

I wanted to go up,  
but the dynamics  
made me go down!

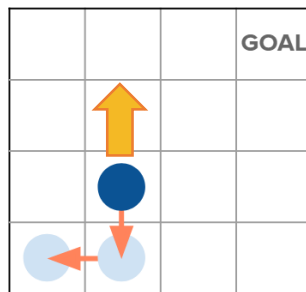
Policy  $\pi_A$



Bad result,  
but right policy

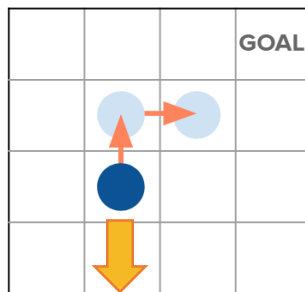
Each execution is not just an action.  
It is a realization of both  $\mathbb{P}$  and  $\pi$ !

Suboptimal  
segment



Equal partial return  
Higher regret

Optimal  
segment



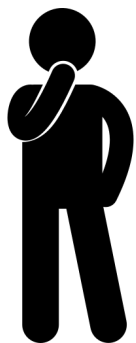
Equal partial return  
Lower regret

Policy  $\pi_B$



Good result,  
but wrong policy

I wanted to go down,  
but the dynamics  
made me go up!



Q. Which way should we prefer?:  
the **result** or the **intended policy**?

A. Evaluate the **result** with **intended policy**

$$Q_*^{\pi}(s_t^{\sigma}, a_t^{\sigma})$$