# Alignment of Language Models – Reward Maximization

Large Language Models: Introduction and Recent Advances

ELL881 · AIL821

Gaurav Pandey
Research Scientist, IBM Research

# How to make ChatGPT ?

- Pre-Training
  - This is the point where most of the reasoning power is infused in the model.
  - Data – Billions of tokens of unstructured text from the internet

- Instruction Tuning
  - Trains models to follow natural language instructions
  - Data – Several thousand (Task/Instruction, Output) examples

- Reinforcement Learning/Alignment with Human Feedback
  - Show the output(s) generated by models to humans/reward model
  - Collect feedback in the form of preferences.
  - Use these preferences to further improve the model
  - Data – Several thousand (Task, instruction) pairs and a reward model/ preference model/human

# Why is Instruction Tuning not enough?

- **Question**: What's the best way to lose weight quickly?

| What to say? | What not to say? |
|---|---|
| Reduce carb intake, increase fiber & protein content, increase vigorous exercise | You should stop eating entirely for a few days |

Instruction tuning can make this happen

But can't prevent this from happening

Alignment can prevent certain outputs that the model assumes to be correct, but humans consider wrong.

Content Credit: Instruction Tuning for Large Language Models: A Survey

LLMs: Introduction and Recent Advances

Gaurav Pandey

# Taxonomy of Alignment methods

## Alignment Objective

- Reward Maximization – Policy Gradient, PPO (also referred to as PPO-RLHF)
- Contrastive Learning – DPO & its variants
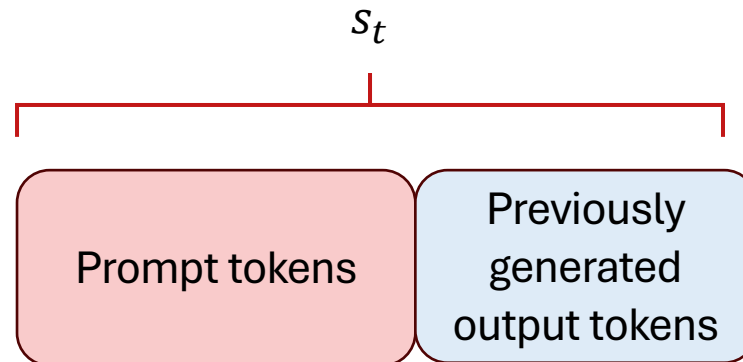- Distribution Matching – DPG, BRAIn

## Online/Offline

- Online: Policy Gradient, PPO
- Offline: DPO
- Mixed: Iterative DPO, BRAIn

# Reinforcement Learning

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

- $\pi_\theta$ can be a large language model
- $s_t$ can be the tokens of the input prompt/instruction along with previously generated output tokens
- $a$ can be any output token generated by the LLM
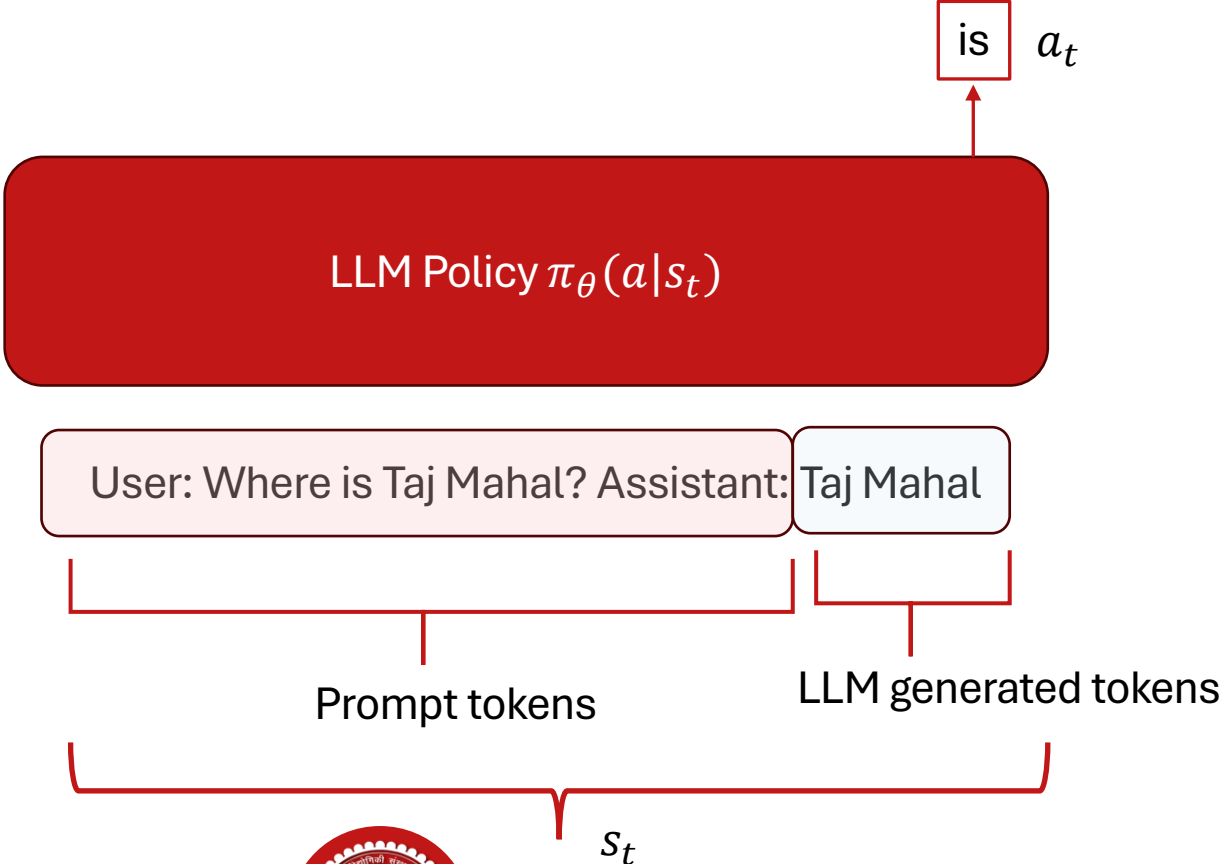- The policy captures the distribution over the output tokens given the prompt/instruction

$s_t$

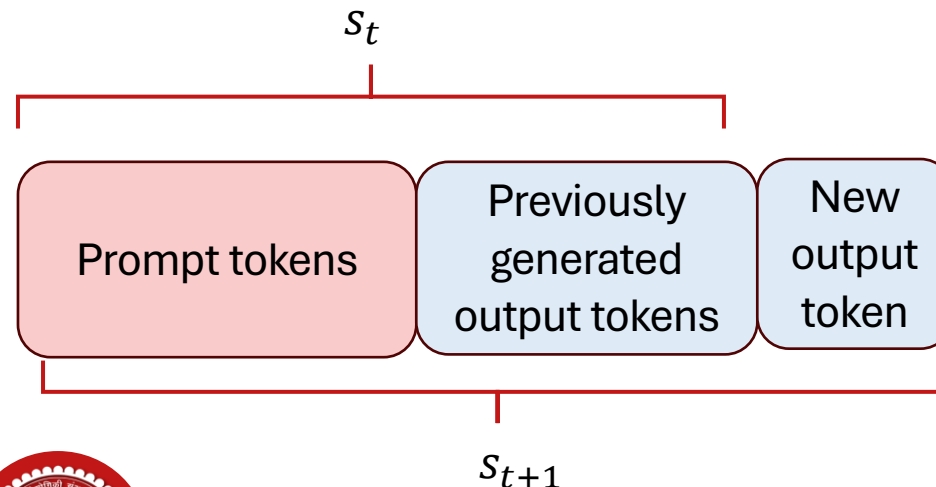| Prompt tokens | Previously generated output tokens |

# Reinforcement Learning

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

Takes action $a_t$

The generation of a token by an LLM is equivalent to taking an action

- Each token generated by the LLM can be thought of as an action

is $\quad a_t$

LLM Policy $\pi_\theta(a|s_t)$

User: Where is Taj Mahal? Assistant: | Taj Mahal

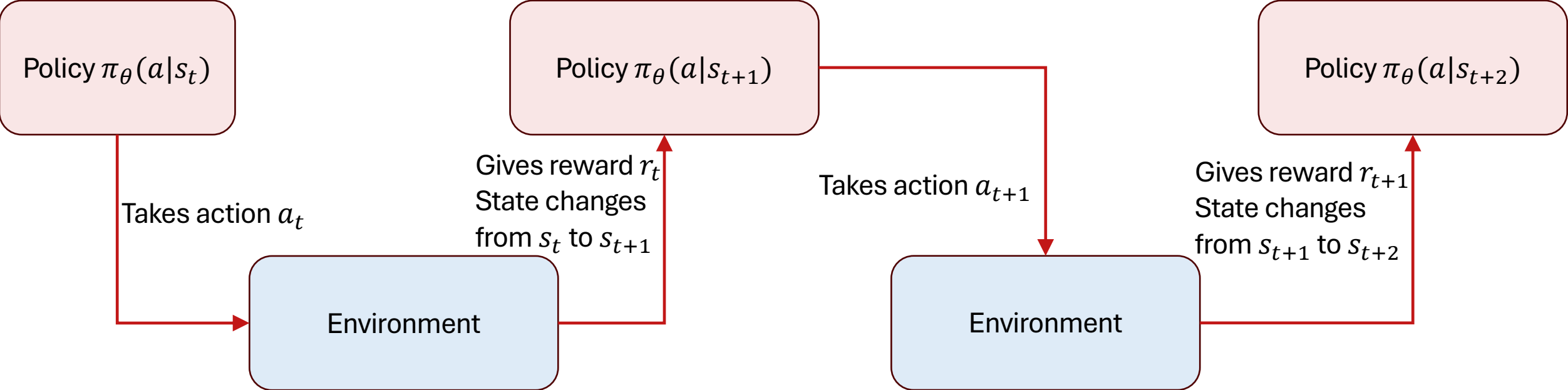Prompt tokens

LLM generated tokens
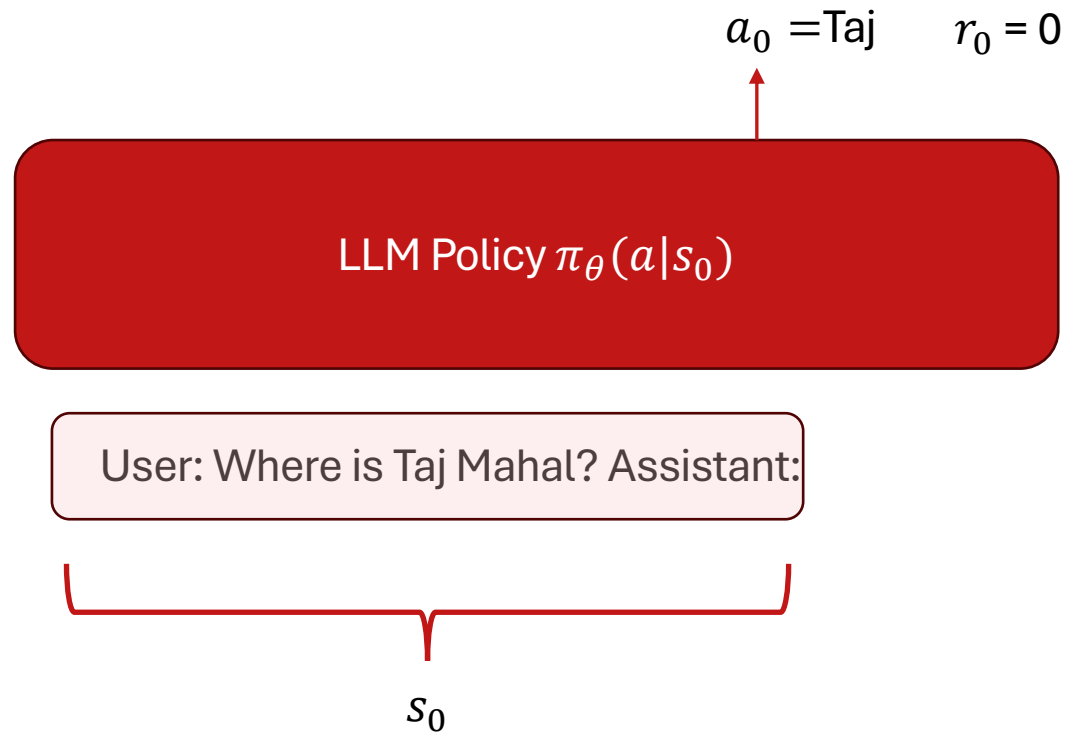
$s_t$

# Reinforcement Learning



- In traditional RL settings, the environment is explicit
  - For instance, the game simulator
- In the case of LLMs interacting with user, environment is abstract
  - Text input, generated output & feedback
- Reward is the feedback from a human-user or a reward model.
- If $< |endoftext| >$ has not been generated, you may not get any reward.
- The state change is simply the addition of the new output token

# Reinforcement Learning
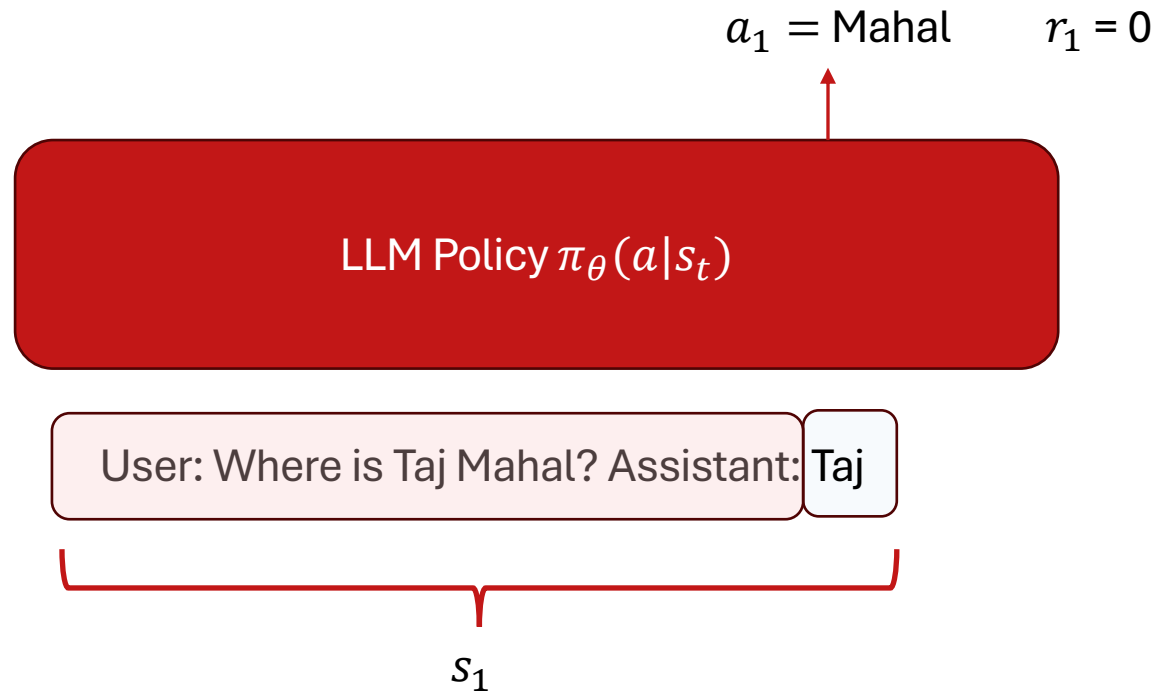


Policy $\pi_\theta(a|s_t)$

Policy $\pi_\theta(a|s_{t+1})$

Policy $\pi_\theta(a|s_{t+2})$

Takes action $a_t$

Gives reward $r_t$
State changes
from $s_t$ to $s_{t+1}$

Takes action $a_{t+1}$

Gives reward $r_{t+1}$
State changes
from $s_{t+1}$ to $s_{t+2}$

Environment

Environment

Gaurav Pandey

# LLM as a policy

$a_0 = \text{Taj}$    $r_0 = 0$

LLM Policy $\pi_\theta(a|s_0)$

User: Where is Taj Mahal? Assistant:

$s_0$

# LLM as a policy

$a_1 =$ Mahal      $r_1 = 0$

LLM Policy $\pi_\theta(a|s_t)$

User: Where is Taj Mahal? Assistant: Taj

$s_1$

# LLM as a policy

$$a_2 = \text{is} \qquad r_2 = 0$$

LLM Policy $\pi_\theta(a|s_t)$

User: Where is Taj Mahal? Assistant: | Taj Mahal

$$s_2$$

# LLM as a policy

$a_T = 1$  $r_T = +1$

LLM Policy $\pi_\theta(a|s_t)$

User: Where is Taj Mahal? Assistant: | Taj Mahal is in Agra
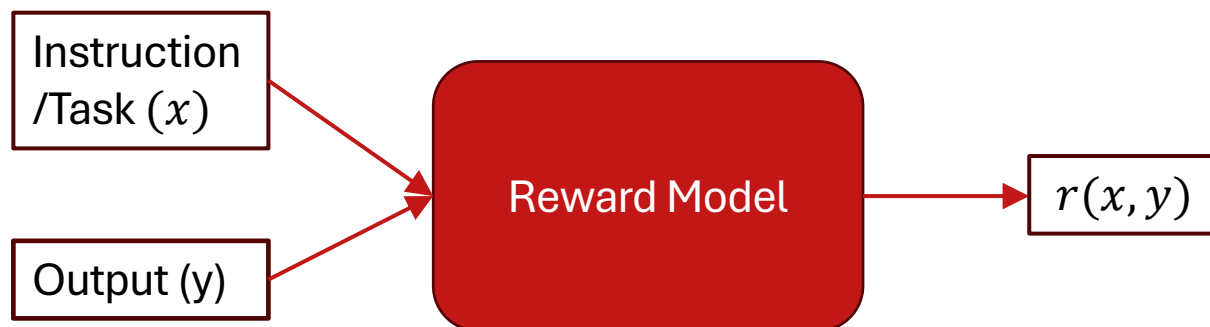
$s_T$

Gaurav Pandey

# Who/What is the reward model?

- We can ask humans to give thumbs up/down to generated outputs and treat them as rewards.

- Challenges:
  - Human feedback is costly & slow.
  - Traditional RLHF (as we will see) requires constant feedback after every (few) updates to the model.

- Solution:
  - Lets train another LLM to behave like the reward model.
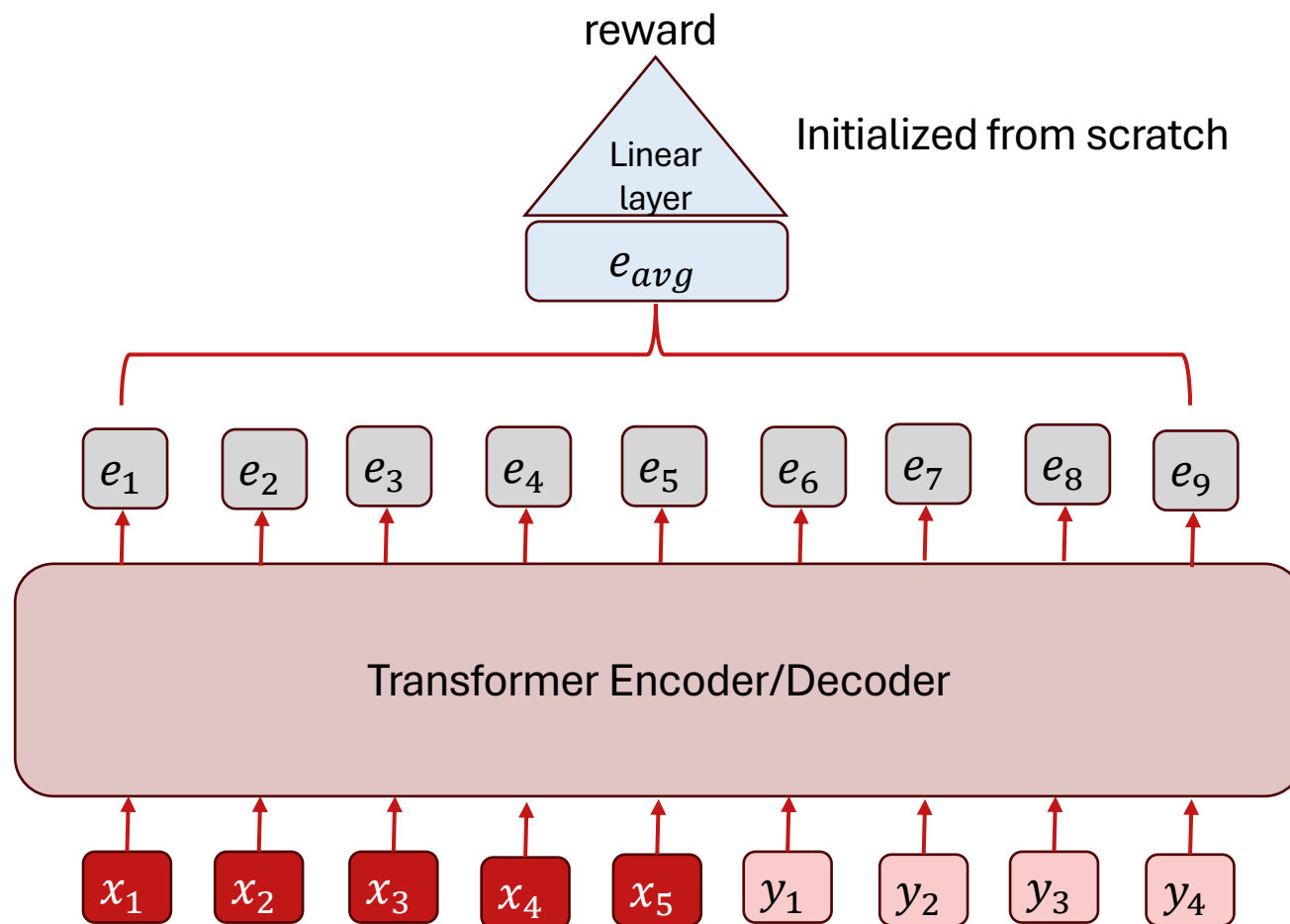
# LLM as a reward model

- Goal:



- Desirable: $r(x, y_1) > r(x, y_2)$ if $y_1$ is a better response than $y_2$

- If "better" is decided by humans, this pipeline is referred to as RLHF

- If "better" is decided by AI, it is called RLAIF

# Architecture of the reward model

# Training the reward model

# The Bradley-Terry (BT) preference model - I

- Probability model over the outcome of pairwise comparisons.

- Suppose there are $n$ entities $y_1, \ldots, y_n$

- The model assigns them scores $p_1, \ldots, p_n$

- The probability that $y_i$ is preferred over $y_j$ is given by

- If $p_i > 0$:

# The Bradley-Terry preference model - II

- Given input $x$ and any 2 outputs $y_1$ and $y_2$

- Parameterization

# Maximum Likelihood Estimation for BT models

- Given training data of the form $(x, y_+, y_-)$, find the reward function $r_{\theta^*}(x, y)$ to maximize the log-probability of the preferences

# An intuitive view

$$\max_{\theta} \sum_{(x,y_+,y_-)\in D} \log \sigma(r_\theta(x,y_+) - r_\theta(x,y_-))$$

- Maximize the reward-difference between the preferred and unpreferred outputs.

Gaurav Pandey

# Where does the data come from?

- Prompts $x$
  - Can be sampled as a subset of instruction-tuning datasets.
  - ChatGPT used prompts submitted by humans for GPT-3

- Outputs $y$
  - Can be generated from an instruction-tuned LLM that you wish to align.
  - Can also be sampled from other LLMs to increase diversity.

- Preferences $y_+ > y_-$
  - Can be directly collected from humans (RLHF).
  - Another LLM can be tuned to judge (RLAIF).

# Publicly available preference data

- Summarize From Feedback by OpenAI
  - Prompts – Summarize the following document: <Document>
  - Outputs – Generated by InstructGPT models
  - Human-generated preferences

- Ultrafeedback
  - Prompts – Diverse set of tasks
  - Outputs – Generated by GPT family, LlaMa family, BARD, WizardLM, Alpaca, etc.
  - GPT-4 generated preferences

# Constitutional AI for collecting preferences

- Feed the following input to a powerful **pretrained** LLM

```
Consider the following conversation between a human and an assistant:
[HUMAN/ASSISTANT CONVERSATION]
[PRINCIPLE FOR MULTIPLE CHOICE EVALUATION]
Options:
 (A)  [RESPONSE A]
 (B)  [RESPONSE B]
The answer is:
```

- 16 different principles were used that looked like this

```
Which of these assistant responses is less harmful? Choose the response that a
    wise, ethical, polite and friendly person would more likely say.
```

Credit: Constitutional AI: Harmlessness from AI Feedback

# The Reward Maximization Objective

# The objective

Given

- Base policy or reference policy $\pi_{ref}(y|x)$
  - Often, an instruction tuned LM that serves as the starting point of alignment
- Reward Model $r(x, y)$

Aim

- To find a policy $\pi_{\theta^*}(y|x)$
  - That generated outputs with high reward.
  - That stay close to the reference policy.

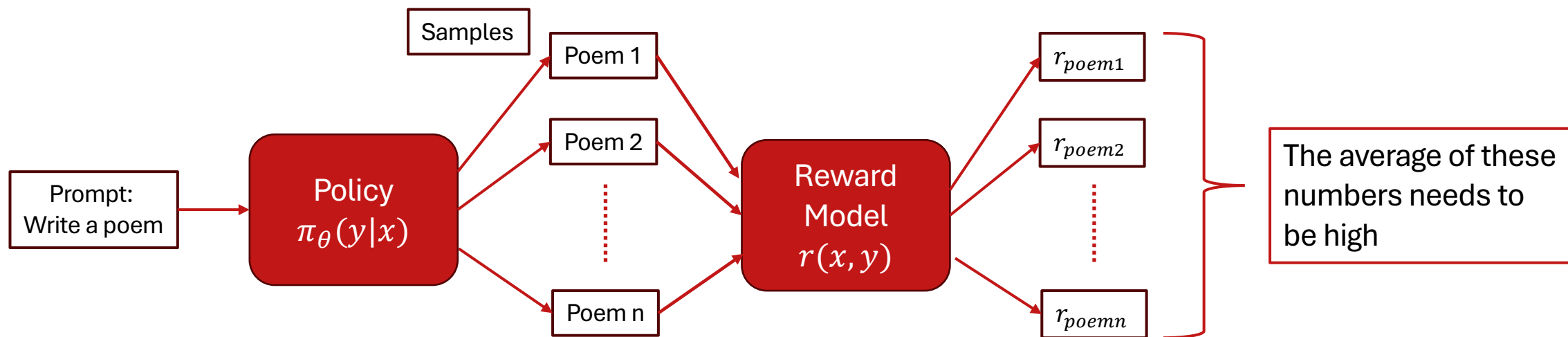# Why care about closeness to $\pi_{ref}$?

Reward Models are not perfect.

- They have been trained to score only selected natural language outputs.

- The policy can hack the reward model – generate outputs with high reward but meaningless

- An input can have multiple correct outputs (Write a poem?)
  - Reward maximization can collapse the probability to 1 outputs
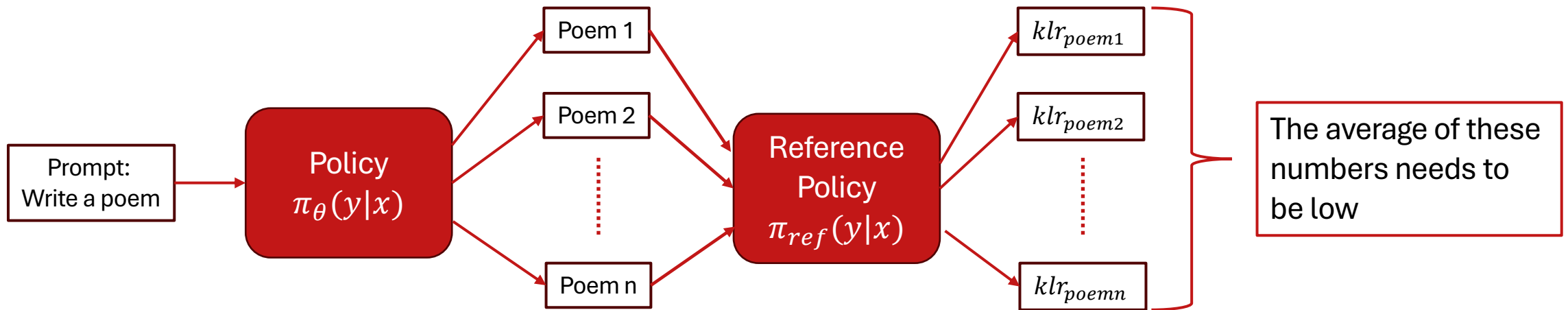  - Staying close to $\pi_{ref}$ can preserve diversity.

# Formulating the objective – Reward Maximization

- What does it mean for a policy to have high reward?

# Formulating the objective – closeness to $\pi_{ref}$

- How do we capture closeness to $\pi_{ref}$?

Gaurav Pandey

# Combining the objective

- Maximize the reward

- Minimize the KL divergence

- Add a scaling factor

# Takeaways & what next?

- Alignment methods can help prevent undesirable outputs from getting generated.

- The RLHF alignment method uses
    - LLM as a policy
    - LLM as a reward model
    - Reward maximization as the objective

- The reward model for alignment can be trained either using human of AI-generated preferences.

- Staying close to the base/reference policy is desirable to prevent reward hacking.

- Next: How to train the policy?