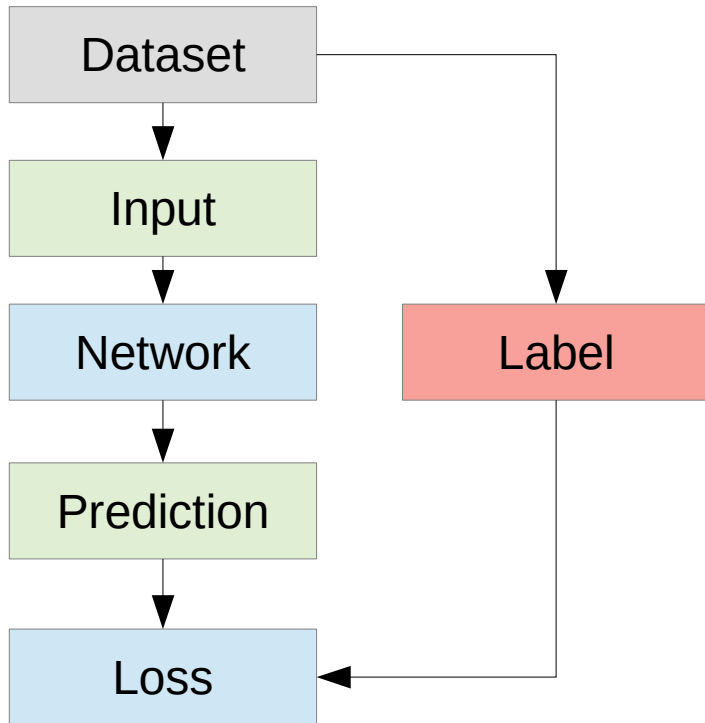




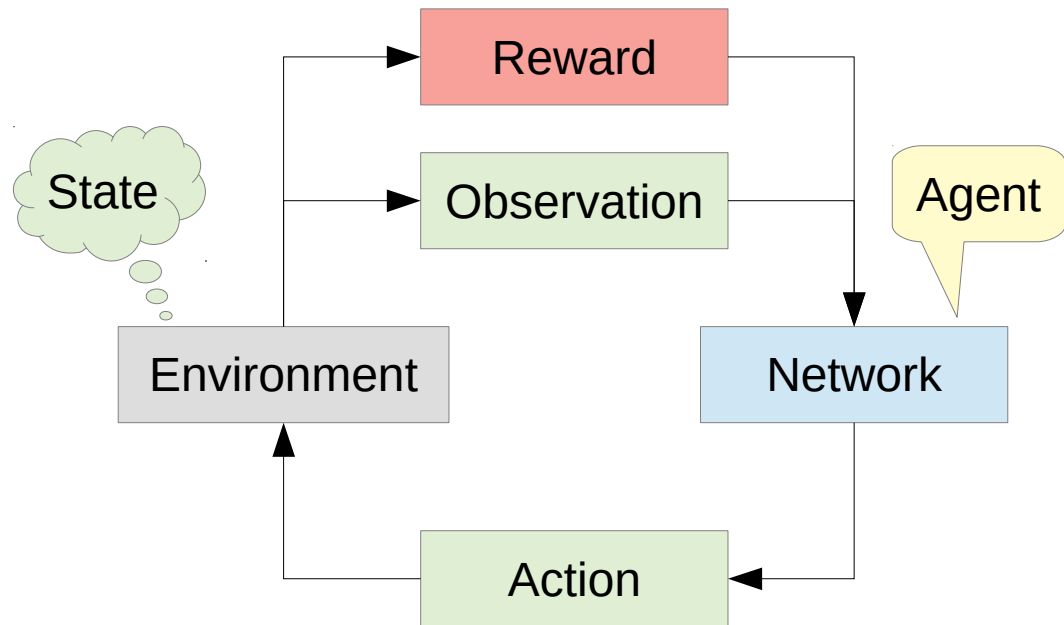
Reinforcement Learning

Acting in an Environment

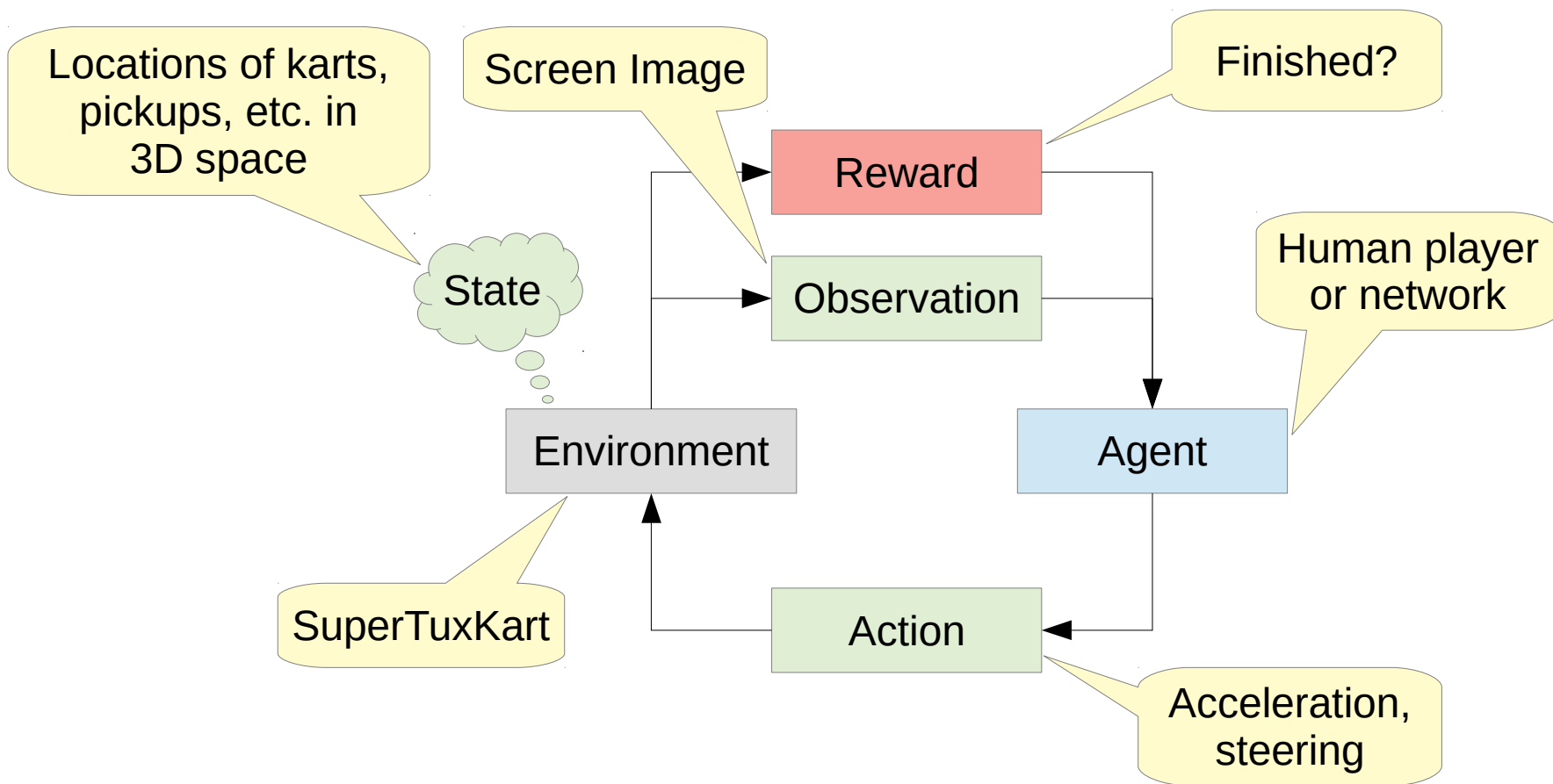
Supervised Learning



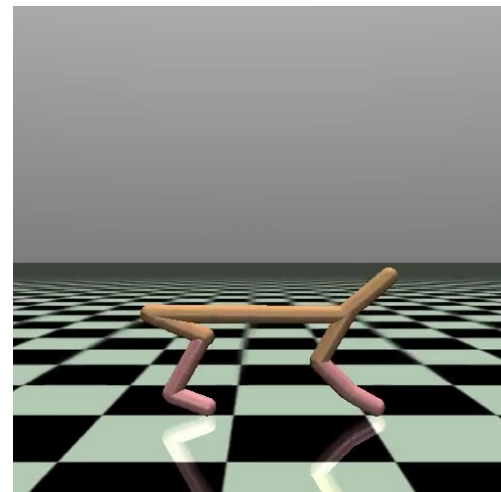
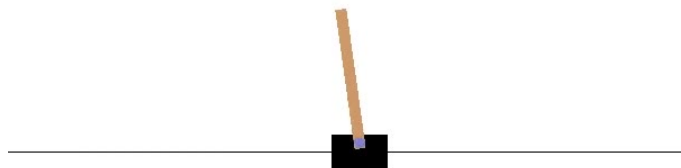
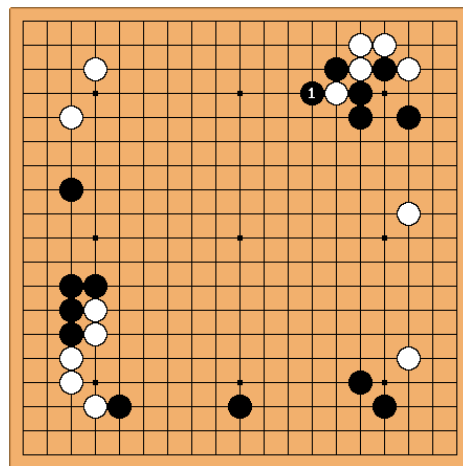
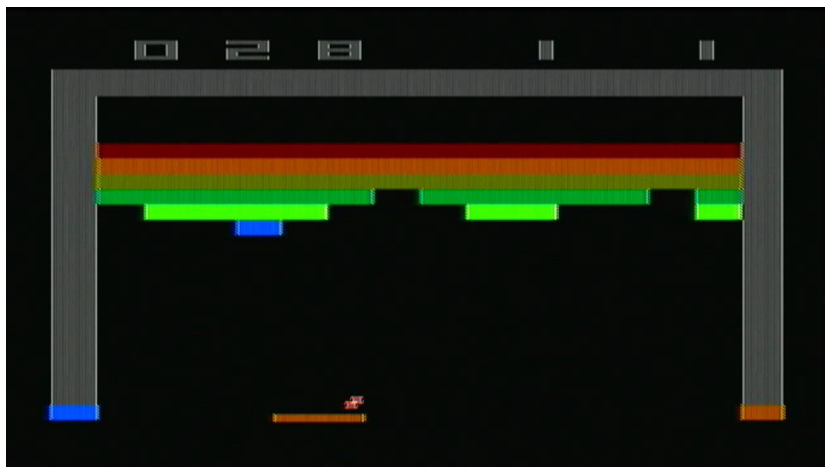
Reinforcement Learning



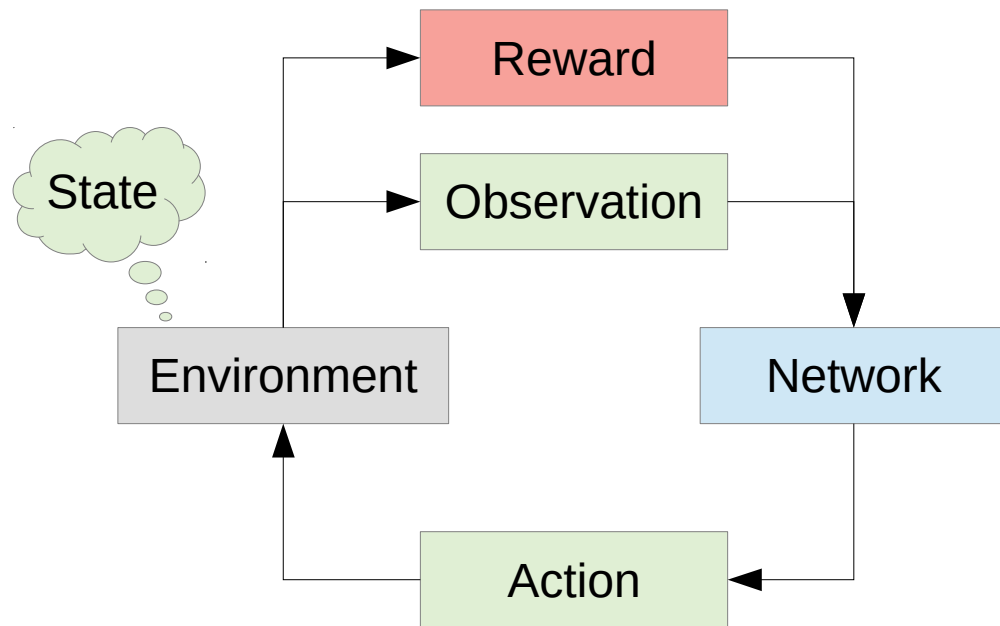
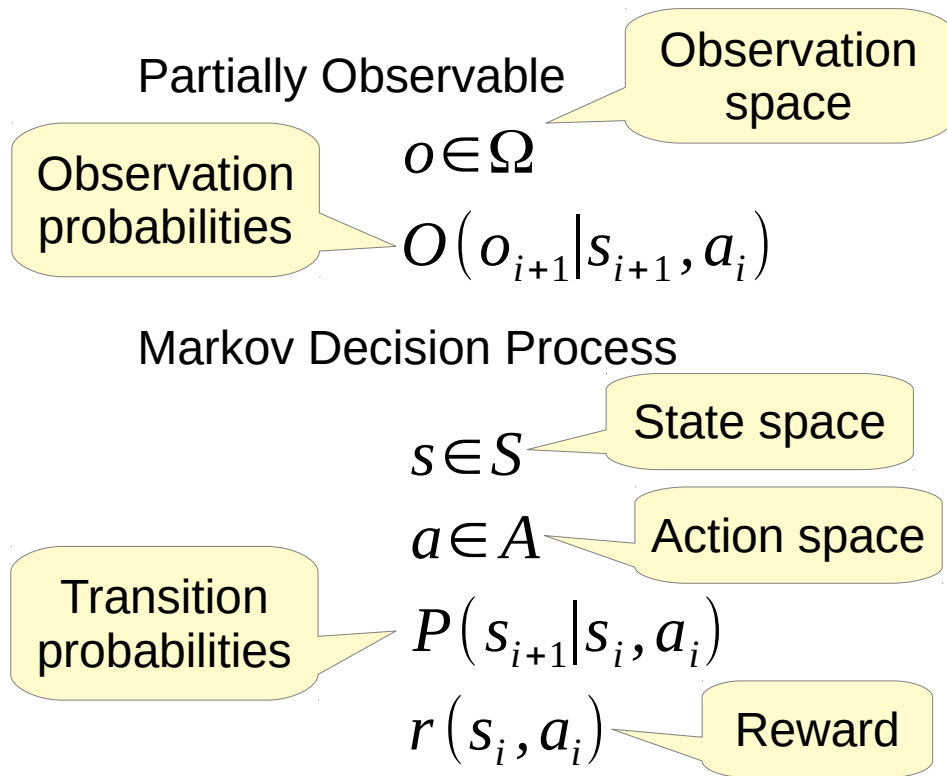
Example: SuperTuxKart



More Examples



Formally: (Partially Observable) Markov Decision Process



Goal

Optimize rewards *in the long term*

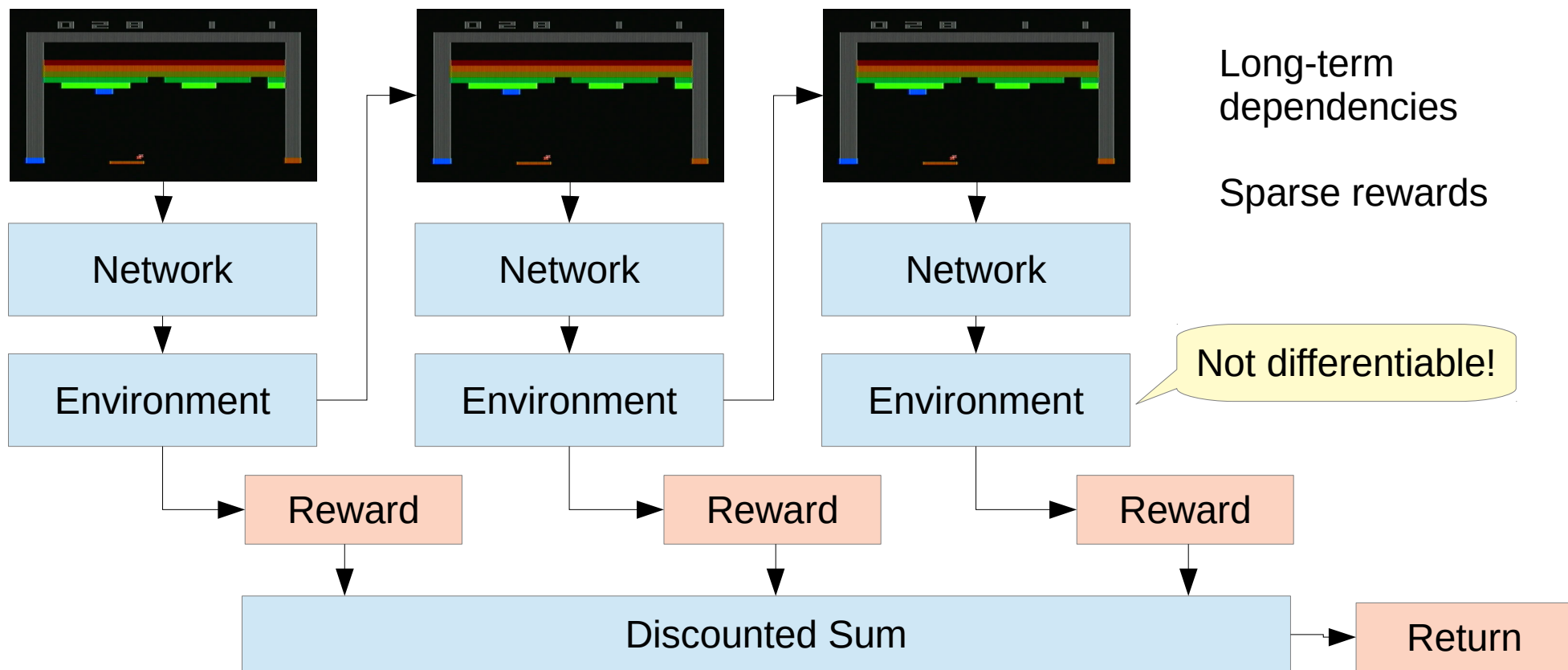
Rollout

Trajectory $\tau = s_0, a_0, s_1, a_1, \dots, s_n, a_n$

Return $R(\tau) = \sum_{i=0}^n \gamma^i r(s_i, a_i)$ “Discounted sum”

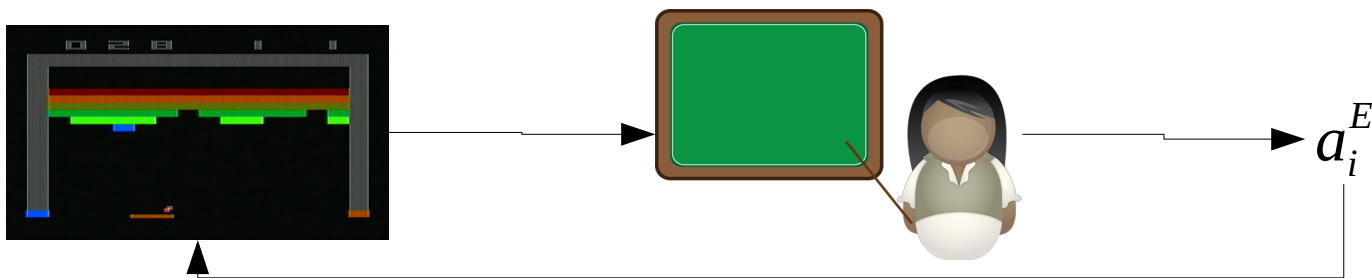
Discount factor

Challenges



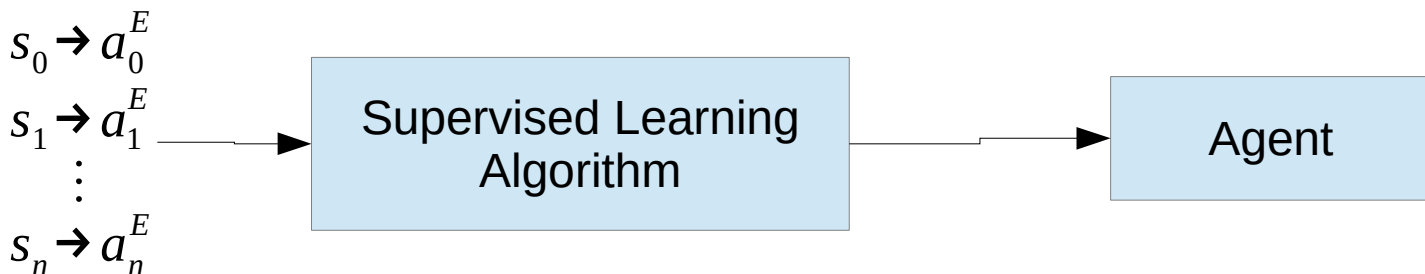
Imitation Learning

- Gather example trajectories from an expert



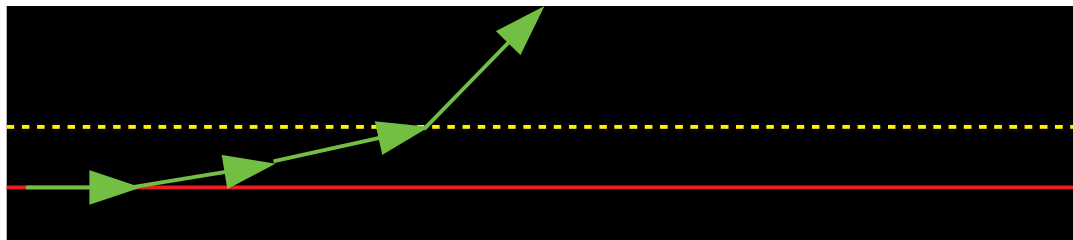
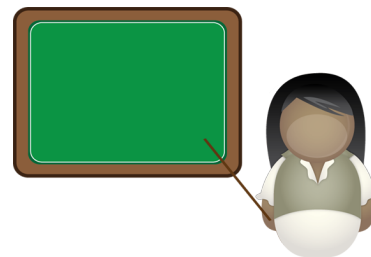
- Use supervised learning

$$\tau^E = s_0, a_0^E, s_1, a_1^E, \dots, s_n, a_n^E$$



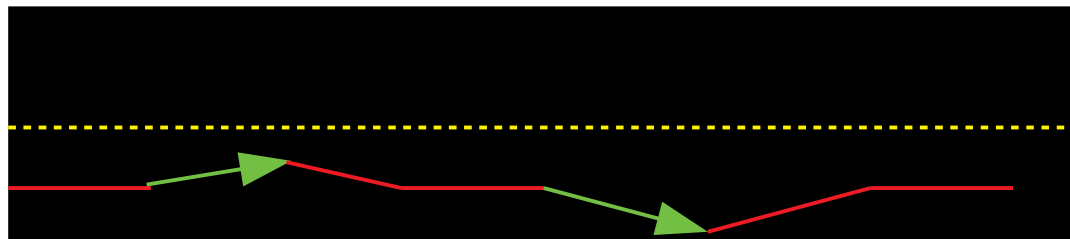
Imitation Learning: Problems

- Expert trajectories can be hard to gather
- Expert limits performance
- Distribution shift

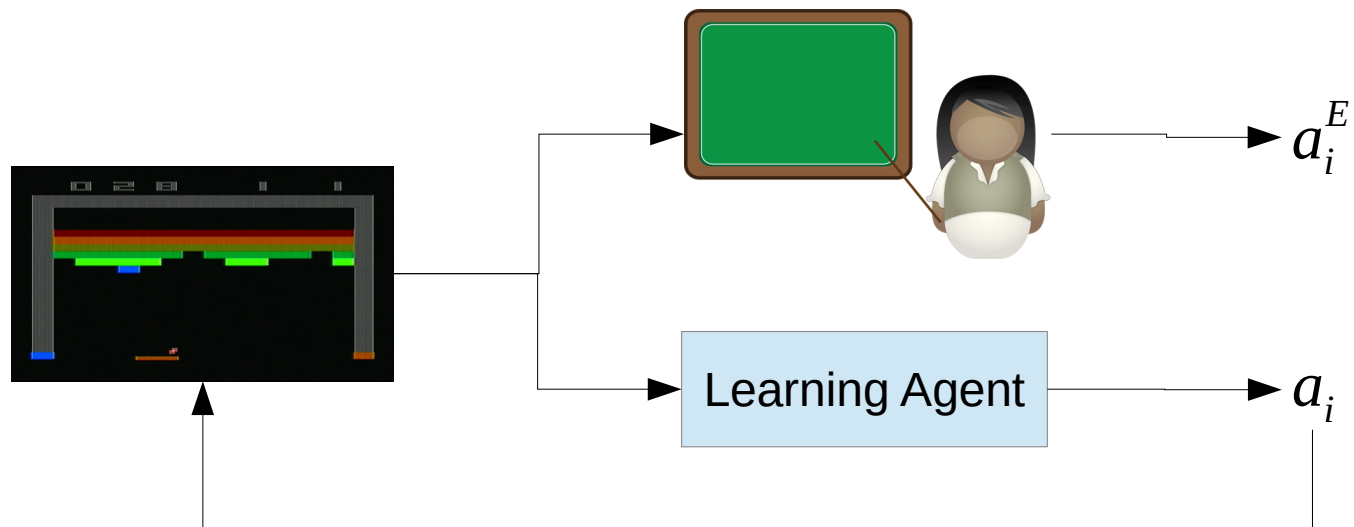


Imitation Learning Tips and Tricks

- Pre-training (for image inputs)
- Data Augmentation



Dagger



“On-policy”