

Deep Reinforcement Learning and Control

Introduction

Fall 2021, CMU 10-703

Instructors:

Katerina Fragkiadaki

Russ Salakhutdinov



Course Logistics

- Course website: https://cmudeepri.github.io/703website_f21/ all you need to know
- Grading:
 - 4 Homework assignments: implementation and question/answering many optional and extra grade questions - 60%
 - 3 quizzes - 40%
- Resources: AWS for those that do not have access to GPUs
- People can audit the course
- The readings on the schedule are required unless noted otherwise

Overview for today

- Goal of the course / why it is important
- What is reinforcement learning
- What is representation learning (and how it helps reinforcement learning and behavior learning in general)
- Reinforcement learning versus supervised learning
- AI's paradox: what is hard and what is easy in behavior learning

Overview for today

- Goal of the course / why it is important
- What is reinforcement learning
- What is representation learning (and how it helps reinforcement learning and behavior learning in general)
- Reinforcement learning versus supervised learning
- AI's paradox: what is hard and what is easy in behavior learning

Goal of the course: Learning to act

Building agents that **learn** to act
and accomplish **goals** in **dynamic**
environments



Goal of the course: Learning to act

Building agents that **learn** to act
and accomplish **goals** in **dynamic**
environments



...as opposed to agents that execute
pre-programmed behaviors in **static**
environments...



Motion and Action are important

“The brain evolved, not to think or feel, but to control movement.”

Daniel Wolpert



Daniel Wolpert: The real reason for brains | TED Talk | TED.com

https://www.ted.com/talks/daniel_wolpert_the_real_reason_for_brains ▼

Motion and Action are important

“The brain evolved, not to think or feel, but to control movement.”

Daniel Wolpert

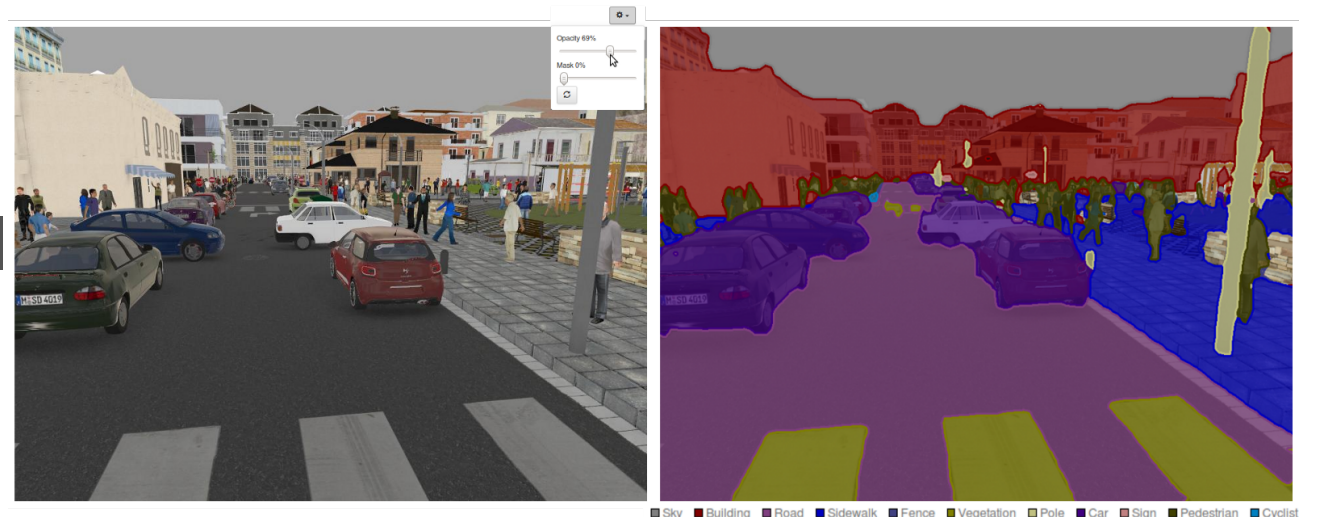


Sea squirts digest their own brain when they decide not to move anymore

Learning to act

- It is considered the most **biologically plausible** objective for learning
- It addresses the **full problem of making artificial agents that act in the world**, so it is driven by the right end goal

...in contrast to, for example,
making artificial agents that label
pixels in images

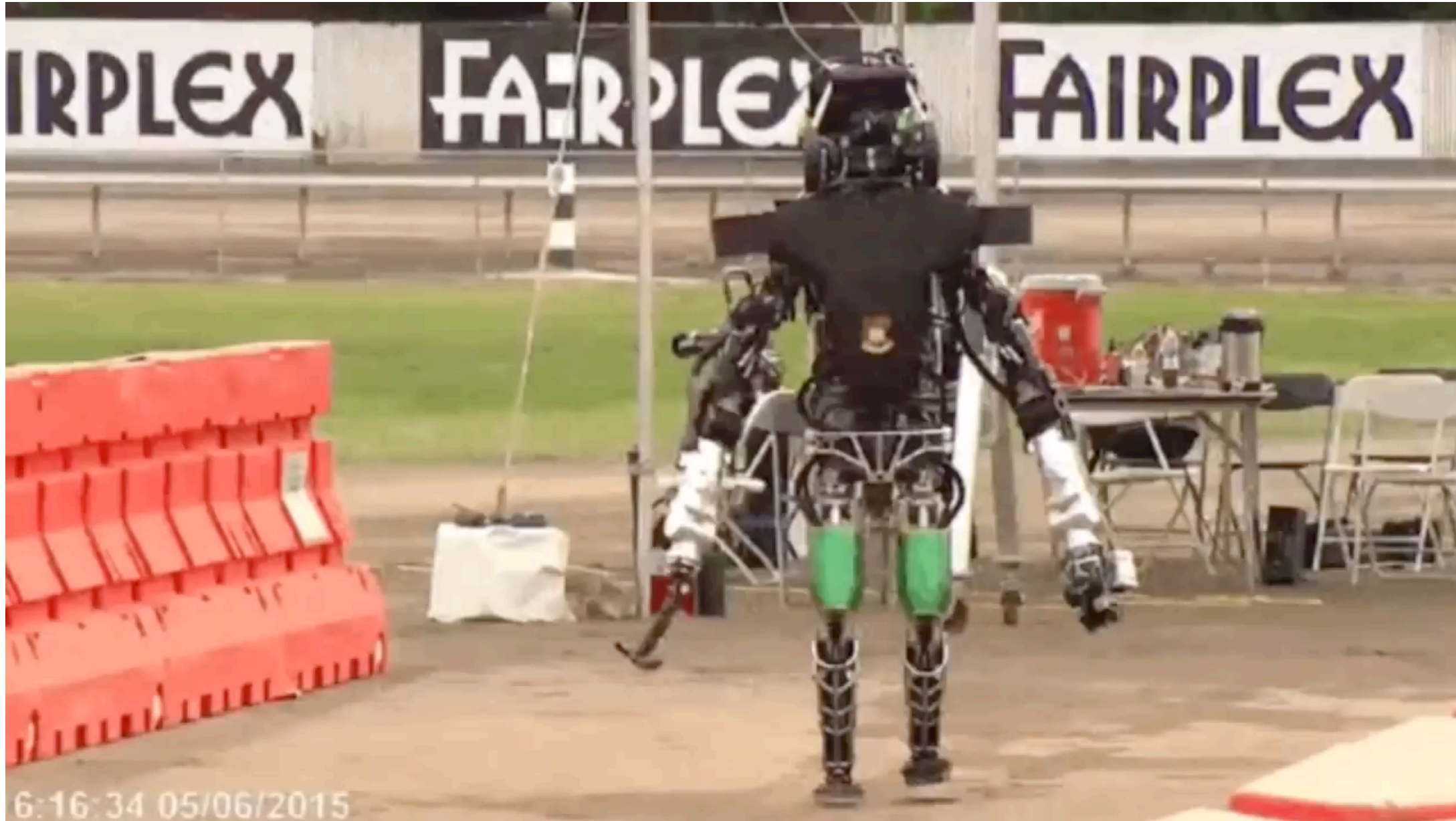


How far are we?



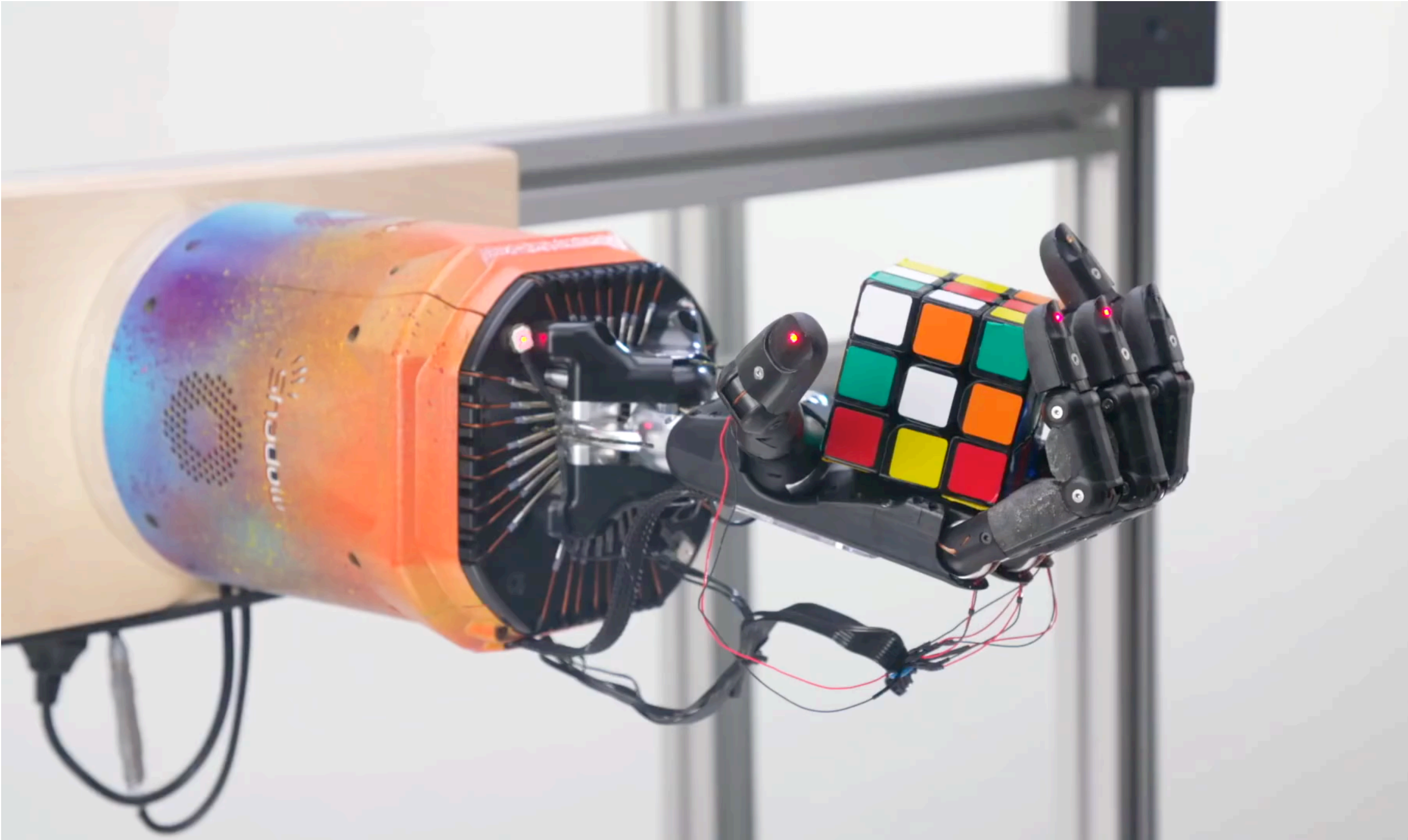
Here the robot is tele-operated: it does not actually operate on its own.

How far are we?



Here the robot operates on its own.

How far are we?



Here the robot operates on its own.

Questions/tasks the course aims to answer/address

- Discovering a behavior through trial-and-error guided by rewards.
- Generalizing/transferring a behavior across different scenarios (camera viewpoints, object identities, objects arrangements) E.g., you show me how to open one door, and I now need to learn how to open other similar doors

Questions/tasks the course aims to answer/address

- Discovering a behavior through trial-and-error guided by rewards.
 - Many algorithm here start tabula rasa: no previous knowledge of anything.
 - Environment doesn't change (camera and objects).
- Generalizing/transferring a behavior across different scenarios (camera viewpoints, object identities, objects arrangements) E.g., you show me how to open one door, and I now need to learn how to open all other doors

Questions/tasks the course aims to answer/address

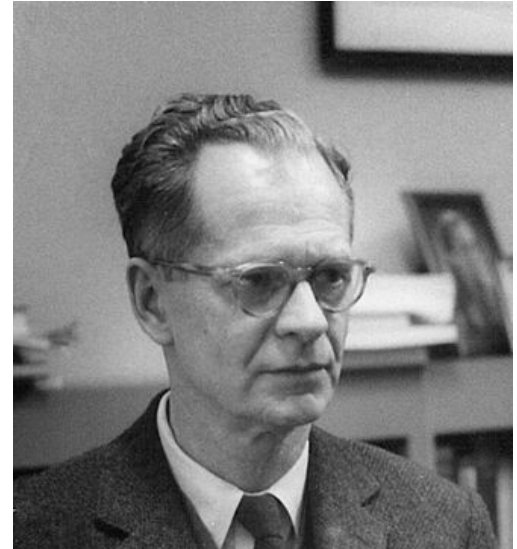
- Discovering a behavior through trial-and-error guided by rewards. E.g., today I discovered how to avoid the ads in y2mate.com, and I also discovered how (many times I need) to turn the key to open the door in the apt.
- Generalizing/transferring a behavior across different scenarios (camera viewpoints, object identities, objects arrangements) E.g., you show me how to open one door, and I now need to learn how to open all other doors
 - We do not start tabula rasa: we have knowledge which we enrich with trial-and-error. Our accomplishments are added to this knowledge with the goal to transfer faster in the future

Overview for today

- Goal of the course / why it is important
- What is reinforcement learning
- What is representation learning (and how it helps reinforcement learning and behavior learning in general)
- Reinforcement learning versus supervised learning
- AI's paradox: what is hard and what is easy in behavior learning

Reinforcement Learning (RL): How behaviors are shaped

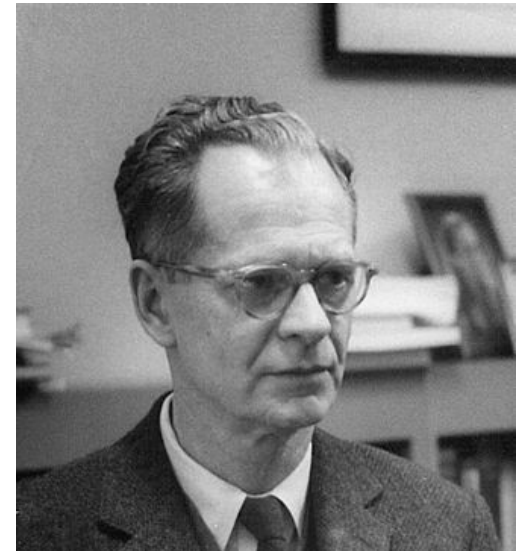
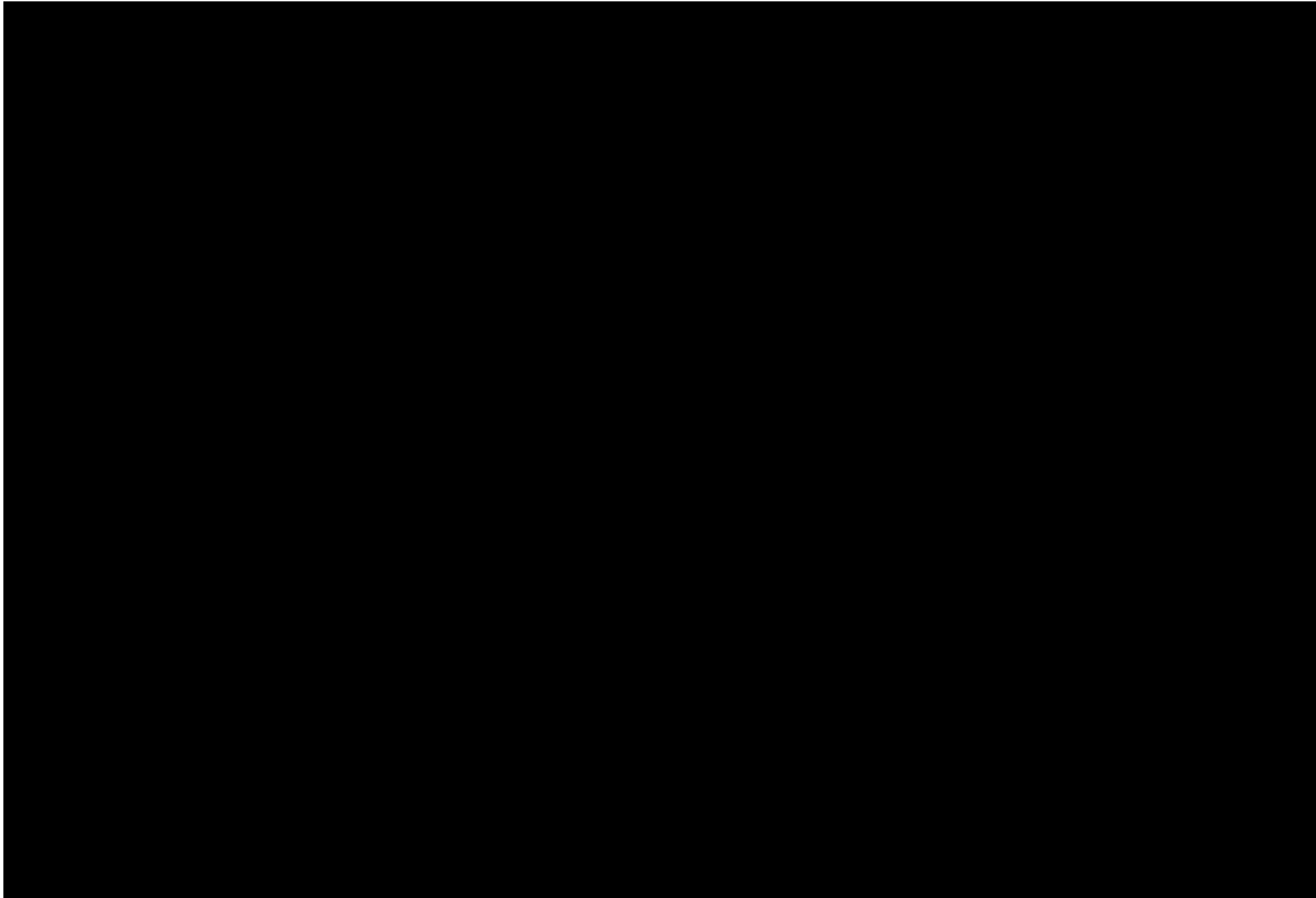
Behavior is primarily shaped by reinforcement rather than free-will.



B.F. Skinner
1904-1990
Harvard psychology

- behaviors that result in praise/pleasure tend to repeat,
- behaviors that result in punishment/pain tend to become extinct.

Reinforcement Learning (RL): How behaviors are shaped



B.F. Skinner
1904-1990
Harvard psychology

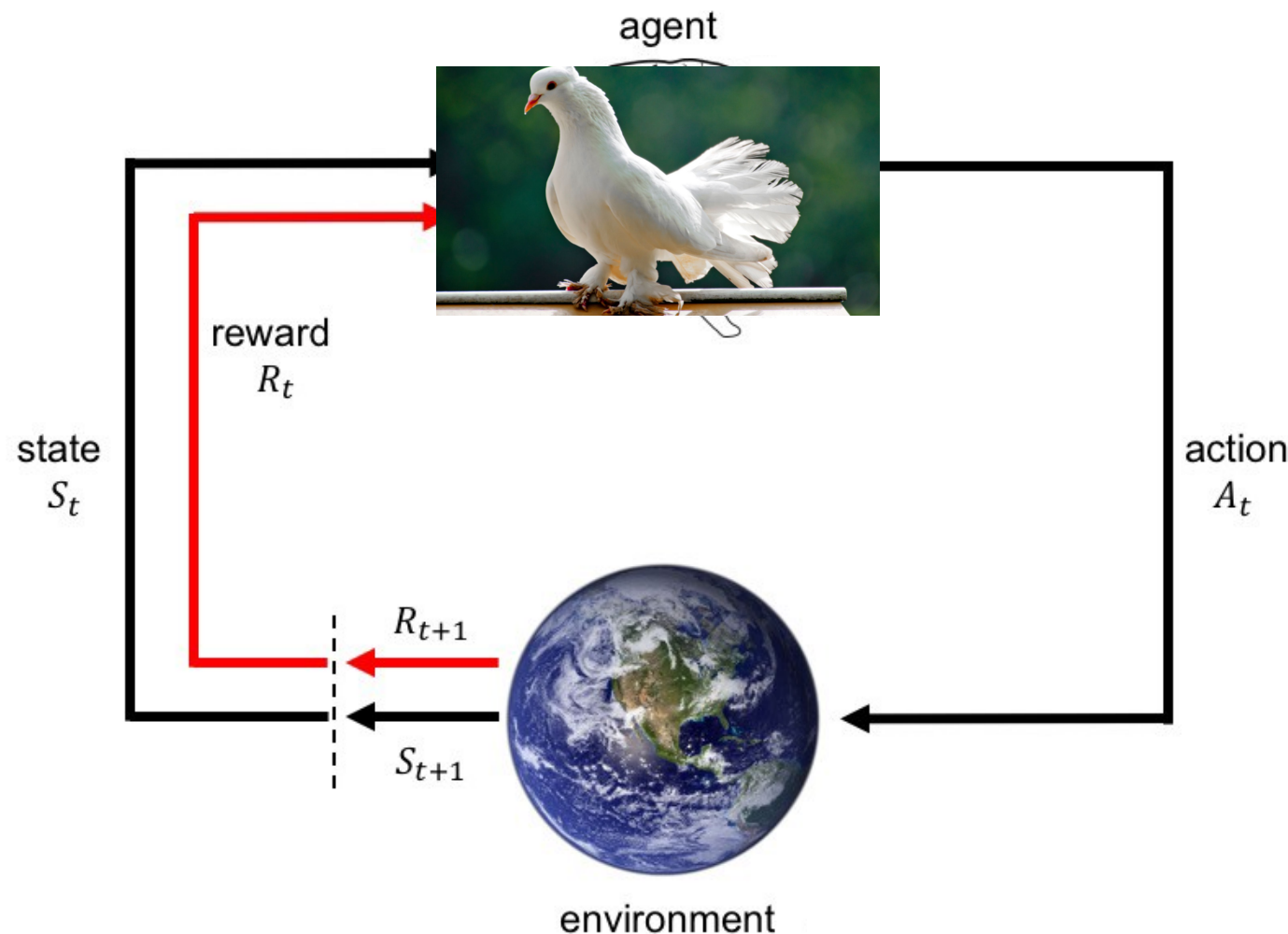
Q: Is the pigeon here
transferring or
discovering?

<https://www.youtube.com/watch?v=yhvaSEJtOV8>

Interesting finding: Pigeons become addicted to pecking under variable (non-consistent) rewarding

Reinforcement learning = trial-and-error learning

Learning policies that maximize a reward function by interacting with the world



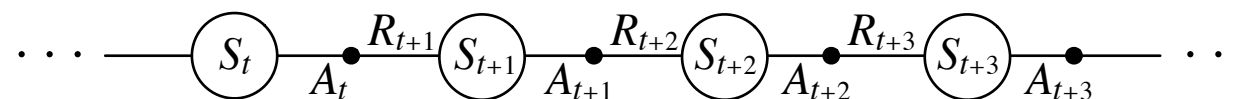
Agent and environment interact at discrete time steps: $t = 0, 1, 2, K$

Agent observes state at step t : $S_t \in \mathcal{S}$

produces action at step t : $A_t \in \mathcal{A}(S_t)$

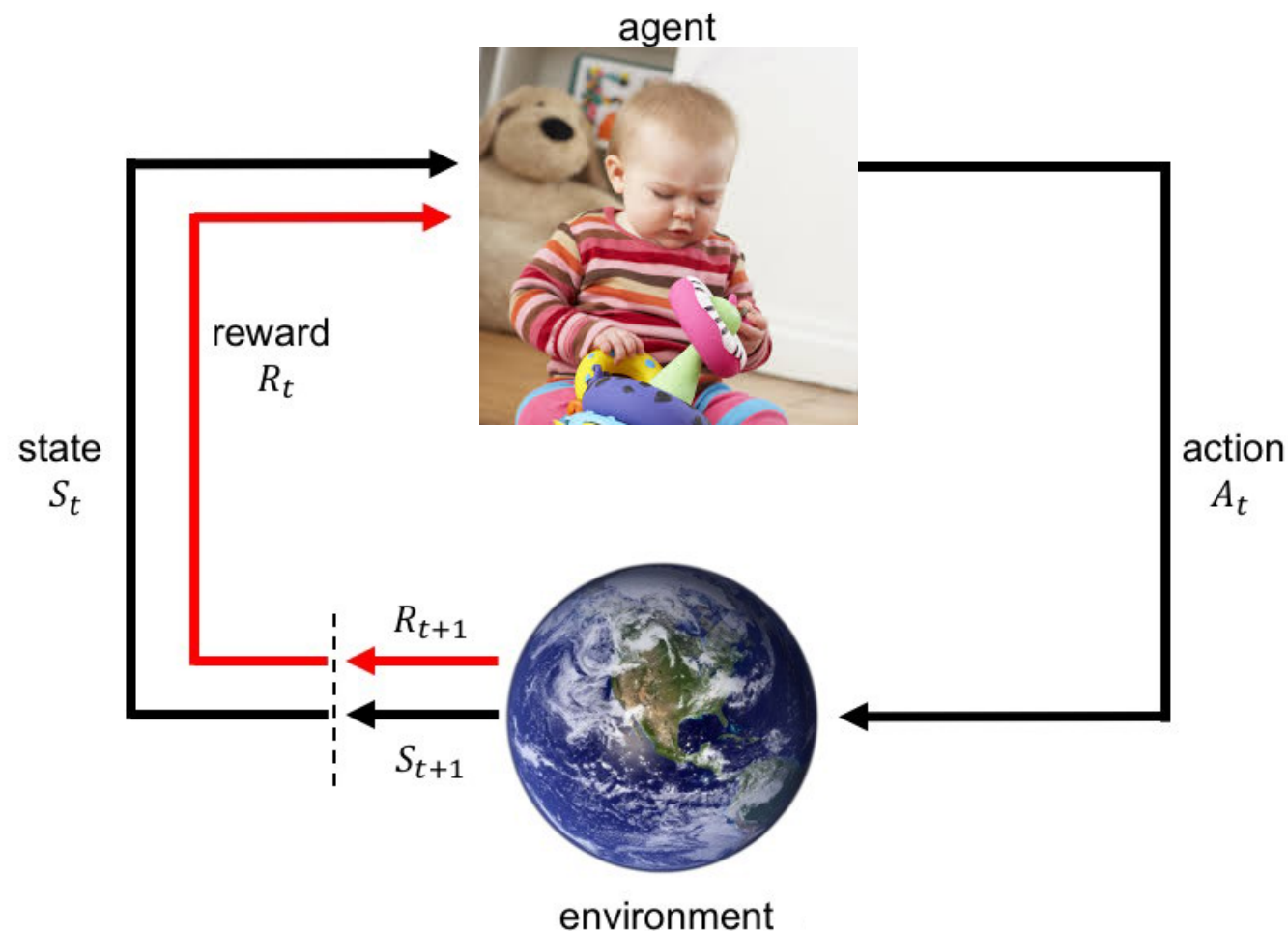
gets resulting reward: $R_{t+1} \in \mathbb{R}$

and resulting next state: $S_{t+1} \in \mathcal{S}^+$



Reinforcement learning

Rewards can be intrinsic, i.e., generated by the agent and guided by its curiosity as opposed to the external environment.



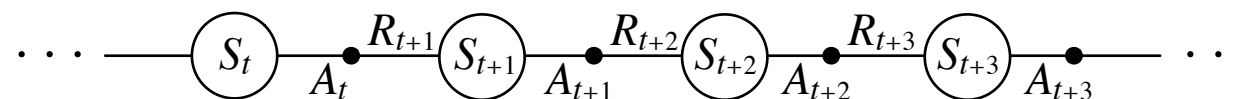
Agent and environment interact at discrete time steps: $t = 0, 1, 2, K$

Agent observes state at step t : $S_t \in \mathcal{S}$

produces action at step t : $A_t \in \mathcal{A}(S_t)$

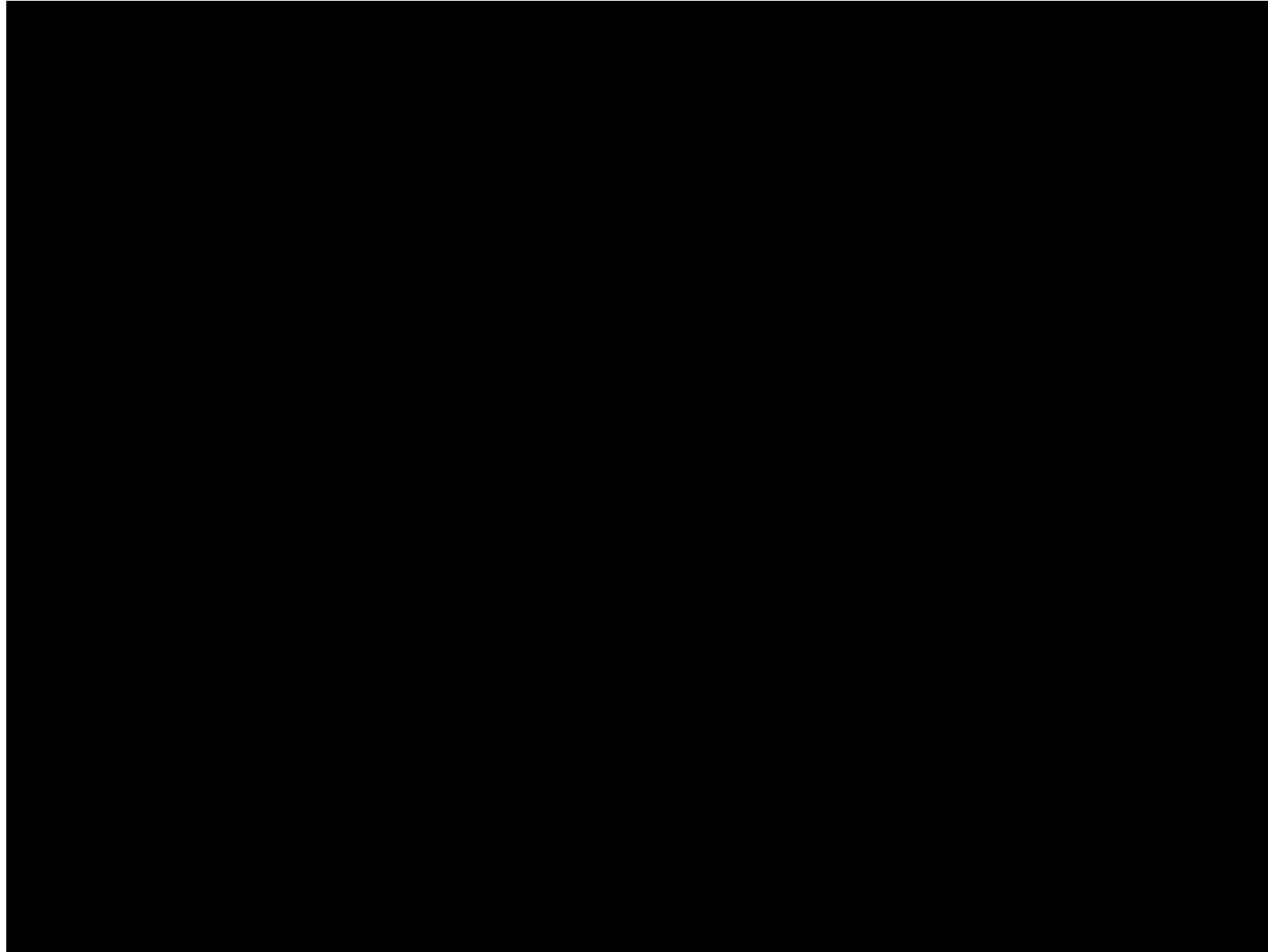
gets resulting reward: $R_{t+1} \in \mathbb{R}$

and resulting next state: $S_{t+1} \in \mathcal{S}^+$



Reinforcement learning

Rewards can be intrinsic, i.e., generated by the agent and guided by its **curiosity** as opposed to the external environment.



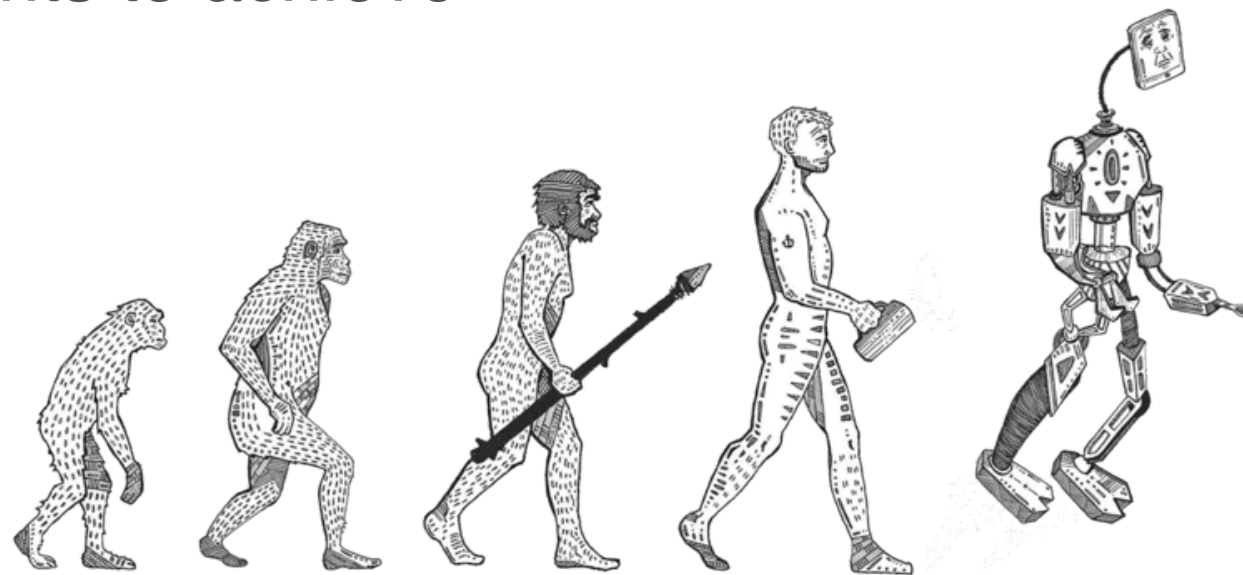
<https://youtu.be/8vNxjwt2AqY>

No food shows up but the baby keeps exploring

Agent

An entity that is equipped with

- sensors, in order to sense the environment,
- end-effectors in order to act in the environment, and
- goals that she wants to achieve



Actions

They are used by the agent to interact with the world:

- Play song with title “Imagine” / lower the lights / increase the volume / call grandma etc..
- Display advertisement , suggest song / movie etc..
- Go straight / turn k degrees / brake etc..
- Robot torques
- Desired gripper translation, rotation, opening

States

- A state captures whatever information is available to the agent at step t about its environment.
- The state can include immediate observations, highly processed observations, and structures built up over time from sequences of sensations, memories etc.

Observations

- An observation a.k.a. sensation: the (raw) input of the agent's sensors, images, tactile signal, waveforms, etc.



Policy

A mapping function from states to actions of the end effectors.

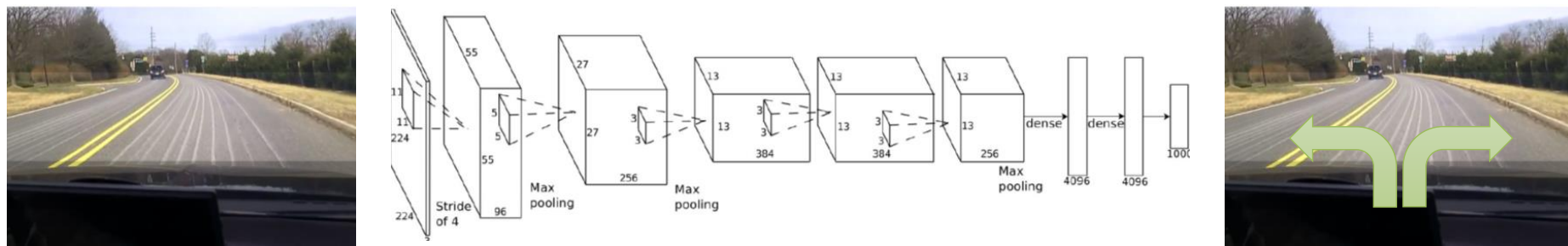
$$\pi(a \mid s) = \mathbb{P}[A_t = a \mid S_t = s]$$

Policy

A mapping function from states to actions of the end effectors.

$$\pi(a \mid s) = \mathbb{P}[A_t = a \mid S_t = s]$$

It can be a shallow or a deep function mapping

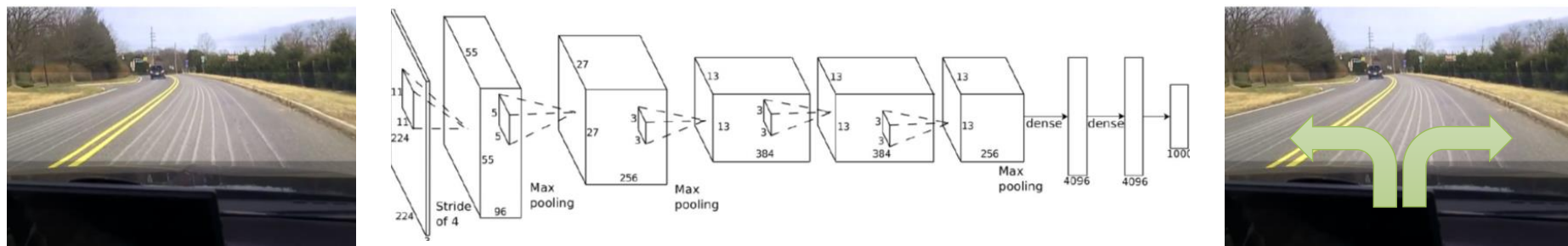


Policy

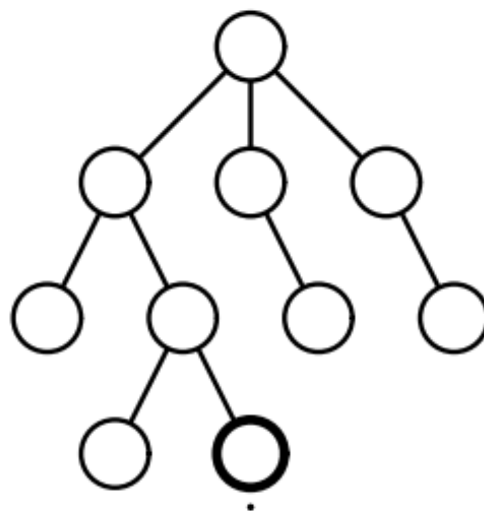
A mapping function from states to actions of the end effectors.

$$\pi(a \mid s) = \mathbb{P}[A_t = a \mid S_t = s]$$

It can be a shallow or a deep function mapping



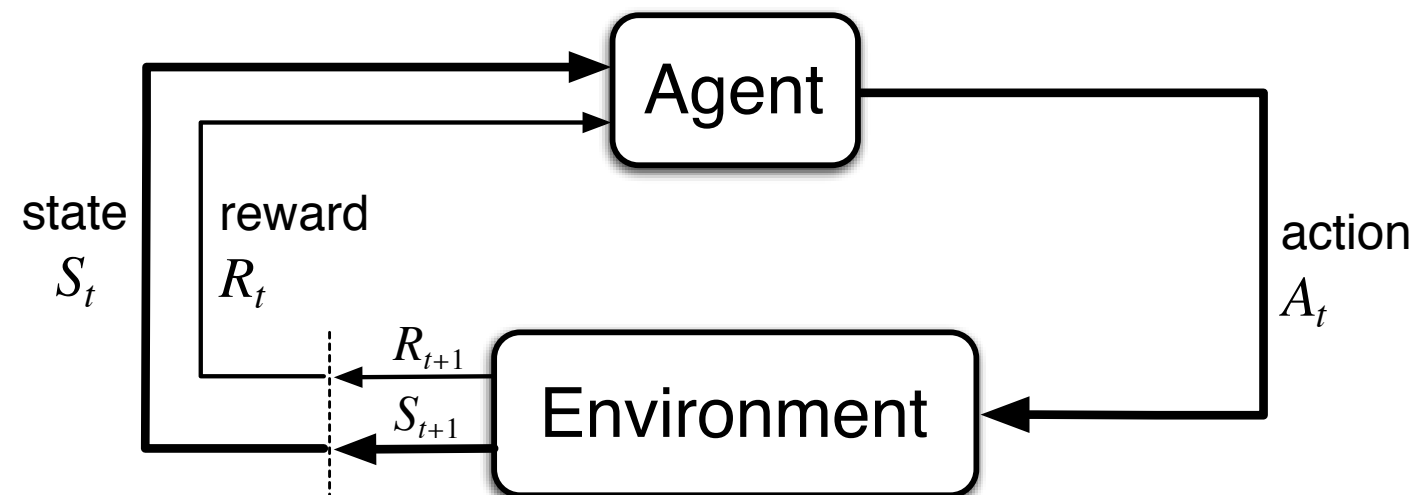
or it can be as complicated as involving a tree look-ahead search



Closed loop sensing and acting

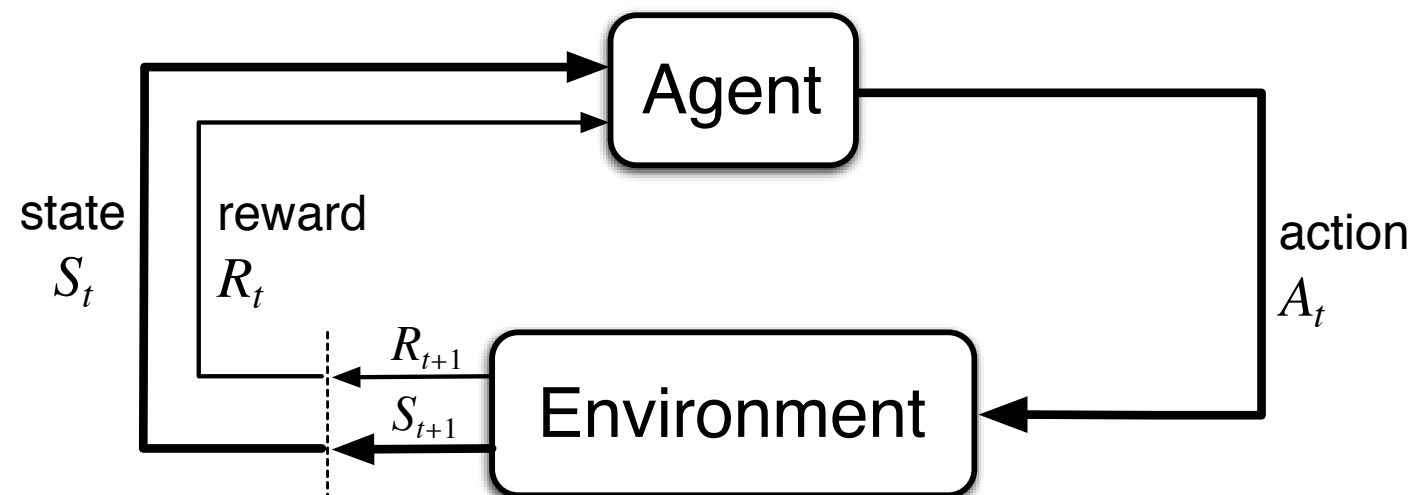
Imagine an agent that wants to pick up an object and has a policy that predicts what the actions should be for the next 2 secs ahead.

This means, for the next 2 secs we switch off the sensors, and just execute the predicted actions. In the next second, due to imperfect sensing, the object is about to fall!



Closed loop sensing and acting

Sensing is always imperfect. Our excellent motor skills are due to continuous sensing and updating of the actions, a.k.a. servoing. So the perception-action loop is in fact extremely short in time.



Rewards

They are scalar values provided provided to the agent that indicate whether goals have been achieved, e.g., 1 if goal is achieved, 0 otherwise, or -1 for overtime step the goal is not achieved

- Rewards specify **what** the agent needs to achieve, not **how** to achieve it.
- The simplest and cheapest form of supervision, and surprisingly general: All of what we mean by goals and purposes can be encoded mathematically as the maximization of the cumulative sum of a received scalar signal (reward)

Returns

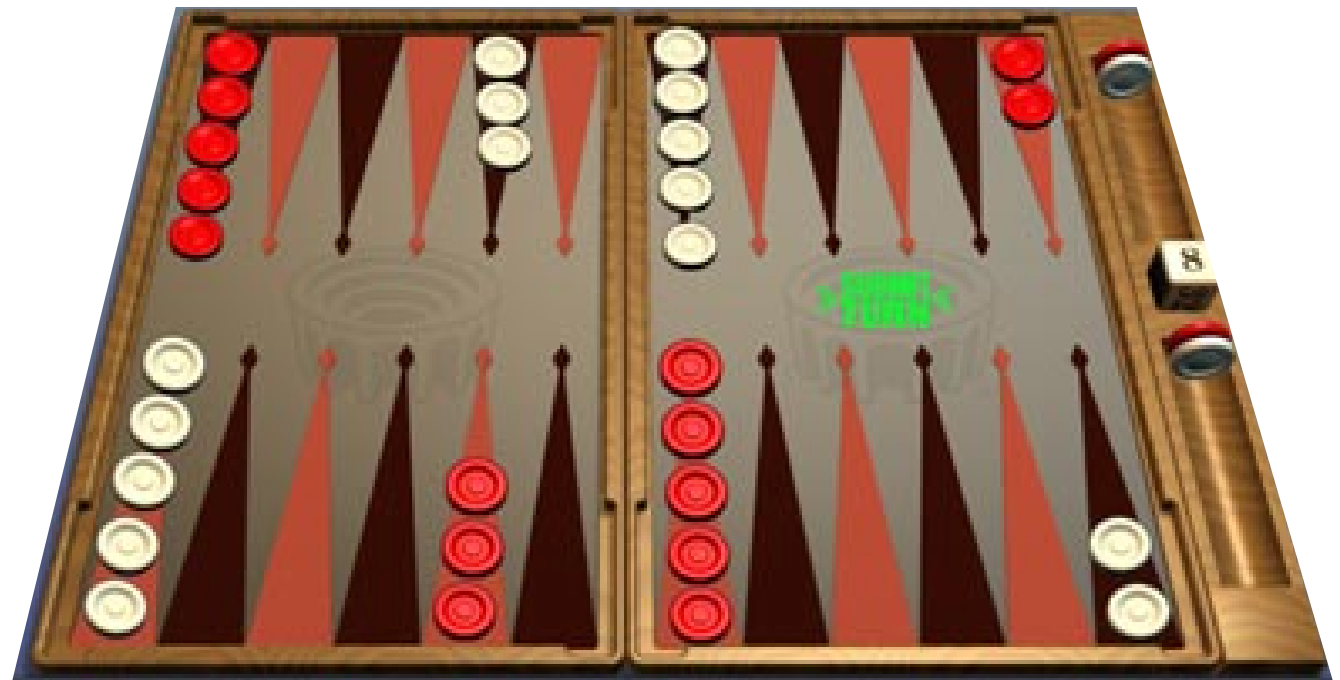
Goal-seeking behavior of an agent can be formalized as the behavior that seeks maximization of the expected value of the **cumulative sum** of (potentially time discounted) rewards, we call it return.

We want to maximize returns.

$$G_t = R_{t+1} + R_{t+2} + \dots + R_T$$

Example: Backgammon

- States: Configurations of the playing board ($\approx 10^{20}$)
- Actions: Moves
- Rewards:
 - win: +1
 - lose: -1
 - else: 0



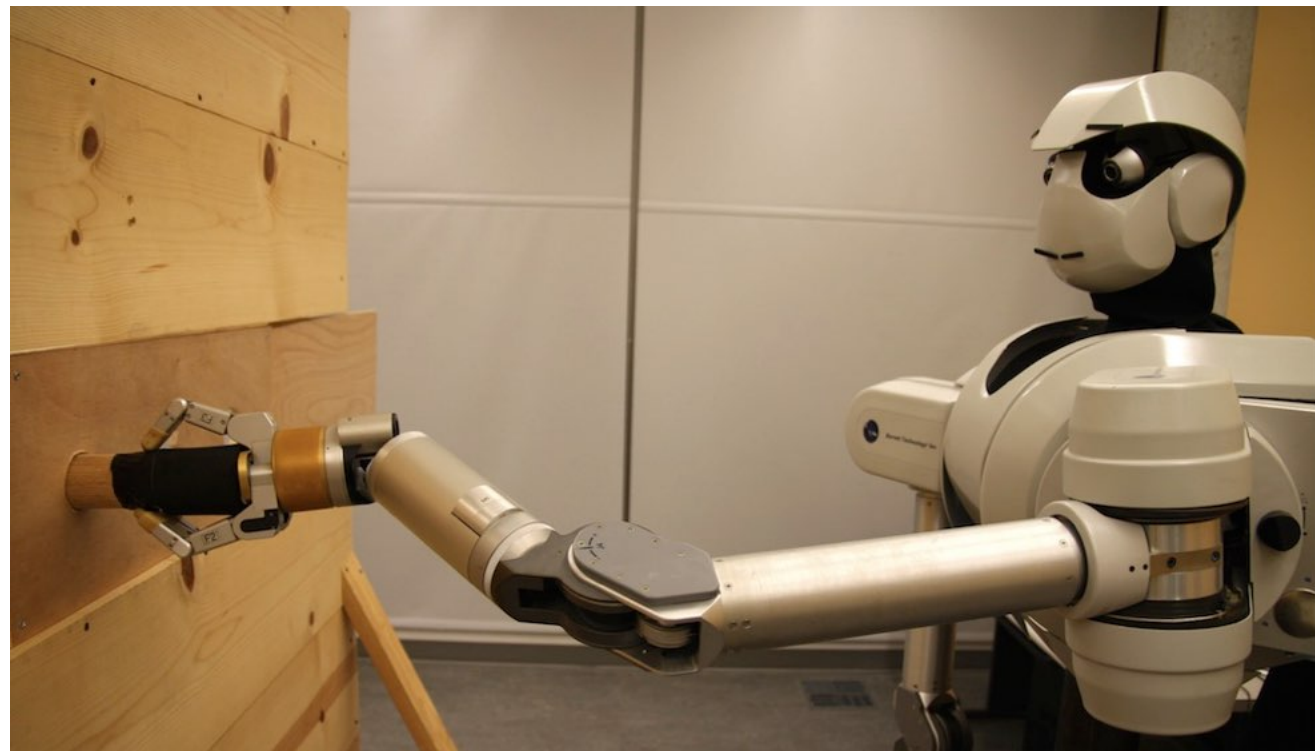
Example: Driving

- States: Road traffic, weather, time of day
- Actions: steering wheel, break
- Rewards:
 - +1 reaching goal not over-tired
 - -1: honking from surrounding drivers
 - -100: collision



Example: Peg in Hole Insertion

- States: Joint configurations ?
- Actions: Torques on joints
- Rewards: Penalize jerky motions, reaching target pose



Example: Peg in Hole Insertion

- States: Joint configurations ?
- Actions: Torques on joints
- Rewards: Penalize jerky motions, reaching target pose

A Framework for Robot Manipulation:
Skill Formalism, Meta Learning and Adaptive Control

Lars Johannessmeier, Malkin Gerchow and Sami Haddadin

Institute of Automatic Control
Gottfried Wilhelm Leibniz Universität Hannover

Dynamics a.k.a. the World Model

- Encodes the results of the actions of the agent.
- How the states and rewards change given the actions of the agent:

$$p(s', r | s, a) = \mathbb{P}\{S_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a\}$$

- Transition function or next step function:

$$T(s' | s, a) = p(s' | s, a) = \mathbb{P}\{S_t = s' | S_{t-1} = s, A_{t-1} = a\} = \sum_{r \in \mathbb{R}} p(s', r | s, a)$$

Dynamics a.k.a. the World Model

“the idea that we predict the consequences of our motor commands has emerged as an important theoretical concept in all aspects of sensorimotor control”

Prediction Precedes Control in Motor Learning

J. Randall Flanagan,^{1*} Philipp Vetter,²
Roland S. Johansson,³ and Daniel M. Wolpert²

Procedures for details). Figure 1 shows, for a single subject, the hand path (top trace) and the grip (middle)

Predicting the Consequences of Our Own Actions: The Role of Sensorimotor Context Estimation

Sarah J. Blakemore, Susan J. Goodbody, and Daniel M. Wolpert

Sobell Department of Neurophysiology, Institute of Neurology, University College London, London WC1N 3BG,

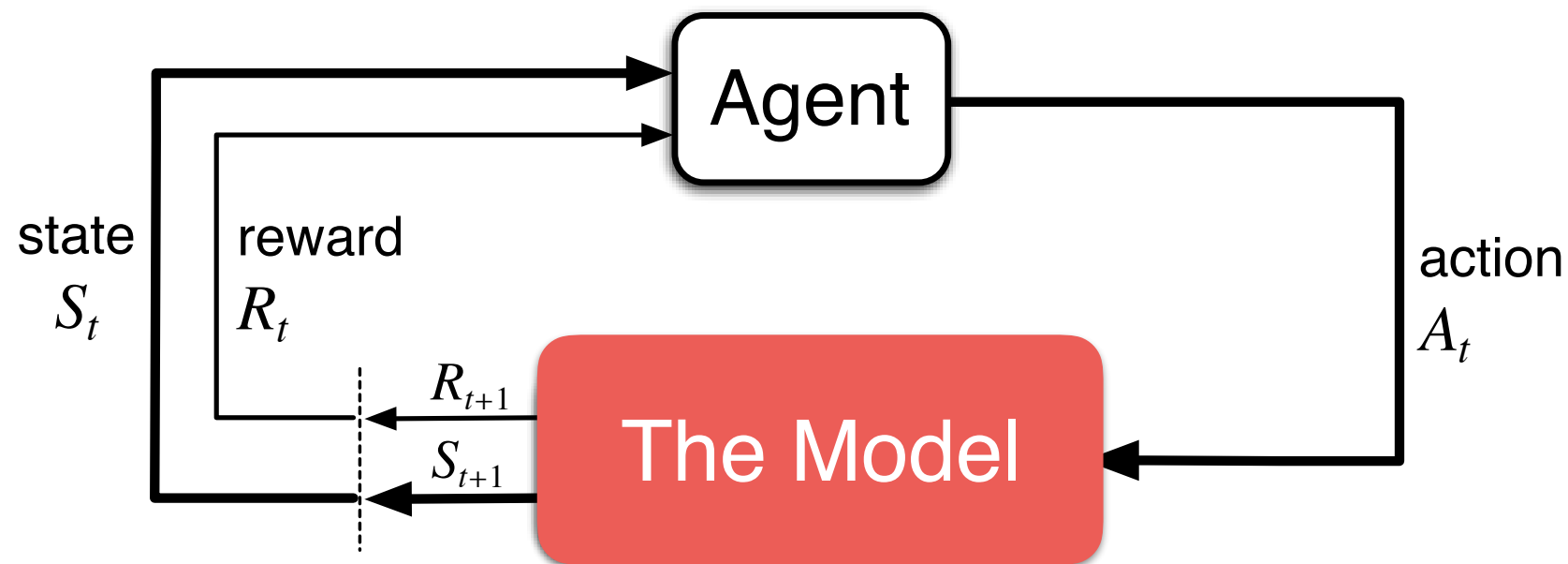
Predictive coding in the visual cortex: a functional interpretation of some extra-classical receptive-field effects

Rajesh P. N. Rao¹ and Dana H. Ballard²

Planning

Planning: unrolling (querying) a model forward in time and selecting the best action sequence that satisfies a specific goal

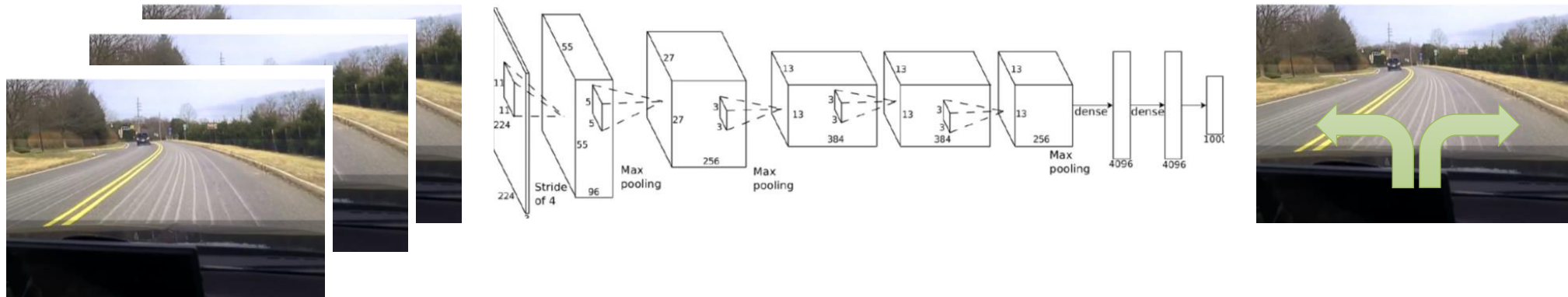
Plan: a sequence of actions



Is planning learning or not?

Why **deep** reinforcement learning?

Because the policy, the model and the value functions (expected returns) will often be represented by some form of a deep neural network.



Limitations of Reinforcement Learning

- Can we think of goal directed behavior learning problems that cannot be modeled or are not meaningful using the trial-and-error reinforcement learning framework?
- The agent should have the chance to try (and fail) enough times
- This is impossible if episode takes too long, e.g., reward=“obtain a great Ph.D.”
- This is impossible when safety is a concern: we can't learn to drive via reinforcement learning in the real world, failure cannot be tolerated

Q: what other forms of supervision humans use to learn to act in the world?

Other forms of supervision for learning behaviors?

1. Learning from rewards
2. Learning from demonstrations
3. Learning from specifications of optimal behavior

Behavior: High Jump

scissors



Fosbury flop



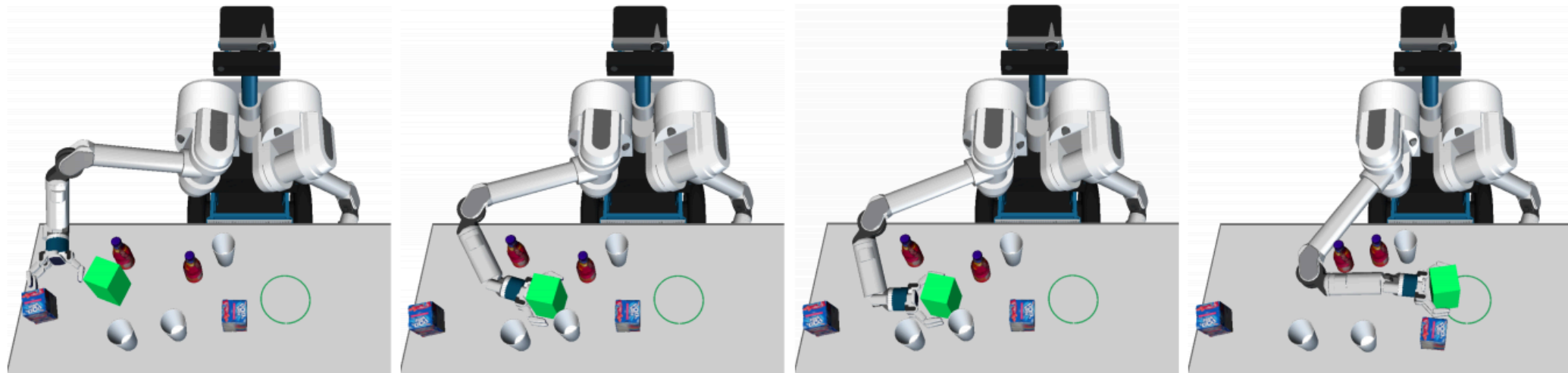
- Learning from **rewards**
 - Reward: jump as high as possible: It took years for athletes to find the right behavior to achieve this
- Learning from **demonstrations**
 - It was way easier for athletes to perfection the jump, once someone showed the right general trajectory
- Learning from **specifications of optimal behavior**
 - For novices, it is much easier to replicate a behavior if additional guidance is provided in natural language: where to place the foot, how to time yourself, etc. .

Overview for today

- Goal of the course / why it is important
- What is reinforcement learning
- What is representation learning (and how it helps reinforcement learning and behavior learning in general)
- Reinforcement learning versus supervised learning
- AI's paradox: what is hard and what is easy in behavior learning

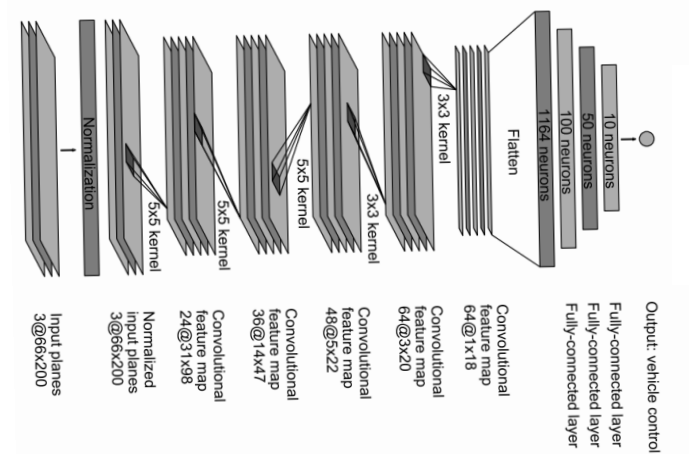
State estimation - Two extremes

- Assuming we know everything about the world (object locations, 3D shapes, physical properties) and world dynamics. Use planners to search for the action sequence to achieve a desired goal.



State estimation - Two extremes

- Assuming we know everything about the world (object locations, 3D shapes, physical properties). Use planners to search for the action sequence to achieve a desired goal.
- Assuming we know nothing about the world. Learn to map pixels directly to actions while optimizing for your end task, i.e., not crashing and obeying the traffic signs, or, imitating human demonstrations.



In practice: A lot of domain knowledge for going from observations to states

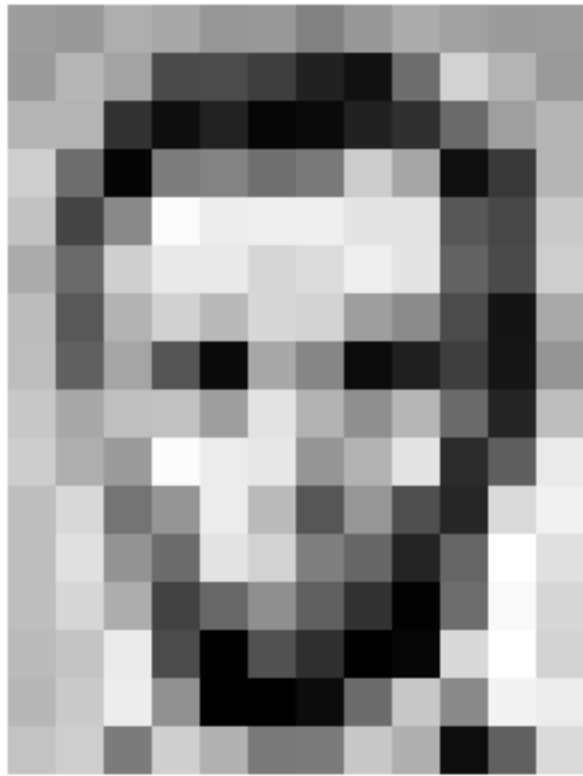


- Q: should the location of the trees and their fruits be part of the state for driving?
- Q: should the location of the trees and their fruits be part of the state for apple picking?

Representation learning helps learning to act

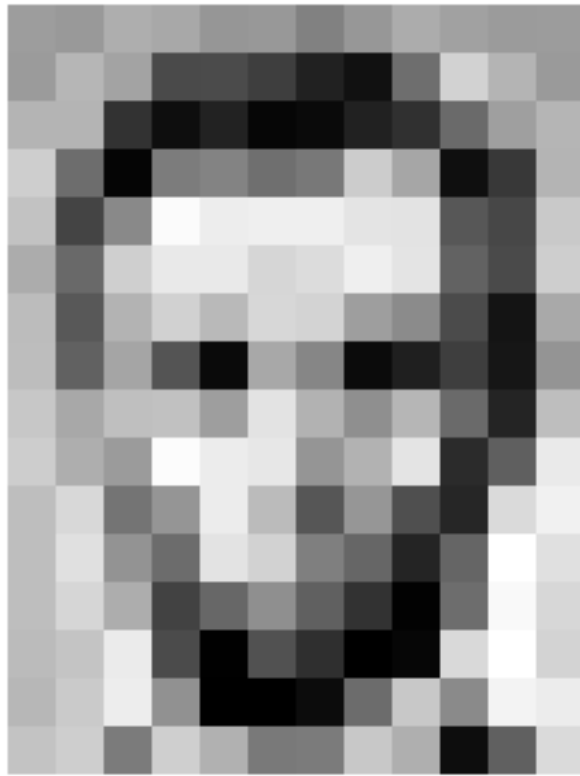
- Representation learning: mapping raw observations to features and structures from which the mapping to actions or to semantic labels is easier to infer.

Representation learning



- Remember what the computer sees

Representation learning



157	153	174	168	150	152	129	151	172	161	155	156
155	182	163	74	75	62	33	17	110	210	180	154
180	180	50	14	34	6	10	33	48	106	159	181
206	109	5	124	131	111	120	204	166	15	56	180
194	68	137	251	237	239	239	228	227	87	71	201
172	106	207	233	233	214	220	239	228	98	74	206
188	88	179	209	185	215	211	158	139	75	20	169
189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	193	158	227	178	143	182	105	36	190
205	174	155	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	108	227	210	127	102	36	101	255	224
190	214	173	66	103	143	95	50	2	109	249	215
187	196	235	75	1	81	47	0	6	217	255	211
183	202	237	145	0	0	12	108	200	138	243	236
195	206	123	207	177	121	123	200	175	13	96	218

- Remember what the computer sees

Representation learning

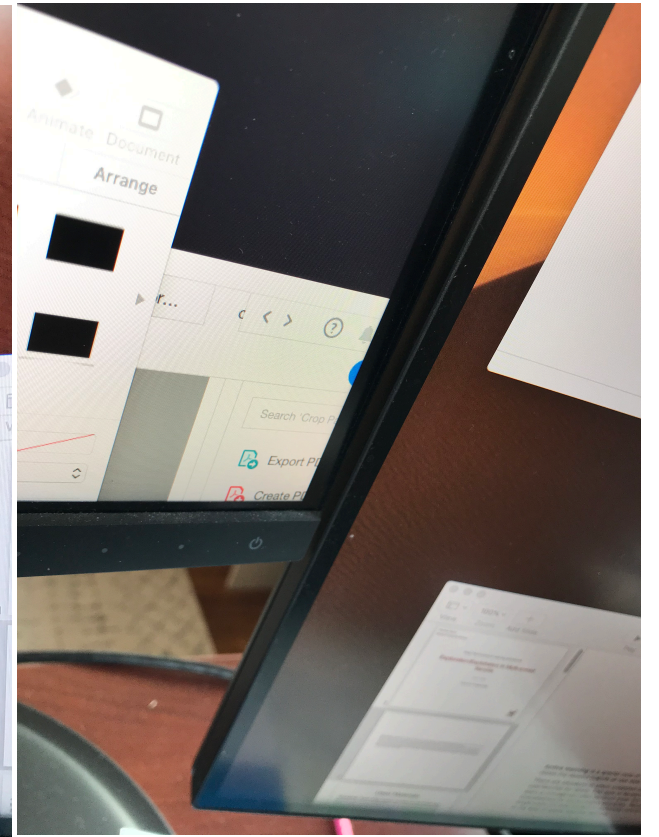
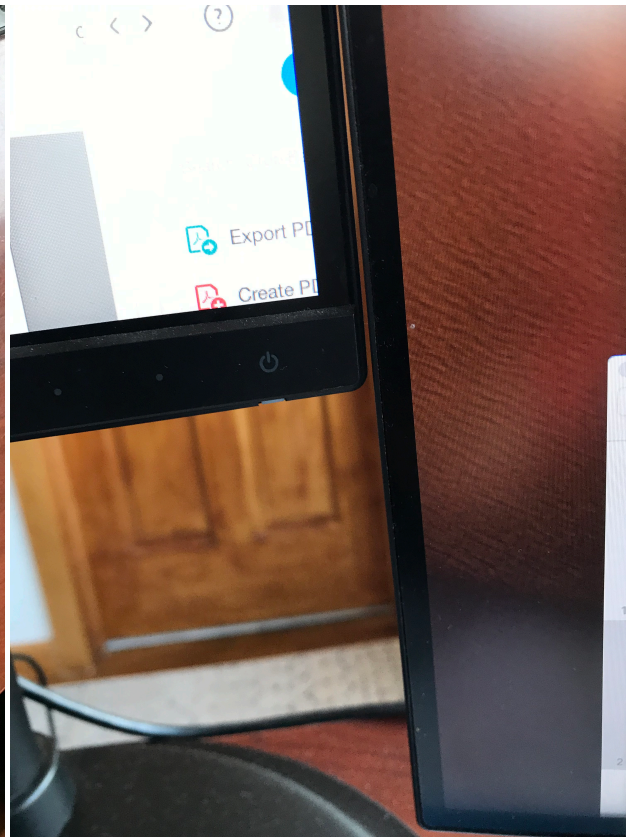
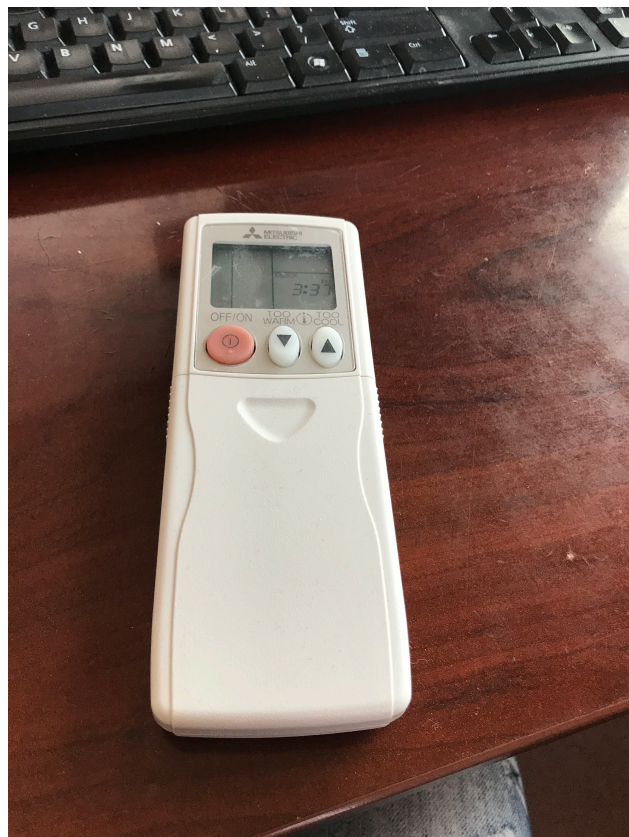


157	153	174	168	150	152	129	151	172	161	155	156
155	182	163	74	75	62	33	17	110	210	180	154
180	180	50	14	34	6	10	33	48	106	159	181
206	109	5	124	131	111	120	204	166	15	56	180
194	68	137	251	237	239	239	228	227	87	71	201
172	106	207	233	233	214	220	239	228	98	74	206
188	88	179	209	185	215	211	158	139	75	20	169
189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	193	158	227	178	143	182	106	36	190
205	174	155	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	86	150	79	38	218	241
190	224	147	108	227	210	127	102	36	101	255	224
190	214	173	66	103	143	96	50	2	109	249	215
187	196	235	75	1	81	47	0	6	217	255	211
183	202	237	145	0	0	12	108	200	138	243	236
195	206	123	207	177	121	123	200	175	13	96	218

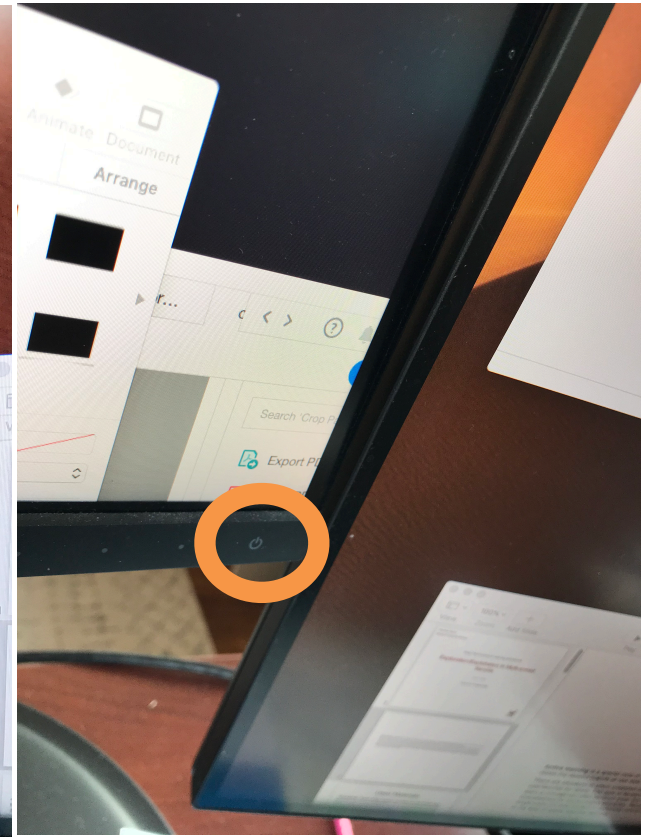
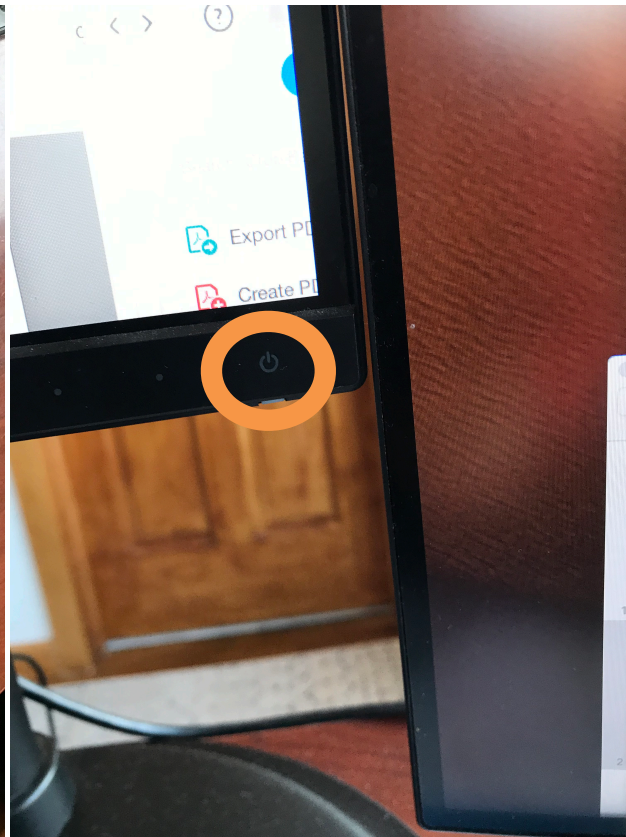
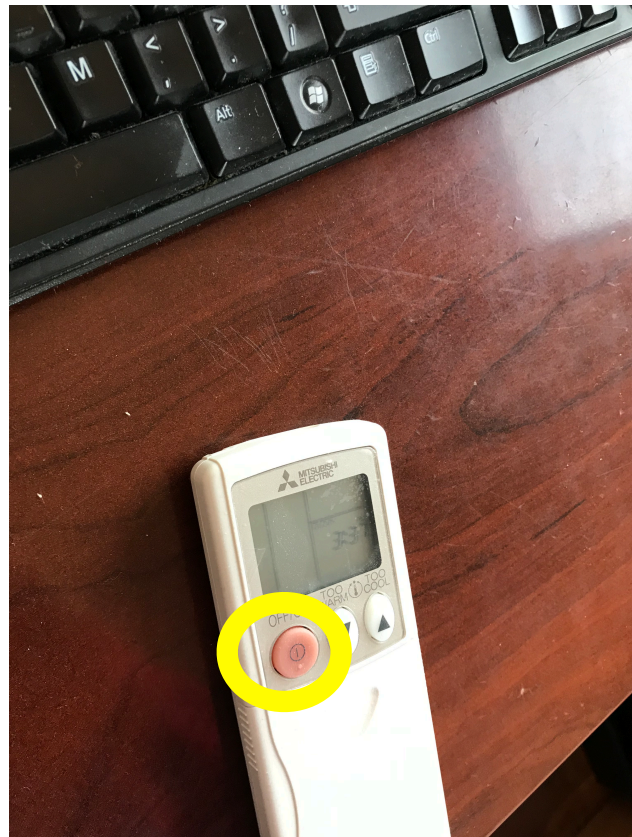
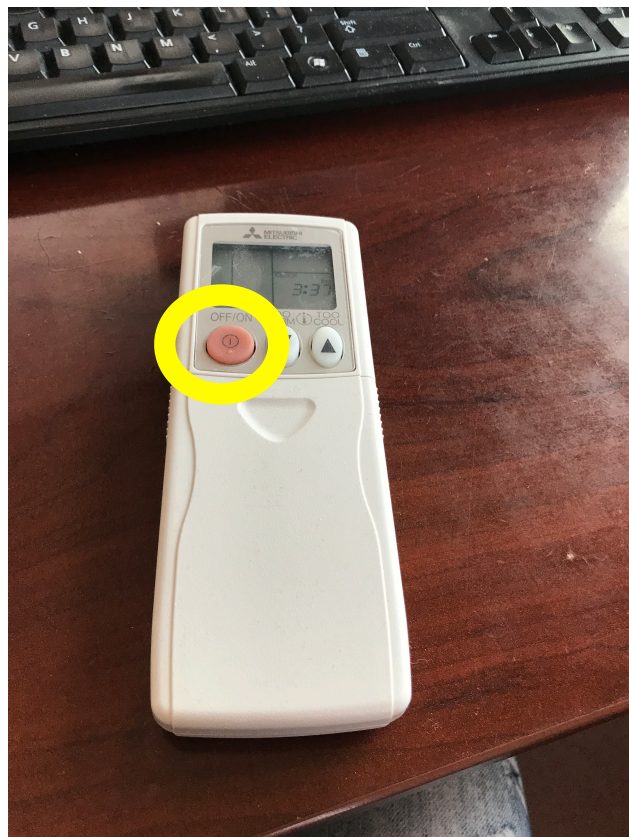
157	153	174	168	150	152	129	151	172	161	155	156
155	182	163	74	75	62	33	17	110	210	180	154
180	180	50	14	34	6	10	33	48	106	159	181
206	109	5	124	131	111	120	204	166	15	56	180
194	68	137	251	237	239	239	228	227	87	71	201
172	106	207	233	233	214	220	239	228	98	74	206
188	88	179	209	185	215	211	158	139	75	20	169
189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	193	158	227	178	143	182	106	36	190
205	174	155	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	86	150	79	38	218	241
190	224	147	108	227	210	127	102	36	101	255	224
190	214	173	66	103	143	96	50	2	109	249	215
187	196	235	75	1	81	47	0	6	217	255	211
183	202	237	145	0	0	12	108	200	138	243	236
195	206	123	207	177	121	123	200	175	13	96	218

- Remember what the computer sees

Representation learning

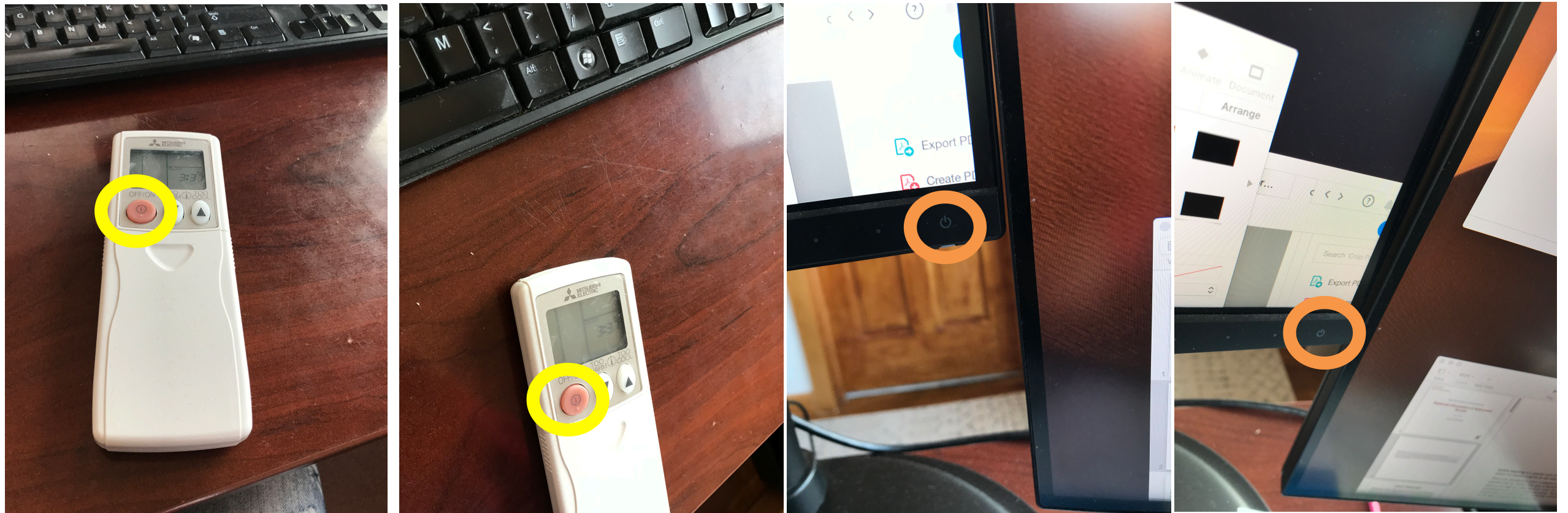


Representation learning



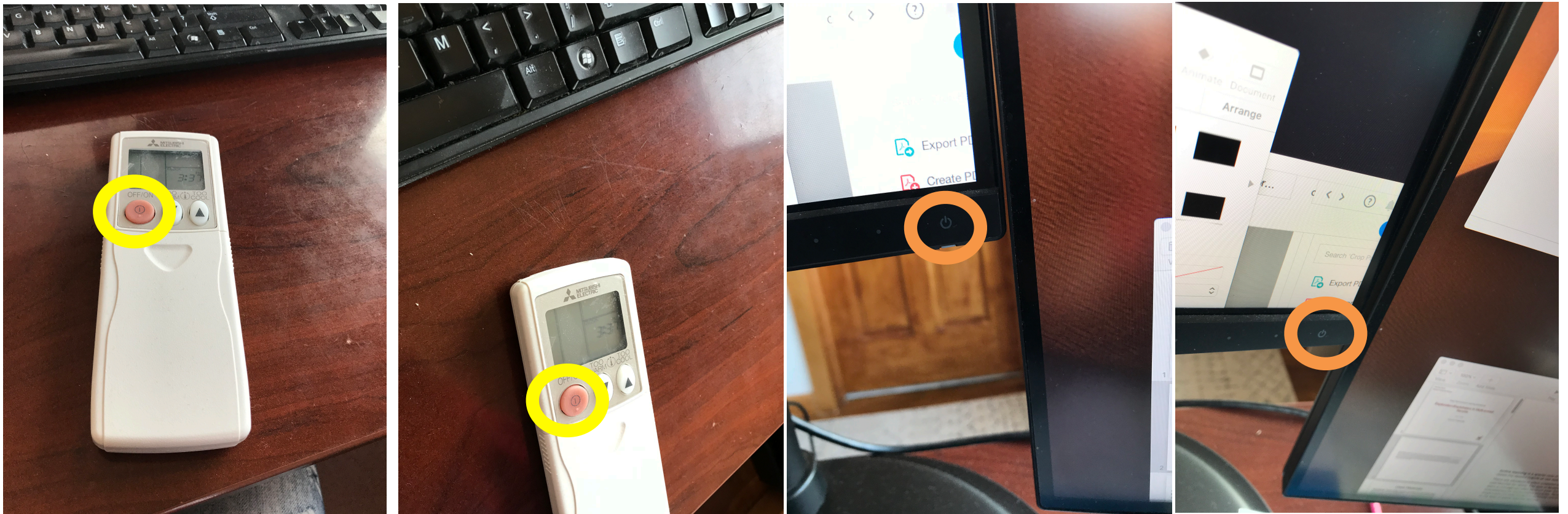
(Visual) Representation learning helps learning to act

- Despite these images have very different pixel values, actions required to achieve the goal of switching on the device are similar.
- Visual perception is instrumental to learning to act, in transforming raw pixels to action-relevant feature vectors and structures.



(Visual) Representation learning helps learning to act

- Having pre-trained our visual representations with auxiliary tasks is likely to dramatically decrease the number of interactions with the environment we need to learn to press buttons.



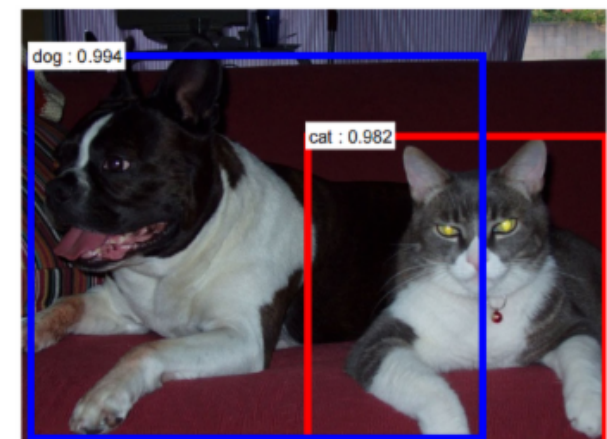
- Q: What are reasonable auxiliary tasks?
 - Supervised: object detection, image classification, pixel labelling.
 - Unsupervised: open research problem

Overview for today

- Goal of the course / why it is important
- What is reinforcement learning
- What is representation learning (and how it helps reinforcement learning and behavior learning in general)
- **Reinforcement learning versus supervised learning**
- AI's paradox: what is hard and what is easy in behavior learning

Reinforcement learning Versus supervised learning

- RL is a form of active learning:
 - the agent gets the chance **to collect her own data** by acting in the world, querying humans, and so on.
 - the data changes over time, it depends on the policy of the agent.
 - To query the environment effectively, the agent needs to keep track of its **uncertainty**: what she knows and what she does not, and thus needs to explore next.
- Supervised learning is a form of passive learning:
 - the data does not depend on the agent in anyway, it is provided by external labellers.
 - the data is static throughout learning.



Reinforcement learning Versus supervised learning

- In RL, we often cannot use gradient-based optimization:
 - e.g., when the agent does not know neither the world model to unroll nor the reward function to maximize.
- In supervised learning, we usually can use gradient-based optimization:
 - E.g., we consider a parametric form for our regressor or classifier and optimize it via stochastic gradient descent (SGD).

Reinforcement learning Versus supervised learning

- RL can be time consuming. Actions take time to carry out in the real world, i.e., each interaction has a non-negligible cost. Our goal is the agent to minimize the amount of interactions with the environment while succeeding in the task.

Reinforcement learning Versus supervised learning

- RL can be time consuming. Actions take time to carry out in the real world, i.e., each interaction has a non-negligible cost. Our goal is the agent to minimize the amount of interactions with the environment while succeeding in the task.
- We can use **simulated experience** and tackle the SIM2REAL (simulation to reality) transfer.

Reinforcement learning Versus supervised learning

- RL can be time consuming. Actions take time to carry out in the real world, i.e., each interaction has a non-negligible cost. Our goal is the agent to minimize the amount of interactions with the environment while succeeding in the task.
- We can use **simulated experience** and tackle the SIM2REAL (simulation to reality) transfer.

MuJoCo physics

Roboti LLC

www.mujoco.org

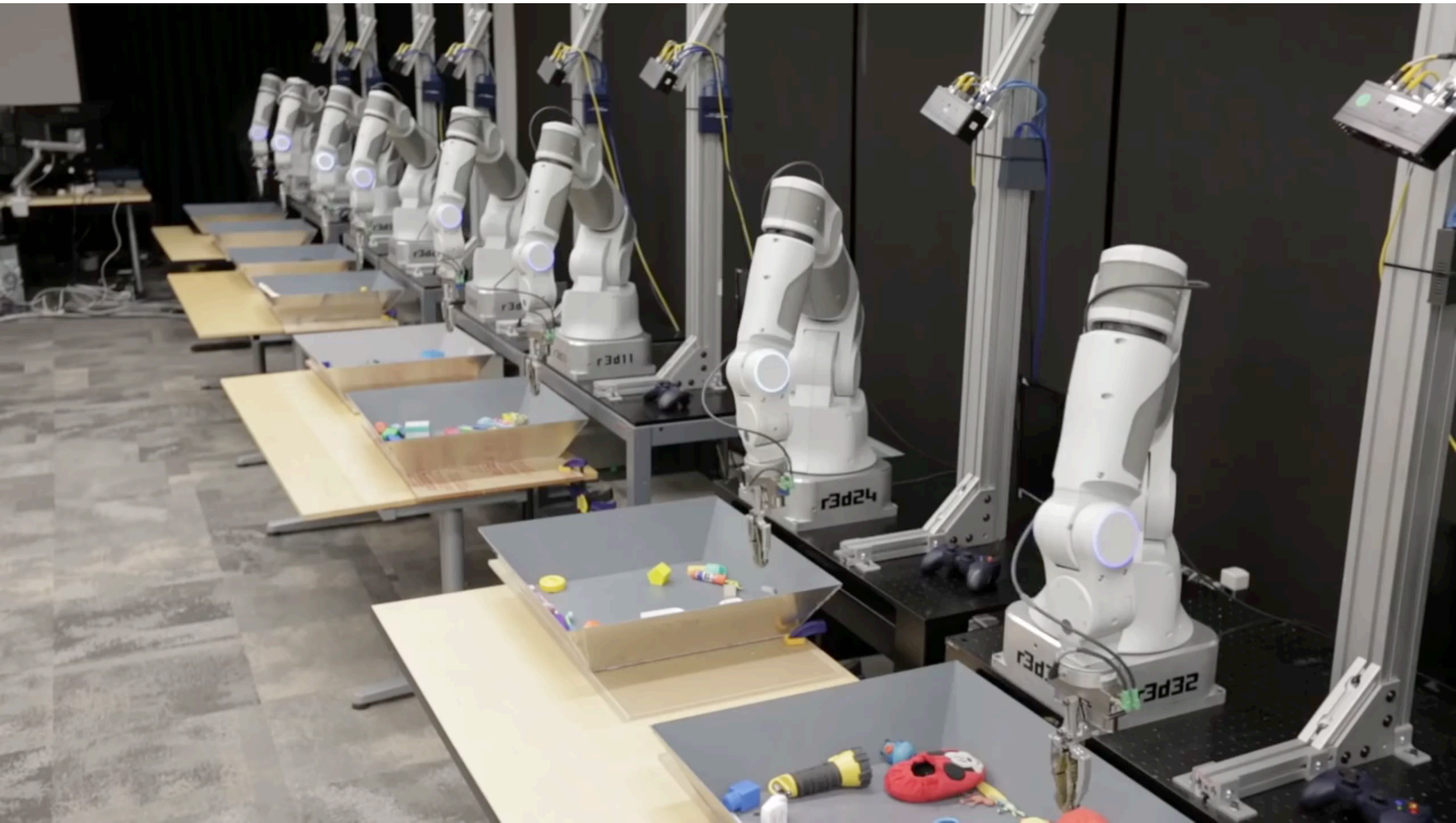
Reinforcement learning Versus supervised learning

- RL can be time consuming. Actions take time to carry out in the real world, i.e., each query has a non-negligible cost. Our goal is the agent to minimize the amount of interactions with the environment while succeeding in the task.
- We can use **simulated experience** and tackle the SIM2REAL (simulation to reality) transfer.
- We can have robots working 24/7

Reinforcement learning Versus supervised learning

- RL can be time consuming. Actions take time to carry out in the real world, i.e., each query has a non-negligible cost. Our goal is the agent to minimize the amount of interactions with the environment while succeeding in the task.
- We can use **simulated experience** and tackle the SIM2REAL (simulation to reality) transfer.
- We can have robots working 24/7
- We can buy many robots

Google's Robot Farm



True or False

Given a dataset of state, action, reward sequences

$(s_1, a_1, r_1, s_2, a_2, r_2, s_3, a_3, r_3, \dots)$:

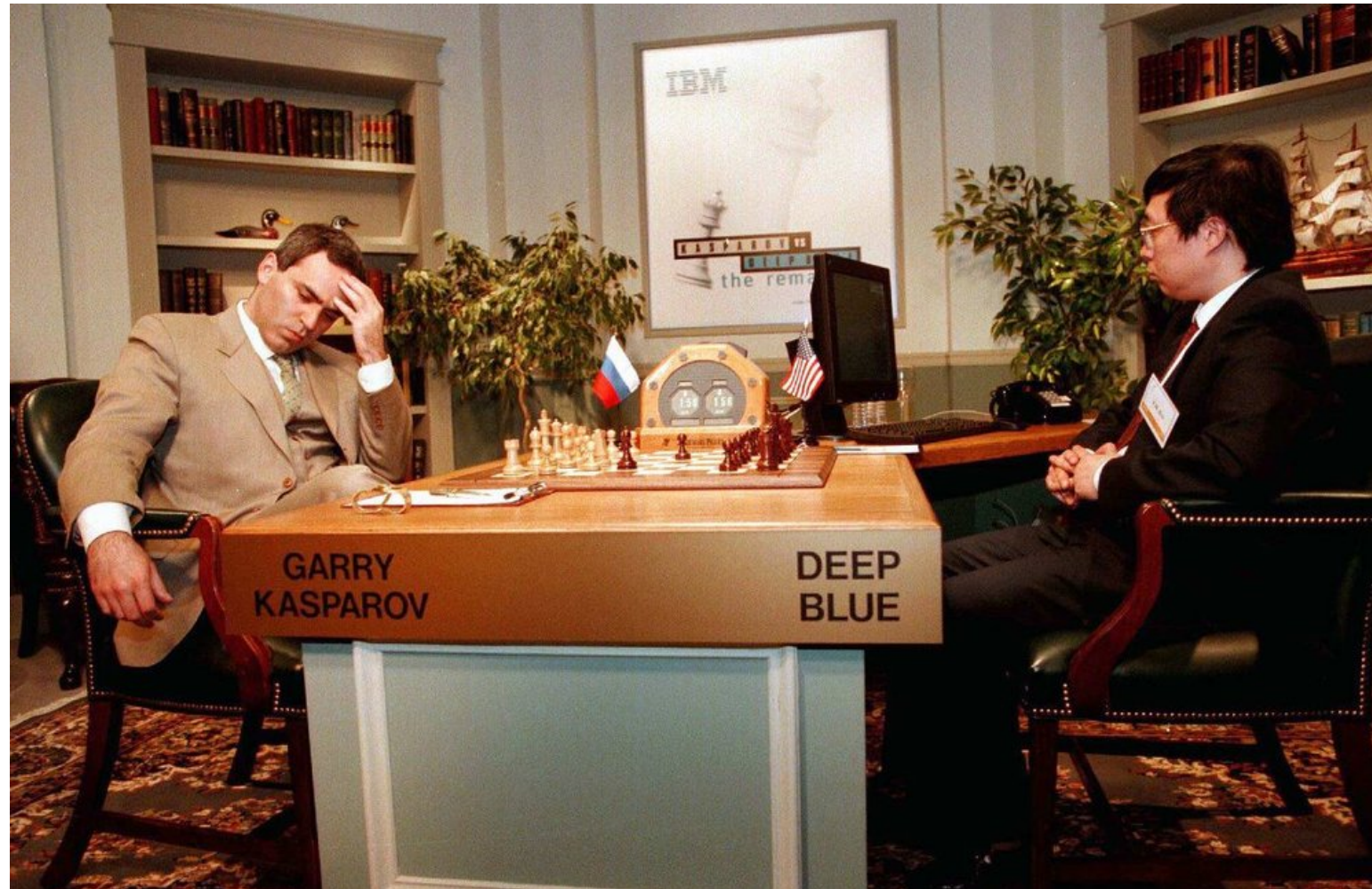
- learning a dynamics model, i.e., mapping of state and actions to next state, is a reinforcement learning problem.
- learning a dynamics model, i.e., mapping of state and actions to next state, is a supervised learning problem.
- for learning a dynamics model, i.e., mapping of state and actions to next state, I can use gradient information.

Given a dataset of state, action, reward sequences

$(s_1, a_1, r_1, s_2, a_2, r_2, s_3, a_3, r_3, \dots)$ from an expert interacting with the environment:

- for learning the expert policy, i.e., mapping of states to expert actions, is a supervised learning problem.
- for learning the expert policy, i.e., mapping of states to expert actions, I do not need to use the rewards.

Deep Blue



A big search with heuristics: manual development of a board evaluation function.

Backgammon



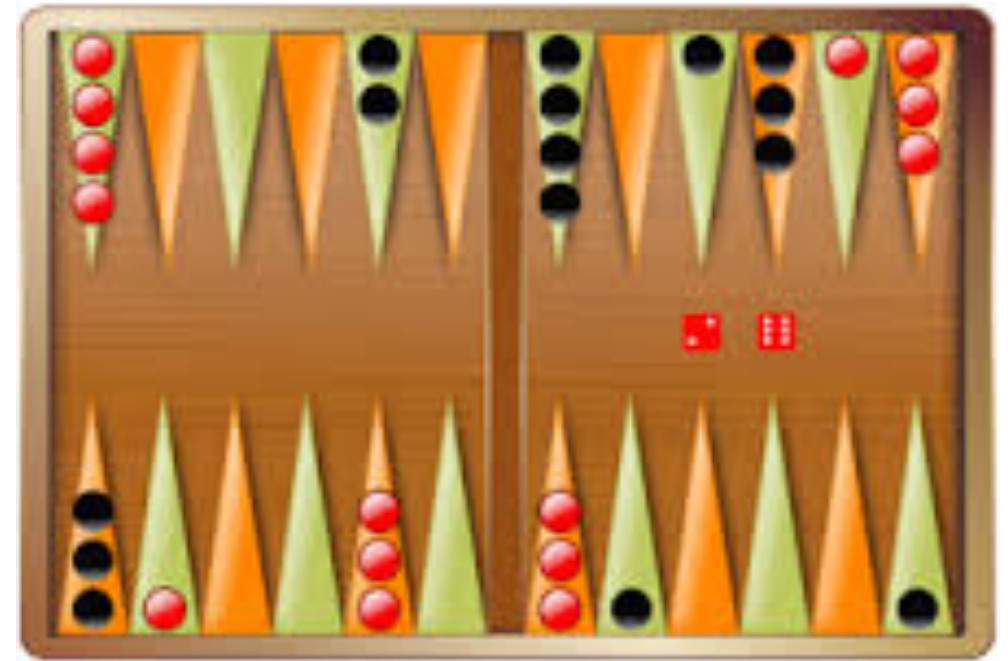
Backgammon



High branching factor due to dice roll prohibits brute force deep searches such as in chess

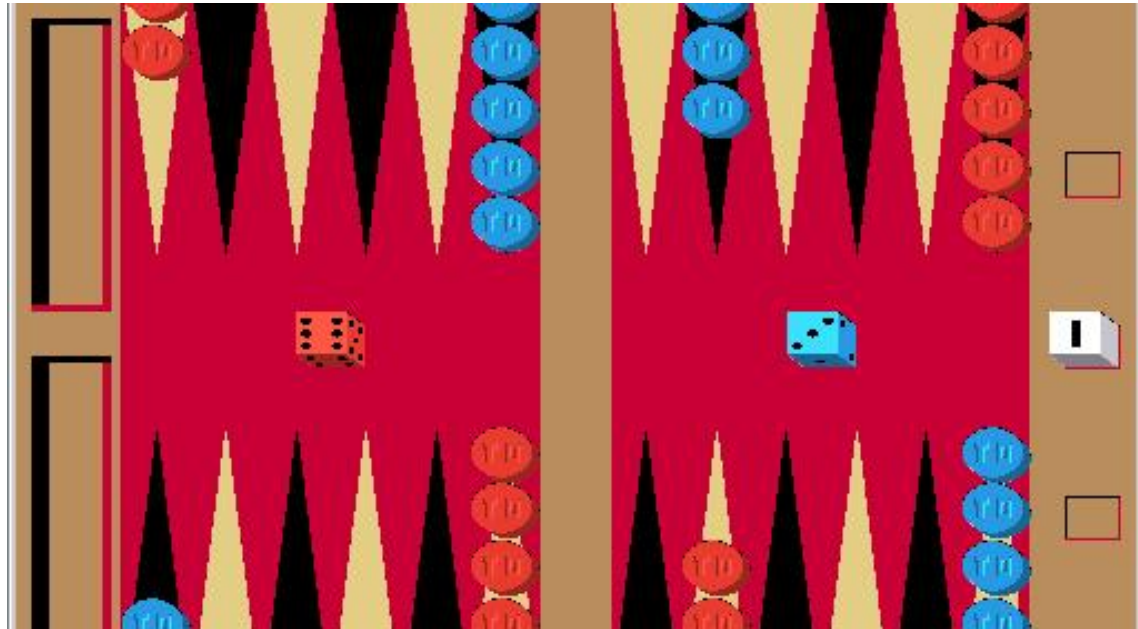


Neuro-Gammon



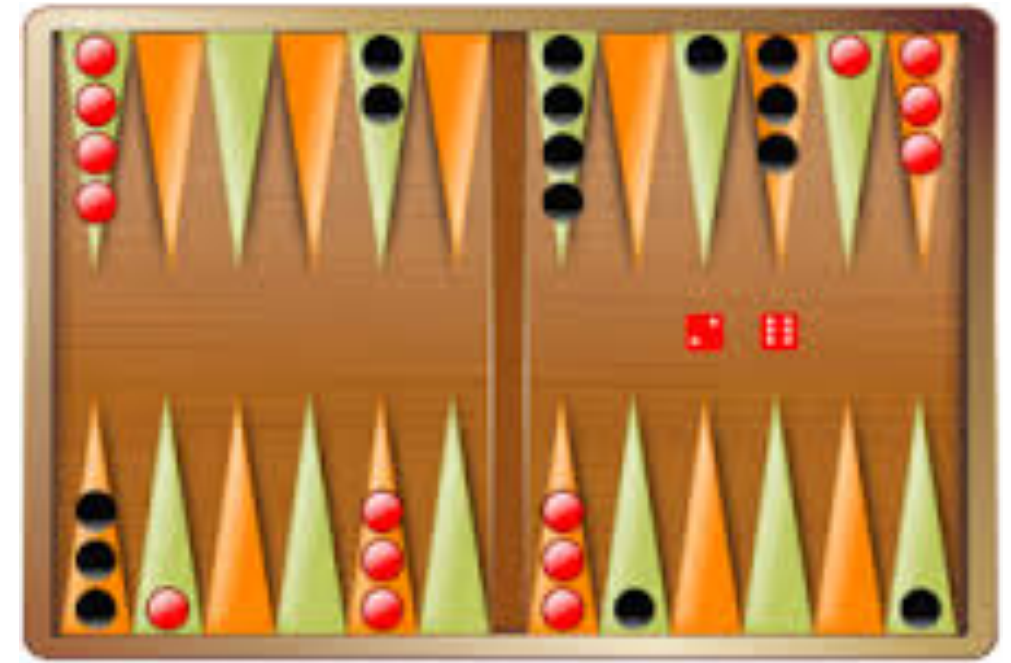
- Developed by Gerald Tesauro in 1989 in IBM's research center
- Trained to mimic expert demonstrations using supervised learning
- Achieved intermediate-level human player

TD-Gammon



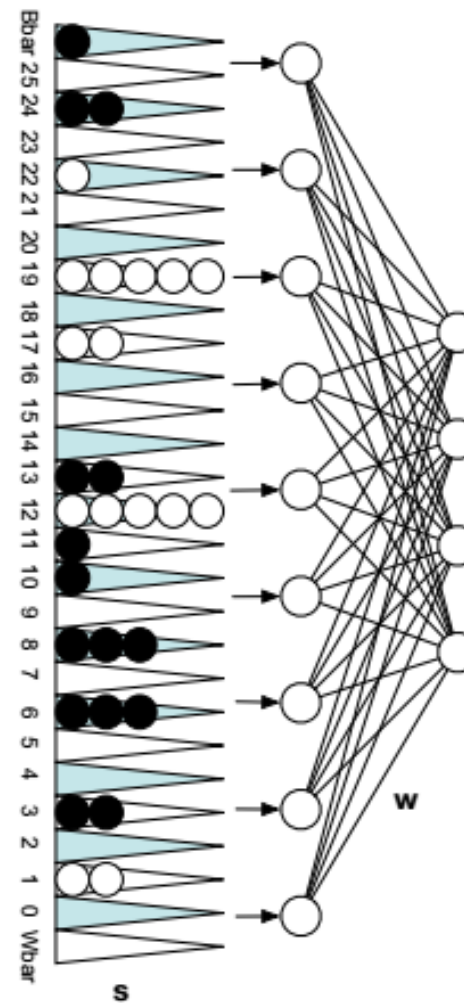
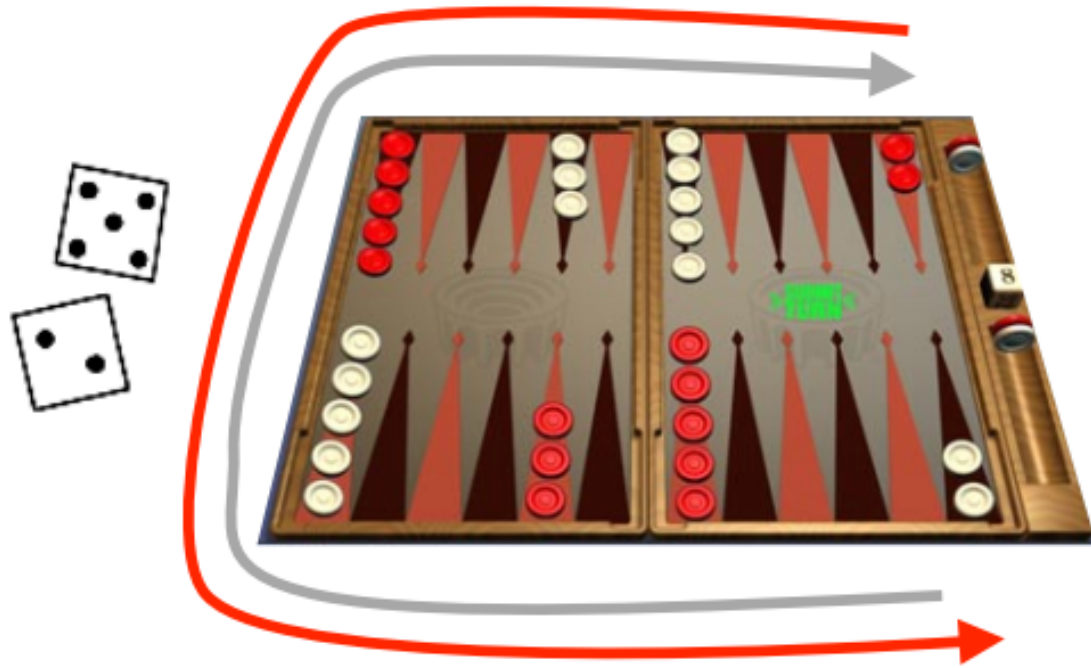
- Developed by Gerald Tesauro in 1992 in IBM's research center
- A neural network that trains itself to be an evaluation function by playing against itself starting from random weights
- Achieved performance close to top human players of its time

Neuro-Gammon



- Developed by Gerald Tesauro in 1989 in IBM's research center
- Trained to mimic expert demonstrations using supervised learning
- Achieved intermediate-level human player

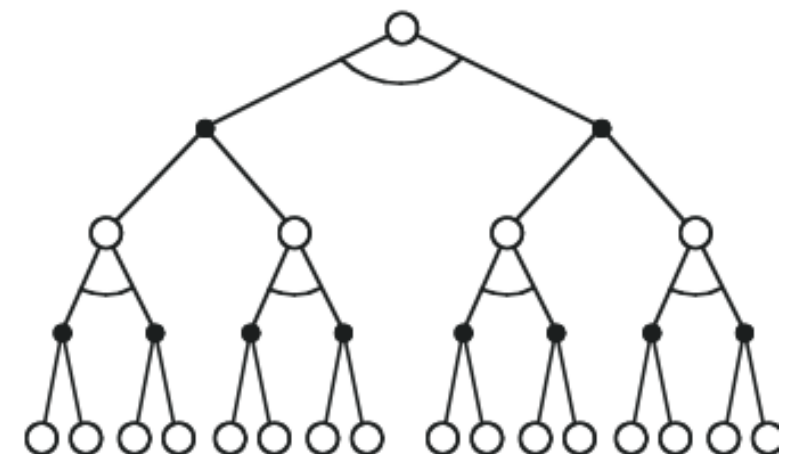
Evaluation function



A neural net with only 80 hidden units..

estimated state value
(\approx prob of winning)

Action selection
by a shallow search



Overview for today

- Goal of the course / why it is important
- What is reinforcement learning
- What is representation learning (and how it helps reinforcement learning and behavior learning in general)
- Reinforcement learning versus supervised learning
- AI's paradox: what is hard and what is easy in behavior learning

GO



AlphaGoZero the program that beat the world champions with only RL



- Monte Carlo Tree Search with neural nets
- self play

Go Versus the real world



Beating the world champion is easier than moving the Go stones.

The difficulty of motor control

What to move where



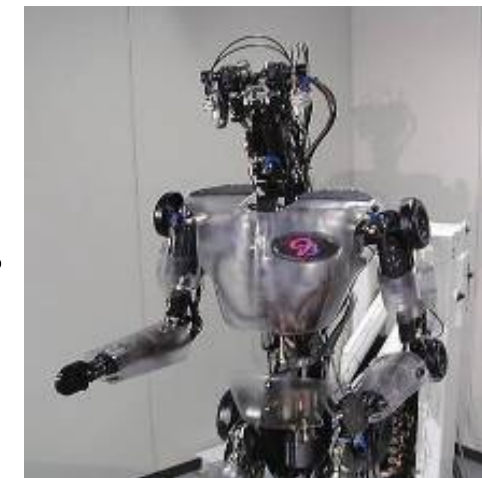
vs.



Moving



vs.



From Dan Wolpert

Reinforcement learning in the real world

How the world of Alpha Go is different than the real world?

1. Known environment (known entities and dynamics) Vs Unknown environment (unknown entities and dynamics).
2. Need for behaviors to transfer across environmental variations since the real world is very diverse
3. Discrete Vs Continuous actions
4. One goal Vs many goals
5. Rewards are provided automatically by an oracle environment VS rewards need themselves to be detected
6. Interactions take time: we really need intelligent exploration

Alpha Go Versus the real world

How the world of Alpha Go is different than the real world?

1. Known environment (known entities and dynamics) Vs Unknown environment (unknown entities and dynamics).
2. Need for behaviors to transfer across environmental variations since the real world is very diverse

Alpha Go Versus the real world

How the world of Alpha Go is different than the real world?

1. Known environment (known entities and dynamics) Vs Unknown environment (unknown entities and dynamics).
2. Need for behaviors to transfer across environmental variations since the real world is very diverse

State estimation: To be able to act you need first to be able to see, detect the objects that you interact with, detect whether you achieved your goal

AI's paradox



Hans Moravec

"it is comparatively easy to make computers exhibit adult level performance on intelligence tests or playing checkers, and difficult or impossible to give them the skills of a one-year-old when it comes to perception and mobility"

AI's paradox



Marvin Minsky

"we're more aware of simple processes that don't work well than of complex ones that work flawlessly"

Evolutionary explanation



Hans Moravec

“We should expect the difficulty of reverse-engineering any human skill to be roughly proportional to the amount of time that skill has been evolving in animals.

The oldest human skills are largely unconscious and so appear to us to be effortless.

Therefore, we should expect skills that appear effortless to be difficult to reverse-engineer, but skills that require effort may not necessarily be difficult to engineer at all.”

AI's paradox

Intelligence was "best characterized as the things that highly educated scientists found challenging", such as chess, **symbolic integration**, proving **mathematical theorems** and solving complicated word algebra problems.



Rodney Brooks

AI's paradox

Intelligence was "best characterized as the things that highly educated scientists found challenging", such as chess, **symbolic integration**, proving **mathematical theorems** and solving complicated word algebra problems.

"The things that children of four or five years could do effortlessly, such as visually distinguishing between a coffee cup and a chair, or walking around on two legs, or finding their way from their bedroom to the living room were not thought of as activities requiring intelligence."



Rodney Brooks

AI's paradox

Intelligence was "best characterized as the things that highly educated scientists found challenging", such as chess, **symbolic integration**, proving **mathematical theorems** and solving complicated word algebra problems.



Rodney Brooks

"The things that children of four or five years could do effortlessly

No cognition. Just sensing and action

coffee cup and a chair, or walking around on two legs, or finding their way from their bedroom to the living room were not thought of as activities requiring intelligence."

Learning from Babies

- Be multi-modal
- Be incremental
- Be physical
- Explore
- Be social
- Learn a language



Take-aways

- Forms of supervision for learning to act: mapping observations to actions for a specific goal
- The reinforcement learning problem, terminology, basic ingredients
- RL vs SL
- Learning to search using evaluation functions
- AI paradox: is hard to learn the abilities of a 2 year old, and easy to learn to beat GO champions, solve theorems and so on: a big search at a kind of small (compared to the real world) state space at the end of the day.