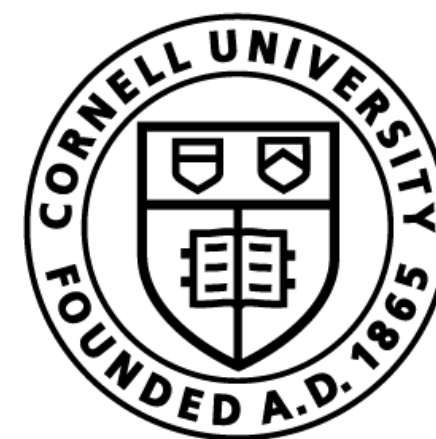


# Reinforcement Learning: From Games to Robotics

Sanjiban Choudhury



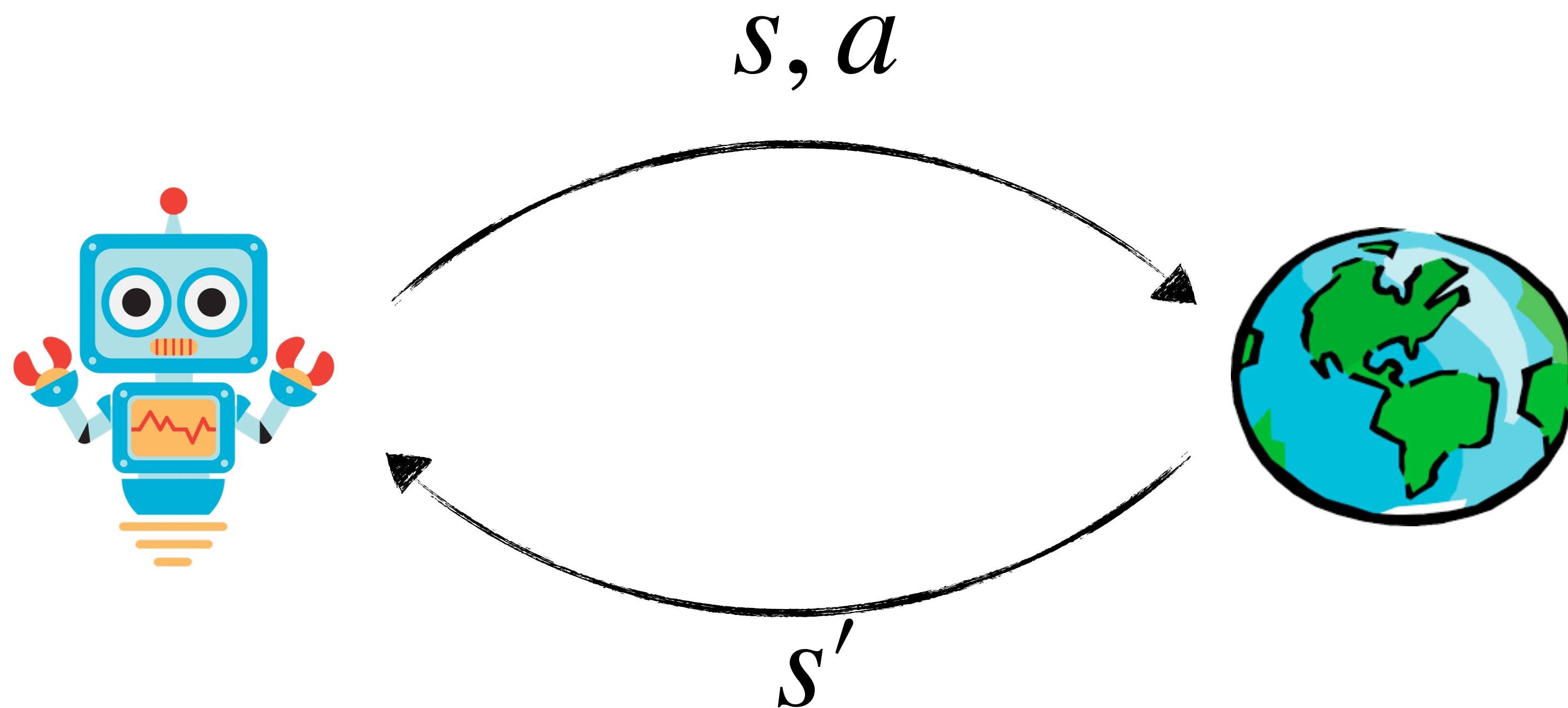
Cornell Bowers CIS  
**Computer Science**

The story thus far ...

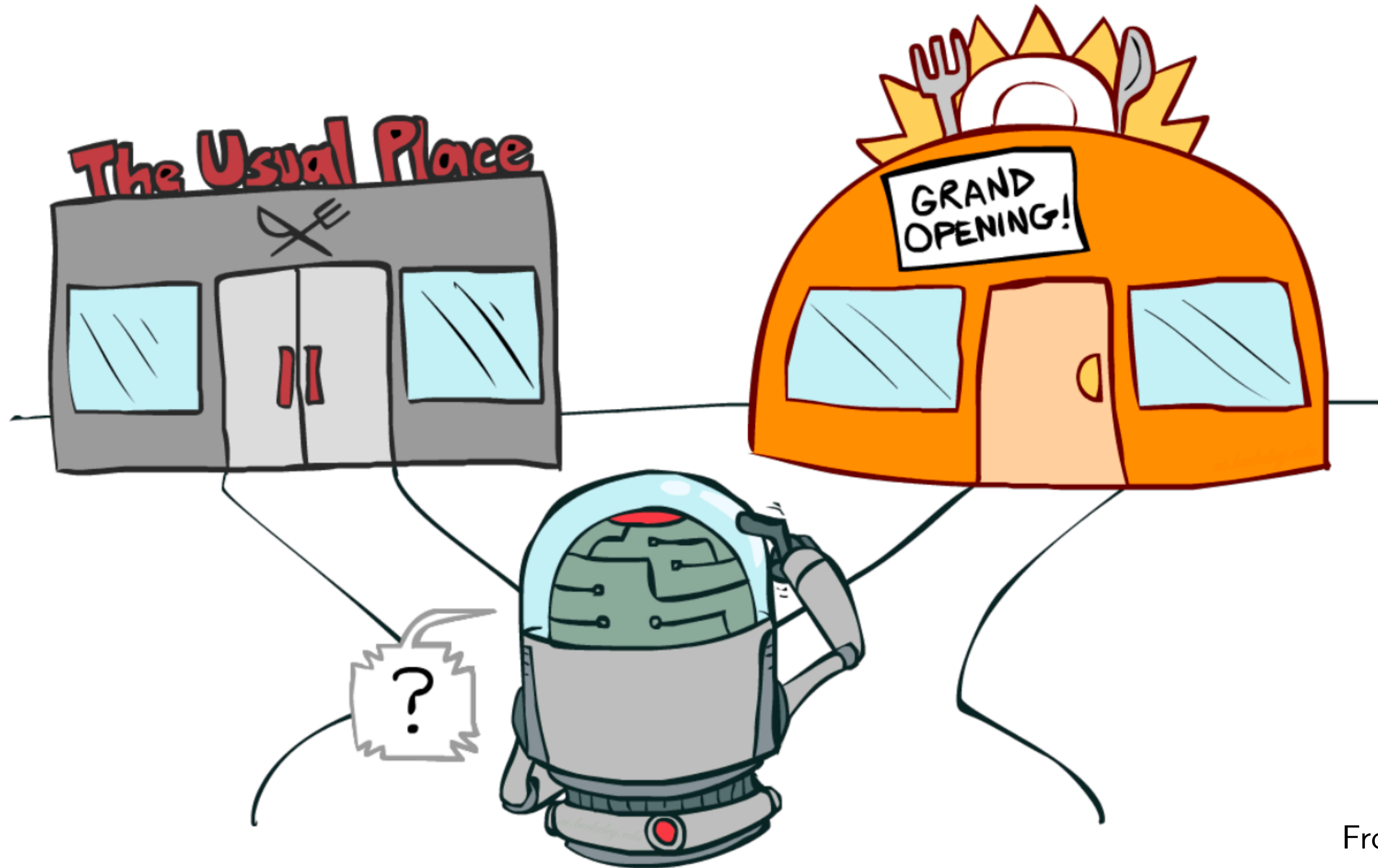


But what if the dynamics are unknown?

$\langle S, A, C, \mathcal{F} \rangle$



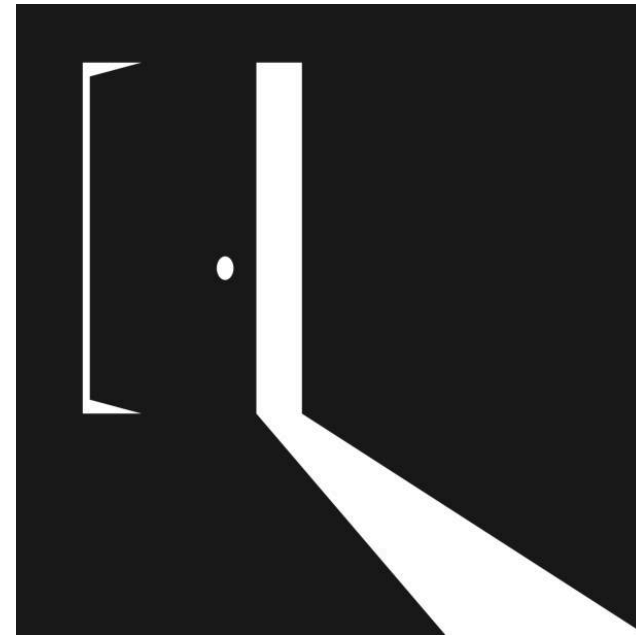
# Exploration vs Exploitation



From Dan Klein

# Doors

$a^1$

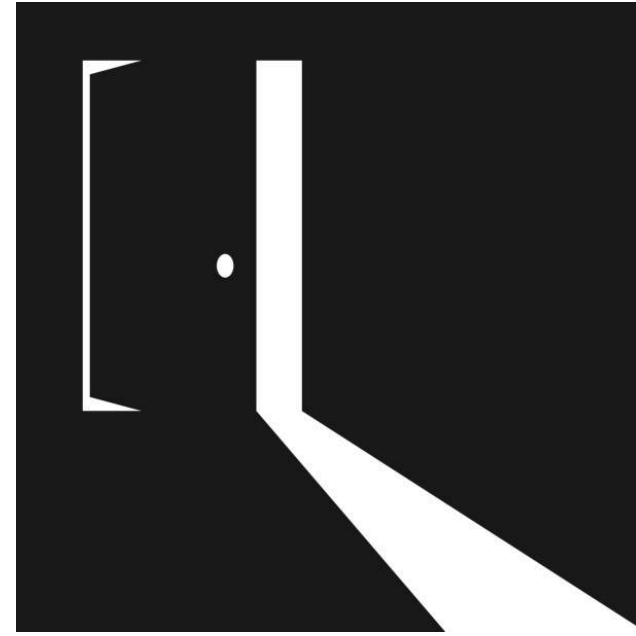


?



+100

$a^2$

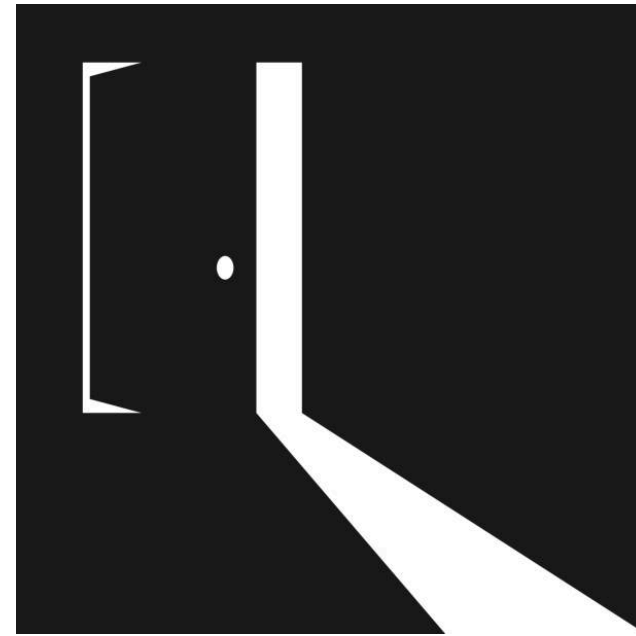


?



+1

$a^3$

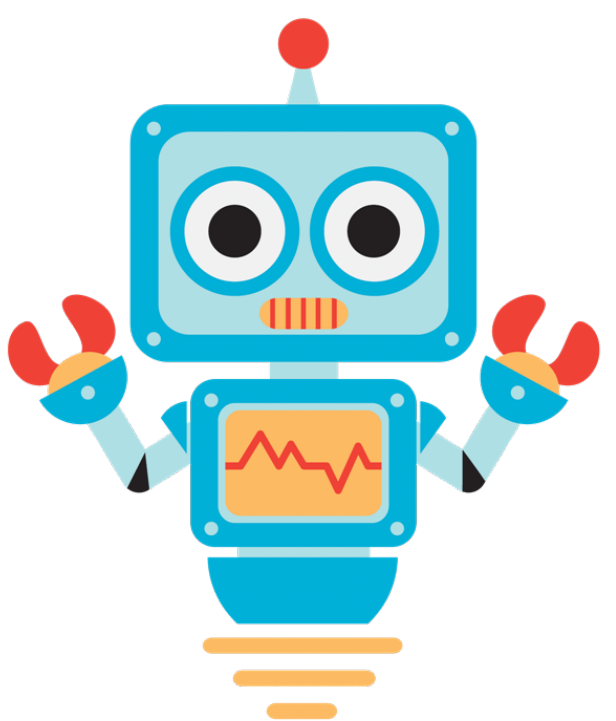


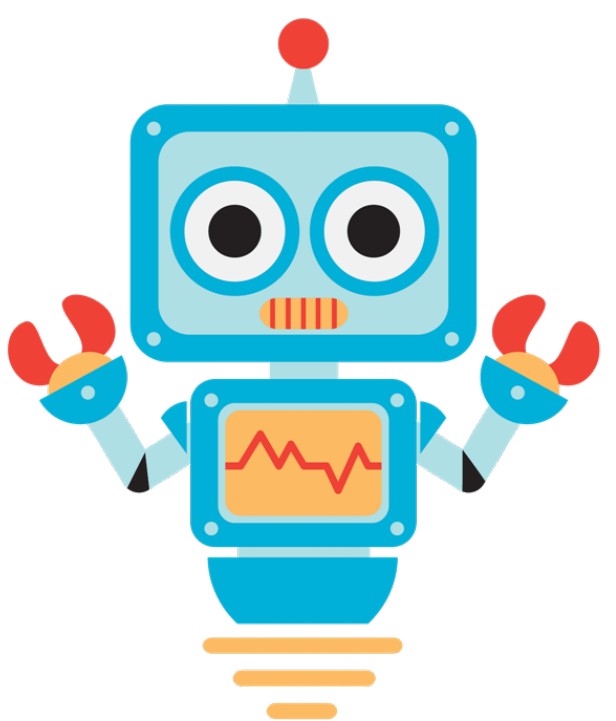
?



-1000

⋮





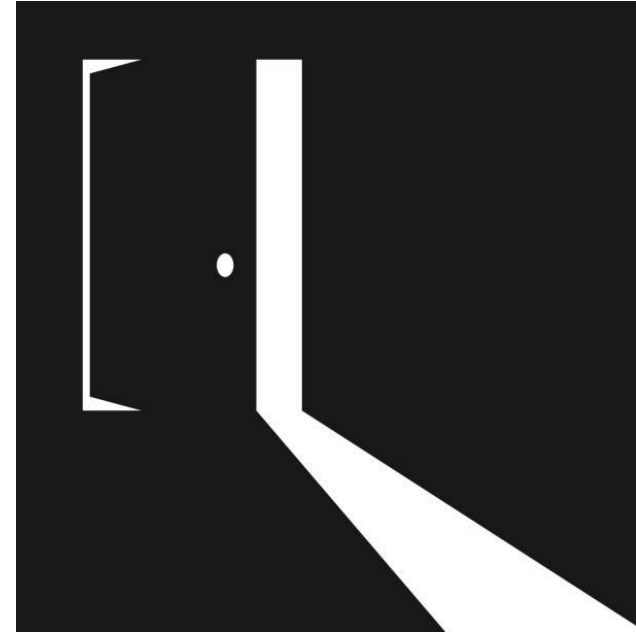
# Doors

Round 1

Round 2

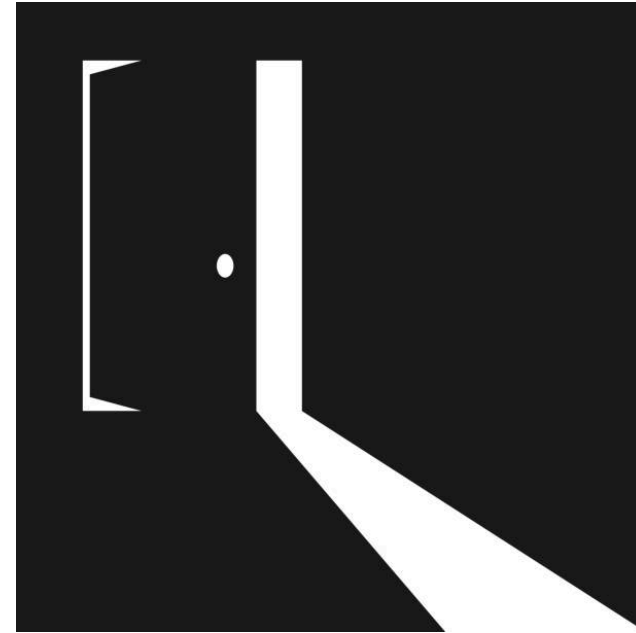
Round 3

$a^1$



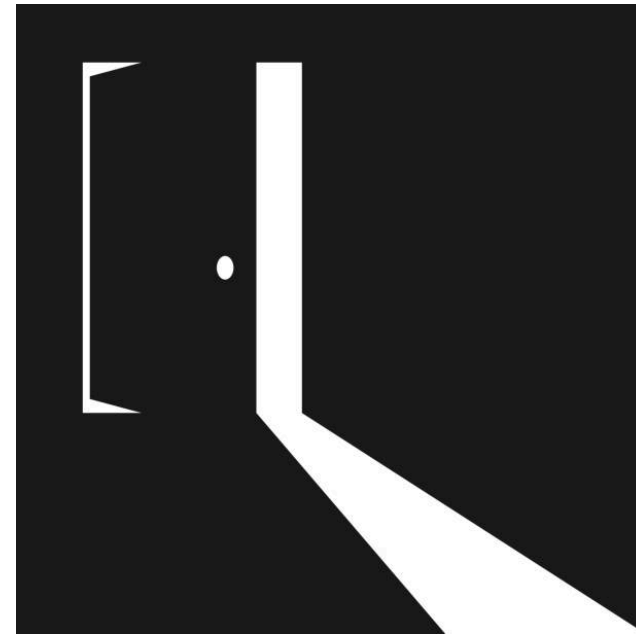
+100

$a^2$



+1

$a^3$



⋮



-1000

Activity!



# Think-Pair-Share

Think (30 sec): What strategy would you pick doors?

Pair: Find a partner

Share (45 sec): Partners exchange ideas

Doors

$a^1$  ?

$a^2$  ?

$a^3$  ?

...

+100

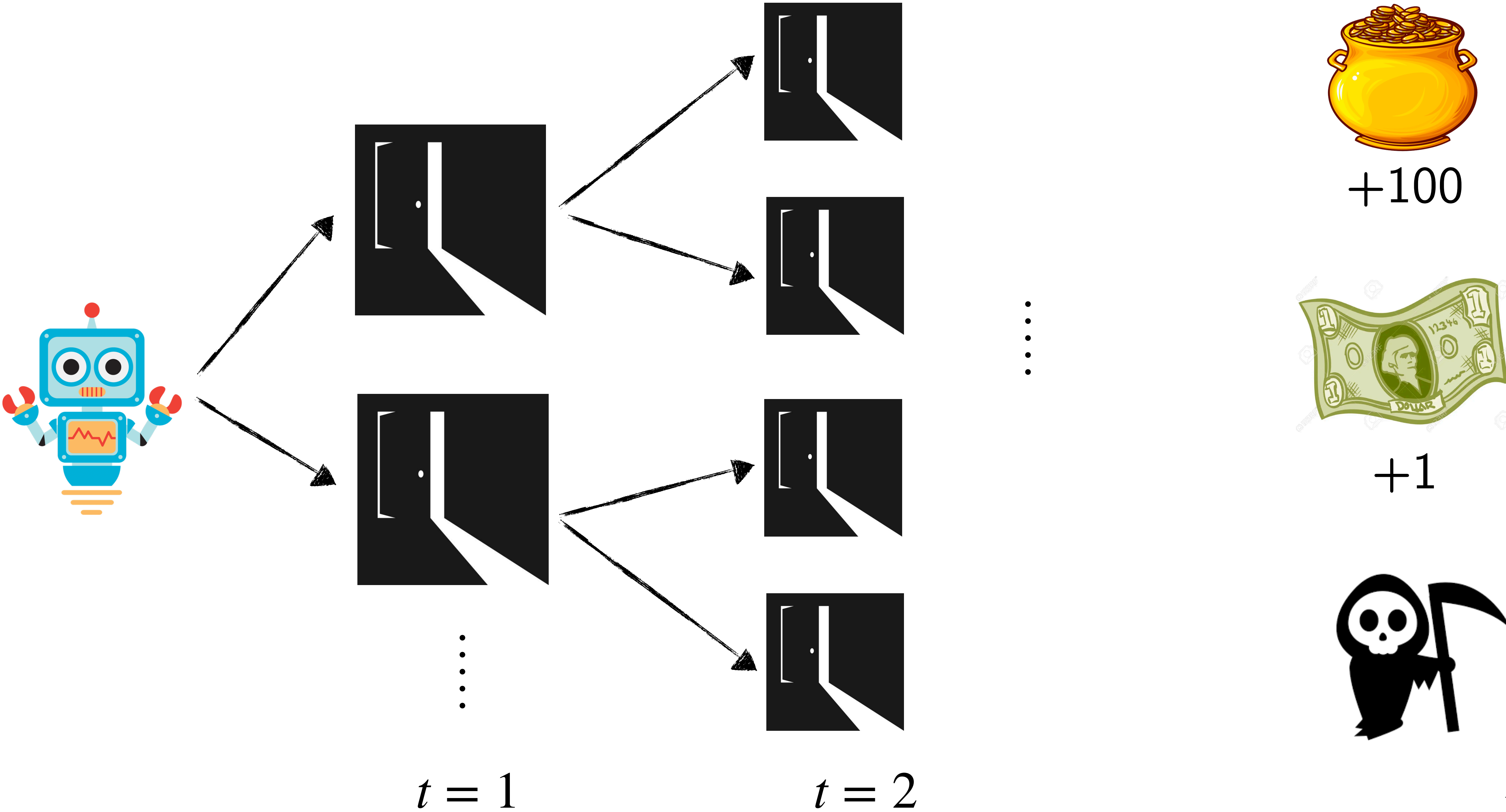
+1

-1000



What if we played the  
game over multiple time  
steps?





$t = 1$

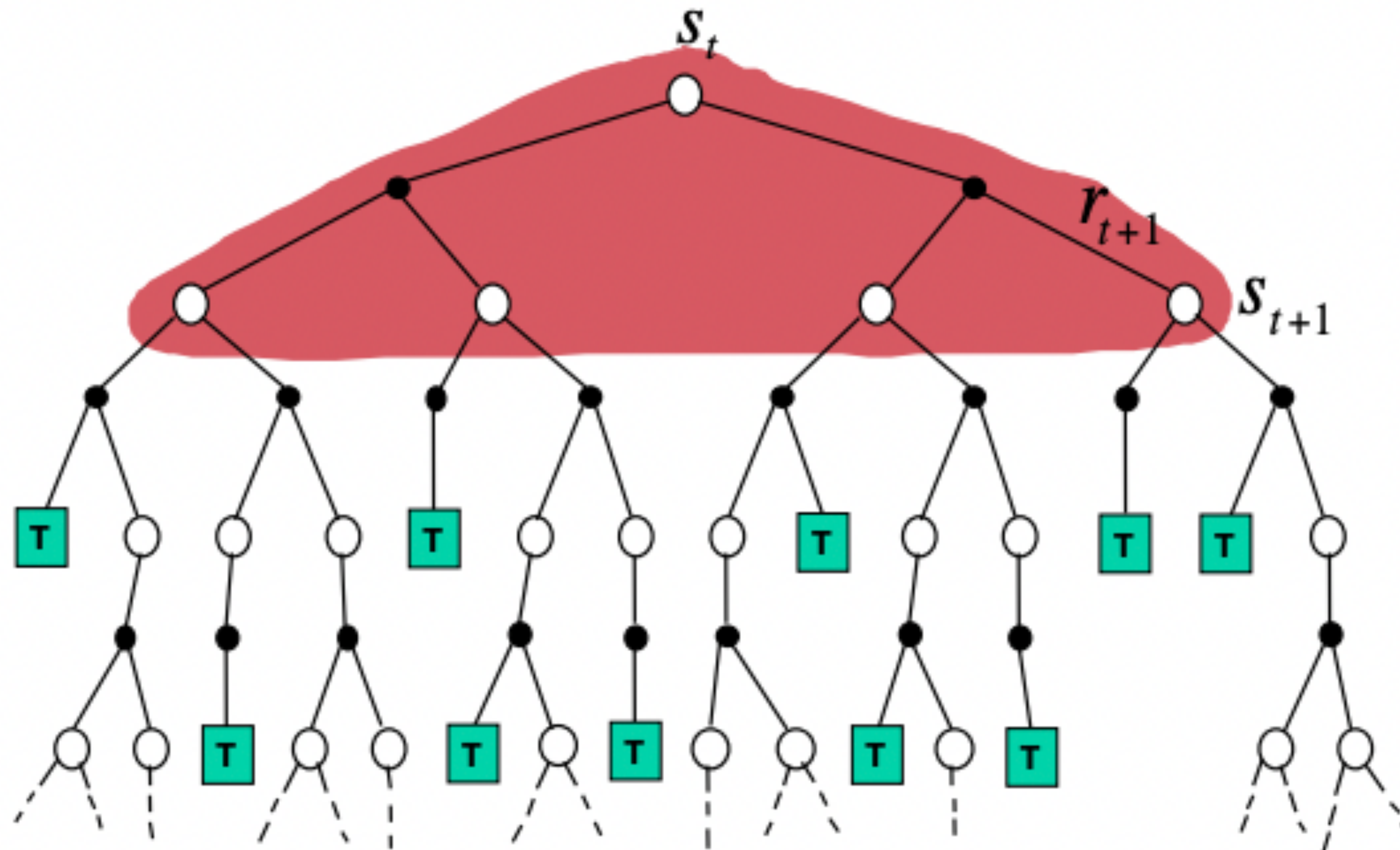
$t = 2$



Don't see how good a  
door is until the end of  
the episode

# When we know the MDP: Dynamic Programming!

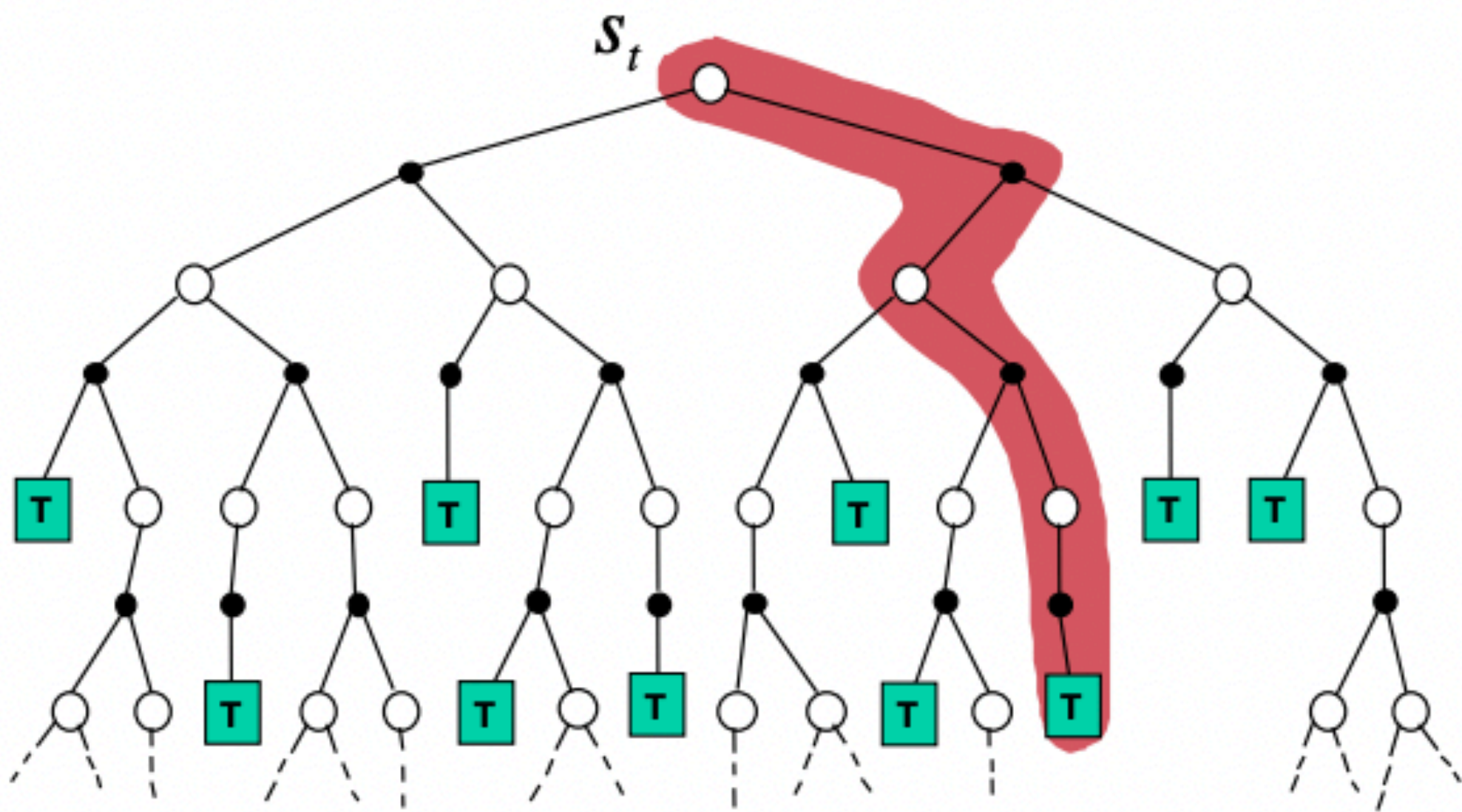
$$V(S_t) \leftarrow \mathbb{E}_\pi [R_{t+1} + \gamma V(S_{t+1})]$$



# When we don't know MDP: *Estimate Values*

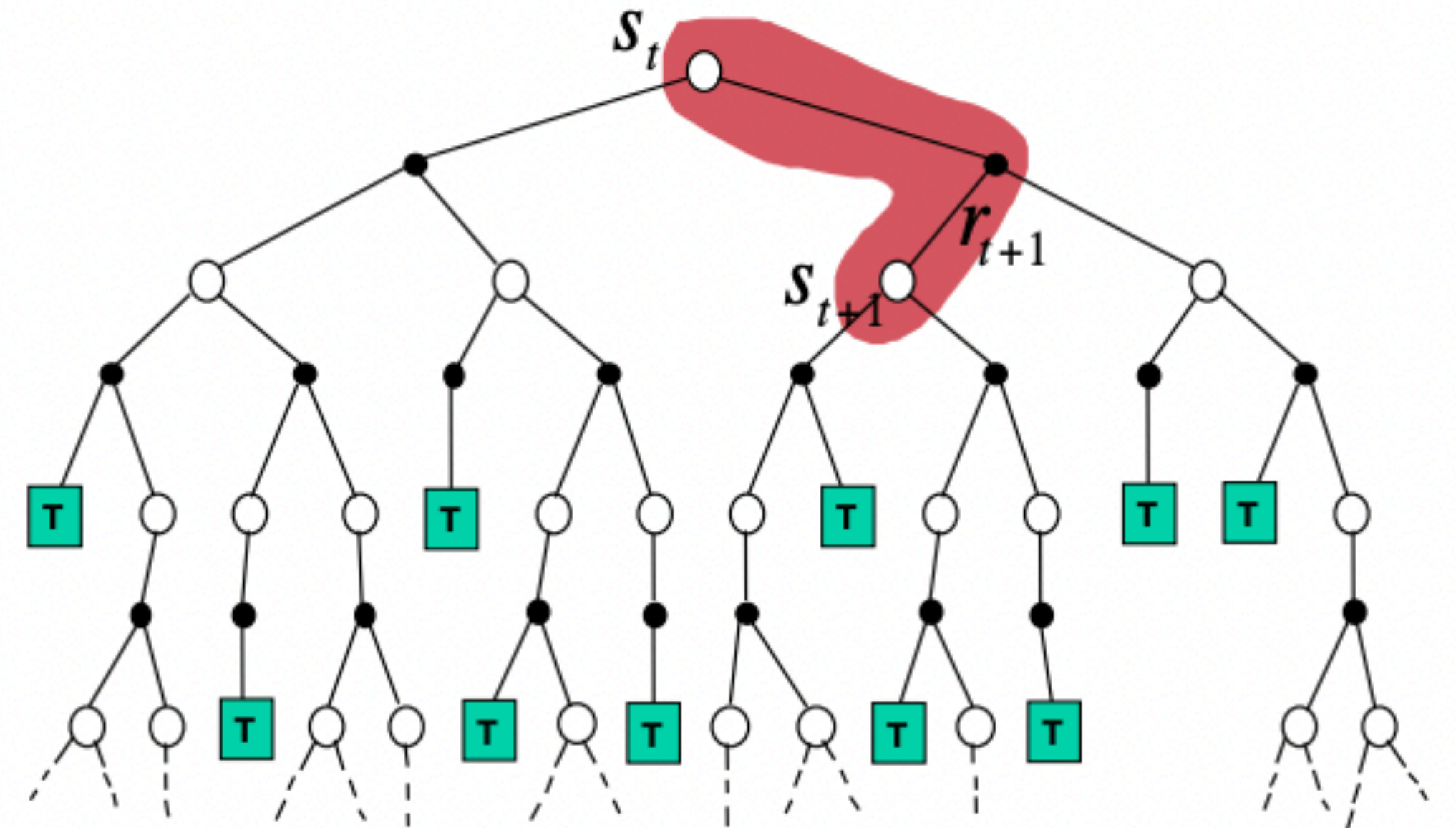
## Monte Carlo

$$V(S_t) \leftarrow V(S_t) + \alpha(G_t - V(S_t))$$

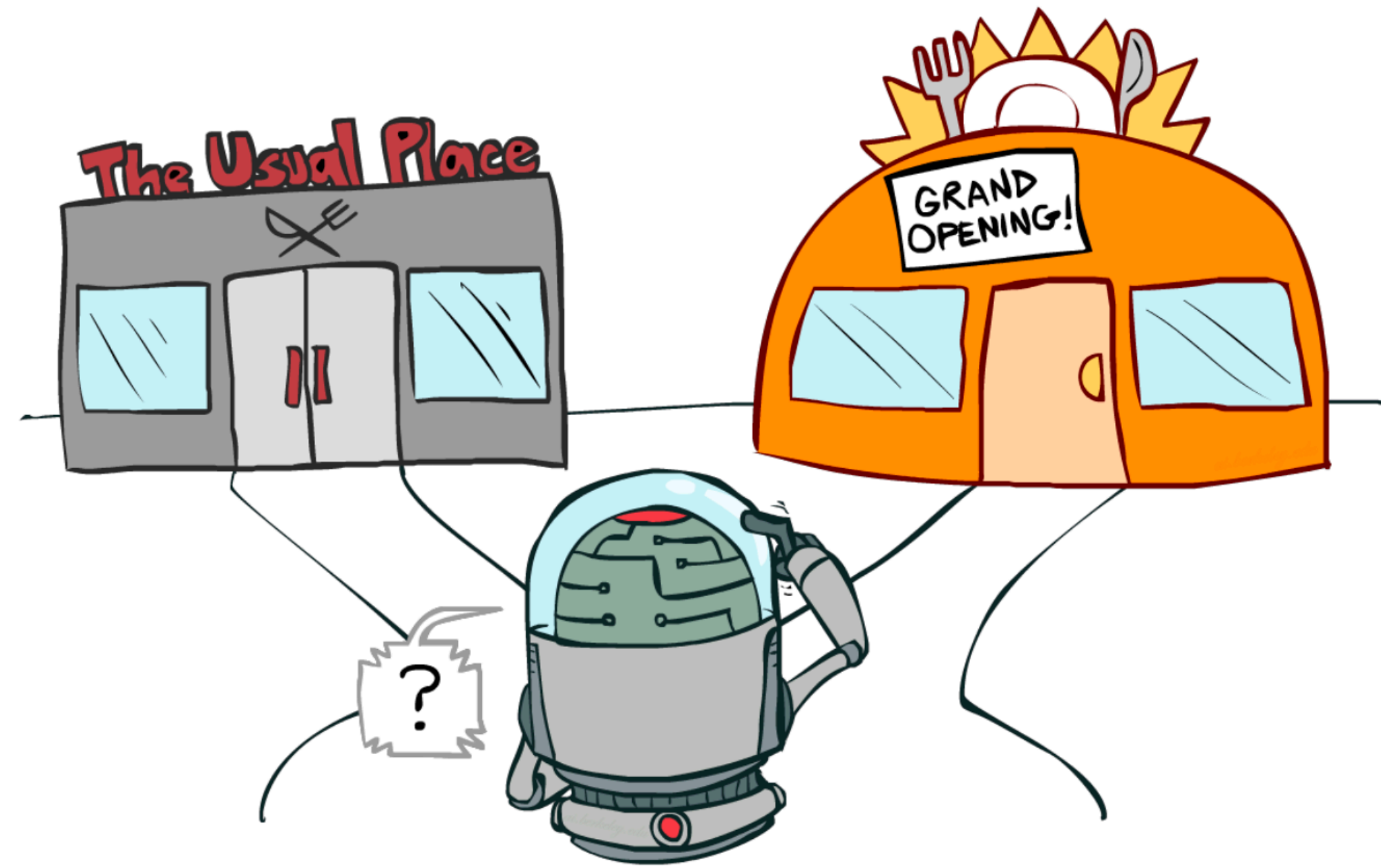


## Temporal Difference

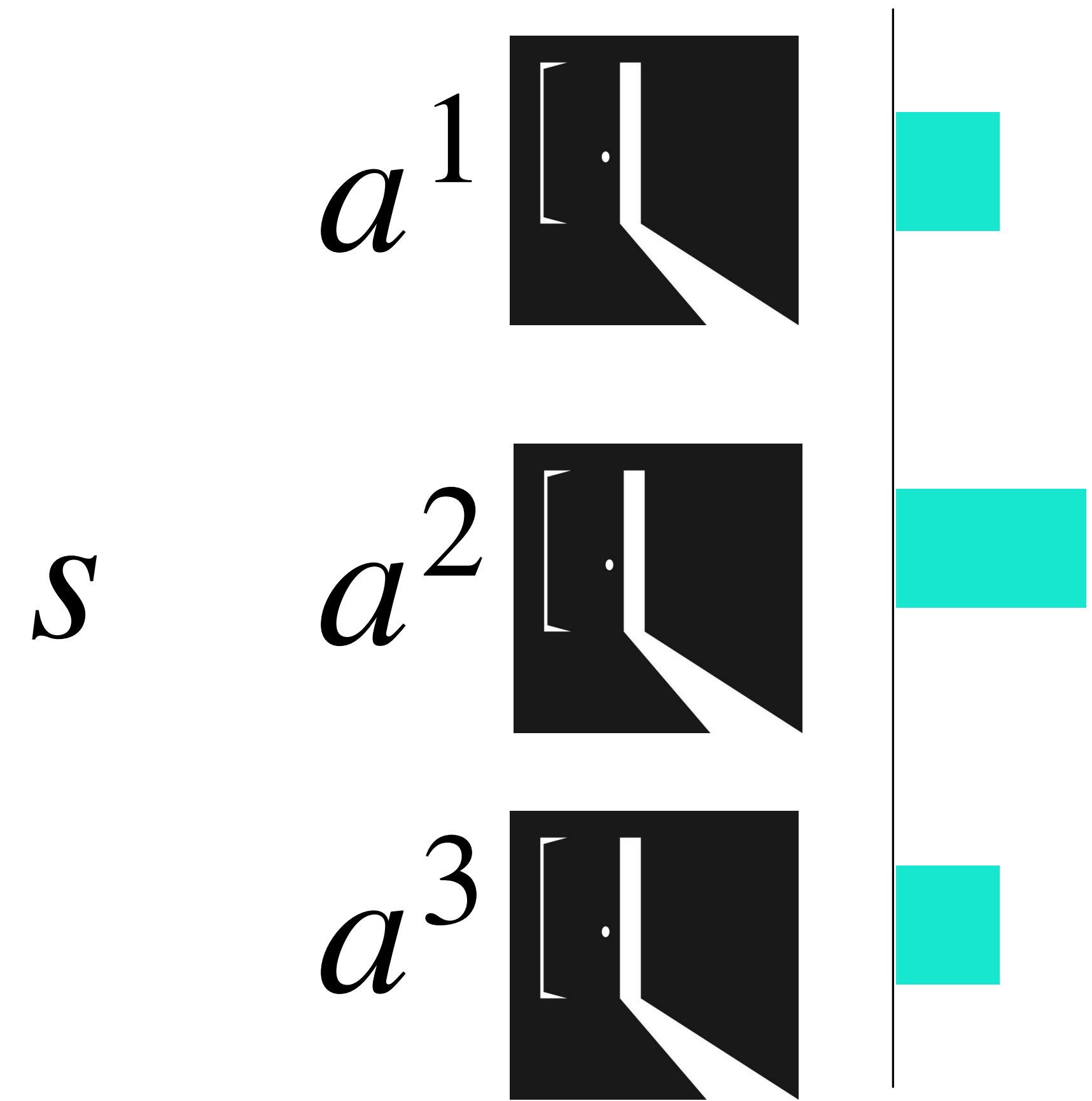
$$V(S_t) \leftarrow V(S_t) + \alpha(R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$



# Two Ingredients of RL



Exploration   Exploitation



Estimate Values  $Q(s, a)$

What is an application  
where we *really* need RL?



A close-up photograph of a hand holding a wooden chess piece, likely a king or queen, over a chessboard. The chessboard is checkered, and several other chess pieces are visible on the board, including a white pawn, a dark knight, a dark rook, a dark pawn, and a white rook. The background is dark and out of focus.

Games!



# Why are Games perfect for RL?

You know the cost function

\* +1 for winning and - 1 for losing

You don't know the transition function

\* Don't know what opponent will do, stochasticity in dice rolls etc.

Perfect for learning from trial and error (by playing itself!)



# Name the Game!

One of the biggest success stories of RL

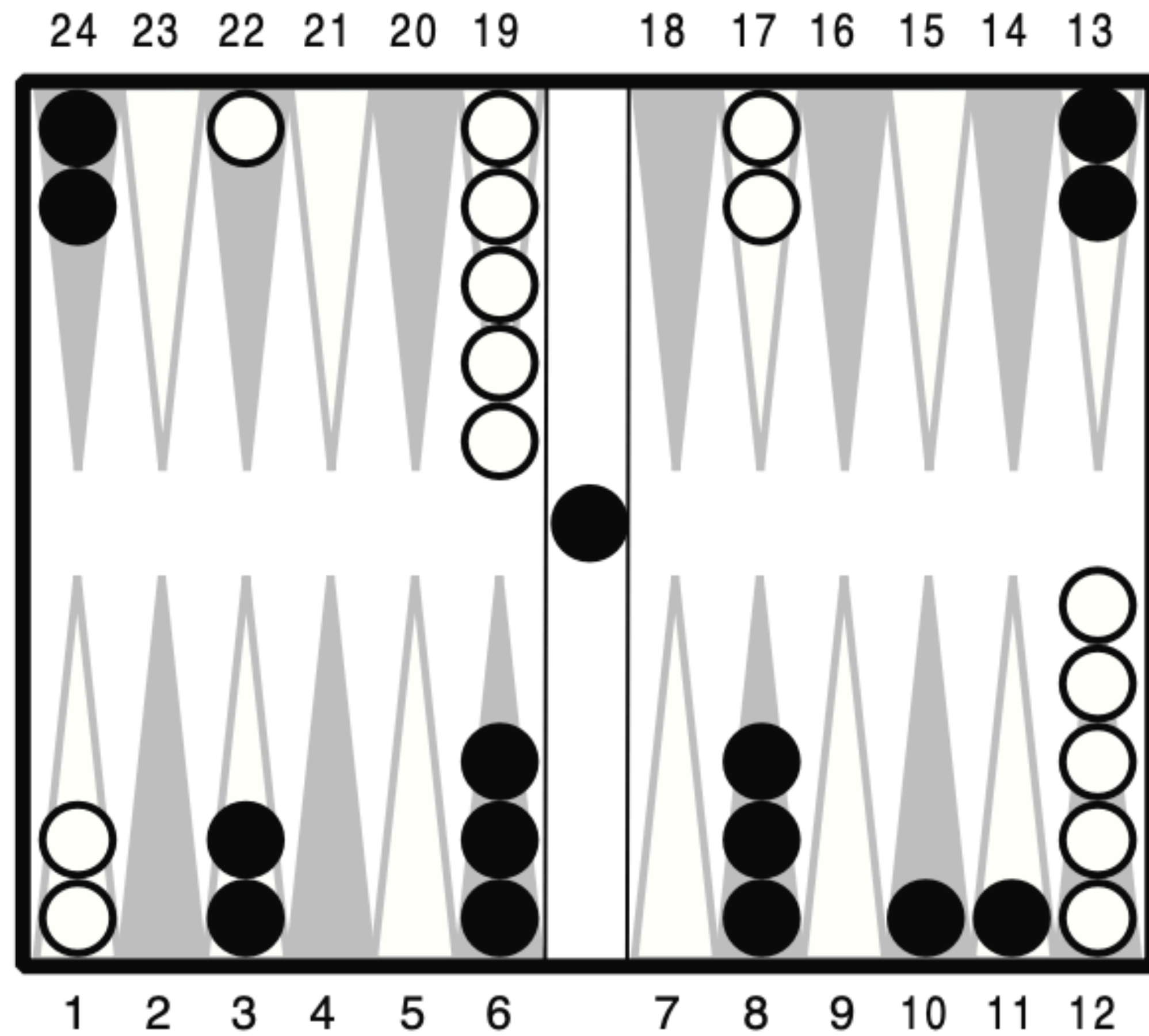
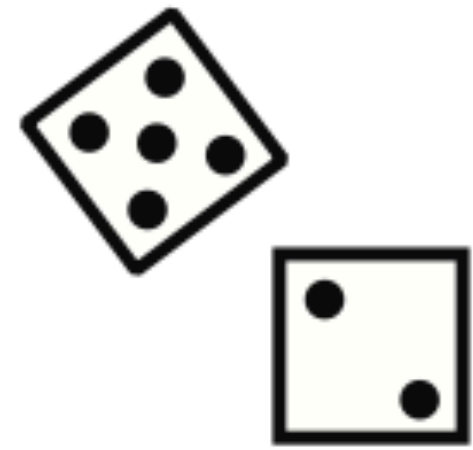
An agent trained via  
self-play and neural networks

Beat the world champion

Discovered totally new moves



# The game of Backgammon



white pieces move  
counterclockwise

black pieces  
move clockwise

A backgammon position

# TD-Gammon

Initialised with random weights

Trained by games of self-play

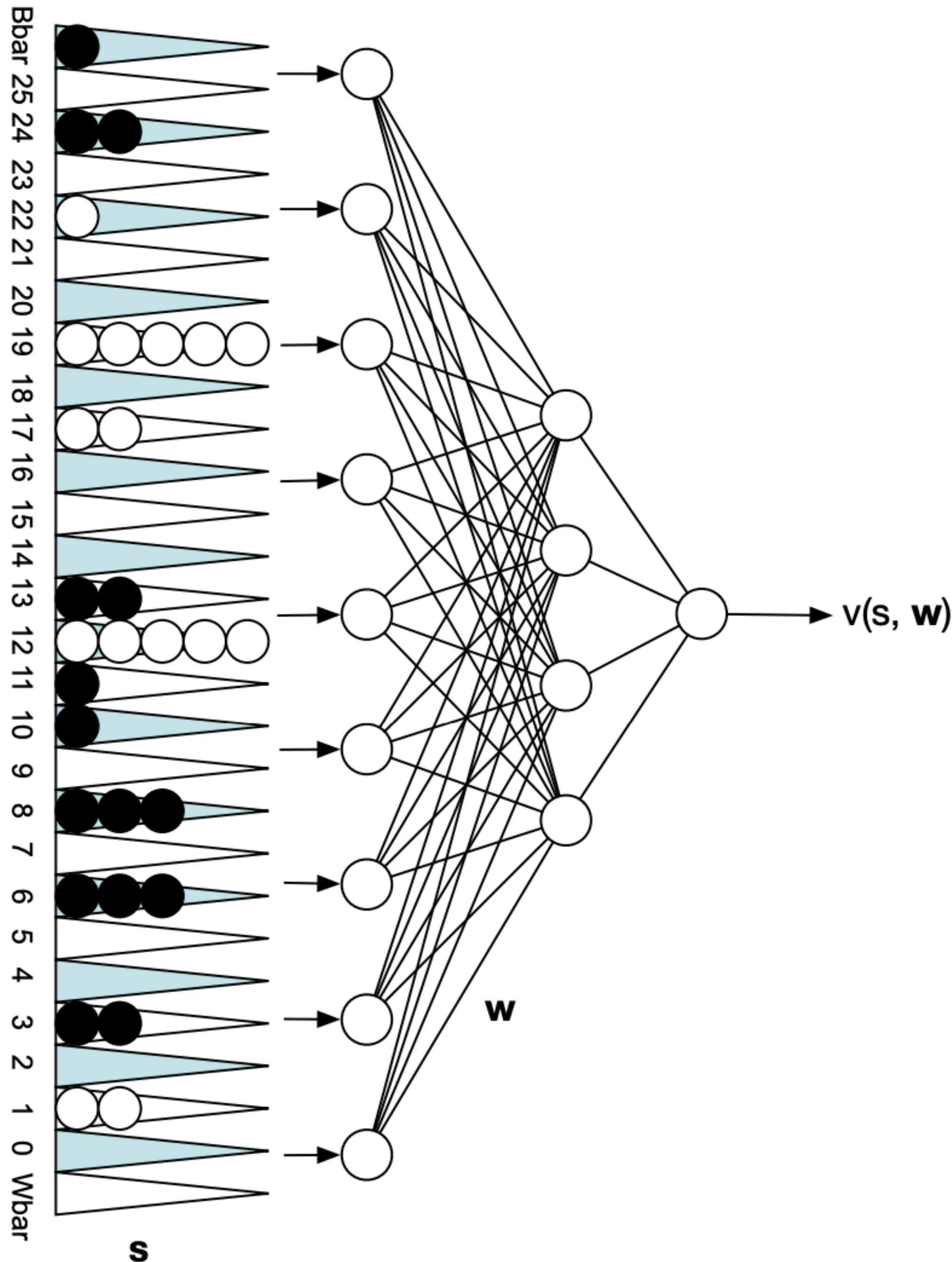
Using non-linear temporal-difference learning

$$\delta_t = v(S_{t+1}, \mathbf{w}) - v(S_t, \mathbf{w})$$

$$\Delta \mathbf{w} = \alpha \delta_t \nabla_{\mathbf{w}} v(S_t, \mathbf{w})$$

Greedy policy improvement (no exploration)

Algorithm always converged in practice

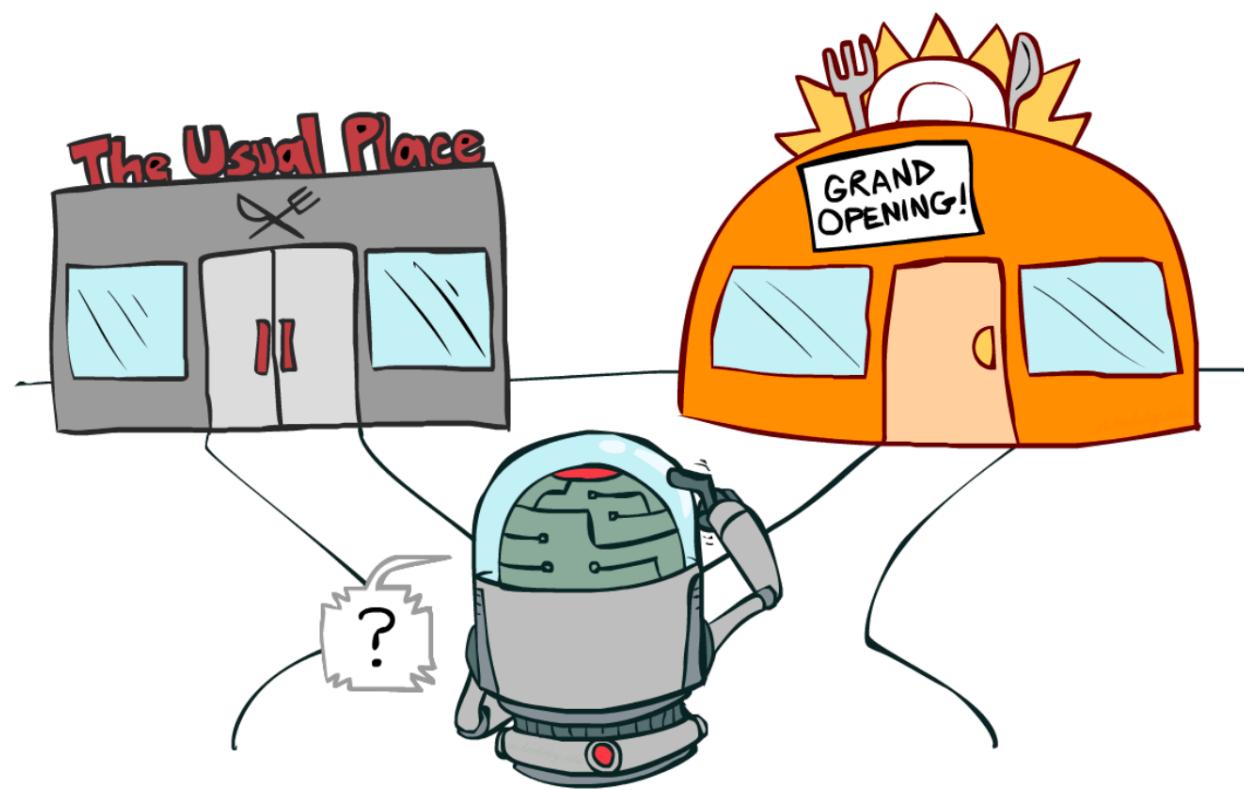


# Why was TD Gammon such a big deal?

- 1) Power of self- play
- 2) First time system could deal with stochastic uncertainty in dynamics (this broke everything in the deep blue style chess engines that were first super human)
- 3) Actually learned a value function: first "human" like behavior rather than just deep search with some heuristic values.
- 4) Changed elements of how backgammon is played because it demonstrated that certain positions were more valuable than self play
- 5) Demonstrated the TD idea could be scaled to super human performance at a game previously unaddressable.
- 6) Use a "deep" (i.e. Neural) representation within RL algorithms
- 7) Saw the benefit of imitation learning, but eventually got better than imitation
- 8) Showed boosting performance by some explicit forward search

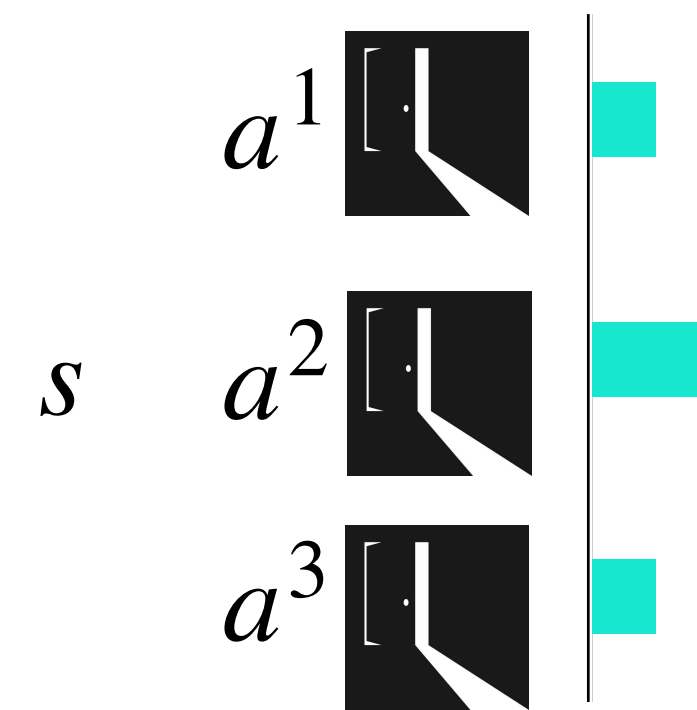
# Two Ingredients of RL: TD-Gammon

No exploration,  
was just greedy



Exploration Exploitation

Used temporal difference  
to estimate values



Estimate Values

Okay, but what about  
games with more  
complex representations?





# Circa 2013

---

## **Playing Atari with Deep Reinforcement Learning**

---

**Volodymyr Mnih   Koray Kavukcuoglu   David Silver   Alex Graves   Ioannis Antonoglou**

**Daan Wierstra   Martin Riedmiller**

DeepMind Technologies

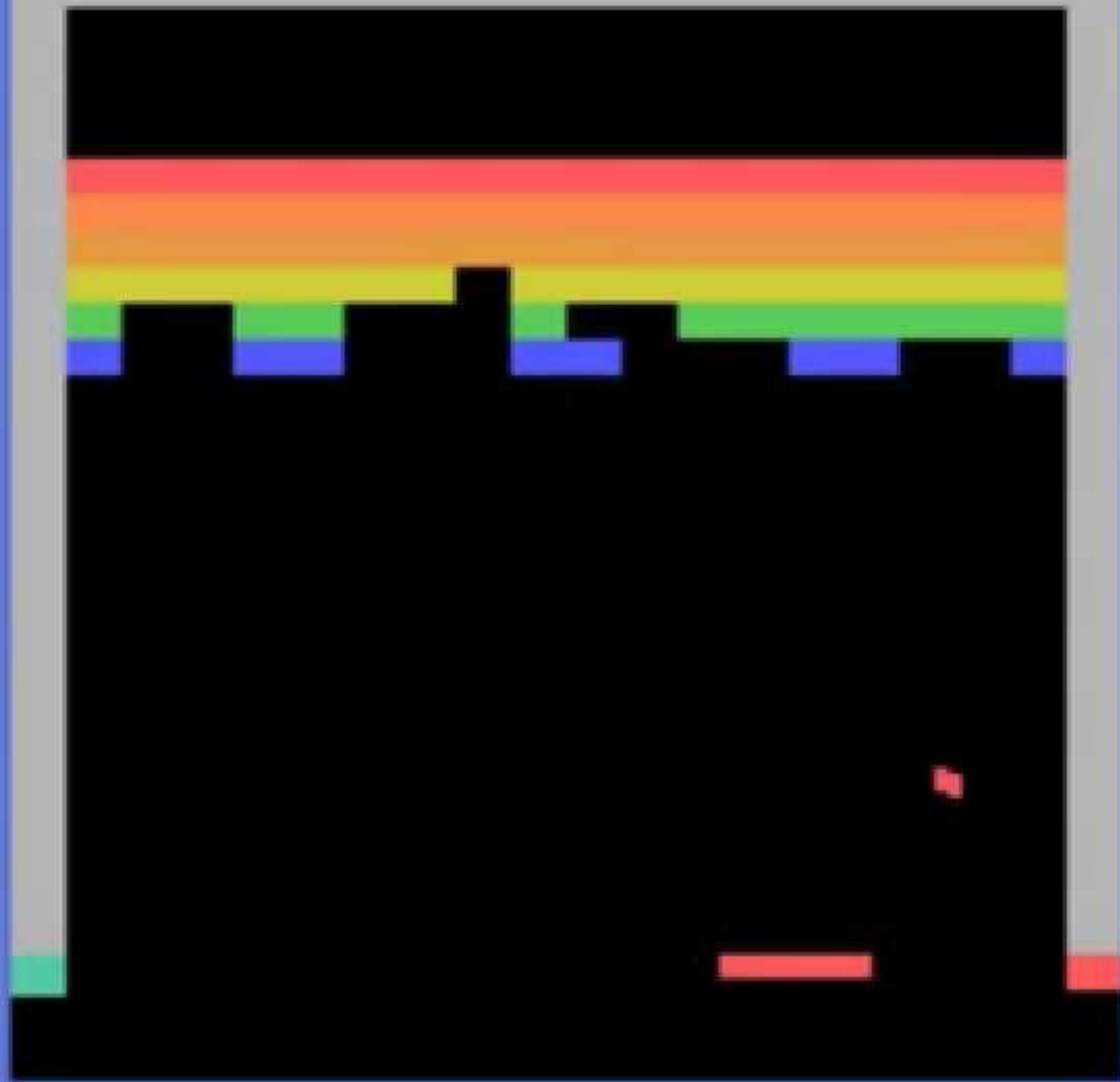
{vlad, koray, david, alex.graves, ioannis, daan, martin.riedmiller} @ deepmind.com



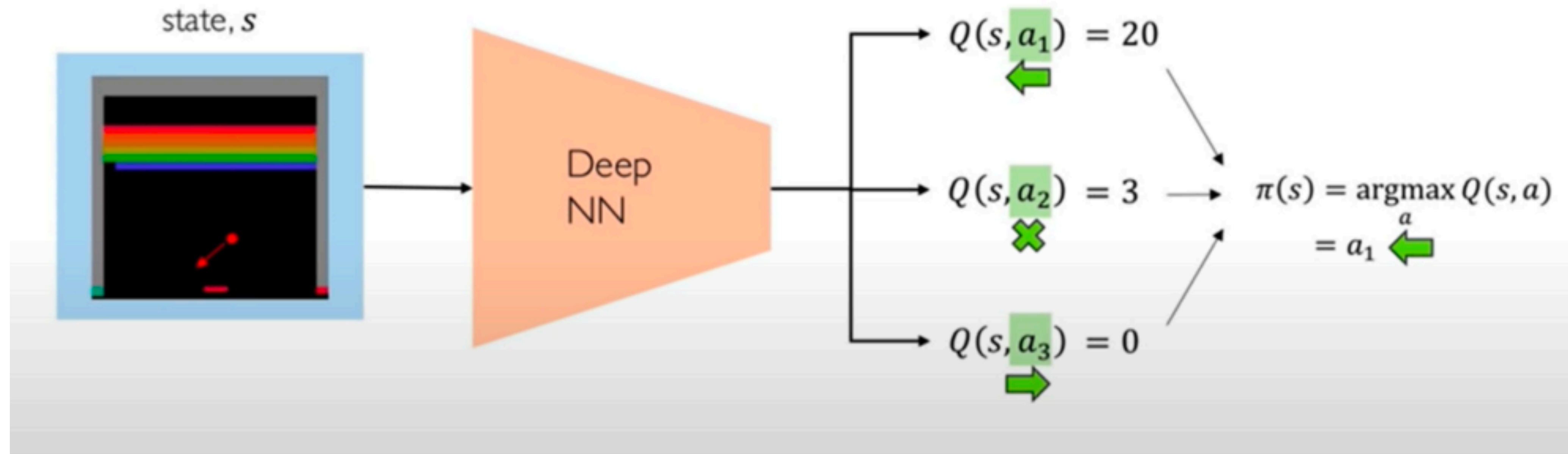
ima...



021 3 1

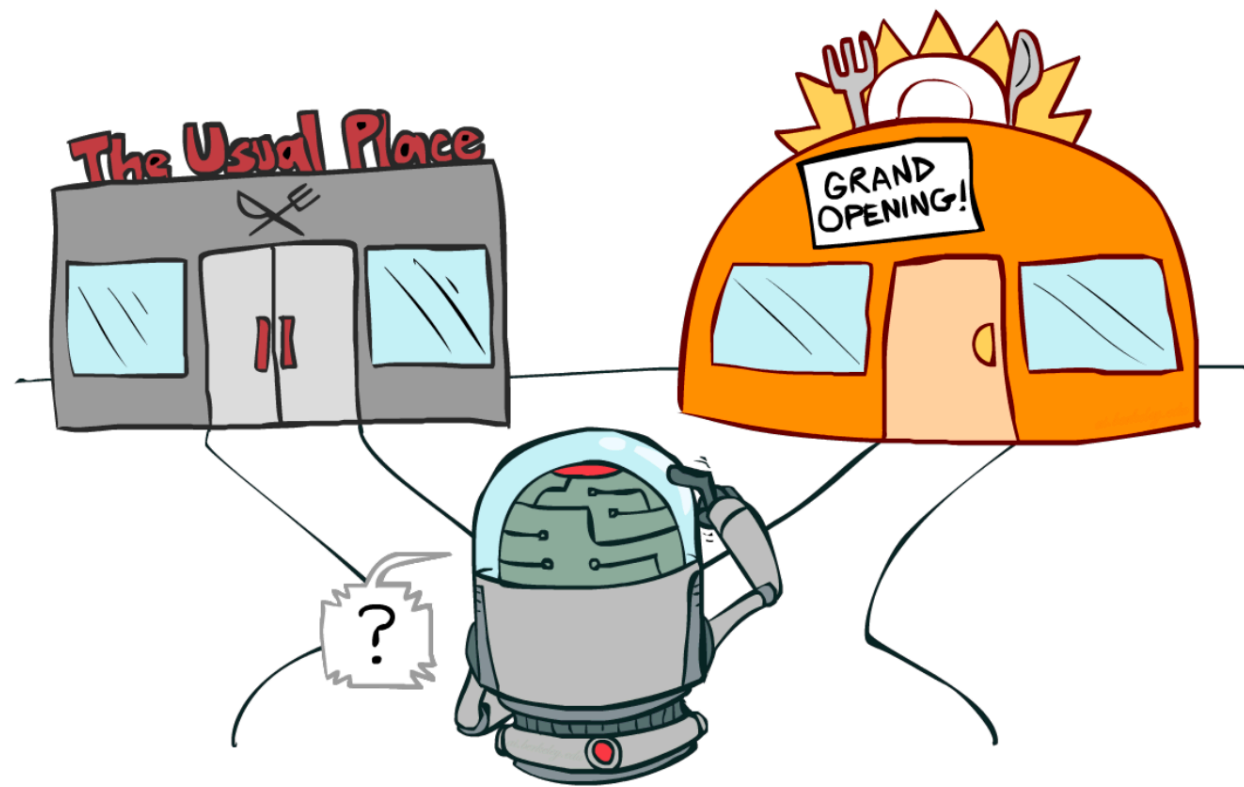


Use NN to learn Q-function and then use to infer the optimal policy,  $\pi(s)$



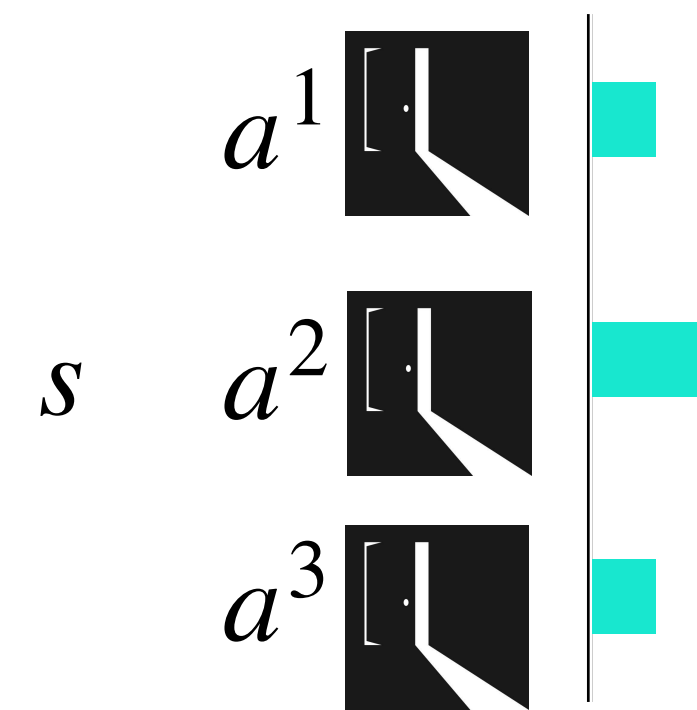
# Two Ingredients of RL: DQN

Epsilon-Greedy



Exploration Exploitation

Used Q-learning with  
Deep CNN



Estimate Values

Okay, but what about games that require some level of planning?



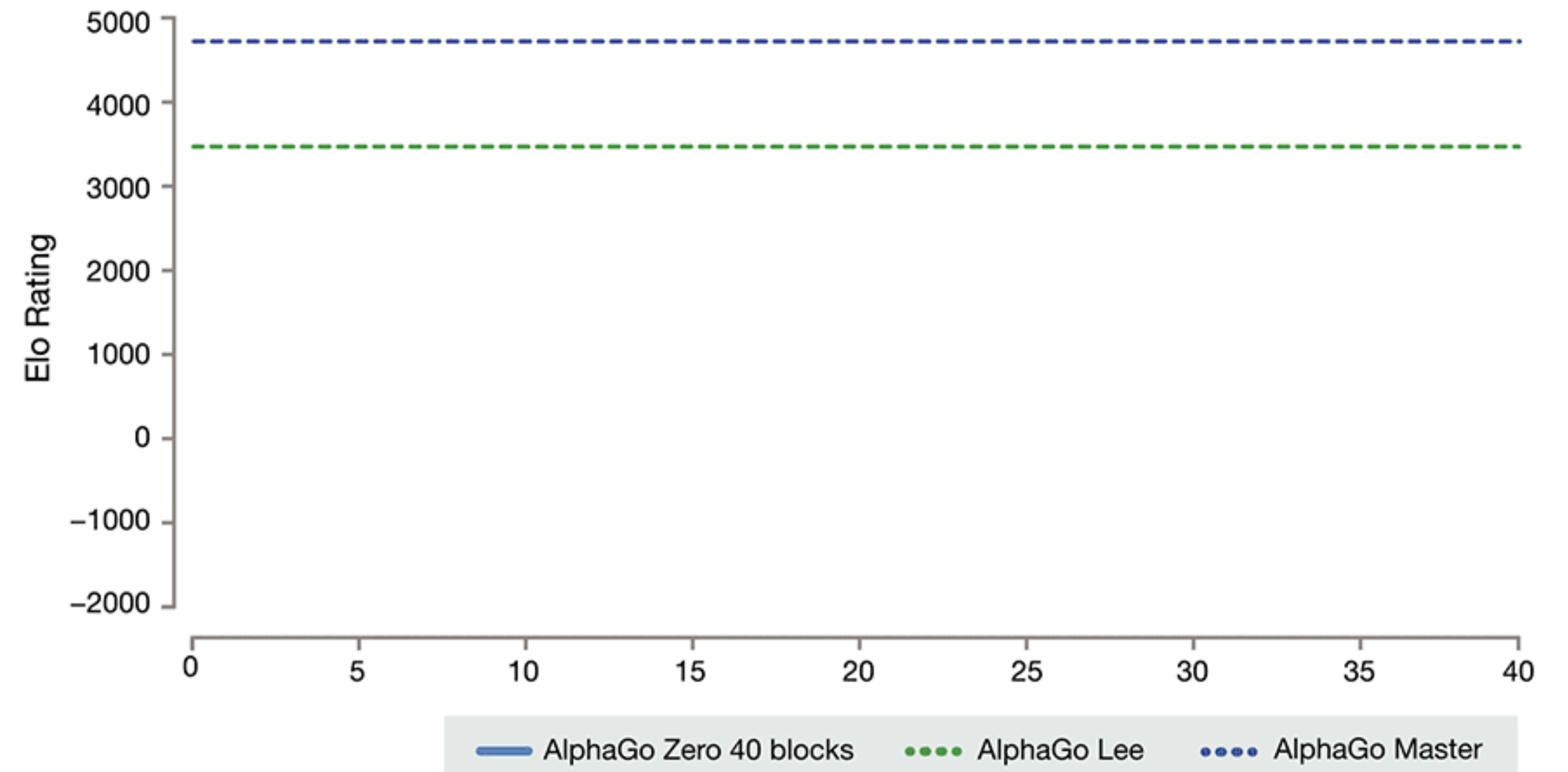
# AlphaGoZero

## Mastering the game of Go without human knowledge

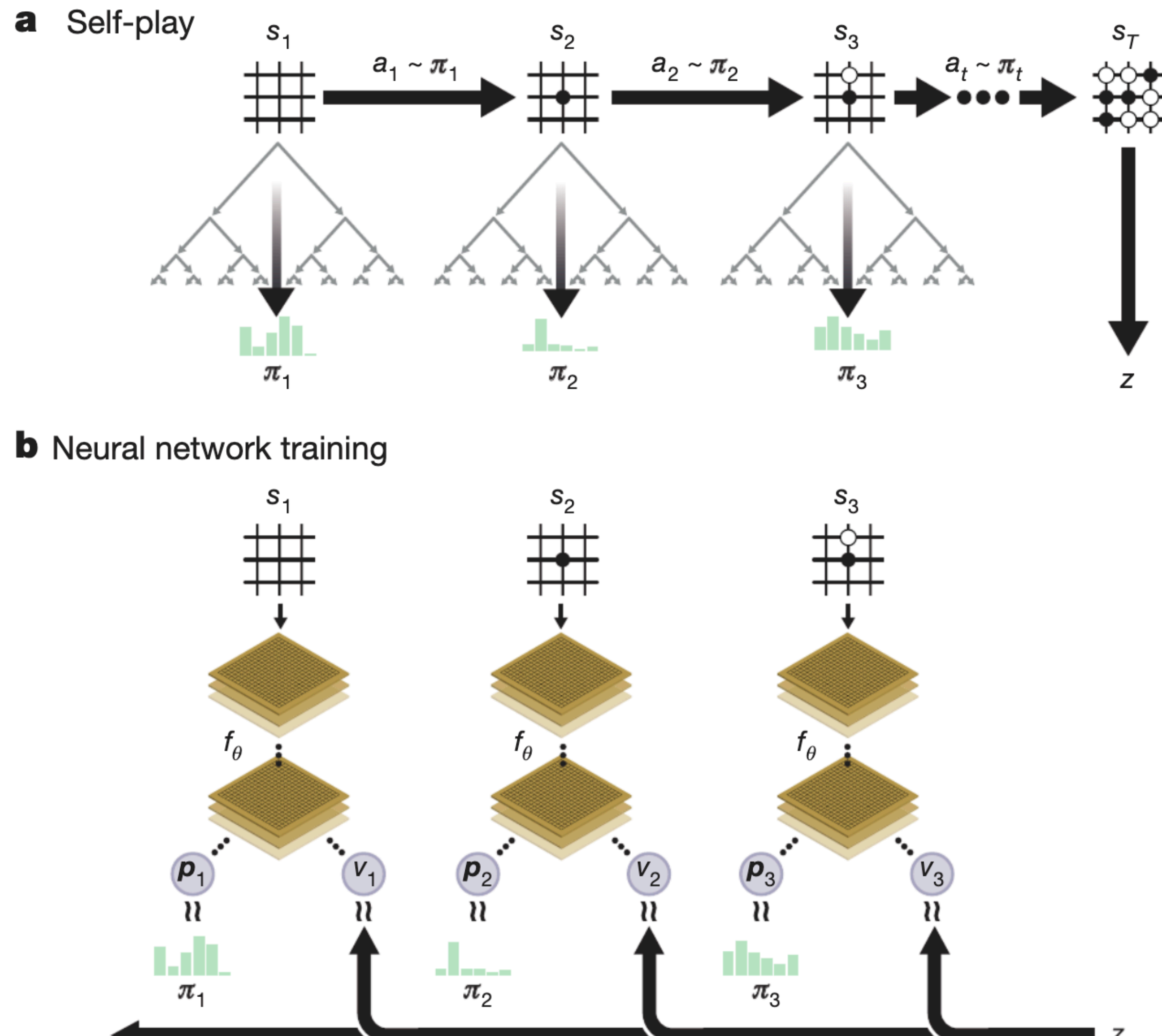
David Silver<sup>1\*</sup>, Julian Schrittwieser<sup>1\*</sup>, Karen Simonyan<sup>1\*</sup>, Ioannis Antonoglou<sup>1</sup>, Aja Huang<sup>1</sup>, Arthur Guez<sup>1</sup>, Thomas Hubert<sup>1</sup>, Lucas Baker<sup>1</sup>, Matthew Lai<sup>1</sup>, Adrian Bolton<sup>1</sup>, Yutian Chen<sup>1</sup>, Timothy Lillicrap<sup>1</sup>, Fan Hui<sup>1</sup>, Laurent Sifre<sup>1</sup>, George van den Driessche<sup>1</sup>, Thore Graepel<sup>1</sup> & Demis Hassabis<sup>1</sup>



AlphaGo



# AlphaGoZero



**Figure 1 | Self-play reinforcement learning in AlphaGo Zero.** **a**, The program plays a game  $s_1, \dots, s_T$  against itself. In each position  $s_t$ , an MCTS  $\alpha_\theta$  is executed (see Fig. 2) using the latest neural network  $f_\theta$ . Moves are selected according to the search probabilities computed by the MCTS,  $a_t \sim \pi_t$ . The terminal position  $s_T$  is scored according to the rules of the game to compute the game winner  $z$ . **b**, Neural network training in AlphaGo Zero. The neural network takes the raw board position  $s_t$  as its input, passes it through many convolutional layers with parameters  $\theta$ , and outputs both a vector  $\mathbf{p}_t$ , representing a probability distribution over moves, and a scalar value  $v_t$ , representing the probability of the current player winning in position  $s_t$ . The neural network parameters  $\theta$  are updated to maximize the similarity of the policy vector  $\mathbf{p}_t$  to the search probabilities  $\pi_t$ , and to minimize the error between the predicted winner  $v_t$  and the game winner  $z$  (see equation (1)). The new parameters are used in the next iteration of self-play as in **a**.

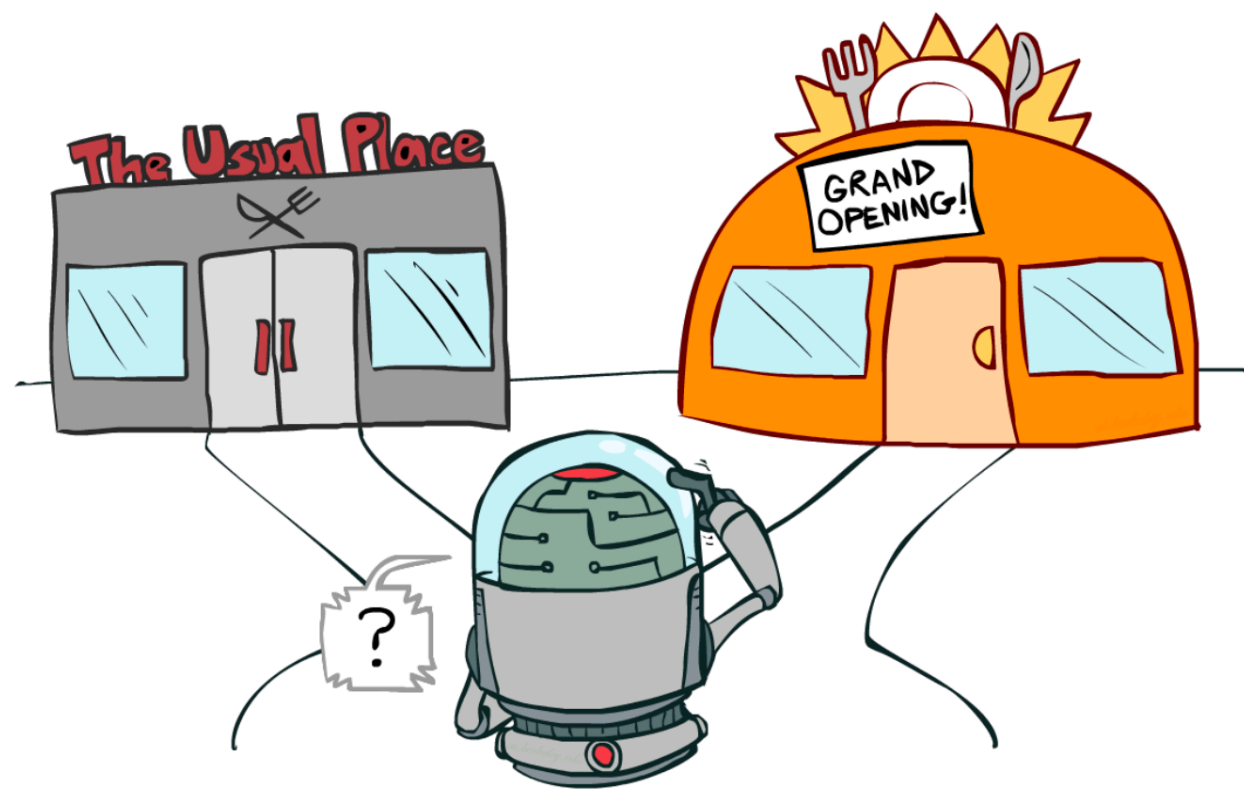
# Two Ingredients of RL: AlphaGoZero

Monte Carlo Tree Search

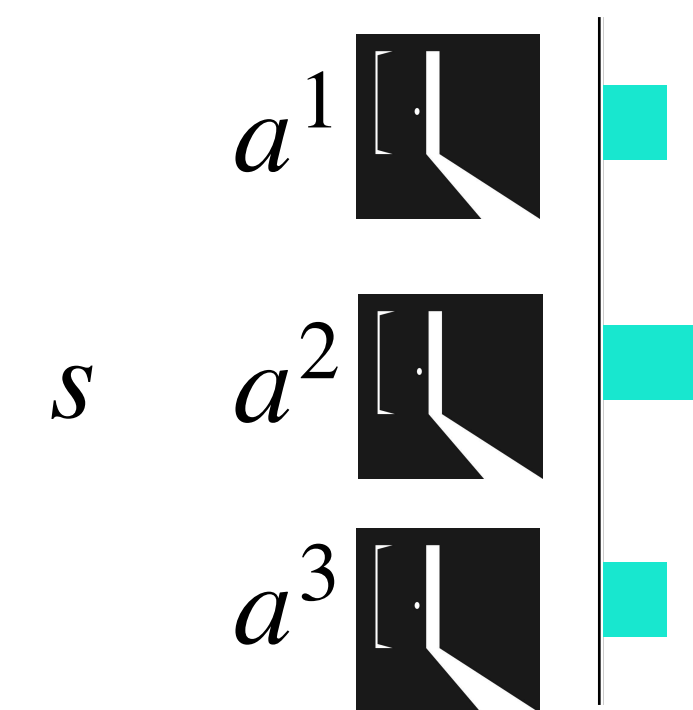
Deep value network

+

Deep policy network



Exploration Exploitation



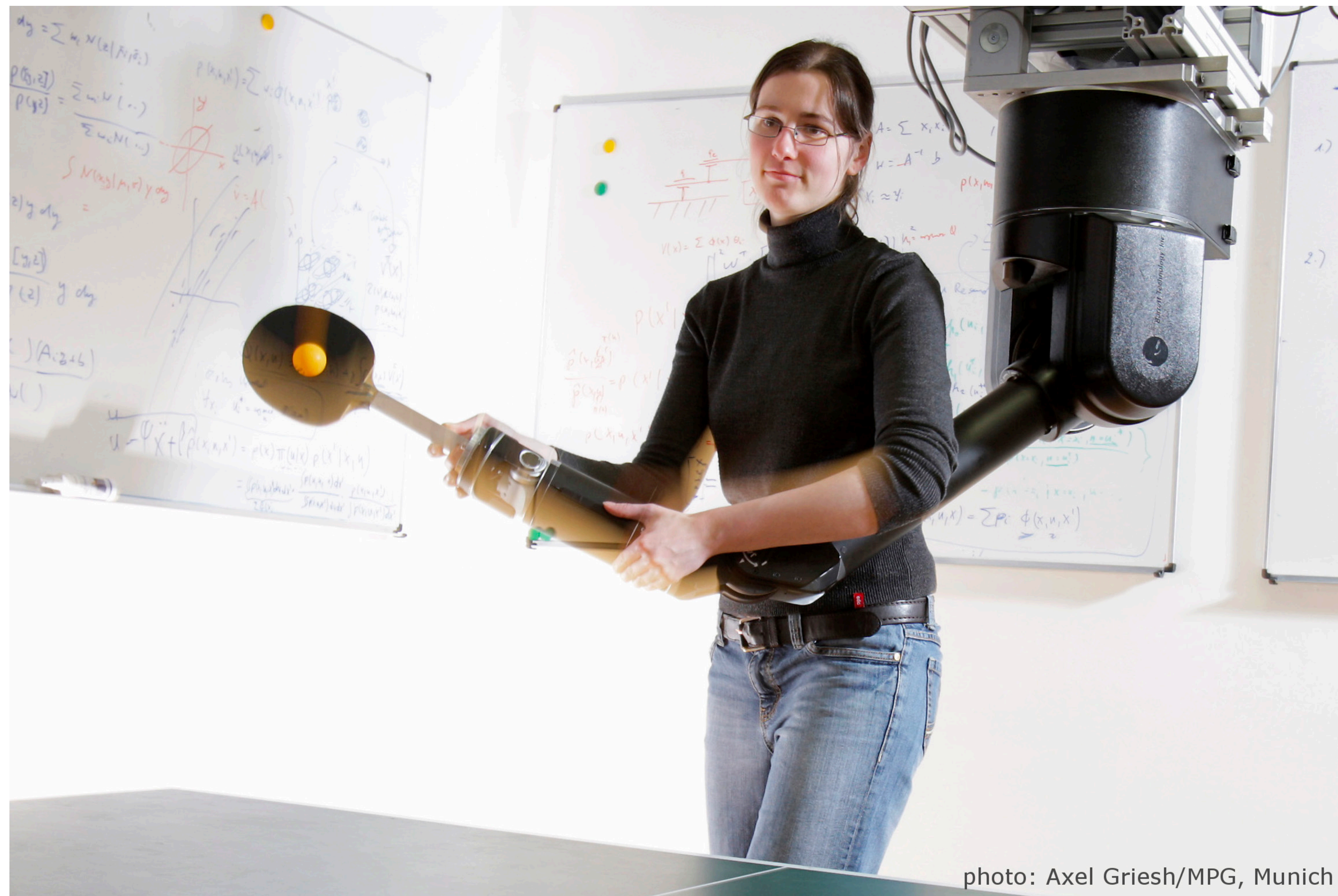
Estimate Values



# So what happens when we replace the human with a robot?







## Learning Strategies in Table Tennis using Inverse Reinforcement Learning

Katharina Muelling · Abdeslam Boularias · Betty Mohler ·  
Bernhard Schölkopf · Jan Peters

# Reinforcement Learning in Robotics: A Survey

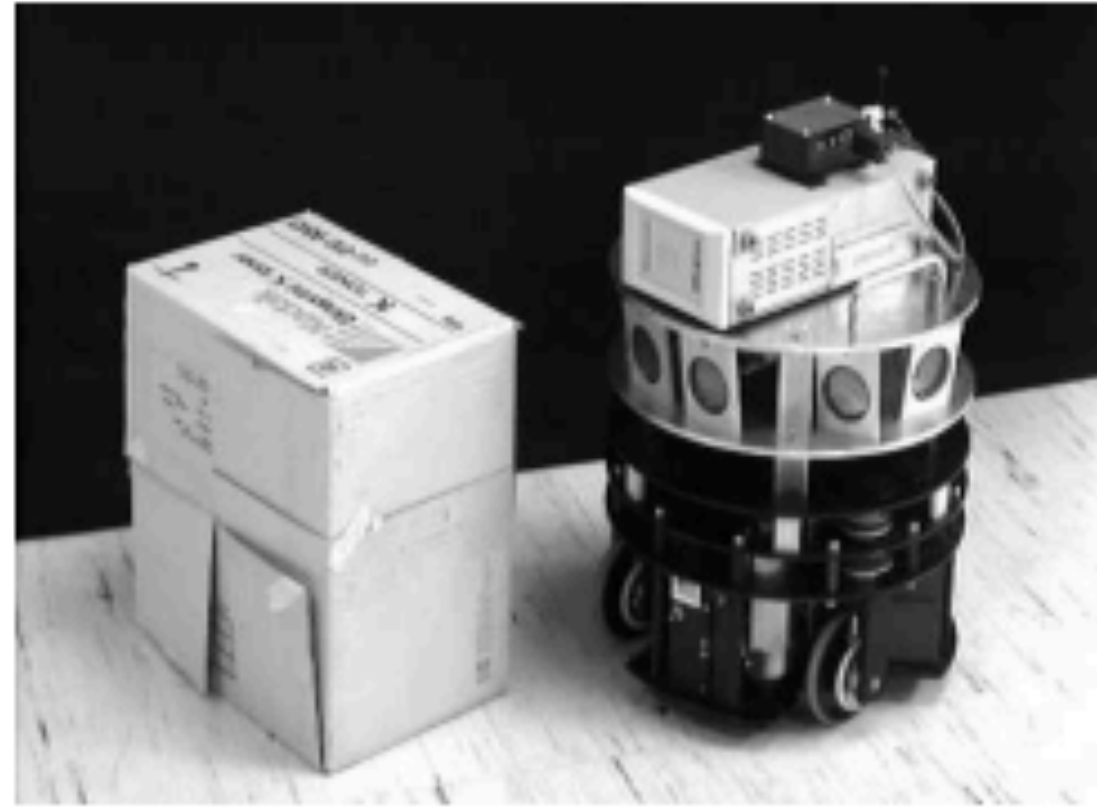
Jens Kober<sup>\*†</sup>

J. Andrew Bagnell<sup>‡</sup>

Jan Peters<sup>§¶</sup>

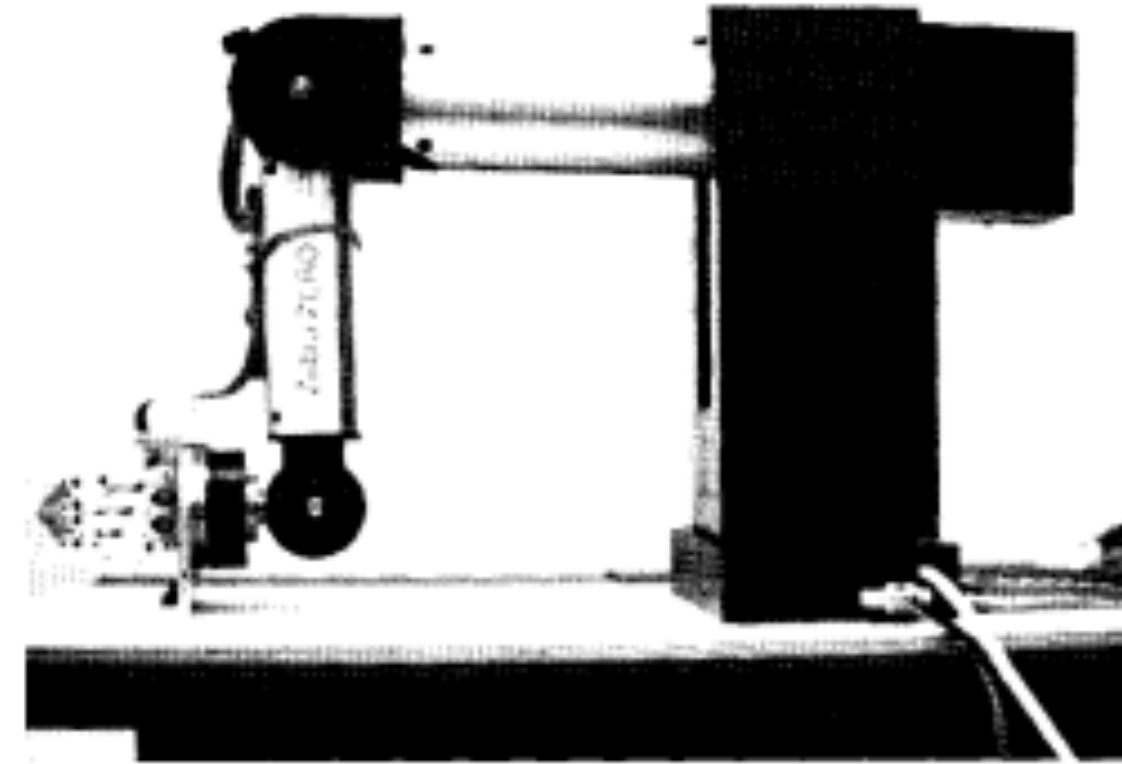
email: jkober@cor-lab.uni-bielefeld.de, dbagnell@ri.cmu.edu, mail@jan-peters.net

1992



(a) OBELIX robot

1994



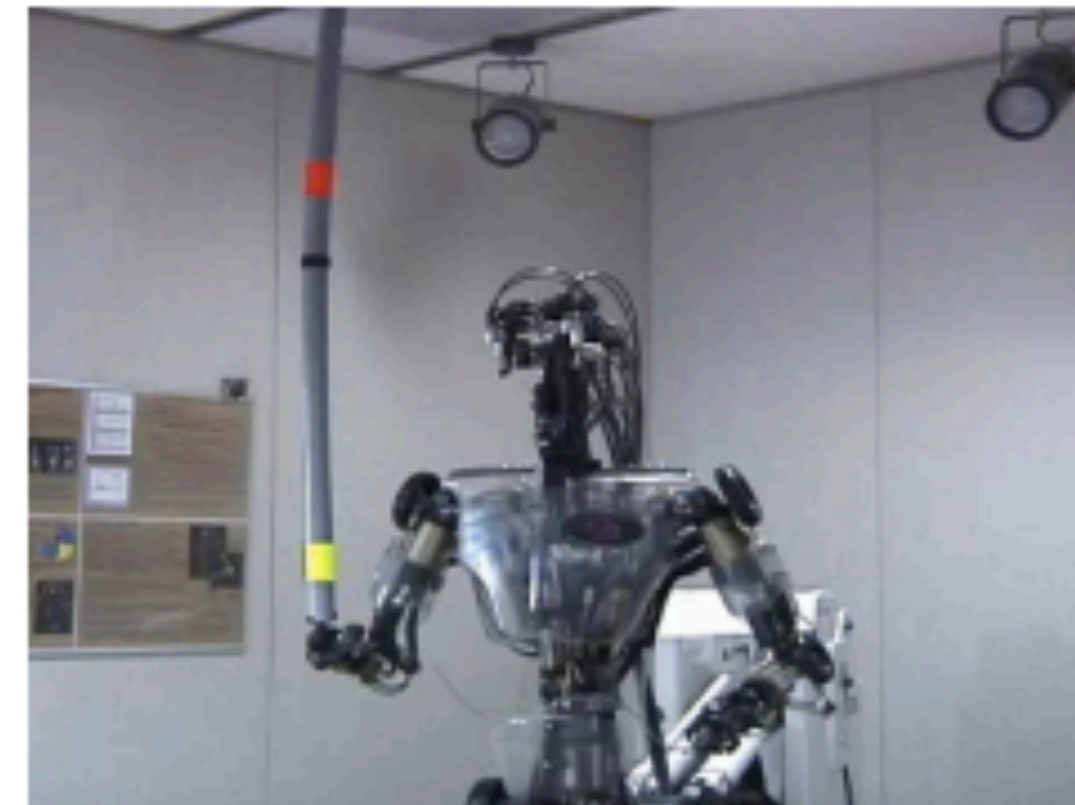
(b) Zebra Zero robot

2001



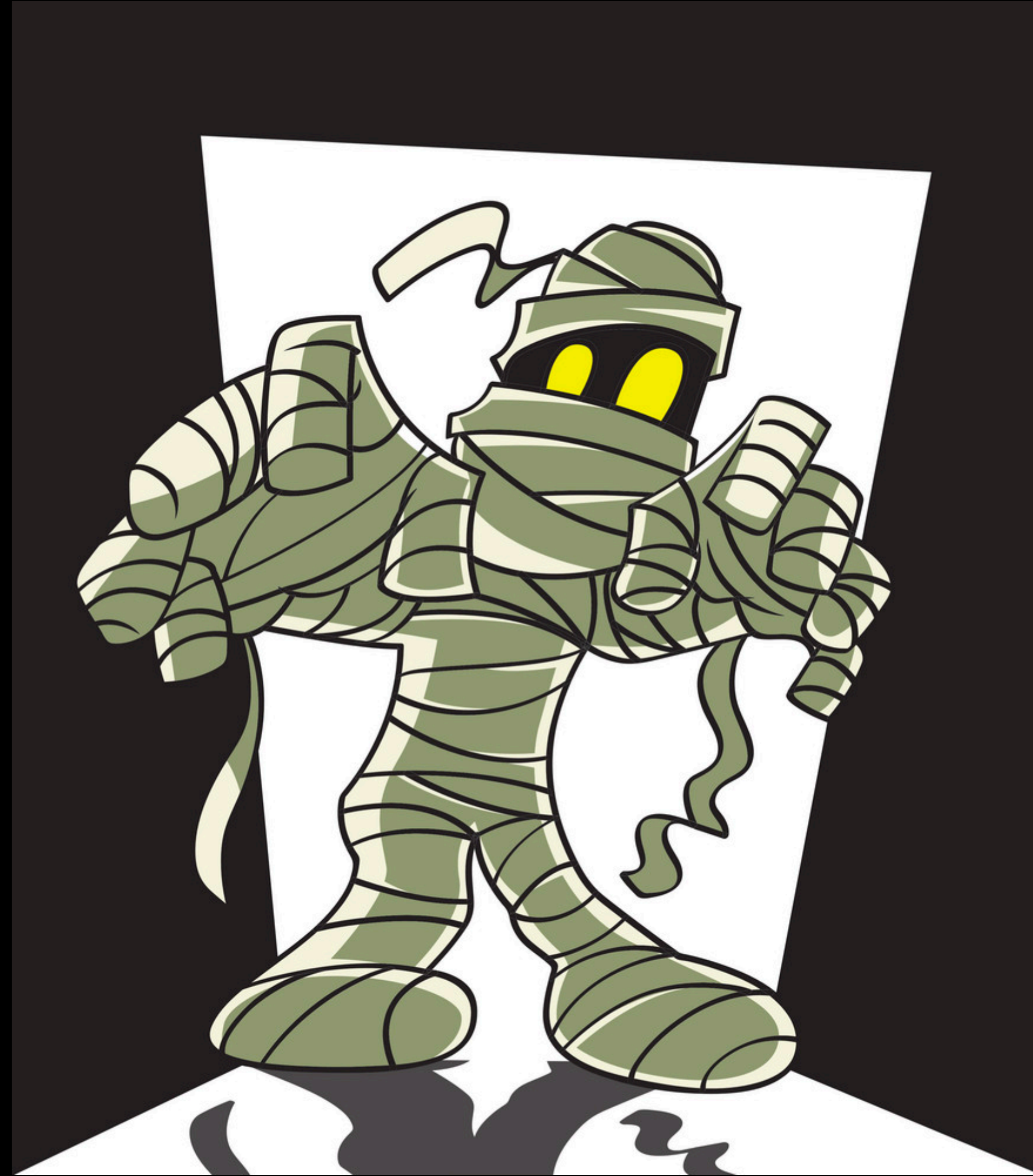
(c) Autonomous helicopter

1996

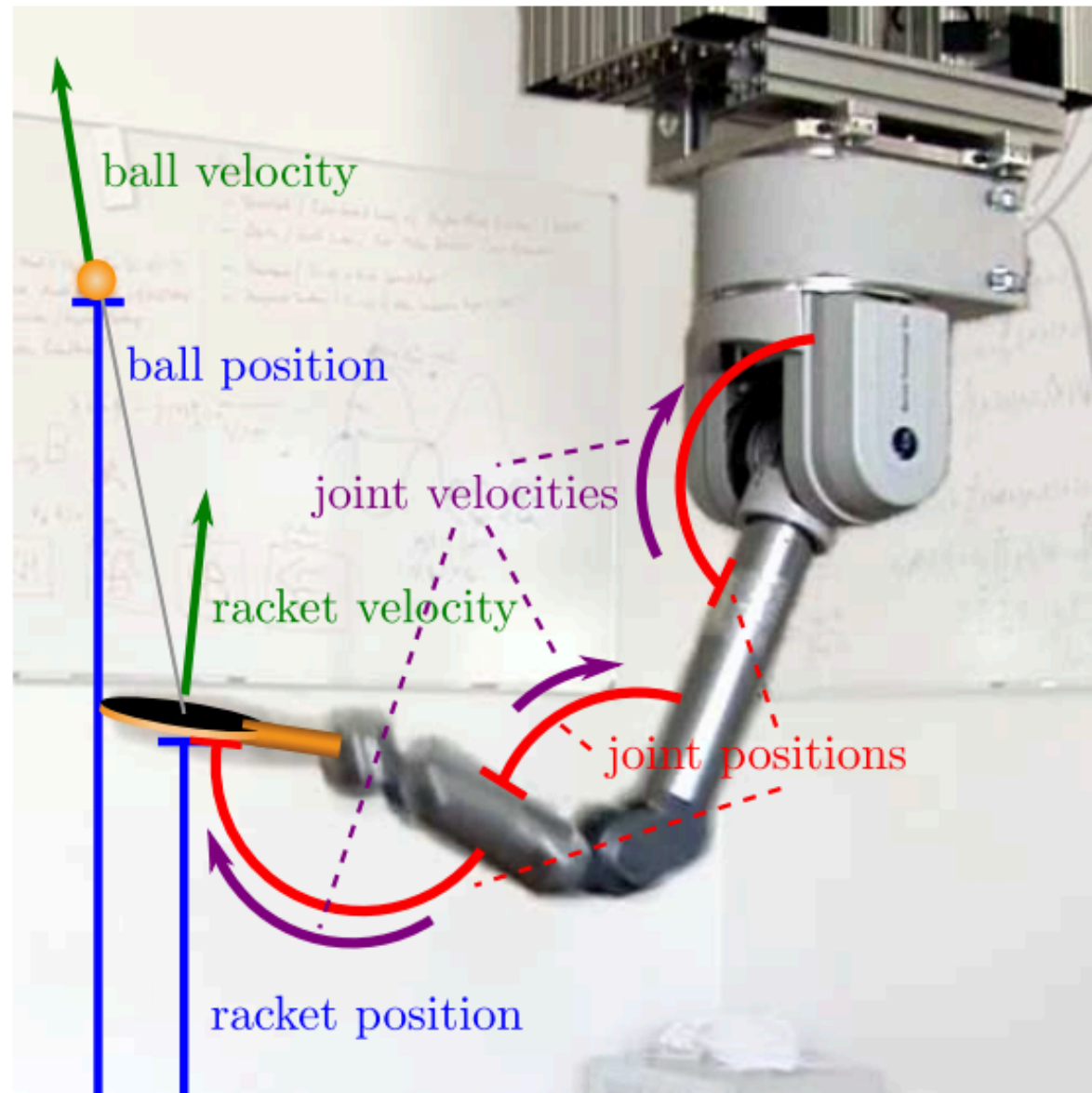


(d) Sarcos humanoid DB

# CURSES IN ROBOTICS



# Curses in RL for Robotics



Curse of Real-World Samples

Curse of Reward Specification

Curse of Dimensionality

Curse of Model Uncertainty

# No Silver Bullet!



# Ingredients for Practical RL



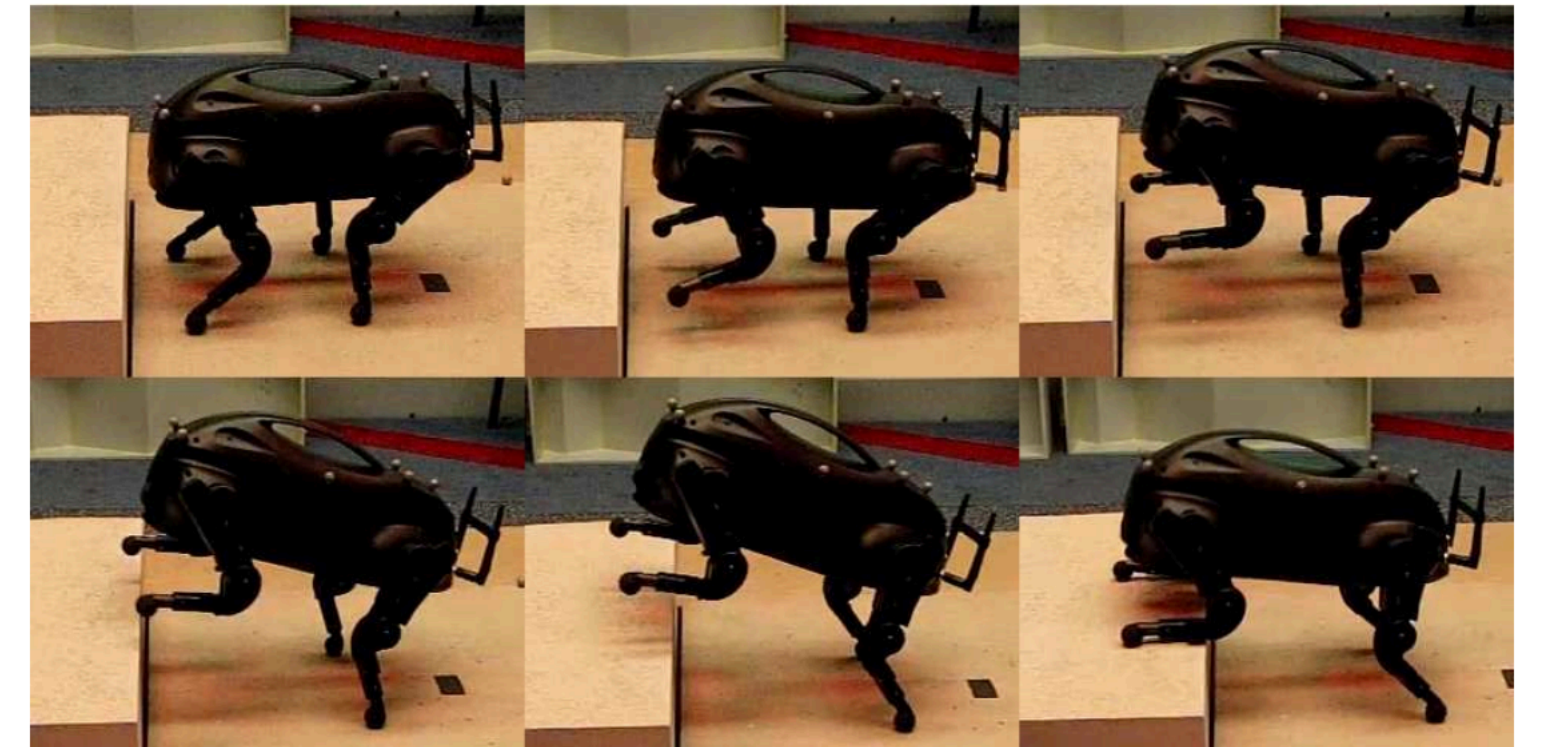
## Leveraging Demonstrations

Kakade et al, Bagnell and Scheider



## Reward Shaping

Ng et al.



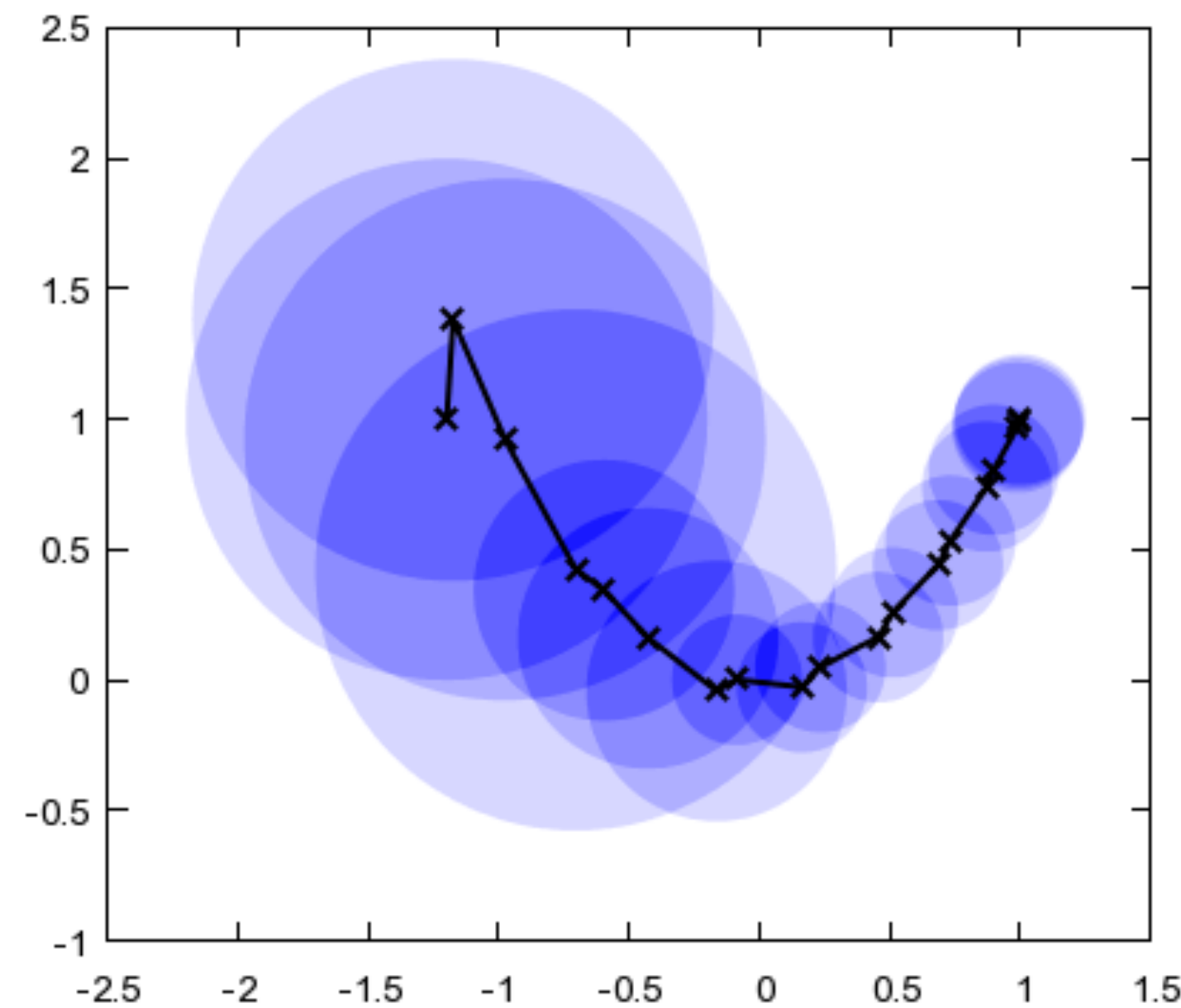
## Appropriate Policy Representation

Locally linear (Kolter and Ng),  
Dynamic Motor Primitives (Schaal)



## Plannable Models

Abbeel and Ng



## Conservatism and Trust

Kakade and Langford, Schulman et al.



What about Deep RL?





**Position 2**

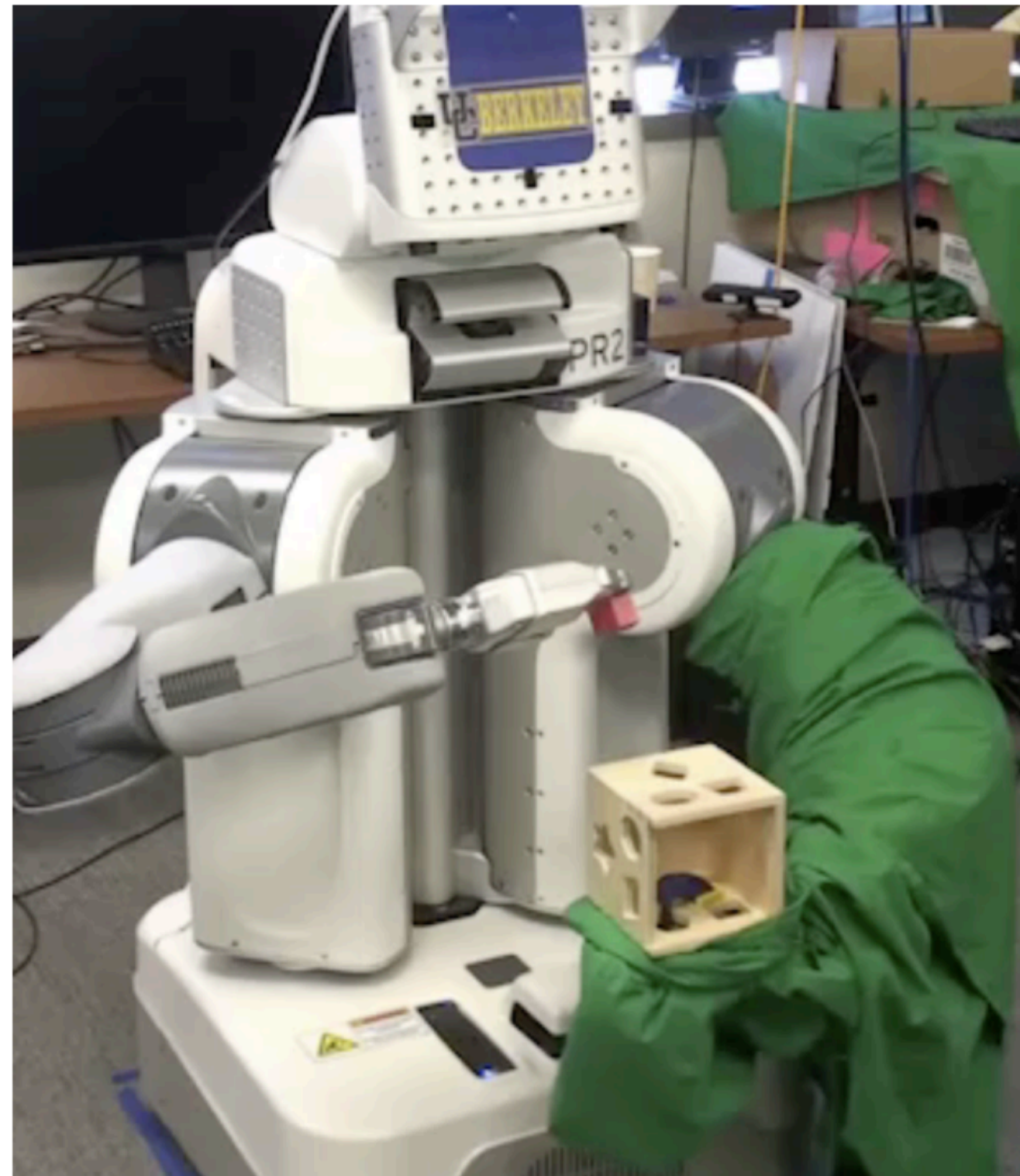
real time

autonomous execution

# How to Train Your Robot with Deep Reinforcement Learning – Lessons We've Learned

Julian Ibarz<sup>1</sup>, Jie Tan<sup>1</sup>, Chelsea Finn<sup>1,3</sup>, Mrinal Kalakrishnan<sup>2</sup>, Peter Pastor<sup>2</sup>, Sergey Levine<sup>1,4</sup>

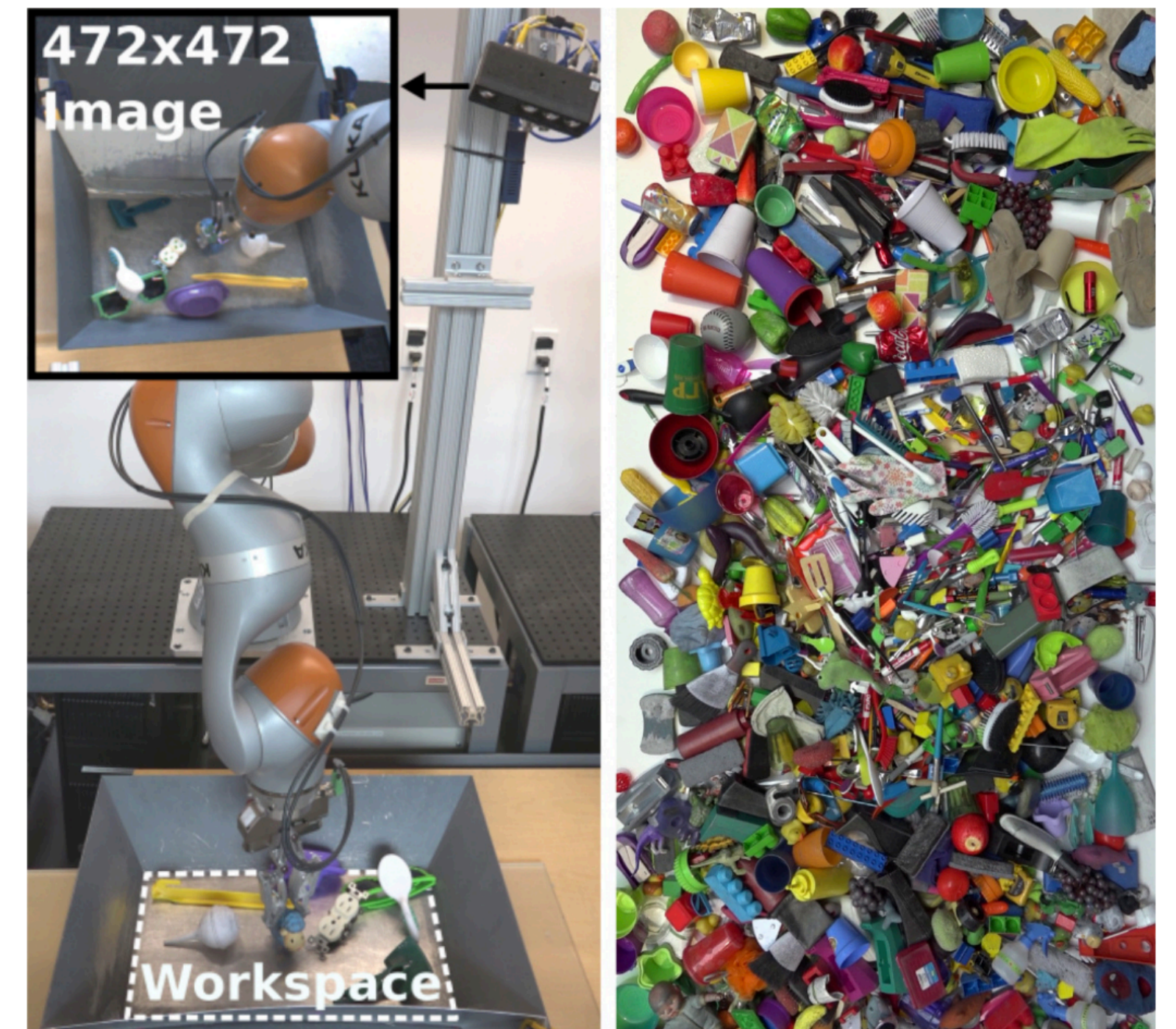
1–22  
©The Author(s) 2020  
Reprints and permission:  
sagepub.co.uk/journalsPermissions.nav  
DOI: 10.1177/ToBeAssigned  
www.sagepub.com/  
SAGE



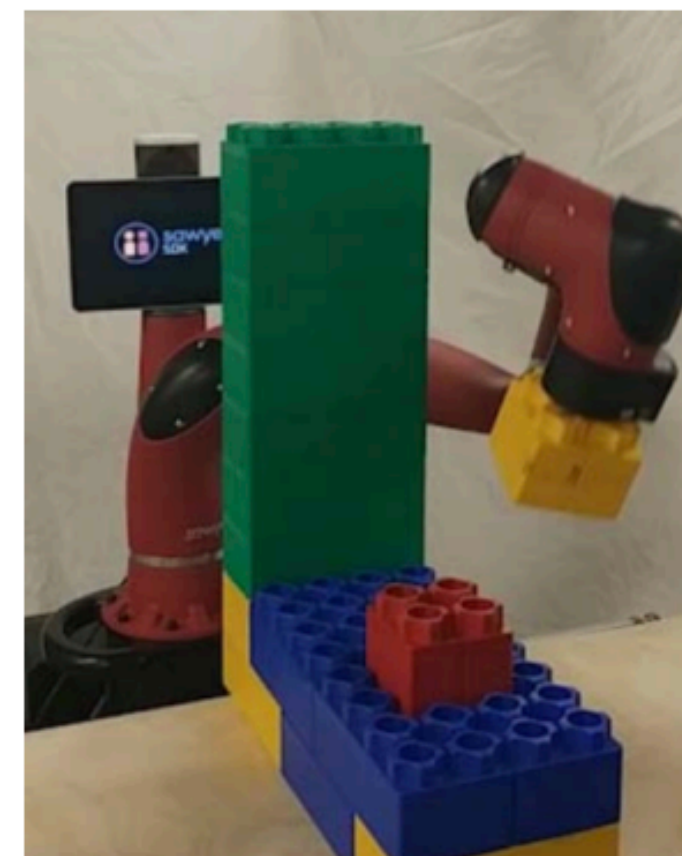
Levine et al 2016



Harnoja et al. 2019

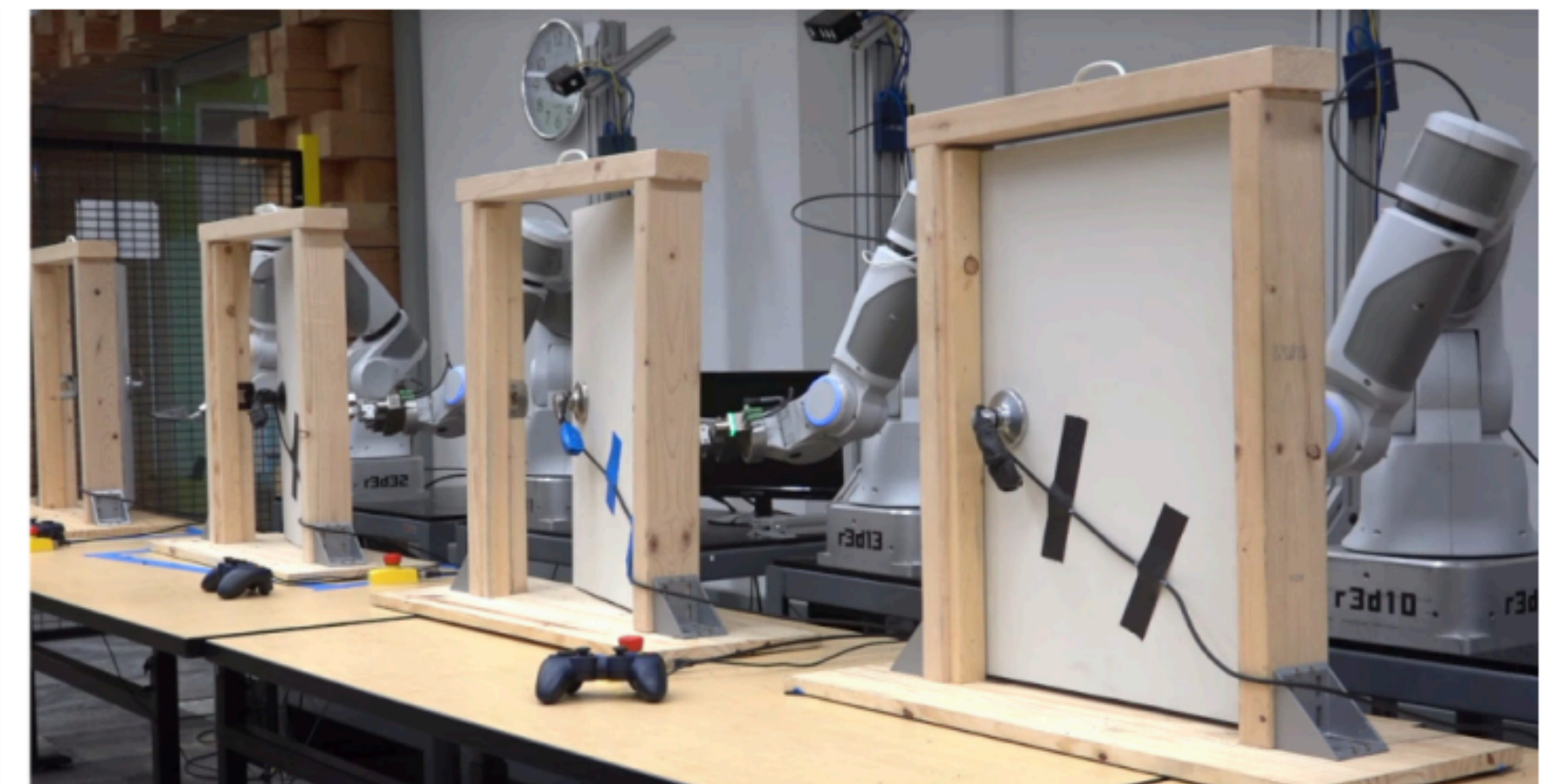


Kalashnikov et al. 2018



(a) block stacking

Harnoja et al 2018



(b) door opening

Gu et al 2017



*There is  
no Deep RL*

# There is no “Deep RL”

There is  
Deep Learning

There is  
Reinforcement Learning



Better representations for state / value function