

# **Interactive Imitation Learning**

# **Announcement**

## **This Thursday:**

lecture will start at 9:50am  
and office hour will end at 11:15am

**Recap**

**Offline IL**

# Recap

## Offline IL

Ground truth reward  $r(s, a) \in [0, 1]$  is unknown;  
assume expert is a near optimal policy  $\pi^\star$

# Recap

## Offline IL

Ground truth reward  $r(s, a) \in [0, 1]$  is unknown;  
assume expert is a near optimal policy  $\pi^\star$

We have a dataset  $\mathcal{D} = (s_i^\star, a_i^\star)_{i=1}^M, s_i^\star \sim d_\mu^{\pi^\star}, a_i^\star \sim \pi^\star(\cdot | s_i^\star)$

# Recap

## The Behavior Cloning algorithm:

Choose regression (for continuous action) or classification loss  $\ell(\pi(s), a)$ , and perform SL:

$$\hat{\pi} = \min_{\pi \in \Pi} \sum_{i=1}^M \ell(\pi^*(s_i^*), a_i^*)$$

# Recap

## The Behavior Cloning algorithm:

Choose regression (for continuous action) or classification loss  $\ell(\pi(s), a)$ , and perform SL:

$$\hat{\pi} = \min_{\pi \in \Pi} \sum_{i=1}^M \ell(\pi^*(s_i^*), a_i^*)$$

### Pros:

Simple, flexible, and often just works reasonably well

# Recap

## The Behavior Cloning algorithm:

Choose regression (for continuous action) or classification loss  $\ell(\pi(s), a)$ , and perform SL:

$$\hat{\pi} = \min_{\pi \in \Pi} \sum_{i=1}^M \ell(\pi^*(s_i^*), a_i^*)$$

### Pros:

Simple, flexible, and often just works reasonably well

### Cons:

Distribution shift issue:  $\hat{\pi}$  does not know what to do outside expert's states



## **Question for today:**

How to mitigate the distribution shift issue?

**Solution:**

**Interactive Imitation Learning Setting**

# **Solution:**

## **Interactive Imitation Learning Setting**

### **Key assumption:**

**we can query expert  $\pi^\star$  at any time and any state during training**

# Solution:

## Interactive Imitation Learning Setting

### Key assumption:

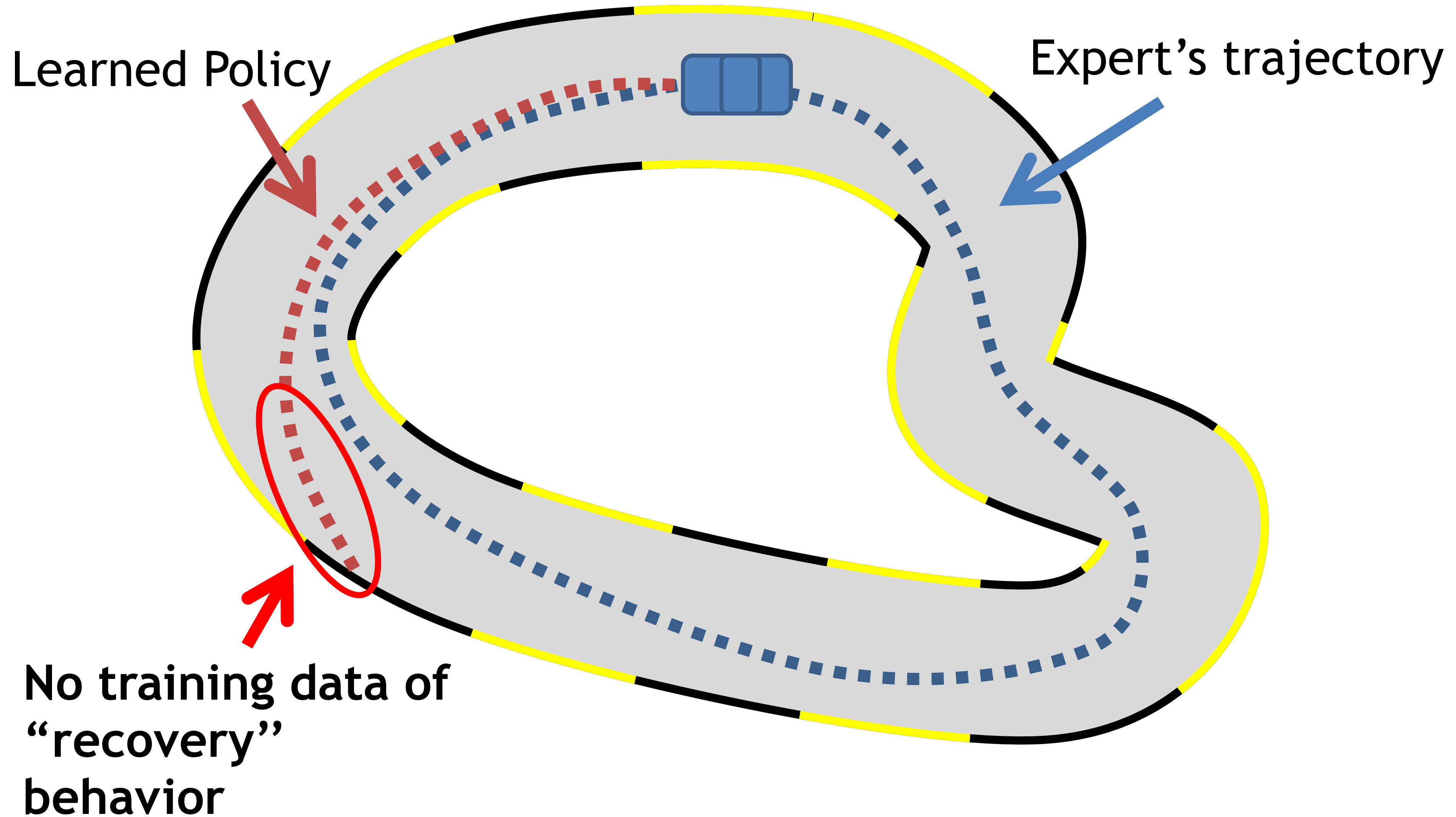
we can query expert  $\pi^\star$  at any time and any state during training

(Recall that previously we only had an offline dataset  $\mathcal{D} = (s_i^\star, a_i^\star)_{i=1}^M \sim d_\mu^{\pi^\star}$ )

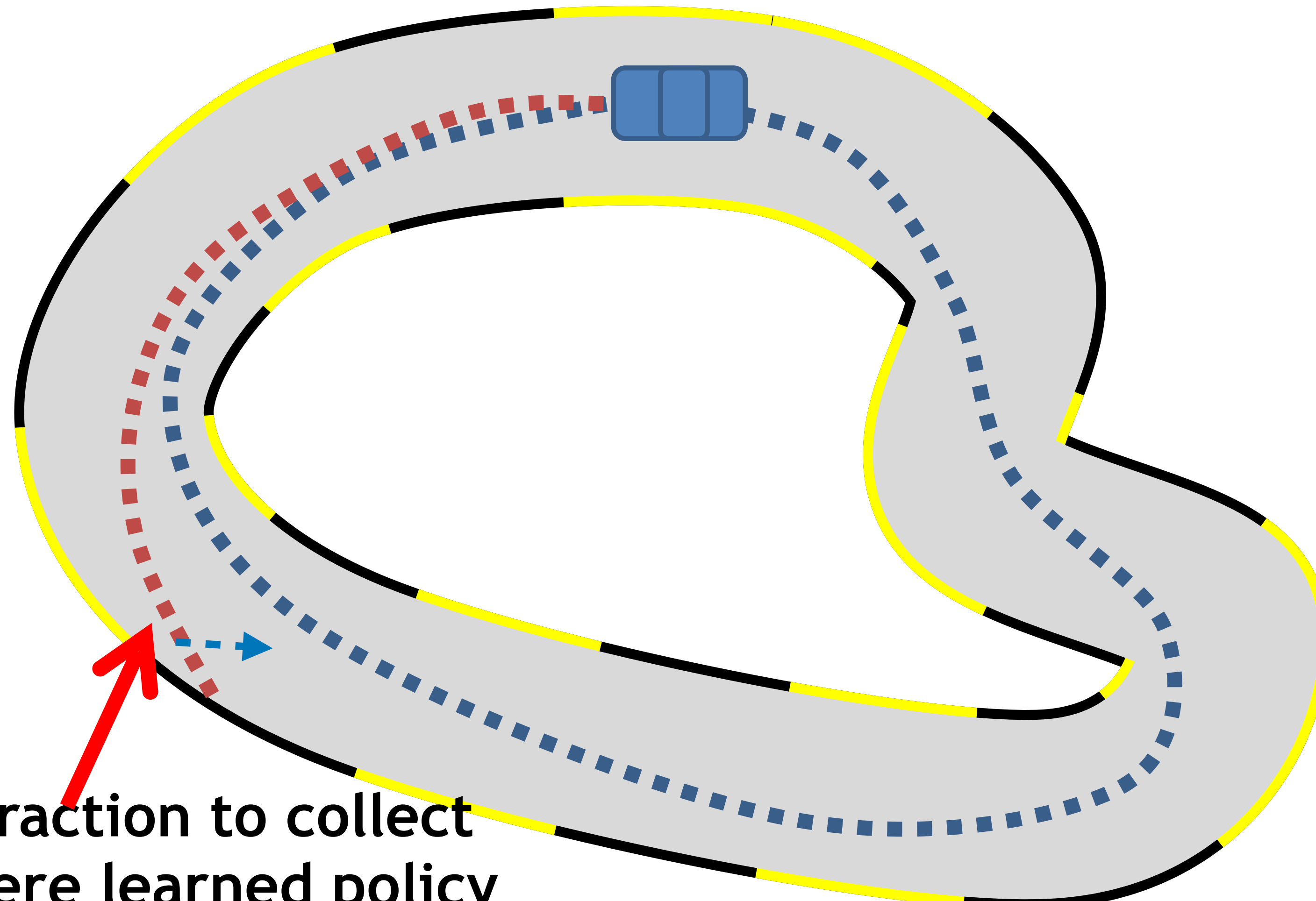
## **Outline for today:**

1. The DAgger (Data Aggregation) Algorithm
2. Analysis of DAgger: a quick intro to Online Learning

# Recall the Main Problem from Behavior Cloning:

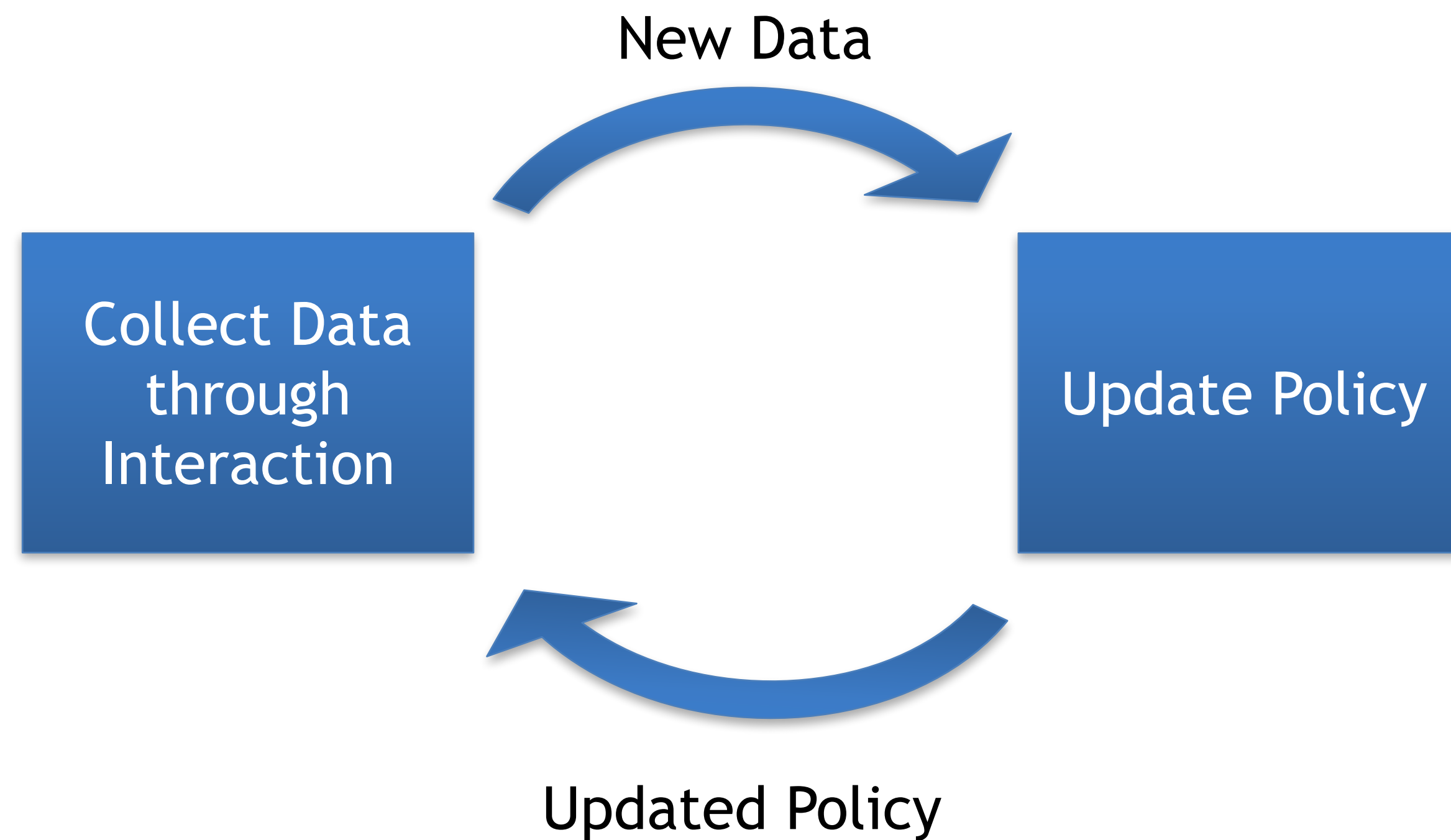


# Intuitive solution: Interaction



Use interaction to collect data where learned policy goes

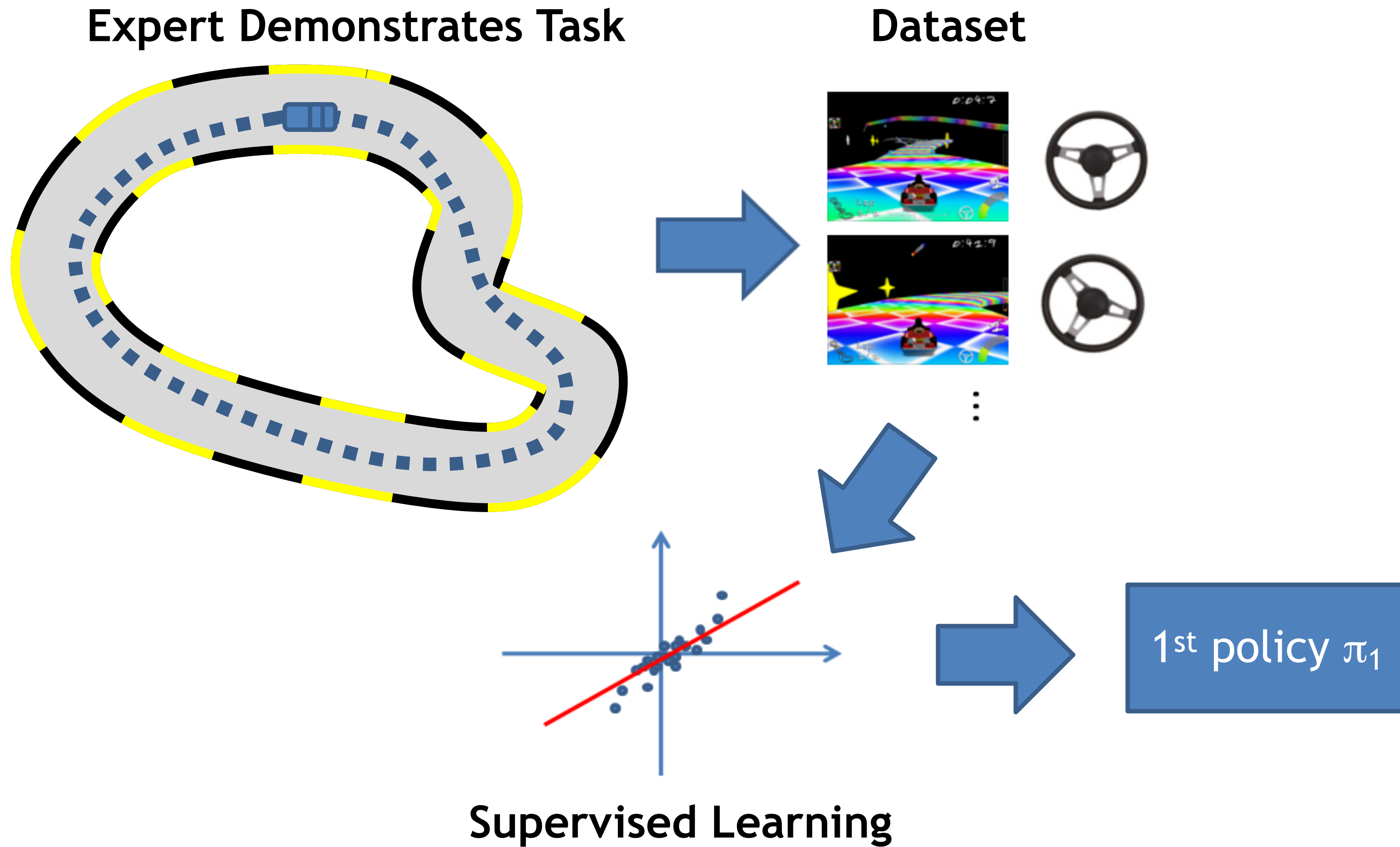
# General Idea: Iterative Interactive Approach





# Dagger: Dataset Aggregation <sup>[Ross11a]</sup>

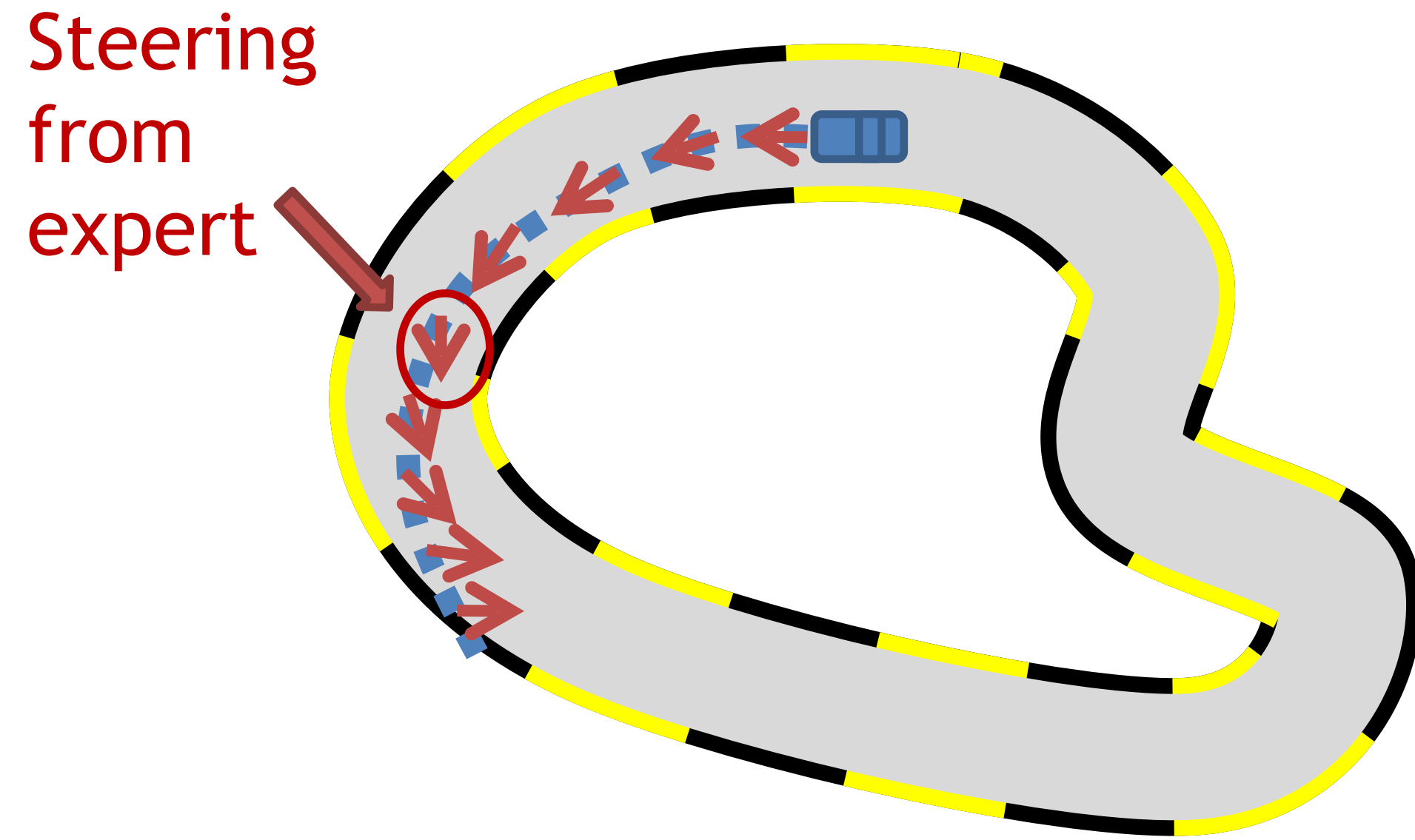
0th iteration



# Dagger: Dataset Aggregation <sup>[Ross11a]</sup>

## 1st iteration

Execute  $\pi_1$  and Query Expert

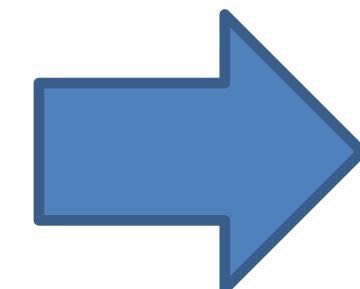
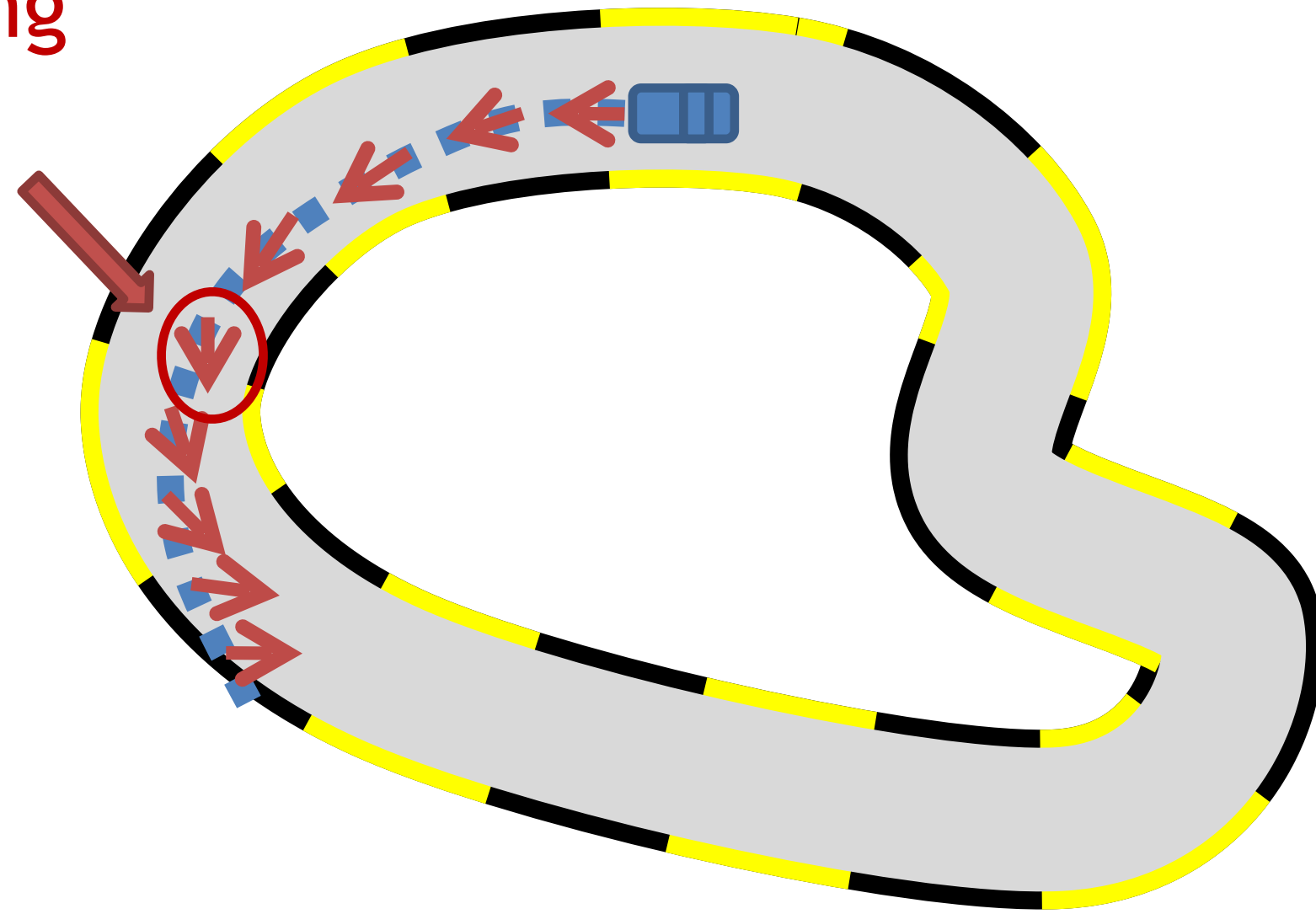


# Dagger: Dataset Aggregation [Ross11a]

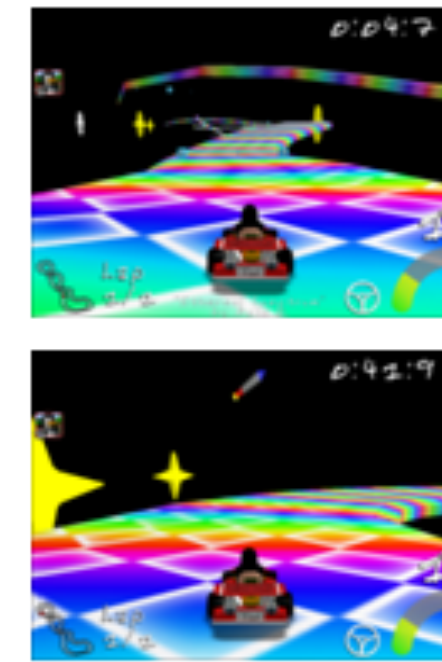
1st iteration

Execute  $\pi_1$  and Query Expert

Steering  
from  
expert



New Data



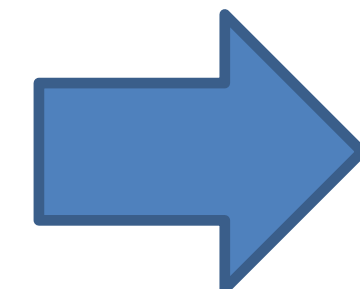
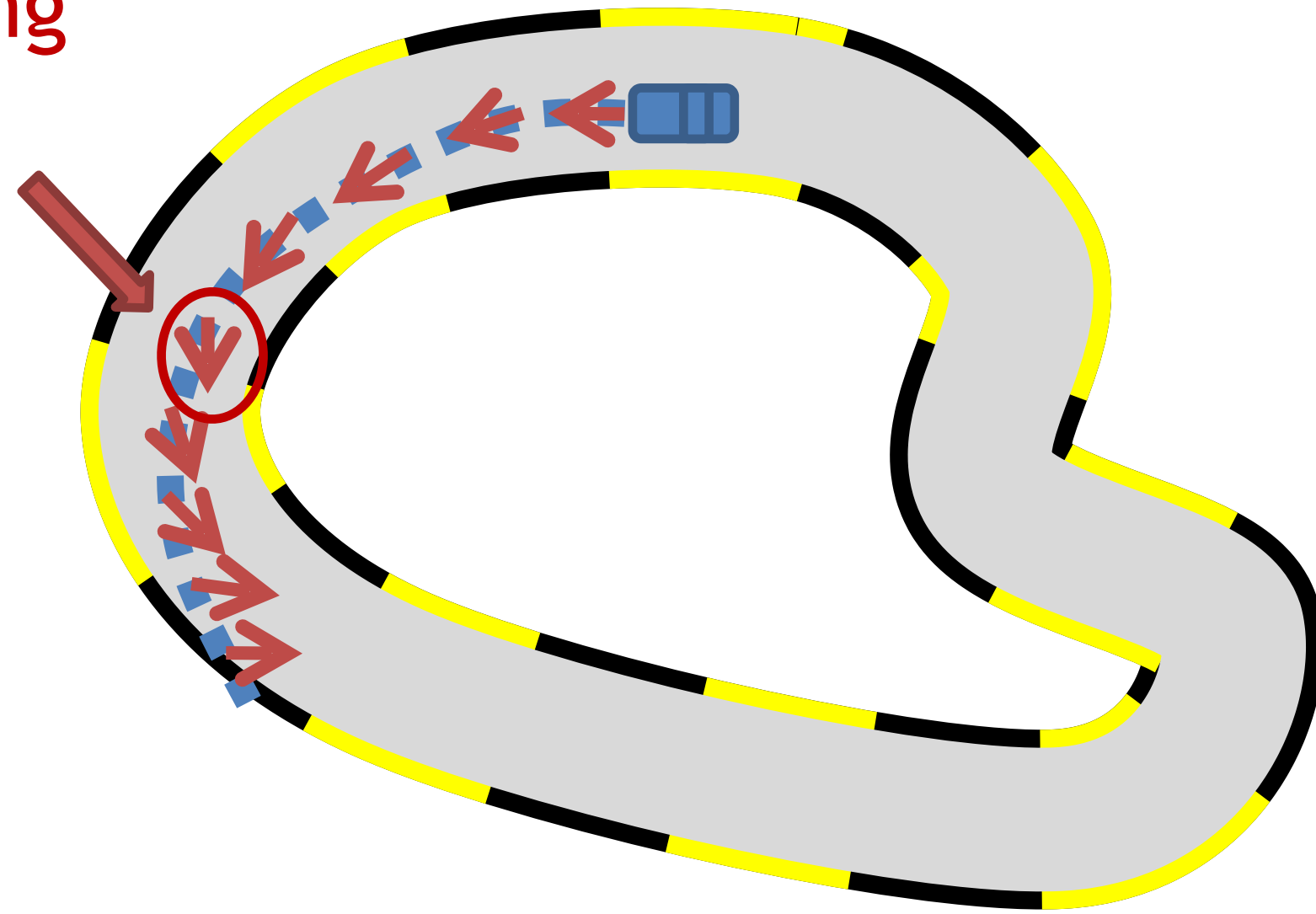
⋮

# Dagger: Dataset Aggregation <sup>[Ross11a]</sup>

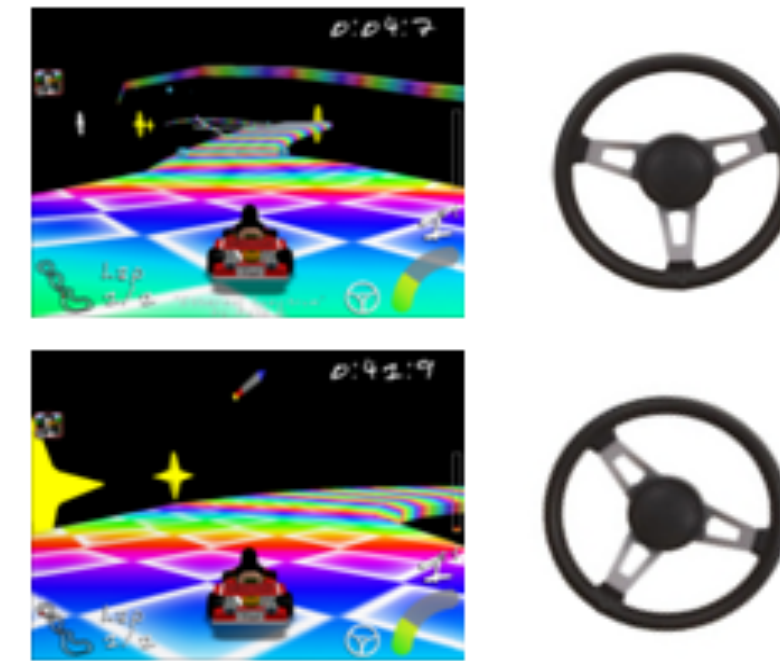
1st iteration

Execute  $\pi_1$  and Query Expert

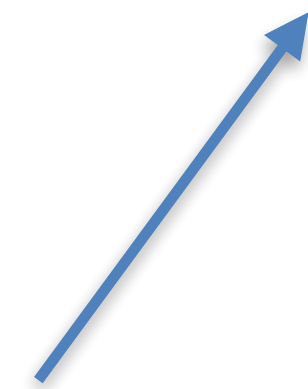
Steering  
from  
expert



New Data



States from  
the learned policy

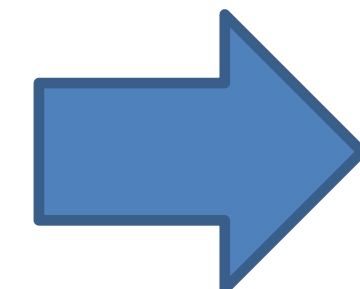
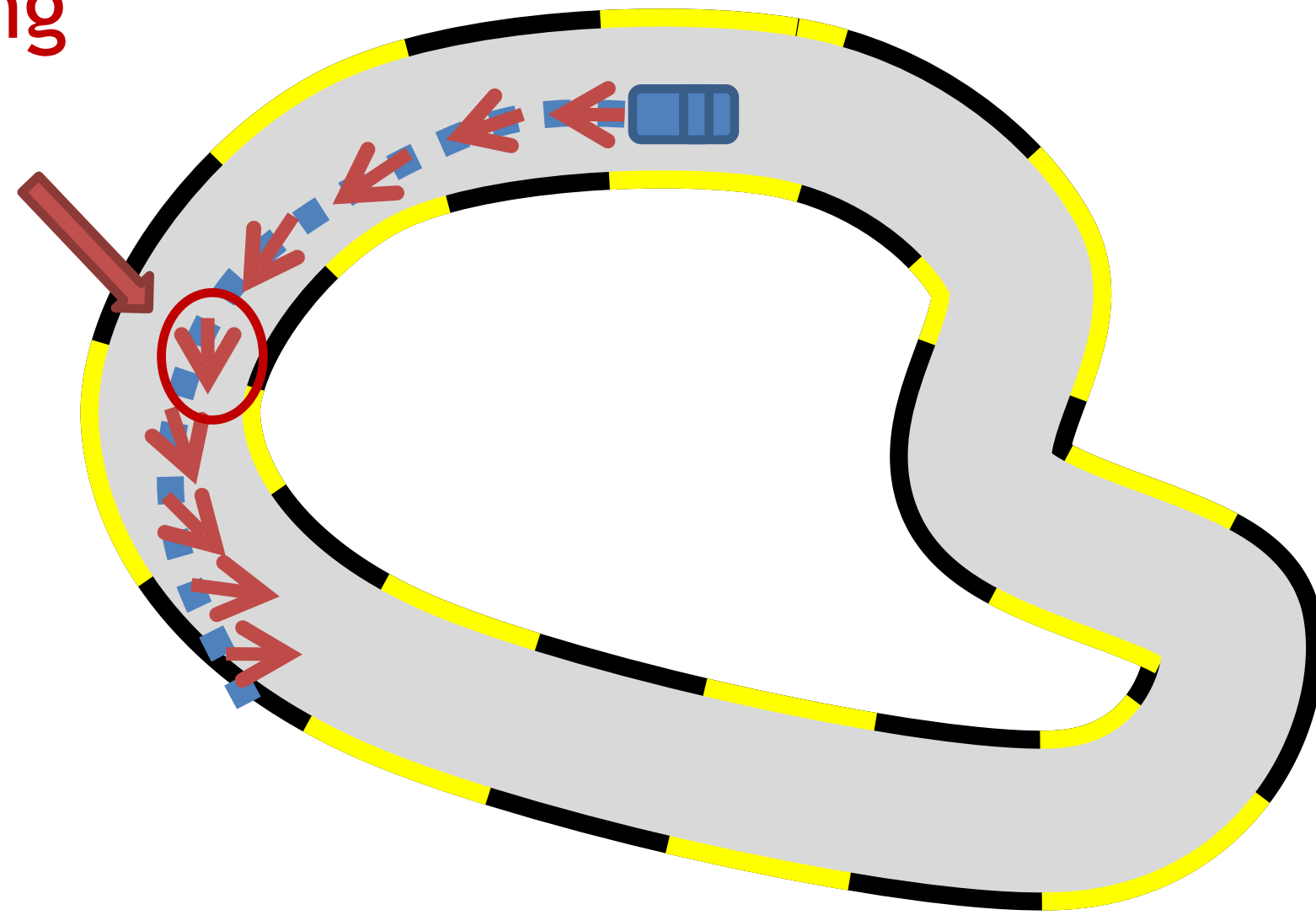


# Dagger: Dataset Aggregation <sup>[Ross11a]</sup>

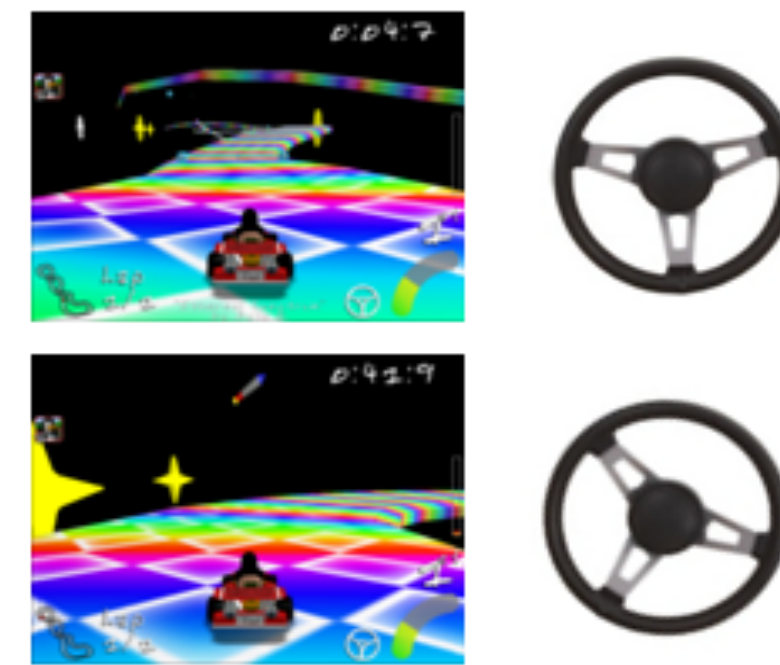
1st iteration

Execute  $\pi_1$  and Query Expert

Steering  
from  
expert



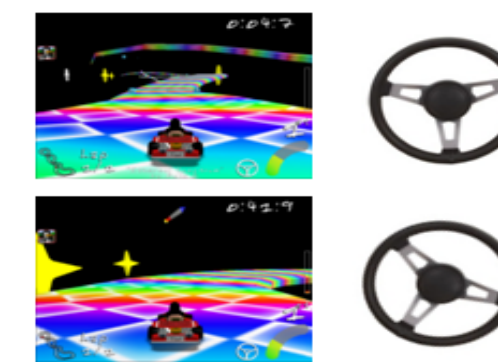
New Data



⋮



All previous data



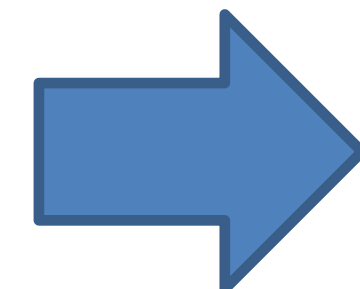
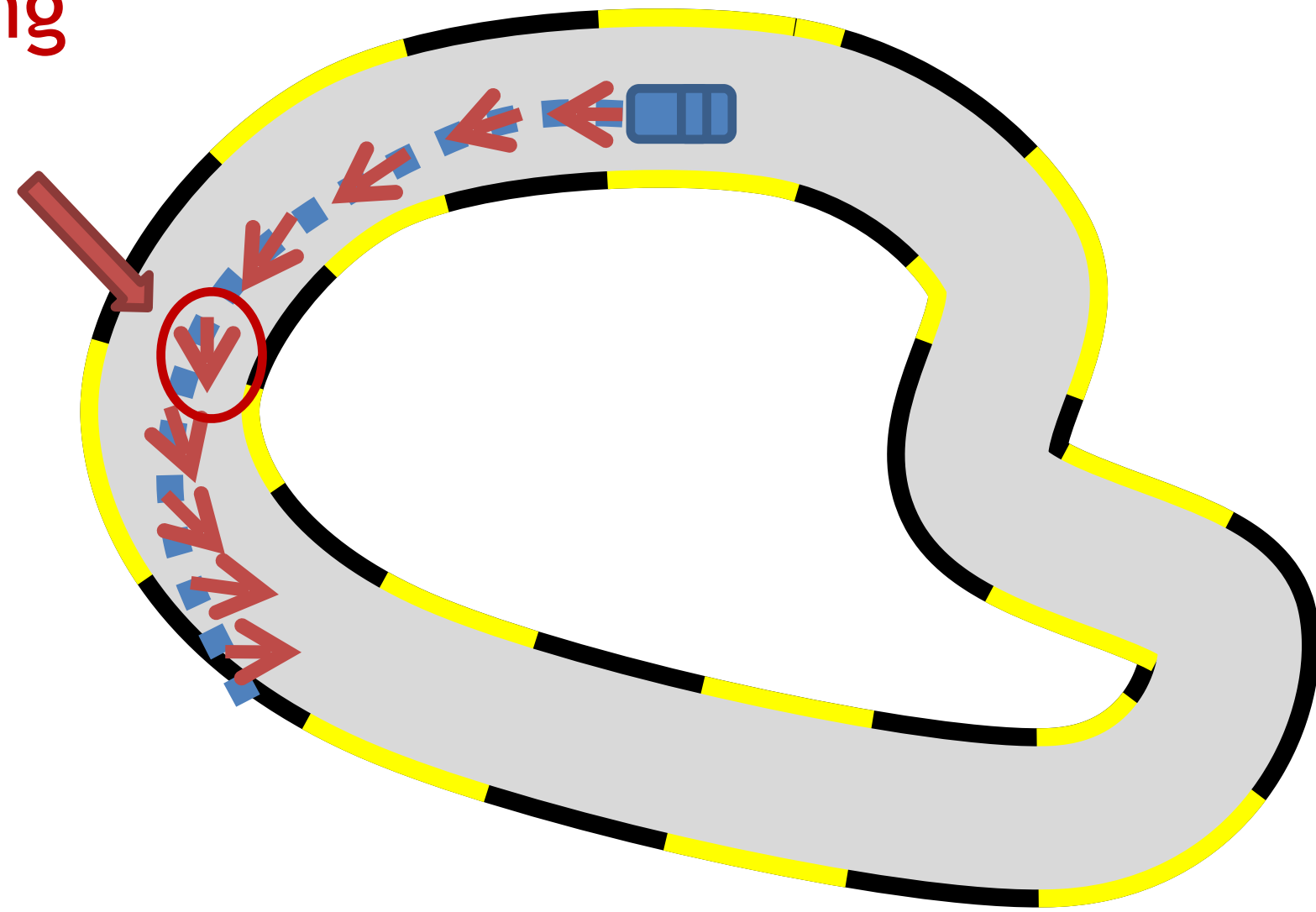
⋮

# Dagger: Dataset Aggregation [Ross11a]

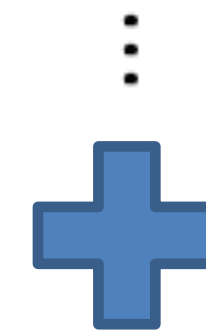
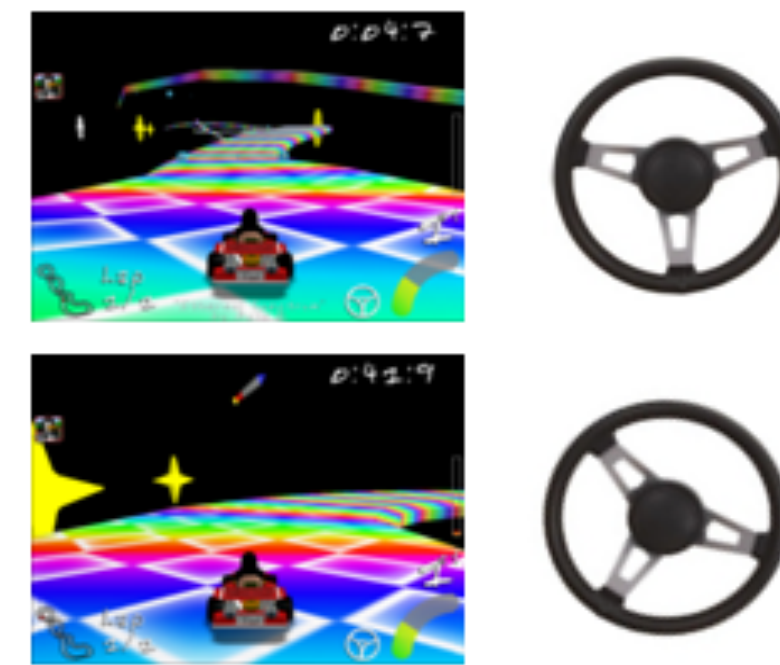
1st iteration

Execute  $\pi_1$  and Query Expert

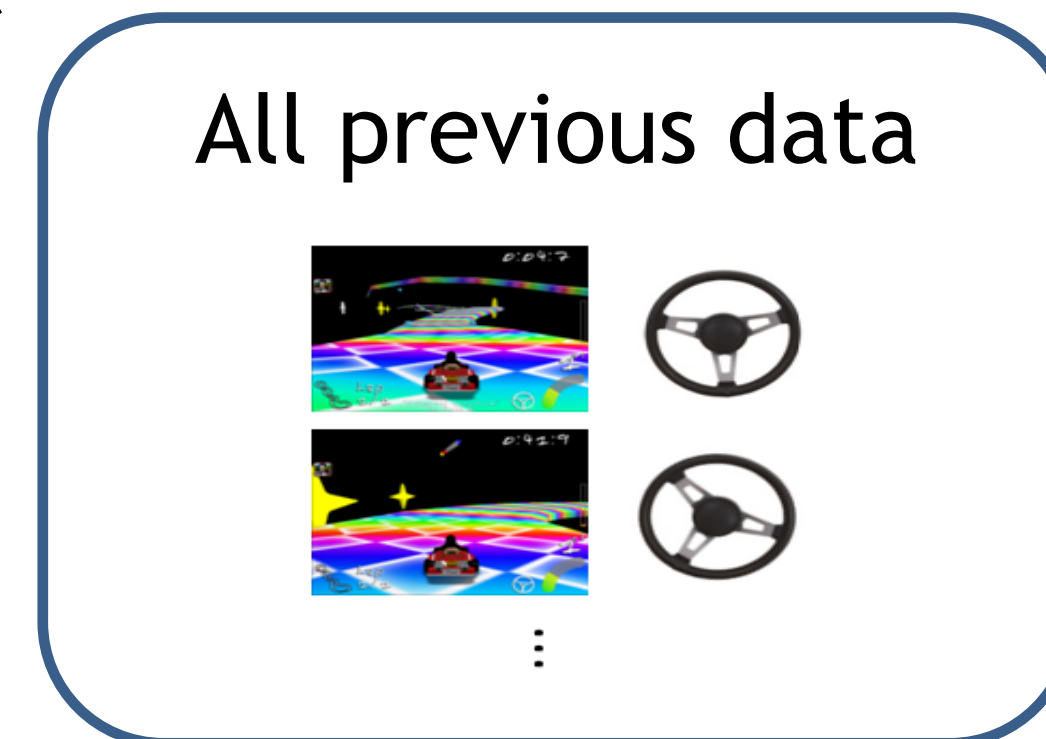
Steering  
from  
expert



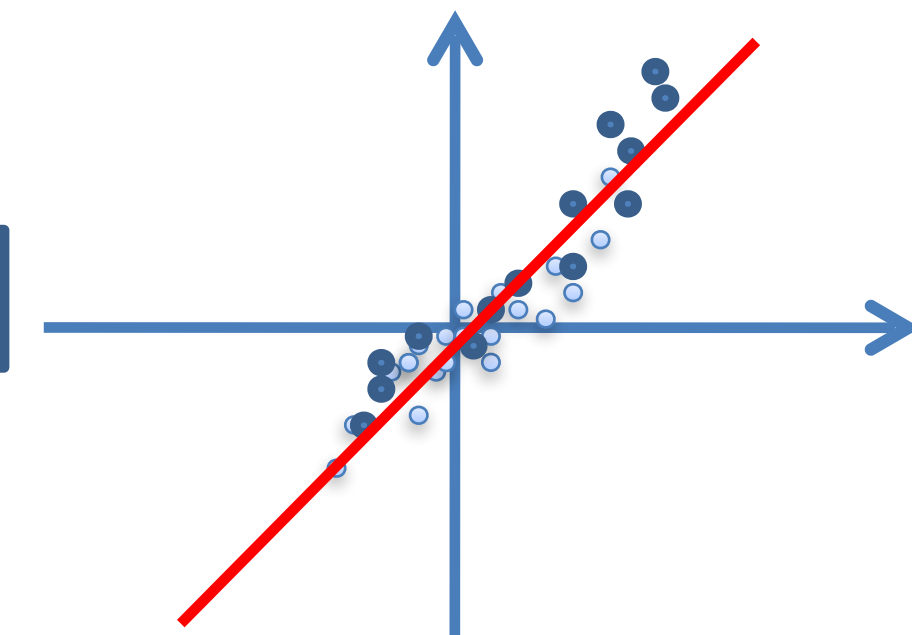
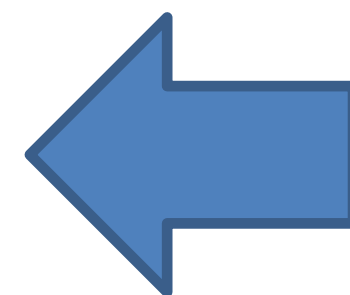
New Data



Aggregate  
Dataset



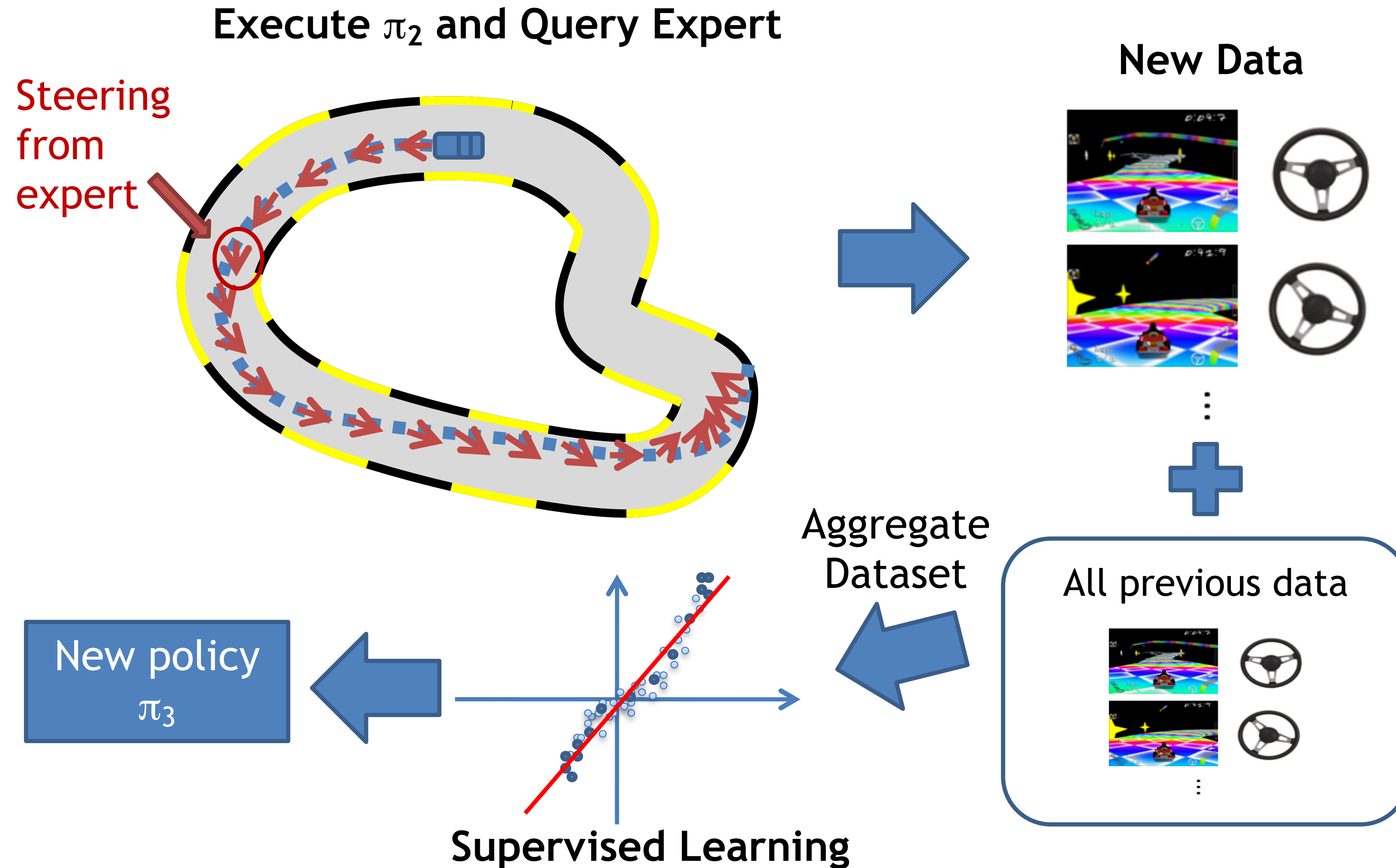
New policy  
 $\pi_2$



Supervised Learning

# Dagger: Dataset Aggregation [Ross11a]

## 2nd iteration

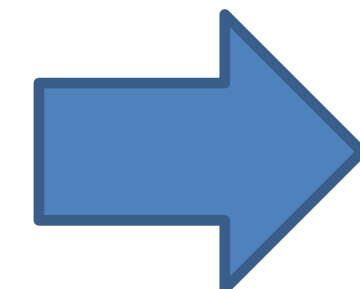
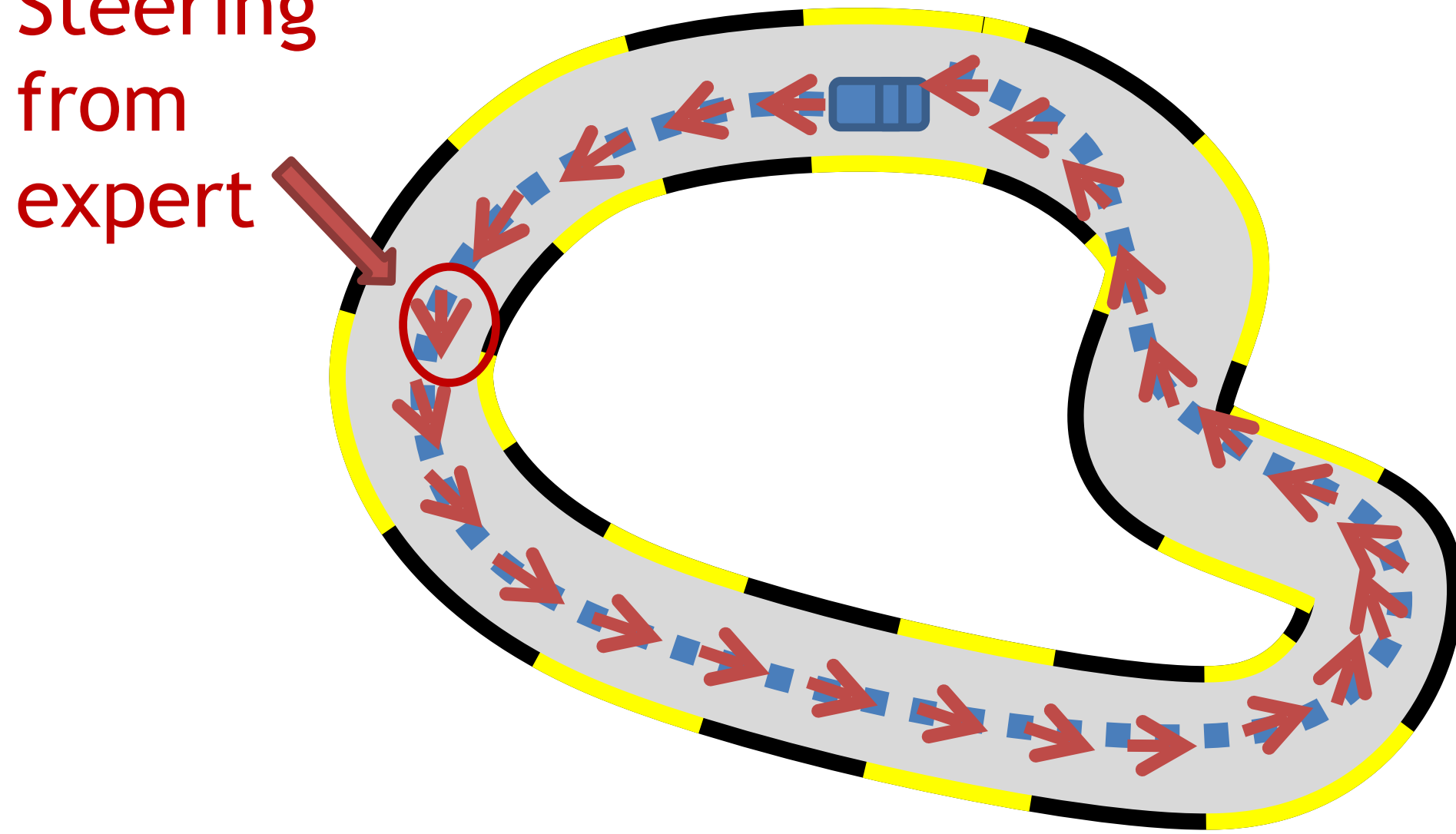


# Dagger: Dataset Aggregation [Ross11a]

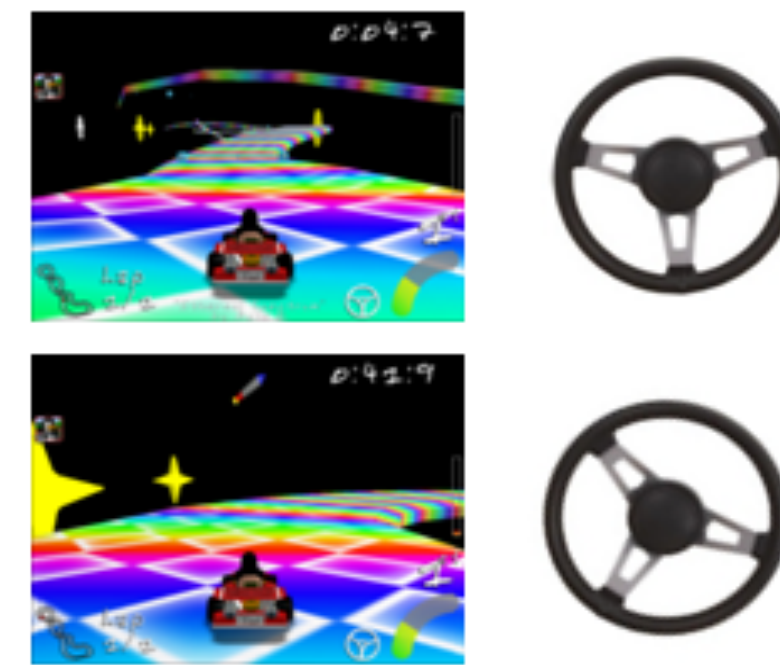
$n^{\text{th}}$  iteration

Execute  $\pi_{n-1}$  and Query Expert

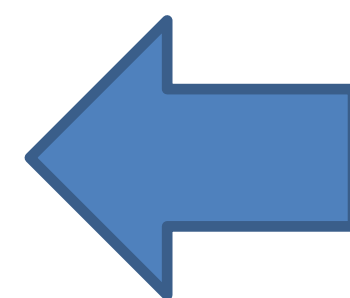
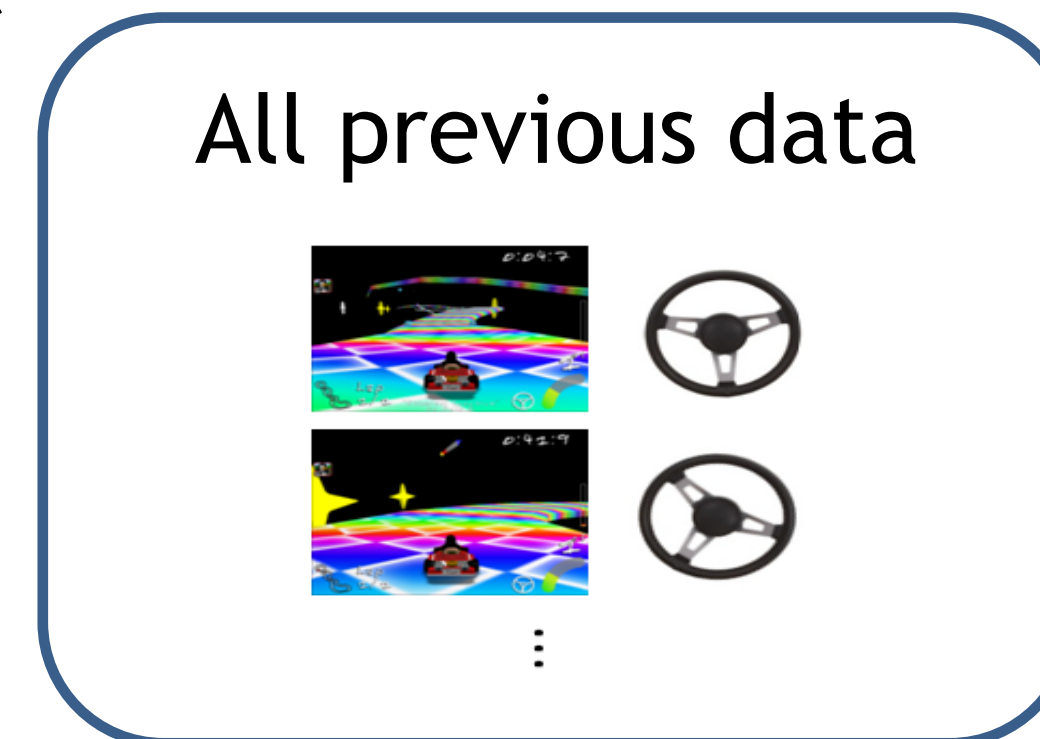
Steering  
from  
expert



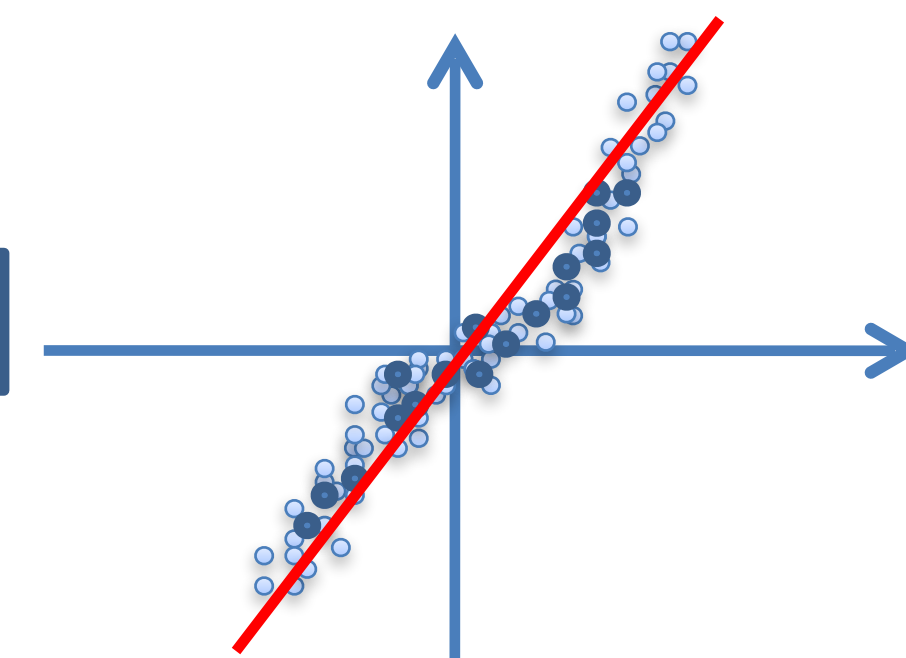
New Data



Aggregate  
Dataset



New policy  
 $\pi_n$



Supervised Learning



# Success!

[Ross AISTATS 2011]



# Success!

[Ross AISTATS 2011]

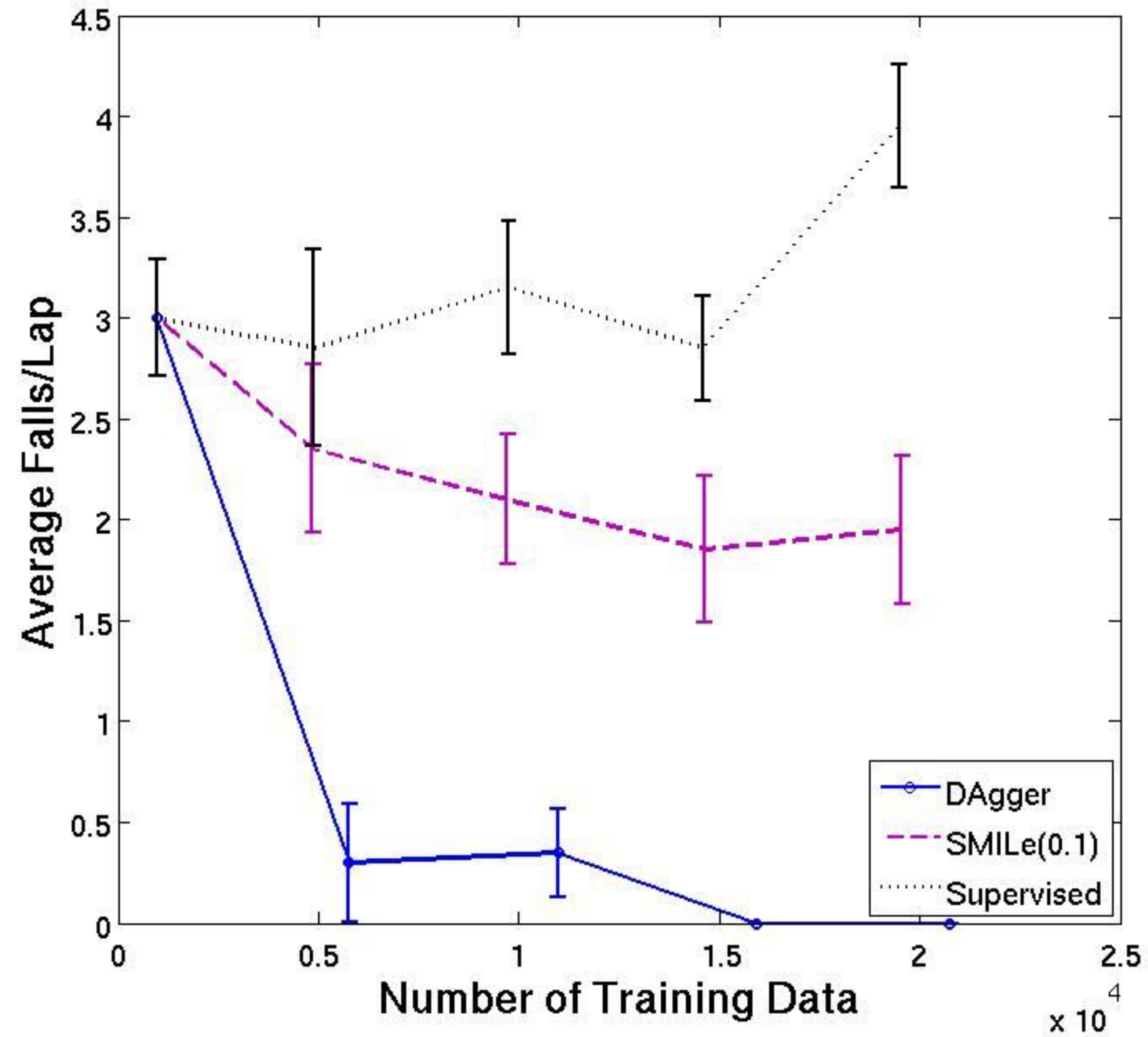
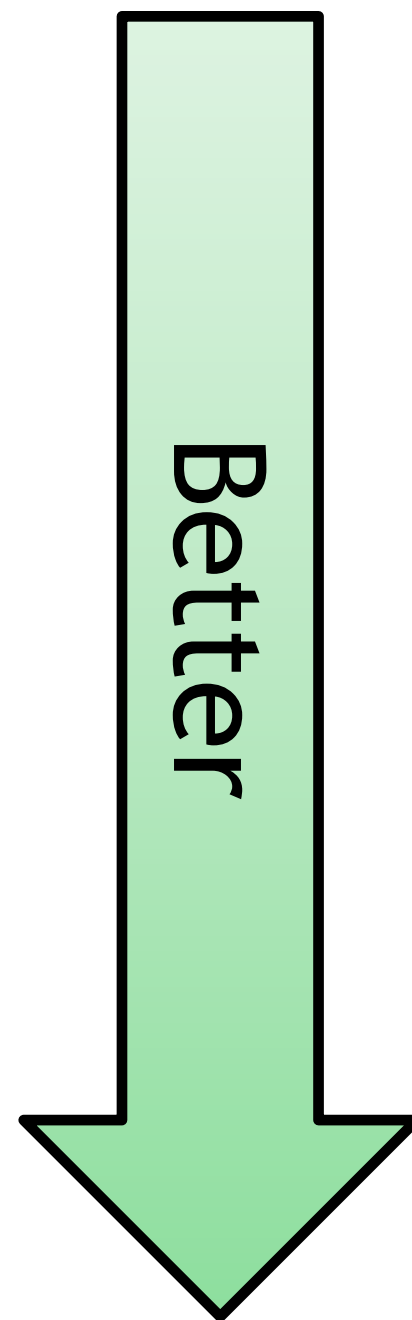


# Success!

[Ross AISTATS 2011]



# Average Falls/Lap



FPS: 24

Attempt: 1 of 1

AgentLinear

Selected Actions:

RIGHT

SPEED

FPS: 24

Attempt: 1 of 1

AgentLinear

Selected Actions:

RIGHT

SPEED

FPS: 24

Attempt: 1 of 1

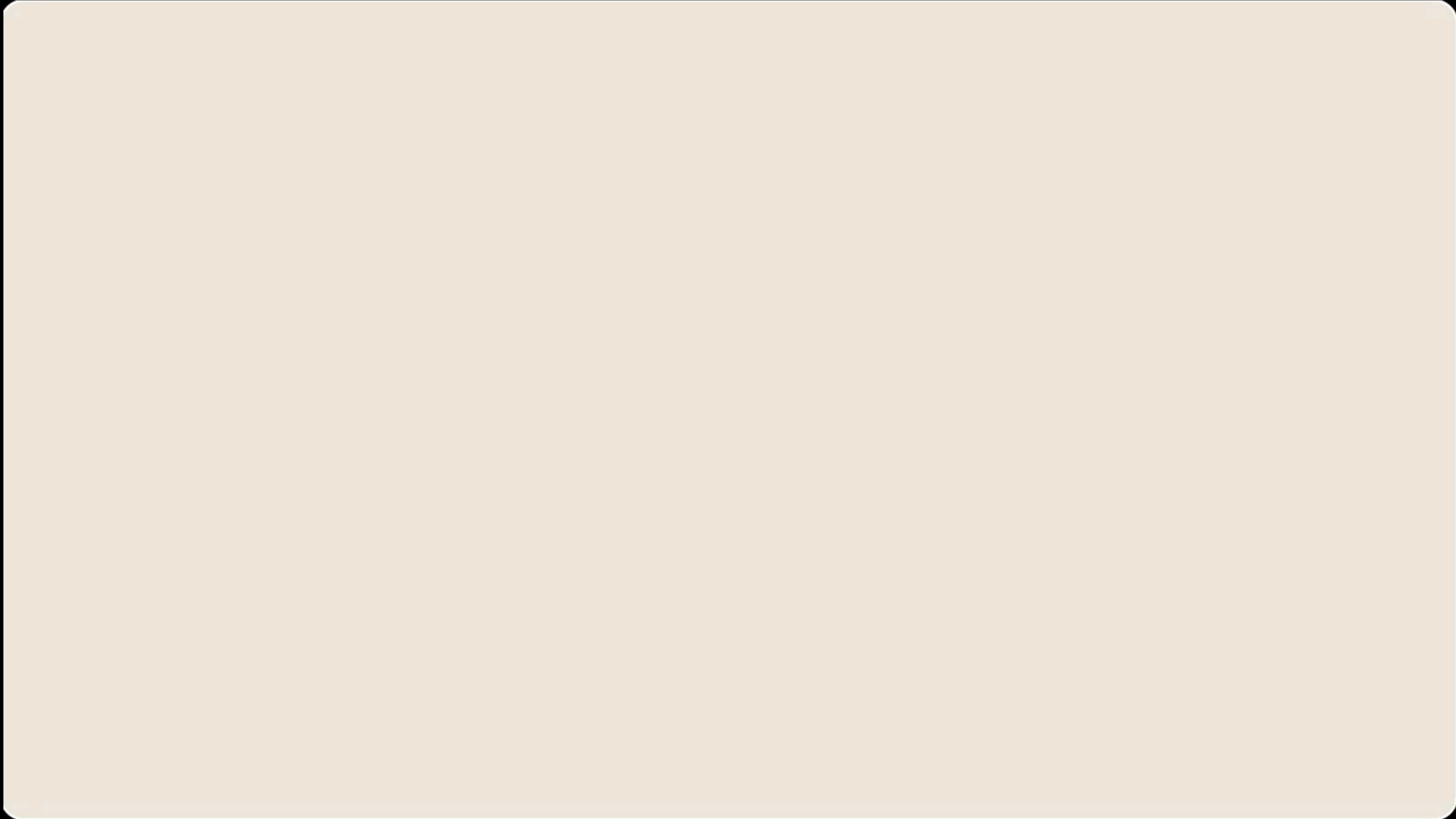
AgentLinear

Selected Actions:

RIGHT

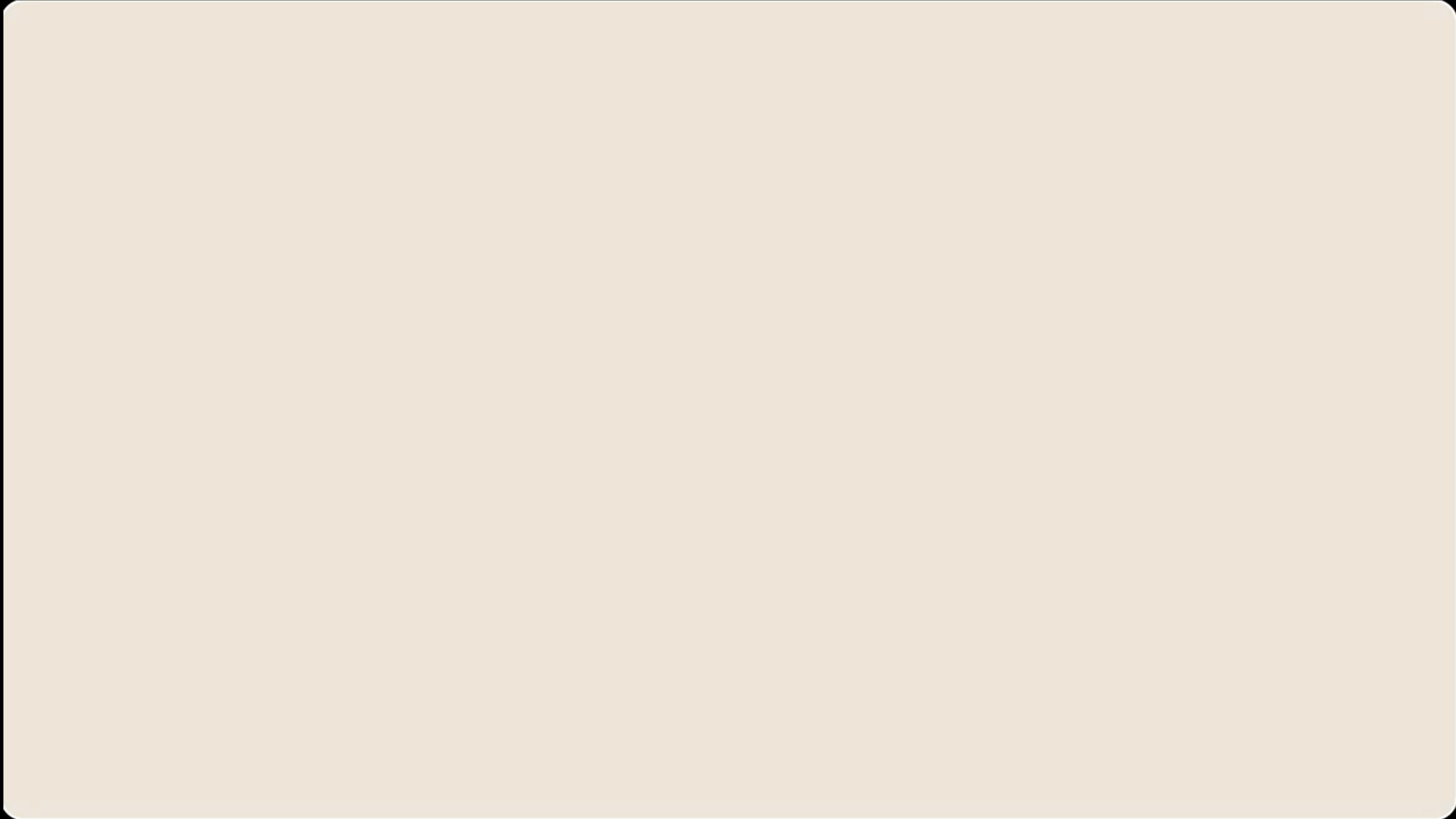
SPEED

# More fun than Video Games...

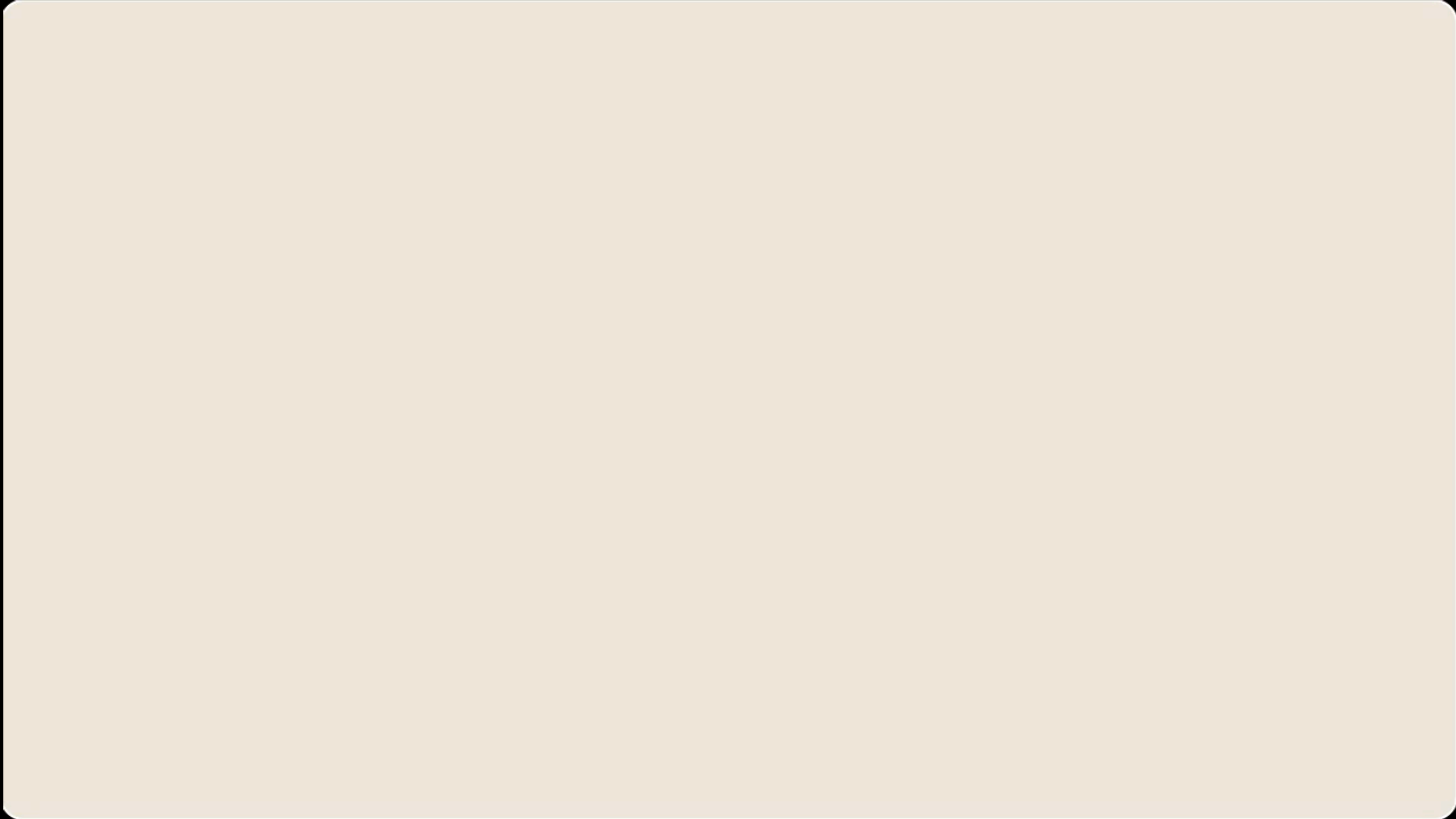




# More fun than Video Games...



# More fun than Video Games...



# Forms of the Interactive Experts

Interactive Expert is expensive, especially when the expert is human...

But expert does not have to be human...

# Forms of the Interactive Experts

Interactive Expert is expensive, especially when the expert is human...

But expert does not have to be human...

## Example: high-speed off-road driving

[Pan et al, RSS 18, Best System Paper]



Fig. 4: The AutoRally car and the test track.

# Forms of the Interactive Experts

Interactive Expert is expensive, especially when the expert is human...

But expert does not have to be human...

## Example: high-speed off-road driving

[Pan et al, RSS 18, Best System Paper]

Goal: learn a racing control policy that maps from data on cheap on-board sensors (raw-pixel image) to low-level control (steer and throttle)



Fig. 4: The AutoRally car and the test track.

# Forms of the Interactive Experts

Interactive Expert is expensive, especially when the expert is human...

But expert does not have to be human...

## Example: high-speed off-road driving

[Pan et al, RSS 18, Best System Paper]



Fig. 4: The AutoRally car and the test track.

Goal: learn a racing control policy that maps from data on cheap on-board sensors (raw-pixel image) to low-level control (steer and throttle)



→ Steering + throttle

(a) raw image

# Forms of the Interactive Experts

**Example: high-speed off-road driving**  
[Pan et al, RSS 18, Best System Paper]

# Forms of the Interactive Experts

**Example: high-speed off-road driving**

[Pan et al, RSS 18, Best System Paper]

Their Setup:

At Training, we have expensive sensors for accurate state estimation and we have computation resources for **MPC** (i.e., high-frequency replanning)



# Forms of the Interactive Experts

**Example: high-speed off-road driving**

[Pan et al, RSS 18, Best System Paper]

Their Setup:

At Training, we have expensive sensors for accurate state estimation and we have computation resources for **MPC** (i.e., high-frequency replanning)

**The MPC is the expert in this case!**

# Forms of the Interactive Experts

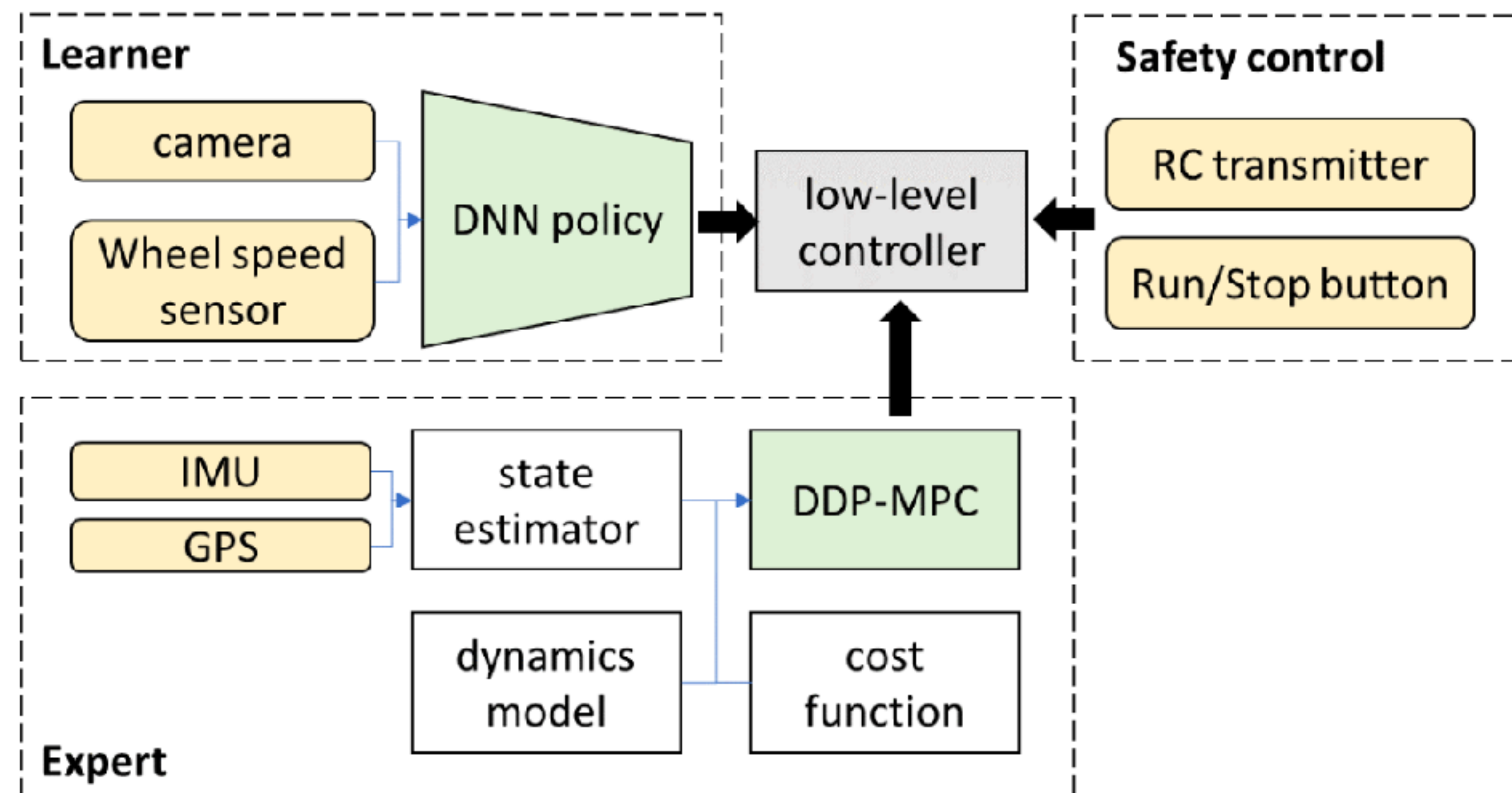
## Example: high-speed off-road driving

[Pan et al, RSS 18, Best System Paper]

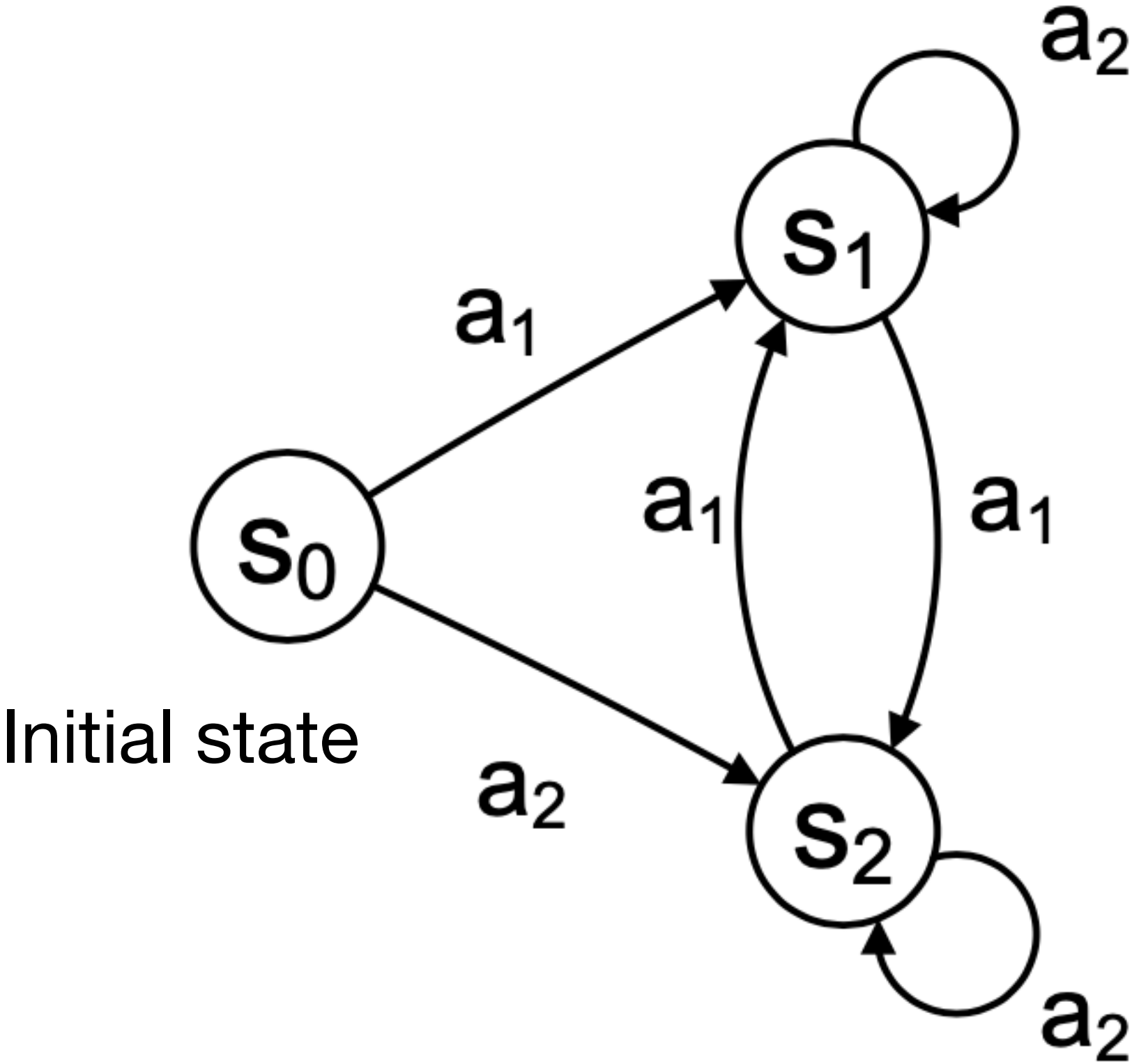
Their Setup:

At Training, we have expensive sensors for accurate state estimation and we have computation resources for **MPC** (i.e., high-frequency replanning)

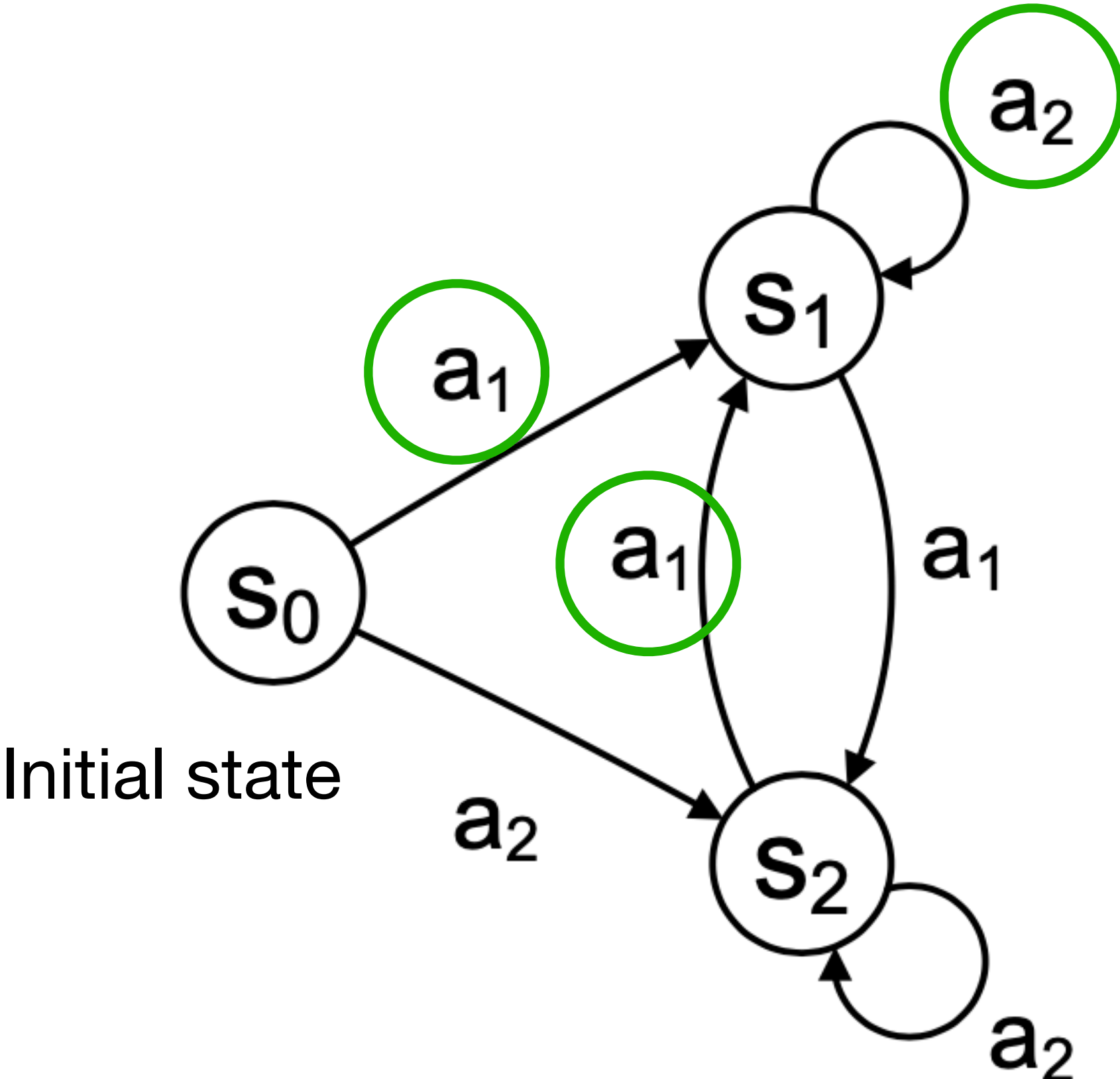
The MPC is the expert in this case!



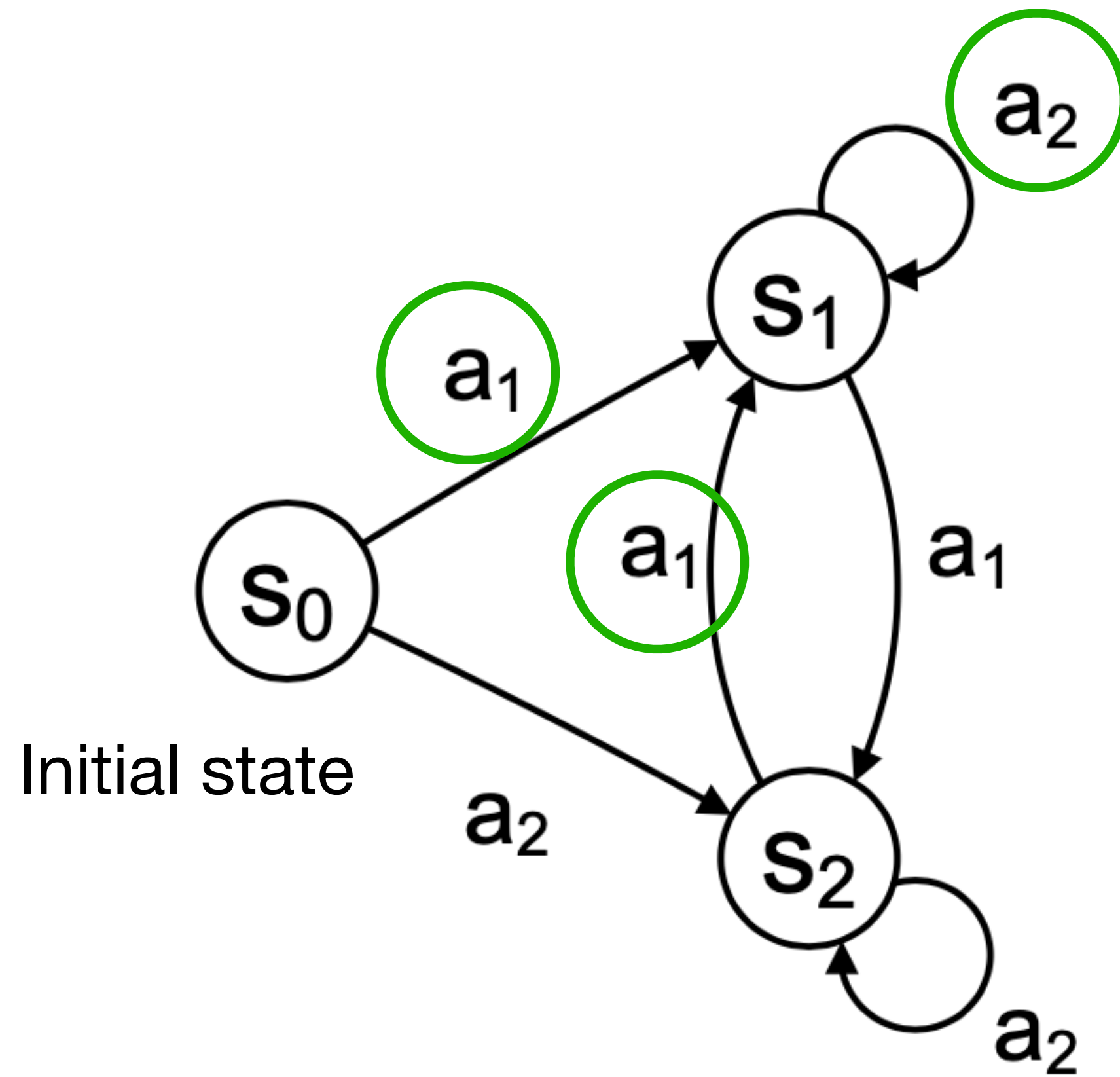
# Distribution Shift: Example



# Distribution Shift: Example

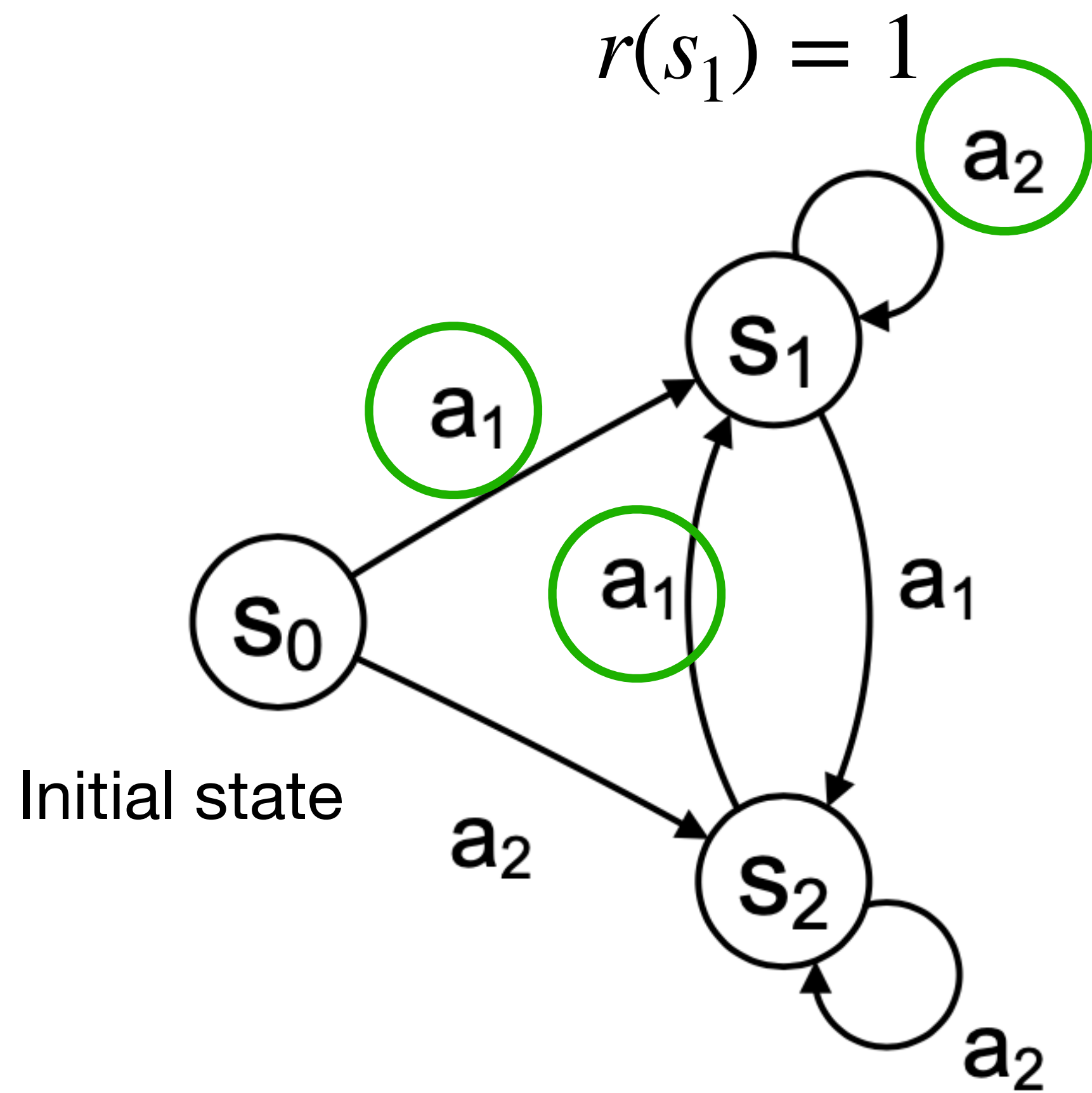


# Distribution Shift: Example



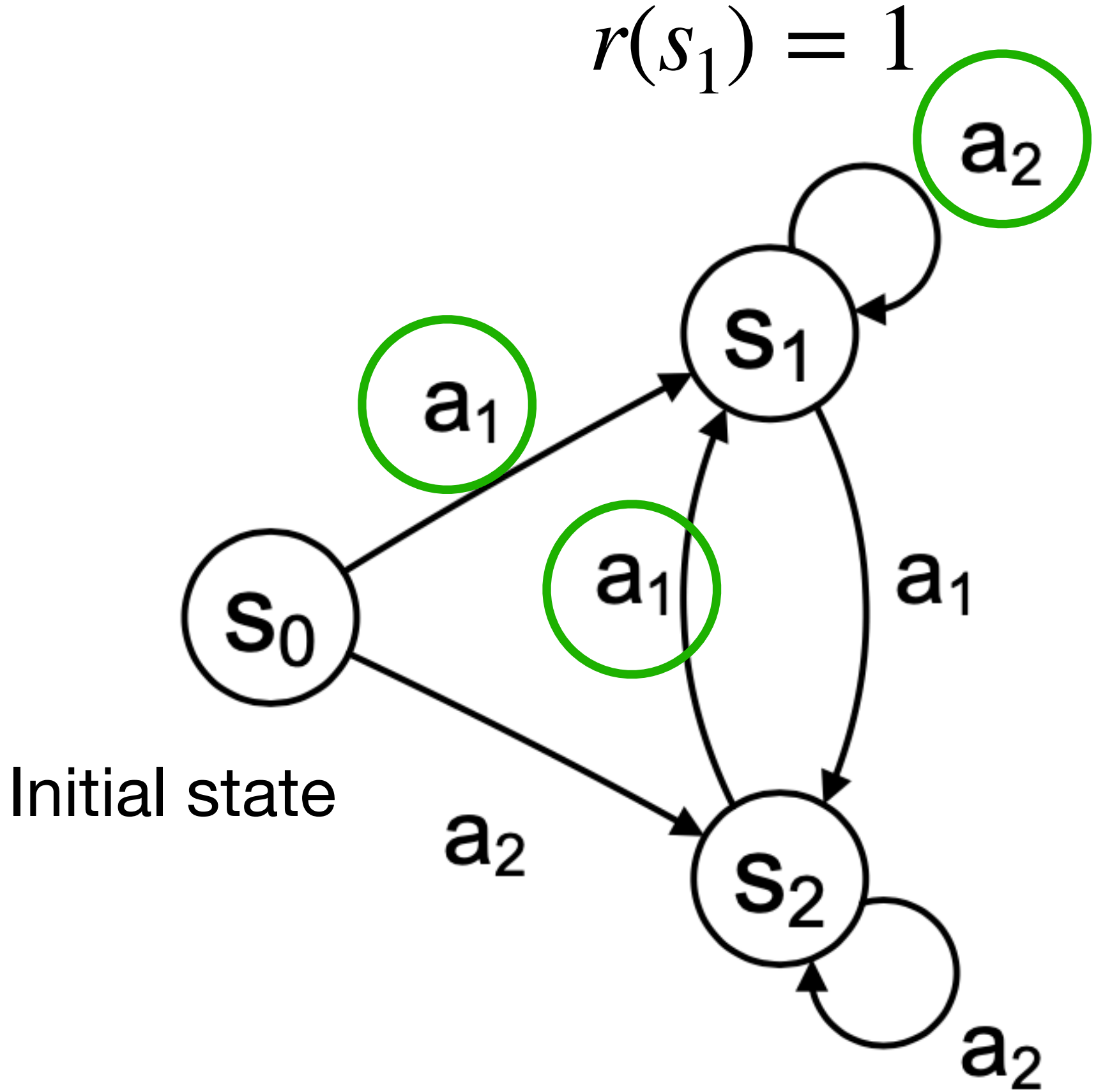
$$d_{s_0}^{\pi^*}(s_0) = 1 - \gamma, d_{s_0}^{\pi^*}(s_1) = \gamma, d_{s_0}^{\pi^*}(s_2) = 0$$

# Distribution Shift: Example



$$d_{s_0}^{\pi^*}(s_0) = 1 - \gamma, d_{s_0}^{\pi^*}(s_1) = \gamma, d_{s_0}^{\pi^*}(s_2) = 0$$

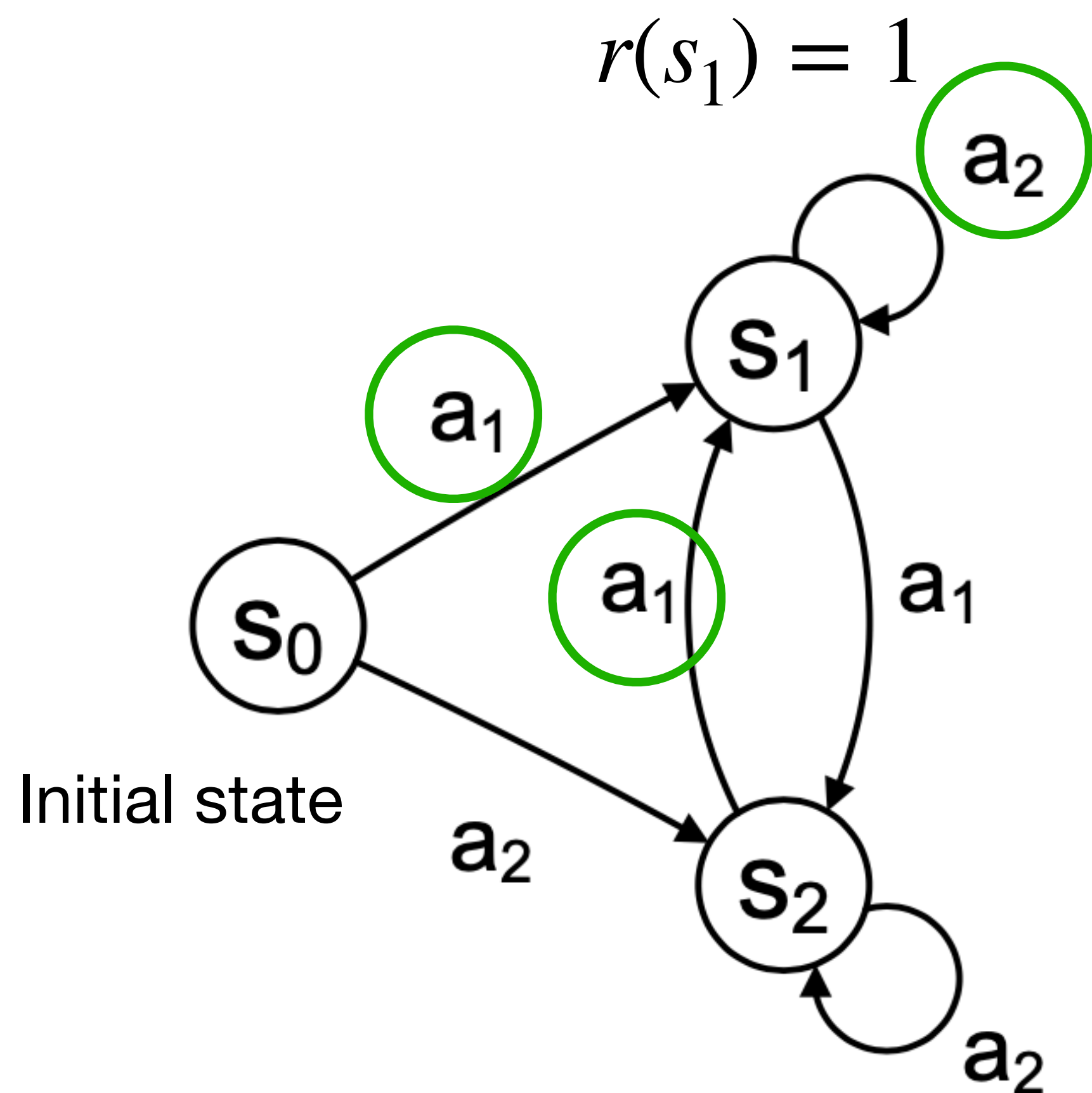
# Distribution Shift: Example



$$d_{s_0}^{\pi^*}(s_0) = 1 - \gamma, d_{s_0}^{\pi^*}(s_1) = \gamma, d_{s_0}^{\pi^*}(s_2) = 0$$

$$V_{s_0}^{\pi^*} = \frac{\gamma}{1 - \gamma}$$

# Distribution Shift: Example



Assume SL returned such policy  $\hat{\pi}$

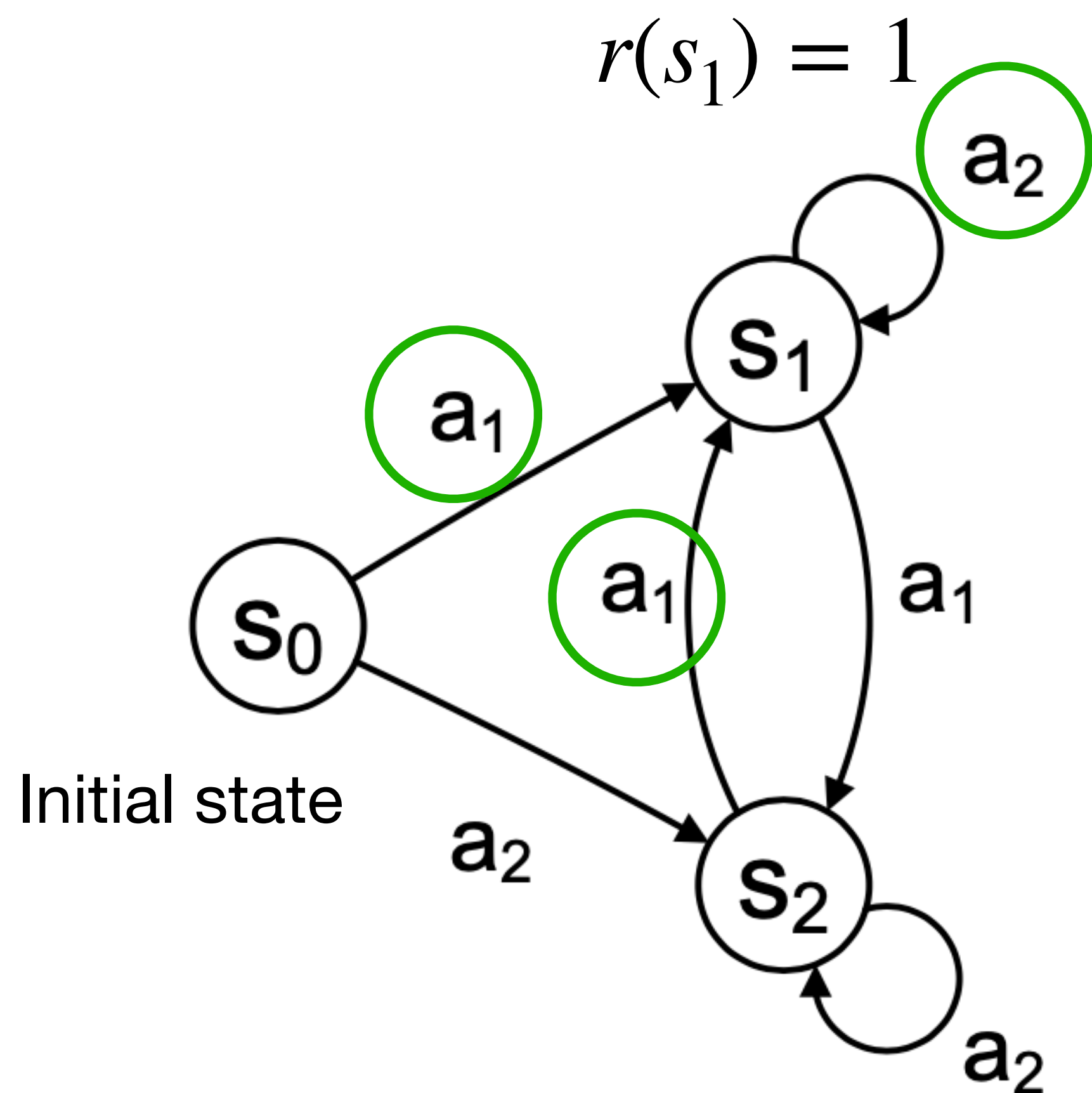
$$\hat{\pi}(s_0) = \begin{cases} a_1 & \text{w/ prob } 1 - \epsilon/(1 - \gamma) \\ a_2 & \text{w/ prob } \epsilon/(1 - \gamma) \end{cases}, \quad \hat{\pi}(s_1) = a_2, \hat{\pi}(s_2) = a_2$$

$$d_{s_0}^{\pi^*}(s_0) = 1 - \gamma, \quad d_{s_0}^{\pi^*}(s_1) = \gamma, \quad d_{s_0}^{\pi^*}(s_2) = 0$$

$$V_{s_0}^{\pi^*} = \frac{\gamma}{1 - \gamma}$$



# Distribution Shift: Example



Assume SL returned such policy  $\hat{\pi}$

$$\hat{\pi}(s_0) = \begin{cases} a_1 & \text{w/ prob } 1 - \epsilon/(1 - \gamma) \\ a_2 & \text{w/ prob } \epsilon/(1 - \gamma) \end{cases}, \quad \hat{\pi}(s_1) = a_2, \hat{\pi}(s_2) = a_2$$

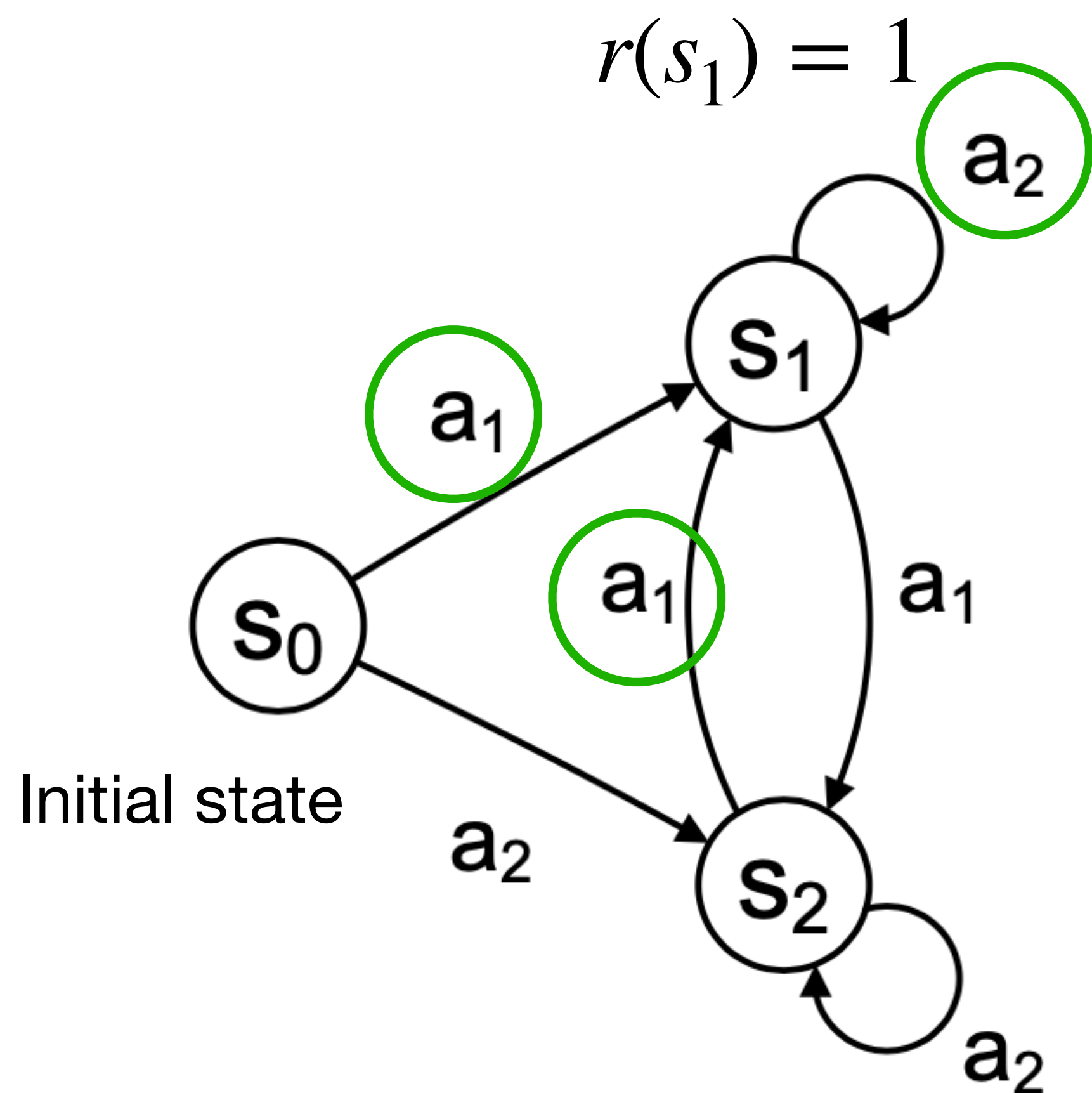
We will have good supervised learning error:

$$\mathbb{E}_{s \sim d_{s_0}^{\pi^*}} \mathbb{E}_{a \sim \hat{\pi}(\cdot|s)} \mathbf{1}(a \neq \pi^*(s)) = \epsilon$$

$$d_{s_0}^{\pi^*}(s_0) = 1 - \gamma, \quad d_{s_0}^{\pi^*}(s_1) = \gamma, \quad d_{s_0}^{\pi^*}(s_2) = 0$$

$$V_{s_0}^{\pi^*} = \frac{\gamma}{1 - \gamma}$$

# Distribution Shift: Example



Assume SL returned such policy  $\hat{\pi}$

$$\hat{\pi}(s_0) = \begin{cases} a_1 & \text{w/ prob } 1 - \epsilon/(1 - \gamma) \\ a_2 & \text{w/ prob } \epsilon/(1 - \gamma) \end{cases}, \quad \hat{\pi}(s_1) = a_2, \hat{\pi}(s_2) = a_2$$

We will have good supervised learning error:

$$\mathbb{E}_{s \sim d_{s_0}^{\pi^*}} \mathbb{E}_{a \sim \hat{\pi}(\cdot|s)} \mathbf{1}(a \neq \pi^*(s)) = \epsilon$$

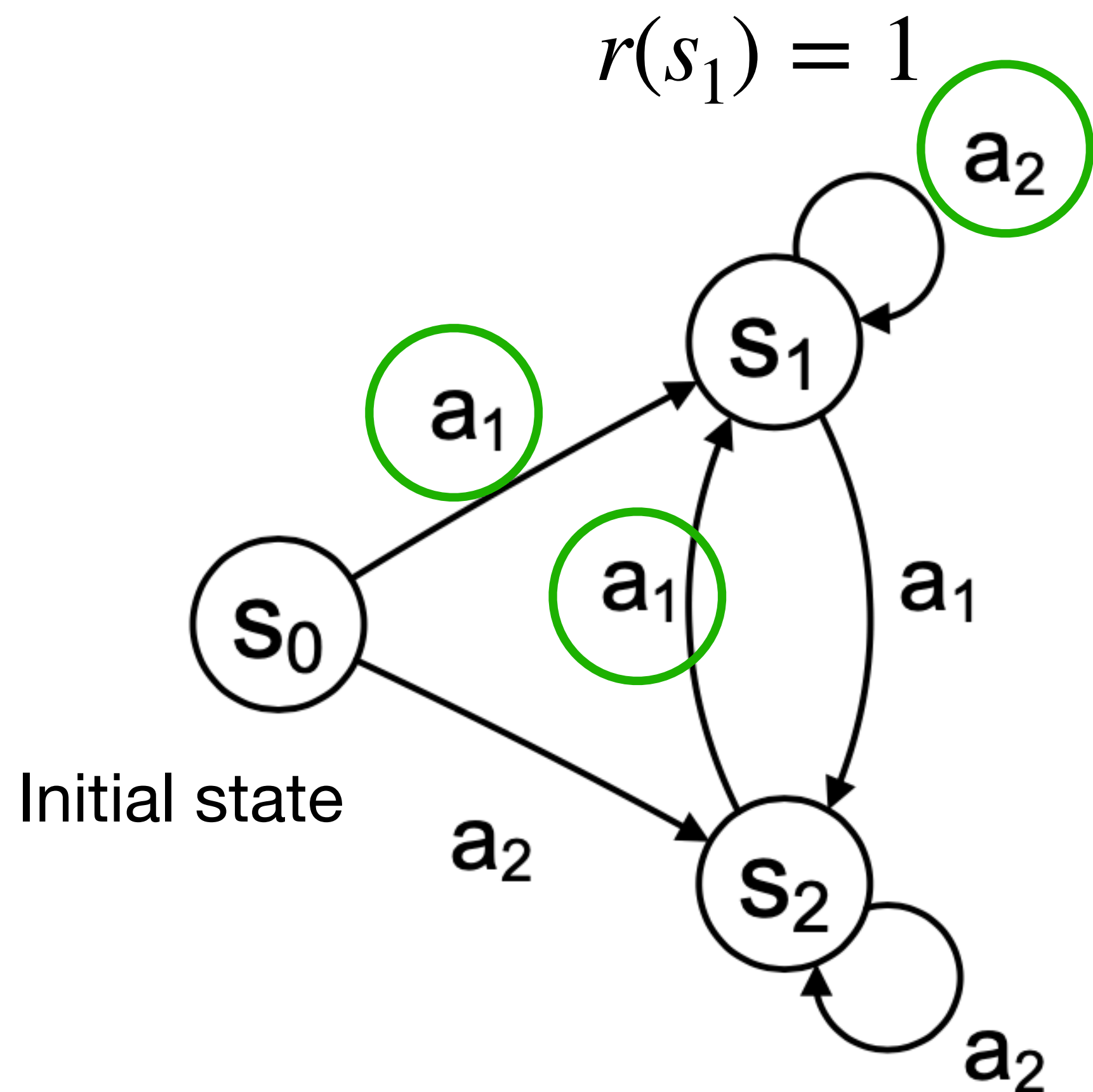
But we have quadratic error in performance:

$$V_{s_0}^{\hat{\pi}} = V_{s_0}^{\pi^*} - \frac{\epsilon\gamma}{(1 - \gamma)^2}$$

$$d_{s_0}^{\pi^*}(s_0) = 1 - \gamma, \quad d_{s_0}^{\pi^*}(s_1) = \gamma, \quad d_{s_0}^{\pi^*}(s_2) = 0$$

$$V_{s_0}^{\pi^*} = \frac{\gamma}{1 - \gamma}$$

# Distribution Shift: Example



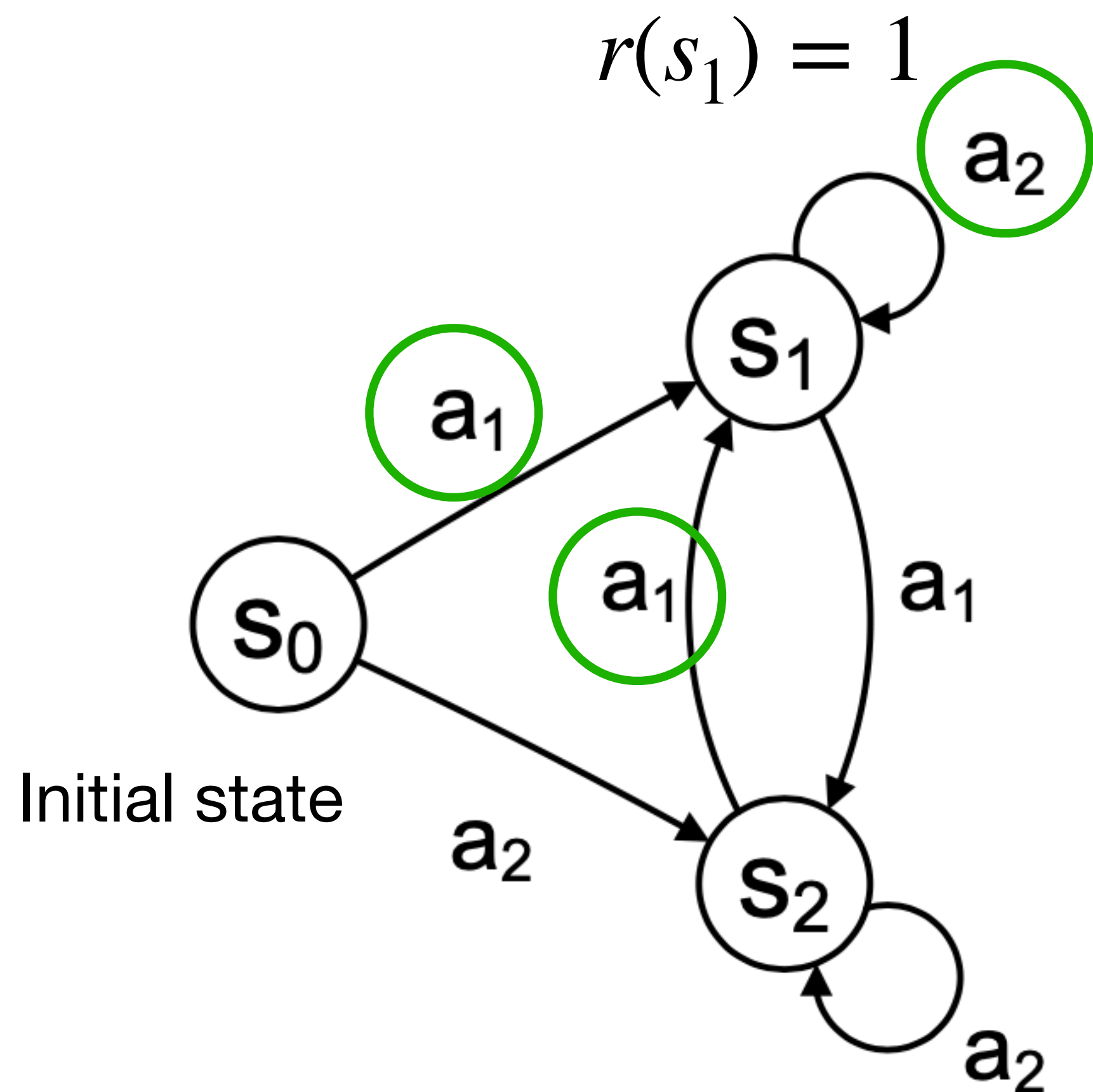
Assume SL returned such policy  $\hat{\pi}$

$$\hat{\pi}(s_0) = \begin{cases} a_1 & \text{w/ prob } 1 - \epsilon/(1 - \gamma) \\ a_2 & \text{w/ prob } \epsilon/(1 - \gamma) \end{cases}, \quad \hat{\pi}(s_1) = a_2, \hat{\pi}(s_2) = a_2$$

$$d_{s_0}^{\pi^*}(s_0) = 1 - \gamma, \quad d_{s_0}^{\pi^*}(s_1) = \gamma, \quad d_{s_0}^{\pi^*}(s_2) = 0$$

$$V_{s_0}^{\pi^*} = \frac{\gamma}{1 - \gamma}$$

# Distribution Shift: Example



Assume SL returned such policy  $\hat{\pi}$

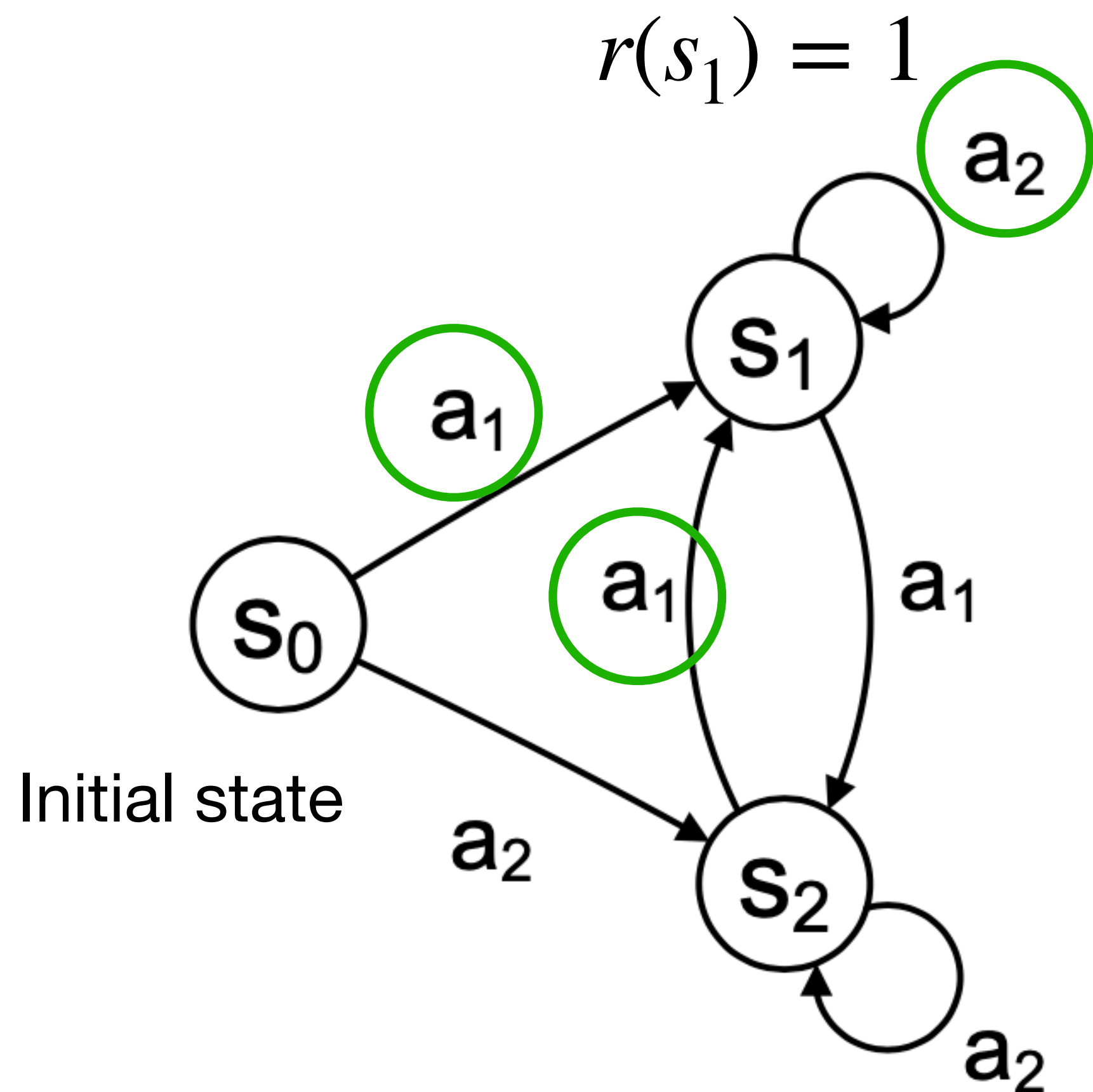
$$\hat{\pi}(s_0) = \begin{cases} a_1 & \text{w/ prob } 1 - \epsilon/(1 - \gamma) \\ a_2 & \text{w/ prob } \epsilon/(1 - \gamma) \end{cases}, \quad \hat{\pi}(s_1) = a_2, \hat{\pi}(s_2) = a_2$$

Why DAgger can fix this problem?

$$d_{s_0}^{\pi^*}(s_0) = 1 - \gamma, \quad d_{s_0}^{\pi^*}(s_1) = \gamma, \quad d_{s_0}^{\pi^*}(s_2) = 0$$

$$V_{s_0}^{\pi^*} = \frac{\gamma}{1 - \gamma}$$

# Distribution Shift: Example



Assume SL returned such policy  $\hat{\pi}$

$$\hat{\pi}(s_0) = \begin{cases} a_1 & \text{w/ prob } 1 - \epsilon/(1 - \gamma) \\ a_2 & \text{w/ prob } \epsilon/(1 - \gamma) \end{cases}, \quad \hat{\pi}(s_1) = a_2, \hat{\pi}(s_2) = a_2$$

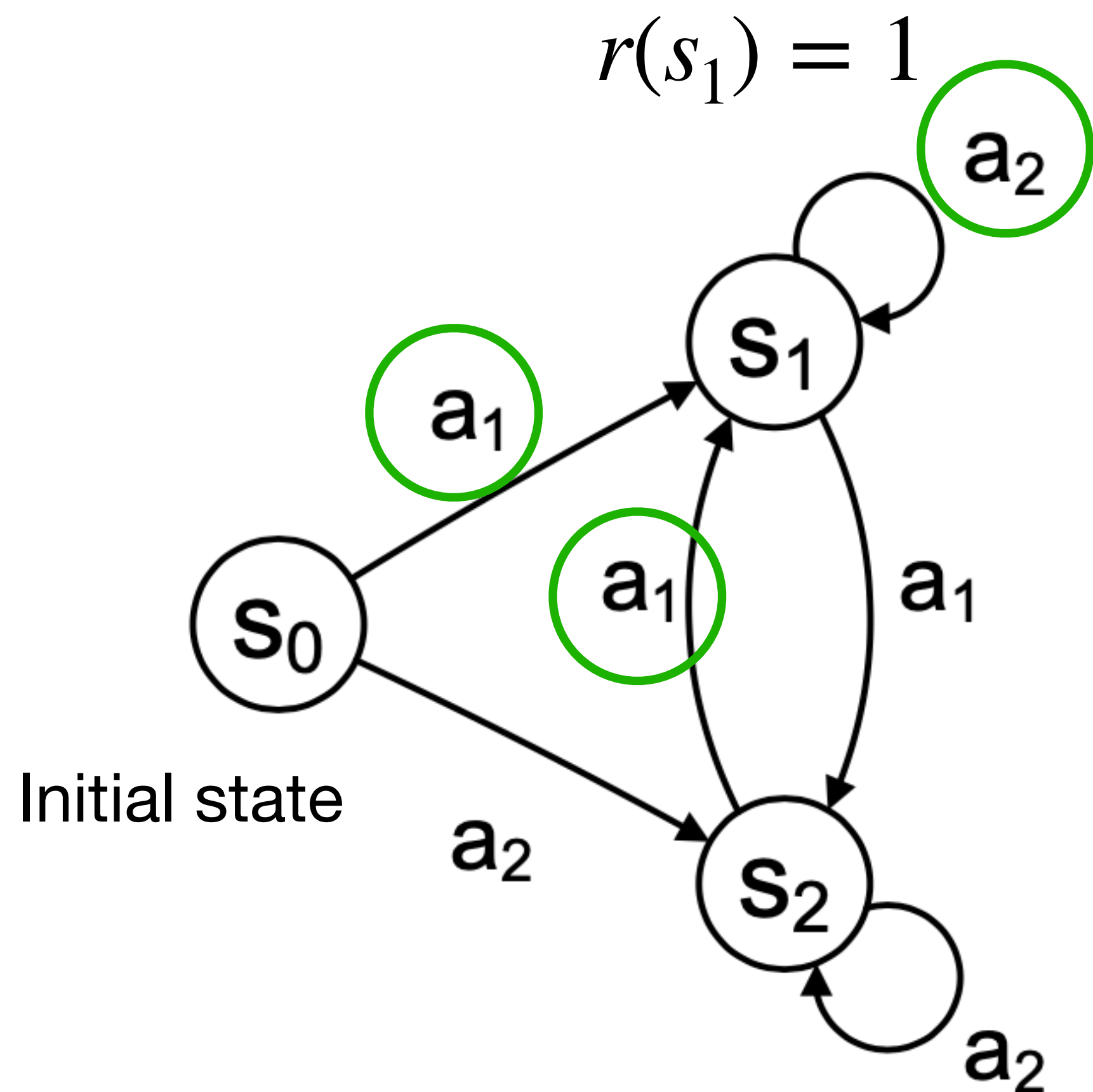
Why DAgger can fix this problem?

$\hat{\pi}$  will visit  $s_2$ , and we collect  $\pi^*(s_2) = a_1$

$$d_{s_0}^{\pi^*}(s_0) = 1 - \gamma, d_{s_0}^{\pi^*}(s_1) = \gamma, d_{s_0}^{\pi^*}(s_2) = 0$$

$$V_{s_0}^{\pi^*} = \frac{\gamma}{1 - \gamma}$$

# Distribution Shift: Example



Assume SL returned such policy  $\hat{\pi}$

$$\hat{\pi}(s_0) = \begin{cases} a_1 & \text{w/ prob } 1 - \epsilon/(1 - \gamma) \\ a_2 & \text{w/ prob } \epsilon/(1 - \gamma) \end{cases}, \quad \hat{\pi}(s_1) = a_2, \hat{\pi}(s_2) = a_2$$

Why DAgger can fix this problem?

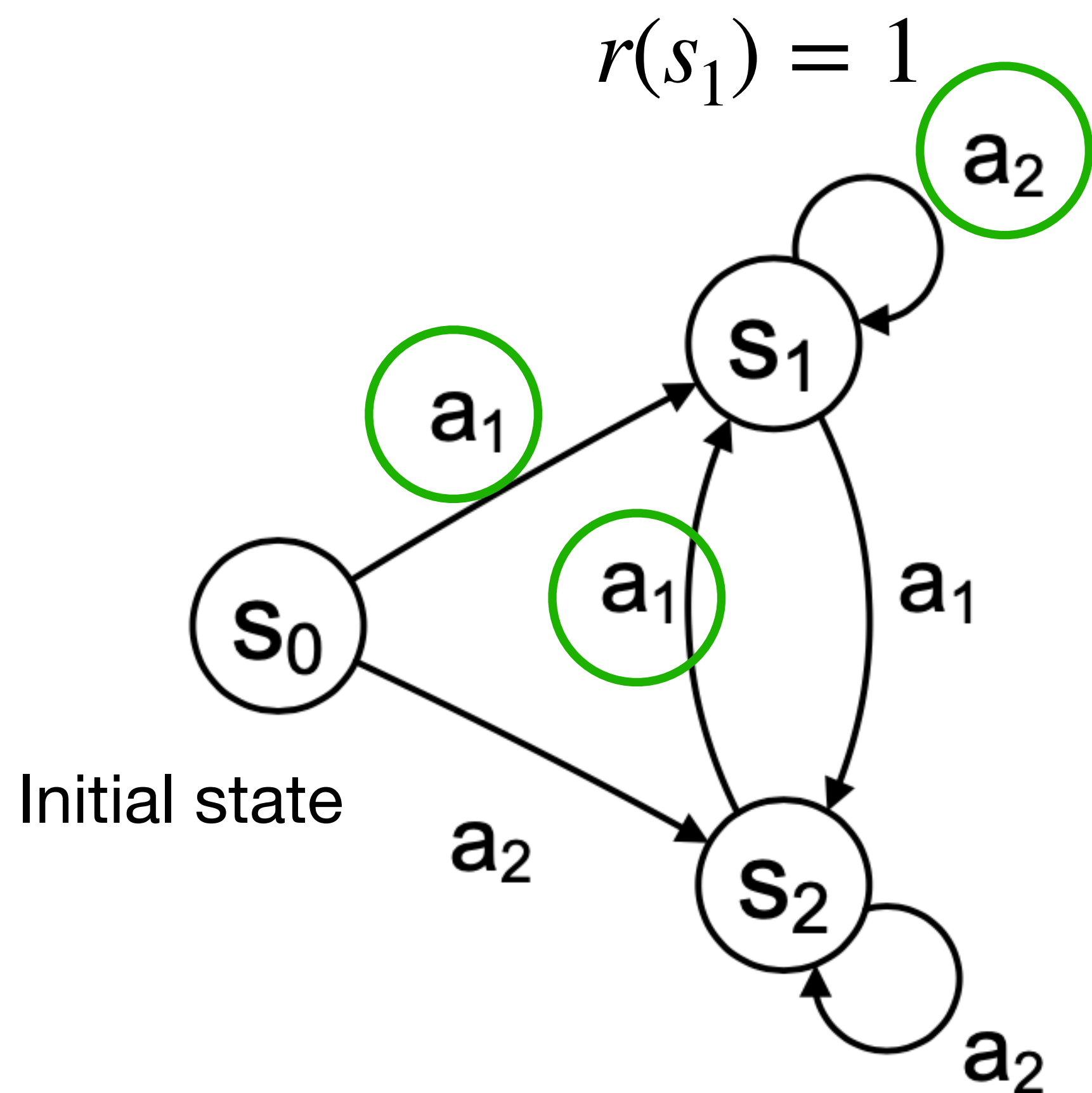
$\hat{\pi}$  will visit  $s_2$ , and we collect  $\pi^*(s_2) = a_1$

By data aggregation, our new dataset will contain  $(s_2, a_1)$  pairs

$$d_{s_0}^{\pi^*}(s_0) = 1 - \gamma, \quad d_{s_0}^{\pi^*}(s_1) = \gamma, \quad d_{s_0}^{\pi^*}(s_2) = 0$$

$$V_{s_0}^{\pi^*} = \frac{\gamma}{1 - \gamma}$$

# Distribution Shift: Example



Assume SL returned such policy  $\hat{\pi}$

$$\hat{\pi}(s_0) = \begin{cases} a_1 & \text{w/ prob } 1 - \epsilon/(1 - \gamma) \\ a_2 & \text{w/ prob } \epsilon/(1 - \gamma) \end{cases}, \quad \hat{\pi}(s_1) = a_2, \hat{\pi}(s_2) = a_2$$

Why DAgger can fix this problem?

$\hat{\pi}$  will visit  $s_2$ , and we collect  $\pi^*(s_2) = a_1$

By data aggregation, our new dataset will contain  $(s_2, a_1)$  pairs

Thus, our new learned policy will know what to do at  $s_2$

$$d_{s_0}^{\pi^*}(s_0) = 1 - \gamma, \quad d_{s_0}^{\pi^*}(s_1) = \gamma, \quad d_{s_0}^{\pi^*}(s_2) = 0$$

$$V_{s_0}^{\pi^*} = \frac{\gamma}{1 - \gamma}$$

## Outline for today:



1. The DAgger (Data Aggregation) Algorithm

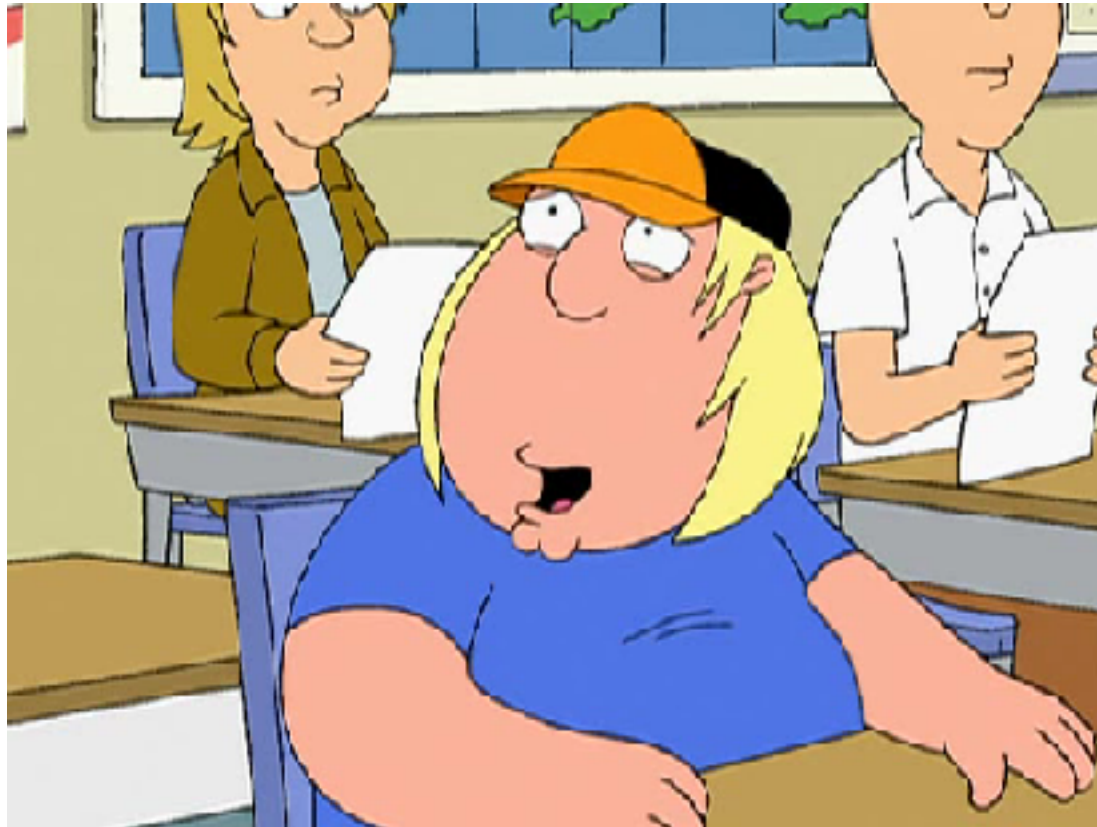
2. Quick intro on Online Learning



[Vovk92, Warmuth94, Freund97, Zinkevich03, Kalai05, Hazan06, Kakade08]

# Online Learning

**Learner**



convex Decision set  $\Theta$

**Adversary**

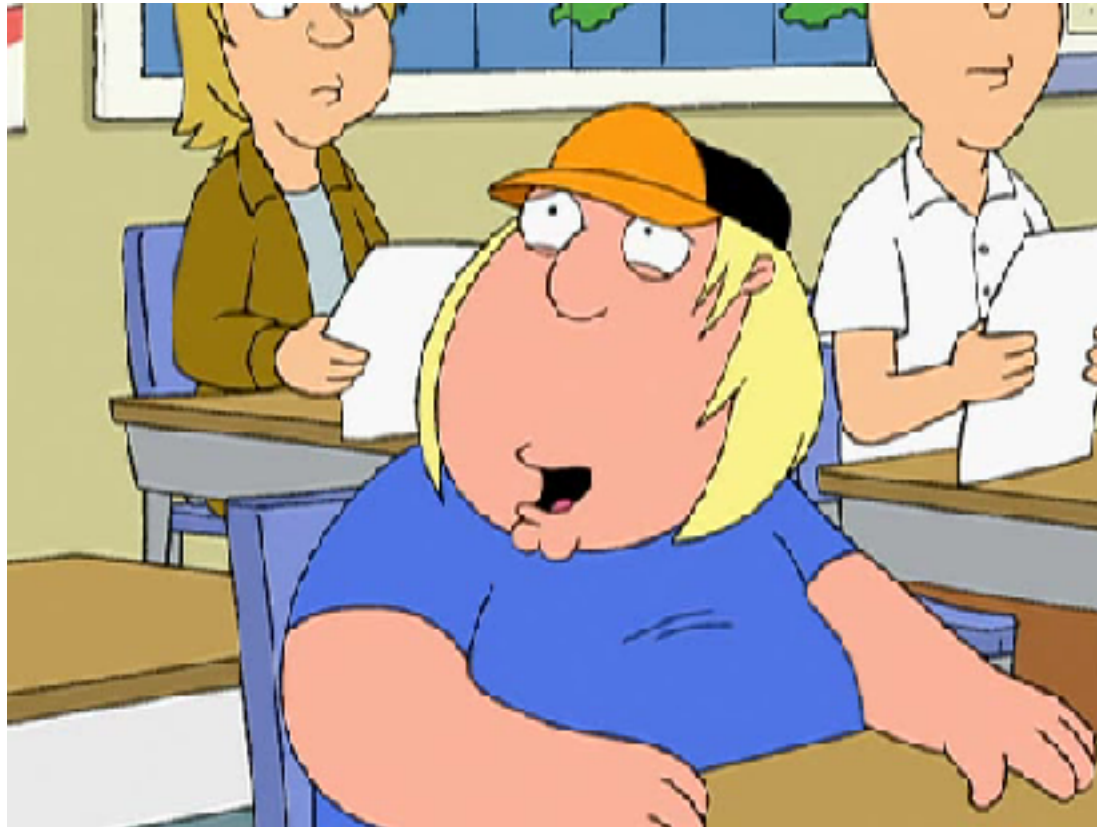


...

[Vovk92, Warmuth94, Freund97, Zinkevich03, Kalai05, Hazan06, Kakade08]

# Online Learning

**Learner**



convex Decision set  $\Theta$

Learner picks a decision  $\theta_0$

A black arrow pointing from the learner towards the adversary.

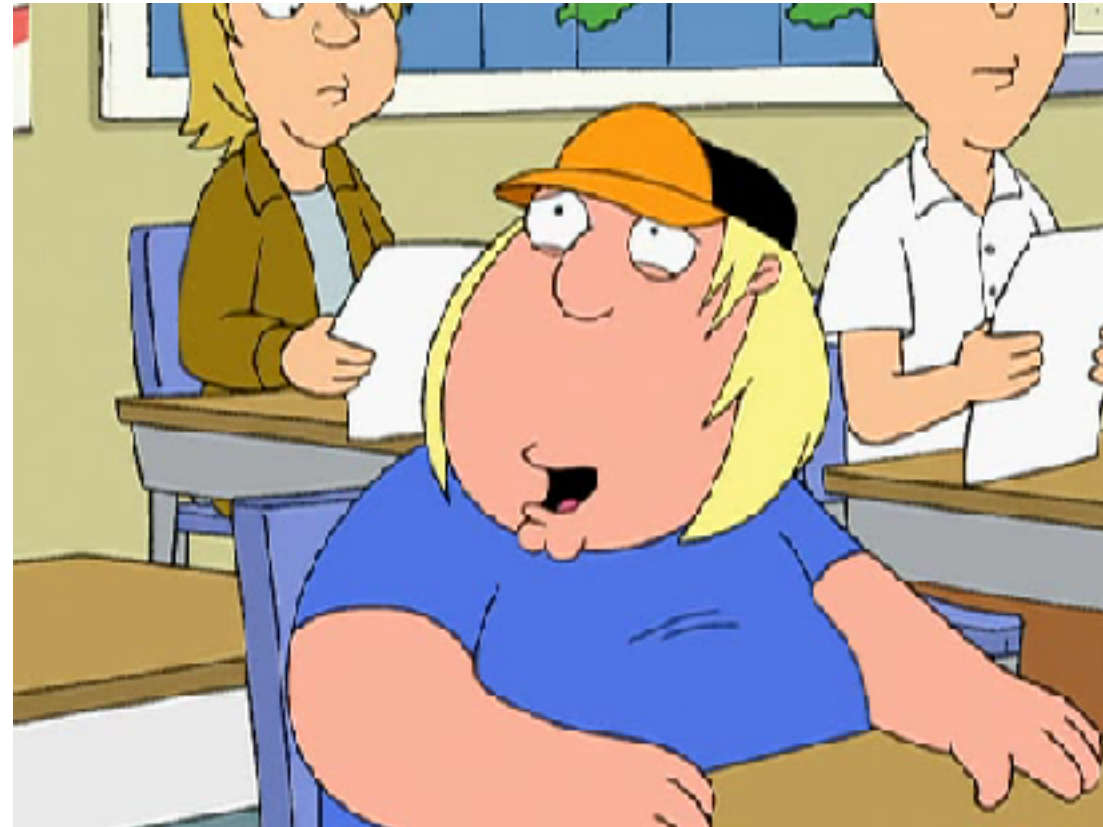
**Adversary**



...

# Online Learning

**Learner**



convex Decision set  $\Theta$

Learner picks a decision  $\theta_0$

→

Adversary picks a loss  $\ell_0 : \Theta \rightarrow \mathbb{R}$

←

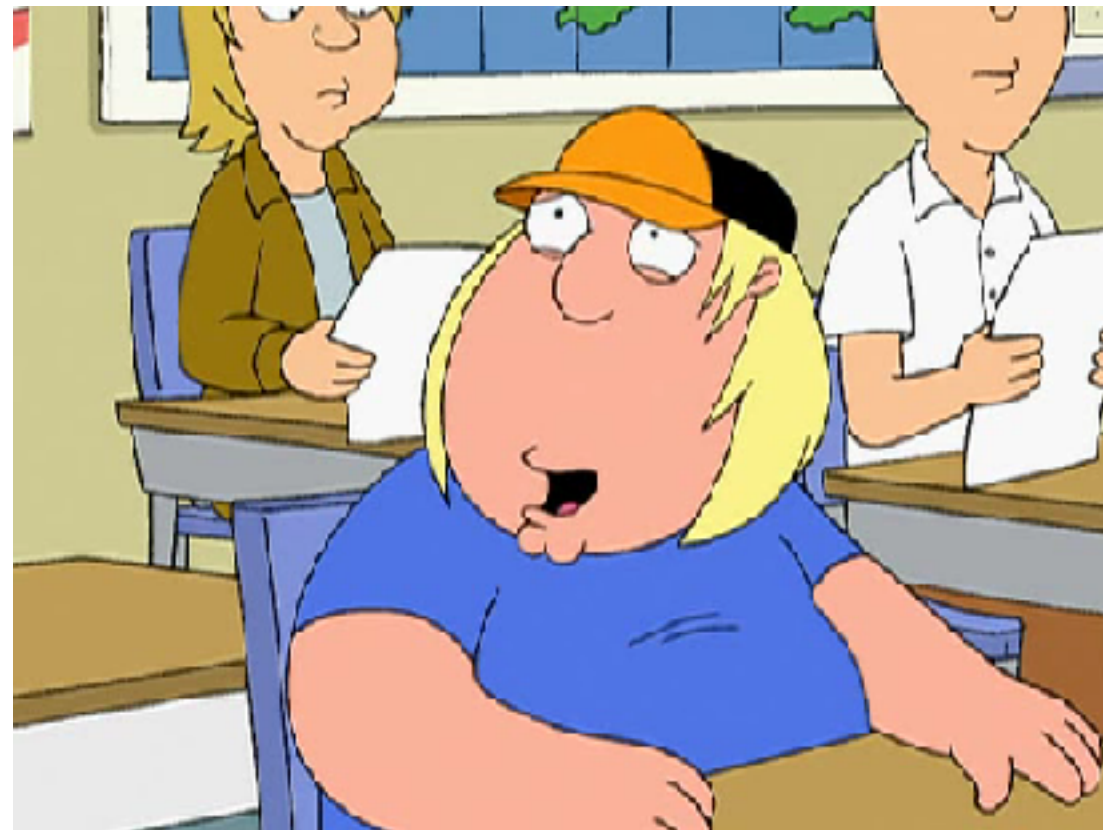
**Adversary**



...

# Online Learning

**Learner**



convex Decision set  $\Theta$

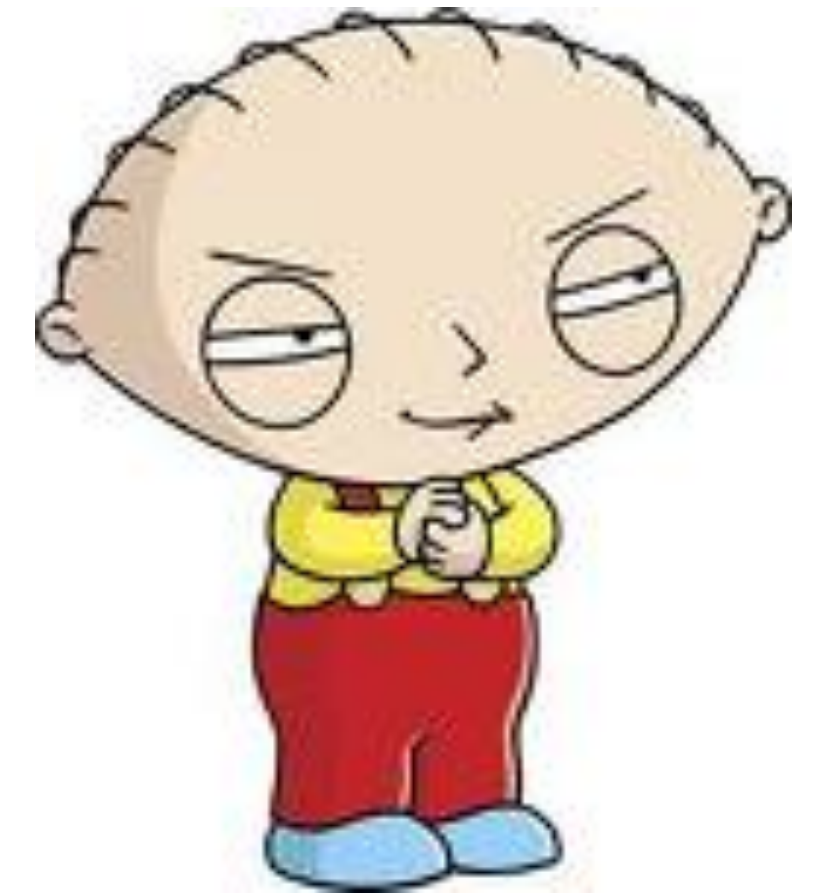
Learner picks a decision  $\theta_0$   
→

← Adversary picks a loss  $\ell_0 : \Theta \rightarrow \mathbb{R}$

Learner picks a new decision  $\theta_1$   
→

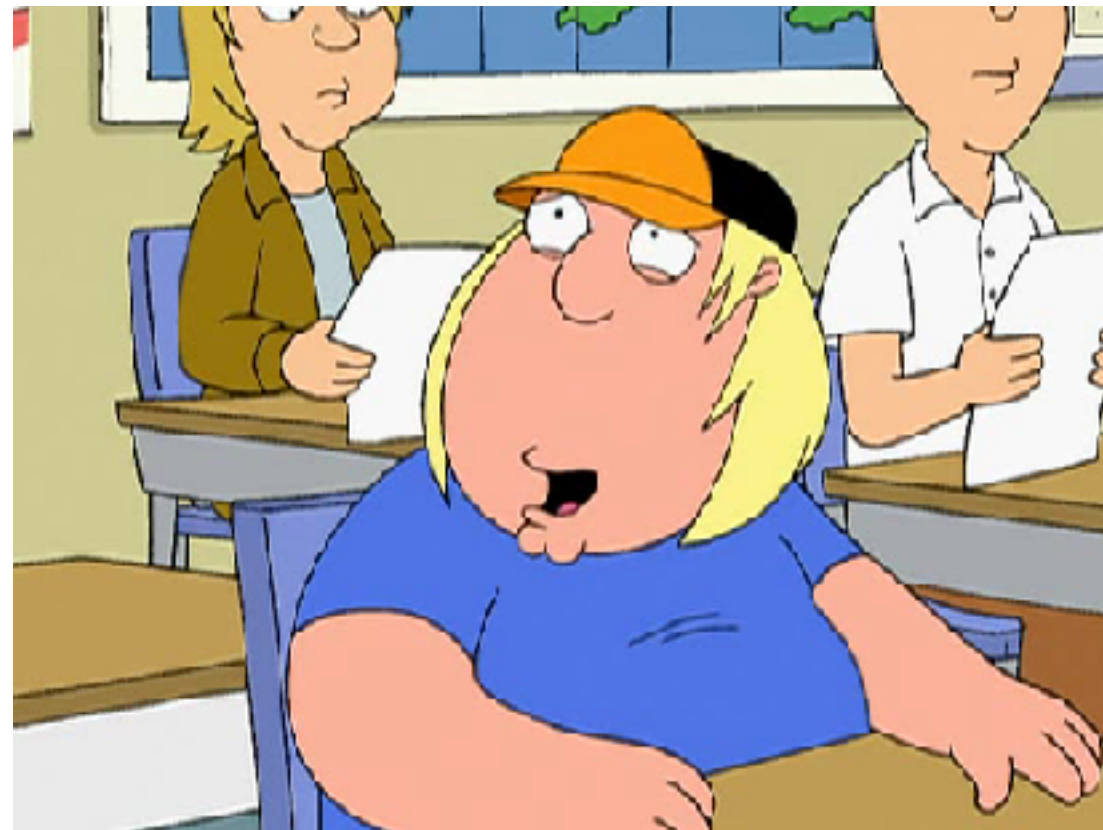
...

**Adversary**



# Online Learning

**Learner**



convex Decision set  $\Theta$

Learner picks a decision  $\theta_0$   
→

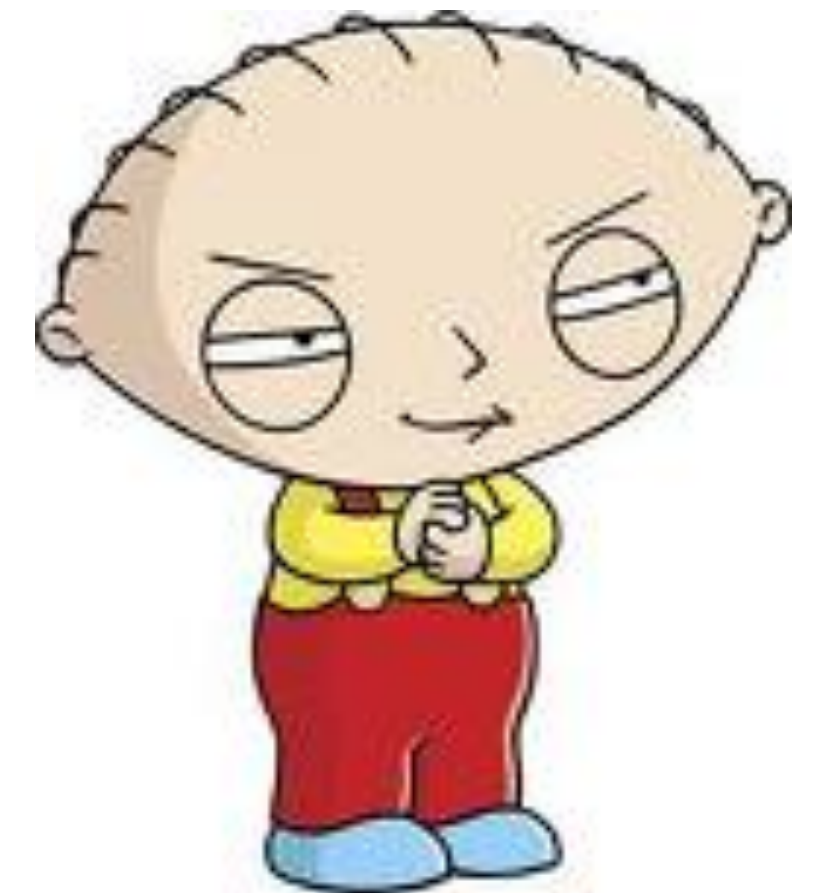
← Adversary picks a loss  $\ell_0 : \Theta \rightarrow \mathbb{R}$

Learner picks a new decision  $\theta_1$   
→

← Adversary picks a loss  $\ell_1 : \Theta \rightarrow \mathbb{R}$

...

**Adversary**



# Online Learning

**Learner**



convex Decision set  $\Theta$

Learner picks a decision  $\theta_0$



Adversary picks a loss  $\ell_0 : \Theta \rightarrow \mathbb{R}$



Learner picks a new decision  $\theta_1$



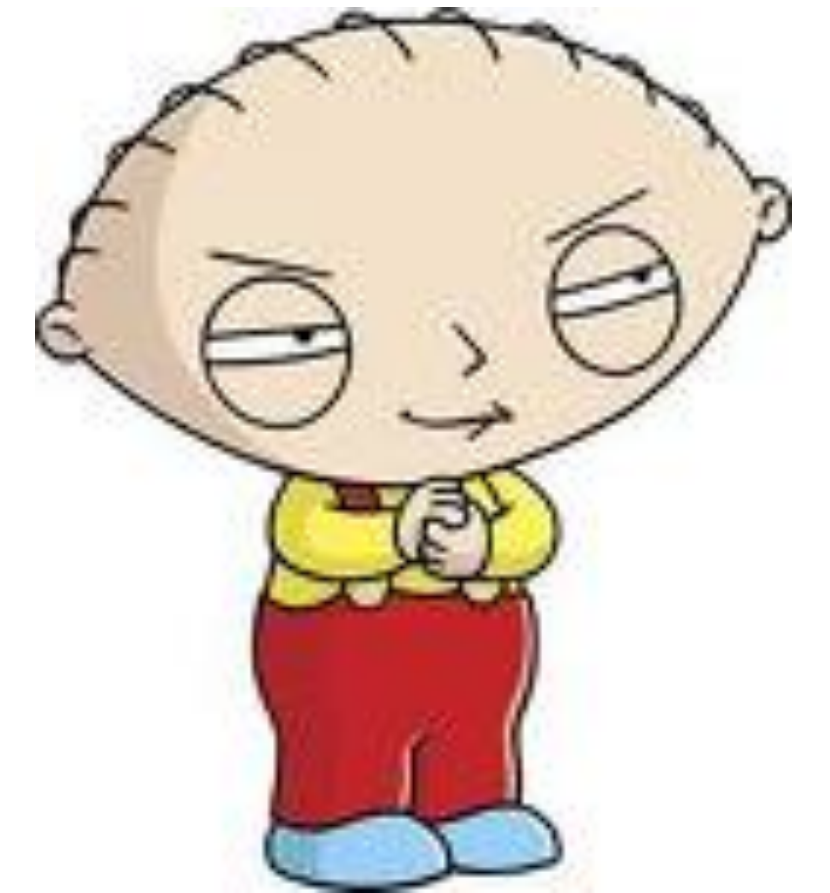
Adversary picks a loss  $\ell_1 : \Theta \rightarrow \mathbb{R}$



...

$$\text{Regret} = \sum_{t=0}^{T-1} \ell_t(\theta_t) - \min_{\theta \in \Theta} \sum_{t=0}^{T-1} \ell_t(\theta)$$

**Adversary**



## **Example: online linear regression**

Can we perform linear regression in online fashion with non i.i.d (or even adversary) data?

## Example: online linear regression

Can we perform linear regression in online fashion with non i.i.d (or even adversary) data?

Every iteration  $t$  :



## Example: online linear regression

Can we perform linear regression in online fashion with non i.i.d (or even adversary) data?

**Every iteration  $t$  :**

1. Learner first picks  $\theta_t \in \text{Ball} \subset \mathbb{R}^d$

## Example: online linear regression

Can we perform linear regression in online fashion with non i.i.d (or even adversary) data?

Every iteration  $t$  :

1. Learner first picks  $\theta_t \in \text{Ball} \subset \mathbb{R}^d$
2. Adversary **then** picks  $x_t \in \mathcal{X} \subset \mathbb{R}^d, y_t \in [a, b]$

## Example: online linear regression

Can we perform linear regression in online fashion with non i.i.d (or even adversary) data?

Every iteration  $t$  :

1. Learner first picks  $\theta_t \in \text{Ball} \subset \mathbb{R}^d$
2. Adversary **then** picks  $x_t \in \mathcal{X} \subset \mathbb{R}^d, y_t \in [a, b]$
3. Learner suffers loss  $\ell_t(\theta_t) = (\theta_t^\top x_t - y_t)^2$

## Example: online linear regression

Can we perform linear regression in online fashion with non i.i.d (or even adversary) data?

Every iteration  $t$  :

1. Learner first picks  $\theta_t \in \text{Ball} \subset \mathbb{R}^d$
2. Adversary **then** picks  $x_t \in \mathcal{X} \subset \mathbb{R}^d, y_t \in [a, b]$
3. Learner suffers loss  $\ell_t(\theta_t) = (\theta_t^\top x_t - y_t)^2$

Learner has to make decision  $\theta_t$  based on history up to  $t - 1$ ,  
while adversary could pick  $(x_t, y_t)$  even after seeing  $\theta_t$

## Example: online linear regression

Can we perform linear regression in online fashion with non i.i.d (or even adversary) data?

Every iteration  $t$  :

1. Learner first picks  $\theta_t \in \text{Ball} \subset \mathbb{R}^d$
2. Adversary **then** picks  $x_t \in \mathcal{X} \subset \mathbb{R}^d, y_t \in [a, b]$
3. Learner suffers loss  $\ell_t(\theta_t) = (\theta_t^\top x_t - y_t)^2$

Learner has to make decision  $\theta_t$  based on history up to  $t - 1$ ,  
while adversary could pick  $(x_t, y_t)$  even after seeing  $\theta_t$

Adversary seems too powerful...

## **Example: online linear regression**

BUT, a very intuitive algorithm actually achieves no-regret property:

## Example: online linear regression

BUT, a very intuitive algorithm actually achieves no-regret property:

**Every iteration  $t$  :**

1. Learner first picks  $\theta_t$  that minimizes the aggregated loss

$$\theta_t = \arg \min_{\theta \in \text{Ball}} \sum_{i=0}^{t-1} (\theta^\top x_i - y_i)^2 + \lambda \|\theta\|_2^2$$

## Example: online linear regression

BUT, a very intuitive algorithm actually achieves no-regret property:

Every iteration  $t$  :

1. Learner first picks  $\theta_t$  that minimizes the aggregated loss

$$\theta_t = \arg \min_{\theta \in \text{Ball}} \sum_{i=0}^{t-1} (\theta^\top x_i - y_i)^2 + \lambda \|\theta\|_2^2$$

This is called Follow-the-Regularized-Leader (FTRL), and it achieves no-regret property:



## Example: online linear regression

BUT, a very intuitive algorithm actually achieves no-regret property:

Every iteration  $t$  :

1. Learner first picks  $\theta_t$  that minimizes the aggregated loss

$$\theta_t = \arg \min_{\theta \in \text{Ball}} \sum_{i=0}^{t-1} (\theta^\top x_i - y_i)^2 + \lambda \|\theta\|_2^2$$

This is called Follow-the-Regularized-Leader (FTRL), and it achieves no-regret property:

$$\sum_{i=0}^{T-1} \ell_i(\theta_i) - \min_{\theta \in \text{Ball}} \sum_{i=0}^{T-1} \ell_i(\theta) = O\left(1/\sqrt{T}\right)$$

## Generally, Follow-the-Regularized-Leader is no-regret

At time step  $t$ , learner has seen  $\ell_0, \dots, \ell_{t-1}$ , which new decision she could pick?

$$\mathbf{FTL: } \theta_t = \min_{\theta \in \Theta} \sum_{i=0}^{t-1} \ell_i(\theta) + \lambda R(\theta)$$

## Generally, Follow-the-Regularized-Leader is no-regret

At time step  $t$ , learner has seen  $\ell_0, \dots, \ell_{t-1}$ , which new decision she could pick?

$$\mathbf{FTL: } \theta_t = \min_{\theta \in \Theta} \sum_{i=0}^{t-1} \ell_i(\theta) + \lambda R(\theta)$$

**Theorem (FTL) (optional):** if  $\Theta$  is convex, and  $\ell_t$  is convex for all  $t$ , and  $R(\theta)$  is strongly convex, then for regret of FTL, we have:

$$\frac{1}{T} \left[ \sum_{t=0}^{T-1} \ell_t(\theta_t) - \min_{\theta \in \Theta} \sum_{t=0}^{T-1} \ell_t(\theta) \right] = O\left(1/\sqrt{T}\right)$$

## **Any questions about no-regret online learning?**

Online learning is a very rich research area — details are out of scope

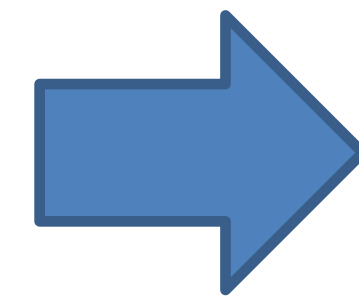
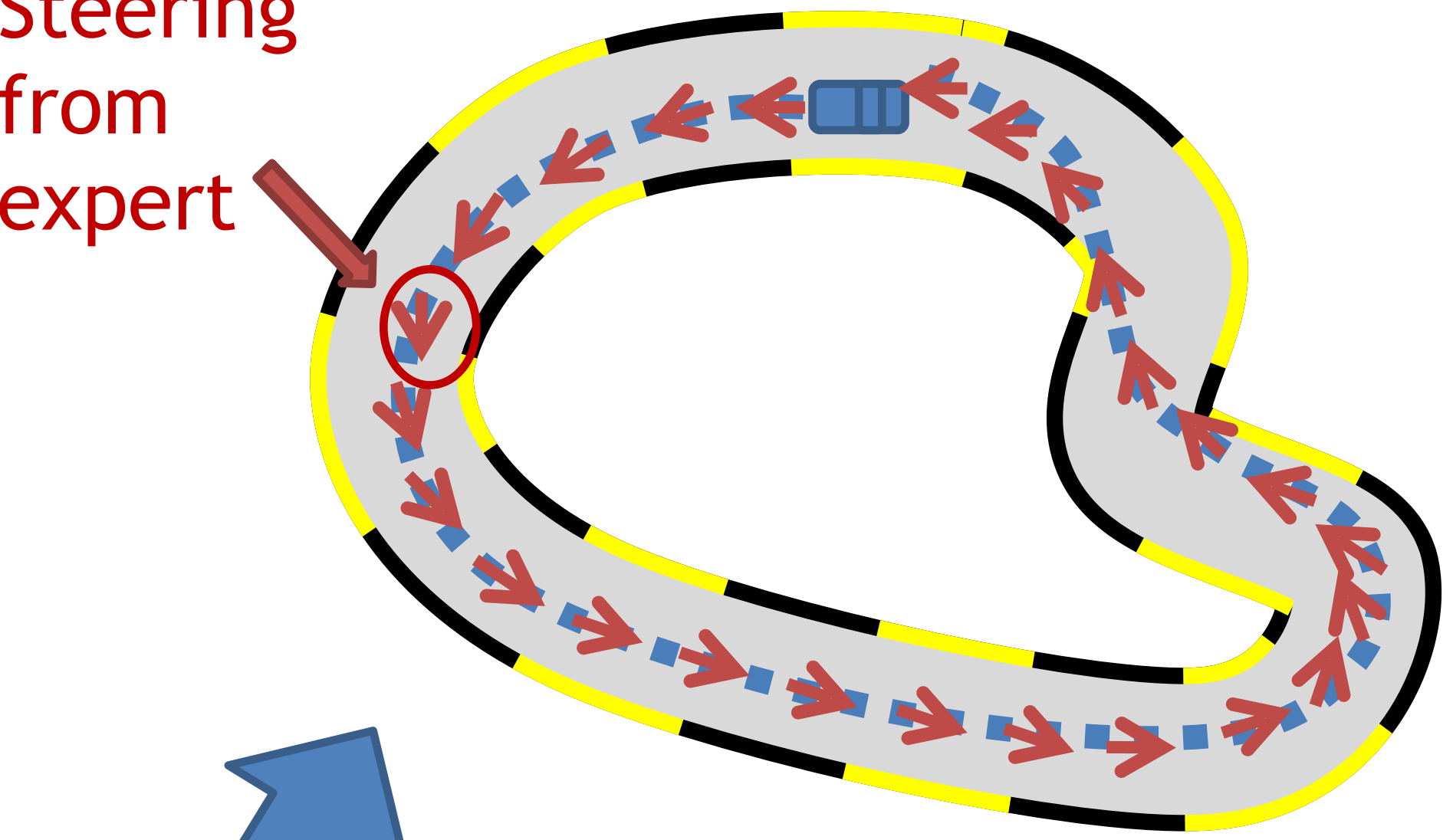
### **Key message:**

Learner has to make a decision before Adversary picks a loss function, yet it is possible to do as well as the best decision in hindsight if we had access to all the loss functions beforehand

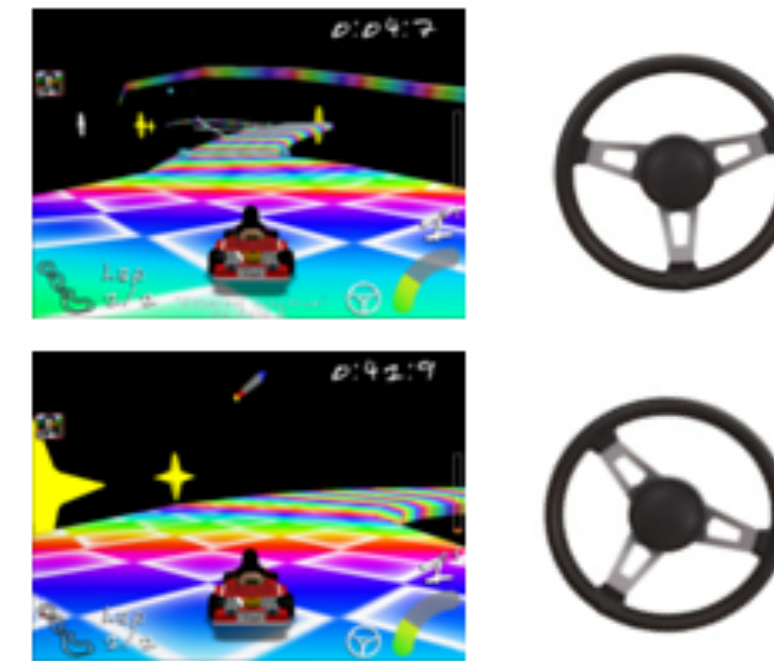
# Dagger Revisit

At iteration t:

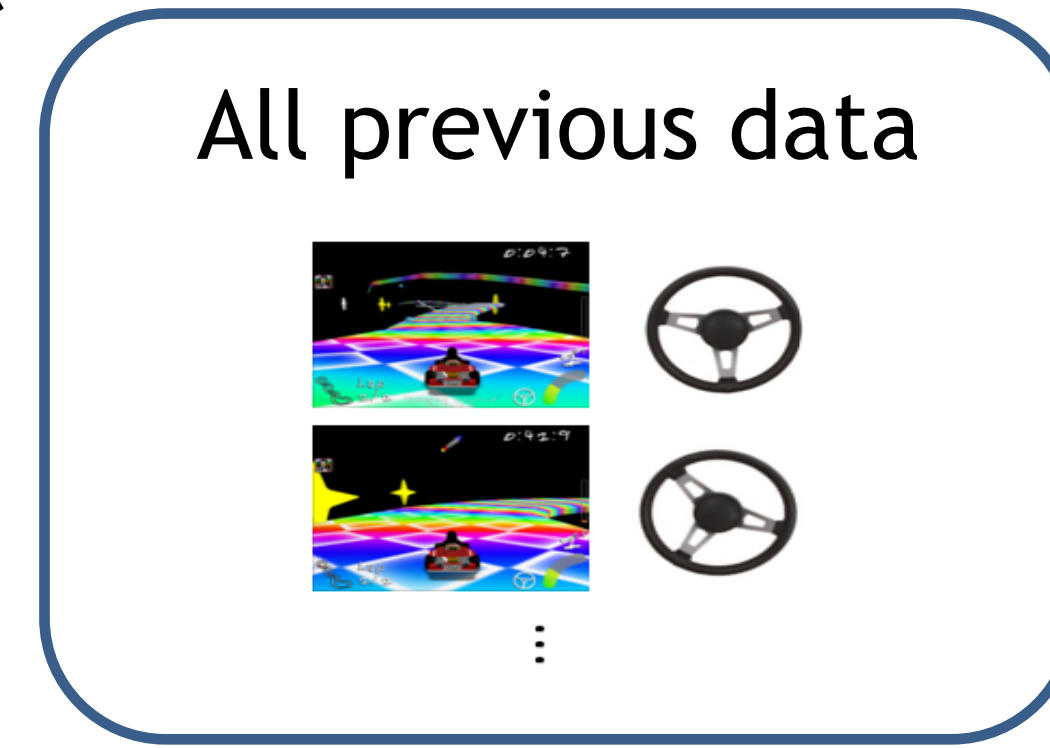
Steering from expert



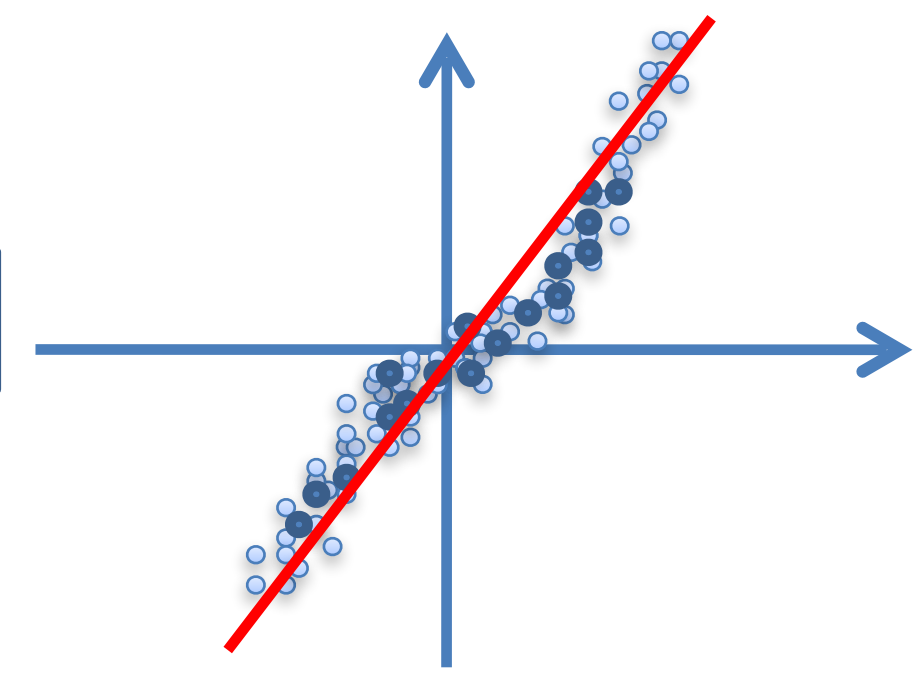
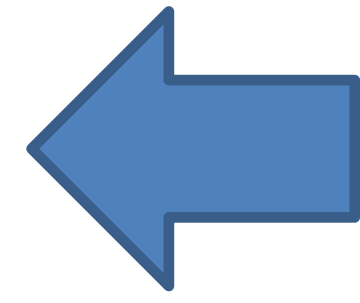
New Data



Aggregate Dataset



New policy  $\pi_n$

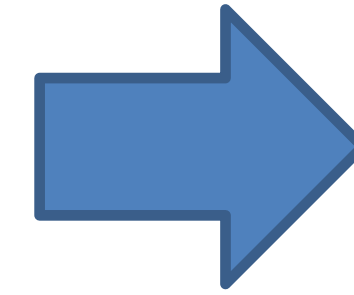
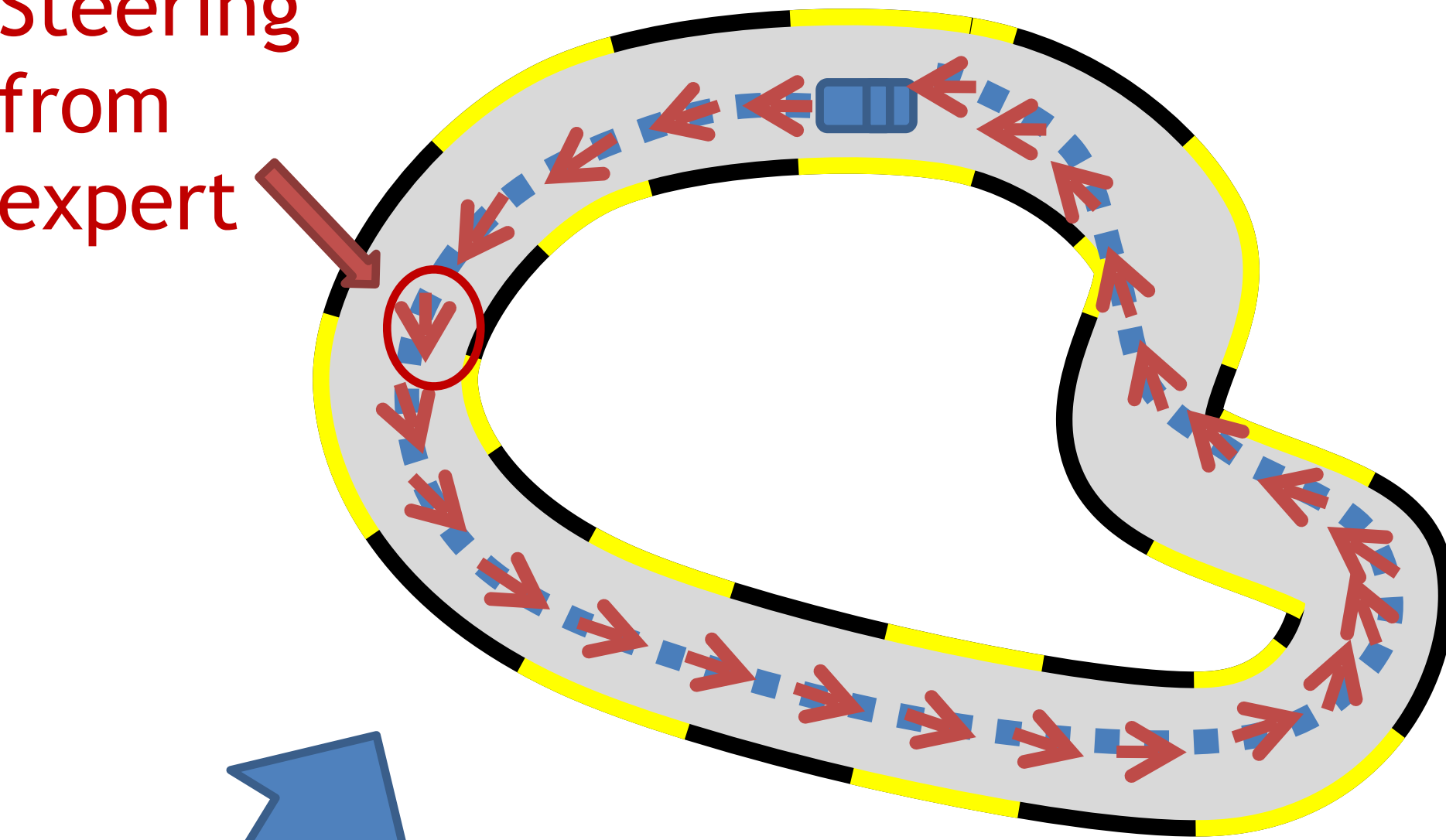


Supervised Learning

# Dagger Revisit

At iteration t:

Steering from expert

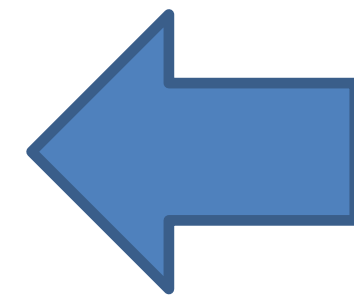
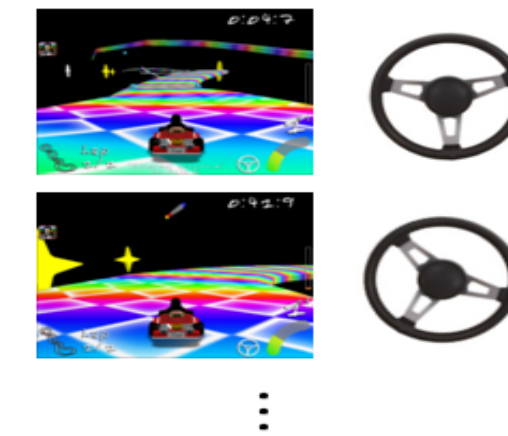


New Data



Aggregate Dataset

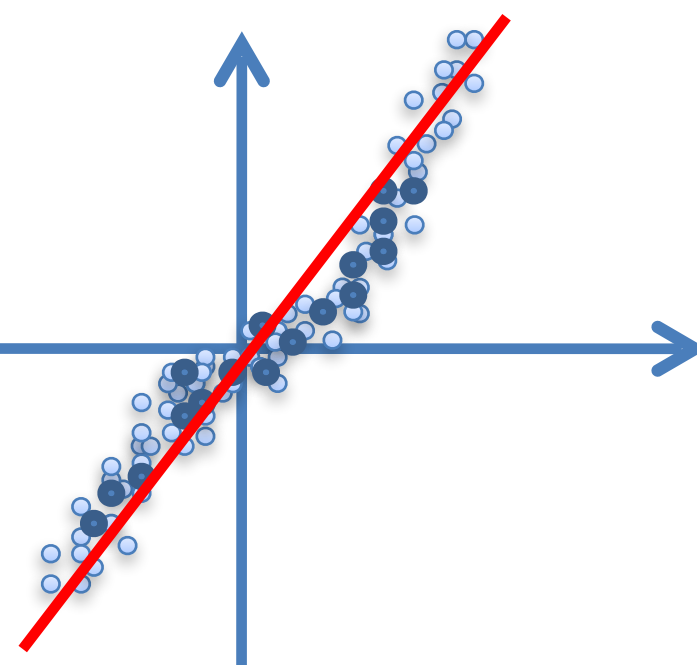
All previous data



New policy

$\pi_n$

Supervised Learning

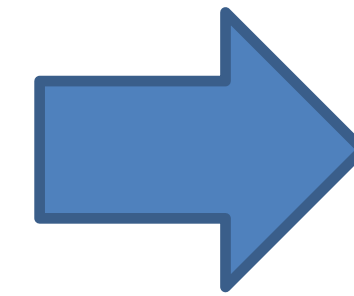
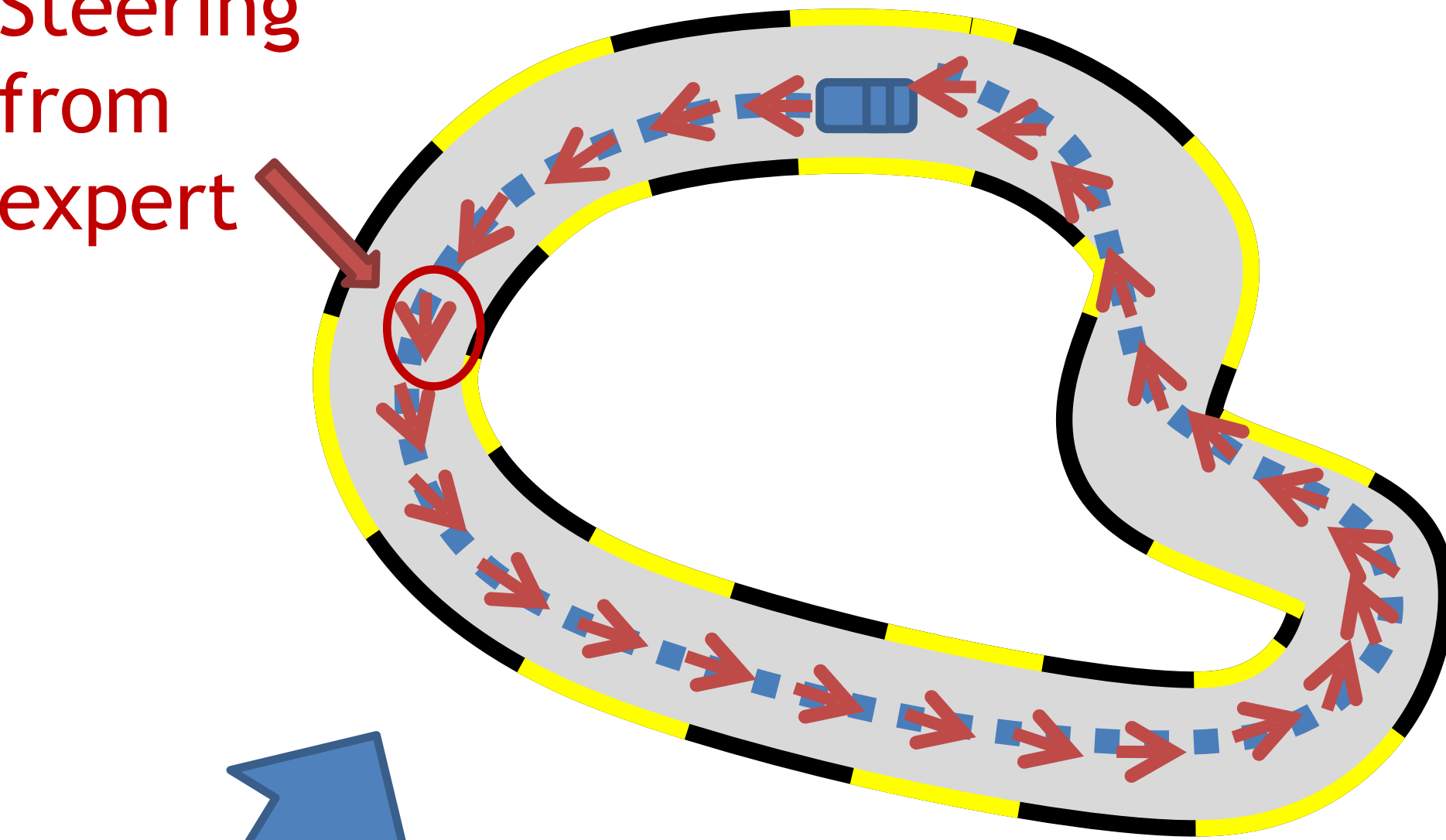


# Dagger Revisit

At iteration t:

New Data

Steering from expert

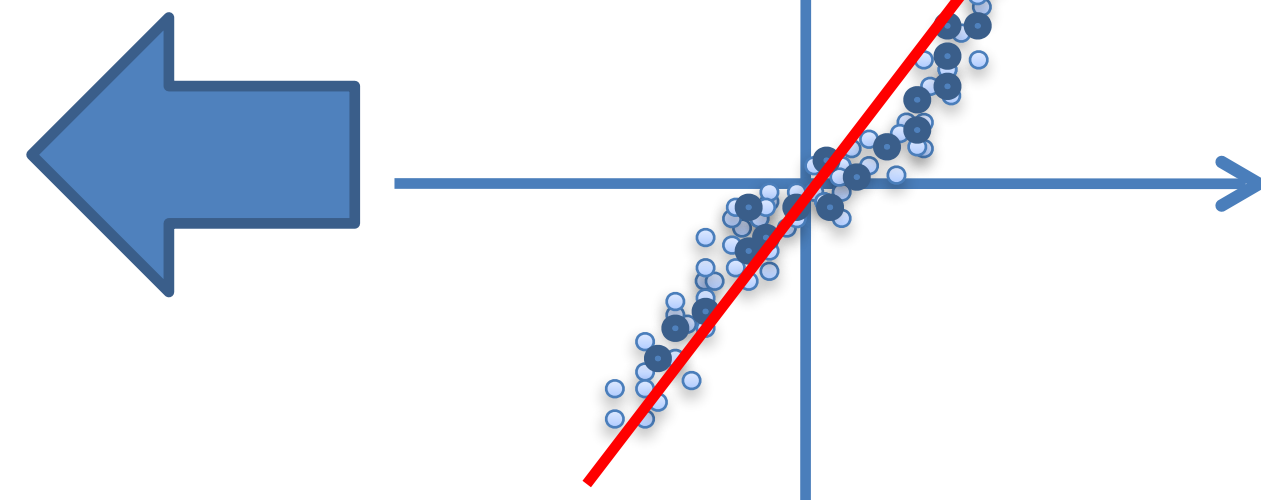
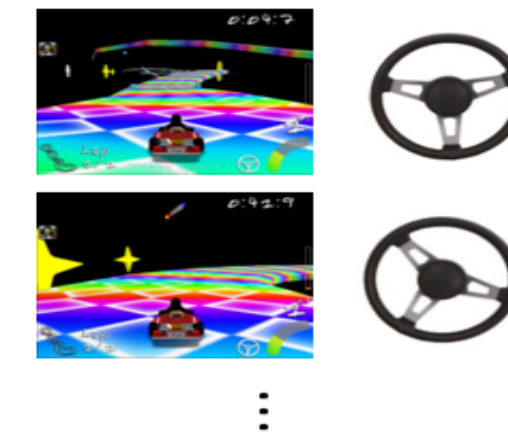


$$\mathcal{L}_t(\pi) = \sum_{i=1}^m \|\pi(s^i) - \pi^*(s^i)\|_2^2$$



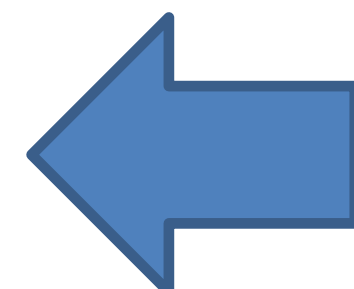
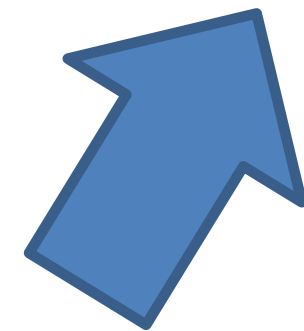
Aggregate Dataset

All previous data



Supervised Learning

New policy  $\pi_n$

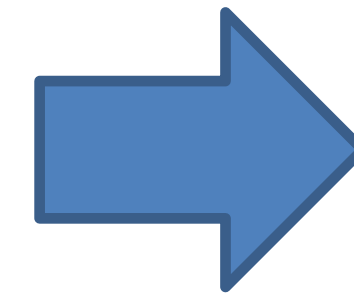
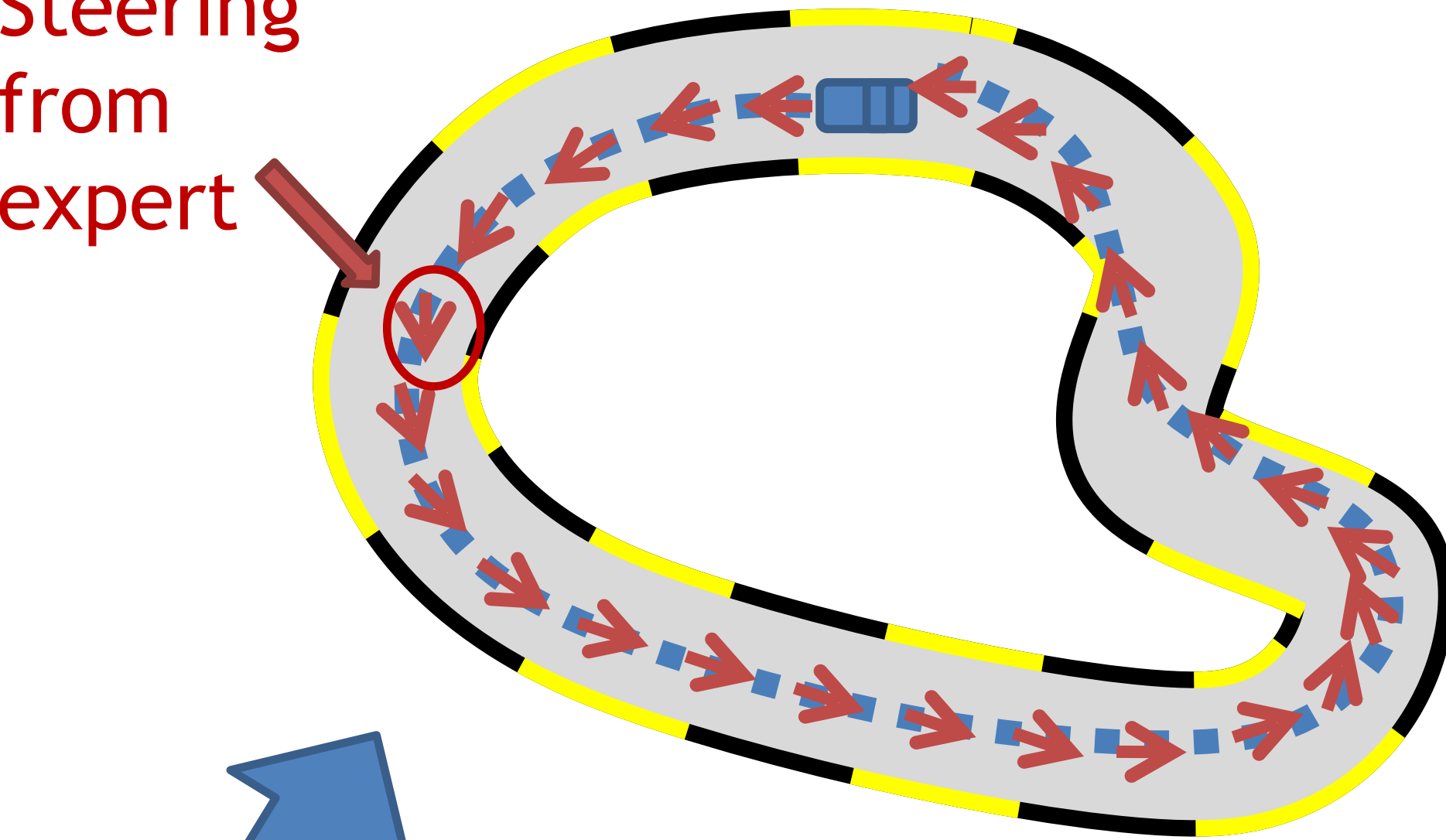


# Dagger Revisit

At iteration t:

New Data

Steering from expert



$$\mathcal{L}_t(\pi) = \sum_{i=1}^m \|\pi(s^i) - \pi^*(s^i)\|_2^2$$

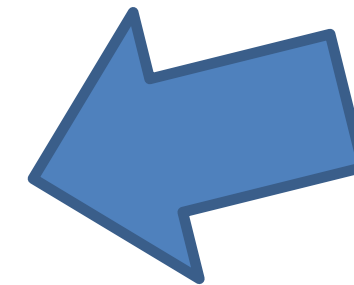
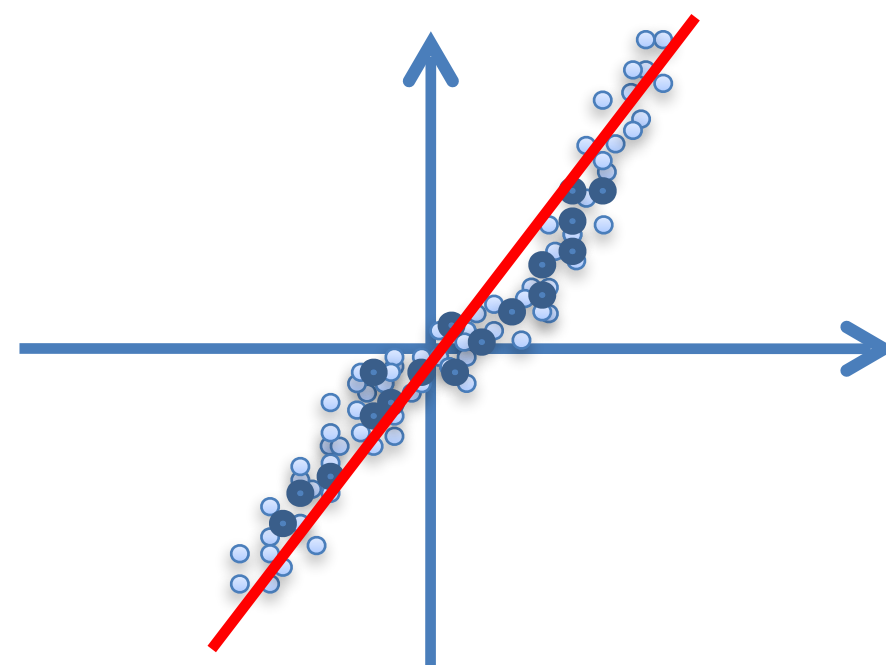
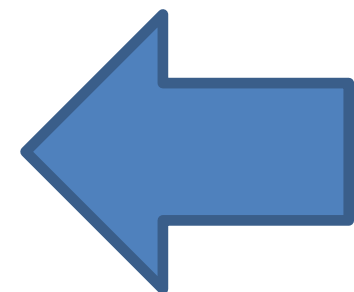
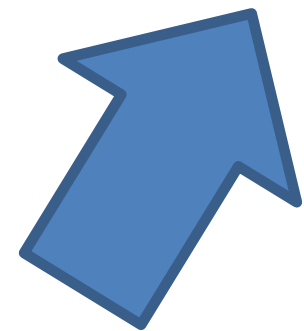


Aggregate Dataset

All previous data

New policy  
 $\pi_n$

Supervised Learning



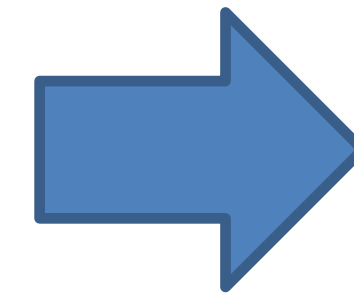
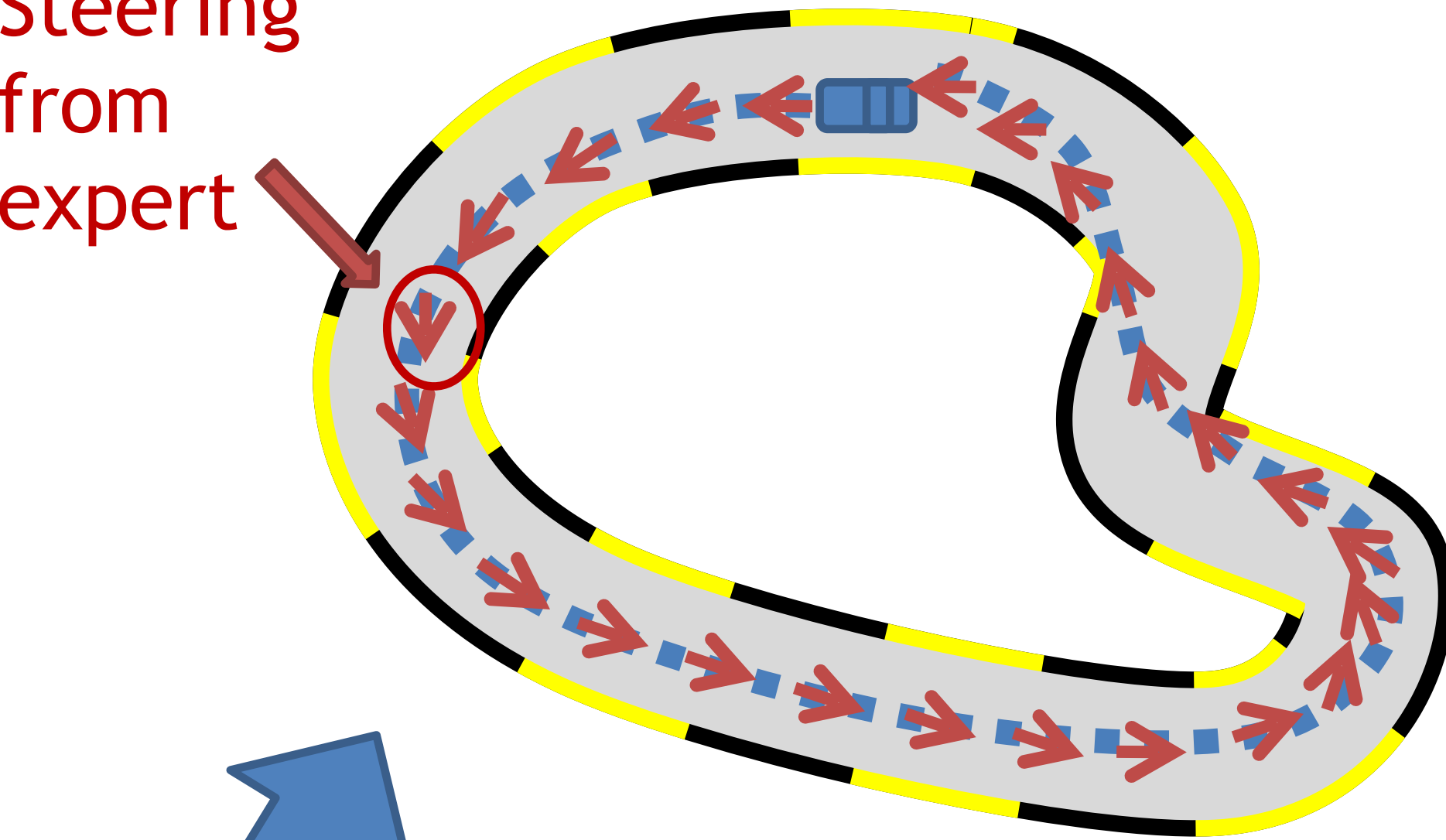


# Dagger Revisit

At iteration t:

New Data

Steering from expert



$$\ell_t(\pi) = \sum_{i=1}^m \|\pi(s^i) - \pi^*(s^i)\|_2^2$$



Aggregate Dataset

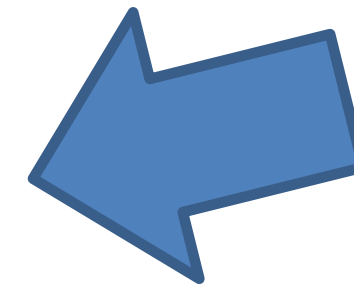
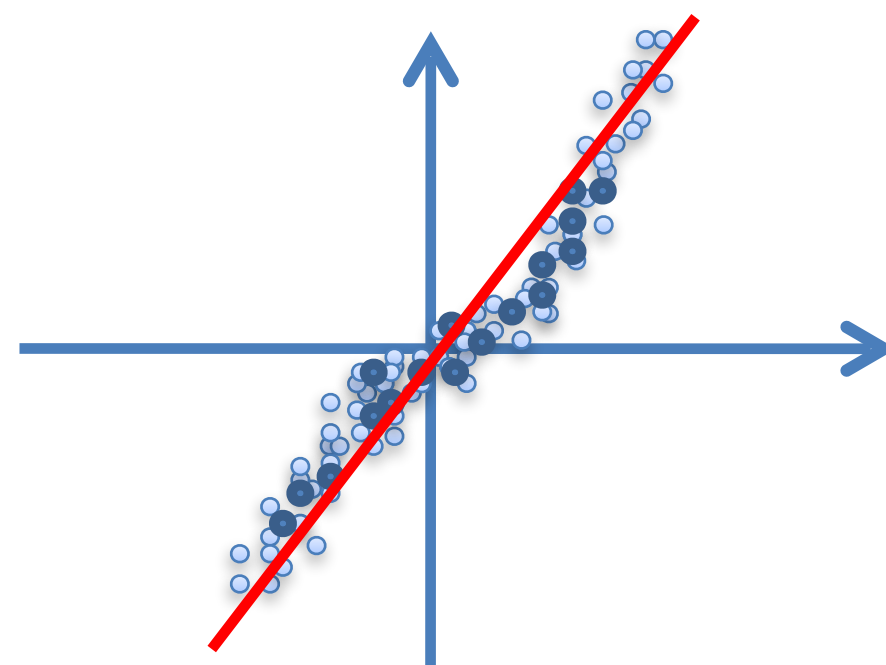
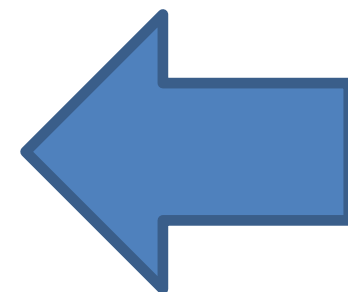
All previous data

$$\sum_{i=0}^{t-1} \ell_i(\pi)$$

New policy

$\pi_n$

Supervised Learning

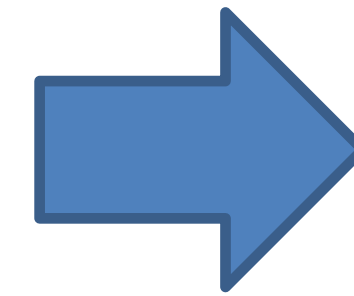
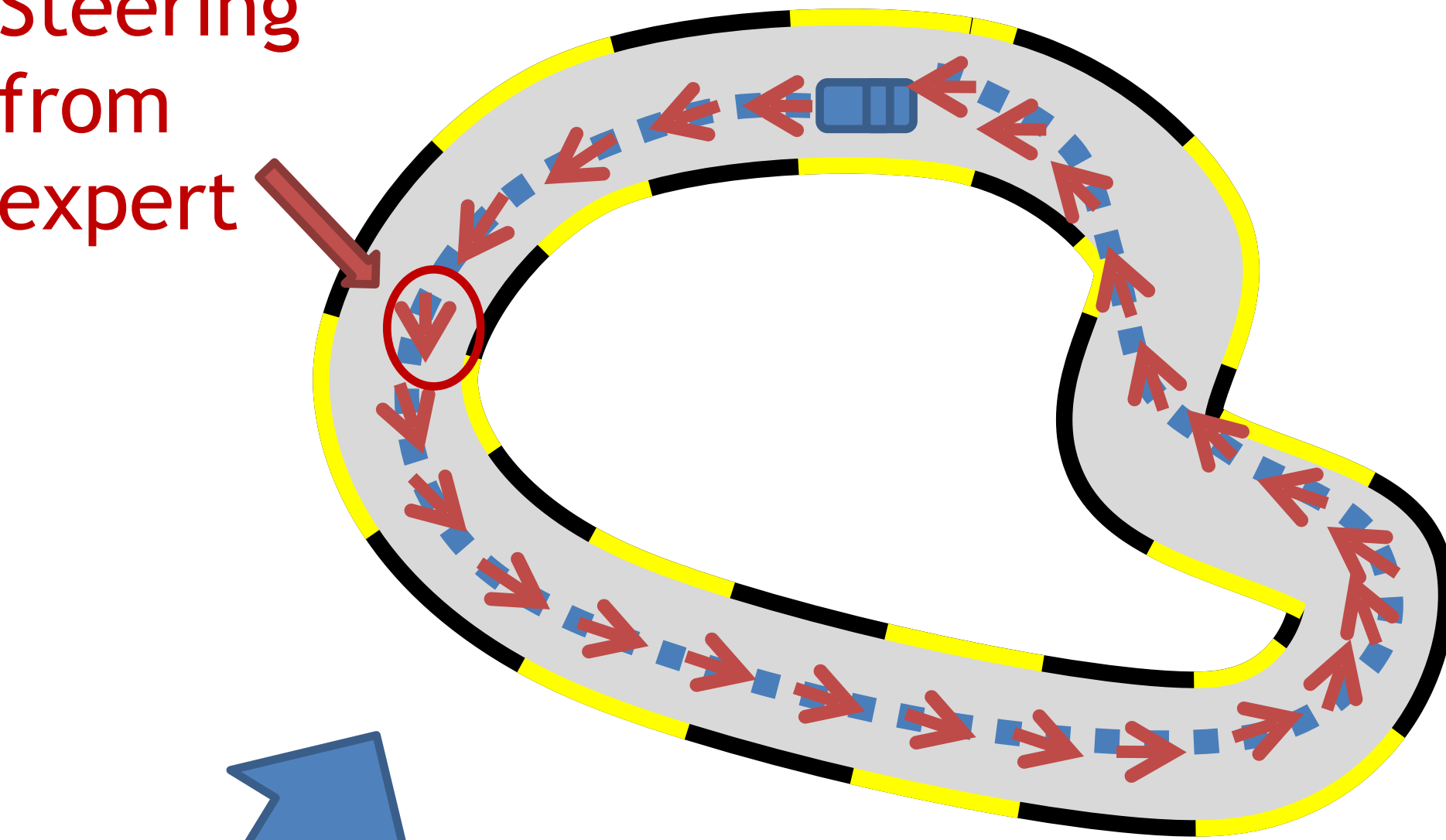


# Dagger Revisit

At iteration t:

New Data

Steering  
from  
expert



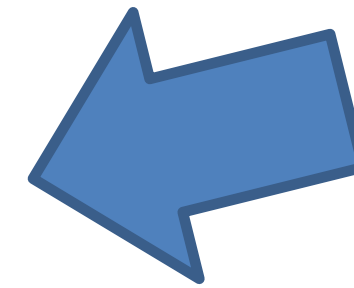
$$\ell_t(\pi) = \sum_{i=1}^m \|\pi(s^i) - \pi^*(s^i)\|_2^2$$



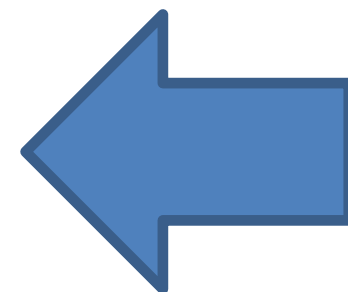
Aggregate  
Dataset

All previous data

$$\sum_{i=0}^{t-1} \ell_i(\pi)$$



New policy  
 $\pi_n$



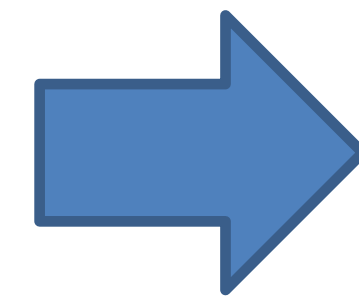
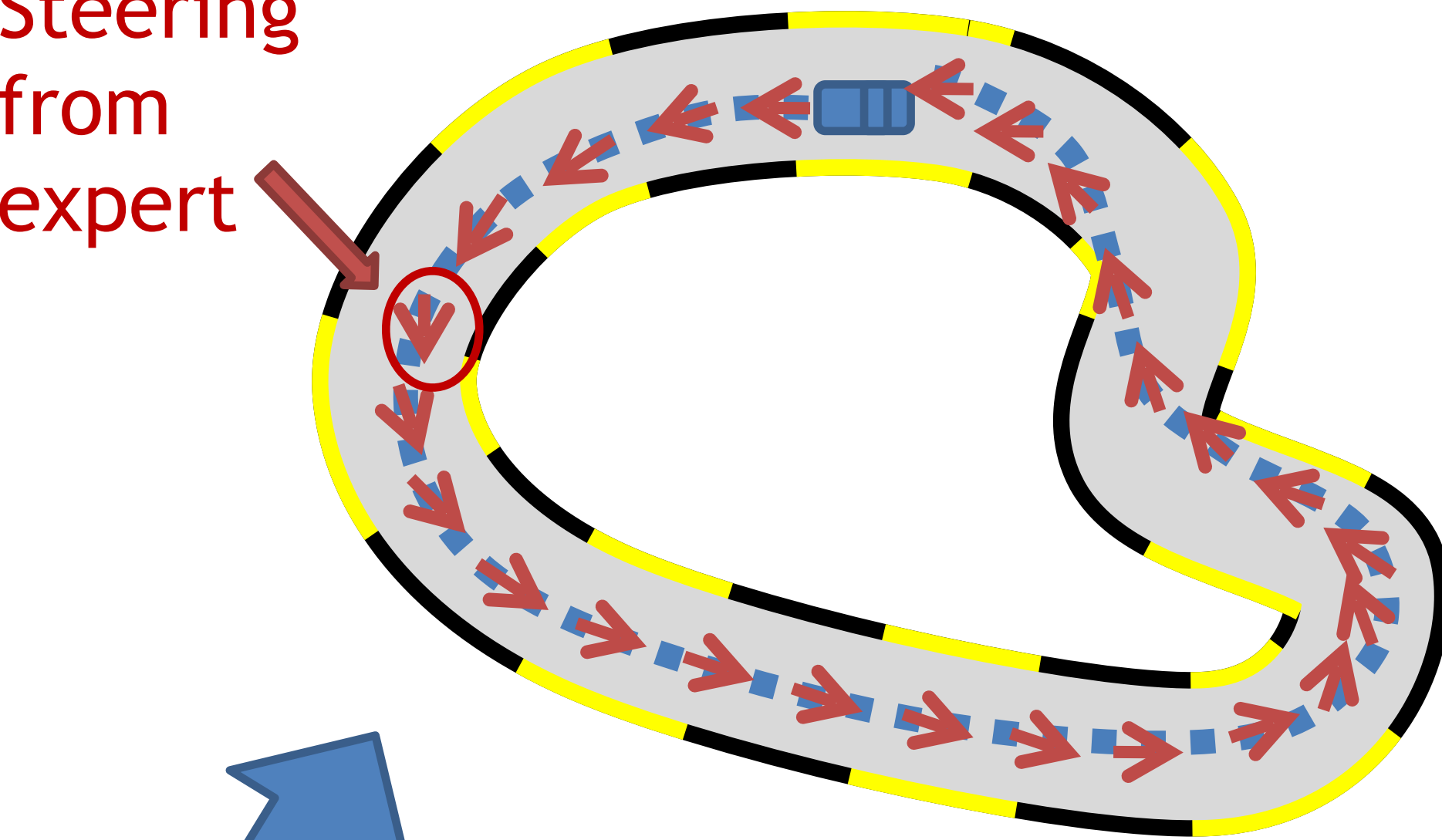
Supervised Learning

# Dagger Revisit

At iteration t:

New Data

Steering from expert



$$\ell_t(\pi) = \sum_{i=1}^m \|\pi(s^i) - \pi^*(s^i)\|_2^2$$



Aggregate Dataset

All previous data

$$\sum_{i=0}^{t-1} \ell_i(\pi)$$

New policy  
 $\pi_n$

$$\pi_t = \arg \min_{\pi} \sum_{i=0}^{t-1} \ell_i(\pi) + \lambda R(\pi)$$

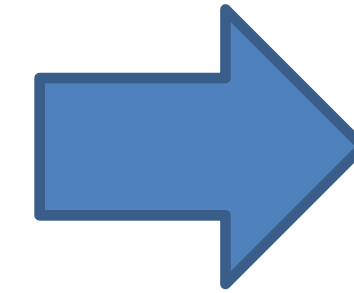
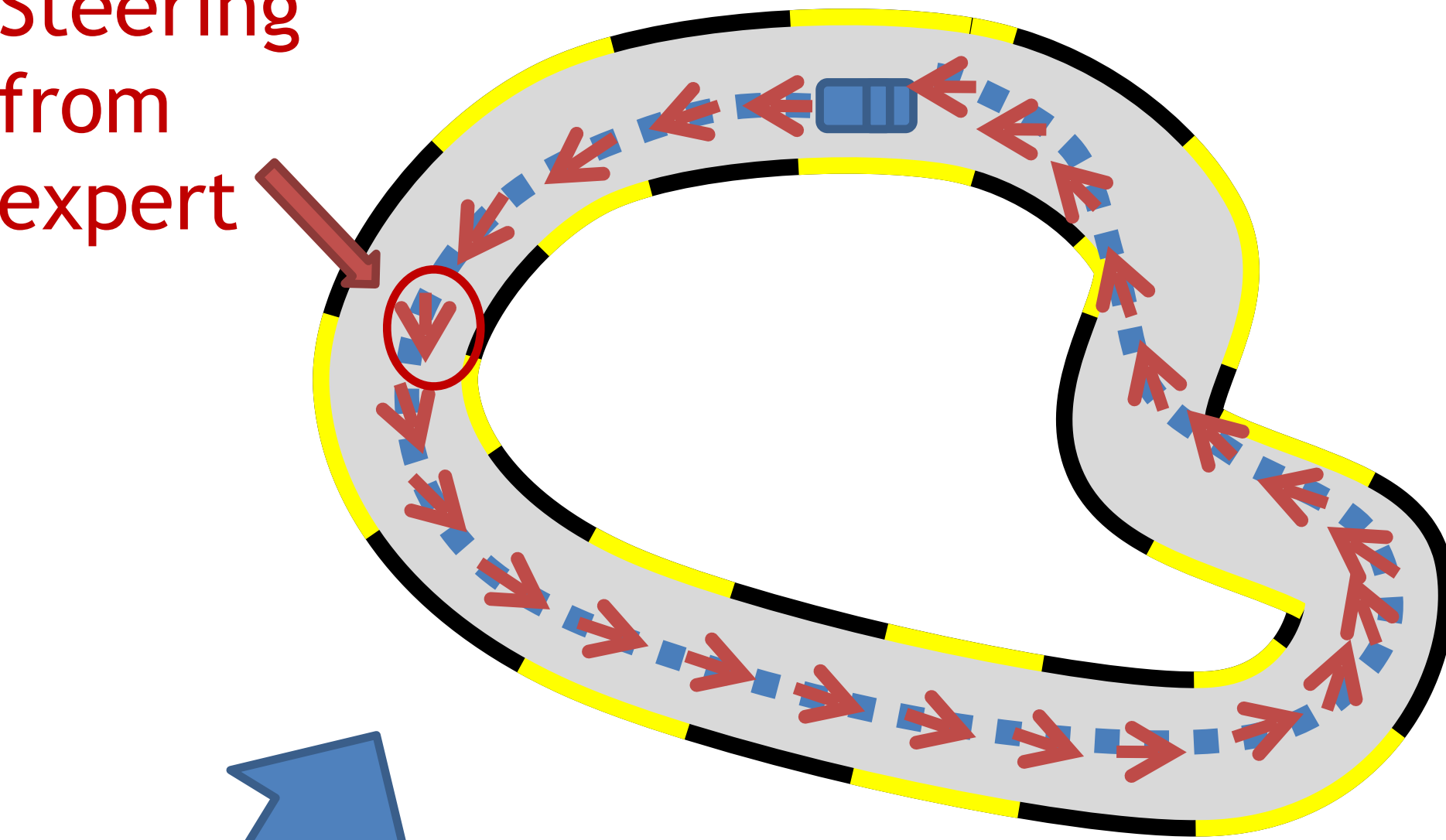
Supervised Learning

# Dagger Revisit

At iteration t:

New Data

Steering from expert



$$\ell_t(\pi) = \sum_{i=1}^m \|\pi(s^i) - \pi^*(s^i)\|_2^2$$



Aggregate Dataset

All previous data

$$\sum_{i=0}^{t-1} \ell_i(\pi)$$

New policy  
 $\pi_n$

$$\pi_t = \arg \min_{\pi} \sum_{i=0}^{t-1} \ell_i(\pi) + \lambda R(\pi)$$

Supervised Learning

Data Aggregation = Follow-the-Regularized-Leader Online Learner

# Summary for Today

## 1. The DAgger algorithm

Initialize  $\pi^0$ , and dataset  $\mathcal{D} = \emptyset$

For  $t = 0 \rightarrow T - 1$ :

1. W/  $\pi^t$ , generate dataset  $\mathcal{D}^t = \{s_i, a_i^\star\}$ ,  $s_i \sim d_\mu^{\pi^t}$ ,  $a_i^\star = \pi^\star(s_i)$

2. **Data aggregation:**  $\mathcal{D} = \mathcal{D} + \mathcal{D}^t$

3. **Update policy via Supervised-Learning:**  $\pi^{t+1} = \text{SL}(\mathcal{D})$

# Summary for Today

## 1. The DAgger algorithm

Initialize  $\pi^0$ , and dataset  $\mathcal{D} = \emptyset$

For  $t = 0 \rightarrow T - 1$ :

1. W/  $\pi^t$ , generate dataset  $\mathcal{D}^t = \{s_i, a_i^\star\}$ ,  $s_i \sim d_\mu^{\pi^t}$ ,  $a_i^\star = \pi^\star(s_i)$

2. **Data aggregation:**  $\mathcal{D} = \mathcal{D} + \mathcal{D}^t$

3. **Update policy via Supervised-Learning:**  $\pi^{t+1} = \text{SL}(\mathcal{D})$

2. We can see that DAgger is essentially an online-learning algorithm (FTRL)