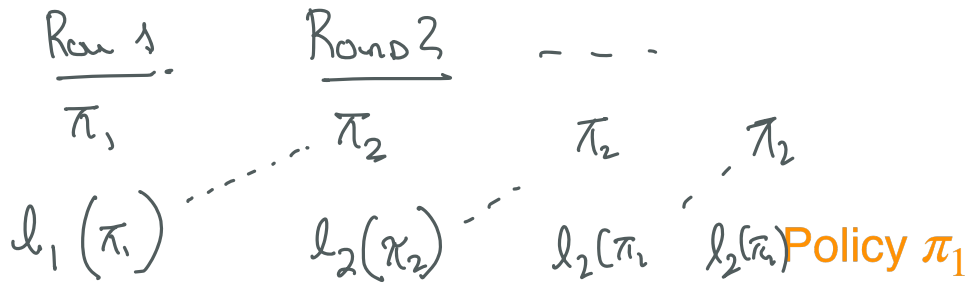
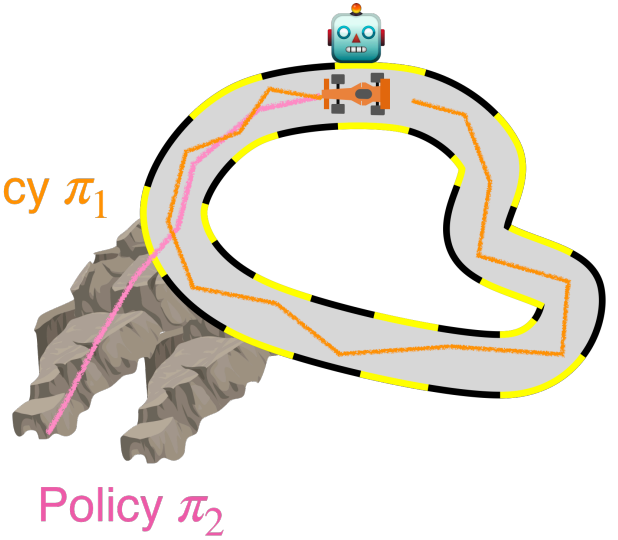


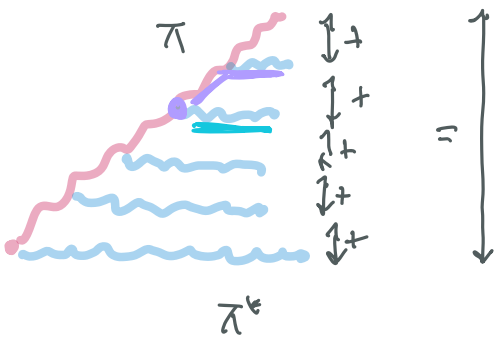
# WHAT WOULD DAGGER RETURN?



$$l_i(\pi)$$
$$= \mathbb{E}_{s \sim d_{\pi_i}} \mathbb{1}(\pi(s) \neq a^*)$$

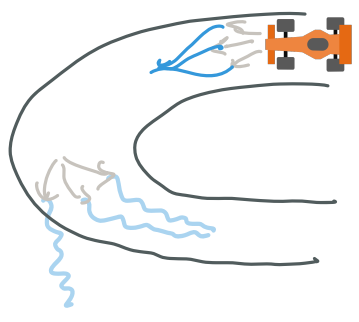


# PERFORMANCE DIFFERENCE LEMMA



$$\begin{aligned}
 & J(\pi) - J(\pi^*) \\
 &= \sum_{t=1}^T \mathbb{E}_{s_t \sim d_{\pi}^t} \left[ \underbrace{Q(s_t, \pi(s_t))}_{\pi^*} - \underbrace{Q(s_t, \pi^*(s_t))}_{\pi^*} \right] \\
 & \quad \underbrace{\hspace{10em}}_{A(s, \pi(s))}
 \end{aligned}$$

$s_t \sim d_{\pi}^t$   
 on-policy



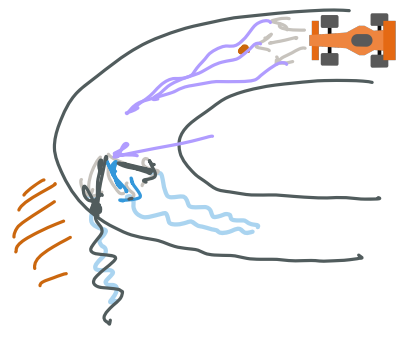
NOT ALL ERRORS ARE EQUAL

FDL

$\pi^*$   
 $A(s, \pi(s))$

~~$O(\epsilon T)$~~   $\Downarrow$   $\mathbb{I}(a^* \neq \pi(s))$

$\begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$



$\begin{bmatrix} 10000.0 \\ 0 \\ \epsilon \end{bmatrix}$

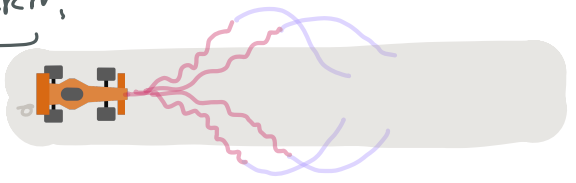
$A^{\pi^*}(s, a_1) = \left[ \frac{Q^{\pi^*}(s, a_1) - Q^{\pi^*}(s, a_2)}{10000.0} \right]$

$A^{\pi^*}(s, \pi(s)) \leq \underbrace{\mathbb{I}(a^* \neq \pi(s))}_{0-1} \cdot \underbrace{\max_a A^{\pi^*}(s, a)}_M$

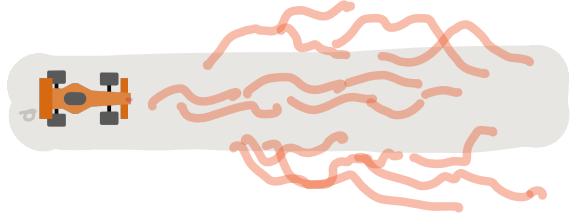
~~$O(\epsilon T)$~~   $\rightarrow O(\epsilon T \cdot M)$   
 $\approx O(\epsilon T^2)$

"RECOVERABILITY"

# WHAT POLICY WOULD HG-DAGGER RETURN?



POLICY 1:  
GOOD: STAY IN IN  
THE TRACK  
OK: RECOVERY



BEST IN  
SIGHT  
N → ∞

POLICY 2:  
BAD: STAYING IN  
TRACK  
AWESOME: RECOVERY

