# CS394R
# Reinforcement Learning: Theory and Practice
# Fall 2007

**Peter Stone**

Department of Computer Sciences
The University of Texas at Austin

# Good Afternoon Colleagues

- Are there any questions?

# Logistics

- How are the final projects coming?

# Helicopter Control

- State: position, orientation, velocity, angular vels
- Actions: Settings of the 4 or 5 controls
- Goal: Hover

# Helicopter Control

- State: position, orientation, velocity, angular vels
- Actions: Settings of the 4 or 5 controls
- Goal: Hover
- How would you formulate the problem "by the book"?

# Helicopter Control

- State: position, orientation, velocity, angular vels
- Actions: Settings of the 4 or 5 controls
- Goal: Hover
- How would you formulate the problem "by the book"?
- Could you implement that? Why or why not?

# Helicopter Control

- State: position, orientation, velocity, angular vels
- Actions: Settings of the 4 or 5 controls
- Goal: Hover
- How would you formulate the problem "by the book"?
- Could you implement that? Why or why not?
- At a high level, what do they do instead?

# Helicopter Control

- State: position, orientation, velocity, angular vels
- Actions: Settings of the 4 or 5 controls
- Goal: Hover
- How would you formulate the problem "by the book"?
- Could you implement that? Why or why not?
- At a high level, what do they do instead?
  - Collect a small amount of human expert data
  - Use that to train a **1-step** model (simulator)
  - Determine the optimal policy in the simulator
  - Fly it!

# Ng paper

- Why hover upside down?

# Ng paper

- Why hover upside down?
- Why quadratic reward (p. 6)?

# Ng paper

- Why hover upside down?
- Why quadratic reward (p. 6)?
- PEGASUS - how does it help policy evaluation?

# Ng paper

- Why hover upside down?
- Why quadratic reward (p. 6)?
- PEGASUS - how does it help policy evaluation?
  - General question: is policy good or lucky?

# Ng paper

- Why hover upside down?
- Why quadratic reward (p. 6)?
- PEGASUS - how does it help policy evaluation?
  - General question: is policy good or lucky?
  - Use same random samples for evaluation of each policy

# Ng paper

- Why hover upside down?
- Why quadratic reward (p. 6)?
- PEGASUS - how does it help policy evaluation?
    - General question: is policy good or lucky?
    - Use same random samples for evaluation of each policy
- How does he do policy optimization?

# Ng paper

- Why hover upside down?
- Why quadratic reward (p. 6)?
- PEGASUS - how does it help policy evaluation?
  - General question: is policy good or lucky?
  - Use same random samples for evaluation of each policy
- How does he do policy optimization?
  - greedy hillclimbing over few parameters (the NNs)!
- Could the approach be used to invert the helicopter? Or is it easier just to hover?
- Can it generalize to adverse conditions?
- Where's the power?  Is it an easy problem or a powerful approach?

# Robot Soccer paper

- Why I selected it. . .