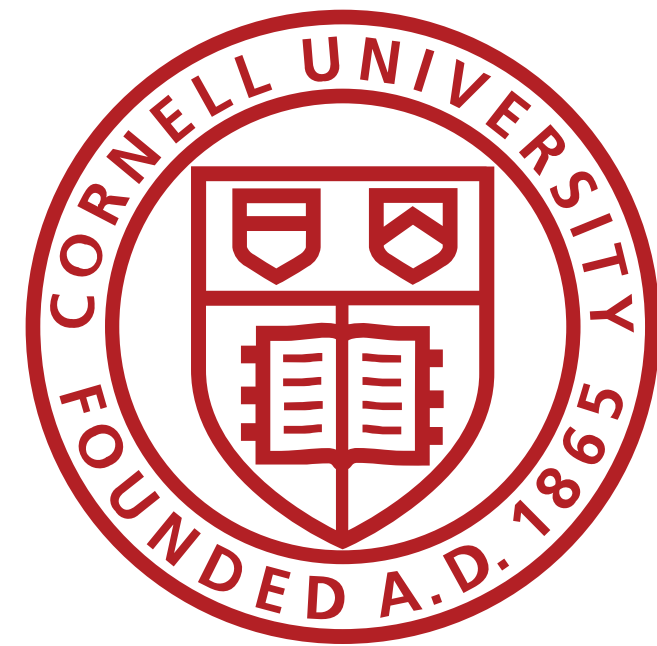


# Learning from Interaction

**Kianté Brantley | Postdoctoral Associate | Cornell University**



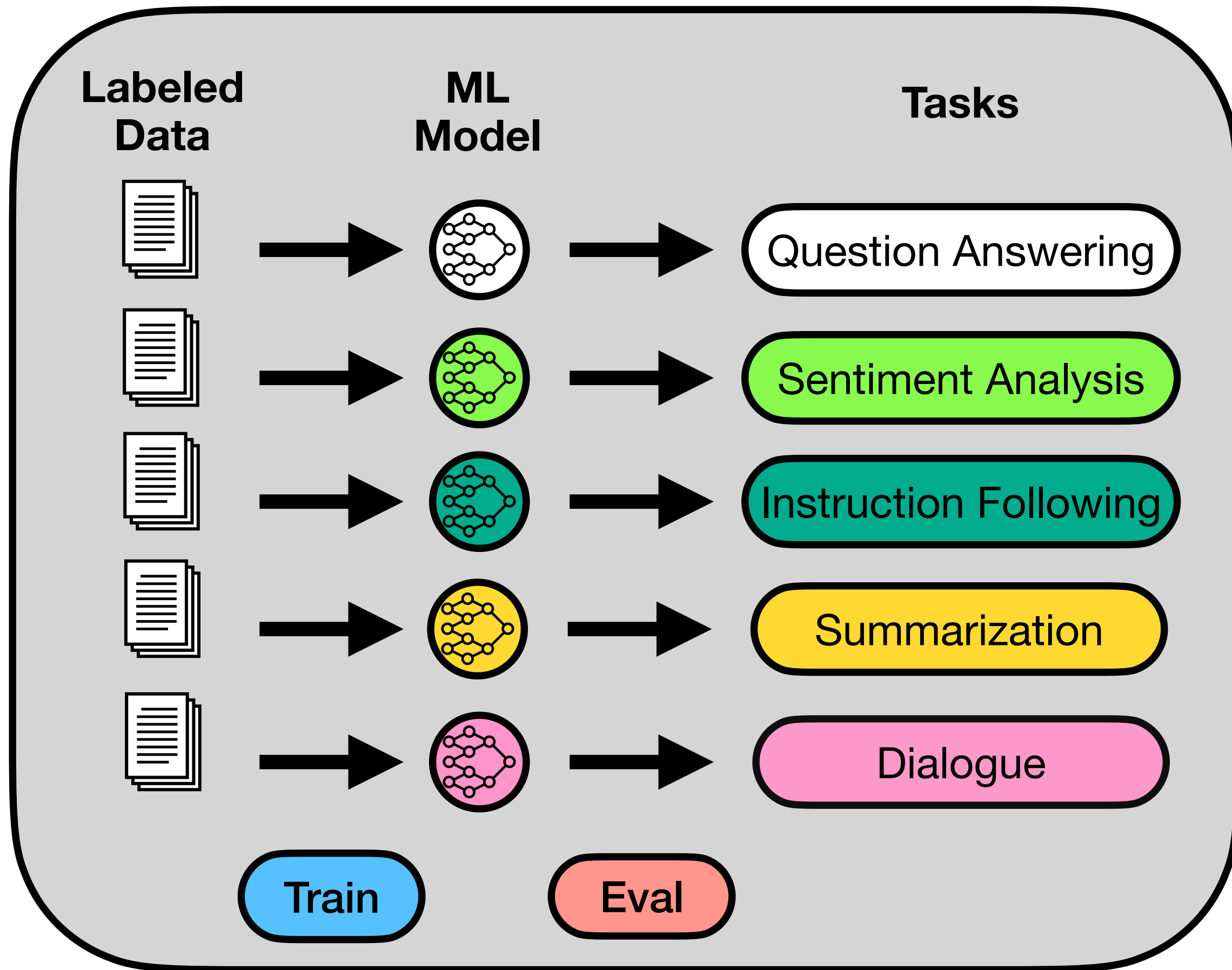
Cornell University®

“A foundation model is any model that is trained on broad data that can be adapted (e.g., fine-tuned) to a wide range of downstream tasks”

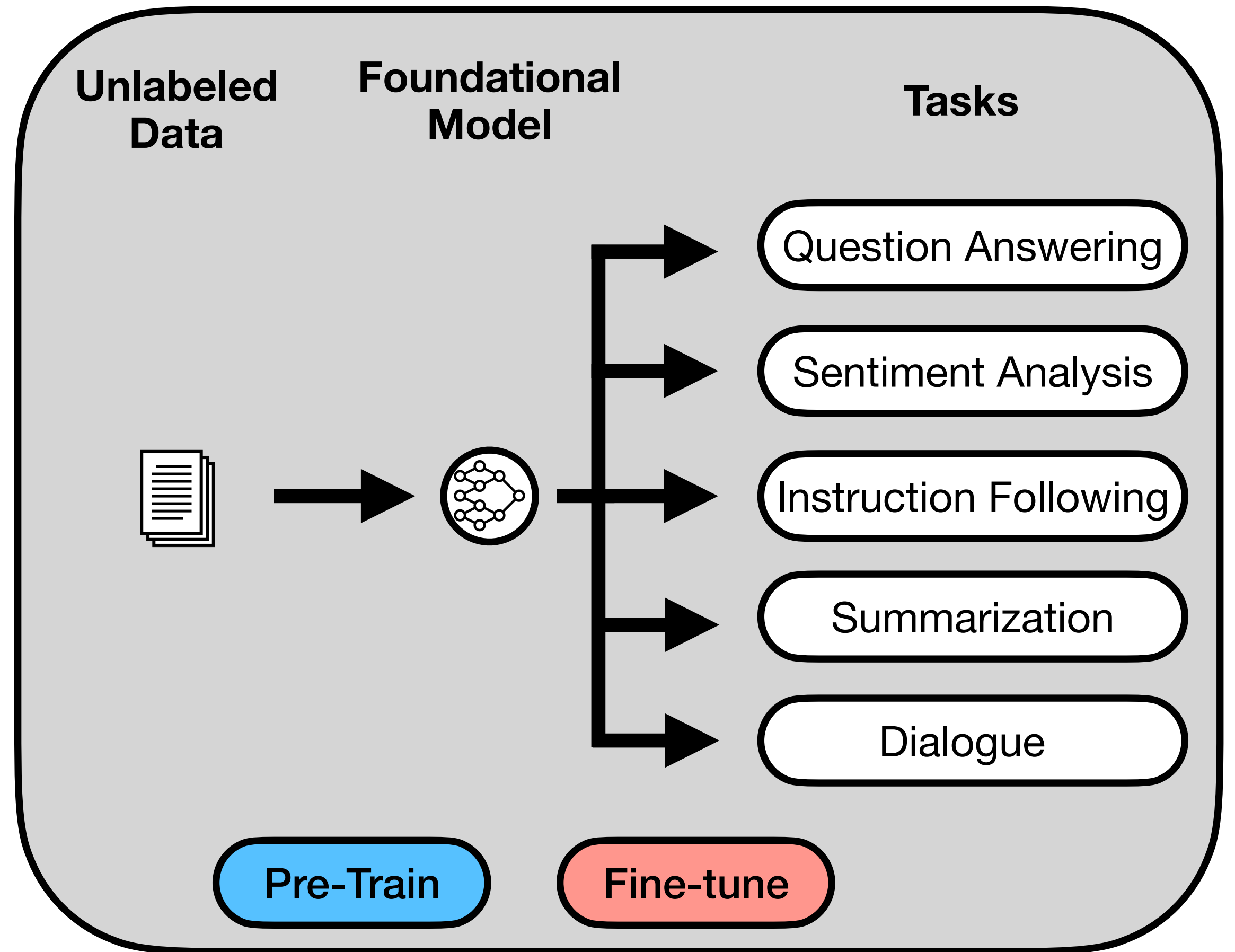


# Foundational Models

# Foundational Models



Traditional Machine Learning Models



Foundational Models


# Foundational Models

> 70B params

Large Language Model

$$p(y_1, y_2, \dots, y_n) = p(y_1)p(y_2 | y_1)p(y_3 | y_1, y_2) \dots = \prod_{i=1}^n p(y_i | y_1, \dots, y_{i-1})$$

for example, context

  
 $p(\cdot | \text{you, only, live})$

  
predict next word



# Foundational Models

## Large Language Model

$$p(y_1, y_2, \dots, y_n) = p(y_1)p(y_2 | y_1)p(y_3 | y_1, y_2) \dots = \prod_{i=1}^n p(y_i | y_1, \dots, y_{i-1})$$

Pre-Training



Large Unlabeled Dataset

# Foundational Models

## Large Language Model

$$L(y) = - \sum_i^n \log p(y_i | y_{<i}) \quad \left. \vphantom{\sum} \right\} \text{next word prediction}$$

Pre-Training



Large Unlabeled Dataset

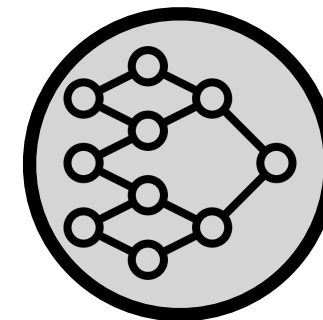
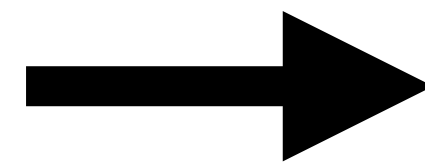
# Foundational Models

## Large Language Model

$$L(x, y) = - \sum_i^n \log p(y_i | y_{<i}, x) \left. \vphantom{\sum} \right\} \begin{array}{l} \text{predicting the next token} \\ \text{conditioned on some} \\ \text{context} \end{array}$$

Fine-Tuning

Input: x



Output: y

What is the capital of France?

Describe the function of a computer motherboard

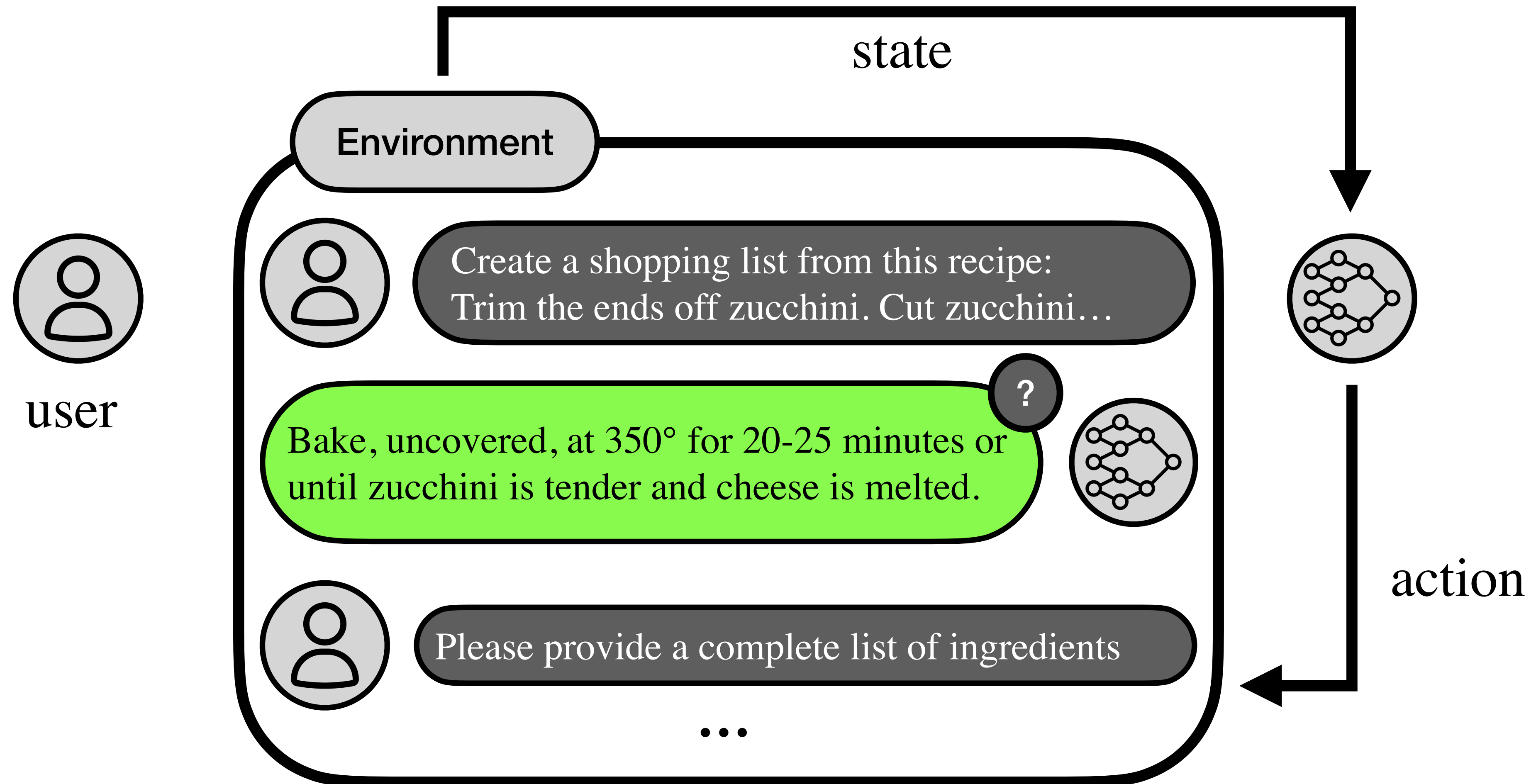
....

The capital of France is Paris.

A computer motherboard is the main circuit board in a computer. It is the backbone of ...

....

# Foundational Models

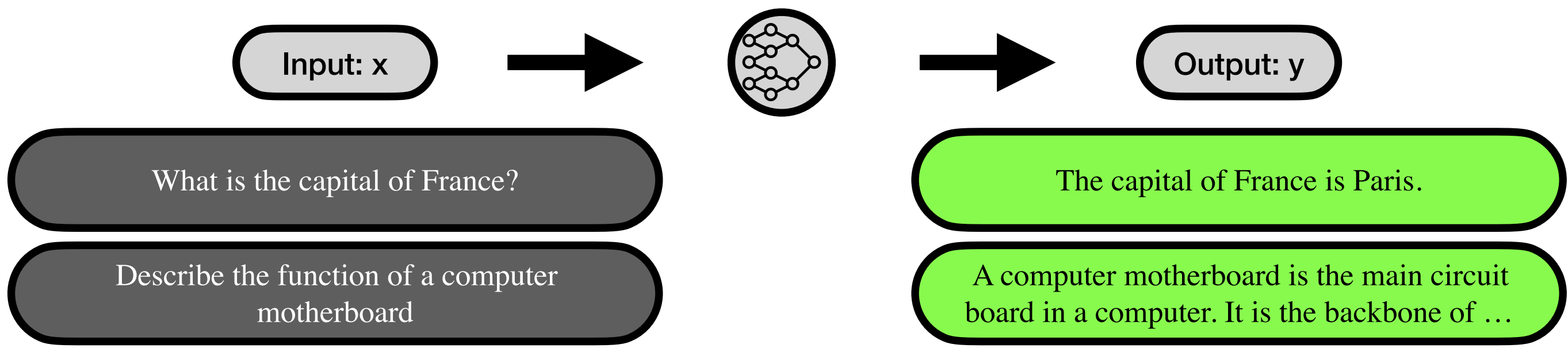


# Foundational Models

objective alignment issue

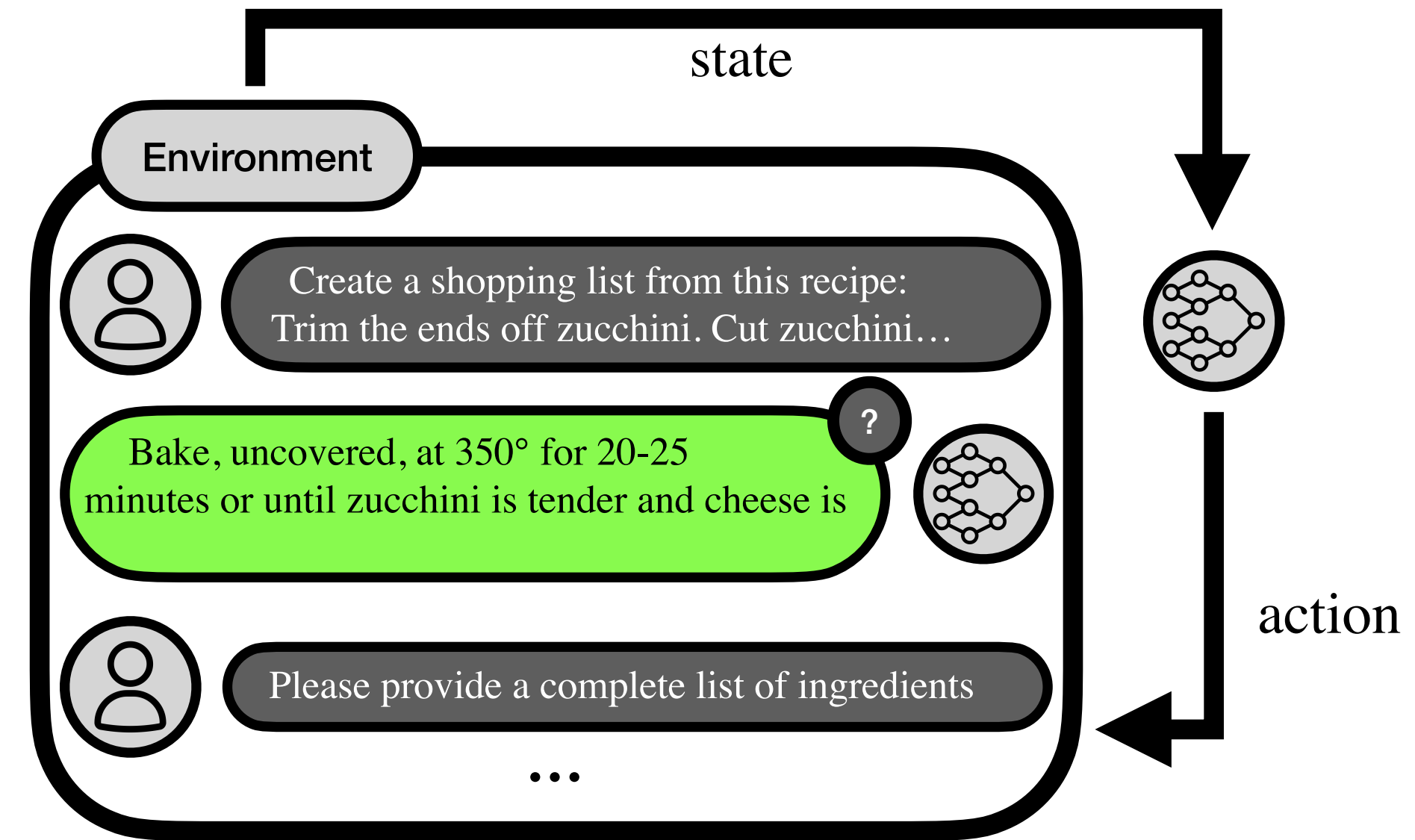
Fine-Tuning

$$L(x, y) = - \sum_i^n \log p(y_i | y_{<i}, x)$$



Next word prediction

Test



Ability to follow instructions

# Foundational Models

“Making language models bigger does not inherently make them better at following a user’s intent. For example, large language models can generate outputs that are untruthful, toxic, or simply not helpful to the user. In other words, these models are not aligned with their users.”

Long Ouyang et al.  
Training language models to follow  
instructions with human feedback  
**OpenAI 2022**



# An Old Problem



Figure 3: NAVLAB, the CMU autonomous navigation test vehicle.

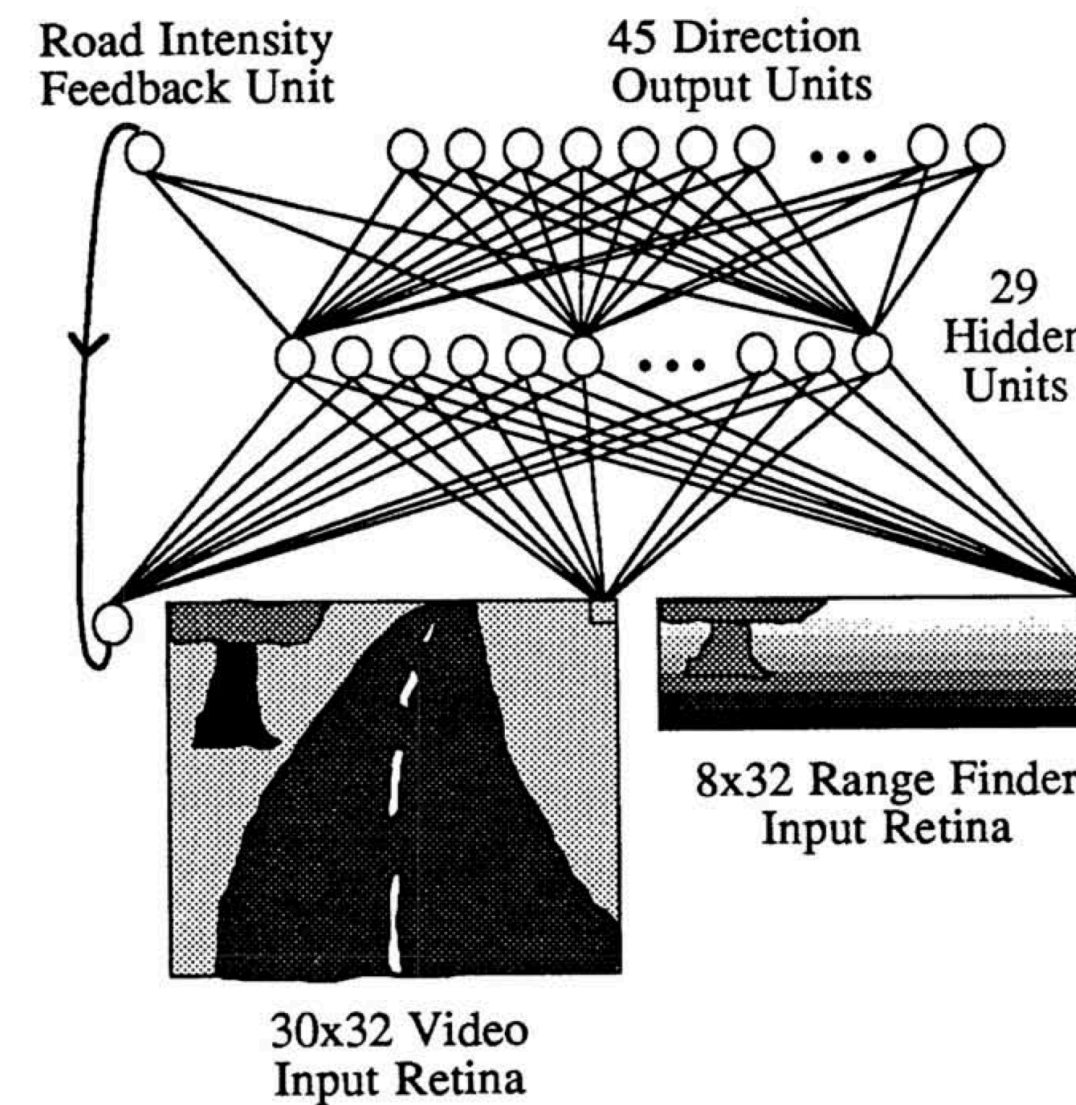


Figure 1: ALVINN Architecture

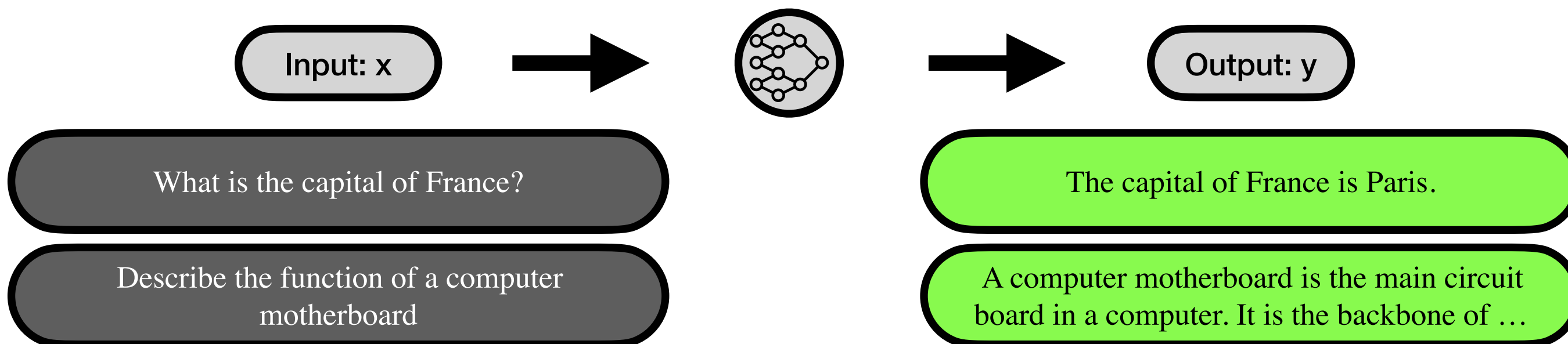
“... the network must not solely be shown examples of accurate driving, **but also how to recover** (i.e. return to the road center) once a mistake has been made.”

Dean Pomerleau  
Alvinn: An Autonomous Land  
Vehicle In A Neural Network  
**NeurIPS 1989**

# Foundational Models

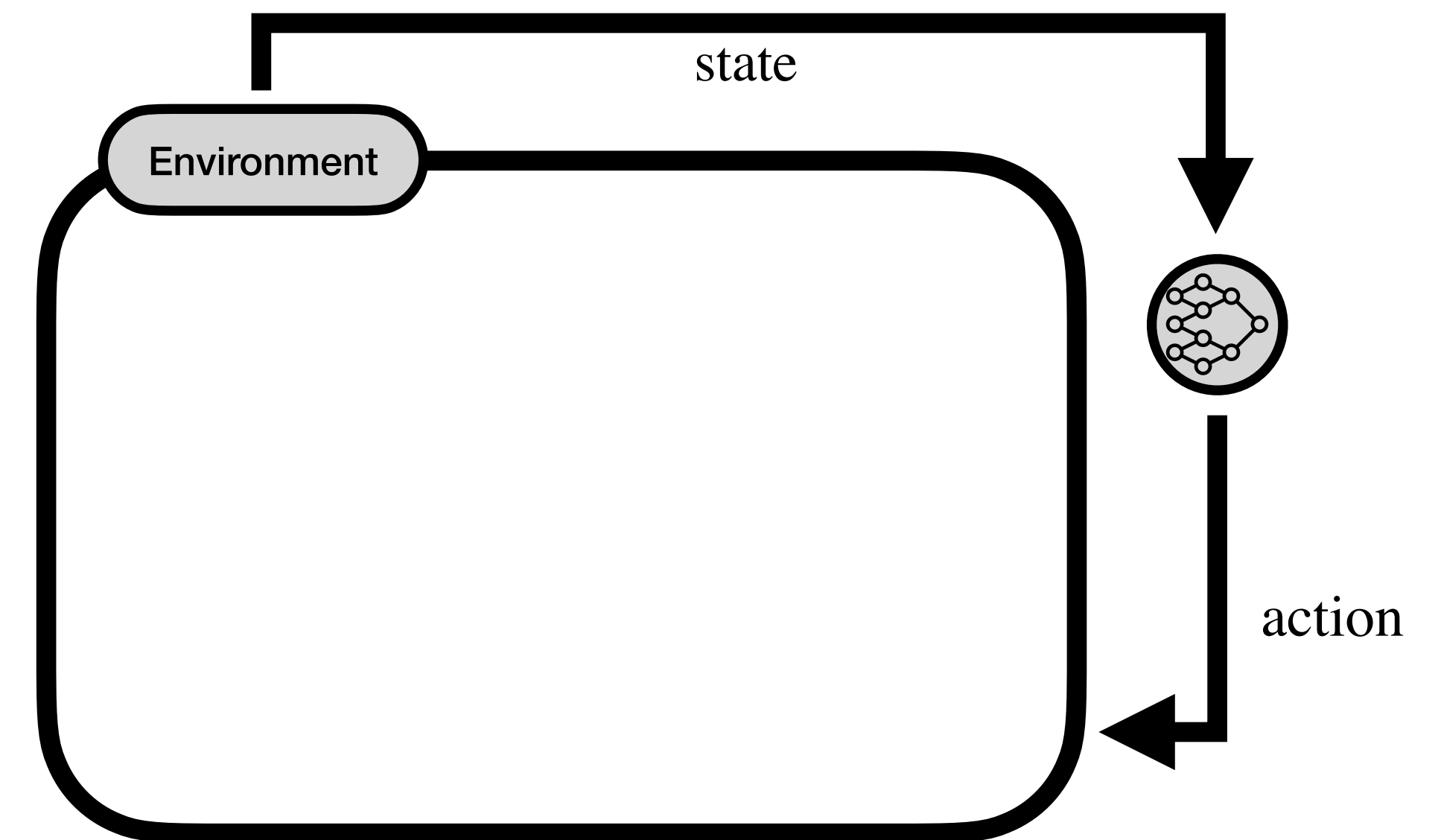
## Fine-Tuning

$$L(x, y) = - \sum_i^n \log p(y_i | y_{<i}, x)$$



## Next word prediction

## Interactive Fine-Tuning





# Foundational Models

Fine-Tuning

Interactive Fine-Tuning

$$L(x, y) = -$$

Input: x

What is the capital of France?

Describe the function of a computer motherboard

## Research Goal:

Develop algorithms that can learn to make automated decisions in the real world from interactions with minimal supervision.

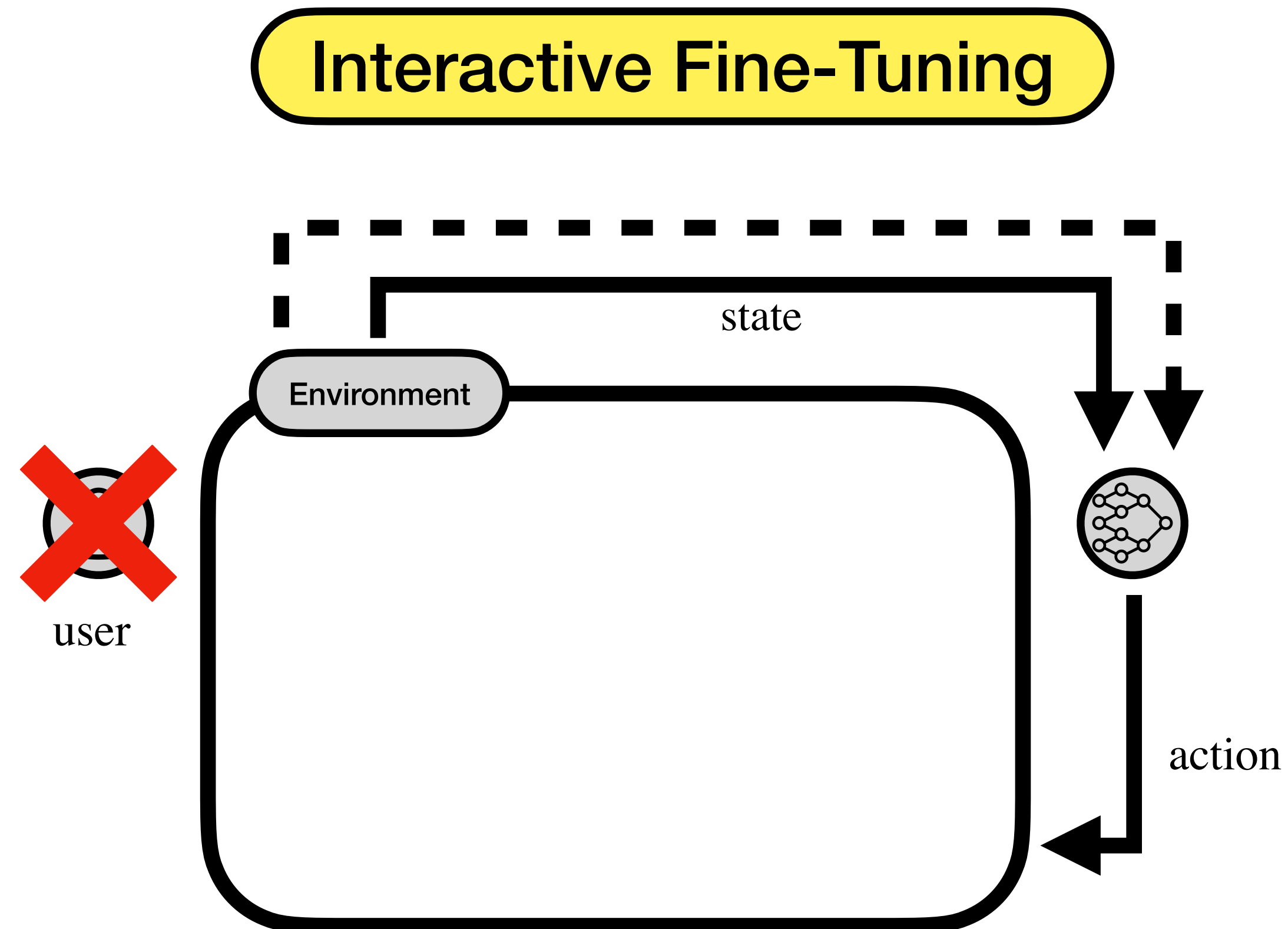
board in a computer. It is the backbone of ...

Next word prediction

state

action

# Foundational Models



How do we learn from interactions in an environment?

# Overview of Past Work

- Imitation Learning [1, 2, 12]
- Imitation Learning with a Computational Oracle [6, 7, 8]
- Preference-Based Reinforcement Learning [3, 4, 5]
- Constrained Reinforcement Learning [9, 10, 11]

1. Disagreement-Regularized Imitation Learning.....	[BSH ICLR 2020]
2. Adversarial Imitation Learning via Boosting.....	[CSHBS ICLR 2024]
3. Is reinforcement learning (not) for natural language processing?.....	[RABHSBHC ICLR 2023]
4. Learning to Generate Better Than Your LLM.....	[CBRMS Instruction Workshop 2023]
5. Policy-Gradient Training of Language Models for Ranking.....	[GCCBJ FMDM Workshop 2023]
6. Non-Monotonic Text Generation.....	[WBDC ICLR 2019]
7. Active Imitation Learning with Noisy Guidance.....	[BSD ACL 2020]
8. Interactive Text Generation.....	[FGPBCZGD EMNLP 2024]
9. Reinforcement Learning with Convex Constraints.....	[MBDDS NeurIPS 2019]
10. Constrained episodic reinforcement learning in concave-convex and knapsack settings.....	[BDLMSSS NeurIPS 2020]
11. Ranking with Long-Term Constraints.....	[BFDJ WSDM 2024]
12. Successor feature sets: Generalizing successor representations across policies.....	[BMG AAI 2021]

# Overview of Past Work **Domains**

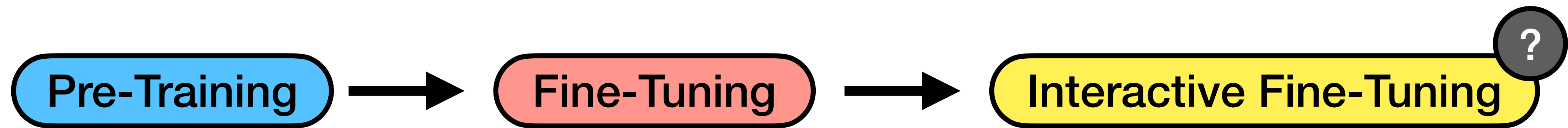
Natural Language Processing [3, 4, 6, 7, 8] •

Classic Control [2, 9, 10, 12] •

Information Retrieval [5, 11] •

Video Games [1] •

1. Disagreement-Regularized Imitation Learning.....	[ <b>BSH</b> ICLR 2020]
2. Adversarial Imitation Learning via Boosting.....	[ <b>CSHBS</b> ICLR 2024]
3. Is reinforcement learning (not) for natural language processing?.....	[ <b>RABHSBHC</b> ICLR 2023]
4. Learning to Generate Better Than Your LLM.....	[ <b>CBRMS</b> Instruction Workshop 2023]
5. Policy-Gradient Training of Language Models for Ranking.....	[ <b>GCCBJ</b> FMDM Workshop 2023]
6. Non-Monotonic Text Generation.....	[ <b>WBDC</b> ICLR 2019]
7. Active Imitation Learning with Noisy Guidance.....	[ <b>BSD</b> ACL 2020]
8. Interactive Text Generation.....	[ <b>FGPBCZGD</b> EMNLP 2024]
9. Reinforcement Learning with Convex Constraints.....	[ <b>MBDDS</b> NeurIPS 2019]
10. Constrained episodic reinforcement learning in concave-convex and knapsack settings.....	[ <b>BDLMSSS</b> NeurIPS 2020]
11. Ranking with Long-Term Constraints.....	[ <b>BFDJ</b> WSDM 2024]
12. Successor feature sets: Generalizing successor representations across policies.....	[ <b>BMG</b> AAI 2021]





Formalize the imitation learning problem

objective alignment issue

How can interaction help?

How can we design interactive learning algorithms specific to LLMs?



Formalize the imitation learning problem

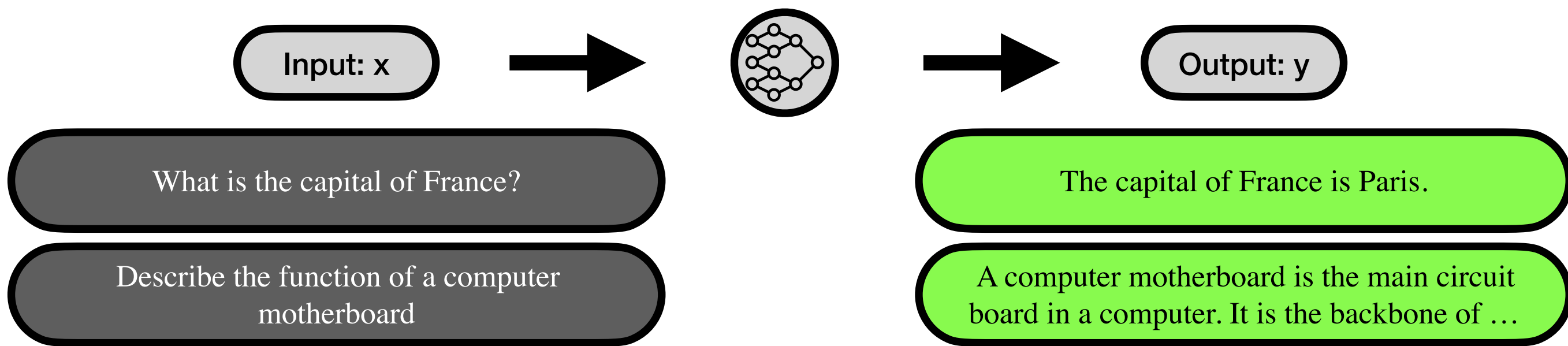
objective alignment issue

# Foundational Models

objective alignment issue

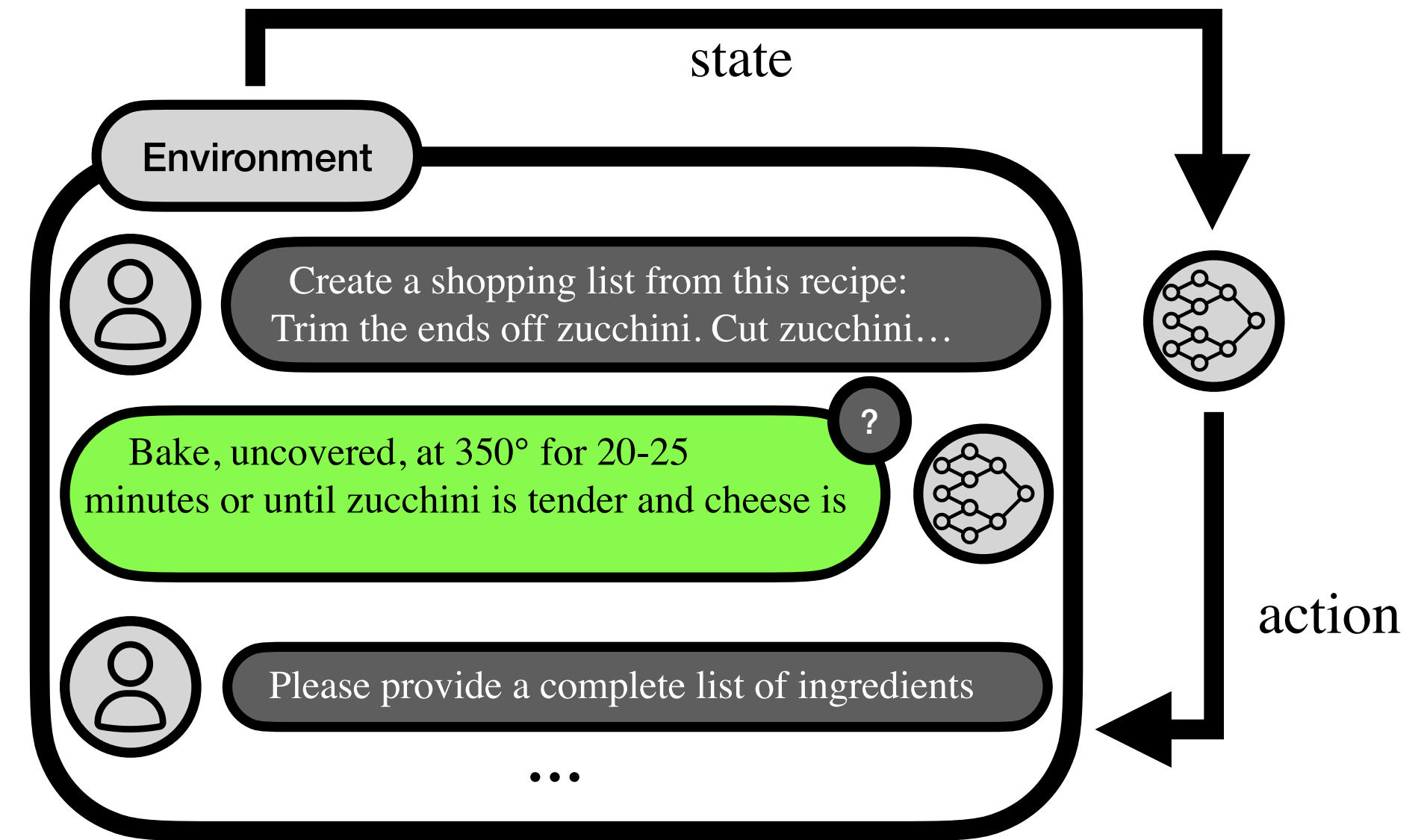
Fine-Tuning

$$L(x, y) = - \sum_i^n \log p(y_i | y_{<i}, x)$$



Next word prediction

Test



Ability to follow instructions



# Reinforcement Learning

## basics of mdp

**States:**

$$\{x, y\}$$

**Actions:**



**Reward:**

$$R(s, a)$$

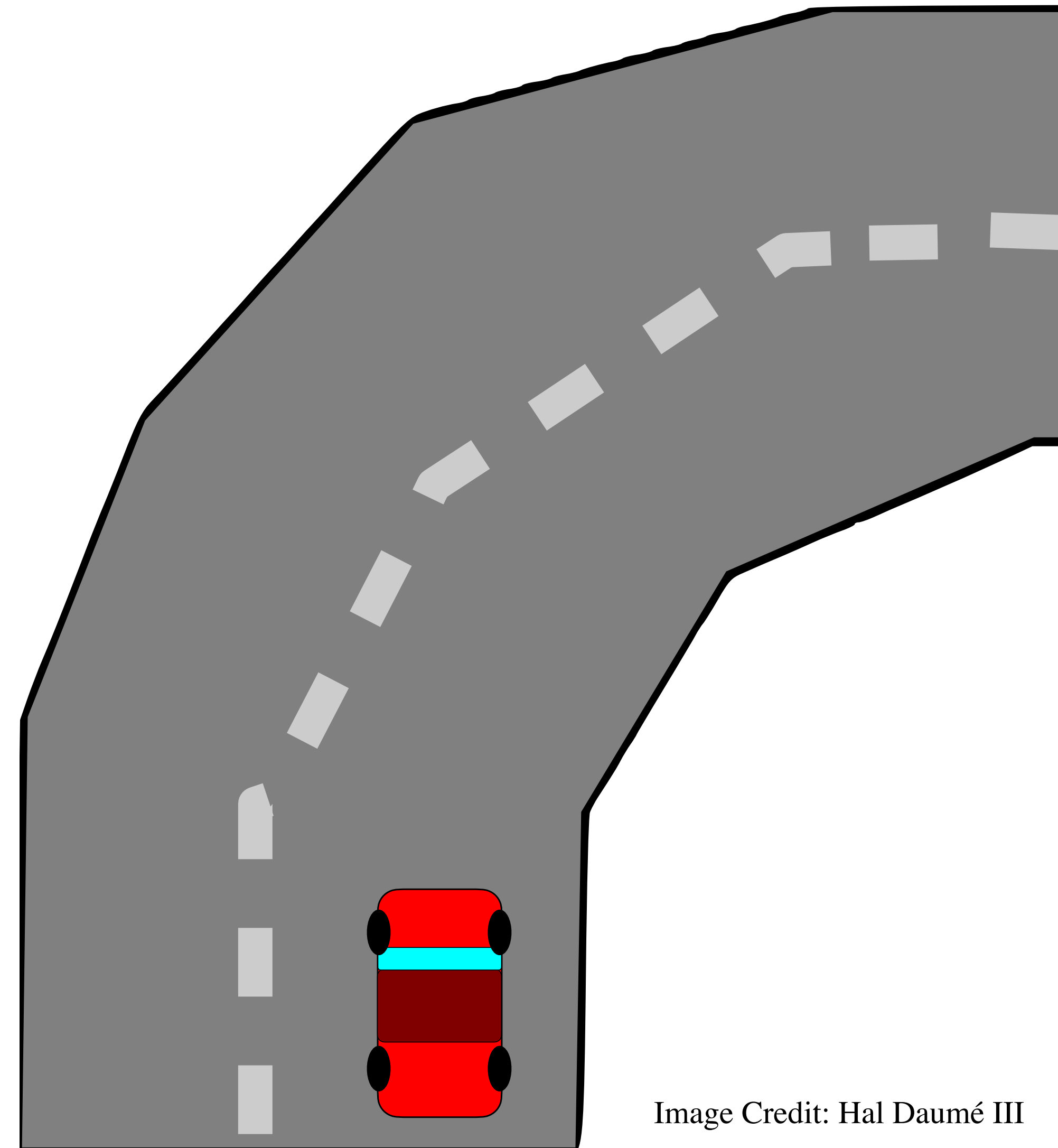
**Transition:**

$$P(s' | s, a)$$

**Policy:**

$$\pi(\cdot | s)$$

**Goal:** Learn policy  $\pi$  to maximize reward



# Imitation Learning

## basics of mdp

States:

$$\{x, y\}$$

Actions:



Reward:

Unknown Reward

Transition:

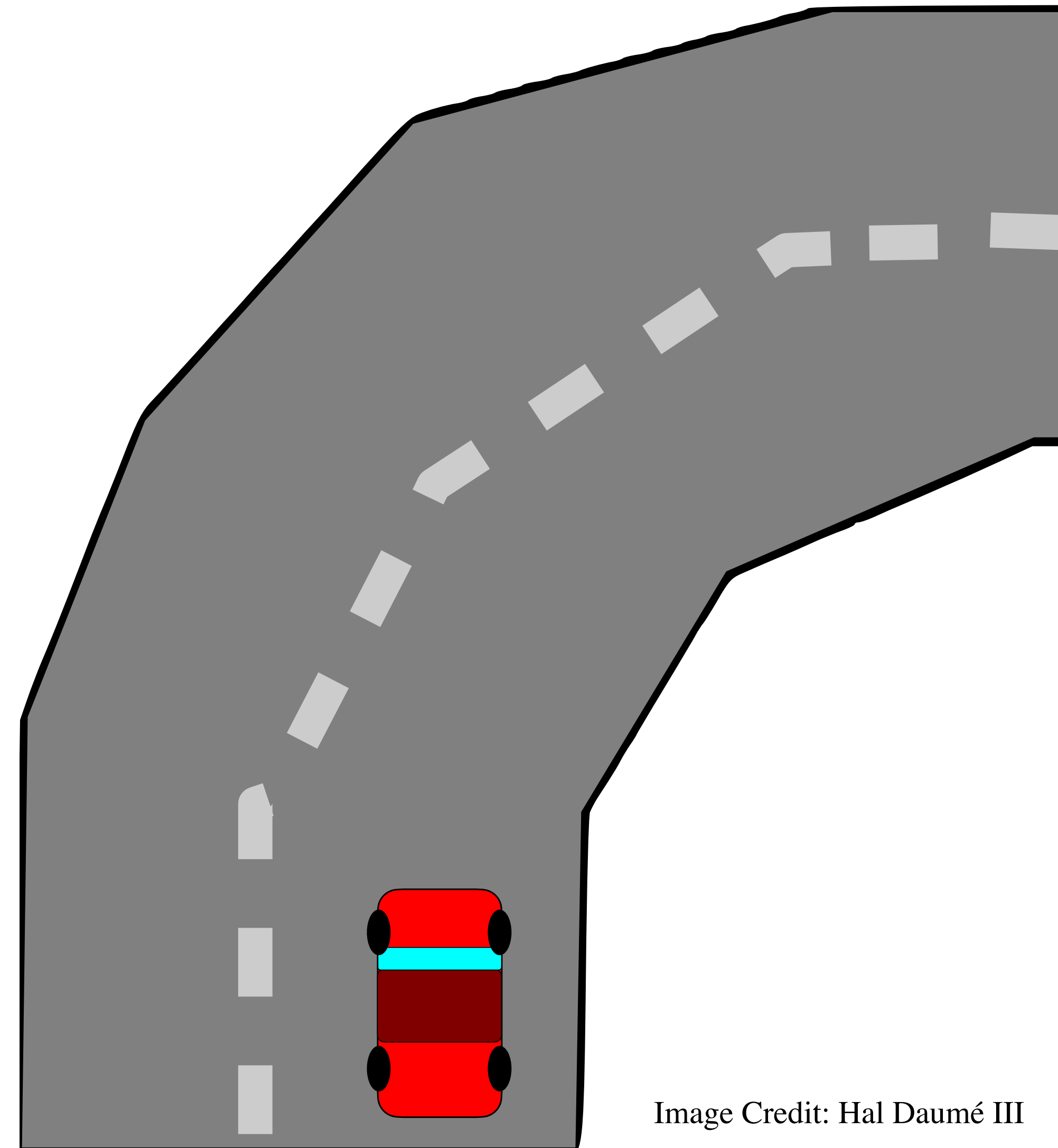
$$P(s' | s, a)$$

Policy:

$$\pi(\cdot | s)$$

**Training dataset:**  $D = \{(state, action)\}$  from an expert  $\pi^*$

**Goal:** Train a policy  $\pi$  to mimic the demonstrations

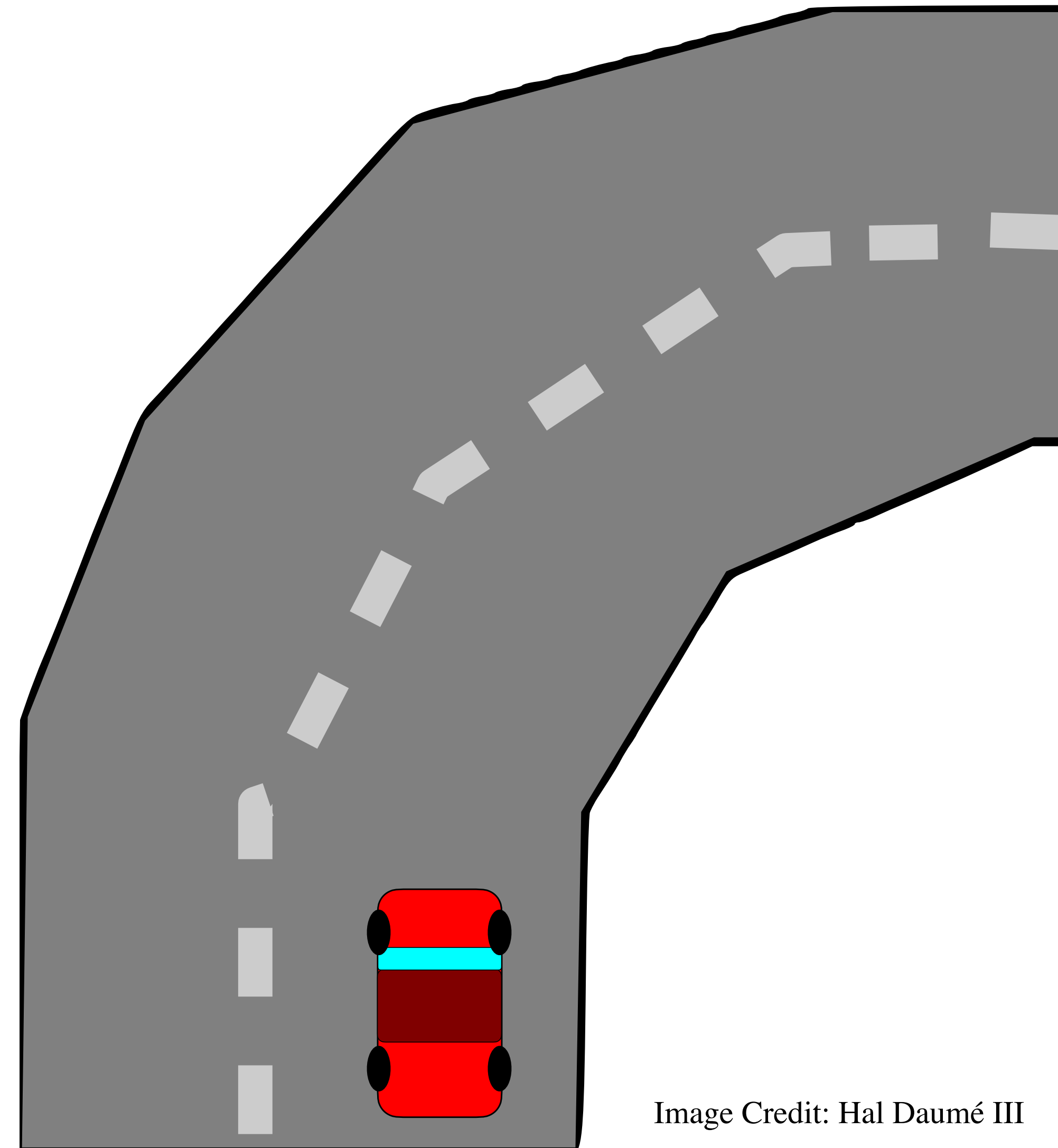


# Behavior Cloning

simplest imitation learning

1. Collect trajectories from expert  $\pi^\star$
2. Create a dataset  $(s_1, a_1), \dots, (s_h, a_h) \sim \rho_{\pi^\star}$
3. Train a policy (classifier)  $\pi$

$$\min_{\pi \in \Pi} \mathbb{E}_{(s,a) \sim \rho_{\pi^\star}} [L(\pi, s, a)]$$

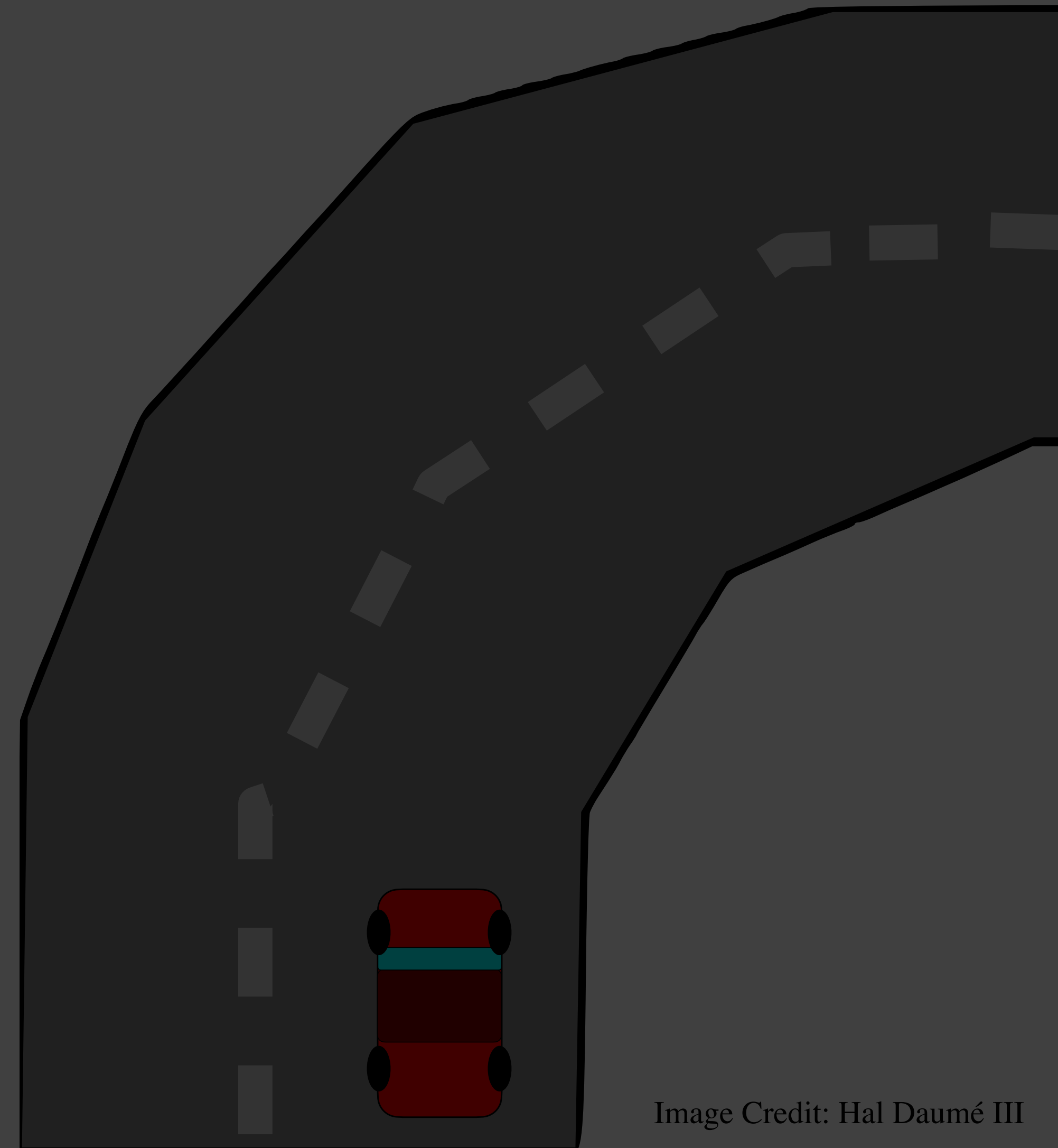


# Behavior Cloning

simplest imitation learning

1. Collect trajectories from expert  $\pi^*$
2. Create a dataset  $(s_1, a_1), \dots, (s_h, a_h) \sim \rho_{\pi^*}$
3. Train a policy (classifier)  $\pi$

$$\min_{\pi \in \Pi} \mathbb{E}_{(s,a) \sim \rho_{\pi^*}} [L(\pi, s, a)]$$

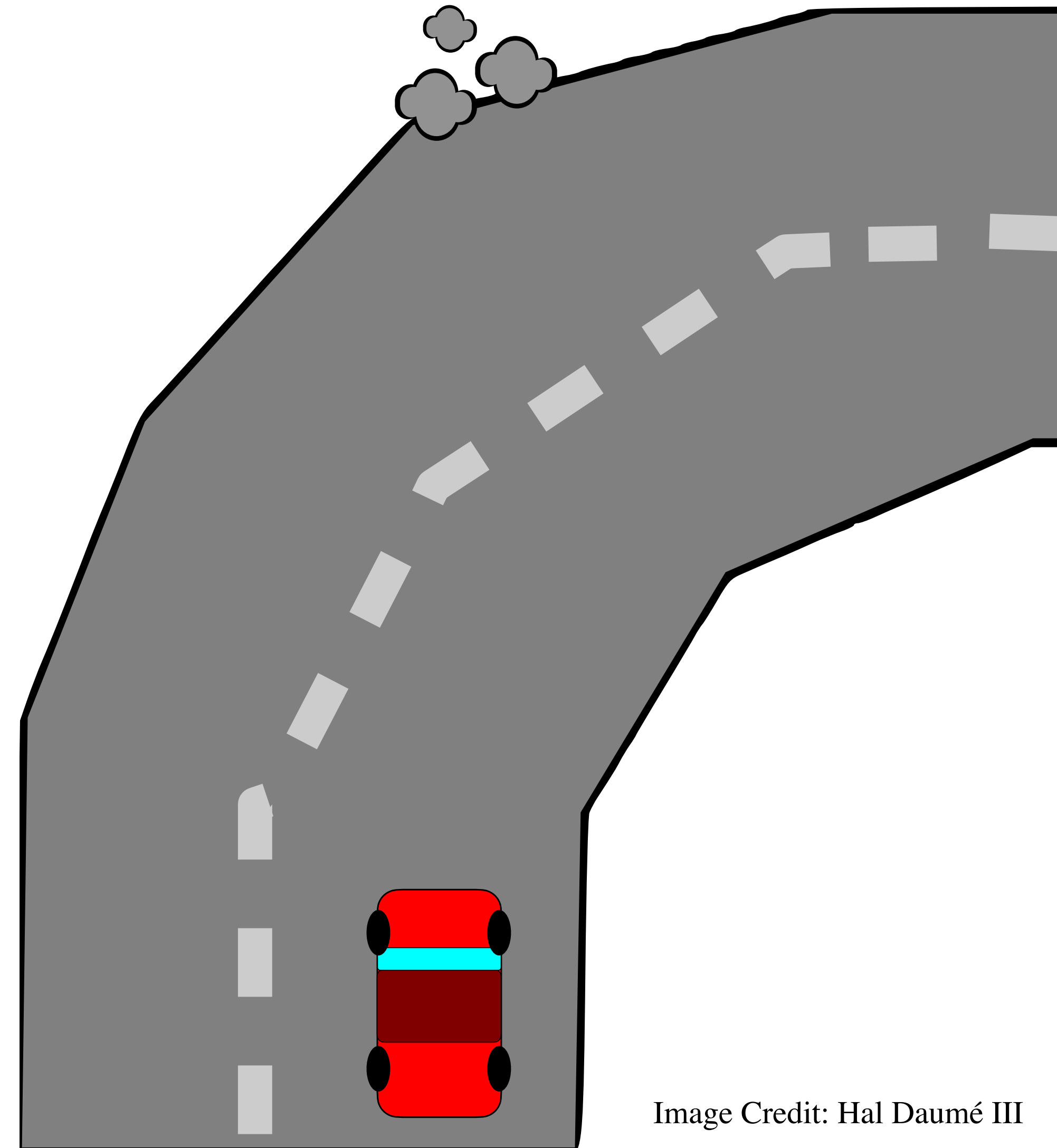


# Behavior Cloning

simplest imitation learning

1. Collect trajectories from expert  $\pi^*$
2. Create a dataset  $(s_1, a_1), \dots, (s_h, a_h) \sim \rho_{\pi^*}$
3. Train a policy (classifier)  $\pi$

$$\min_{\pi \in \Pi} \mathbb{E}_{(s,a) \sim \rho_{\pi^*}} [L(\pi, s, a)] \leq \epsilon$$

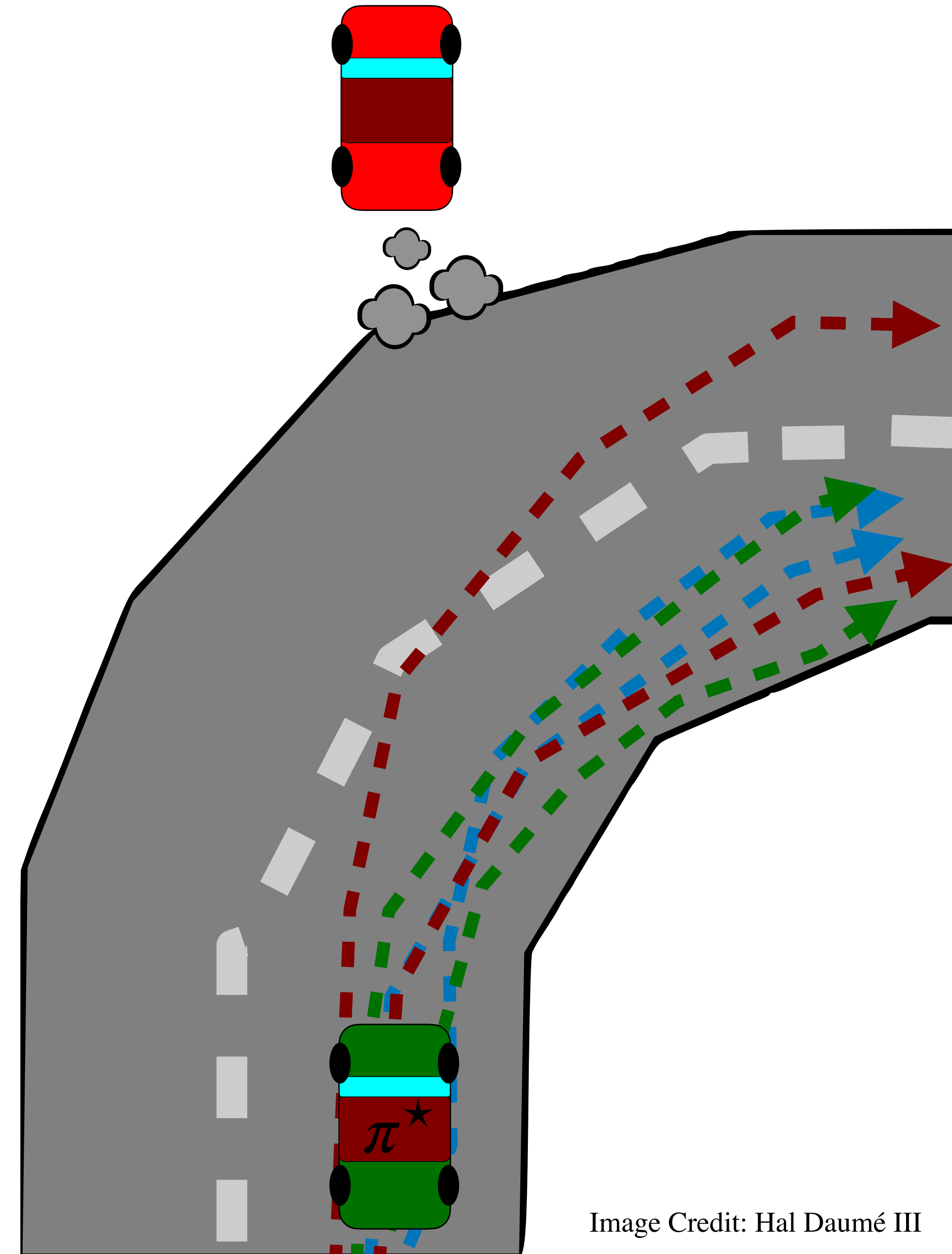


# Behavior Cloning

simplest imitation learning

1. Collect trajectories from expert  $\pi^*$
2. Create a dataset  $(s_1, a_1), \dots, (s_h, a_h) \sim \rho_{\pi^*}$
3. Train a policy (classifier)  $\pi$

$$\min_{\pi \in \Pi} \mathbb{E}_{(s,a) \sim \rho_{\pi^*}} [L(\pi, s, a)] \leq \epsilon$$



# Behavior Cloning

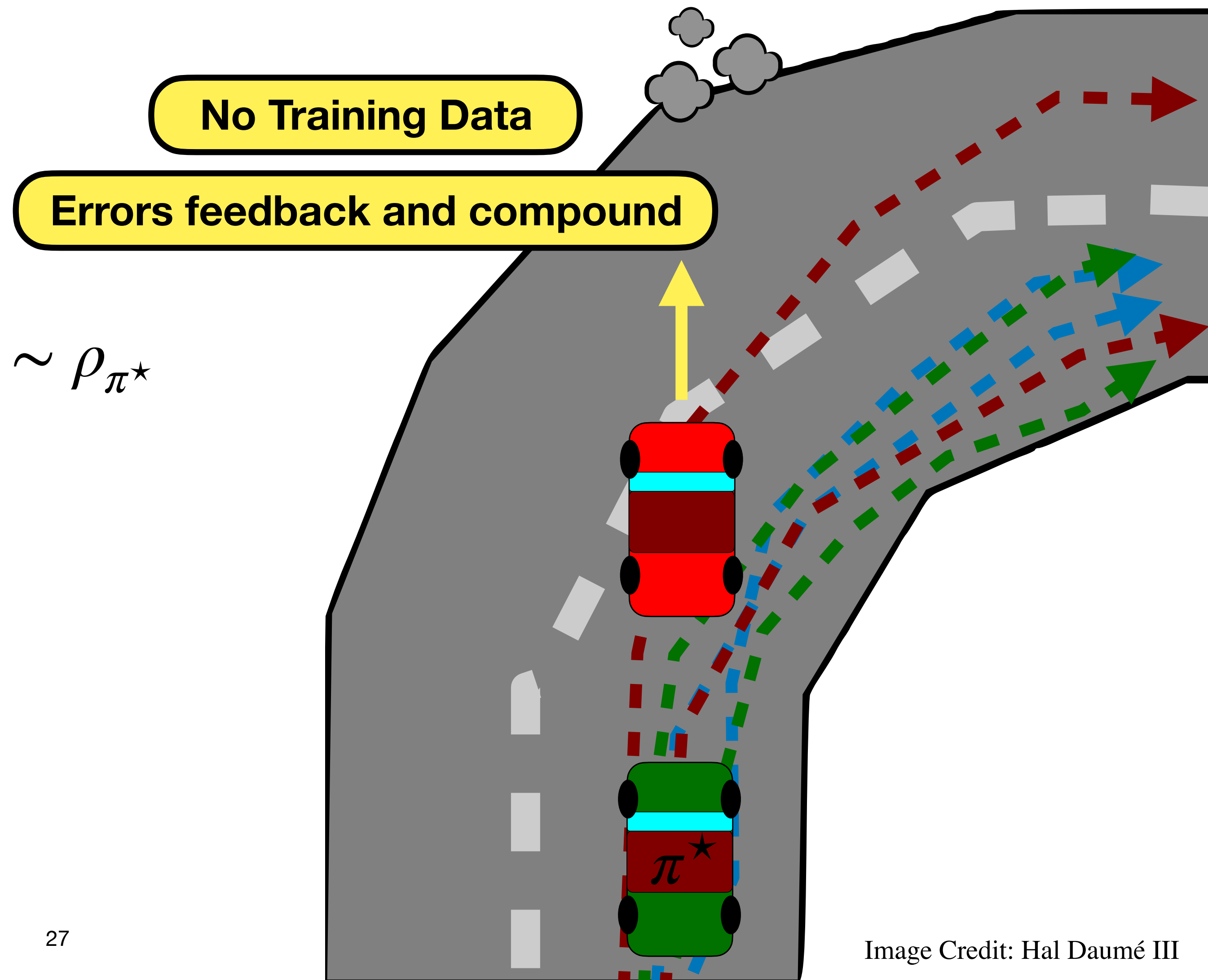
simplest imitation learning

1. Collect trajectories from expert  $\pi^*$

2. Create a dataset  $(s_1, a_1), \dots, (s_h, a_h) \sim \rho_{\pi^*}$

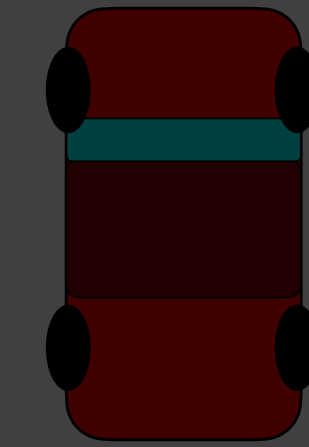
3. Train a policy (classifier)  $\pi$

$$\min_{\pi \in \Pi} \mathbb{E}_{(s,a) \sim \rho_{\pi^*}} [L(\pi, s, a)] \leq \epsilon$$



# Behavior Cloning

simplest imitation learning



1. Collect trajectories

2. Create a dataset

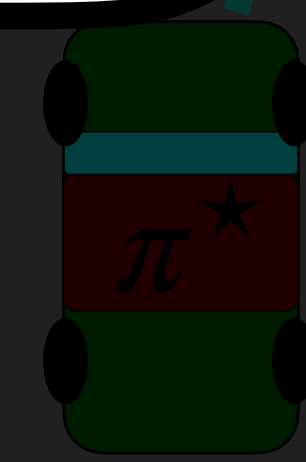
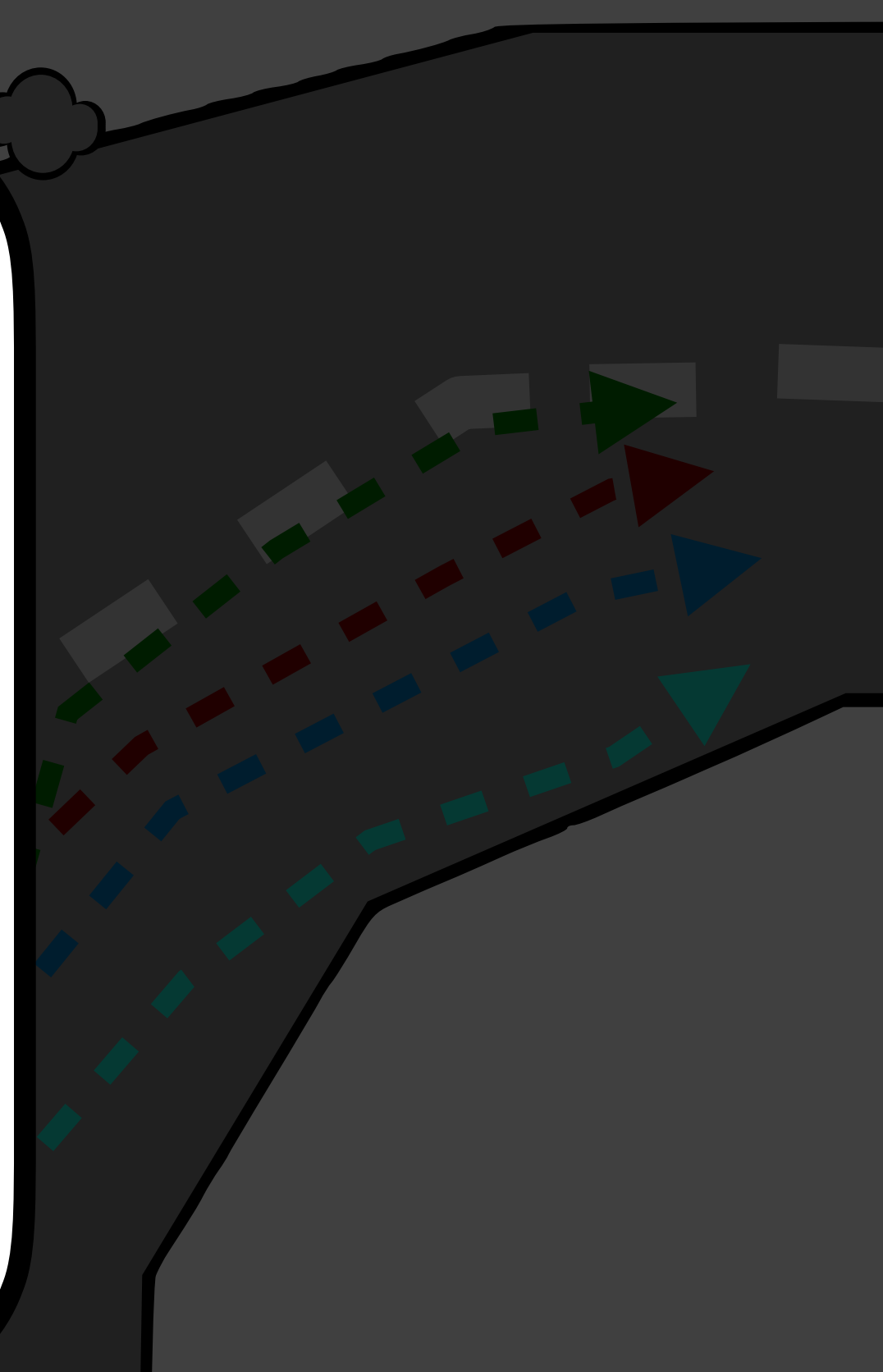
3. Train a policy

$$\min_{\pi \in \Pi} \mathbb{E}_{(s,a) \sim \rho_{\pi^*}}$$

## Problem:

The training distribution is different than the test distribution resulting in covariate shift issues.

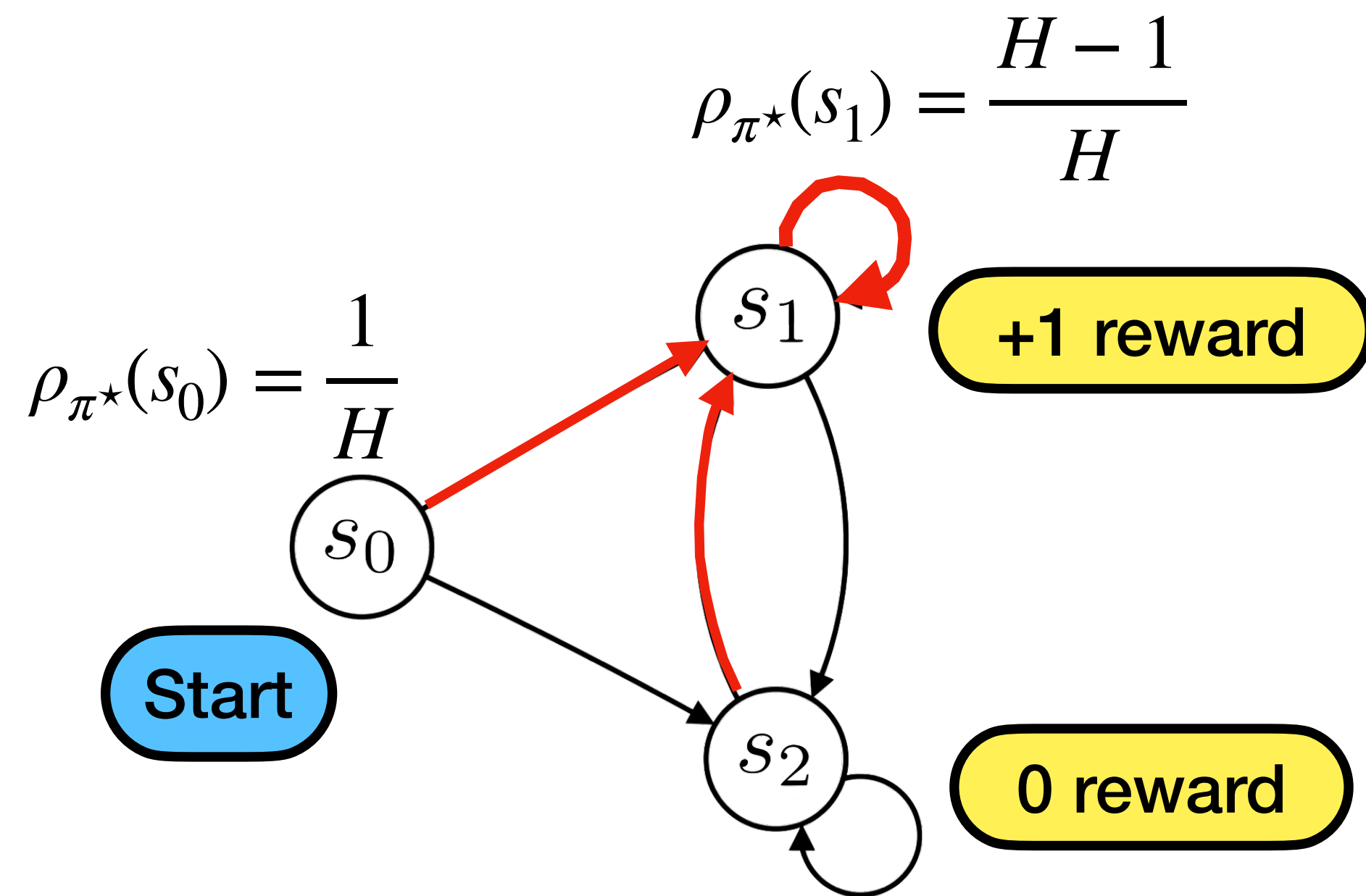
	Supervised Learning	Behavior Cloning
Train	$(x, y) \sim \rho$	$(s, a) \sim \rho_{\pi^*}$
Test	$(x, y) \sim \rho$	$(s, a) \sim \rho_{\pi}$





# Formalizing the covariate shift issue

Given an expert policy:  $\pi^*$



**Behavior Cloning Loss:**

$$\text{Loss}(\hat{\pi}) = \epsilon$$

**(loss is small)**

**Performance**

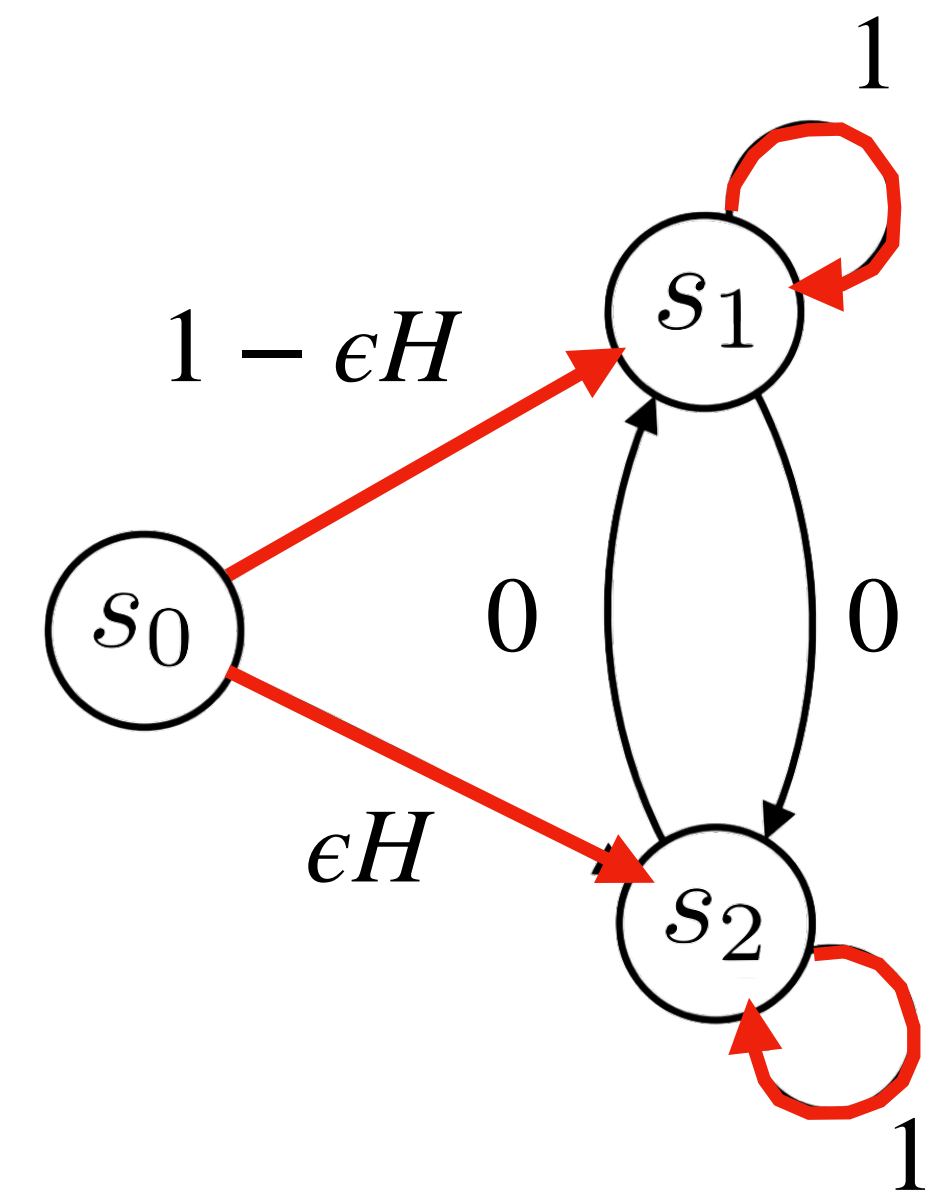
$$J(\pi) := \mathbb{E}_{\pi} \left[ \sum_{h=1}^H R_h \right]$$

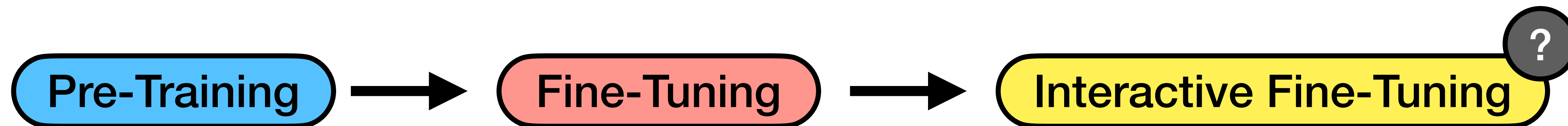
**Behavior Cloning Performance Gap:**

$$J(\pi^*) - J(\hat{\pi}) \leq \mathcal{O}(\epsilon H^2)$$

**(gap scales quadratic)**

Consider a policy:  $\hat{\pi}$





Formalize the imitation learning problem

How can interaction help?

objective alignment issue

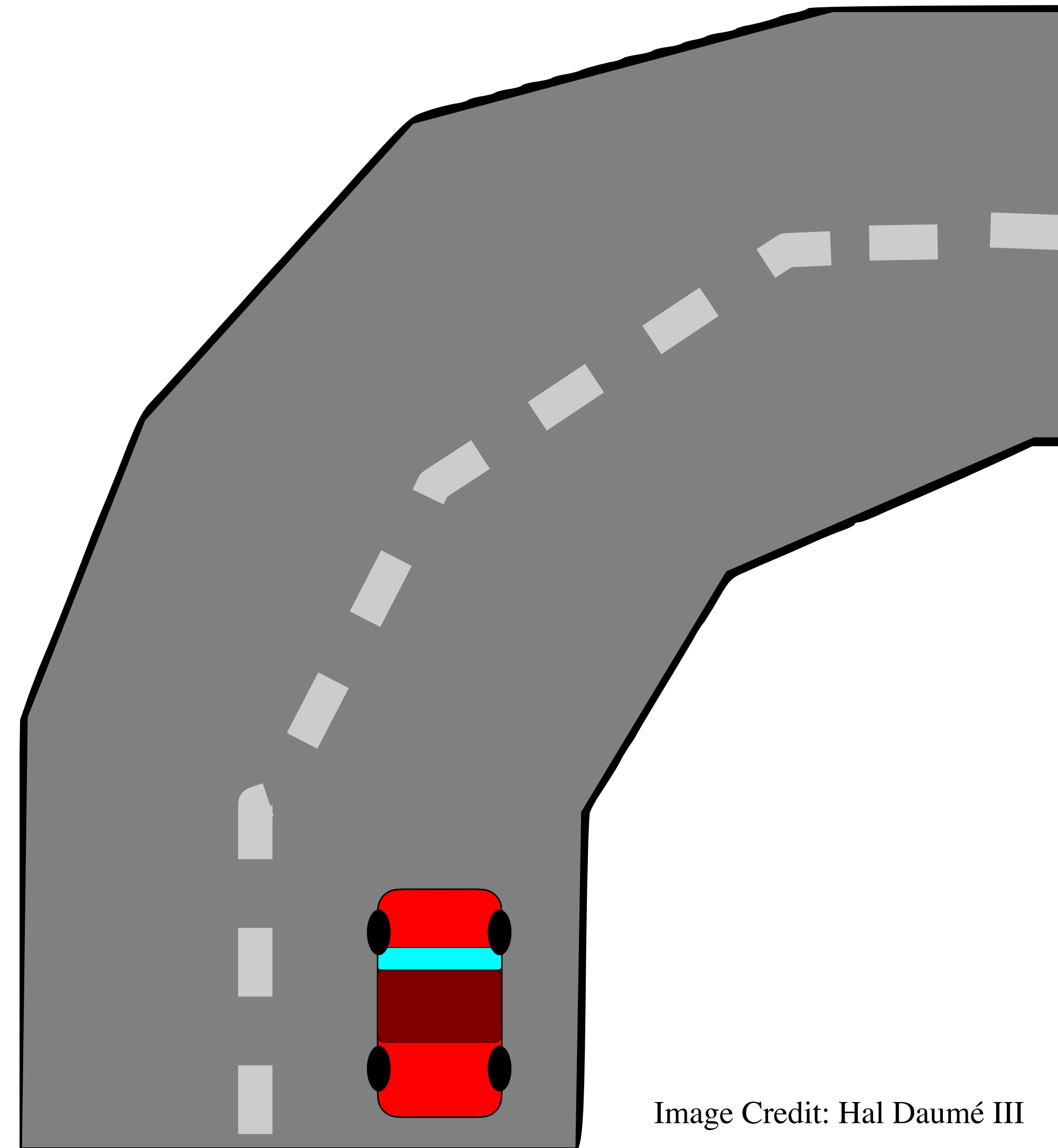
# Interactive Imitation Learning with DAgger

Let  $\pi^\star$  be an expert user 

Loop:

- Interact in an environment  $(s, a) \sim \rho_\pi$
- Get Dataset  $D_i = \{(s, \pi^\star(s))\}$
- Aggregate dataset  $D \leftarrow D \cup D_i$
- Train a policy  $\pi$  on dataset

$$\min_{\pi \in \Pi} \mathbb{E}_{(s,a) \sim \rho_\pi} [L(\pi, s, \pi^\star(s))]$$



# Interactive Imitation Learning with DAgger

Let  $\pi^*$  be an expert user 

Loop:

- Interact in an environment
- Get Dataset  $D_i = \{(s, \pi^*(s))\}$
- Aggregate dataset  $D \leftarrow D_i$
- Train a policy  $\pi$  on dataset

$$\min_{\pi \in \Pi} \mathbb{E}_{(s,a) \sim \rho_\pi} [L(\pi, s, \pi^*(s))]$$

**Behavior Cloning  
Performance Gap:**

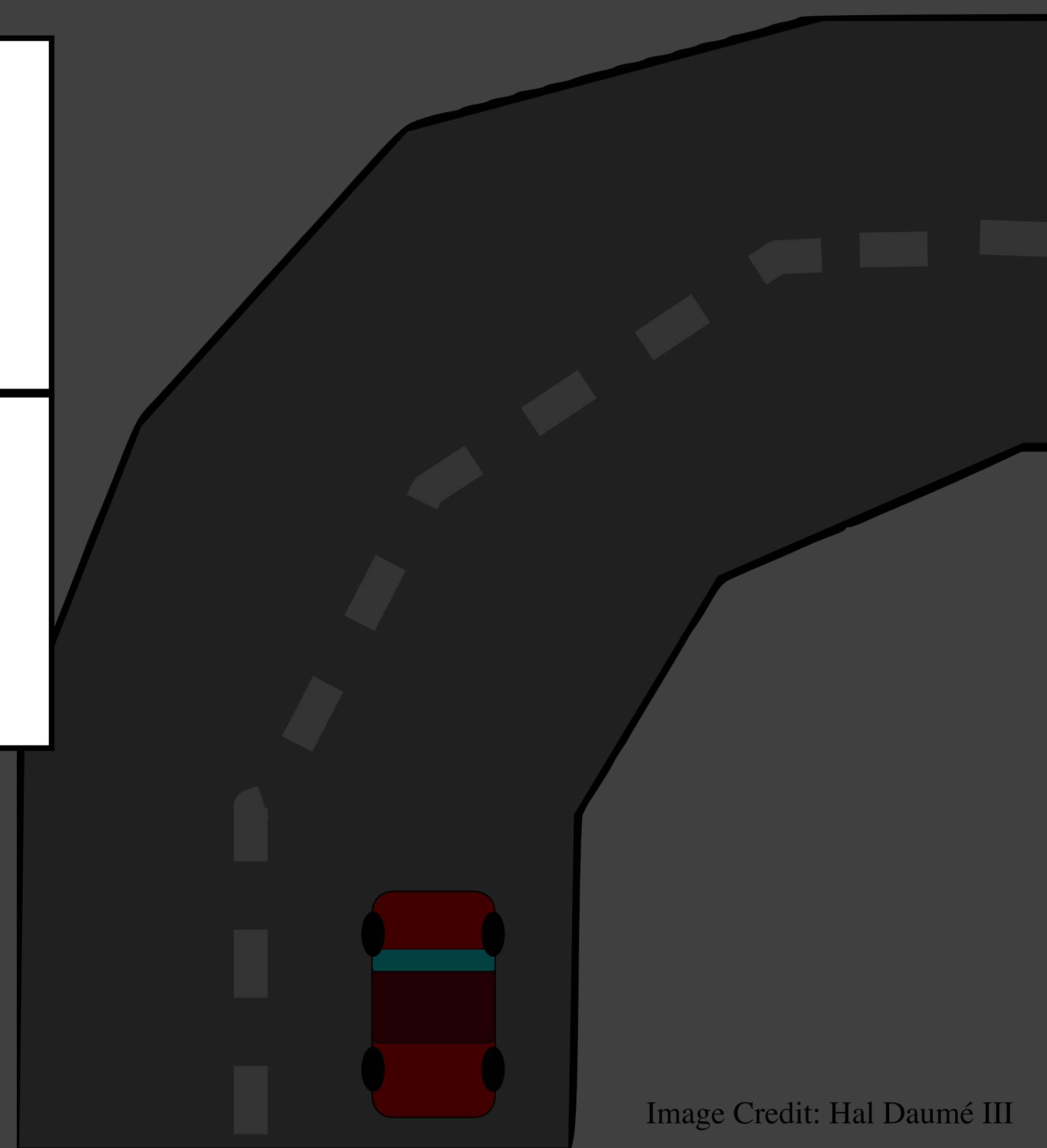
$$\mathcal{O}(\epsilon H^2)$$

**(gap scales quadratic)**

**DAgger  
Performance Gap:**

$$\mathcal{O}(\epsilon H)$$

**(gap scales linear)**



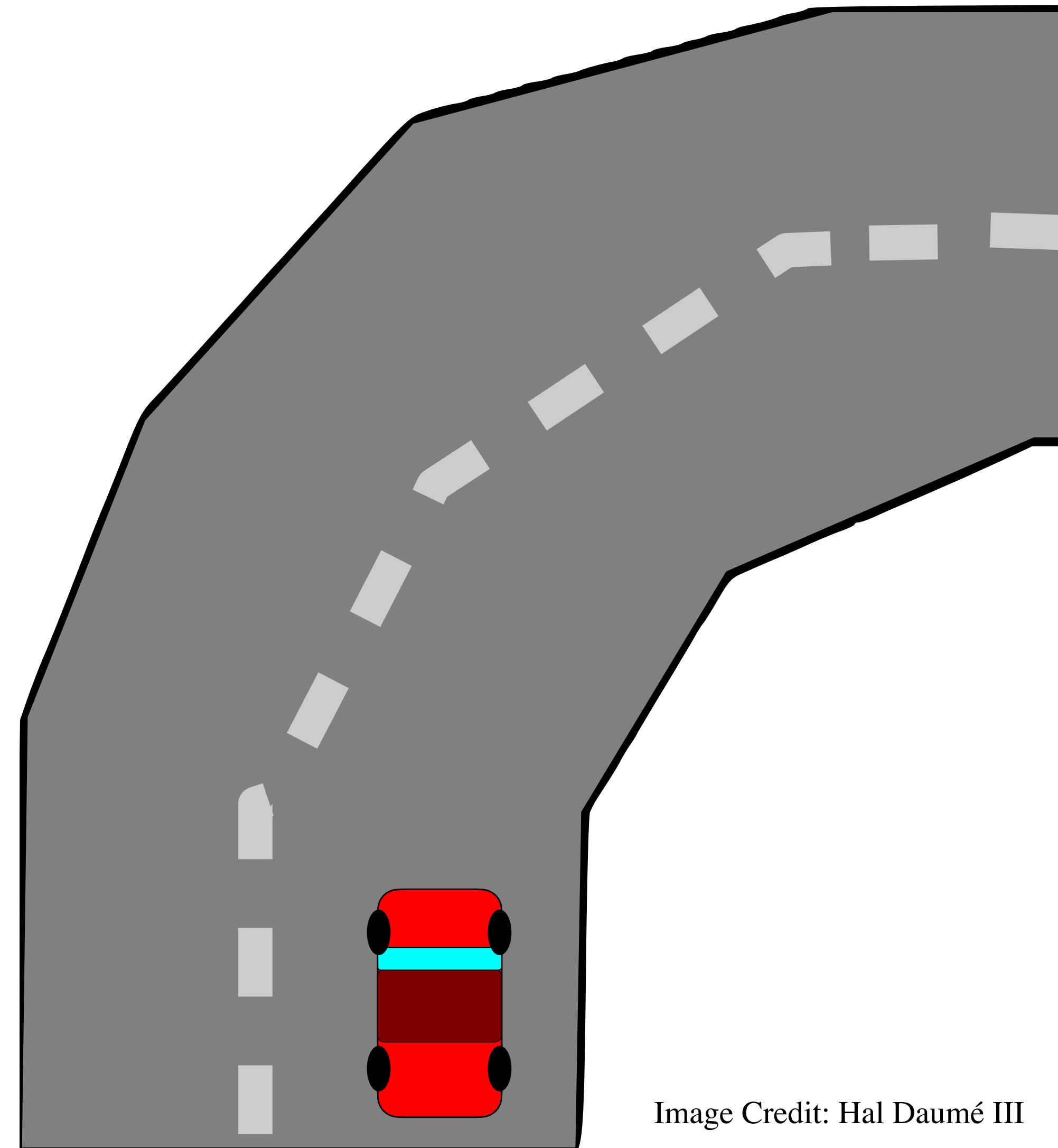
# Interactive Imitation Learning with DAgger

Let  $\pi^\star$  be an expert user 

Loop:

- Interact in an environment  $(s, a) \sim \rho_\pi$
- Get Dataset  $D_i = \{(s, \pi^\star(s))\}$
- Aggregate dataset  $D \leftarrow D \cup D_i$
- Train a policy  $\pi$  on dataset

$$\min_{\pi \in \Pi} \mathbb{E}_{(s,a) \sim \rho_\pi} [L(\pi, s, \pi^\star(s))]$$



# Interactive Imitation Learning with DAgger

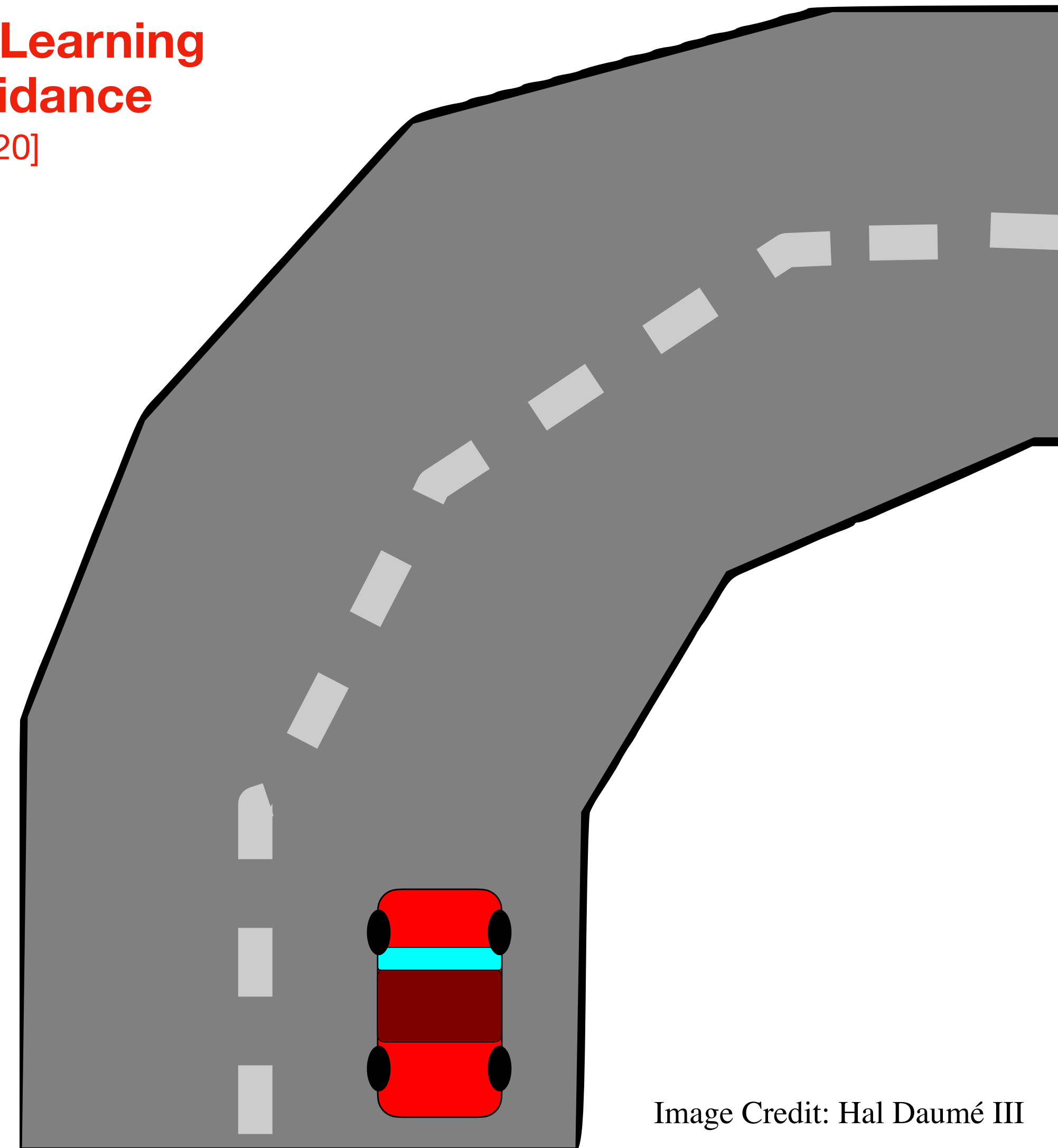
Let  $\pi^*$  be an expert user 

Active Imitation Learning  
with Noisy Guidance  
[BSD ACL 2020]

Loop:

- Interact in an environment  $(s, a) \sim \rho_\pi$
- **If  $\pi$  is not confident:**
  - **Get Dataset**  $D_i = \{(s, \pi^*(s))\}$
- **Aggregate dataset**  $D \leftarrow D \cup D_i$
- **Train a policy  $\pi$  on dataset**

$$\min_{\pi \in \Pi} \mathbb{E}_{(s,a) \sim \rho_\pi} [L(\pi, s, \pi^*(s))]$$



# Interactive Imitation Learning with DAgger

Let  $\pi^*$  be an expert

Loop:

- Interact in an environment  $(s, a) \sim \rho_\pi$

- If  $\pi$  is not optimal

- Get Data

- Aggregate

- Train a policy  $\pi$  on dataset

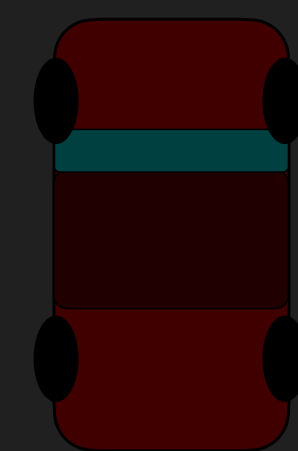
$$\min_{\pi \in \Pi} \mathbb{E}_{(s,a) \sim \rho_\pi} [L(\pi, s, \pi^*(s))]$$

## Problem:

We need online access to an expert to label each state we visit.

## Research Question:

How can we remove the need for online access to an expert?



# Our Approach: DRIL [BSH ICLR 2020]

## disagreement regularized imitation learning

Mimic expert within the expert distribution

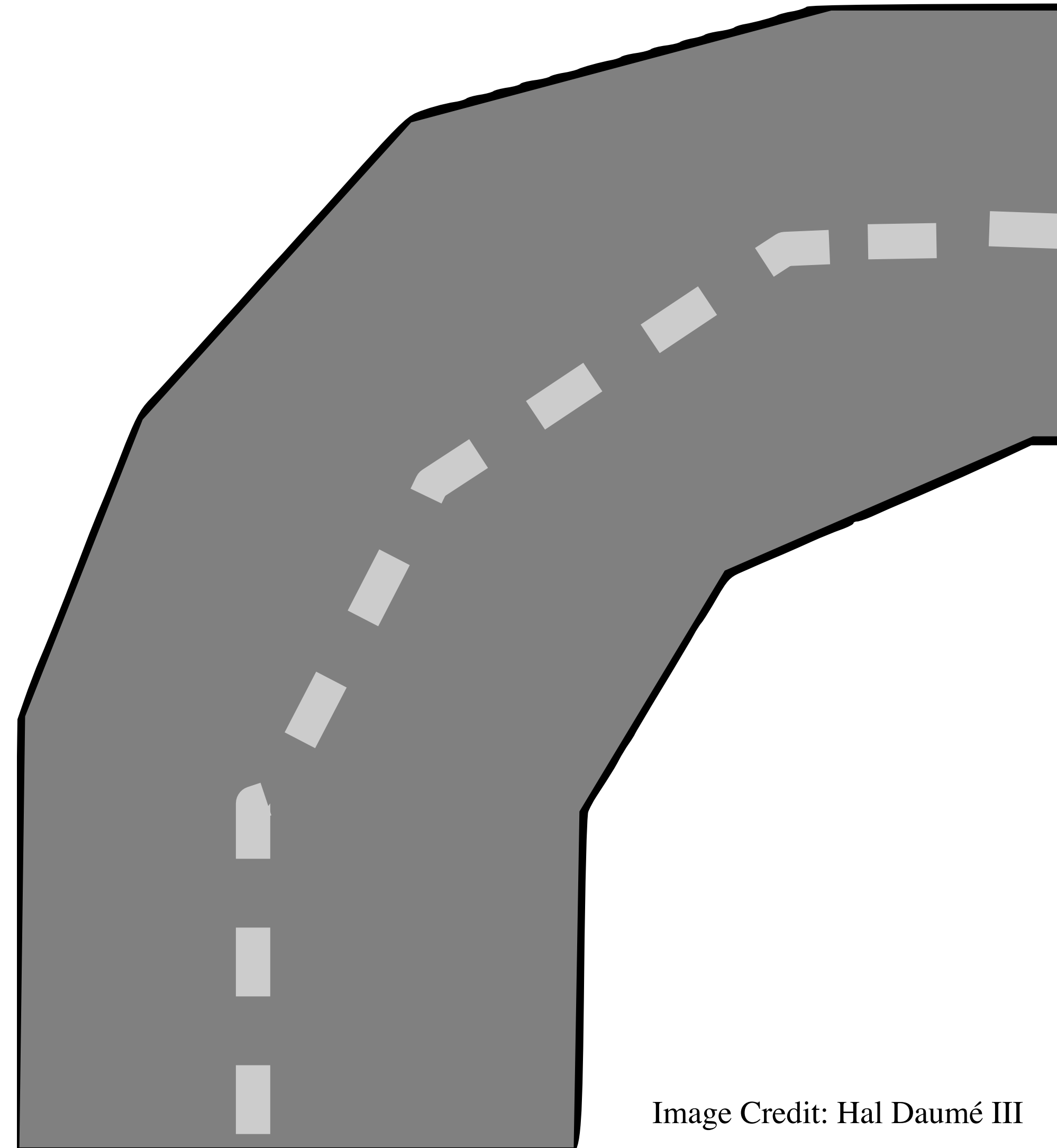
$$\underbrace{\mathbb{E}_{(s,a) \sim \rho_\pi} [C_U(s, a)]}_{\text{Stay within the expert distribution}} + \lambda \overbrace{D_{\text{KL}}(\rho_{\pi^*} || \rho_\pi)}^{\text{Mimic expert within the expert distribution}}$$

Stay within the expert distribution

Train ensemble of policies  $\Pi_E = \{\pi_1, \dots, \pi_E\}$  on expert data

**Uncertainty Cost:**  $C_U(s, a) = \text{Var}_{\pi \sim \Pi_E}(\pi(a | s))$

**DRIL cost can be optimized using any RL algorithm**





# Our Approach: DRIL [BSH ICLR 2020]

disagreement regularized imitation learning

Mimic expert within the expert distribution

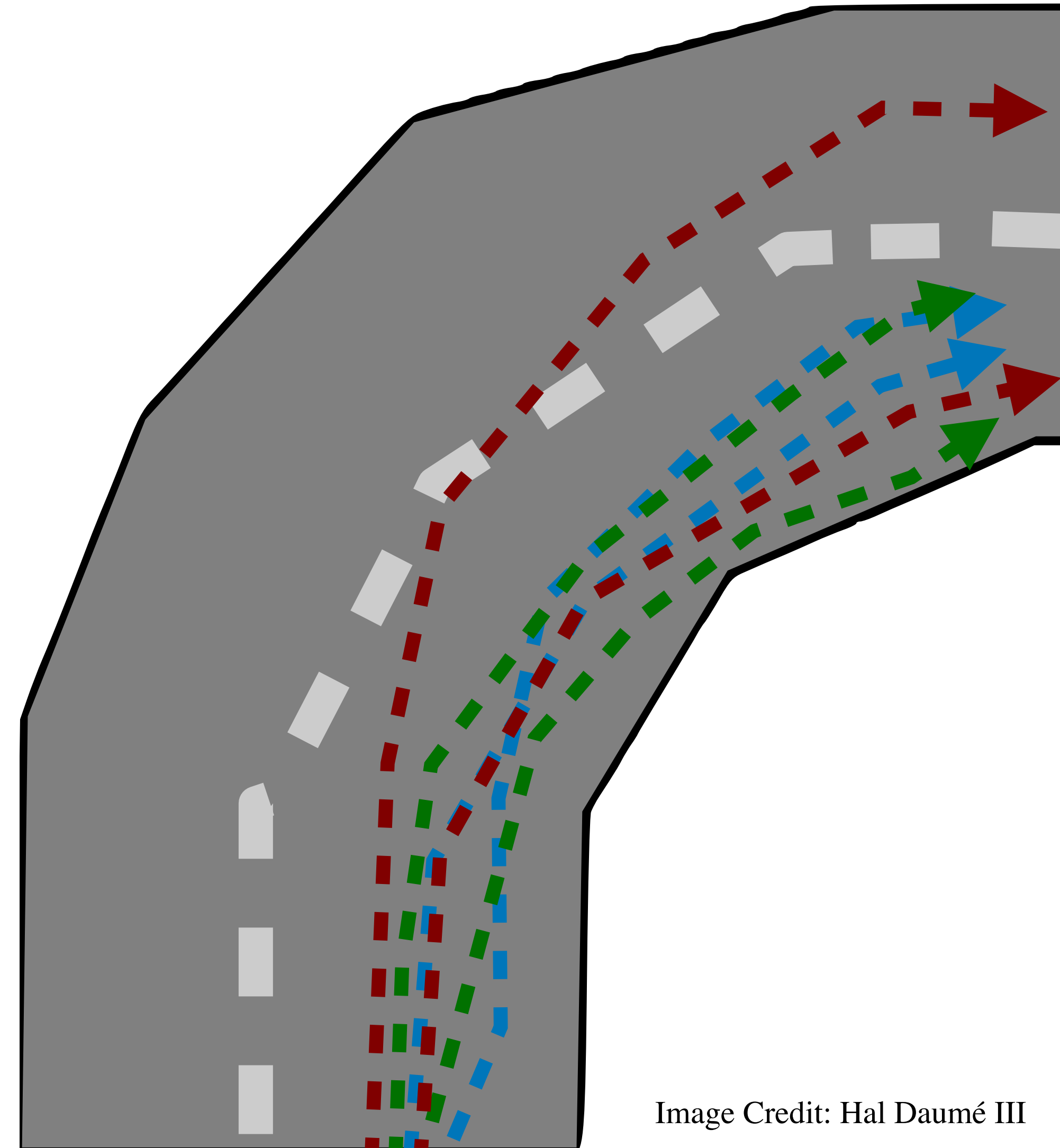
$$\underbrace{\mathbb{E}_{(s,a) \sim \rho_{\pi}} [C_U(s, a)]}_{\text{Uncertainty Cost}} + \lambda \overbrace{D_{\text{KL}}(\rho_{\pi^*} || \rho_{\pi})}_{\text{Disagreement Regularization}}$$

Stay within the expert distribution

Train ensemble of policies  $\Pi_E = \{\pi_1, \pi_2, \pi_3\}$  on expert data

**Uncertainty Cost:**  $C_U(s, a) = \text{Var}_{\pi \sim \Pi_E}(\pi(a | s))$

**DRIL cost can be optimized using any RL algorithm**



# Our Approach: DRIL [BSH ICLR 2020]

disagreement regularized imitation learning

Mimic expert within the expert distribution

$$\mathbb{E}_{(s,a) \sim \rho_\pi} [C_U(s, a)] + \lambda D_{\text{KL}}(\rho_{\pi^*} || \rho_\pi)$$

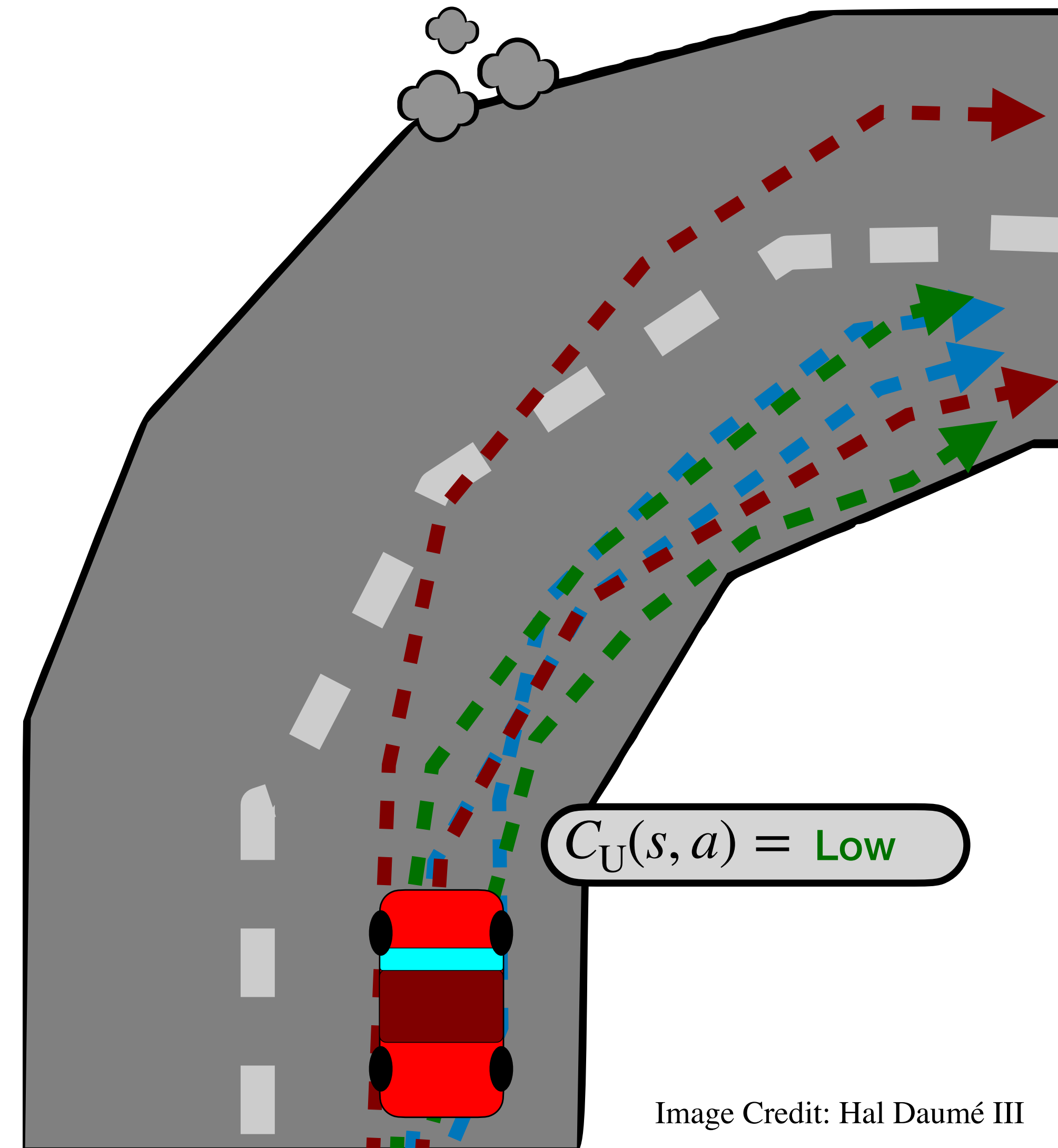
Stay within the expert distribution

Train ensemble of policies  $\Pi_E = \{\pi_1, \pi_2, \pi_3\}$  on expert data

Uncertainty Cost:  $C_U(s, a) = \text{Var}_{\pi \sim \Pi_E}(\pi(a | s))$

DRIL cost can be optimized using any RL algorithm

$C_U(s, a) = \text{High}$



$C_U(s, a) = \text{Low}$

# DRIL

## algorithm

**Train Policy Ensemble  $\Pi_E = \{\pi_1, \dots, \pi_E\}$  using demonstration data  $D$**

**Initialize  $\pi$  using behavior cloning**

**Loop:**

- **Interact in an environment  $(s, a) \sim \rho_\pi$**
- **Minimize  $D_{\text{KL}}(\rho_{\pi^*} || \rho_\pi)$  using mini-batch of expert data**
- **Minimize  $C_U(s, a)$  with the collected samples using RL**

# Theory of DRIL

Performance gap compared to globally optimal policy  $\pi^\star$

Behavior cloning  
 $Loss(\pi) = \epsilon$

$$J(\pi^\star) - J(\hat{\pi}) \leq \mathcal{O}(\kappa \epsilon H)$$

$$\kappa = \min_{\mathcal{U} \subseteq \mathcal{S}} \frac{\max_{\pi \in \Pi} \sup_{s \in \mathcal{U}} \frac{d_\pi(s)}{d_{\pi^\star}(s)}}{\min_{s \notin \mathcal{U}, a \in A} \text{VAR}_{\pi \sim \Pi_E}[\pi(a | s)]}$$

Concentrability inside of  $\mathcal{U}$

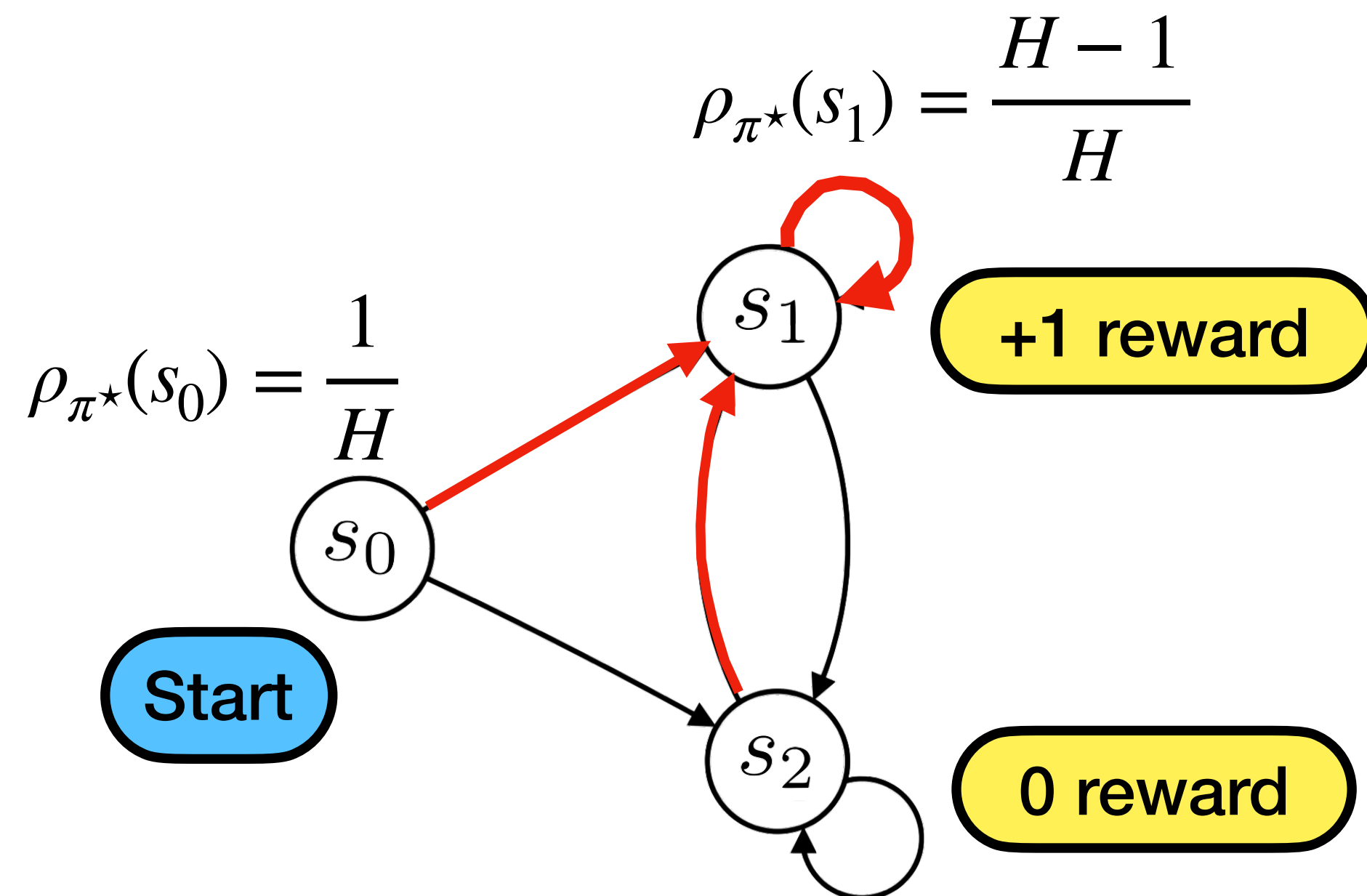
Minimum ensemble variance outside of  $\mathcal{U}$

$\kappa$  problem dependent.

Depends on environment dynamics, expert distribution, and uncertainty of ensemble

# Revisiting the covariate shift problem

Given an expert policy:  $\pi^*$



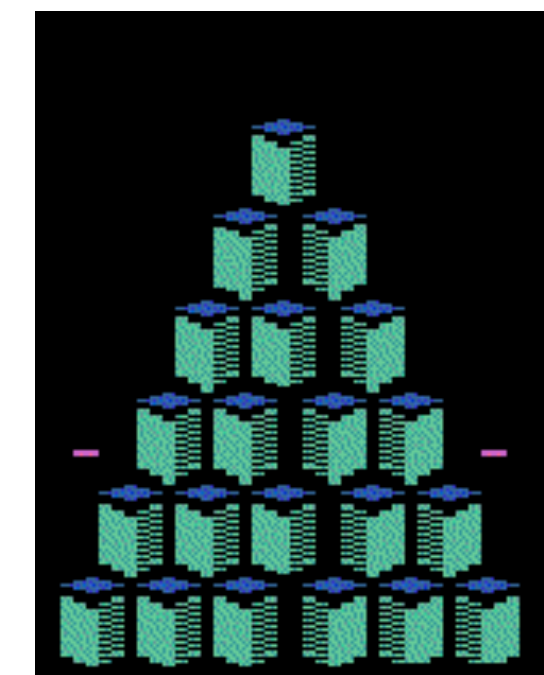
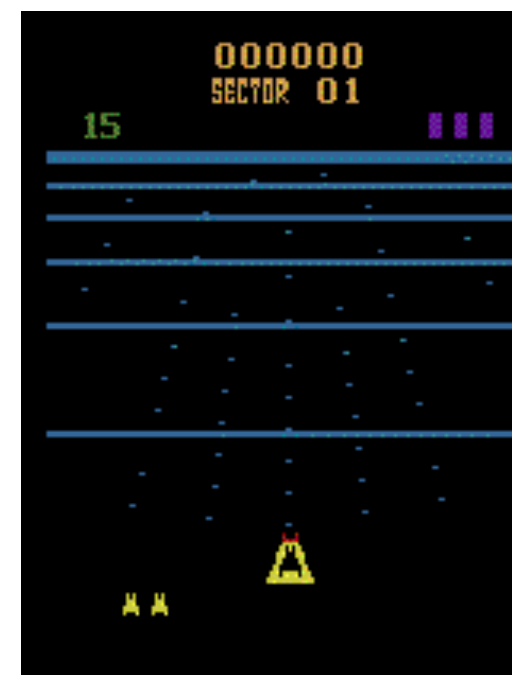
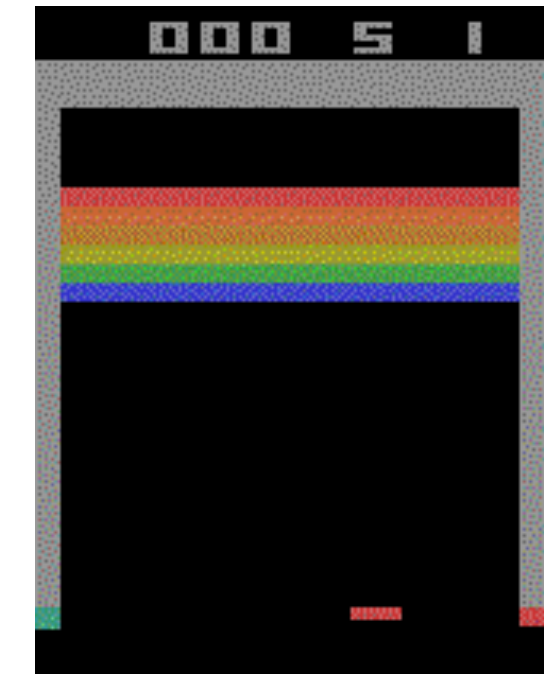
<p><b>Behavior Cloning</b> Performance Gap: <math>\mathcal{O}(\epsilon H^2)</math> <b>(gap scales quadratic)</b></p>
<p><b>Dagger</b> Performance Gap: <math>\mathcal{O}(\epsilon H)</math> <b>(gap scales linear)</b></p>
<p><b>DRIL</b> Performance Gap: <math>\kappa = 12 + \mathcal{O}\left(\frac{1}{\sqrt{ \text{ensemble} }}\right)</math> <b>(gap scales linear)</b></p>



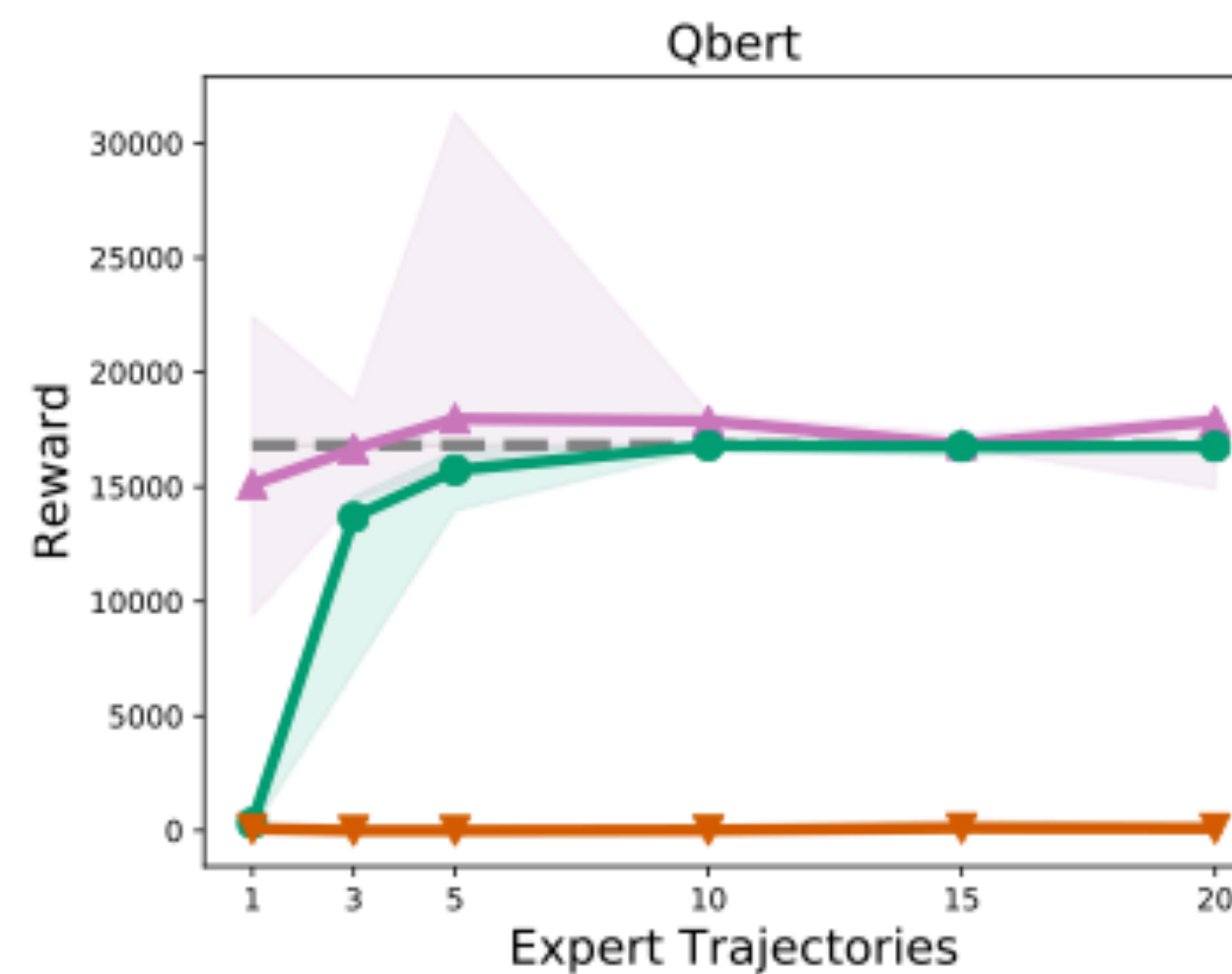
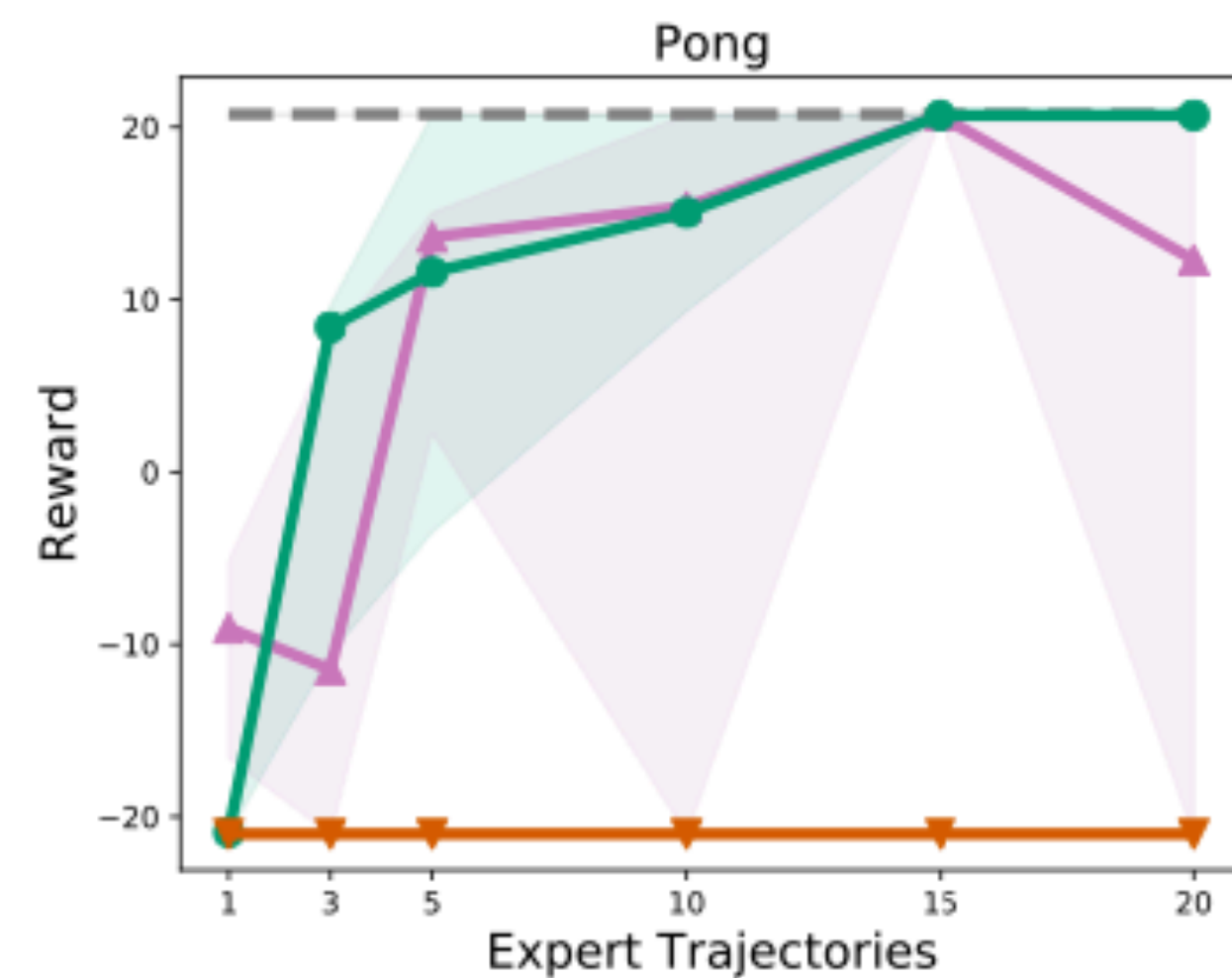
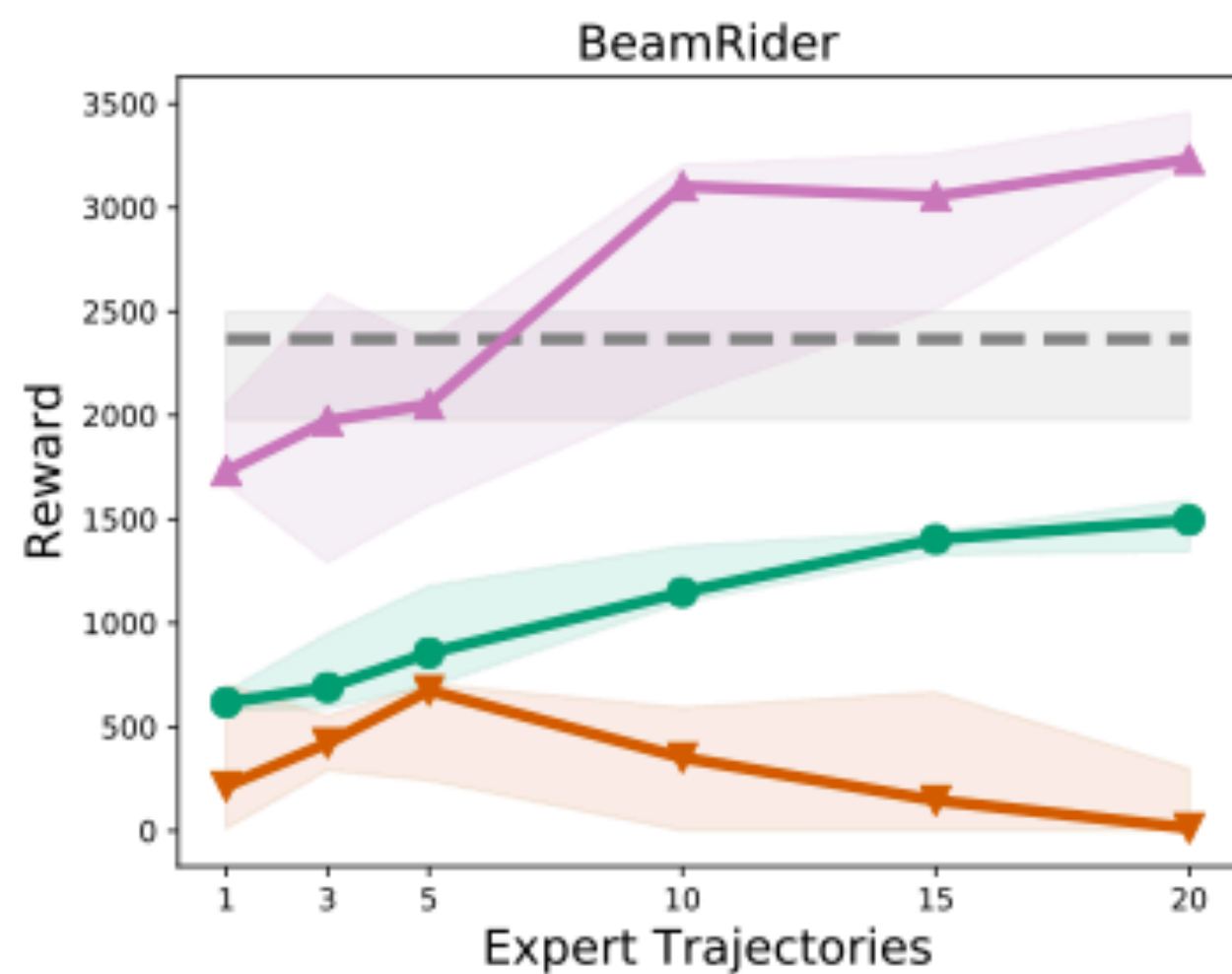
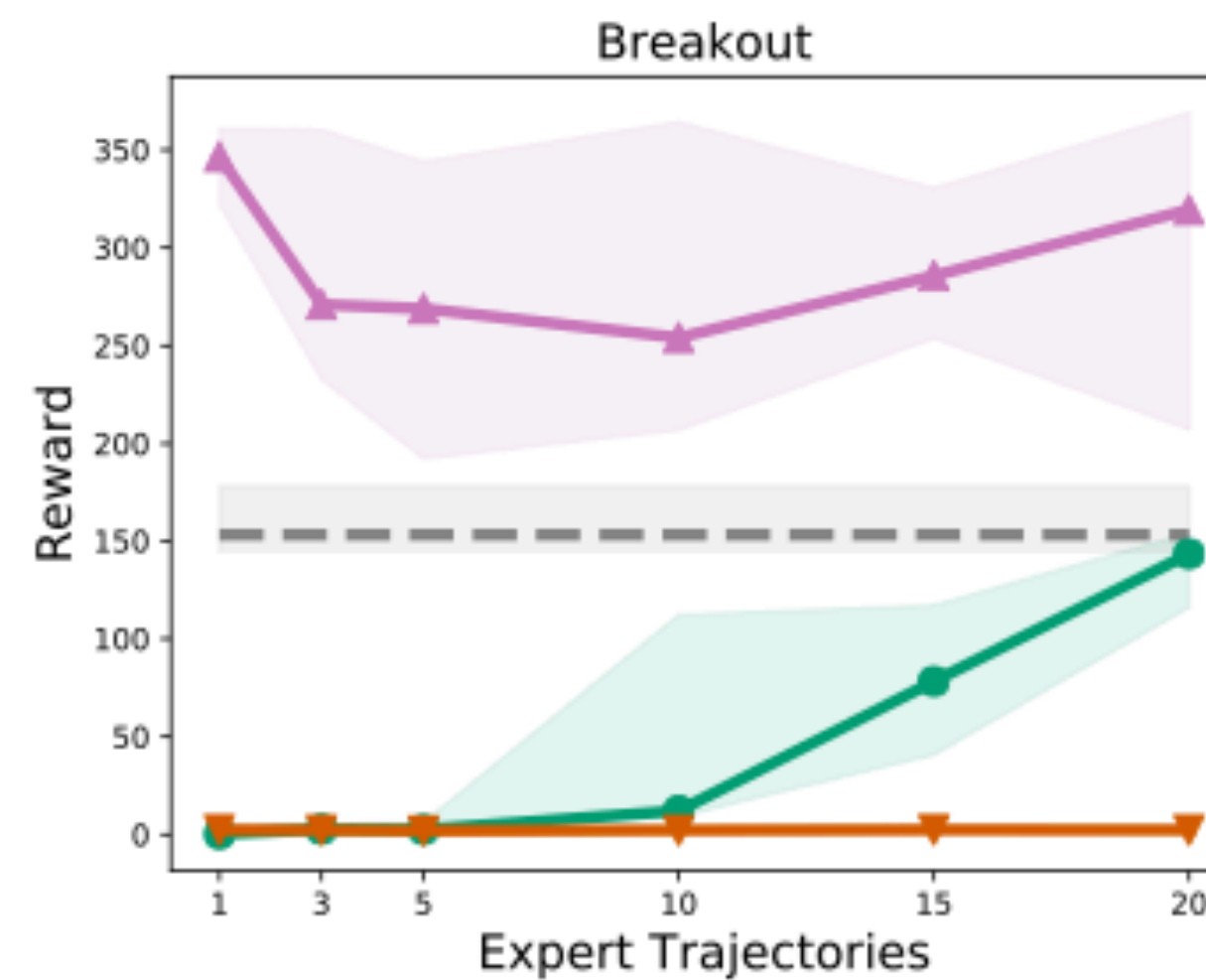
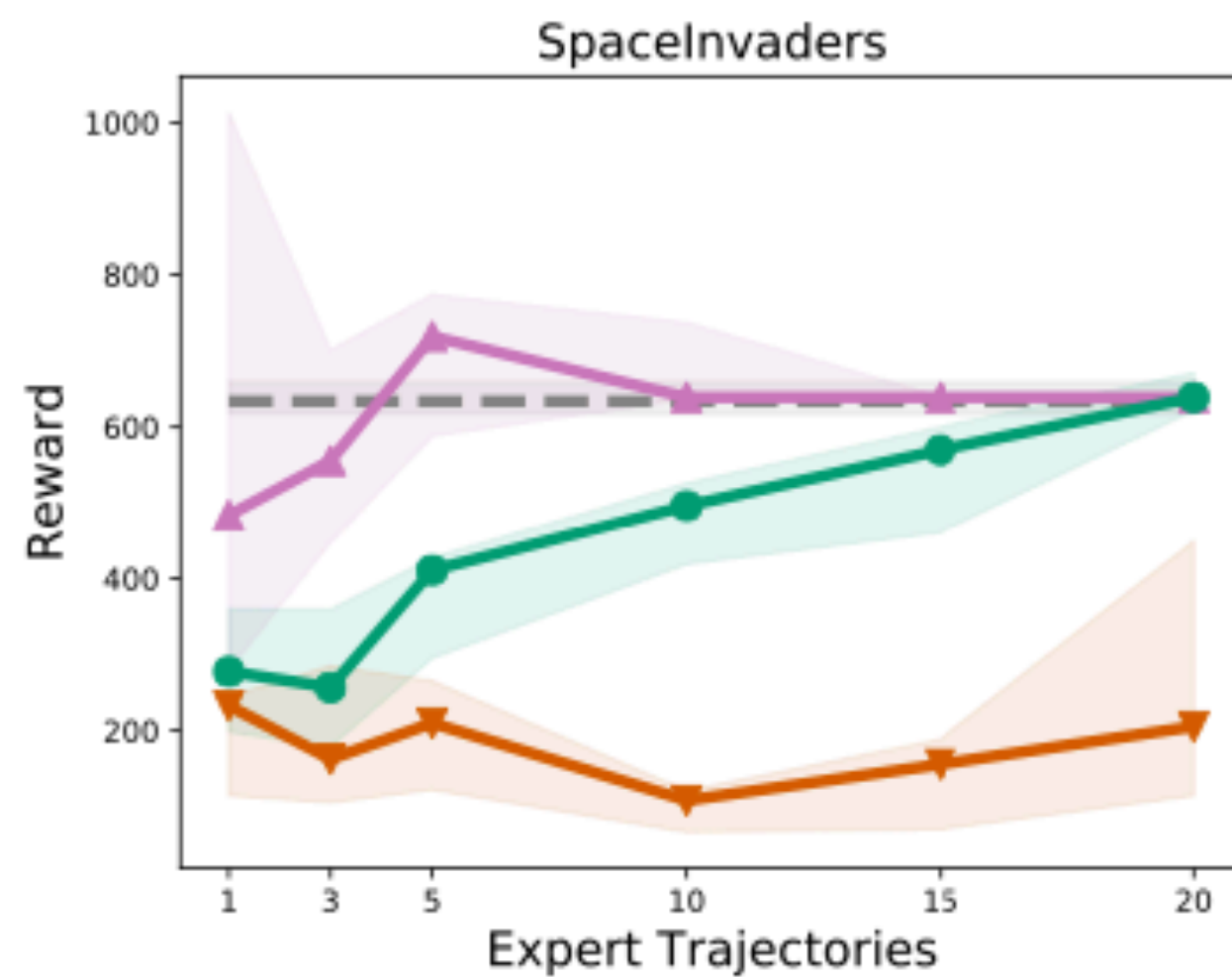
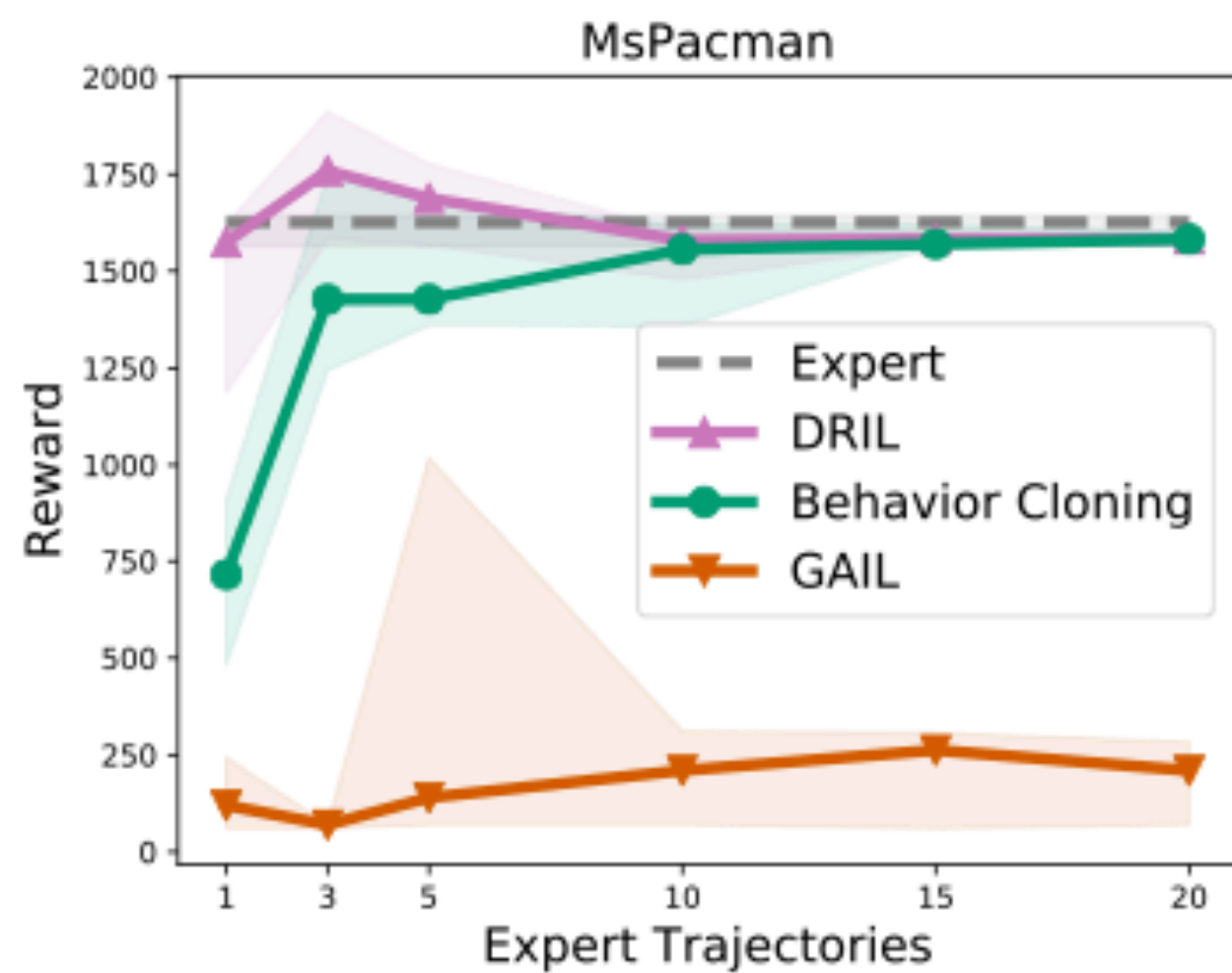
# Experimental Setup

- Generated demonstration data from pre-trained experts
- Then trained DRIL on the demonstration data and evaluated on true environment reward

## Atari Imitation Learning Benchmark



# Experimental Results



# Impact of our work

**extended to LLMs** [Coste et al. ICLR 2024]

## Task

Alpaca datasets: General instruction fine-tuning dataset incorporating wide range of tasks

## Dataset Composition

52,000 instructions and demonstrations generated from GPT 3.5

## Example of Instructions

What is the capital of France?

The capital of France is Paris.

Describe the function of a computer  
motherboard

A computer motherboard is the main circuit  
board in a computer. It is the backbone of ...

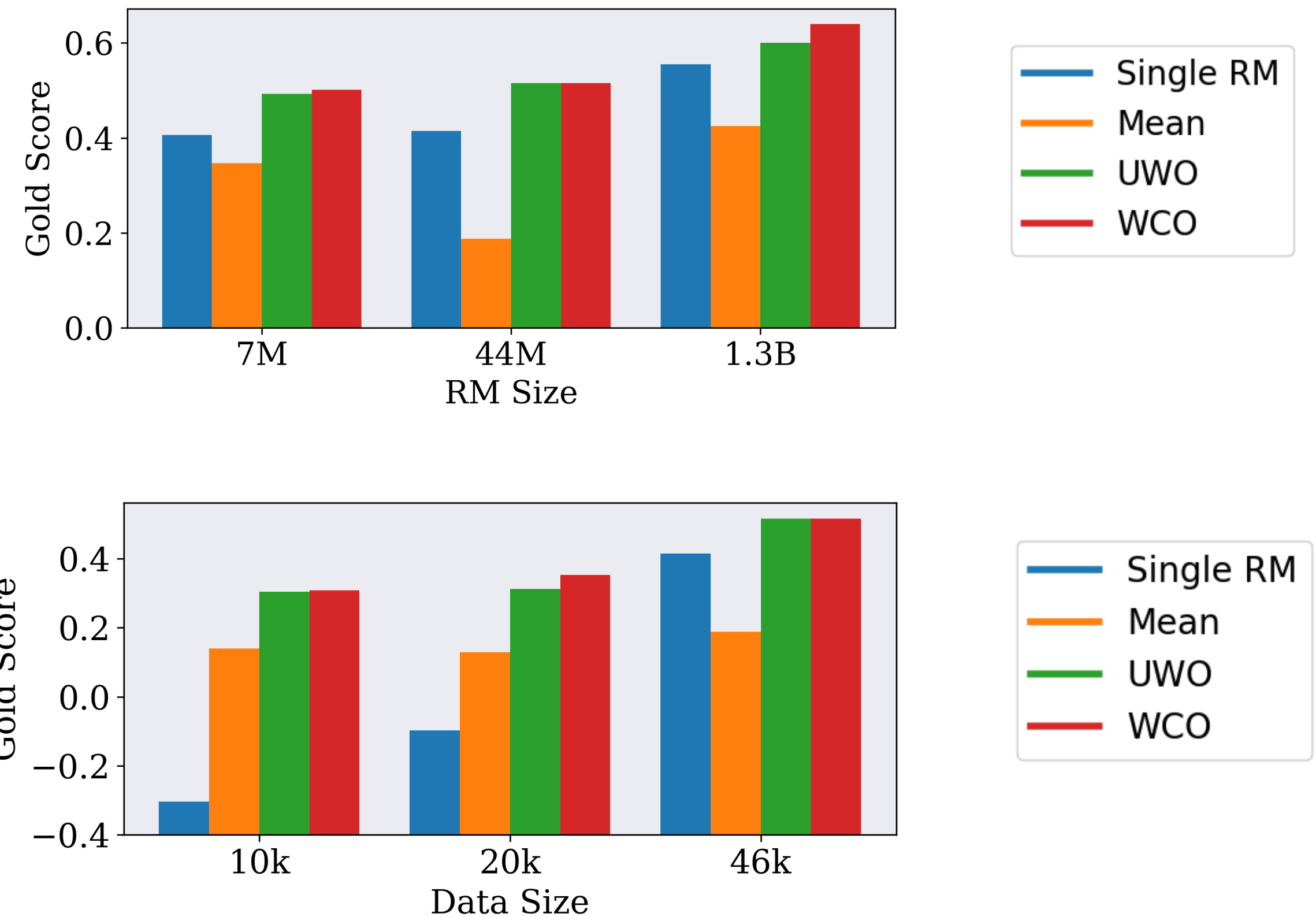


# Impact of our work

extended to LLMs [Coste et al. ICLR 2024]

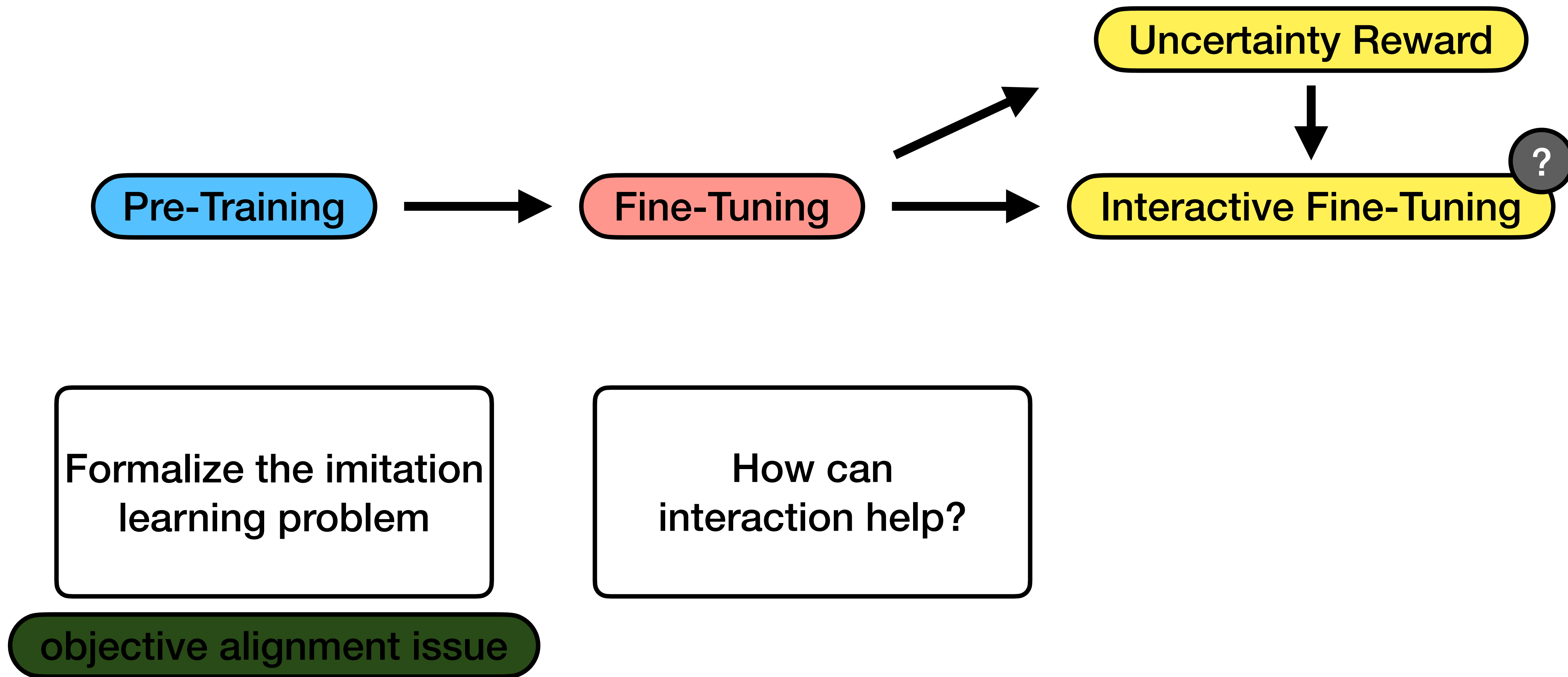
## Uncertainty-Weighted Optimization (UWO):

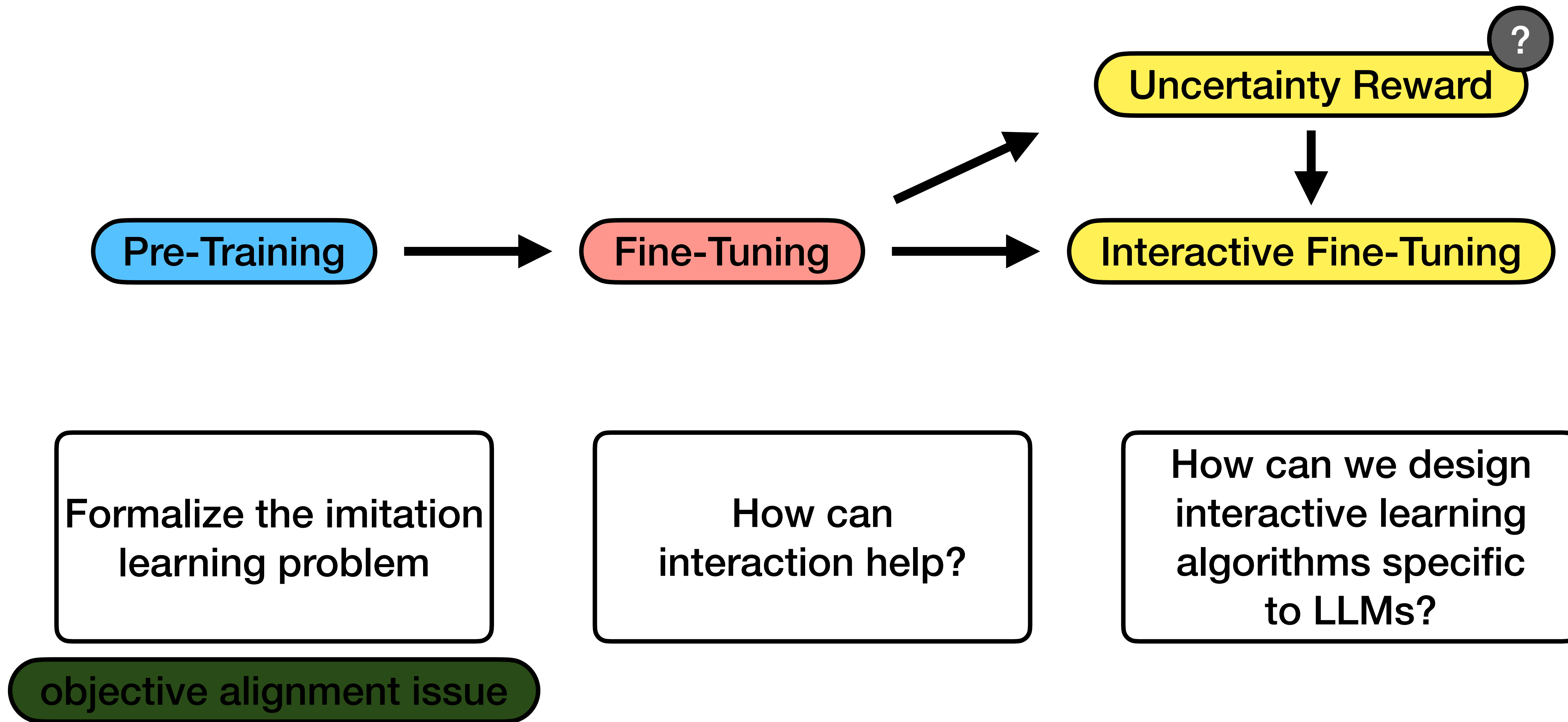
UWO works by penalizing the policy for generating responses for which there is high disagreement among reward models within the ensemble



# Summary

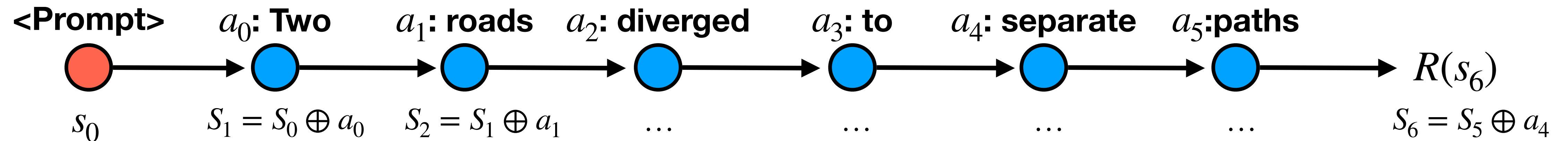
- Interacting in an environment can help address covariate shift issues
- DRIL has theoretical guarantees in some settings
- DRIL empirically performs better than Behavior Cloning





# Reinforcement Learning

## for text-generation



**Initial State:**

$s_0 \sim D$  a prompt

**Actions:**

words from our vocabulary

**States:**

$s_t = s_0 \oplus \langle a_0, a_1, a_2, a_3, a_4, a_5 \rangle$

**Transition:**

$P(s' | s, a)$

**Reward:**

$R(s)$  score

**Horizon:**

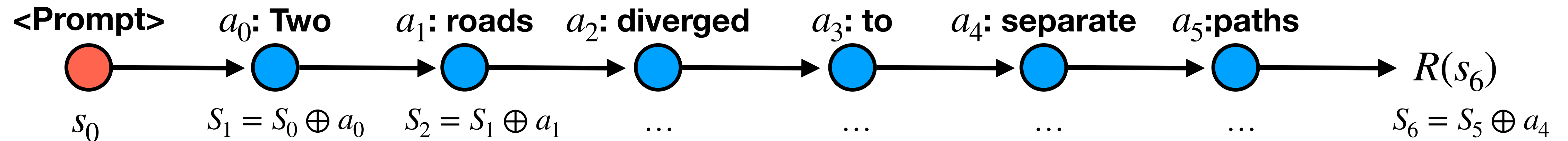
$H$  max generation length

**Policy:**

$\pi(\cdot | s)$

**Goal:** Learn policy  $\pi$  to maximize reward

# Reinforcement Learning for text-generation



**Initial State:**

**Actions:**

**States:**

**Transition:**

**Reward:**

**Horizon:**

**Policy:**

$s_0 \sim D$  a prompt

words from our vocabulary

$s_6 = s_0 \oplus \langle a_0, a_1, a_2, a_3, a_4, a_5 \rangle$

$P(s' | s, a)$

$R(s)$  score **NLP Metric**

$H$  max generation length

$\pi(\cdot | s)$

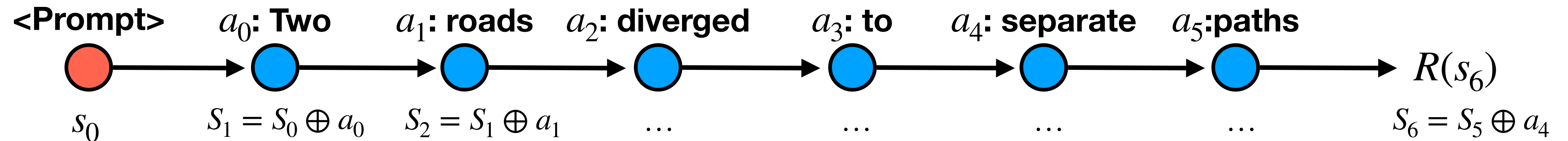
**Is reinforcement learning (not) for natural language processing?:  
Benchmarks, baselines, and building blocks for natural language policy optimization**

[RABHSBHC ICLR 2023]

The first benchmark comparing interactive fine-tuning with RL and supervised fine-tuning across a range of NLP tasks.

# Reinforcement Learning

## for text-generation



**Initial State:**

**Actions:**

**States:**

**Transition:**

**Unknown Reward:**

**Horizon:**

**Policy:**

$s_0 \sim D$  a prompt

words from our vocabulary

$S_6 = S_0 \oplus \langle a_0, a_1, a_2, a_3, a_4, a_5 \rangle$

$P(s' | s, a)$

$R(s)$  score

$H$  max generation length

$\pi(\cdot | s)$

**Given Preference dataset** from an expert  $\pi^\star$

**Goal:** Learn policy  $\pi$  using this dataset

# Reinforcement Learning From Human Feedback

## 1. Supervised Fine-Tuning (SFT)

$$\pi^{\text{SFT}} = \arg \min_{\pi \in \Pi} \mathbb{E}_{(s,a) \sim \rho_{\pi^*}} [L(\pi, s, a)]$$

## 3. Learn a preference reward function

for example using logistic regression

$$\hat{r} = \arg \min_r \sum_D \log \sigma(r(\tau_1) - r(\tau_2))$$

## 4. Interactive Fine-Tuning with RL

$$\mathbb{E}_{\tau \sim \rho_{\pi}} [\hat{r}(\tau)] + \lambda D_{\text{KL}}(\pi || \pi^{\text{SFT}})$$

## 2. Collect preferences data

Prompt  $s_0 \sim D$

Calculate the median number [2, 3, 7, 8, 10]

Responses from  $\pi^{\text{SFT}}$

0 7

Human Labeler rank outputs

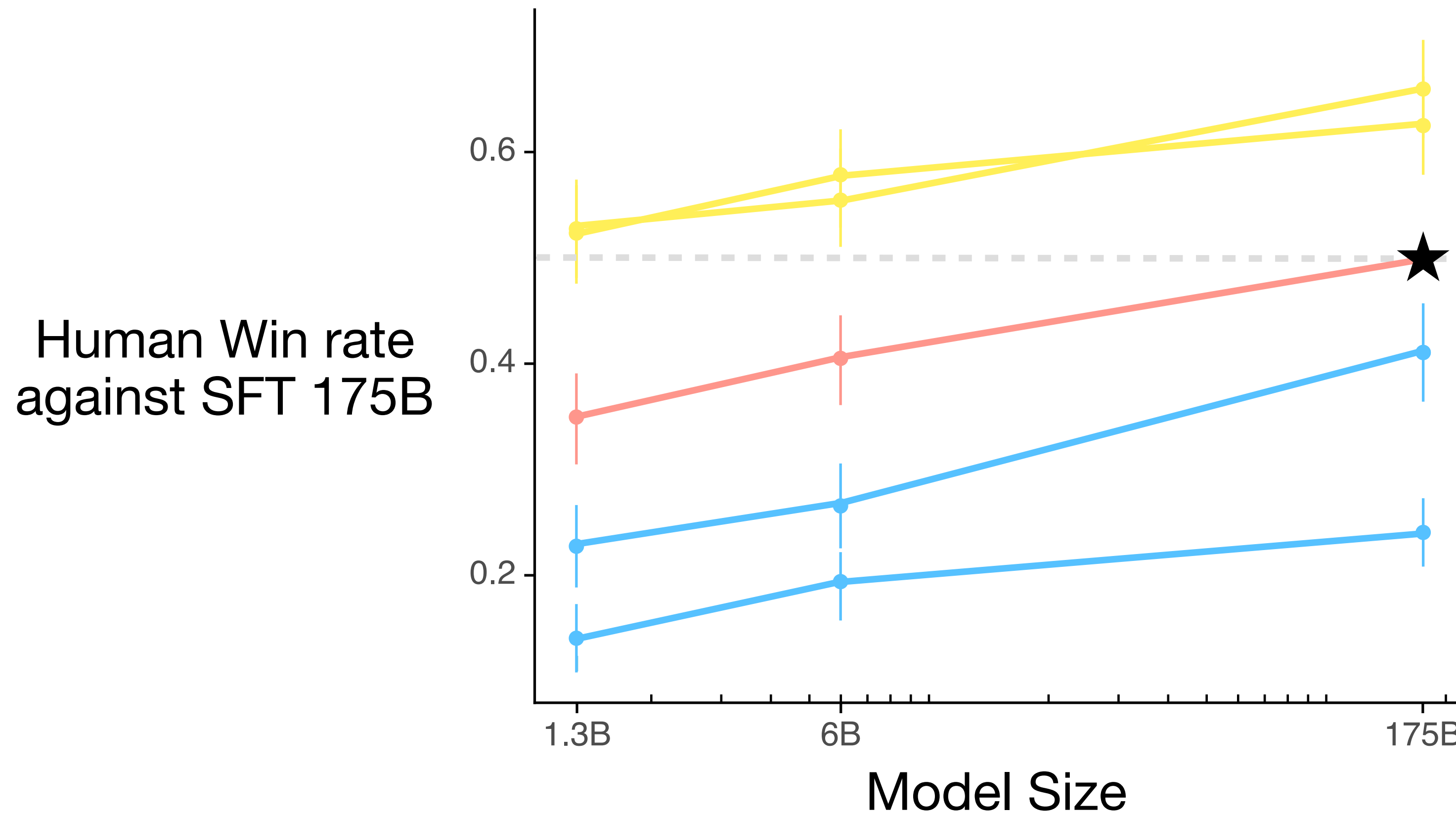
7 > 0

Aggregate

$$D = D \cup \{(\tau_1, \tau_2, s_0)\}$$



# Reinforcement Learning From Human Feedback



Long Ouyang et al.  
Training language models to follow instructions with human feedback  
**OpenAI 2022**

# Reinforcement Learning From Human Feedback

Pre-Training

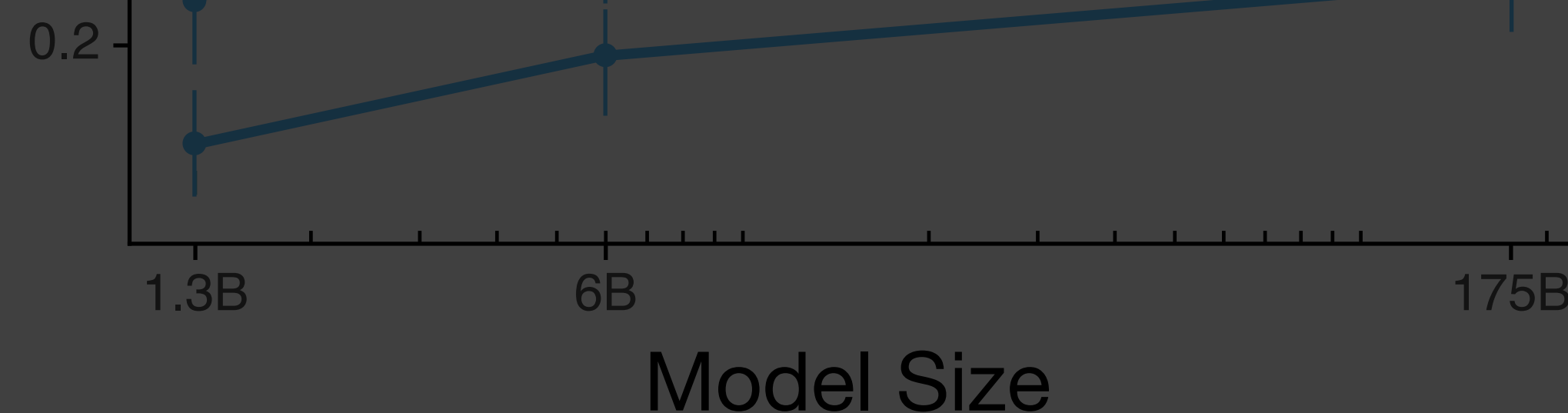
Human Feedback

Reinforcement Learning From Human Feedback

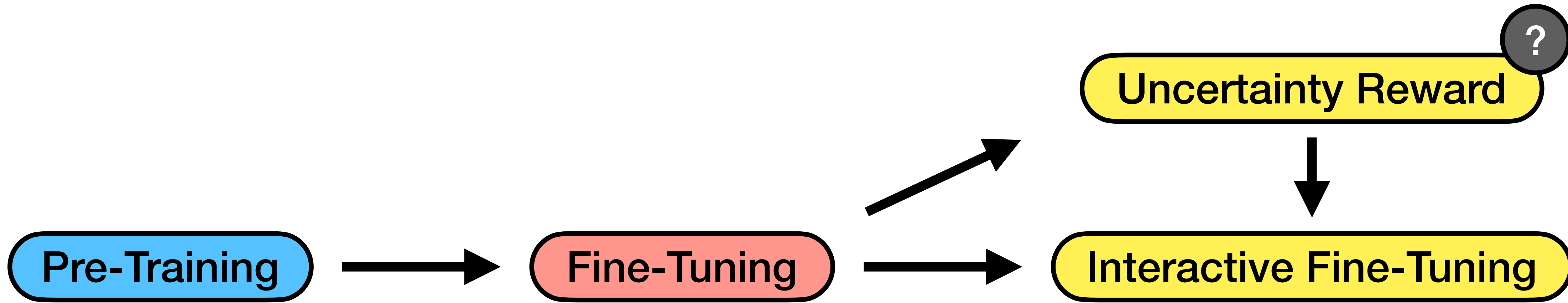
## Problem:

- Text generation is a **LARGE** combinatorial search problem.
- RL algorithms are complex and **unstable**, making them challenging to **tune effectively**.

Human Win rate  
against SFT 175B



Long Ouyang et al.  
Training language models to follow  
instructions with human feedback  
**OpenAI 2022**

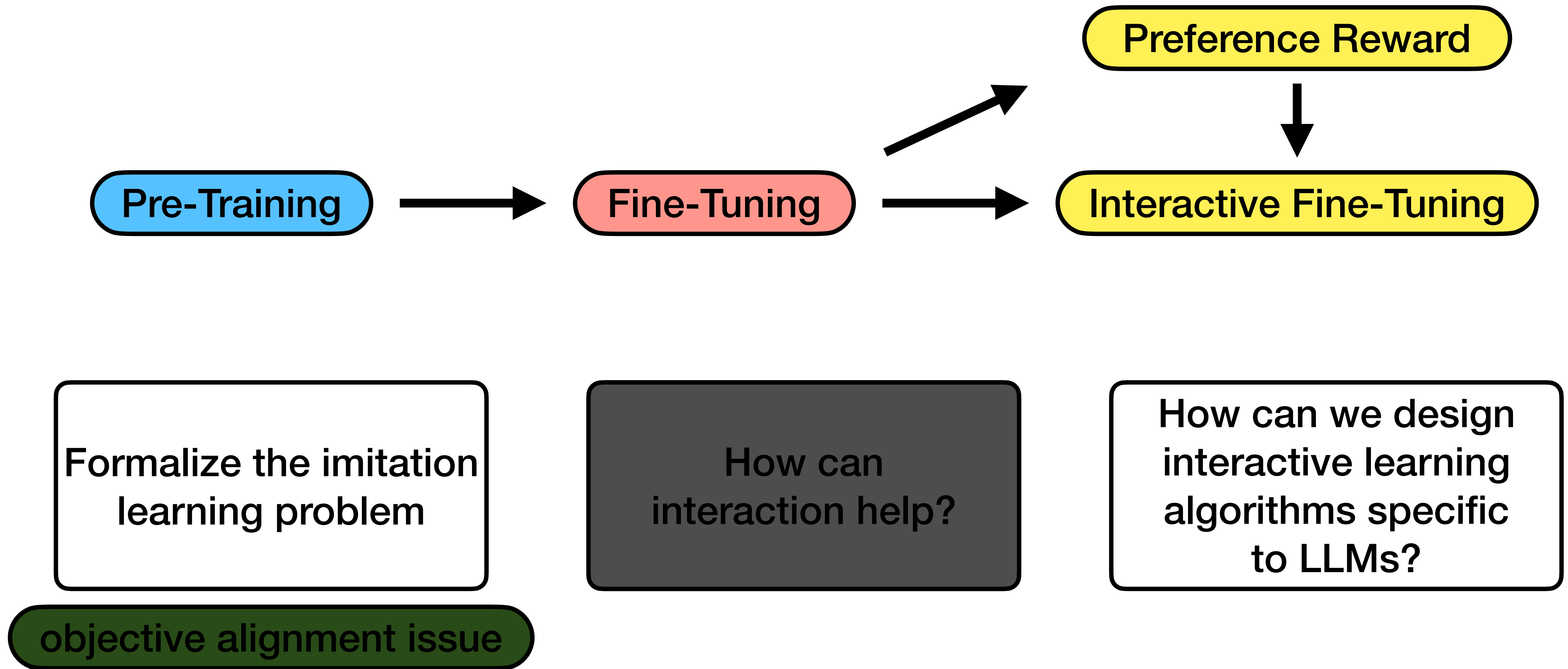


Formalize the imitation learning problem

objective alignment issue

How can interaction help?

How can we design interactive learning algorithms specific to LLMs?

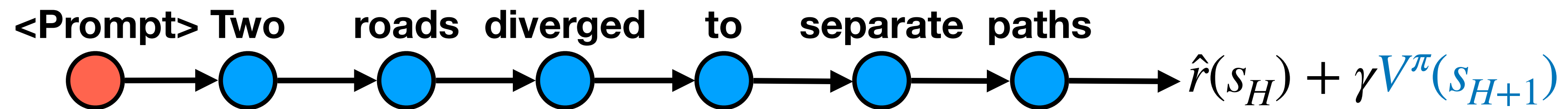


# Reset Property

Reset allows us to rollout a policy from partial sentences

?

1. Sample a prompt from  $D$



2. Sample a generation from  $\pi$

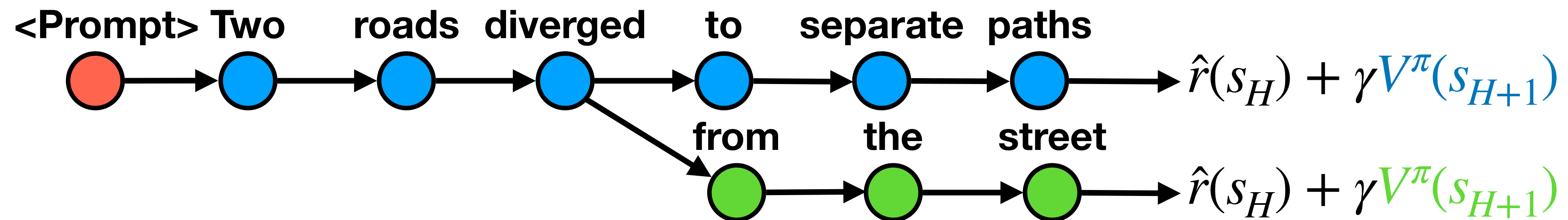
# Reset Property

Reset allows us to rollout a policy from partial sentences

Transition:  $P(s' | s, a)$

Deterministic

1. Sample a prompt from  $D$



2. Sample a generation from  $\pi$

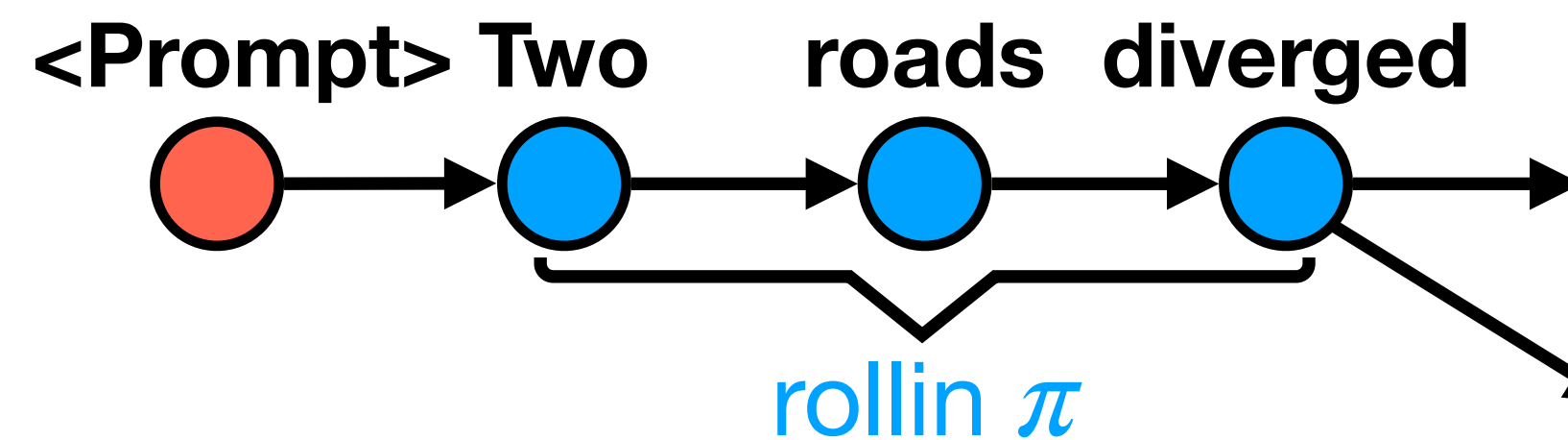
# Reset Property

Reset allows us to rollout a policy from partial sentences

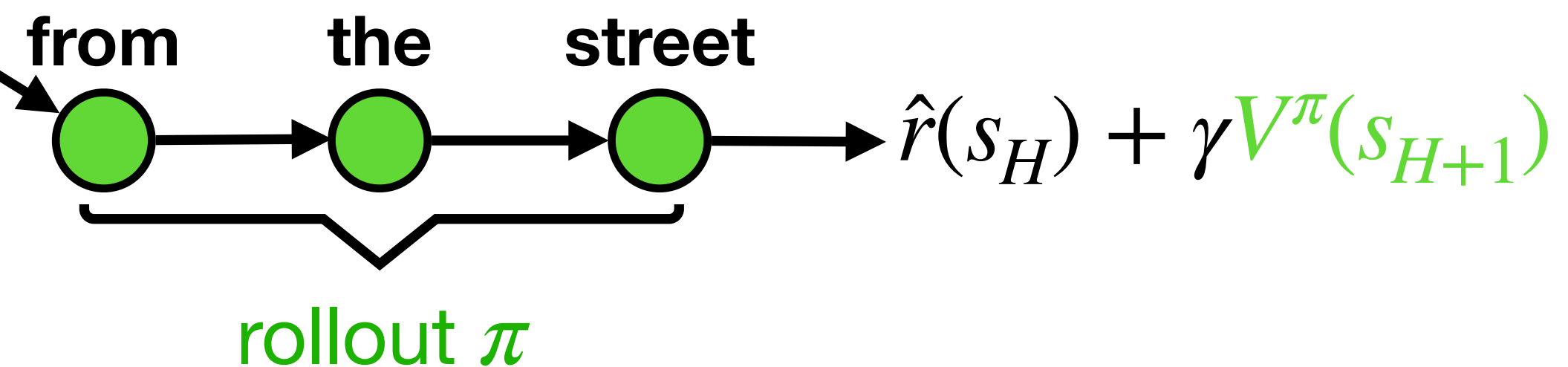
Transition:  $P(s' | s, a)$

Deterministic

1. Sample a prompt from  $D$



2. Sample a generation from  $\pi$



# Reset Property

Reset allows us to rollout a policy from partial sentences

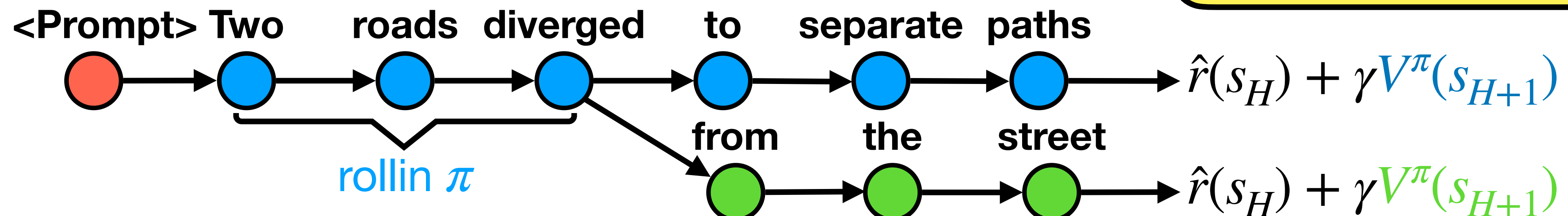
?

Transition:  $P(s' | s, a)$  **Deterministic**

$$\mathbb{E}_{\tau \sim \rho_{\pi}} [\hat{r}(\tau)] + \underbrace{\lambda D_{\text{KL}}(\pi || \pi^{\text{SFT}})}$$

we do not want to move too far from  $\pi^{\text{SFT}}$

1. Sample a prompt from  $D$

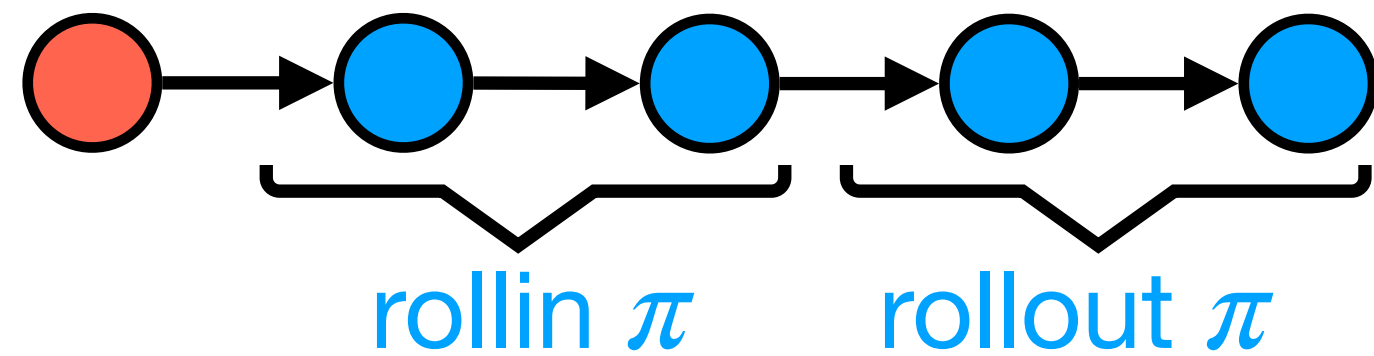


2. Sample a generation from  $\pi$



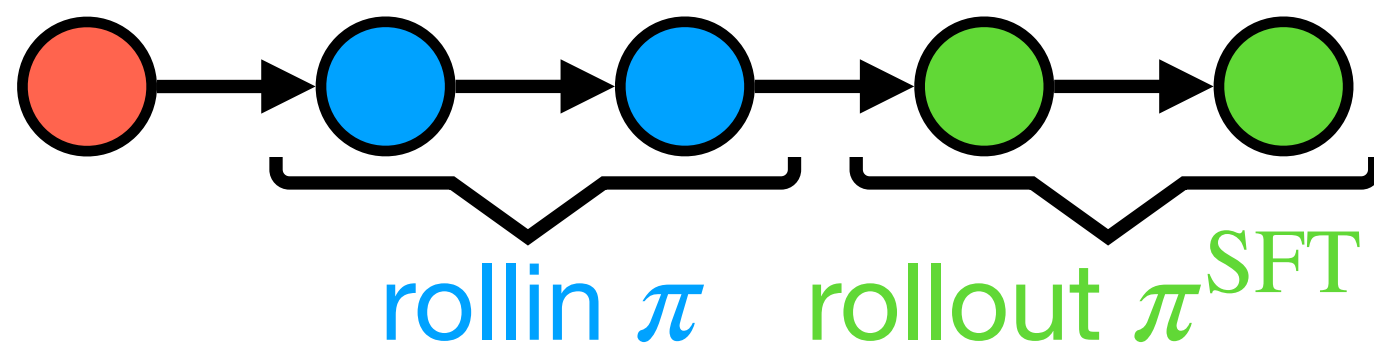
# Approaches

## PPO (RL algorithm)



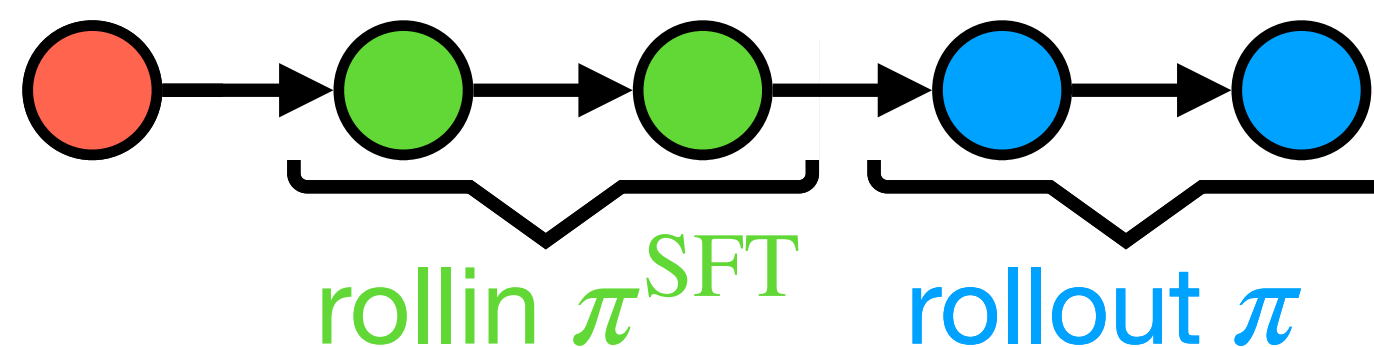
- Does not take advantage of problem specific structure
- Samples prompts  $s_0 \sim D$
- Scores action with  $\hat{r} + \gamma V^\pi$

## AggreVaTeD



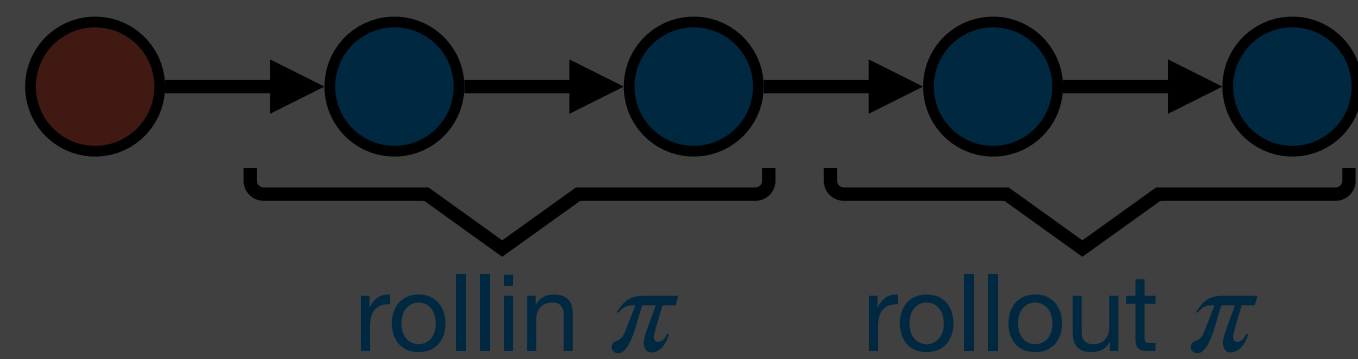
- Samples prompts  $s_0 \sim D$
- Scores action with  $\hat{r} + \gamma V^{\pi^{\text{SFT}}}$
- Intuition: Receive feedback on partial generations

## PPO++ [CBRMS Instruction Workshop 2023]



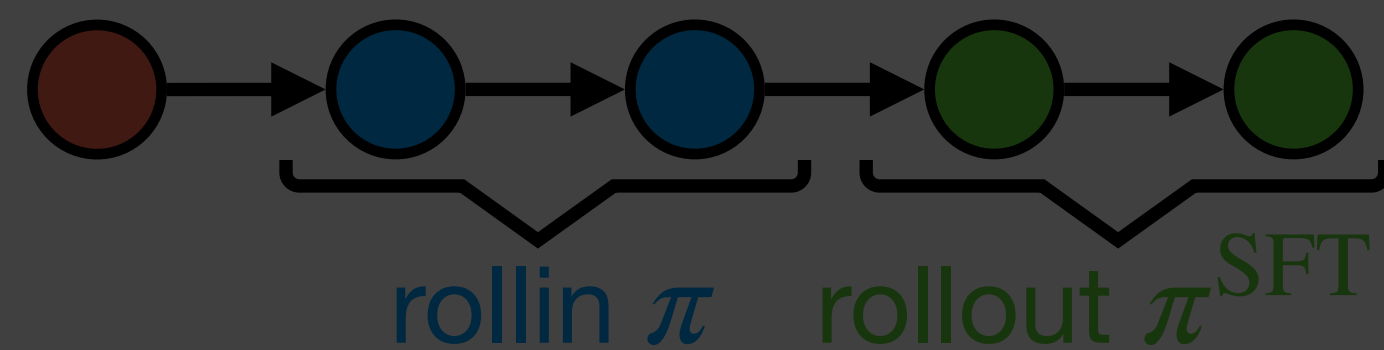
# Approaches

## PPO (RL algorithm)



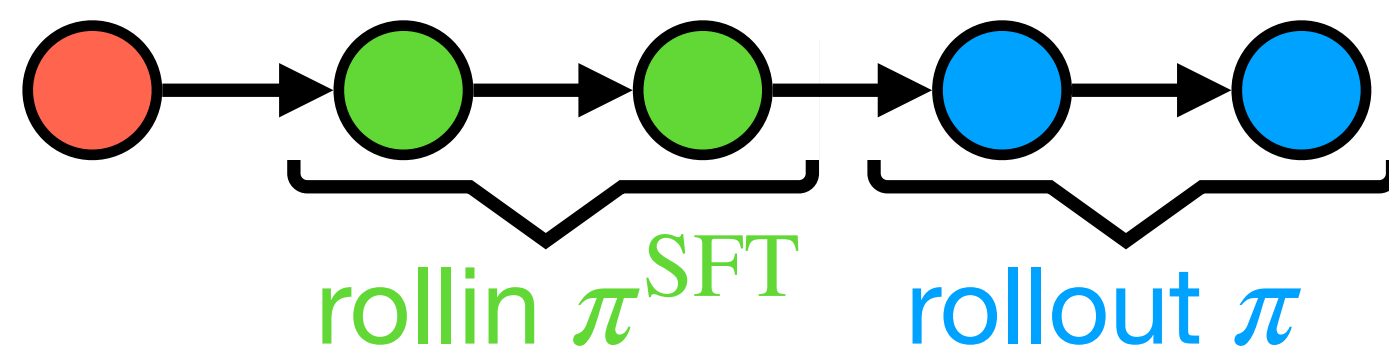
- Does not take advantage of problem specific structure
- Samples prompts  $s_0 \sim D$
- Scores action with  $\hat{r} + \gamma V^\pi$

## AggreVaTeD



- Samples prompts  $s_0 \sim D$
- Scores action with  $\hat{r} + \gamma V^{\pi^{\text{SFT}}}$
- Intuition: Receive feedback on partial generations

## PPO<sup>++</sup> [CBRMS Instruction Workshop 2023]



- Samples prompts  $s_0 \sim \beta \rho_{\pi^{\text{SFT}}} + (1 - \beta)D$
- Scores actions with  $\hat{r} + \gamma V^\pi$
- Intuition: Richer initial state distribution

# Theory of PPO++

Let  $\pi^\star$  be a high quality policy covered by  $\pi^{\text{SFT}}$

$$\underbrace{J(\pi^\star) - J(\pi^t)}_{\text{Performance gap}} \leq O \left( H^2 \max_s \underbrace{\left( \frac{d^{\pi^\star}(s)}{d^{\pi^{\text{SFT}}}(s)} \right)}_{\substack{\text{Assume bound density ratio and} \\ \pi^{\text{SFT}} \text{ provides coverage for } \pi^\star}} \underbrace{\epsilon}_{\substack{\mathbb{E}_{s \sim \beta \rho^{\pi^{\text{SFT}}} + (1-\beta)D} \left[ \max_a A^{\pi^t}(s, a) \right] \leq \epsilon \\ \text{Assume that one-step local} \\ \text{improvement over } \pi^t \text{ is small}}} \right)$$

# Experimental Setup

## Task Statement


Given a reddit post, write a TL;DR (short summary).

## Dataset Composition

- 210K Prompts total
- 117K Prompts with *Human Labels*
- 93K Prompts with *Human Preference Labels*

# Experimental Results


Uncertainty of a model at predicting  
the ground truth summaries



Algorithms	Perplexity (↓)	GPT4 Win Rate (↑)
SFT	14.09	29.5%
SFT+PPO	14.87	60.7%
SFT+PPO <sup>++</sup>	<b>13.42</b>	<b>64.4%</b>
SFT+AggreVaTeD	13.53	54.12%

# Experimental Results

Rafailov et al. 2023 conducted human study showing alignment between GPT4 and Human preference



Algorithms	Perplexity (↓)	GPT4 Win Rate (↑)
SFT	14.09	29.5%
SFT+PPO	14.87	60.7%
SFT+PPO <sup>++</sup>	<b>13.42</b>	<b>64.4%</b>
SFT+AggreVaTeD	13.53	54.12%

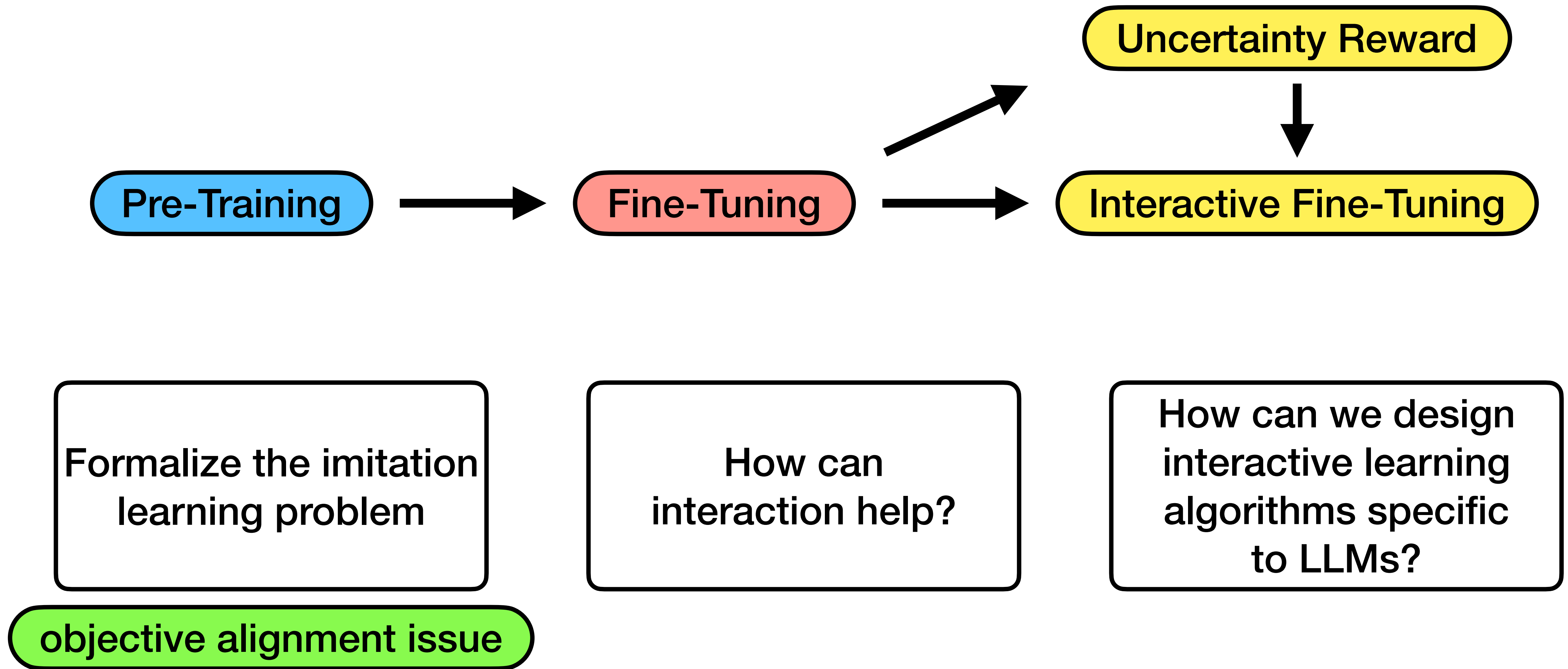
## Winrate Prompt Template:

Which of the following summaries does a better job of summarizing the most important points in the given forum Post? FIRST provide a one-sentence comparison of the two summaries, explaining which you prefer and why. SECOND, on a new line, state only "A" or "B" to indicate your choice.

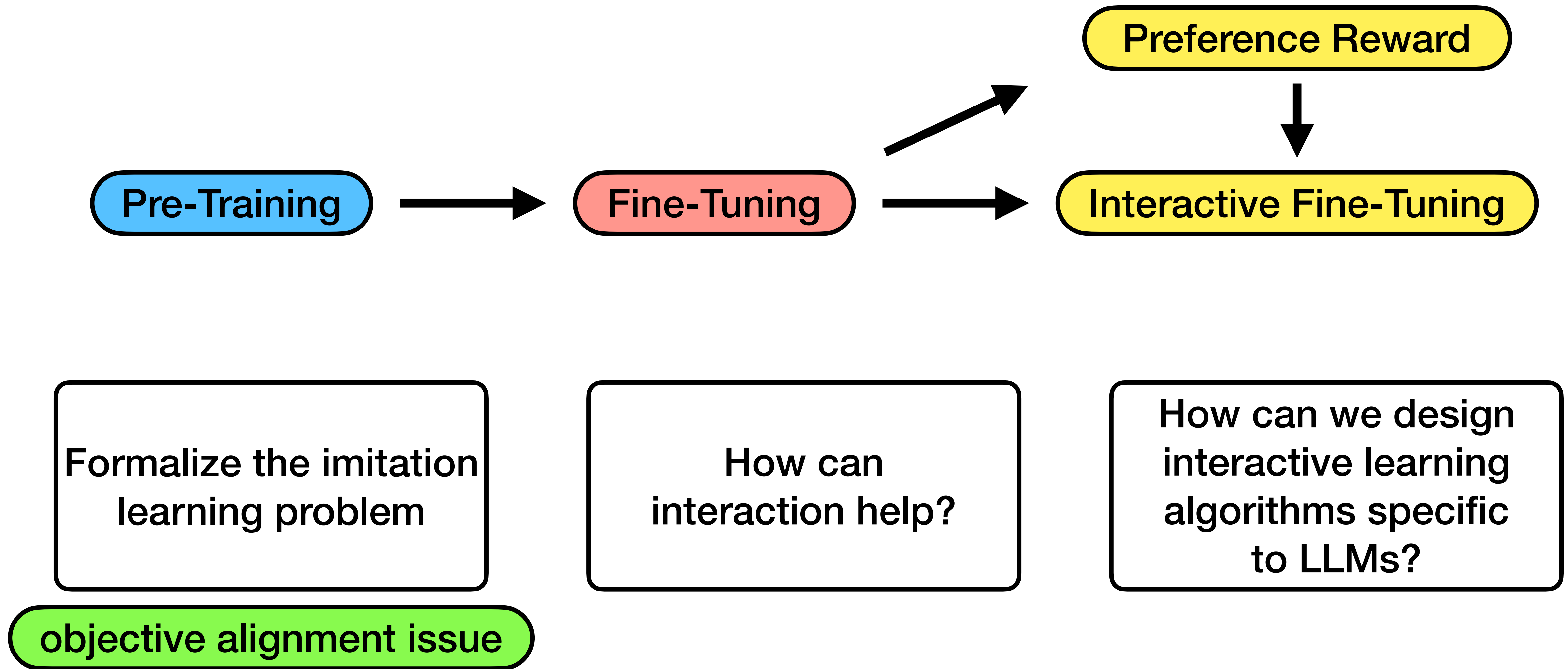
Post: <Post>  
A: <TLDR A>  
B: <TLDR B>

# Summary

- Designing RL algorithms specific for LLMs can improve performance
- Resetting is a special property of MDPs for LLMs
- PPO++ simple algorithm that use reset







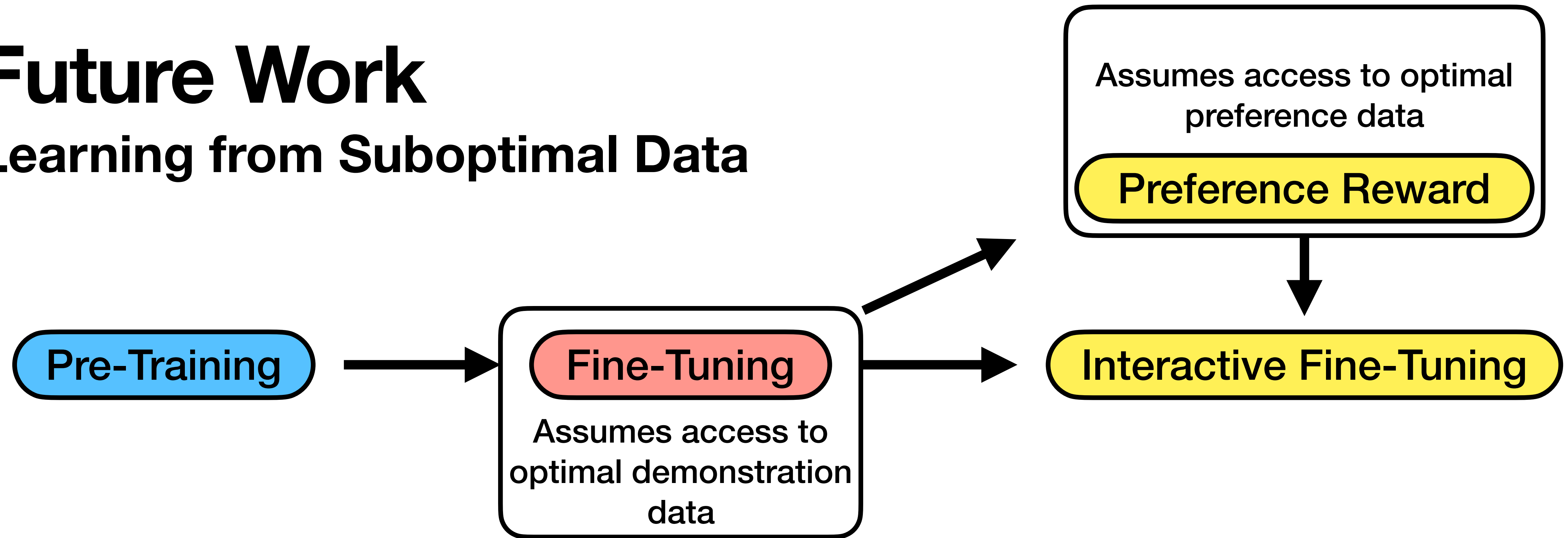
# Future Work

- Imitation Learning [1, 2]
- Imitation Learning with a Computational Oracle [6, 7, 8, 12]
- Imitation Learning with Suboptimal Data

1. Disagreement-Regularized Imitation Learning.....	[BSH ICLR 2020]
2. Adversarial Imitation Learning via Boosting.....	[CSHBS ICLR 2024]
3. Is reinforcement learning (not) for natural language processing?.....	[RABHSBHC ICLR 2023]
4. Learning to Generate Better Than Your LLM.....	[CBRMS Instruction Workshop 2023]
5. Policy-Gradient Training of Language Models for Ranking.....	[GCCBJ FMDM Workshop 2023]
6. Non-Monotonic Text Generation.....	[WBDC ICLR 2019]
7. Active Imitation Learning with Noisy Guidance.....	[BSD ACL 2020]
8. Interactive Text Generation.....	[FGPBCZGD EMNLP 2024]
9. Reinforcement Learning with Convex Constraints.....	[MBDDS NeurIPS 2019]
10. Constrained episodic reinforcement learning in concave-convex and knapsack settings.....	[BDLMSSS NeurIPS 2020]
11. Ranking with Long-Term Constraints.....	[BFDJ WSDM 2024]
12. lilGym: Natural Language Visual Reasoning with Reinforcement Learning.....	[WBKA ACL 2022]

# Future Work

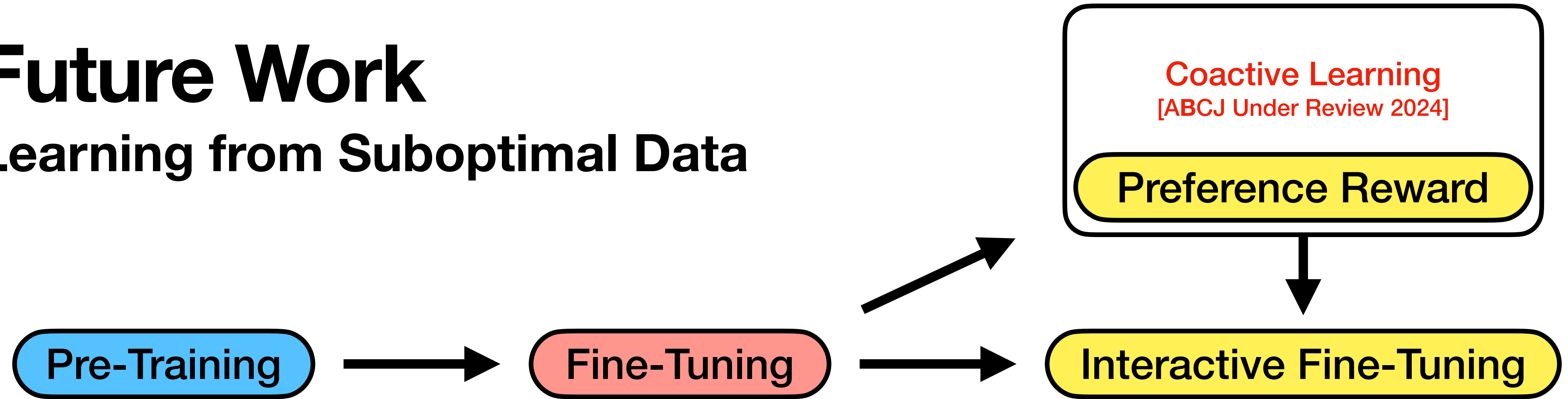
## Learning from Suboptimal Data



- Develop **fine-tuning algorithms** that can learn from suboptimal data.
- Develop **reward functions** that can learn from suboptimal preferences.

# Future Work

## Learning from Suboptimal Data



- Develop **fine-tuning algorithms** that can learn from suboptimal data.
- Develop **reward functions** that can learn from suboptimal preferences.

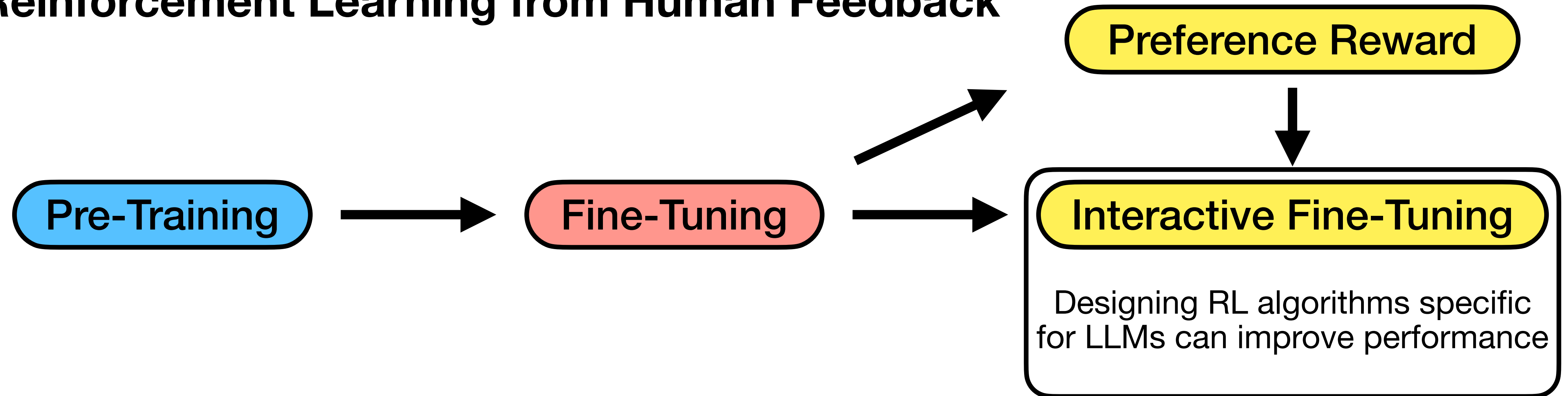
# Future Work

- Reinforcement Learning from Human Feedback [3, 4]
- Constrained Reinforcement Learning [9, 10, 11]

1. Disagreement-Regularized Imitation Learning.....	[BSH ICLR 2020]
2. Adversarial Imitation Learning via Boosting.....	[CSHBS ICLR 2024]
3. Is reinforcement learning (not) for natural language processing?.....	[RABHSBHC ICLR 2023]
4. Learning to Generate Better Than Your LLM.....	[CBRMS Instruction Workshop 2023]
5. Policy-Gradient Training of Language Models for Ranking.....	[GCCBJ FMDM Workshop 2023]
6. Non-Monotonic Text Generation.....	[WBDC ICLR 2019]
7. Active Imitation Learning with Noisy Guidance.....	[BSD ACL 2020]
8. Interactive Text Generation.....	[FGPBCZGD EMNLP 2024]
9. Reinforcement Learning with Convex Constraints.....	[MBDDS NeurIPS 2019]
10. Constrained episodic reinforcement learning in concave-convex and knapsack settings.....	[BDLMSSS NeurIPS 2020]
11. Ranking with Long-Term Constraints.....	[BFDJ WSDM 2024]
12. lilGym: Natural Language Visual Reasoning with Reinforcement Learning.....	[WBKA ACL 2022]

# Future Work

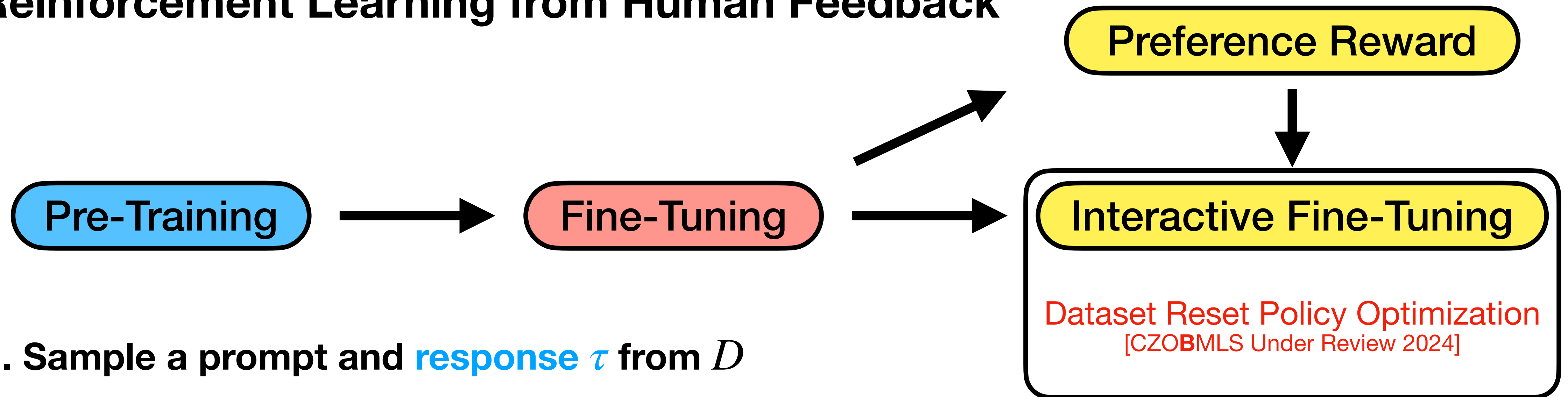
## Reinforcement Learning from Human Feedback



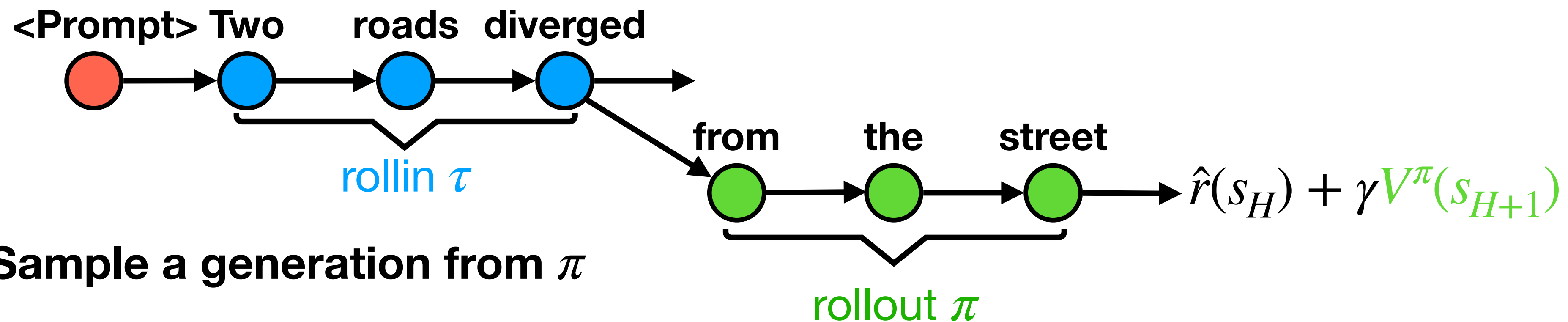
- Continue to **explore RL algorithms** specialized for NLP systems.
  - Fast Simulation
  - Deterministic Environment
  - Demonstration Data
  - Trajectory-Level Rewards

# Future Work

## Reinforcement Learning from Human Feedback



1. Sample a prompt and response  $\tau$  from  $D$



2. Sample a generation from  $\pi$



# Future Work

Natural Language Processing [3, 4, 6, 7, 8, 12] •

Information Retrieval [5, 11] •

Classic Control [2, 9, 10] •

Video Games [1] •

• Imitation Learning [1, 2, 6, 7, 8, 12]

• Reinforcement Learning [3, 4, 5, 9, 10, 11]

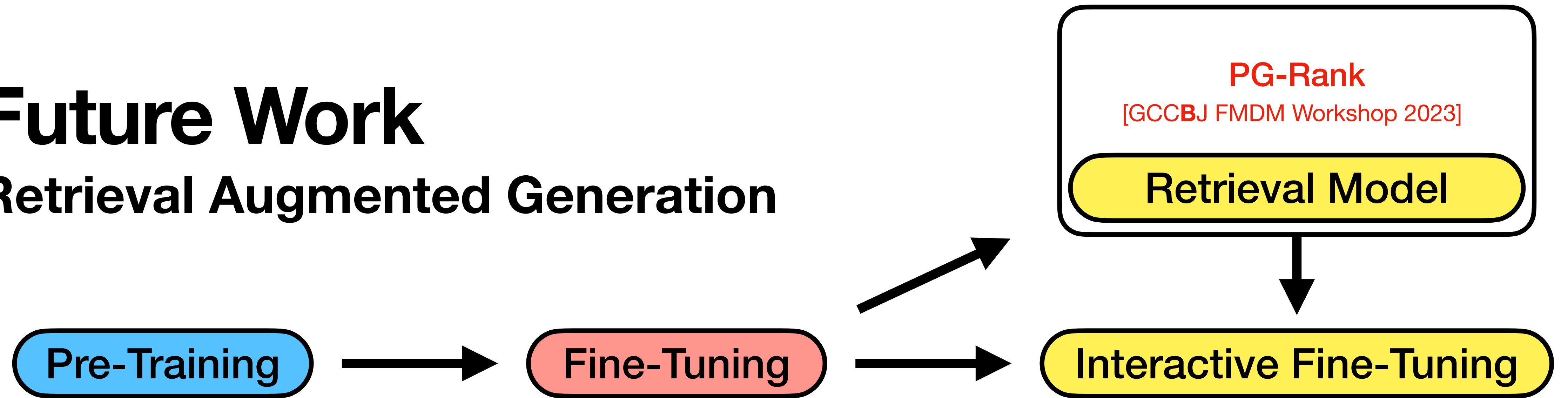
1. Disagreement-Regularized Imitation Learning.....	[BSH ICLR 2020]
2. Adversarial Imitation Learning via Boosting.....	[CSHBS ICLR 2024]
3. Is reinforcement learning (not) for natural language processing?.....	[RABHSBHC ICLR 2023]
4. Learning to Generate Better Than Your LLM.....	[CBRMS Instruction Workshop 2023]
5. Policy-Gradient Training of Language Models for Ranking.....	[GCCBJ FMDM Workshop 2023]
6. Non-Monotonic Text Generation.....	[WBDC ICLR 2019]
7. Active Imitation Learning with Noisy Guidance.....	[BSD ACL 2020]
8. Interactive Text Generation.....	[FGPBCZGD EMNLP 2024]
9. Reinforcement Learning with Convex Constraints.....	[MBDDS NeurIPS 2019]
10. Constrained episodic reinforcement learning in concave-convex and knapsack settings.....	[BDLMSSS NeurIPS 2020]
11. Ranking with Long-Term Constraints.....	[BFDJ WSDM 2024]
12. Successor feature sets: Generalizing successor representations across policies.....	[BMG AAI 2021]

? •

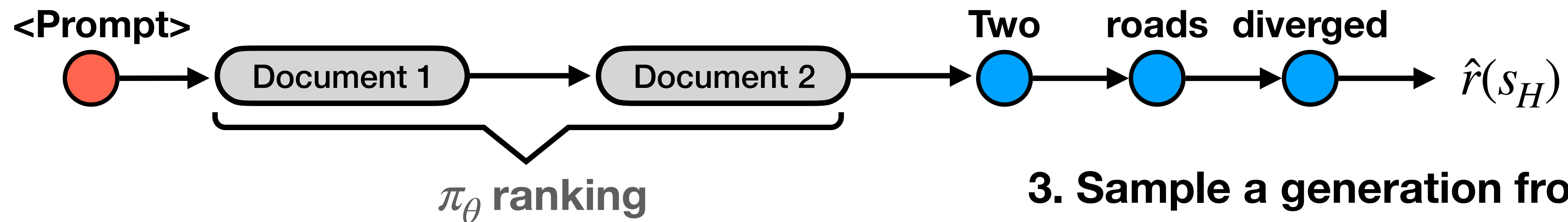


# Future Work

## Retrieval Augmented Generation



### 1. Sample a prompt from $D$

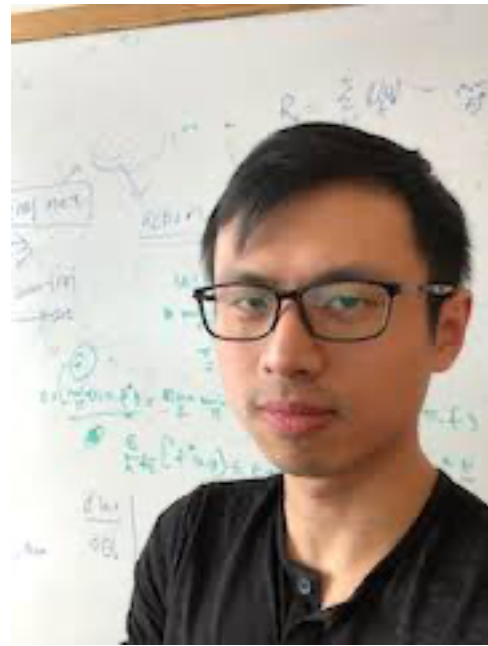


### 3. Sample a generation from $\pi$

### 2. Query a retrieval model



# Thank You!



Wen Sun



Mikael Henaff



Rajkumar Ramamurthy



Dipendra Misra



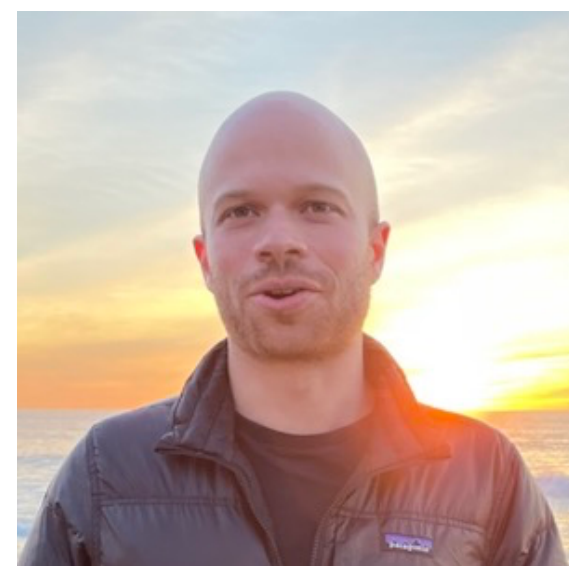
Jonathan D. Chang



Wenhao Zhan



Thorsten Joachims



Aaron Tucker



Adam Cahall



Yoav Artzi



Anne Wu



Jason Lee

# Questions?