

"A beginning is the time for taking the most delicate care that the balances are correct."

Frank Herbert, *Dune*

CMPUT 365

Introduction to RL

Plan

- Introduction
- Course logistics
 - Instruction team
 - Pre-requisites
 - Textbook
 - Coursera
 - Academic integrity
 - Evaluation
- What is reinforcement learning?

Note

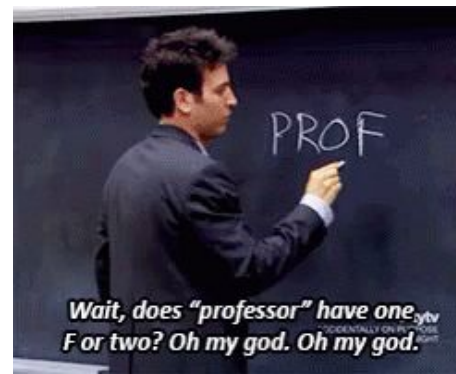
Lectures may be audio recorded for the purpose of a student's individual study as part of an approved academic accommodation.

Please, interrupt me at any time!

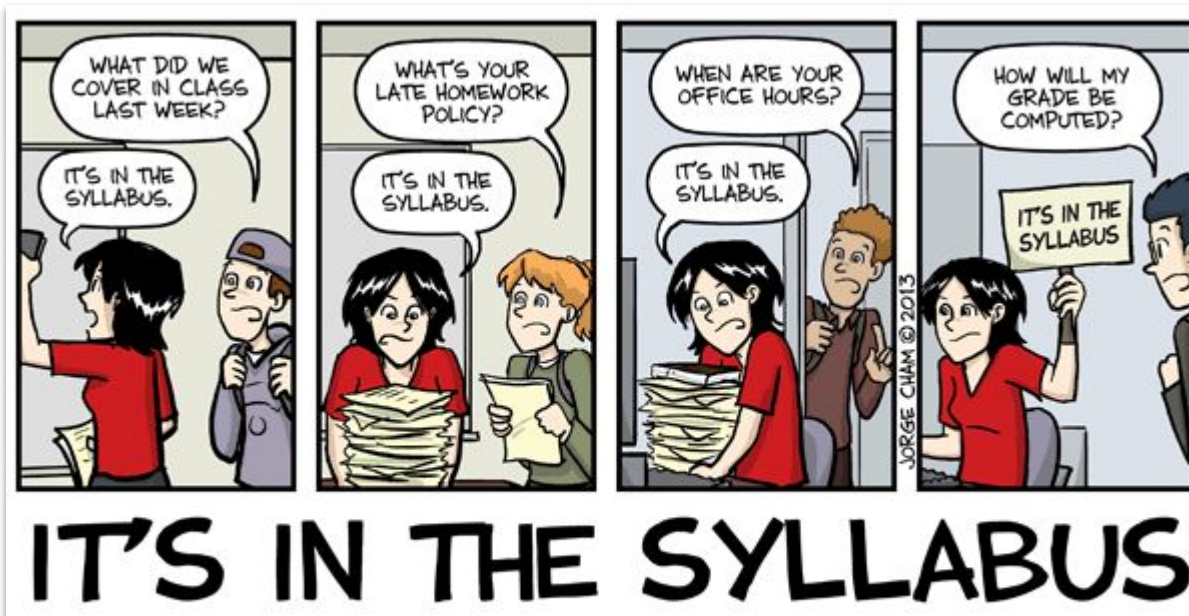


About myself

- Name: Marlos C. Machado
- I was born in Brazil
- I have been living in Edmonton for 10+ years
- I have 2 kids
- Ph.D. working on reinforcement learning
 - Interned at Microsoft Research, IBM Research, and DeepMind
- Worked 4 years at Google Brain and DeepMind
 - Among several other things, we deployed RL to fly balloons in the stratosphere



Course overview and logistics



Canvas: [link](#)

Slack: [link](#)

• My website: [link](#)

• Google drive: [link](#)

Start here!

University of Alberta

CMPUT 365: Introduction to Reinforcement Learning
LEC A1
Fall 2025

Instructor: Marlos C. Machado
Teaching Assistants: A. Naitpour, B. Kikar, D. Hart, L. Ouz, S. Chandrasekar, Ba G., T. Tan

Office: UCCMH 7041
E-mail: machado@ualberta.ca
Web Page: <https://courses.ualberta.ca/courses/27863>

Office hours: The location and time (p) which the TAs will hold office hours will be available on Canvas. Slack and Canvas: asynchronously

TA email address: comput365@ualberta.ca
Do not personally email the TAs. They will only respond via comput365@ualberta.ca.

Lecture room & time: ESB 3-27, MW 13:00 - 13:50
Attendance isn't mandatory, although strongly encouraged.

Slack invitation link: We will use Slack as an optional alternative to Canvas for communication and question-answering. The invitation link will be provided to the students on Canvas.

TERRITORIAL ACKNOWLEDGEMENT

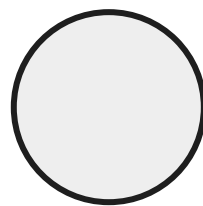
The University of Alberta respectfully acknowledges that we are situated on Treaty 6 territory, traditional lands of First Nations and Métis people.

COURSE CONTENT

Course Description: This course provides an introduction to reinforcement learning, which focuses on the study and design of learning agents that interact with a complex, uncertain world to achieve a goal. The course will cover multi-armed bandits, Markov decision processes, reinforcement learning, planning, and function approximation (deep supervised learning). The course will take an information-processing approach to the study of intelligence and briefly touch on perspectives from psychology, neuroscience, and philosophy.

Key resources

- Syllabus
 - Canvas, Slack, my website, Google Drive.
- Teaching assistants



Amirhossein • Bavish • Dikshant • Lucas • Siddarth • Shashank • Tian

- TA email address: cmput365@ualberta.ca
- My email address: machado@ualberta.ca
- Slack invitation link: [link](#)

I want to make this course a **safe** and **inclusive** environment, for everyone.

It is ok to make mistakes.

We should all strive to be **respectful** to each other.

If you want me to address you by a **different name**, or if you want to tell me your **pronouns**, I'm more than happy to hear!

Office hours

- Slack and Canvas: Asynchronous
- Marlos: *After class* @ *here*
- Lucas Cruz: *Monday* 15:00 – 17:00 @ UCOMM 2-138
- Siddarth Chandrasekar: *Tuesday* 13:00 – 15:00 @ UCOMM 3-162
- Tian Tian: *Wednesday* 15:00 – 17:00 @ UCOMM 2-138
- Dikshant: *Thursday* 09:00 – 11:00 @ UCOMM 2-138
- Amirhossein Rajabpour: *Thursday* 11:00 – 13:00 @ UCOMM 2-138
- Sai Shashank Gunuputi: *Friday* 10:00 – 12:00 @ UCOMM 2-138
- Bavish Kulur: *Friday* 15:00 – 17:00 @ UCOMM 2-138

Syllabus [[Canvas](#), [Slack](#), [website](#), [Google Drive](#)]

Pre-requisites

- CMPUT 175 or CMPUT 275
- CMPUT 267 or 466, or STAT 265
- Python
- Probability (e.g., expectations of random variables, conditional expectations)
- Calculus (e.g., partial derivatives)
- Linear algebra (e.g., vectors and matrices)

You should either be familiar with these topics or be ready to pick them up quickly as needed by consulting outside resources.

This will **not** be a flipped classroom!

- In the past, this course used to be taught in a flipped classroom
 - Roughly, you are initially introduced to **new topics outside** the classroom, using classroom time to explore topics in greater depth
- The number of students in this class has been steadily increasing, though
 - I don't know how to scale a flipped classroom without relying more and more on TAs to teach you
- Some of the feedback I received revolved around it feeling too repetitive
 - First read the textbook, watch the recorded lectures, do exercises, and then come to class
 - You can (and **should**) still do some of that before coming to class
- All this to say this will be a regular class, for better or for worse 😊

Required textbook

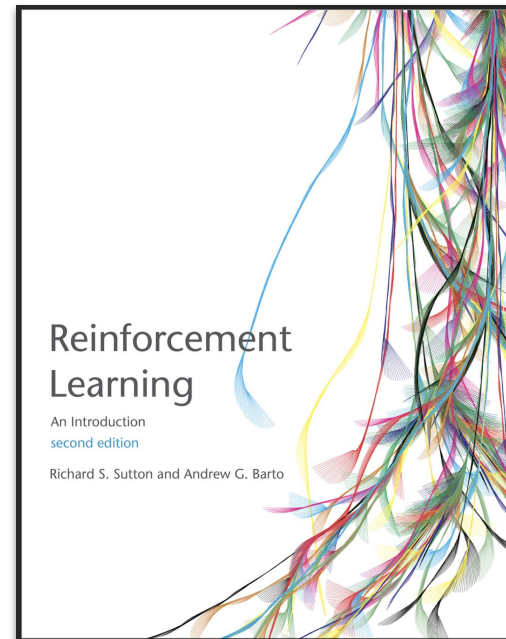
Reinforcement Learning: An Introduction

Richard S. Sutton & Andrew G. Barto

MIT Press. 2nd Edition.

<http://www.incompleteideas.net/book/the-book-2nd.html>

- You will need to read the book!
That's how you study for this course!
- The book is really good!



GRADE EVALUATION		
Assessment	Weight	Date
Practice quizzes (80% pass)	$9 \times 1\% = 9\%$	Day of the last class on the topic of the week at 23:59:59 (see Course schedule at the end for details)
Assessments (graded quizzes/notebooks on Coursera)	$9 \times 2.5\% = 22.5\%$	Day of the last class on the topic of the week at 23:59:59 (see Course schedule at the end for details)
Midterm 1 exam	20 %	October 3, 2025
Midterm 2 exam	20%	October 31, 2025
Final exam	30%	December 15, 2025*

GRADE EVALUATION

Assessment	Weight	Date
Practice quizzes (80% pass)	$9 \times 1\% = 9\%$	Day of the last class on the topic of the week at 23:59:59 (see Course schedule at the end for details)
Assessments (graded quizzes/notebooks on Coursera)	$9 \times 2.5\% = 22.5\%$	Day of the last class on the topic of the week at 23:59:59 (see Course schedule at the end for details)
Midterm 1 exam	20 %	October 3, 2025
Midterm 2 exam	20%	October 31, 2025
Final exam	30%	December 15, 2025*

GRADE EVALUATION

Assessment	Weight	Date
Practice quizzes (80% pass)	$9 \times 1\% = 9\%$	Day of the last class on the topic of the week at 23:59:59 (see Course schedule at the end for details)
Assessments (graded quizzes/notebooks on Coursera)	$9 \times 2.5\% = 22.5\%$	Day of the last class on the topic of the week at 23:59:59 (see Course schedule at the end for details)
Coursera, almost every week (<u>starting next week, Sep 12</u>): 31.5%		
Final exam	30%	December 15, 2025*

GRADE EVALUATION

Assessment	Weight	Date
Practice quizzes (80% pass)	$9 \times 1\% = 9\%$	Day of the last class on the topic of the week at 23:59:59 (see Course schedule at the end for details)
Assessments (graded quizzes/notebooks on Coursera)	$9 \times 2.5\% = 22.5\%$	Day of the last class on the topic of the week at 23:59:59 (see Course schedule at the end for details)

Coursera, almost every week (starting next week, Sep 12): 31.5%

Late submissions will not be accepted. There are 11 quizzes and 11 graded assignments. You're expected to do all of them, but s**t happens, so you can miss 2 of each and still get full marks.

GRADE EVALUATION		
Assessment	Weight	Date
Practice quizzes (80% pass)	$9 \times 1\% = 9\%$	Day of the last class on the topic of the week at 23:59:59 (see Course schedule at the end for details)
Assessments (graded quizzes/notebooks on Coursera)	$9 \times 2.5\% = 22.5\%$	Day of the last class on the topic of the week at 23:59:59 (see Course schedule at the end for details)
Midterm 1 exam	20 %	October 3, 2025
Midterm 2 exam	20%	October 31, 2025

Two midterms, summing to 40%. Closed book.

GRADE EVALUATION		
Assessment	Weight	Date
Practice quizzes (80% pass)	$9 \times 1\% = 9\%$	Day of the last class on the topic of the week at 23:59:59 (see Course schedule at the end for details)
Assessments (graded quizzes/notebooks on Coursera)	$9 \times 2.5\% = 22.5\%$	Day of the last class on the topic of the week at 23:59:59 (see Course schedule at the end for details)
Midterm 1 exam	20 %	October 3, 2025
Midterm 2 exam	20%	October 31, 2025

Two midterms, summing to 40%. Closed book.
If you miss the midterm, you can can apply for an excused absence.
If granted, the weight of the missed midterm will be deferred to the final.

GRADE EVALUATION		
Assessment	Weight	Date
Practice quizzes (80% pass)	$9 \times 1\% = 9\%$	Day of the last class on the topic of the week at 23:59:59 (see Course schedule at the end for details)
Assessments (graded quizzes/notebooks on Coursera)	$9 \times 2.5\% = 22.5\%$	Day of the last class on the topic of the week at 23:59:59 (see Course schedule at the end for details)
Midterm 1 exam	20 %	October 3, 2025
Midterm 2 exam	20%	October 31, 2025
Final exam	30%	December 15, 2025*

GRADE EVALUATION		
Assessment	Weight	Date
Practice quizzes (80% pass)	9 x 1% = 9%	Day of the last class on the topic of the week at 23:59:59 (see Course schedule at the end for details)
Assessments (graded quizzes/notebooks on Coursera)	9 x 2.5% = 22.5%	Day of the last class on the topic of the week at 23:59:59 (see Course schedule at the end for details)
The final, worth 30%, will be about the whole course. If you miss the final, you can apply to a deferred final examination.		
Final exam	30%	December 15, 2025*

GRADE EVALUATION		
Assessment	Weight	Date
Practice quizzes (80% pass)	$9 \times 1\% = 9\%$	Day of the last class on the topic of the week at 23:59:59 (see Course schedule at the end for details)
Assessments (graded quizzes/notebooks on Coursera)	$9 \times 2.5\% = 22.5\%$	Day of the last class on the topic of the week at 23:59:59 (see Course schedule at the end for details)
Midterm 1 exam	20 %	October 3, 2025
Midterm 2 exam	20%	October 31, 2025
Final exam	30%	December 15, 2025*

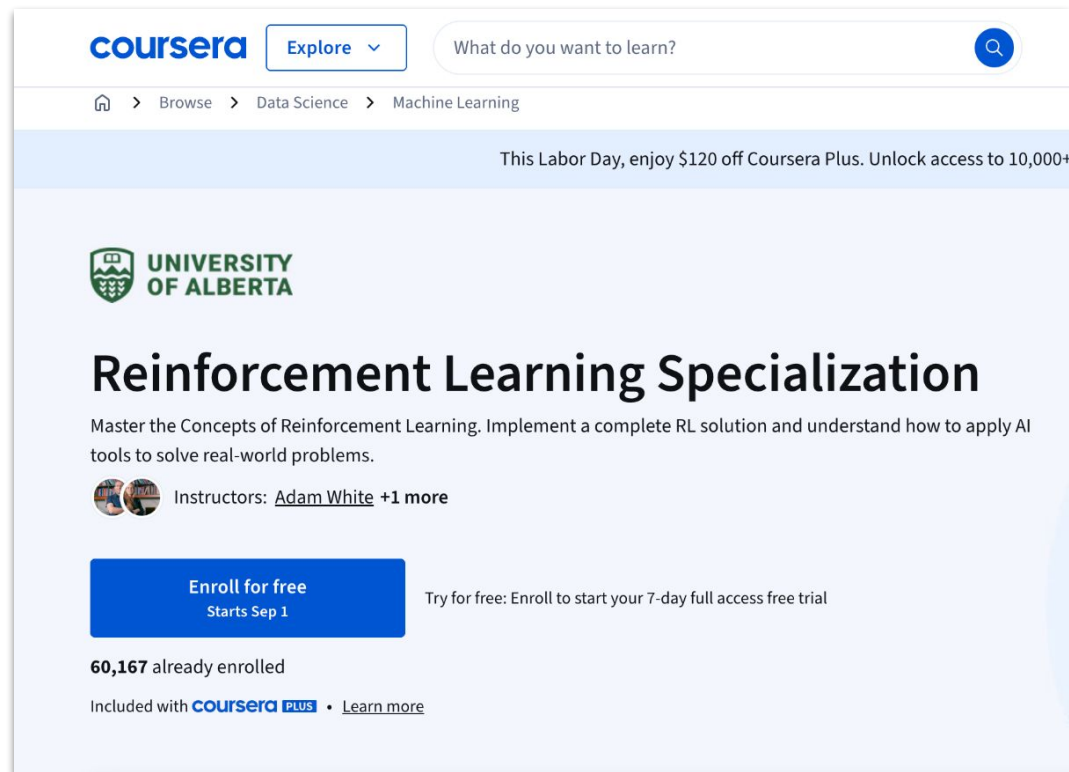
GRADE EVALUATION	
Assessment	Weight
Practice quizzes (80% pass)	9 x 1% = 9%
Assessments (graded quizzes/notebooks on Coursera)	9 x 2.5% = 22.5%
Midterm 1 exam	20 %
Midterm 2 exam	20%
Final exam	30%

Total: **101.5%.**
You can **not-submit 2 quizzes** and **2 assessments.**

Grades **will not** be rounded at the end, and **no more extra marks** will be given.
No exceptions.

Coursera

- Coursera will be essential to CMPUT 365
- You should have been added to a private session of the RL courses (we used your university's email)
 - If you don't have access you should let me know!
 - **IMPORTANT: If you don't use the private session you won't get credit for submitted work!**



The screenshot shows the Coursera website interface. At the top, the Coursera logo is on the left, followed by an 'Explore' button with a dropdown arrow. To the right is a search bar with the placeholder text 'What do you want to learn?' and a magnifying glass icon. Below the header, a breadcrumb trail reads 'Home > Browse > Data Science > Machine Learning'. A light blue banner across the page states: 'This Labor Day, enjoy \$120 off Coursera Plus. Unlock access to 10,000+'. The main content area features the University of Alberta logo and the course title 'Reinforcement Learning Specialization' in large, bold black text. Below the title, a description reads: 'Master the Concepts of Reinforcement Learning. Implement a complete RL solution and understand how to apply AI tools to solve real-world problems.' Underneath the description is a row of instructor profile pictures followed by the text 'Instructors: Adam White +1 more'. A prominent blue button with white text says 'Enroll for free' and 'Starts Sep 1'. To the right of this button, smaller text says 'Try for free: Enroll to start your 7-day full access free trial'. Below the button, it states '60,167 already enrolled'. At the bottom, it says 'Included with coursera PLUS' with a blue 'PLUS' badge, followed by a link to 'Learn more'.

Coursera

Coursera

RL

Search in course

Search

Viewing:

CMPUT 365-Fall 2025

Private

Live

September 1, 2025 - December 22, 2025

Fundamentals of Reinforcement Learning

Course Material

Module 1

Module 2

Module 3

Module 4

Module 5

Grades

Notes

Discussion Forums

Messages

Live Events

Classmates

An Introduction to Sequential Decision-Making

46 min of videos left

1h 10m of readings left

2 graded assessments left

For the first week of this course, you will learn how to understand the exploration-exploitation trade-off in sequential decision-making, implement incremental algorithms for estimating action-values, and compare the strengths and weaknesses to...

Show Learning Objectives

The K-Armed Bandit Problem

Module 1 Learning Objectives
Reading • 10 min

Weekly Reading
Reading • 30 min

Let's play a game!
Ungraded Plugin • 15 min

Sequential Decision Making with Evaluative Feedback
Video • 5 min

Compare bandits to supervised learning
Discussion Prompt • 10 min

What to Learn? Estimating Action Values

Learning Action Values
Video • 4 min

Academic integrity

- [Code of Student Behaviour](#)
- [Student Conduct Policy](#)
- [Academic Integrity website](#)
- **Appropriate collaboration:** You are allowed to discuss the quizzes and assignments with your classmates. Note, however, that you are not allowed to exchange any written text, code, or to give and/or receive detailed step-by-step instructions on how to solve the proposed problems.
- **Cell phones:** Cell phones are to be turned off during lectures, labs and seminars.
- **Recording and/or Distribution of Course Materials:** Audio or video recording, digital or otherwise, by students is allowed only with my prior written consent as a part of an approved accommodation plan.

Academic integrity – **Expectations for AI use**

The primary goal of this course is to foster *individual* critical, creative thinking, and problem-solving skills related to reinforcement learning. Thus, in order to achieve such learning outcomes, you can submit each practice quiz and graded assignment multiple times, which allows for many learning opportunities.

The use of advanced AI-tools based on large-language models such as ChatGPT is **strictly prohibited** for all quizzes and graded assignments. The only exception is their use for Python-related queries (but the use of such tools to help with the programming assignments themselves is still strictly prohibited).

As stated in the university's [AI-Squared - Artificial Intelligence and Academic Integrity](#) webpage, “*learning is not only about the product; learning is also about the process of acquiring new knowledge or learning ways to think and reason.*”

Schedule

- The course will be structured in “weeks”. **Not every week starts on Monday**
- We have 12 weeks of content classes and we’ll cover 13 weeks of the MOOC
 - This corresponds to 9 chapters of the textbook

Schedule

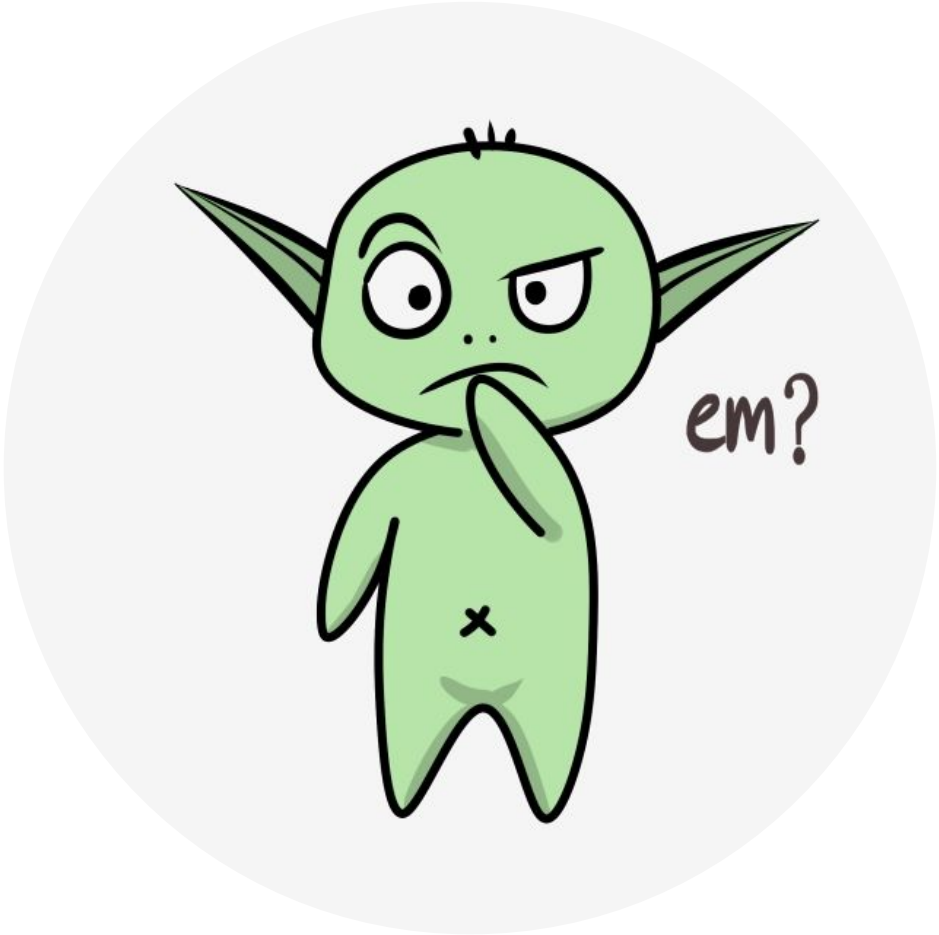
- The course will be structured in “weeks”. **Not every week starts on Monday**
- We have 12 weeks of content classes and we’ll cover 13 weeks of the MOOC
 - This corresponds to 9 chapters of the textbook
- A practice quiz and a graded assignment are due at the end of each “week” in terms of content – You should look at the syllabus / schedule all the time
- The deadline for submitting assignments and quizzes is 23:59:59

Schedule

Course Schedule & Assigned Readings

Week	Date	Topic	Deadlines (all due at 23:59:59)	Readings
0	Wed, Sep 3	Course overview Discussion about what is reinf. learning		
1	Fri, Sep 5	Fundamentals of RL: An introduction to sequential decision-making		Chapter 2, up to §2.7 (pp. 25-36), and §2.10 (pp. 42-44)
0	Mon, Sep 8	[NO CLASS – RECORDED LECTURE] Background review: Probability, statistics, linear algebra, and calculus		
0	Wed, Sep 10	Guest Lecture: Richard Sutton		
1	Fri, Sep 12	Fundamentals of RL: An introduction to sequential decision-making	Practice quiz and Progr. assignment (Bandits & exploration / exploitation)	
2	Mon, Sep 15	Fundamentals of RL: Markov decision processes (MDPs)		Chapter 3, up to §3.3 (pp. 47-56)
2	Wed, Sep 17	Fundamentals of RL: Markov decision processes (MDPs)	Practice quiz (MDPs)	
3	Fri, Sep 19	Fundamentals of RL: Value functions & Bellman equations		Chapter 3, §3.5-§3.8 (pp. 58-69)
3	Mon, Sep 22	Fundamentals of RL: Value functions & Bellman equations		
3	Wed, Sep 24	Fundamentals of RL: Value functions & Bellman equations	Practice and Graded quiz (Value functions & Bellman equations)	
4	Fri, Sep 26	Fundamentals of RL: Dynamic programming		Chapter 4, §4.1-§4.4 (pp. 73-84); §4.6-§4.7 (pp. 86-89)
4	Mon, Sep 29	Fundamentals of RL: Dynamic programming	Practice quiz and Progr. assignment (Optimal policies with dyn. progr.)	
-	Wed, Oct 1	General Overview and Q&A		

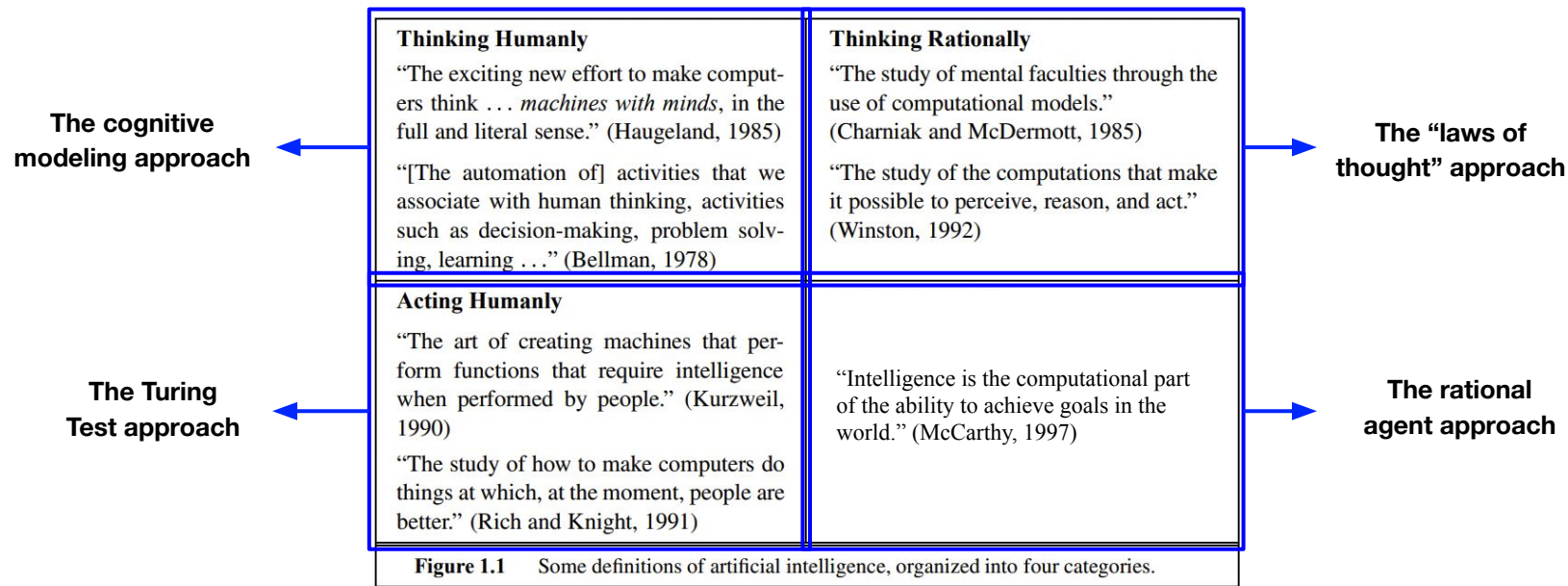
Syllabus [[Canvas](#), [Slack](#), [website](#), [Google Drive](#)]



What is reinforcement learning?

Artificial intelligence

“AI is the ability of machines to perform tasks that are typically associated with human intelligence, such as learning and problem-solving.” –Wikipedia



(Russell & Norvig; 2010)

Artificial intelligence

“AI is the ability of machines to perform tasks that are typically associated with human intelligence, such as learning and problem-solving.” –Wikipedia

The less a science has advanced, the more its terminology tends to rest on an uncritical assumption of mutual understanding.

– W. V. Quine



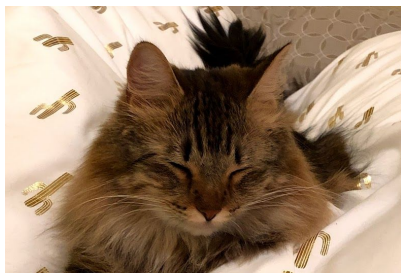
Machine learning

Machine learning is a subfield of AI in which the system's desired behavior is not explicitly programmed, instead it is *learned* from data

Machine learning

Machine learning is a subfield of AI in which the system's desired behavior is not explicitly programmed, instead it is *learned* from data

- “*Supervised learning* is learning from a training set of labeled examples provided by a knowledgeable external supervisor” (Sutton & Barto; 2018)



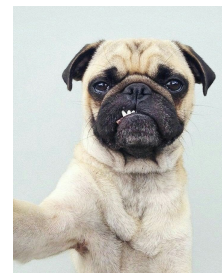
Cat



Cat



Not cat



Cat or not cat?

Machine learning

Machine learning is a subfield of AI in which the system's desired behavior is not explicitly programmed, instead it is *learned* from data

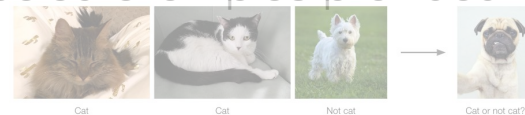
- “*Supervised learning* is learning from a training set of labeled examples provided by a knowledgeable external supervisor” (Sutton & Barto; 2018)
- “*Unsupervised learning* is typically about finding structure hidden in collections of unlabeled data” (Sutton & Barto; 2018)



Machine learning

Machine learning is a subfield of AI in which the system's desired behavior is not explicitly programmed, instead it is *learned* from data

- “*Supervised learning* is learning from a training set of labeled examples provided by a knowledgeable external supervisor” (Sutton & Barto; 2018)
- “*Unsupervised learning* is typically about finding structure hidden in collections of unlabeled data” (Sutton & Barto; 2018)



... and *reinforcement learning*!

Artificial intelligence

Machine learning

Reinforcement learning

Reinforcement learning

Reinforcement learning is a computational approach to learning from interaction to maximize a numerical reward signal (Sutton & Barto; 2018)



Reinforcement learning

Reinforcement learning is a computational approach to learning from interaction to maximize a numerical reward signal (Sutton & Barto; 2018)

- The idea of learning by interacting with our environment is very natural
- It is based on the idea of a learning system that wants something, and that adapts its behavior to get that



Reinforcement learning

Reinforcement learning is a computational approach to learning from interaction to maximize a numerical reward signal (Sutton & Barto; 2018)

- The idea of learning by interacting with our environment is very natural
- It is based on the idea of a learning system that wants something, and that adapts its behavior to get that



Some features are unique to reinforcement learning:

- Trial-and-error
- The trade-off between exploration and exploitation
- The delayed credit assignment / delayed reward problem

Artificial intelligence

Machine learning

Reinforcement learning

Reinforcement learning

Reinforcement learning is a computational paradigm for learning from interaction to maximize a numerical reward signal (Sutton & Barto, 2018)

- The idea of learning by interacting with an environment is very natural
- It is based on the idea of a *learning agent* that wants something, and that *action* to get that

Some features are unique to reinforcement learning:

- Trial-and-error learning
- The trade-off between exploration and exploitation
- The delayed reward / delayed reward problem

Problem or solution?



RL is now commonly deployed in the real-world

- **Recommendation systems**
 - Ads, news articles, videos, etc
- **General game playing**
 - Go, Chess, Shogi, Atari 2600, Starcraft, Minecraft, Gran Turismo
- **Industrial automation**
 - Cooling commercial buildings
 - Inventory management
 - Gas turbine optimization
 - Optimizing combustion in coal-fired power plants
- **Algorithms**
 - Video compression on YouTube
 - Faster matrix multiplication
 - Faster sorting algorithms
- **Control / Robotics**
 - Navigating stratospheric balloons
 - Plasm control for nuclear fusion
- **And more (see Csaba's [slides](#))**
 - COVID-19 border testing
 - Conversational agents
 - ...

The 2024 ACM A.M. Turing Award Winners “Created” RL



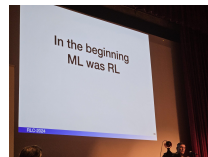
Association for
Computing Machinery

On intelligence, AGI, ASI, etc etc...

- People in the field have different, non-competing, perspectives and motivations
 - Some study RL to learn about / develop tools for solving sequential decision-making problems
 - Some look at RL as a computational model of intelligence

On intelligence, AGI, etc etc...

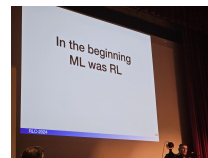
- People in the field have different, non-competing, perspectives and motivations
 - Some study RL to learn about / develop tools for solving sequential decision-making problems
 - Some look at RL as a computational model of intelligence
- RL was originally developed to understand intelligence/the brain
 - We should develop a critical view around these topics, and an ability to recognize hype / PR pieces



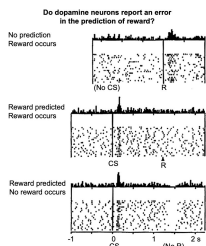
A. Barto (2024)

On intelligence, AGI, etc etc...

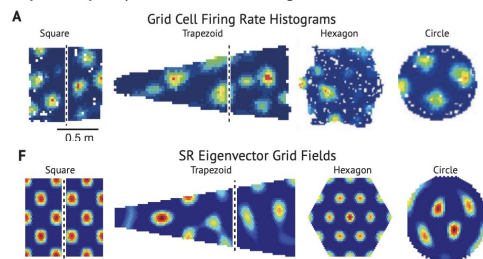
- People in the field have different, non-competing, perspectives and motivations
 - Some study RL to learn about / develop tools for solving sequential decision-making problems
 - Some look at RL as a computational model of intelligence
- RL was originally developed to understand intelligence/the brain
 - We should develop a critical view around these topics, and an ability to recognize hype / PR pieces
- Both perspectives are valid and both had had successes in the past **But they are different!!**



A. Barto (2024)



(Schultz, Dayan, & Montague; 1997)

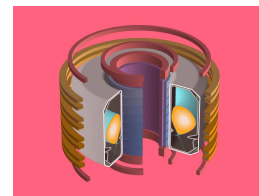


(Stachenfeld, Botvinick, & Gershman; 2017)

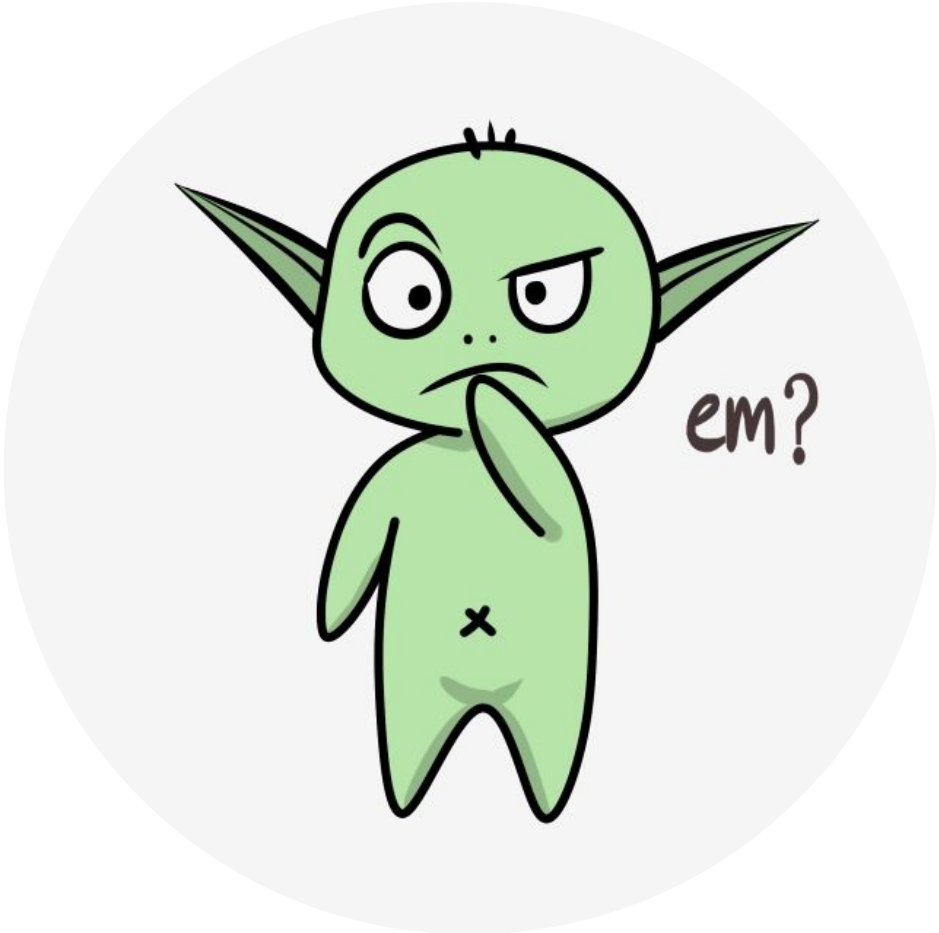
(Silver et al.; 2016)



(Degraeve et al.; 2022)



(Bellemare et al.; 2020)



Next class

- What **I** plan to do:
 - Fundamentals of RL: An introduction to sequential decision-making (Bandits)
 - I will have to miss class on Monday and Wednesday next week
 - I will record a background review so you can watch (if you want)
 - ***Rich Sutton will give a guest lecture on Wednesday***
- What I recommend **YOU** to do for next class:
 - Make sure you have access to Coursera, Canvas, and Slack
 - Read Chapter 1 of the textbook (not mandatory)
 - Read Chapter 2 of the textbook up to §2.7 (inclusive)