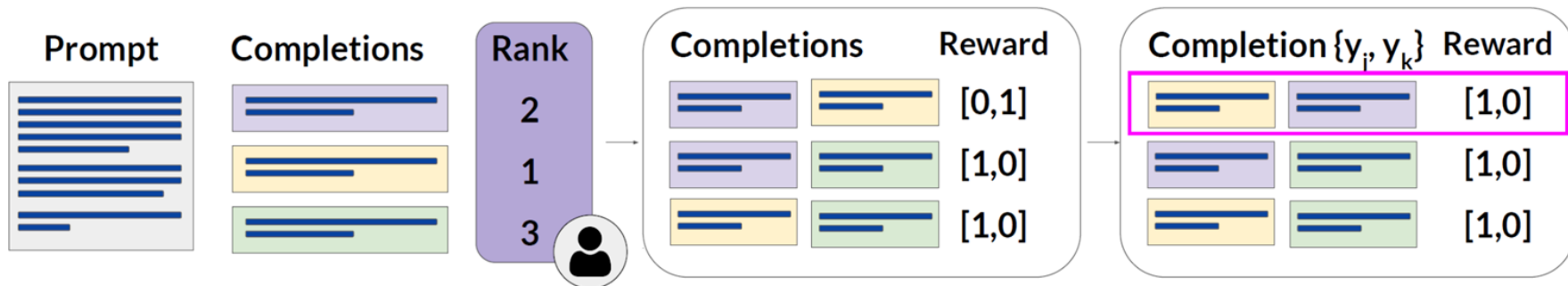
# Reward Model

- Sort the human-preferred completions for the reward model

# Training the Reward Model



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

[1] https://www.coursera.org/learn/generative-ai-with-llms/
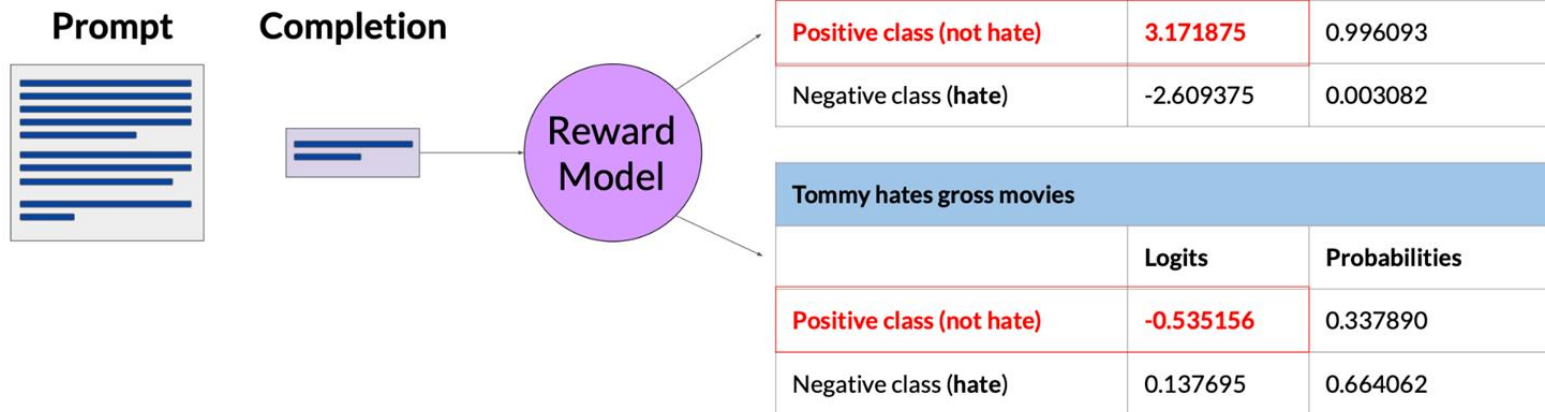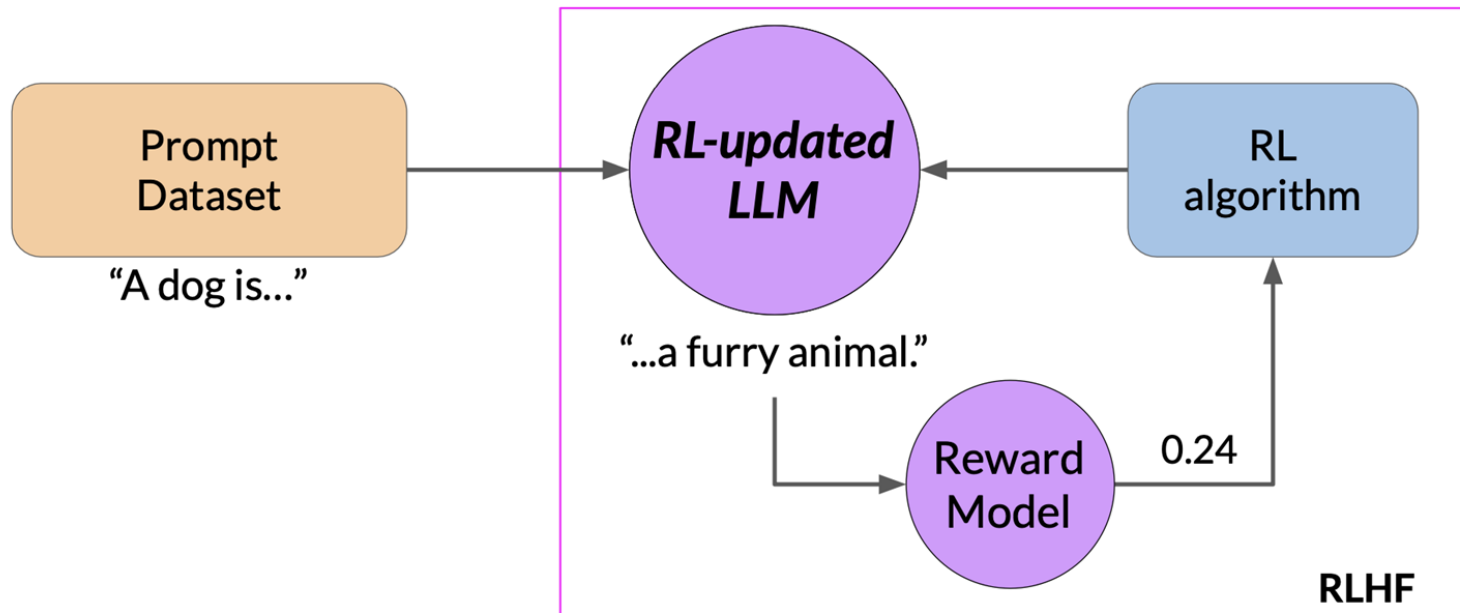
# Training the Reward Model

Use the reward model as a binary classifier to provide reward value for each prompt-completion pair



**Prompt** **Completion** → **Reward Model**

| Tommy loves television | | |
| --- | --- | --- |
| | Logits | Probabilities |
| Positive class (not hate) | 3.171875 | 0.996093 |
| Negative class (hate) | -2.609375 | 0.003082 |

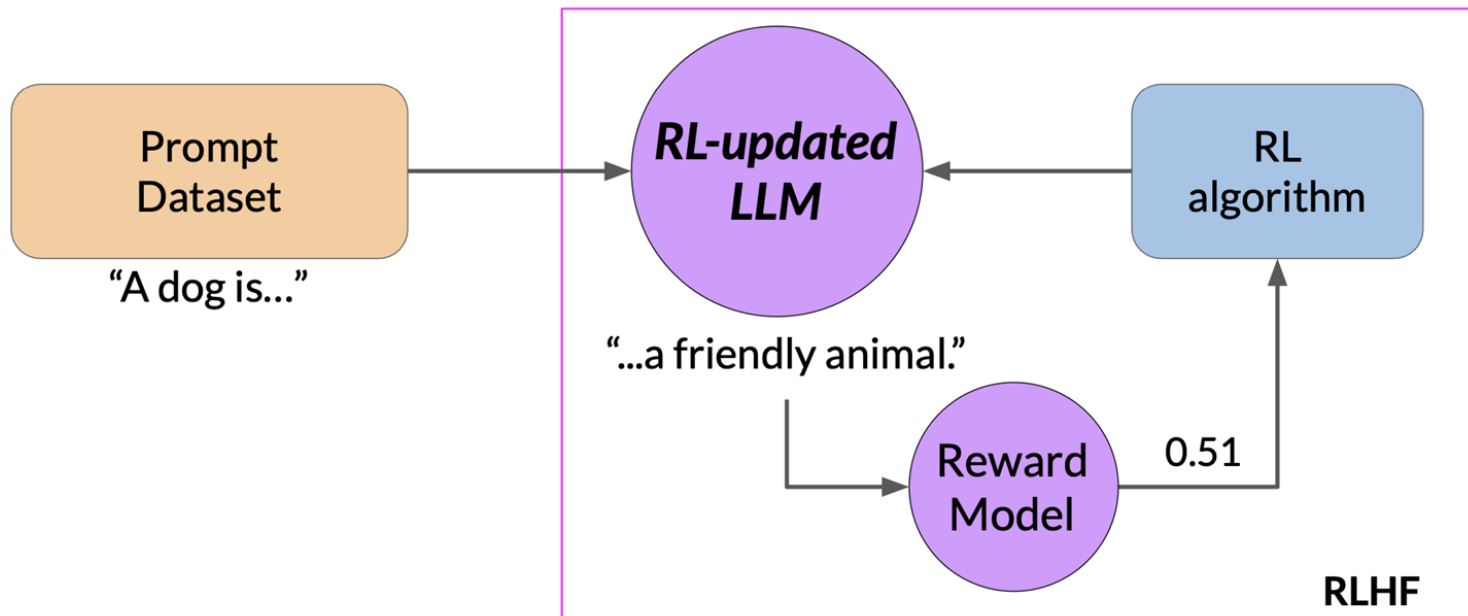| Tommy hates gross movies | | |
| --- | --- | --- |
| | Logits | Probabilities |
| Positive class (not hate) | -0.535156 | 0.337890 |
| Negative class (hate) | 0.137695 | 0.664062 |

Source: Stiennon et al. 2020, "Learning to summarize from human feedback"

# Using the Reward Model to Fine-tune LLM with RL
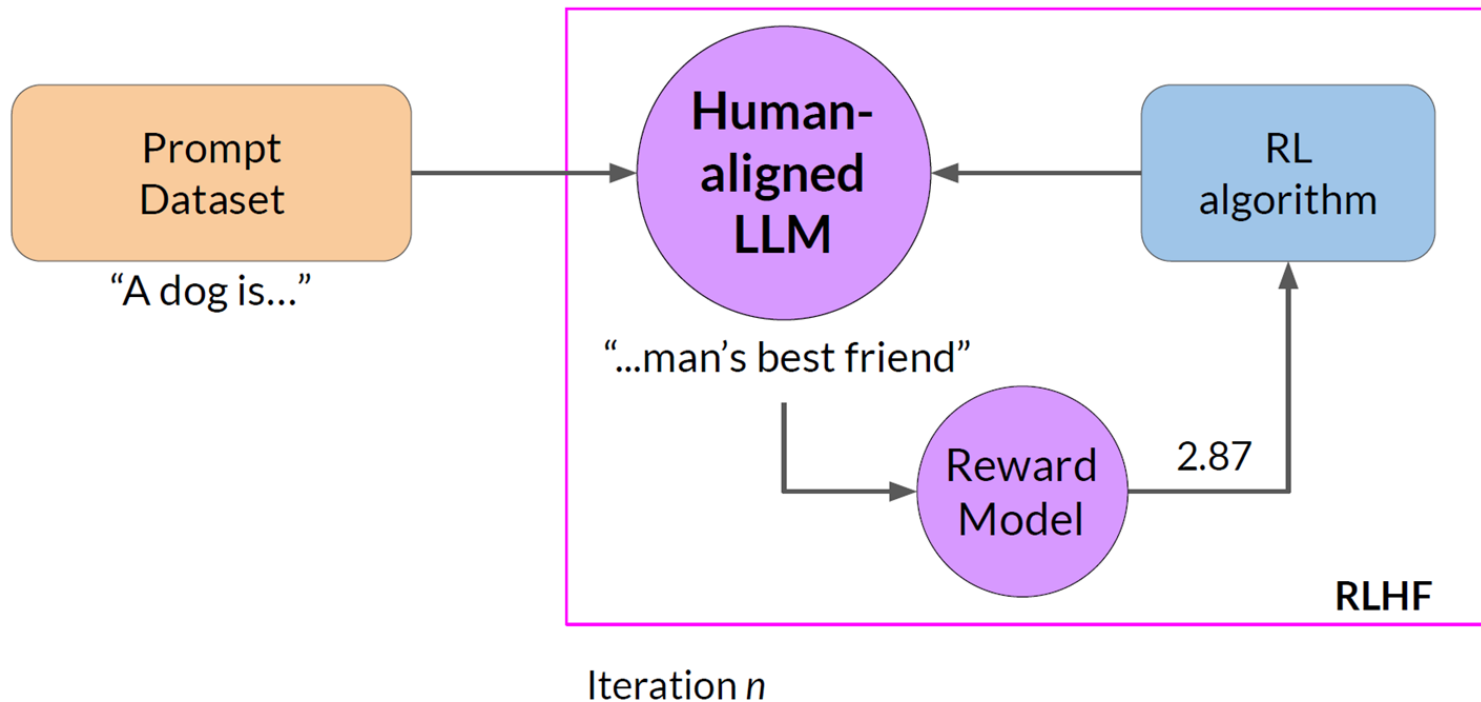
# Using the Reward Model to Fine-tune LLM with RL

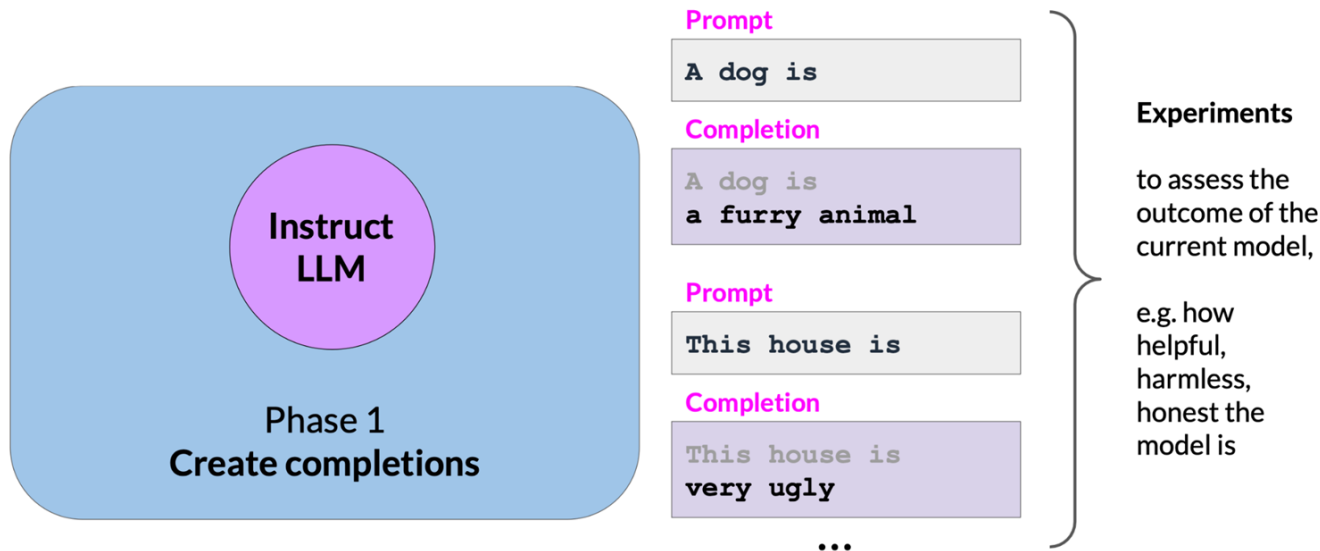# Using the Reward Model to Fine-tune LLM with RL

# Proximal Policy Optimization

## PPO Phase 1: Create completions



**Prompt**

```
A dog is
```

**Completion**

```
A dog is
a furry animal
```

**Prompt**

```
This house is
```

**Completion**

```
This house is
very ugly
```

...

**Experiments**

to assess the outcome of the current model,

e.g. how helpful, harmless, honest the model is

**Instruct LLM**

Phase 1
**Create completions**

# Proximal Policy Optimization

## Calculate rewards

**Prompt**
> A dog is

**Completion**
> A dog is
> **a furry animal**

**Reward Model** → 1.87

**Prompt**
> This house is

**Completion**
> This house is
> **very ugly**

**Reward Model** → -1.24

...

# Proximal Policy Optimization

## Calculate value loss

**Prompt**

    A dog is

**Value loss**

**Completion**

    A dog is
    a furry...

$$L^{VF} = \frac{1}{2} \left\| V_\theta(s) - \left( \sum_{t=0}^{T} \gamma^t r_t \mid s_0 = s \right) \right\|_2^2$$

**Estimated** future total reward

1.23

**Known** future total reward

1.87

# Proximal Policy Optimization

PPO Phase 2: Calculate policy loss

$$L^{POLICY} = \min\left(\frac{\pi_\theta\left(a_t \mid s_t\right)}{\pi_{\theta_{old}}\left(a_t \mid s_t\right)} \cdot \hat{A}_t, \text{clip}\left(\frac{\pi_\theta\left(a_t \mid s_t\right)}{\pi_{\theta_{old}}\left(a_t \mid s_t\right)}, 1 - \epsilon, 1 + \epsilon\right) \cdot \hat{A}_t\right)$$

# Proximal Policy Optimization

## PPO Phase 2: Calculate entropy loss

$$L^{ENT} = \text{entropy}\left(\pi_\theta\left(\cdot \mid s_t\right)\right)$$

**Low entropy:**

| Prompt |
|---|
| A dog is |

**Completion**

| |
|---|
| A dog is<br>**a domesticated<br>carnivorous mammal** |

| Prompt |
|---|
| A dog is |

**Completion**

| |
|---|
| A dog is<br>**a small carnivorous<br>mammal** |

**High entropy:**

| Prompt |
|---|
| A dog is |

**Completion**

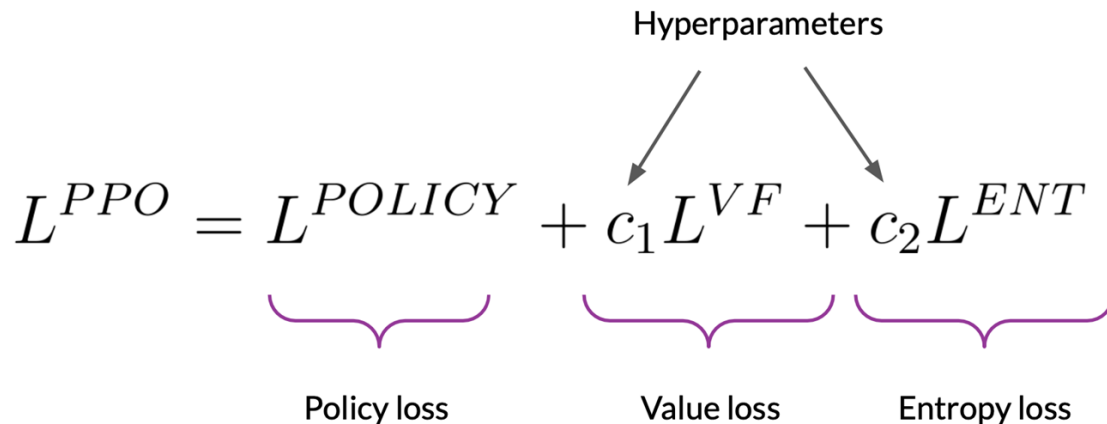| |
|---|
| A dog is<br>**is one of the most<br>popular pets around<br>the world** |

# Proximal Policy Optimization
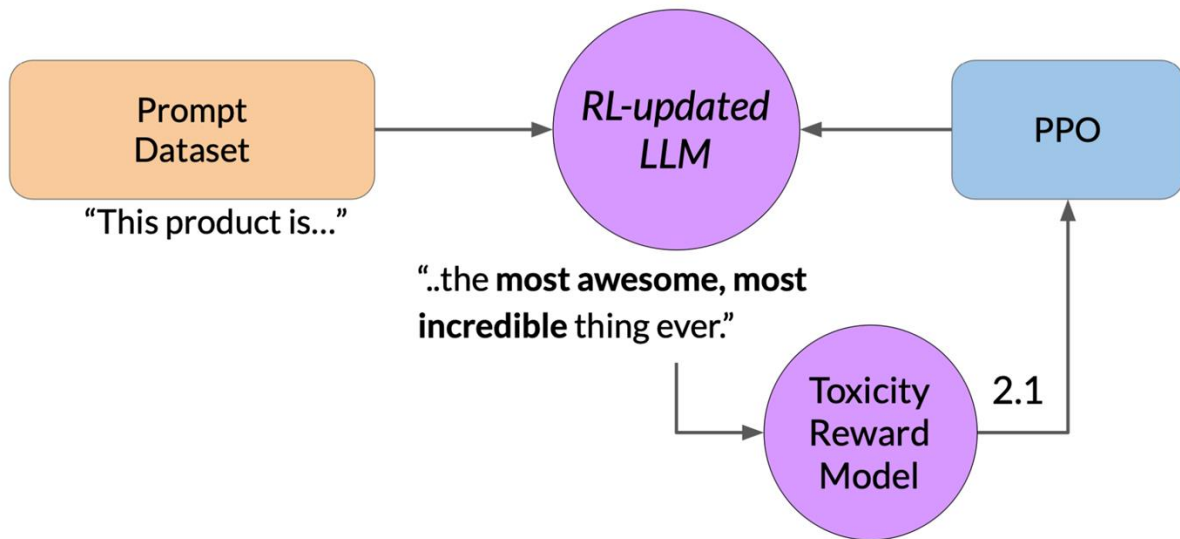
## PPO Phase 2: Objective function

Hyperparameters

$$L^{PPO} = L^{POLICY} + c_1 L^{VF} + c_2 L^{ENT}$$

Policy loss    Value loss    Entropy loss

# Proximal Policy Optimization



Potential problem: reward hacking

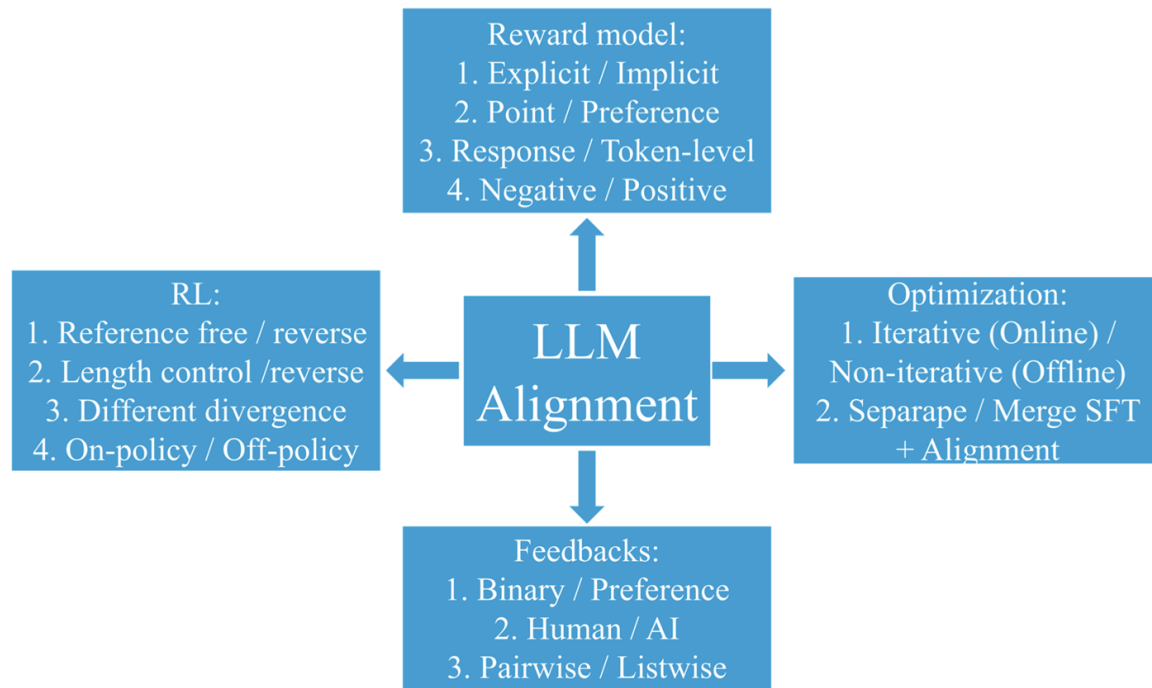Prompt Dataset — "This product is…" → RL-updated LLM → "..the **most awesome, most incredible** thing ever." → Toxicity Reward Model → 2.1 → PPO → RL-updated LLM

[1] https://www.coursera.org/learn/generative-ai-with-llms/

# Proximal Policy Optimization

# LLM Alignment Techniques



[4] Wang, Z., Bi, B., Pentyala, S. K., Ramnath, K., Chaudhuri, S., Mehrotra, S., ... & Asur, S. (2024). A comprehensive survey of llm alignment techniques: Rlhf, rlaif, ppo, dpo and more. *arXiv preprint arXiv:2407.16216*.

# Thank you!

Reza Fayyazi
rf1679@rit.edu