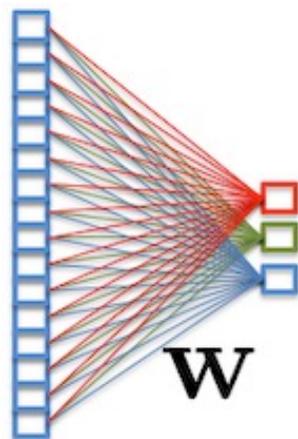
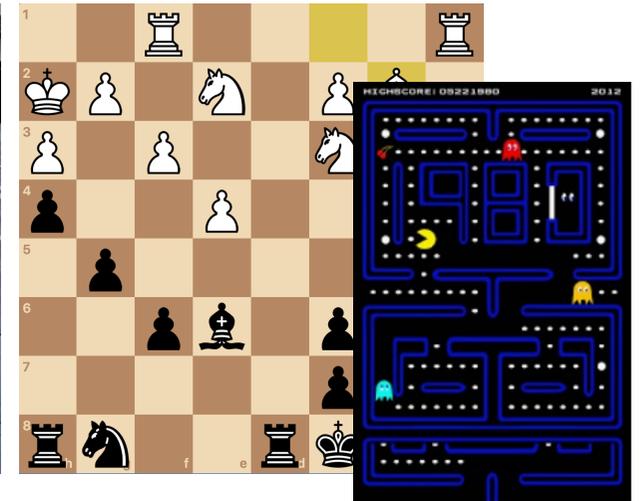
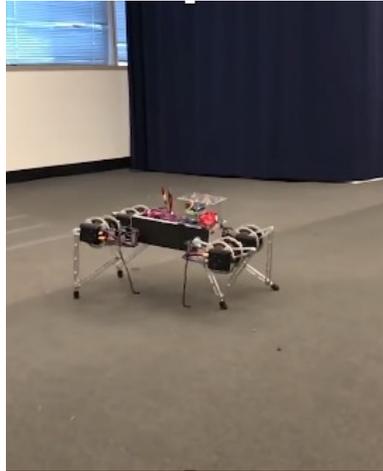


# AI and Data Analytics

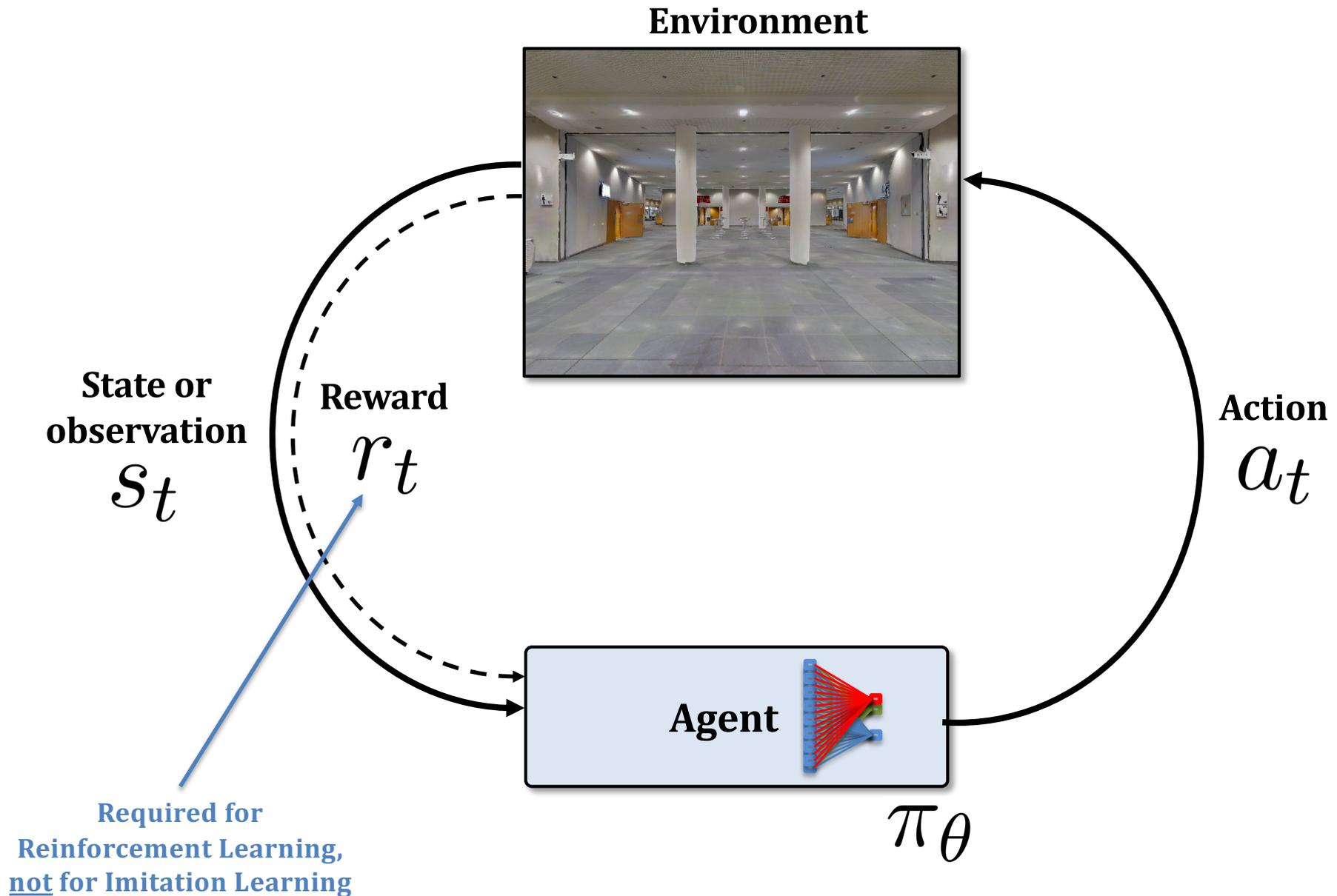
## 3.1 Imitation Learning



# Learning to control



# Interacting with an "environment"



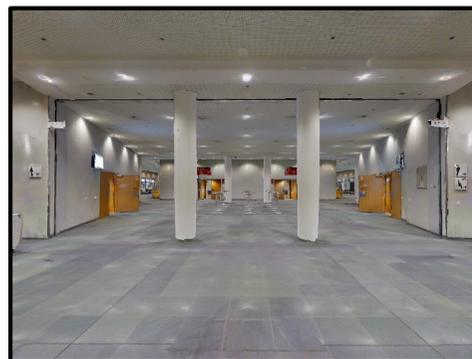
# States

An environment is in a given state  $s_t$  at any instant  $t$ .

**Fully observable problems** — the state is observable to the agent at any time, e.g. board games (chess, Go), 2D computer games with visible full screen (Pacman).



**Partially observable problems** — only an observation is returned to the agent at each instant, not the full state. Examples: robot navigation; 3D computer games with ego-centric view.



# Actions

At each time  $t$ , the agent takes an action  $a_t$  from an action space  $\mathcal{A}$ .

## Discrete action spaces

- Chess:  $\mathcal{A} = \{a1, Qf3, Ra1, bf3, \dots\}$ .
- Go:  $\mathcal{A} = [1 \dots 19] \times [1 \dots 19]$ .
- Dialogue:  $\mathcal{A} =$  words or characters.

## Continuous action spaces

- Robotic arm control:  $\mathcal{A} =$  motor commands
- Drone (UAV) control:  $\mathcal{A} =$  motor commands

# Policies

The policy  $\pi$  of an agent is a function which

- chooses an action given that the agent is in a current state;
- or, in the stochastic case, it assigns a probability to an action:

$$\pi(a_t | s_t)$$

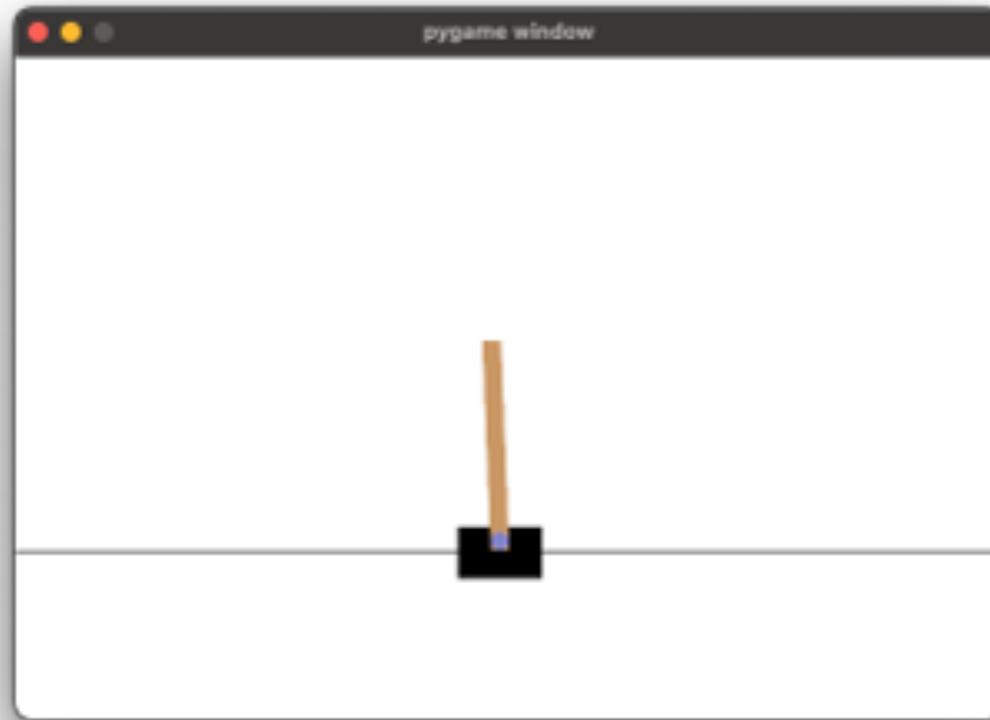
Policies are generally parametrized by a vector  $\theta$  and implemented as a neural network, denoted as  $\pi_\theta$ .

# The Gym-Cartpole environment

A cart with an inverted pendulum is controlled with two different actions:

- 0 → move to the left
- 1 → move to the right

The goal is to keep the pendulum/pole stable upright.

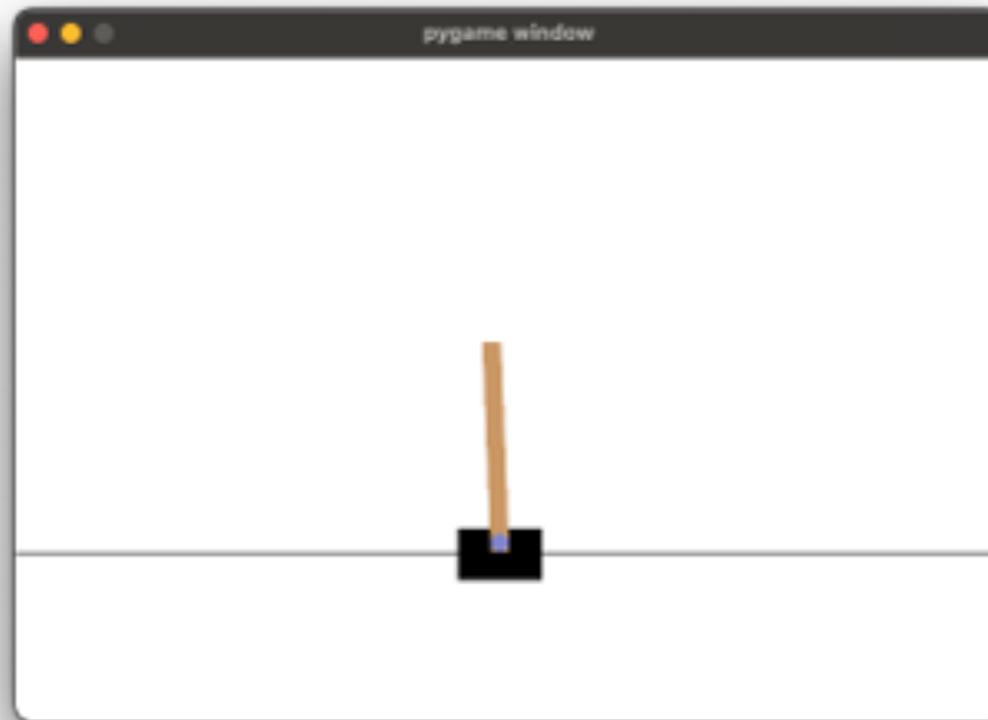


# The Gym-Cartpole environment

Termination criteria:

- Angle too steep: failure
- Cart too far off the center: failure
- Time limit reached: success

Lecture demo source code: `cartpole_test.py`



# The Gym-Cartpole environment

States space: 4 values

Action space: 2 values

$$f\left(\begin{bmatrix} x \\ \dot{x} \\ \theta \\ \dot{\theta} \end{bmatrix}\right) = \begin{bmatrix} \text{Score for Left} \\ \text{Score for Right} \end{bmatrix}$$

		min	max
$x$	Cart Position	-4.8	4.8
$\dot{x}$	Cart Velocity	$-\infty$	$\infty$
$\theta$	Pole Angle	$\approx -0.418 = -24^\circ$	$\approx 0.418 = 24^\circ$
$\dot{\theta}$	Pole Angular Velocity	$-\infty$	$\infty$

# Gym “Cartpol”: programming interface

```
1 # Initialize environment
2 env = gym.make("CartPole-v1", render_mode="human")
3 env.reset()
4
5 print("Left/Right arrows; 'ESC' to quit.")
6 running = True
7 while running:
8     action = None
9     keys = pygame.key.get_pressed()
10
11     if keys[pygame.K_LEFT]:
12         action = 0
13     elif keys[pygame.K_RIGHT]:
14         action = 1
15     if action is not None:
16         obs, reward, term, succ, info = env.step(action)
17
18 env.close()
19 pygame.quit()
```

# Adding success/failure checking

```
1  if action is not None:
2      obs, reward, term, succ, info = env.step(action)
3
4      if term:
5          print("\n" + "="*20 + "\n  EPISODE FAILED!\n" +
6              "="*20 + "\n" )
7
8          # Keep the window open for a second, then reset
9          time.sleep(1.0)
10         env.reset()
11
12         elif succ:
13             print("Time Limit Reached (Success!)")
14             env.reset()
15
16         # User wants to quit
17         for event in pygame.event.get():
18             if event.type == pygame.QUIT or (event.type ==
19                 pygame.KEYDOWN and event.key == pygame.K_ESCAPE):
20                 running = False
```

# Behavior cloning

“**Behavior cloning**” is a special case of “**Imitation Learning**”:

- Use a dataset  $\mathcal{D} = \{(s_i, a_i^*)\}$  of pairs of input states  $s_i$  and expert actions  $a_i^*$ .
- Use supervised learning: train a policy  $\pi$  to predict  $a_i^*$  when being in state  $s_i$ :

$$\min_{\theta} \sum_i \mathcal{L}(\pi_{\theta}(s_i), a_i^*)$$

- The loss function  $\mathcal{L}$  depends on the action space. For instance, for discrete actions (Gym-Cartpole), a cross-entropy loss could be used.

# Implementing the policy

```
1 class StudentNet(nn.Module):
2     def __init__(self, state_dim, action_dim):
3         super().__init__()
4         self.net = nn.Sequential(
5             nn.Linear(state_dim, 32),
6             nn.ReLU(),
7             nn.Linear(32, action_dim)
8         )
9     def forward(self, x):
10        return self.net(x)
```

# Behavior Cloning in PyTorch

```
1 def behavior_cloning(exp_states, exp_act, num_epochs):
2     States=4
3     Actions=2
4     model = StudentNet(States, Actions)
5     optimizer = optim.Adam(model.parameters(), lr=1e-3)
6     criterion = nn.CrossEntropyLoss()
7
8     # For simplicity, we do not batch:
9     # do one forward pass with the full data
10    for epoch in range(num_epochs):
11        pred_act = model(exp_states)
12        loss = criterion(pred_act, exp_act)
13
14        optimizer.zero_grad()
15        loss.backward()
16        optimizer.step()
17
18    return model
```

# Demo



```
pytorch_control -- -bash -- 130x40
(coursia) (base) cwolf pytorch_control $ █
```

The image shows a terminal window with a dark background. At the top, the window title is "pytorch\_control -- -bash -- 130x40". The prompt is "(coursia) (base) cwolf pytorch\_control \$ █". The rest of the terminal is mostly blank with some faint, illegible text visible in the background.

Lecture demo full code: [beh\\_cloning.py](#)