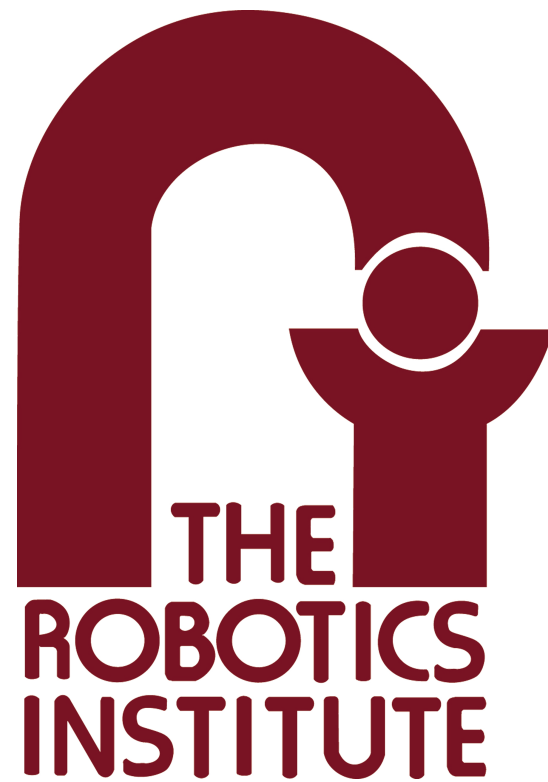


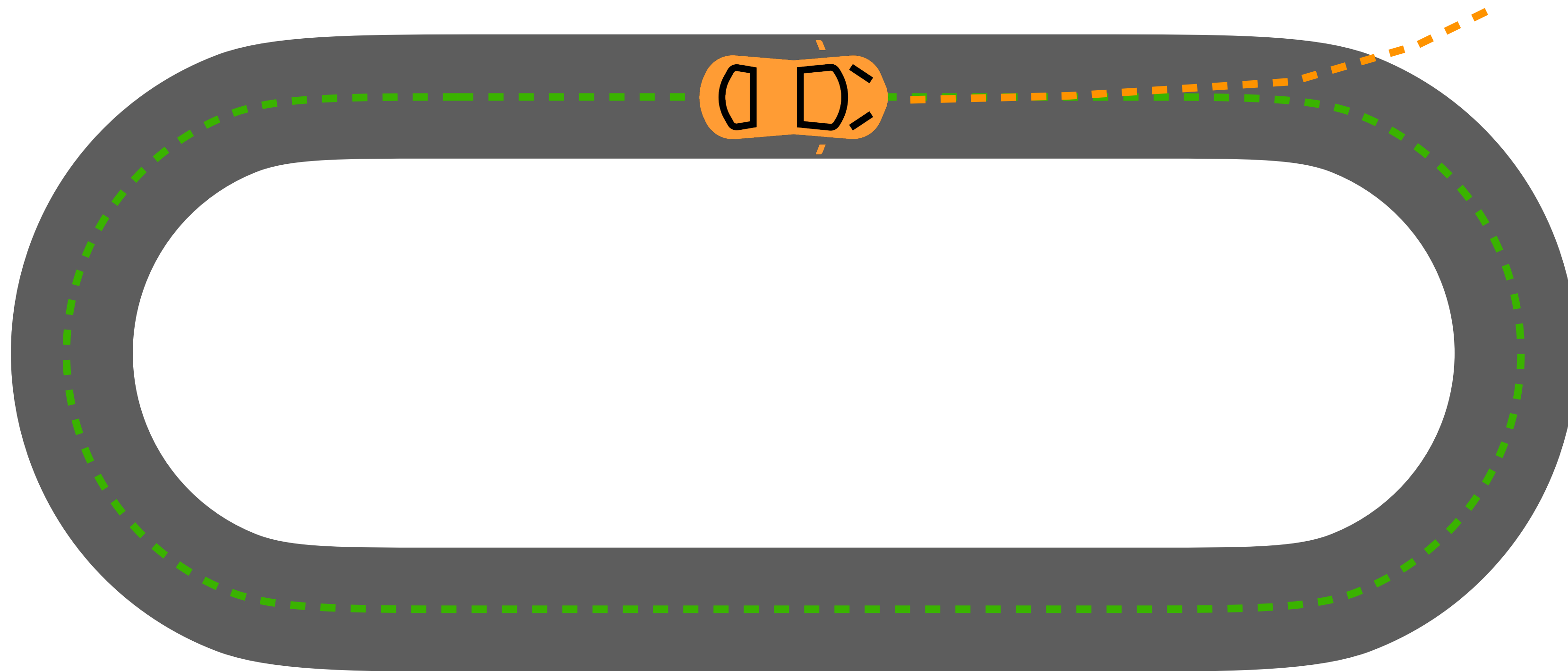
# Of Moments and Matching: A Game-Theoretic Framework for Closing the Imitation Gap

*Gokul Swamy, Sanjiban Choudhury, Drew Bagnell, Steven Wu*

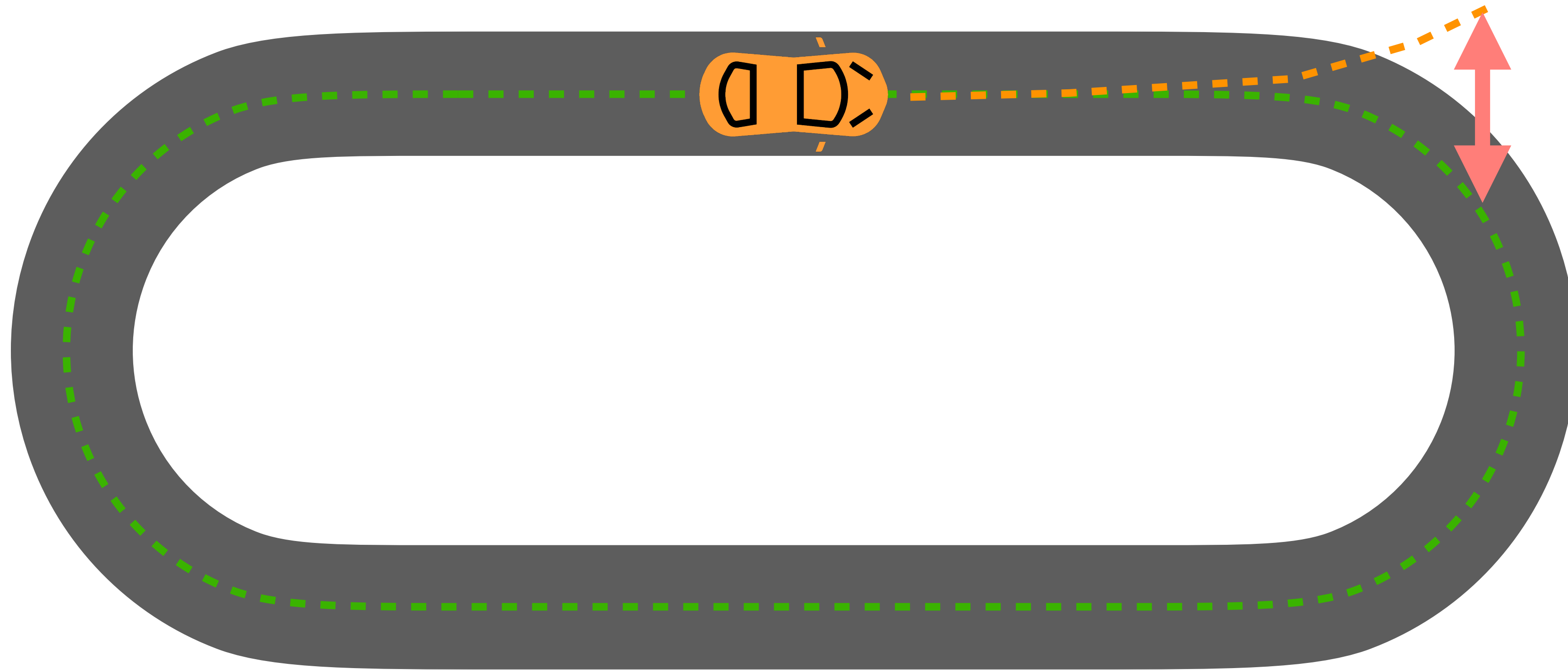


Aurora

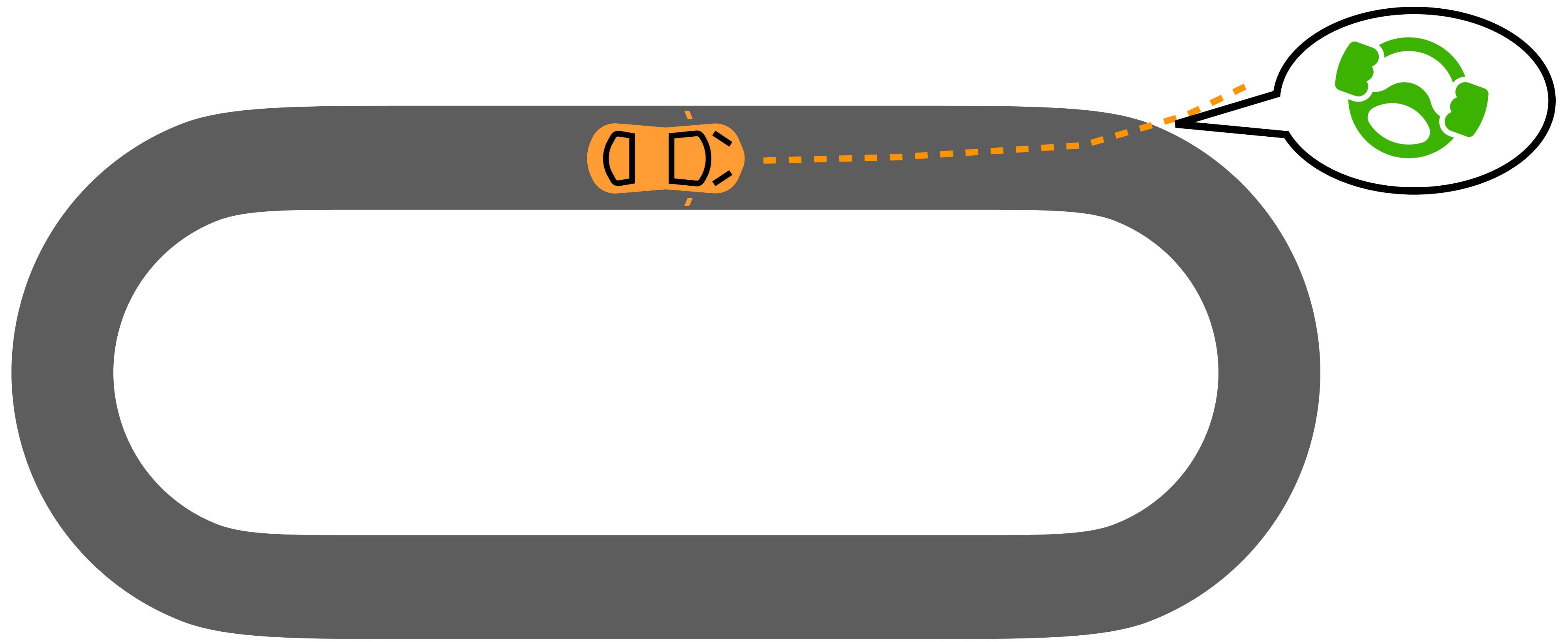




$$\{s_1 \dots s_n\} \mapsto \{a_1 \dots a_n\}$$



$$\begin{array}{ccc} \{s_1 \dots s_n\} & \longleftrightarrow & \{s_1 \dots s_n\} \\ \{a_1 \dots a_n\} & & \{a_1 \dots a_n\} \end{array}$$



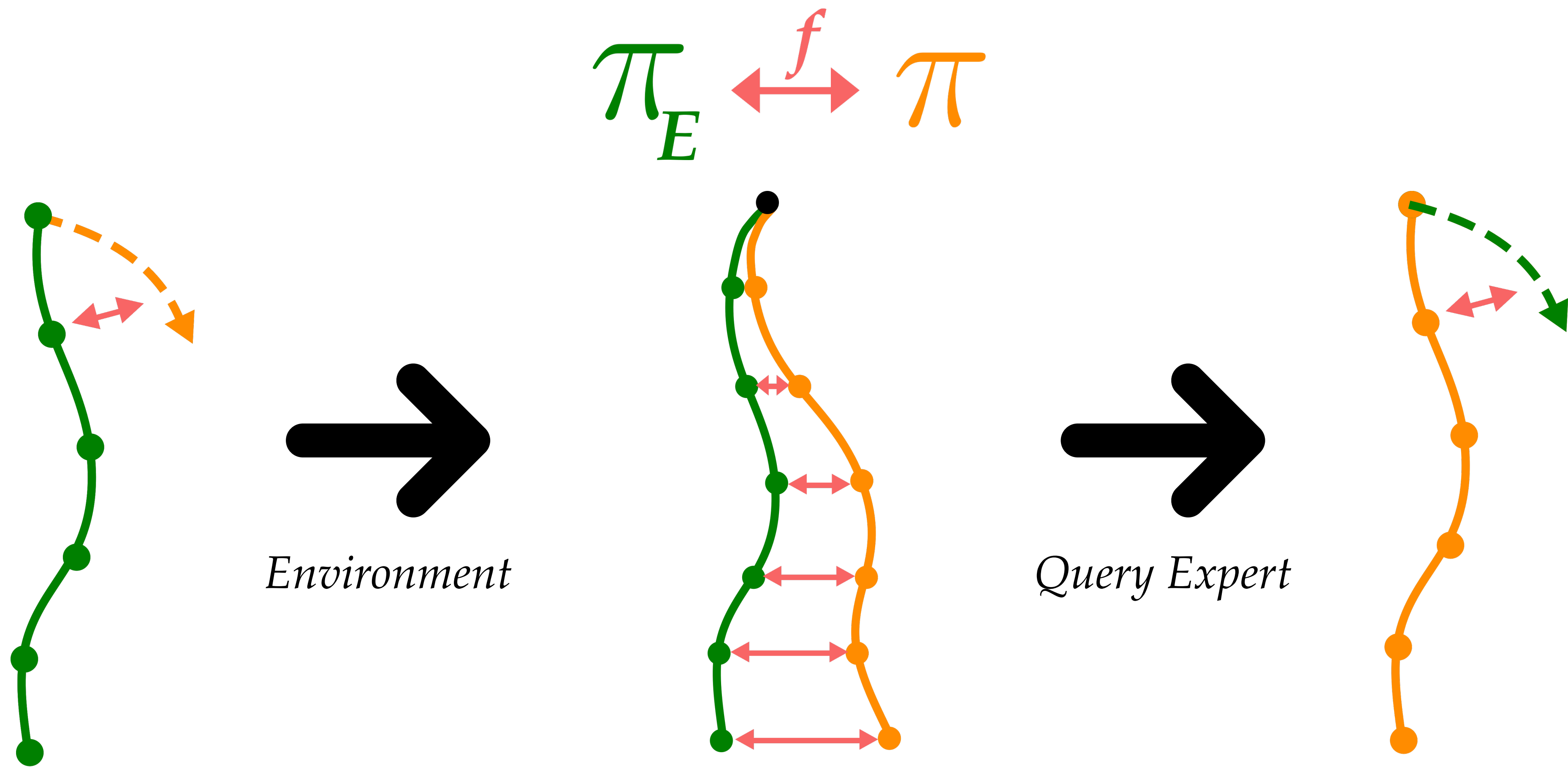
$$\{s_1 \dots s_n\} \mapsto \{a_1 \dots a_n\}$$

*Q1:* How well can we expect a learned policy to do in each of these three settings?

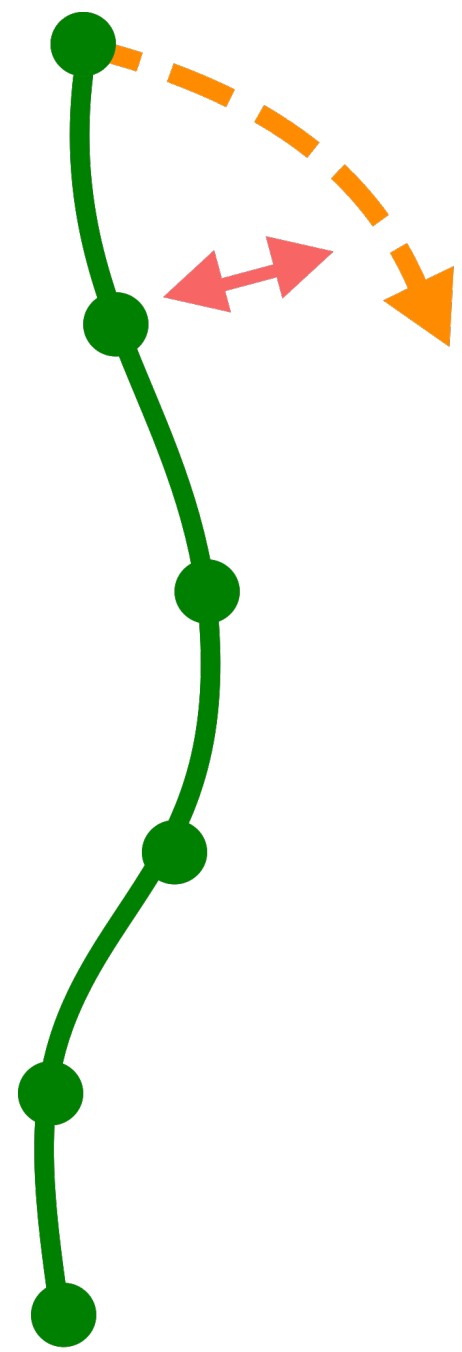
$\min_{\pi}$

$$J(\pi_E) - J(\pi)$$

$$\min_{\pi} \max_f f(\pi_E) - f(\pi)$$

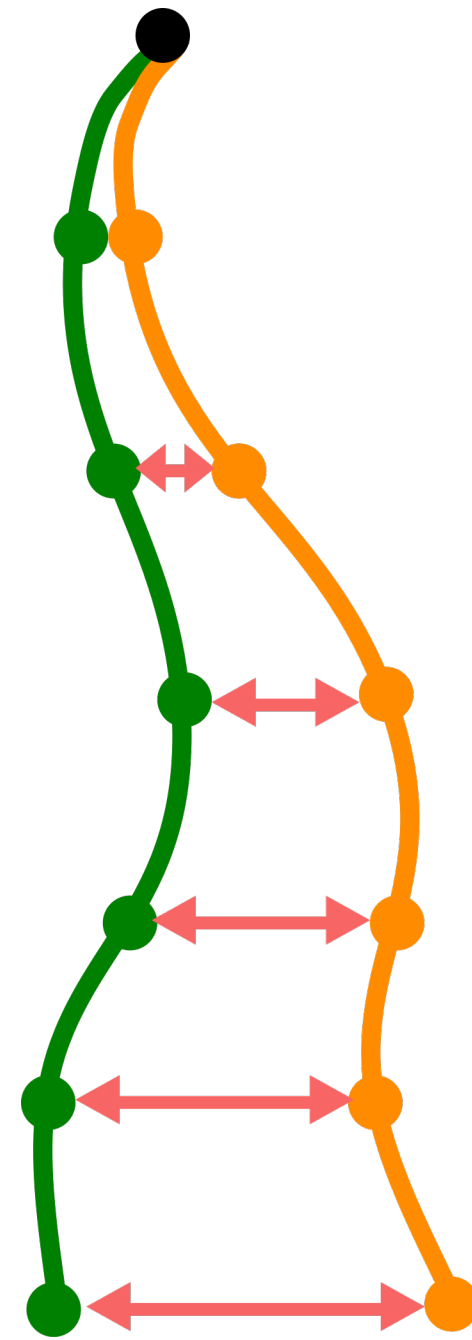


*Key Insight: Each of these classes of approaches corresponds to solving a game with a different class of discriminators. Stronger feedback leads to more powerful discriminators and a tighter performance bound.*

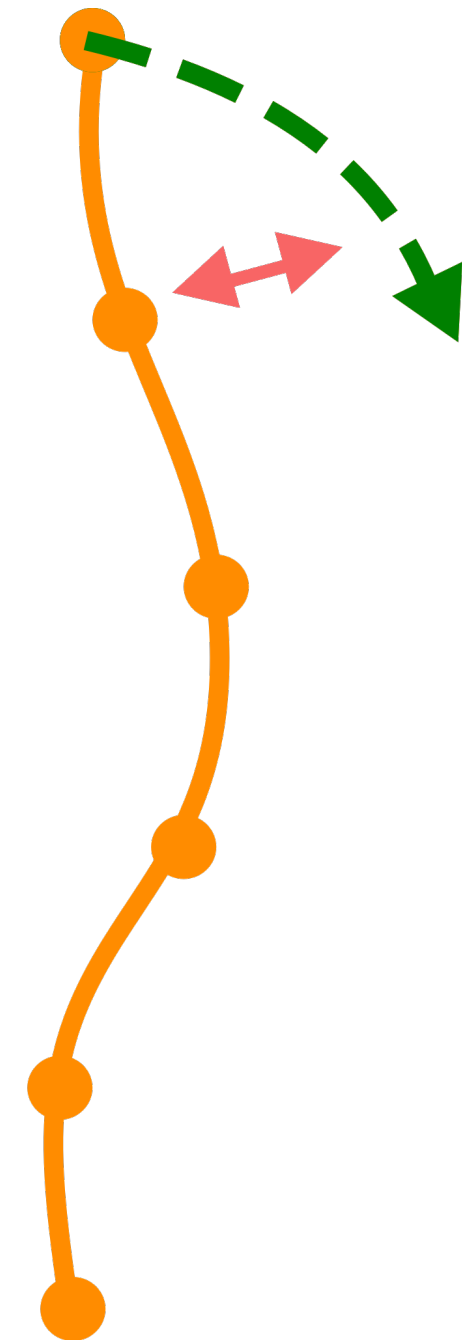


Environment

$$\pi_E \xleftrightarrow{f} \pi$$



Query Expert



$$J(\pi_E) - J(\pi) \leq O(\epsilon T^2)$$

*Behavioral Cloning,  
ValueDICE,*

...

$$J(\pi_E) - J(\pi) \leq O(\epsilon T)$$

*GAIL, SQIL, MaxEnt  
IRL, LEARCH, Max  
Margin Planning*

...

$$J(\pi_E) - J(\pi) \leq O(\epsilon HT)$$

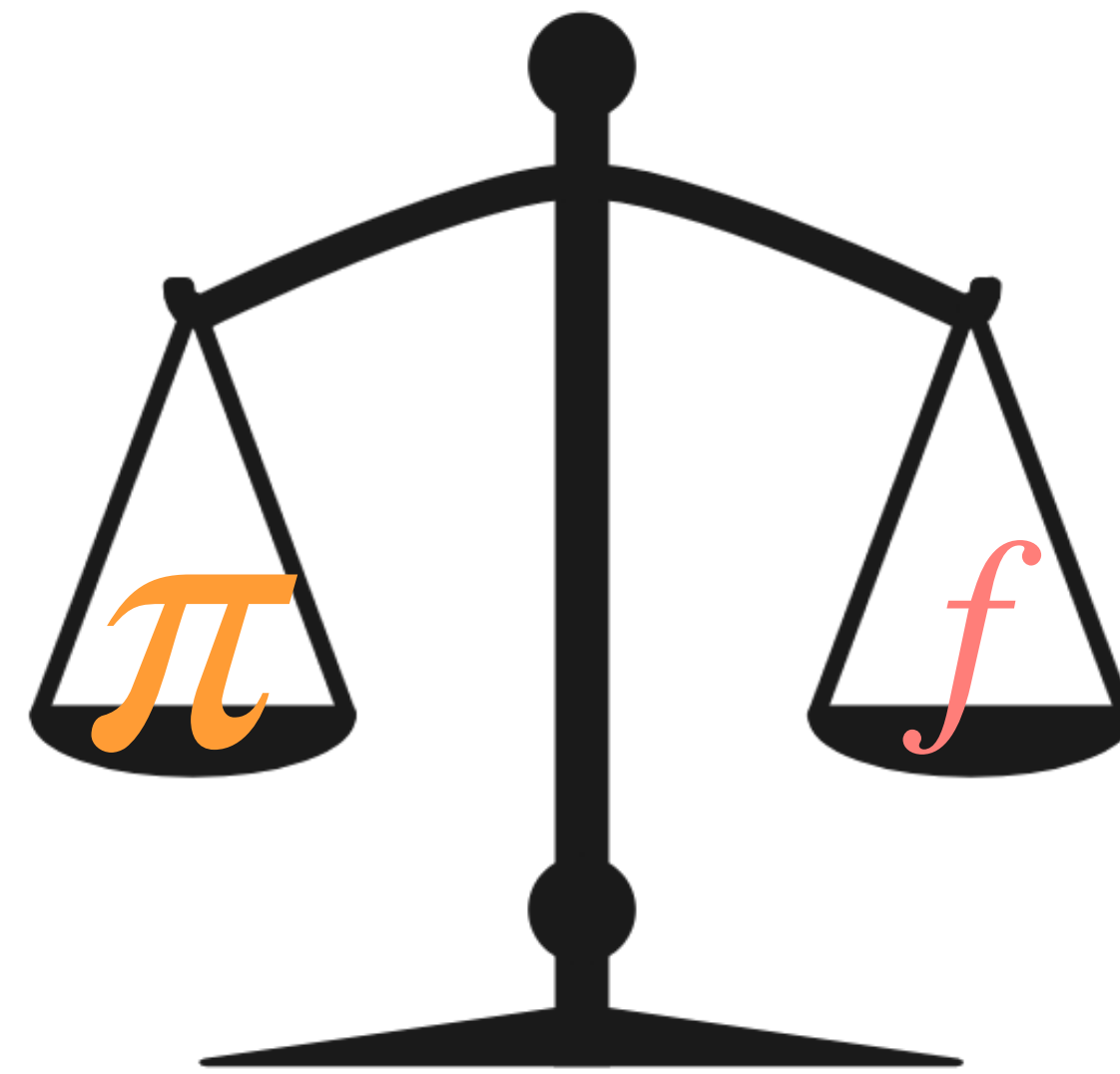
*DAgger, Guided Policy  
Search, iFAIL*

...



**Q2:** How can we efficiently find a performant policy in each of these settings?

**A:**

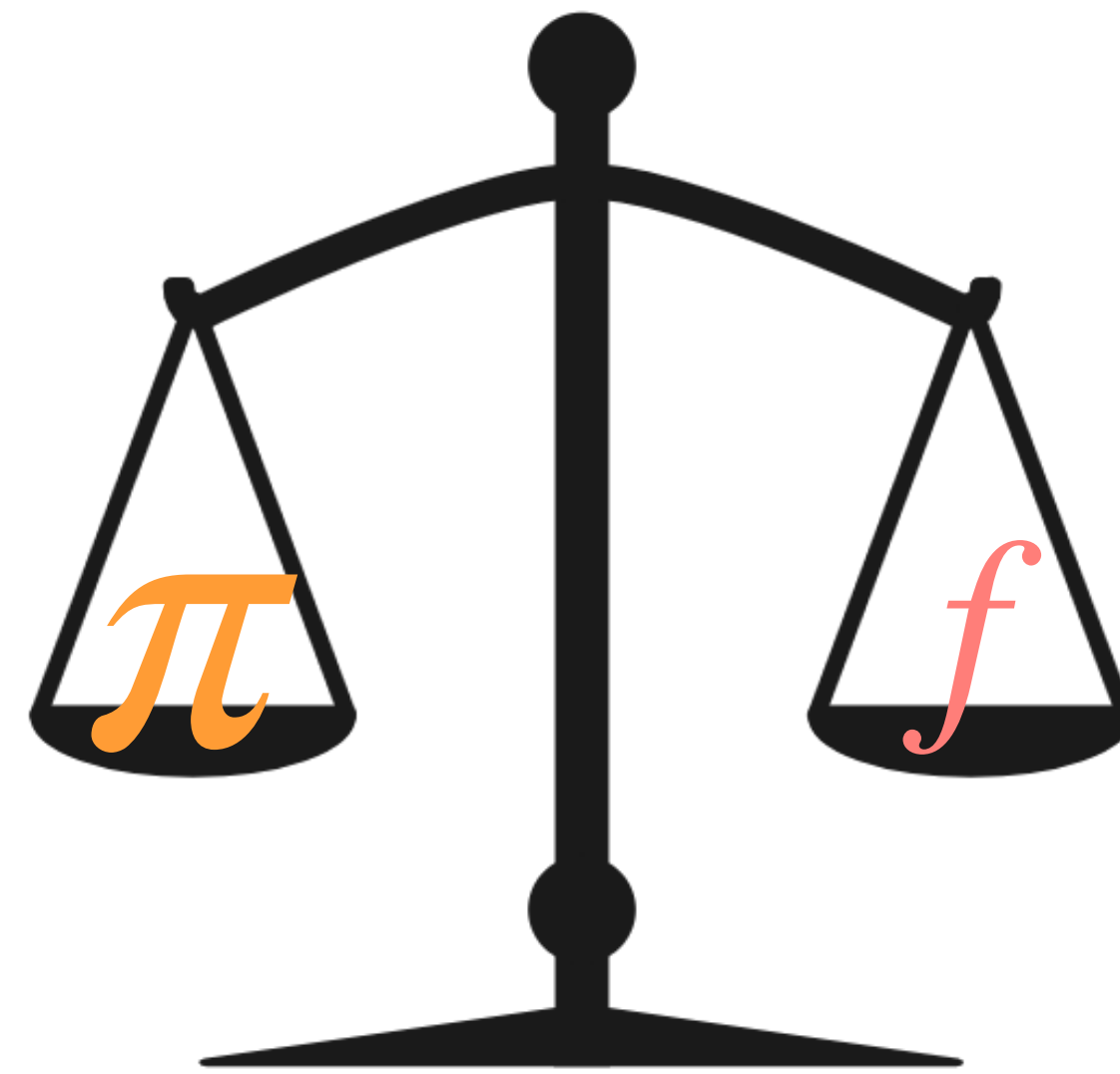


= *No Regret on  $\pi$*  vs. *Best Response on  $f$*

= *Best Response on  $\pi$*  vs. *No Regret on  $f$*

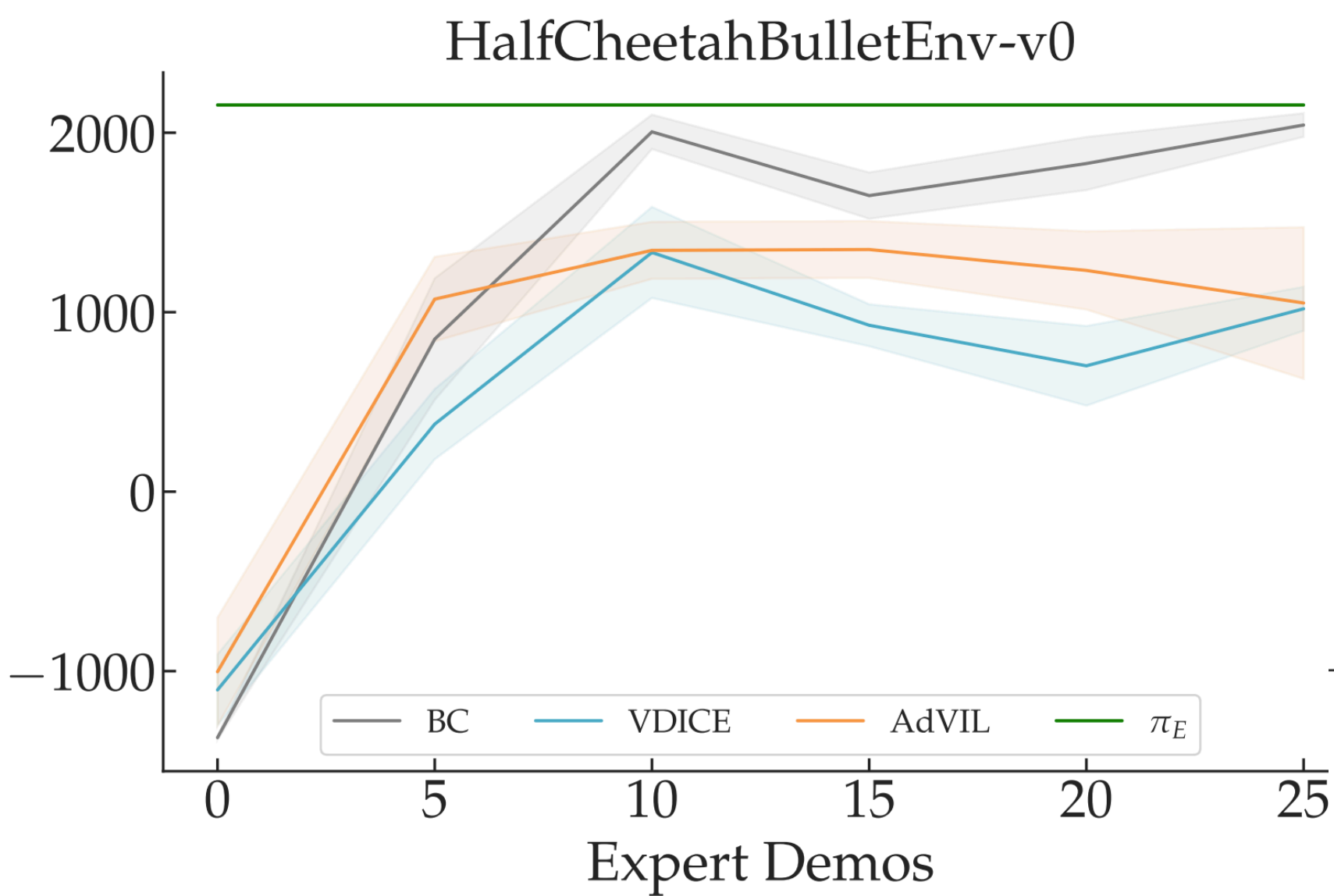
**Q2:** How can we efficiently find a performant policy in each of these settings?

**A:**

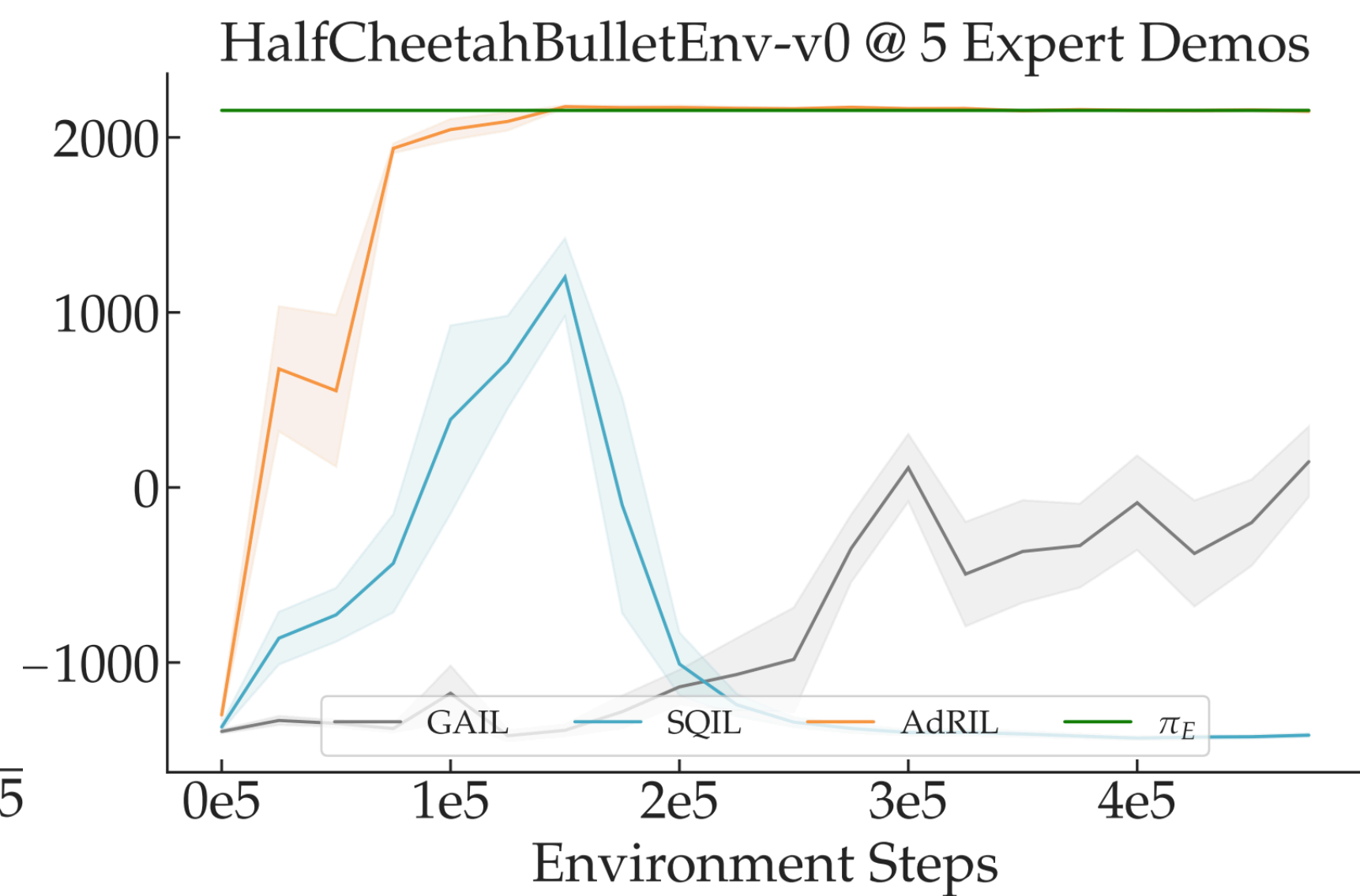


= \_\_\_\_\_ on  $\pi$  vs. *Best Response on  $f$*   
= *Best Response on  $\pi$*  vs. \_\_\_\_\_ on  $f$

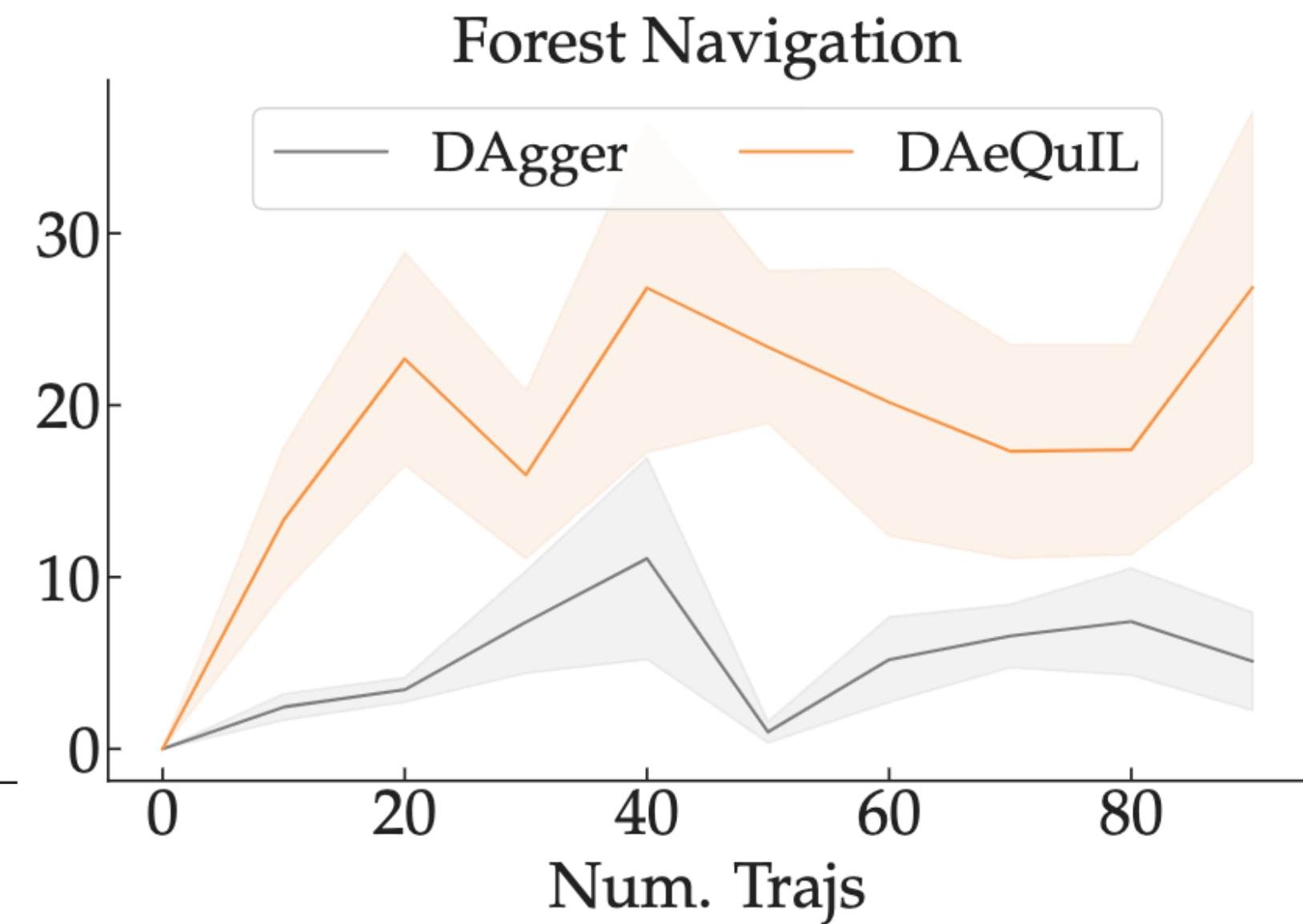
# AdVIL



# AdRIL

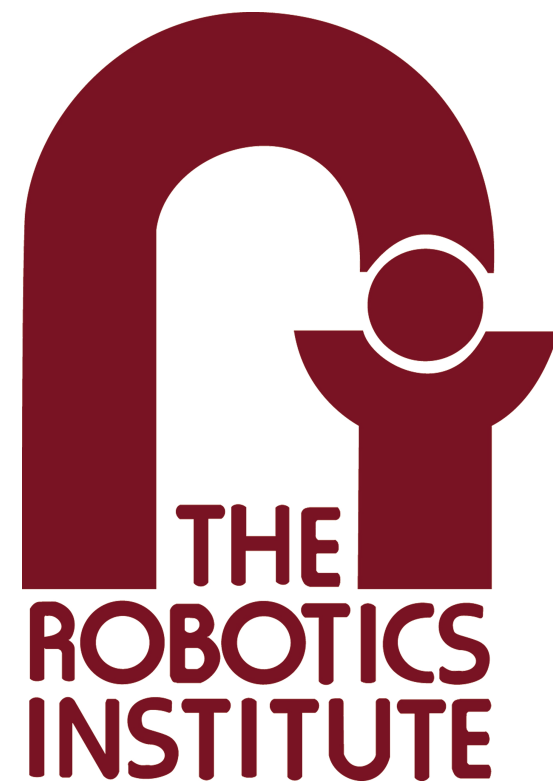


# DAeQuIL



# Of Moments and Matching: A Game-Theoretic Framework for Closing the Imitation Gap

*Gokul Swamy, Sanjiban Choudhury, Drew Bagnell, Steven Wu*



Aurora



<https://gokul.dev/mmil/>