

*CSCE 580: Introduction to AI*  
*CSCE 581: Trusted AI*

## Lecture 25: Planning and Reinforcement Learning

---

PROF. BIPLAV SRIVASTAVA, AI INSTITUTE

21<sup>ST</sup> NOV, 2023

**Carolinian Creed: “I will practice personal and academic integrity.”**

**Credits: Copyrights of all material reused acknowledged**

# Organization of Lecture 25

---

- Introduction Segment
  - Recap of Lectures 23 and 24
- Main Segment
  - Making Sequential Decisions
  - Planning
  - Reinforcement Learning
- Concluding Segment
  - Course Project Discussion
  - About Next Lecture – Lecture 26
  - Ask me anything

# Introduction Section

---

# Recap of Lecture 23 and 24

---

- Topic discussed
  - Making Decisions
  - Simple Decisions
  - Complex Decisions
  - Quiz 4

# Graduate Paper Presentation

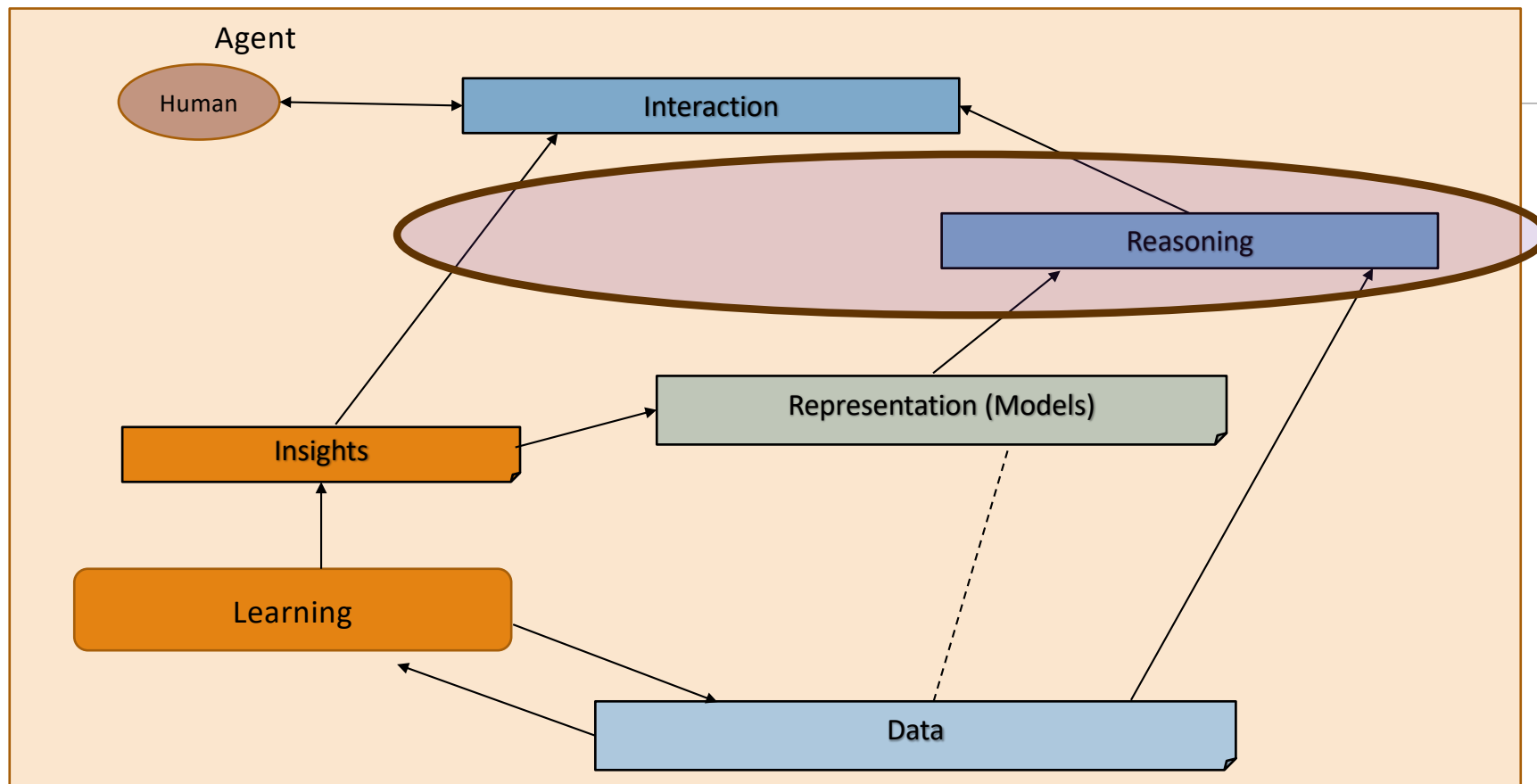
---

- Papers between 2021-2023 (last 3 years)
- At top AI venues: AAAI, Neurips, IJCAI, ICML, ICLR, or discuss with instructor
- Guideline on presentation
  - See template format shared in Google drive
- More in the concluding section of lecture

# Intelligent Agent Model



# Relationship Between Main AI Topics



# Where We Are in the Course

## CSCE 580/ 581 – In This Course

- Week 1: Introduction, Aim: Chatbot / Intelligence Agent
- Weeks 2-3: Data: Formats, Representation and the Trust Problem
- Week 4-5: Search, Heuristics - Decision Making
- Week 6: Constraints, Optimization – Decision Making
- Week 7: Classical Machine Learning – Decision Making, Explanation
- Week 8: Machine Learning - Classification
- Week 9: Machine Learning - Classification – Trust Issues and Mitigation Methods
- Topic 10: Learning neural network, deep learning, Adversarial attacks
- Week 11: Large Language Models – Representation, Issues
- Topic 12: Markov Decision Processes, Hidden Markov models - Decision making
- Topic 13: Planning, Reinforcement Learning – Sequential decision making
- Week 14: AI for Real World: Tools, Emerging Standards and Laws; Safe AI/ Chatbots



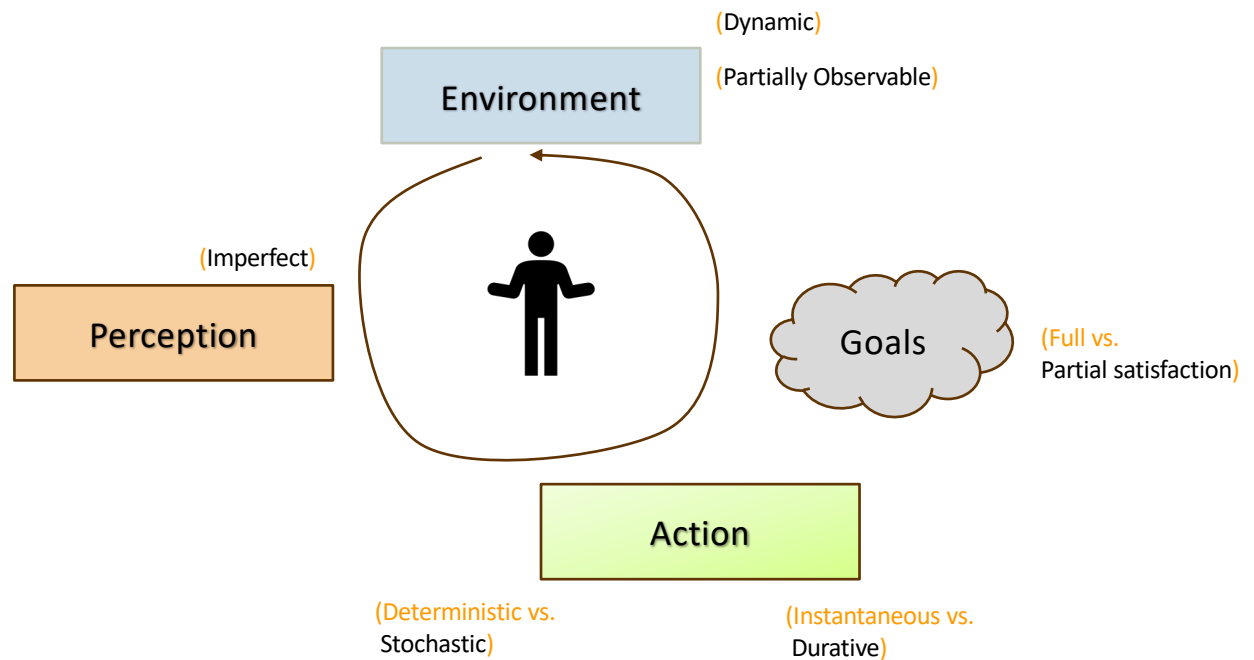
# Main Section

---

**Credit:** Retrieved from internet

# Complex Decisions

- Making a sequence of decisions
- Making a single decision but with
  - Environment changing
  - Actions not being deterministic
  - Perception not being perfect
  - ...



# Goal-Based Agents

## Generating Sequence of Actions

---

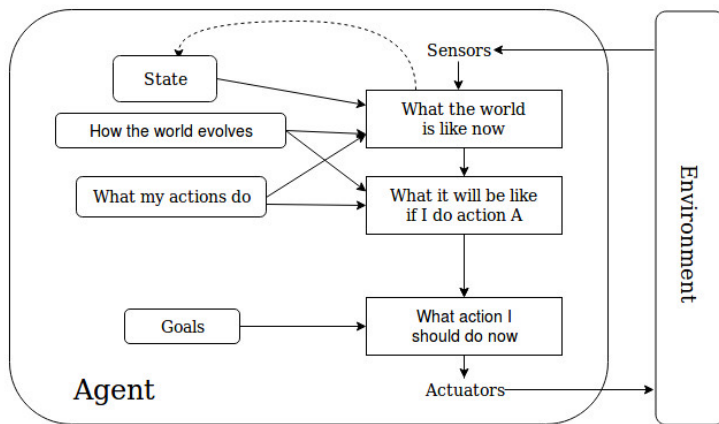
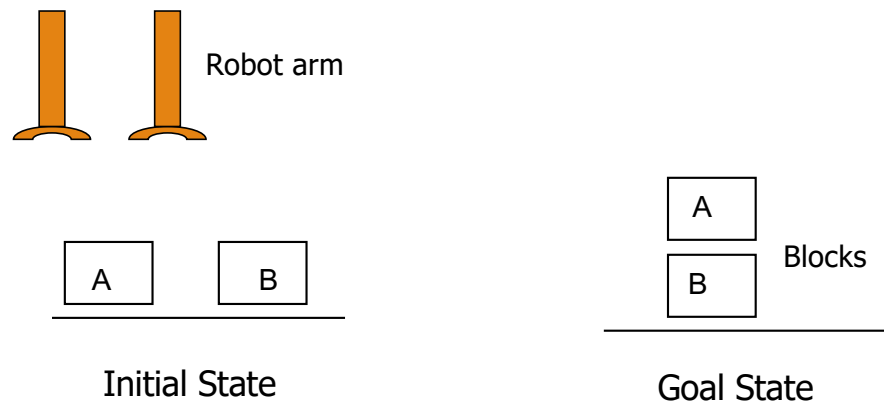


Figure Source: Russell & Norvig, AI: A Modern Approach

# Reasoning Illustration - Planning Example

---

## *Blocks World*

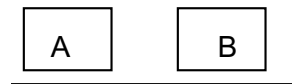


All robots are equivalent

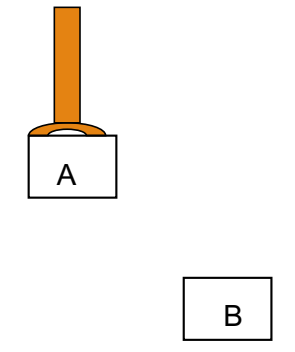
# Reasoning Illustration - Representation

---

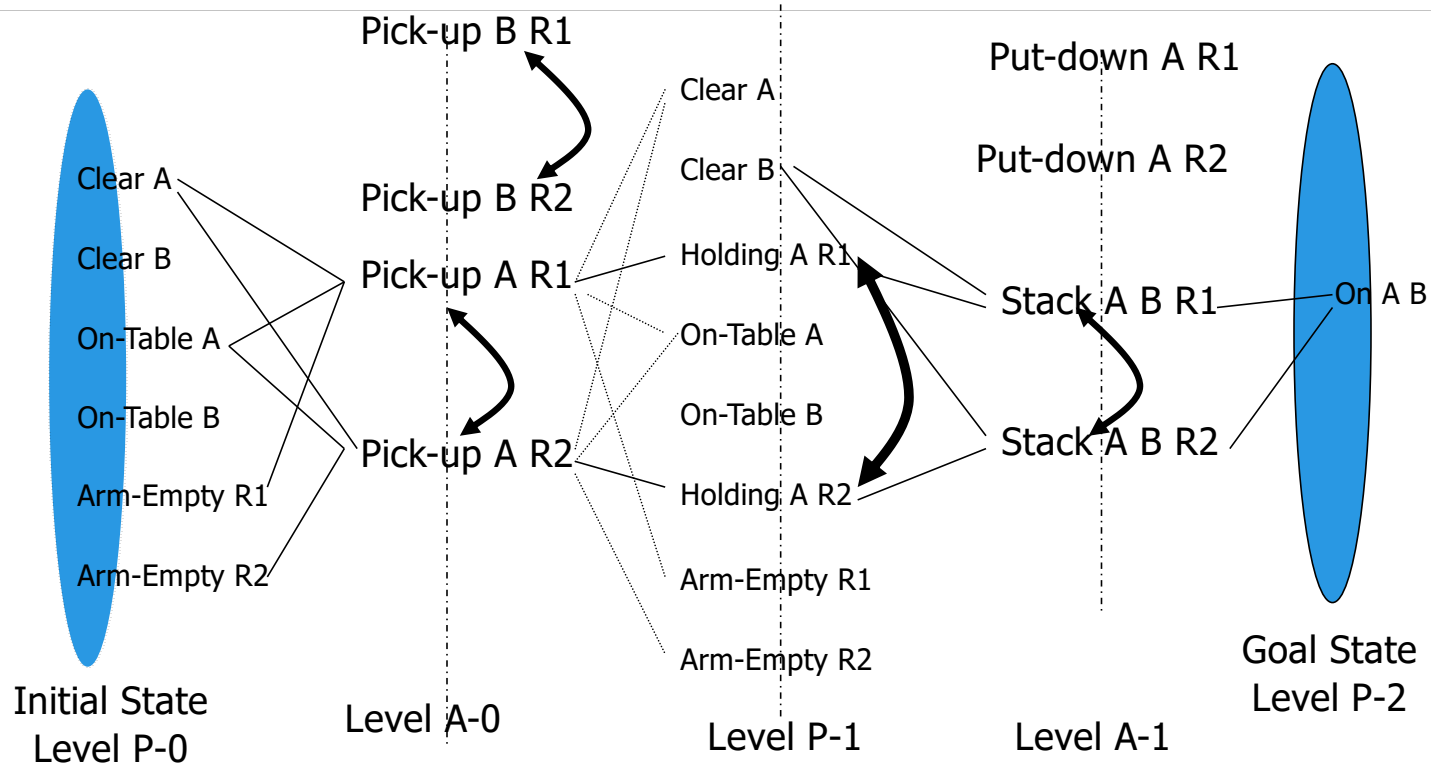
States: ((On-Table A) (On-Table B) ...)



Actions: ((Name: (Pickup ?block ?robot)  
Precondition: ((Clear ?block)  
(Arm-Empty ?robot)  
(On-Table ?block))  
Add: ((Holding ?block ?robot))  
Delete: ((Clear ?block)  
(Arm-Empty ?robot)))...)



# Reasoning Illustration - Planning Process



# Active Area of Research

---

## Considerations

- What to find:
  - Any workable plan
  - Optimal plan – but then what is the criteria
  - All plans
  - Diverse plans
- How to find
  - Plan at the end
  - Plan anytime
- How to represent problem
- How to explain solution

# Hand's On With Planning

---

- Site: <http://planning.domains/>
  - Try the editor: <https://editor.planning.domains/#>
- Code example with API: <https://github.com/biplav-s/course-ai-tai-f23/blob/main/sample-code/Class25-Planning/PlannerInvokerWithAPIs.ipynb>



# Exercise: 10 mins

---

- Try any domain from domain.pddl or classical planning repo:  
<https://github.com/AI-Planning/classical-domains/tree/main/classical>
- Change sample code with domain and problem files
- Run the sample code

# Forms of Uncertainty and Planning

---

- Uncertain knowledge, caused by
  - Incomplete knowledge
  - Incorrect knowledge
- Uncertain actions, caused by
  - Physics of the domain
  - External events

# Forms of Uncertainty

- Uncertain knowledge, caused by
  - Incomplete knowledge
  - Incorrect knowledge
- Uncertain actions, caused by
  - Physics of the domain
  - External events

Alternative approaches to represent

- Degree of belief: Probability. The sentence still is true or false
- Degree of truth: Fuzzy logic

Language	Ontological Commitment (What exists in the world)	Epistemological Commitment (What an agent believes about facts)
Propositional logic	facts	true/false/unknown
First-order logic	facts, objects, relations	true/false/unknown
Temporal logic	facts, objects, relations, times	true/false/unknown
Probability theory	facts	degree of belief 0...1
Fuzzy logic	degree of truth	degree of belief 0...1

## Credits:

- Russell & Norvig, AI - A Modern Approach
- Deepak Khemani - A First Course in AI

# Forms of Uncertainty

---

- Uncertain knowledge, caused by
  - Incomplete knowledge
  - Incorrect knowledge
- Uncertain actions, caused by
  - Physics of the domain
  - External events



Use Probability Theory  
Infer using probabilities

Decision Processes = create  
situational policies (state-action based)

# Decision-theoretic Agent

Probability theory: degree of belief in sentences

- Summarizes the uncertainty  $t$

Utility theory: represent and reason with preferences

**function** DT-AGENT(*percept*)**returns** an *action*

**static:** a set probabilistic beliefs about the state of the world

calculate updated probabilities for current state based on  
available evidence including current percept and previous action

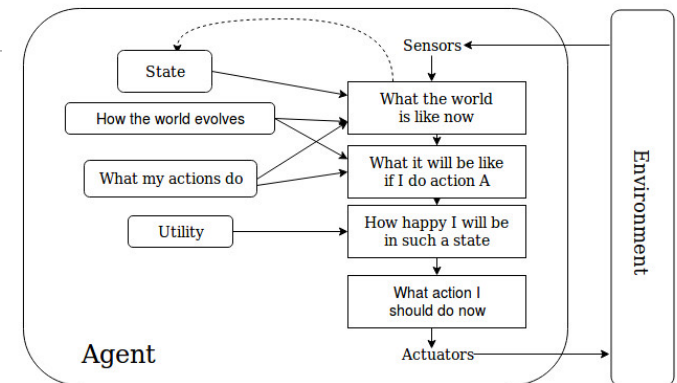
calculate outcome probabilities for actions,

given action descriptions and probabilities of current states

select *action* with highest expected utility

given probabilities of outcomes and utility information

**return** *action*



Source: Russell & Norvig, AI - A Modern Approach

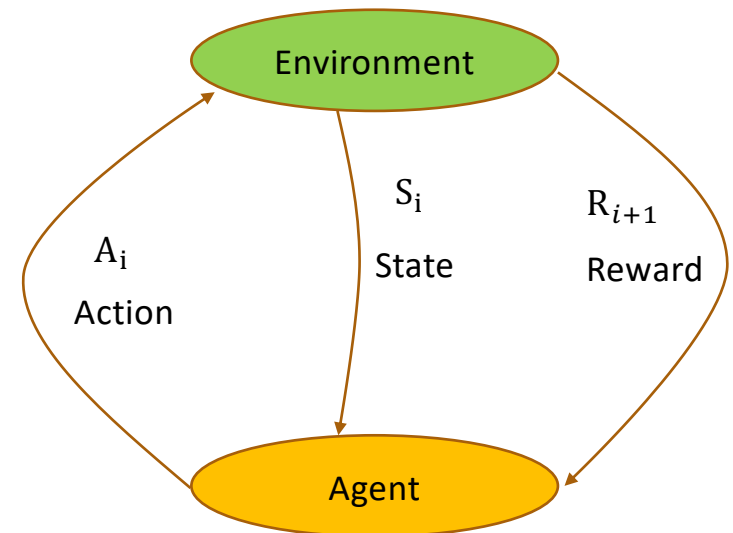
# Reinforcement Learning

---



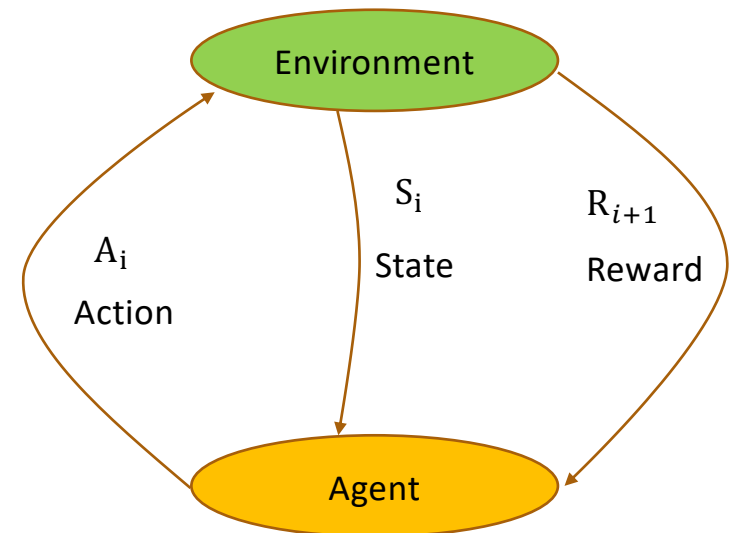
# Reinforcement Learning Setting

- An agent in an environment
- Agent
  - Can see **state**
  - Can take **action**
  - Will get **rewards**
- Precisely, at each time step  $i$ 
  - In state  $S_i$ , agent takes action  $A_i$
  - Based on state  $s_i$  and action  $a_i$ , the environment transitions to state  $S_{i+1}$  and outputs reward  $R_{i+1}$
- **Objective:** learn mapping of **states** to **actions** so that the agent maximizes the **reward** from the **environment**.



# Reinforcement Learning

- **Objective:** learn mapping of **states** to **actions** so that the agent maximizes the **reward** from the **environment**.
- **Output**
  - Deterministic:  $a = \pi(s)$
  - Stochastic:  $\pi(a|s) = P(A_i = a|S_i = s)$





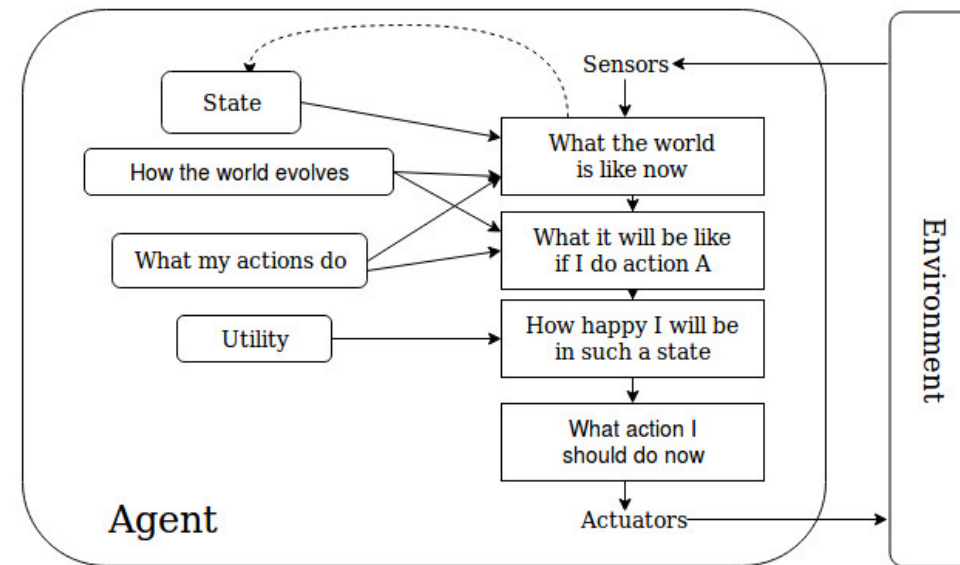
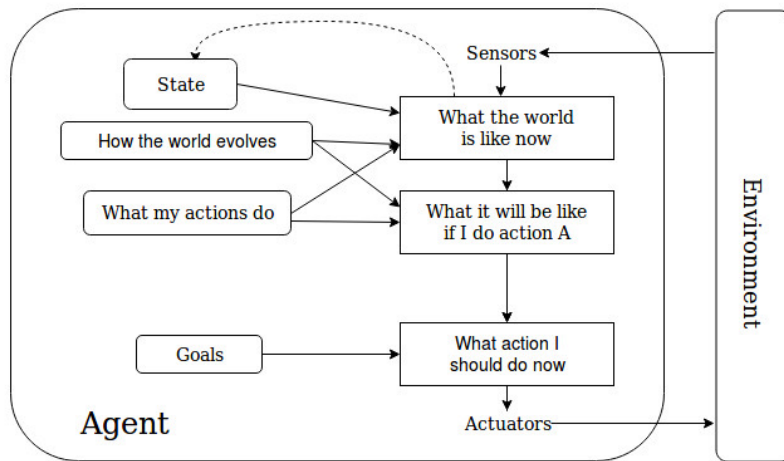
# Comparison With Other Learning

---

- Supervised learning
  - Training information: labels
  - Objective: learn (input-label) mapping
  - Goodness criteria: Reduce error = (Predicted label – Actual label)
- Reinforcement learning
  - Training information: reward functions
  - Objective: learn policy
  - Goodness criteria: maximal reward
- These two forms of learning are orthogonal – for different tasks

# RL as a Learning-Based Agent

A general, alternative way of solving goal-based problems from just execution traces

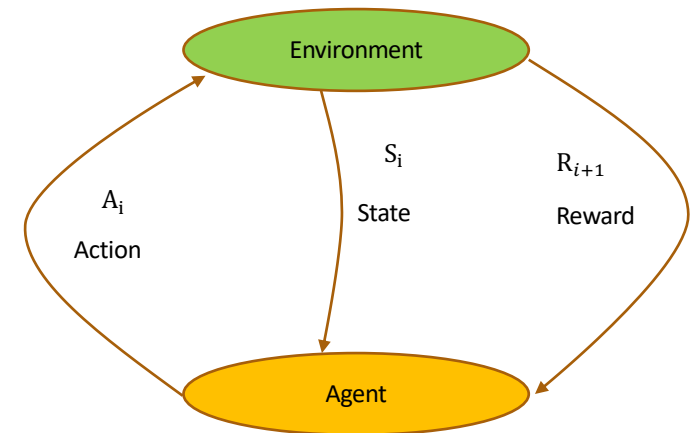
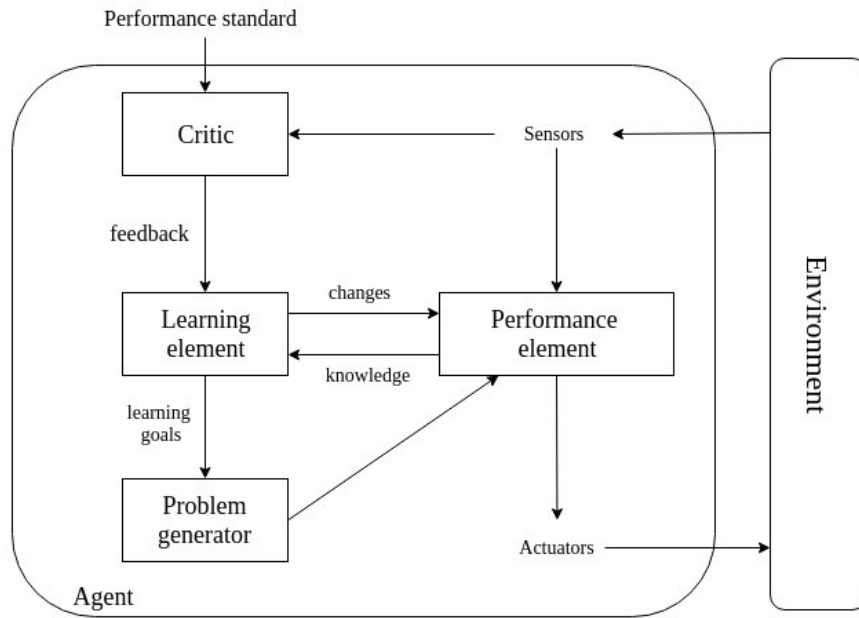


Goal- and Utility-based Intelligent Agent



# RL as a Learning-Based Agent

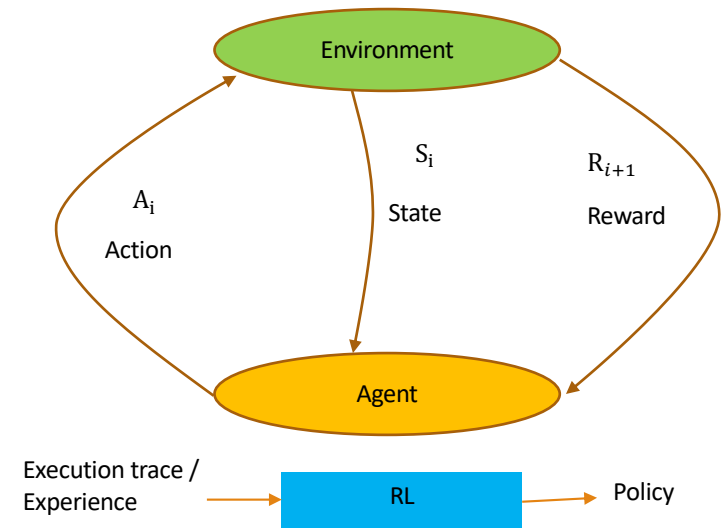
A general, alternative way of solving goal-based problems from just execution traces



# RL as a Learning-Based Agent

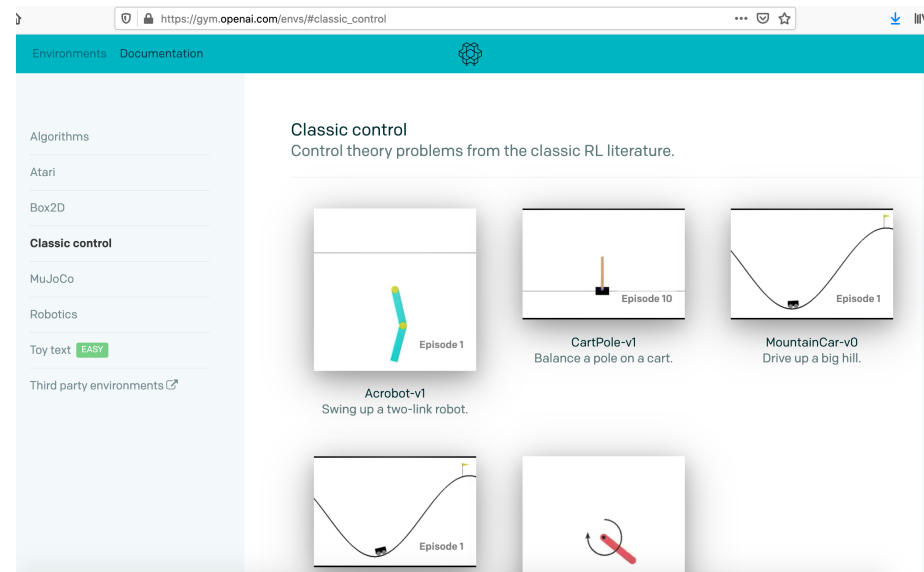
A general, alternative way of solving goal-based problems from just execution traces

Goal- and Utility-based Intelligent Agent



# Exercise and Code – Gym RL

- RL using Open AI's Gym
  - <https://gym.openai.com/>
  - Environments: [https://gym.openai.com/envs/#classic\\_control](https://gym.openai.com/envs/#classic_control)
- Exercise (5 mins):
  - Look at the various categories
  - Explore the videos



# Exercise and Code – Gym RL

---

- RL using Open AI's Gym
  - <https://gymnasium.farama.org/>
  - Old: <https://gym.openai.com/>
- Code: <https://github.com/biplav-s/course-d2d-ai/blob/main/sample-code/l18-learning-agent/RL%20using%20Gym.ipynb>

Source: Russell & Norvig, AI: A Modern Approach

# Diversity in RL Problems

---

- Environment - accessible or inaccessible
  - Accessible: states can be identified with percepts
  - Inaccessible environment: agent has to learn and maintain representation of state to track environment
- Knowledge of effects of action and utility, or learn
- Rewards
  - Available for all states or only terminal states
  - Actual utility or hints of increase/ decrease
- Ability to execute actions - Active learner or passive learner
  - A passive learner simply watches the world going by, and tries to learn the utility of being in various states
  - An active learner can actions to explore unknown environment

Source: Russell & Norvig, AI - A Modern Approach

# Passive RL

---

- **Input**

- policy:  $\pi_i$
- // Has no knowledge Reward  $R(s)$  and Transition function  $P(s' | s, a)$

- **Output**

- Expected utility for each state,  $U(s)$

- **Procedure:**

- Execute a sequence of runs
- At any instant, the agent knows only its current state and current reward, and the action it must take next. This action may lead it to more than one state, with different probabilities.

- **Expected Utility**

$$U^\pi(s) = E(\sum_{t=0}^{\infty} \gamma^t R^t(s'))$$

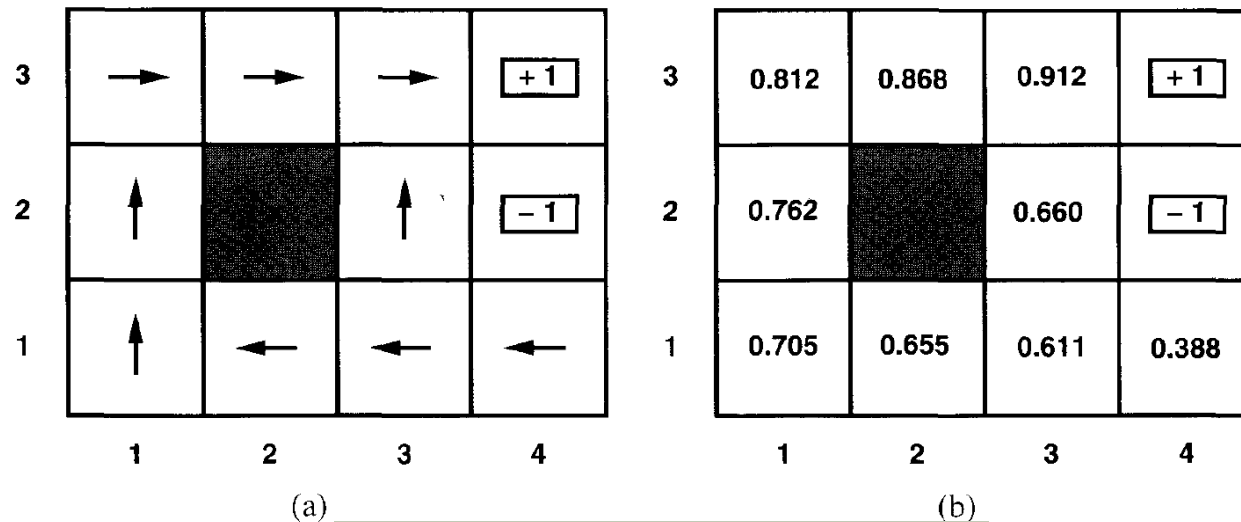


# Illustration

```
# Action Directions
north = (0, 1)
south = (0, -1)
west = (-1, 0)
east = (1, 0)

policy = {
    (0, 2): east, (1, 2): east, (2, 2): east, (3, 2): None,
    (0, 1): north, (2, 1): north, (3, 1): None,
    (0, 0): north, (1, 0): west, (2, 0): west, (3, 0): west,
}
```

Policy: [https://github.com/biplav-s/course-d2d-ai/blob/main/sample-code/l15-l16-l17-l18-agents/reinforcement\\_learning.ipynb](https://github.com/biplav-s/course-d2d-ai/blob/main/sample-code/l15-l16-l17-l18-agents/reinforcement_learning.ipynb)



Input Policy and Output Optimal Utility

Source: Russell & Norvig, AI - A Modern Approach

# The Markov Property – True of Many Domains

---

- **Our policy at timepoint  $t$  is only dependent on the current state  $s$** 
  - $\pi(a|s) = P(A_t = a|S_t = s)$
- Although the agent has a history up until  $S_t$ 
  - $H_t = S_0, A_0, R_1 S_1, A_1, R_2 \dots S_{t-1}, A_{t-1}, R_t, S_t$
- One may assume that all relevant information about the future is contained in the current state and action
  - $P(S_{t+1} = s', R_{t+1} = r|S_t = s, A_t = a) = P(S_{t+1} = s', R_{t+1} = r|H_t = h_{t+1}, A_t = a)$
- This is a generalization of the Markov property to sequential decision problems
  - $P(S_{t+1}|S_t) = P(S_{t+1}|S_t, S_{t-1}, \dots S_0)$

# RL with Finite States

---

## *Solving a Finite MDP*

- **States:** A discrete and finite set  $\mathcal{S}$
- **Actions:** A discrete and finite set  $\mathcal{A}$
- **Transition Probabilities:**  $P(S_{t+1} = s', R_{t+1} = r | S_t = s, A_t = a)$ 
  - Defines the dynamics of the MDP
- The state-transition probabilities can be obtained from the transition probabilities
  - $p(s' | s, a) = \sum_{r \in \mathcal{R}} p(s', r | s, a)$  // Estimating state-transition by looking at reward of samples
- The **expected reward** can be obtained from the transition probabilities
  - $r(s, a) = \sum_{r \in \mathcal{R}} r \sum_{s' \in \mathcal{S}} p(s', r | s, a) = \mathbb{E}[R_{t+1} | S_t = s, A_t = a]$  // Estimating reward from transitions seen

Adapted from: Forest A.'s RL Course

# Model-free RL: Q-learning

---

- Learning action-value functions
- $Q(a,i)$ : value of doing action  $a$  in state  $i$
- Relationship between utility  $U$  of state and  $Q$  value
  - $U(i) = \max_a Q(a, i)$
- Finding  $Q$  value based on whether transition probability is known
  - When  $M$  (transition is known)

$$Q(a, i) = R(i) + \sum_j M_{ij}^a \max_{a'} Q(a', j)$$

- Estimating with TD method

$$Q(a, i) \leftarrow Q(a, i) + \alpha (R(i) + \max_{a'} Q(a', j) - Q(a, i))$$

Source: Russell & Norvig, AI - A Modern Approach

# RL with Deep Learning

---

- For small problems, like games, state-value function (U), action- utility value (Q), and transition functions (M), and policy functions are represented using a table
- But for large and realistic problems, number of states are countably large/ practically infinite
- Deep learning are excellent function approximators
  - Estimate Q-value i.e., action-value
- Not covered in this class

# Exercise and Code – RL

---

- RL settings and solution methods
- Code: <https://github.com/biplav-s/course-d2d-ai/blob/main/sample-code/l18-learning-agent/RL%20using%20Gym.ipynb>

Source: Russell & Norvig, AI: A Modern Approach

# Inverse Reinforcement Learning

---

- Given  $\pi^*$  and transition function  $M$ ,
  - can we recover  $R$
- Or, given execution traces corresponding to  $\pi^*$ 
  - can we recover  $R$ ?
- Applications
  - Path planning
  - Automated-driving
- Reference: Pieter Abbel's course slides: <https://people.eecs.berkeley.edu/~pabbeel/cs287-fa12/slides/inverseRL.pdf>

# RL References

---

- Sutton and Barto's Book: <http://incompleteideas.net/book/the-book.html>
- Russell and Norvig, AI – A modern Approach
- David Silver's RL course, <https://www.davidsilver.uk/teaching/>
- Inverse RL
  - A Survey of Inverse Reinforcement Learning: Challenges, Methods and Progress, <https://arxiv.org/abs/1806.06877>, 2018
  - Pieter Abbel's course slides: <https://people.eecs.berkeley.edu/~pabbeel/cs287-fa12/slides/inverseRL.pdf>



# Course Project

---

# Project Discussion: What Problem Fascinates You ?

---

- Data
  - Water
  - Finance
  - ...
- Analytics
  - Search, Optimization, Learning, Planning, ...
- Application
  - Building chatbot
- Users
  - Diverse demographics
  - Diverse abilities
  - Multiple human languages

## Project execution in sprints

- Sprint 1: (Sep 12 – Oct 5)
  - **Solving**: Choose a decision problem, identify data, work on solution methods
  - **Human interaction**: Develop a basic chatbot (no AI), no problem focus
- Sprint 2: (Oct 10 – Nov 9)
  - **Solving**: Evaluate your solution on problem
  - **Human interaction**: Integrated your choice of chatbot (rule-based or learning-based) and methods
- Sprint 3: (Nov 14 – 30)
  - **Evaluation**: Comparison of your solver chatbot with an LLM-based alternative, like ChatGPT

# Project Discussion: Dates and Deliverables

---

## Project execution in sprints

- Sprint 1: (Sep 12 – Oct 5)
  - **Solving**: Choose a decision problem, identify data, work on solution methods
  - **Human interaction**: Develop a basic chatbot (no AI), no problem focus
- Sprint 2: (Oct 10 – Nov 9)
  - **Solving**: Evaluate your solution on problem
  - **Human interaction**: Integrated your choice of chatbot (rule-based or learning-based) and methods
- Sprint 3: (Nov 14 – 30)
  - **Evaluation**: Comparison of your solver chatbot with an LLM-based alternative, like ChatGPT

- Oct 12, 2023
  - Project checkpoint
  - In-class presentation
- Nov 30, 2023
  - Project report due
- Dec 5 / 7, 2023
  - In-class presentation

# Skeleton: A Basic Chatbot

- Run in an infinite loop until the user wants to quit
- Handle any user response
  - User can quit by typing “Quit” or “quit” or just “q”
  - User can enter any other text and the program has to handle it. The program should write back what the user entered and say – “I do not know this information”.
- Handle known user query types // Depends on your project
  - “Tell me about N-queens”, “What is N ?”
  - “Solve for N=4?”
  - “Why is this a solution? ”
- Handle chitchat // Support at least 5, extensible from a file
  - “Hi” => “Hello”
  - ...
- *Store session details in a file*

## Illustrative Project

1. **Title:** Solve and explain solving of n-queens puzzle
2. **Key idea:** Show students how a course project will look like
3. **Who will care when done:** students of the course, prospective AI students and teachers
4. **Data need:** n: the size of game; interaction
5. **Methods:** search
6. **Evaluation:** correctness of solution, quality of explanation, appropriateness of chat
7. **Users:** with and without AI background; with and without chess background
8. **Trust issue:** user may not believe in the solution, may find interaction offensive (why queens, not kings? ...)

# Project Discussion: Illustration

1. Create a private Github repository called “CSCE58x-Fall2023-<studentname>-Repo”. Share with Instructor (biplav-s) and TA (kausik-l)
2. Create Google folder called “CSCE58x-Fall2023-<studentname>-SharedInfo”. Share with Instructor ([prof.biplav@gmail.com](mailto:prof.biplav@gmail.com)) and TA ([lakkarajukausik90@gmail.com](mailto:lakkarajukausik90@gmail.com))
3. Create a Google doc in your Google repo called “Project Plan” and have the following by next class (Sep 5, 2023)

1. **Title:** Solve and explain solving of n-queens puzzle
2. **Key idea:** Show students how a course project will look like
3. **Who will care when done:** students of the course, prospective AI students and teachers
4. **Data need:** n: the size of game; interaction
5. **Methods:** search
6. **Evaluation:** correctness of solution, quality of explanation, appropriateness of chat
7. **Users:** with and without AI background; with and without chess background
8. **Trust issue:** user may not believe in the solution, may find interaction offensive (why queens, not kings? ...)

# Project Illustration: N-Queens

---

- Sprint 1: (Sep 12 – Oct 5)
  - **Solving**: Choose a decision problem, identify data, work on solution methods
    - Method 1: Random solution
    - Method 2: Search – BFS
    - Method 3: Search - ...
  - **Human interaction**: Develop a basic chatbot (no AI) as outlined
  - Deliverable
    - Code structure in Github
      - ./data
      - ./code
      - ./docs
      - ./test
    - Presentation: Make sprint presentation on Oct 12, 2023

# Reference: Project Rubric - NEW

- **Project report – 60%**
  - Project description: problem, related work, approach, evaluation – 40%
  - Working system demo/ video – 10%
    - Well organized Github with code (./data, ./code, ./docs, ./test) – 10%
- **Project presentation – 40%**
  - Evaluation by peers, instructor and TA
- **Bonus**
  - Instructor discretion – 10%
- **Penalty**
  - Lack of timeliness as per announced policy (right) - up to 30%

## Milestones and Penalties

- Oct 12, 2023
  - Project checkpoint
  - In-class presentation
  - **Penalty: presentation not ready by Oct 10, 2023 [-10%]**
- Nov 30, 2023
  - Project report due
  - **Project report not ready by date [-10%]**
- Dec 5 / 7, 2023
  - In-class presentation
  - **Project presentations not ready by Dec 4, 2023 [-10%]**

# Evaluation of Presentation

---

1. An online form will be available during presentation
2. During a presentation, three students will be assigned to review along with instructor and TA
3. They will enter following survey questions:
  1. Their name
  2. Presentation number
  3. How useful is the system – will you use it? [1-5 scale]
  4. How well have you understood the project from the presentation? [1-5 scale]
4. Top and bottom scores will be removed. Average of remaining three will be used for final presentation marks



# Lecture 5: Summary

---

- We talked about
  - Planning
  - Uncertainty
  - Reinforcement Learning

# Concluding Section

---

# Quiz 4

---

- November 14-21, 2023
  - Due today

# About Next Lecture – Lecture 26

---

# Student Assessment

A = [900-1000]  
B+ = [870-899]  
B = [800-869]  
C+ = [770-799]  
C = [700-769]  
D+ = [670-699]  
D = [600-669]  
F = [0-599]

Tests	Undergrad	Grad
Course Project – report, in-class presentation	600	600
Quiz – best of 3 from 4	200	200
Final Exam	200	100
Additional Final Exam – Paper summary, in-class presentation		100
Total	1000 points	1000 points

# Final Exam

---

- Graduate students
  - Paper presentation [100 points]
  - Write about their paper presented [100 points]
- Undergraduate students
  - Write about 2 papers presented in class by graduate students [150 points]
  - Vote for the papers presented [50 points]
- Paper reports due by Dec 5, 2023 (Tuesday)

Final Exam	200	100
Additional Final Exam – Paper summary, in-class presentation		100

# Lecture 26: Graduate Student Presentations

- 5 presentations
  - Sample template in drive (folder shared via Piazza); make a copy and edit
- Evaluation
  - By undergrads as well as instructor and TA
  - All undergraduates to attend and give survey response; link to be shared
  - Those undergrads not giving inputs will be given negative marks as part of the final score [-10 point per presentation]
- What to have in the report – minimum 1 page per paper (<500 words).
  - Paper summary
  - Key contributions
  - Your critique about the paper.

Nov 21 (Tu)	Sequential Decision Making: Planning, RL	Quiz 4- end [Week 14]
Nov 23 (Th)		Holiday - Thanksgiving
Nov 28 (Tu)	Paper presentation (grad students only)	
Nov 30 (Th)	AI for the Real World – Bringing All Together	Project – Sprint 3 - end
Dec 5 (Tu)	Project presentation	
Dec 7 (Th)	Project presentation	Last day of class
Dec 9 (Sat)		Reading Day
Dec 12 (Tu)	4pm – Final Overview	Optional, information shared