

# Two-Player Games

**Lucas Janson**

**CS/Stat 184(0): Introduction to Reinforcement Learning**  
**Fall 2024**

# Today

- Feedback from last lecture
- Recap
- Game Playing: AlphaBeta Search/Rule Based Systems
- MCTS
- AlphaZero and Self-Play

# Feedback from feedback forms

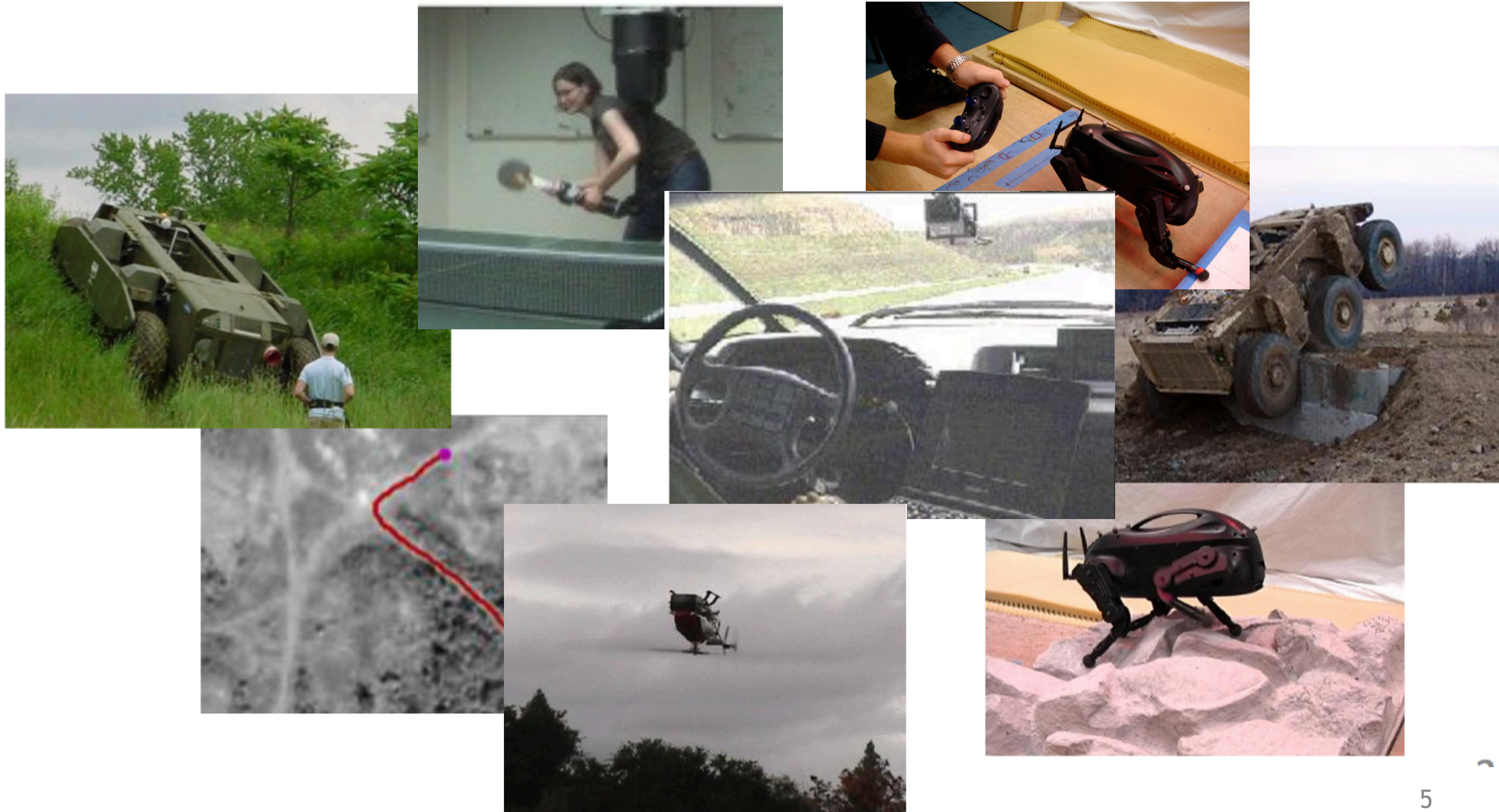
1. Thank you to everyone who filled out the forms!

# Today

- ✓ • Feedback from last lecture
- Recap
- Game Playing: AlphaBeta Search/Rule Based Systems
- MCTS
- AlphaZero and Self-Play

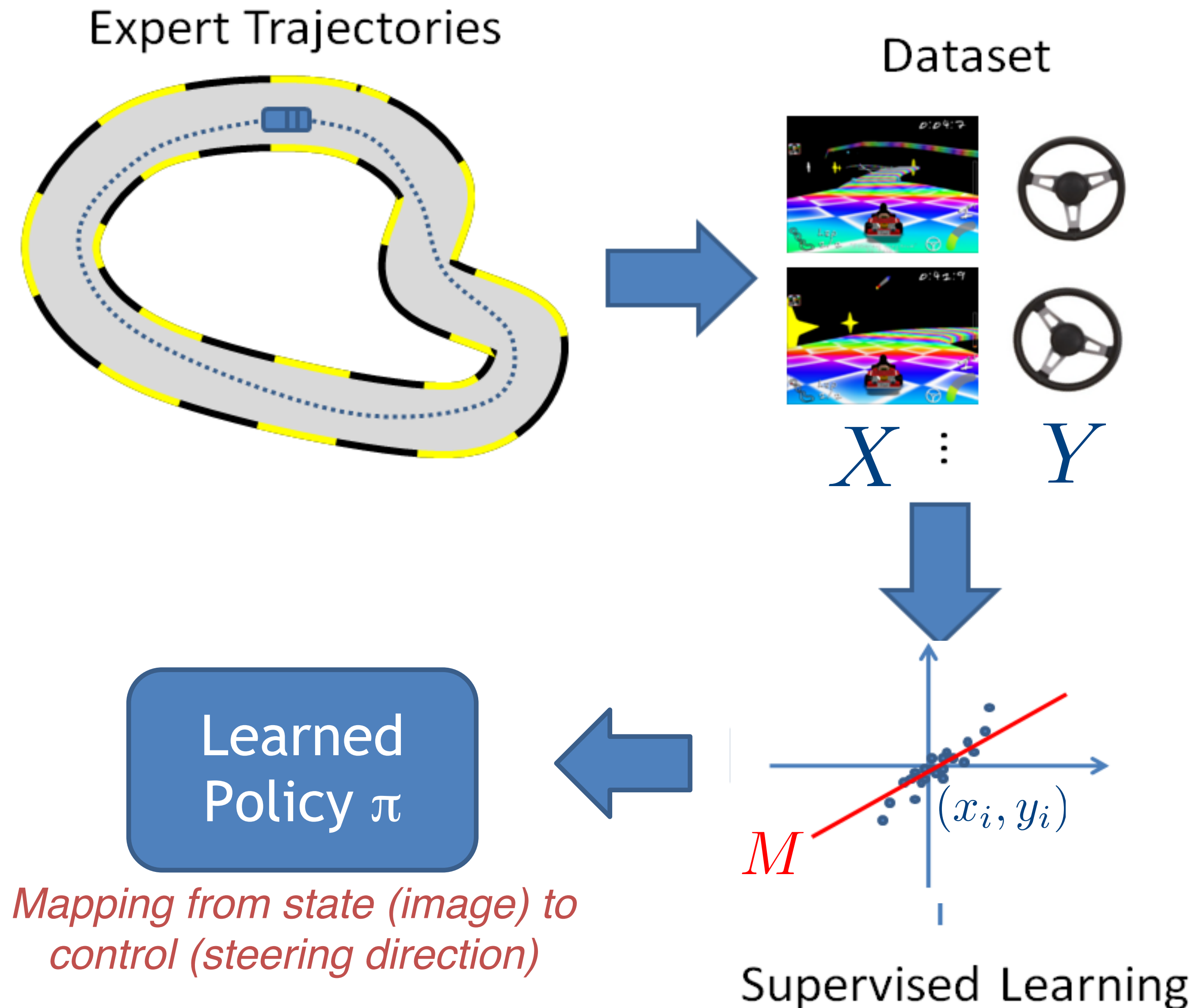


# Imitation Learning



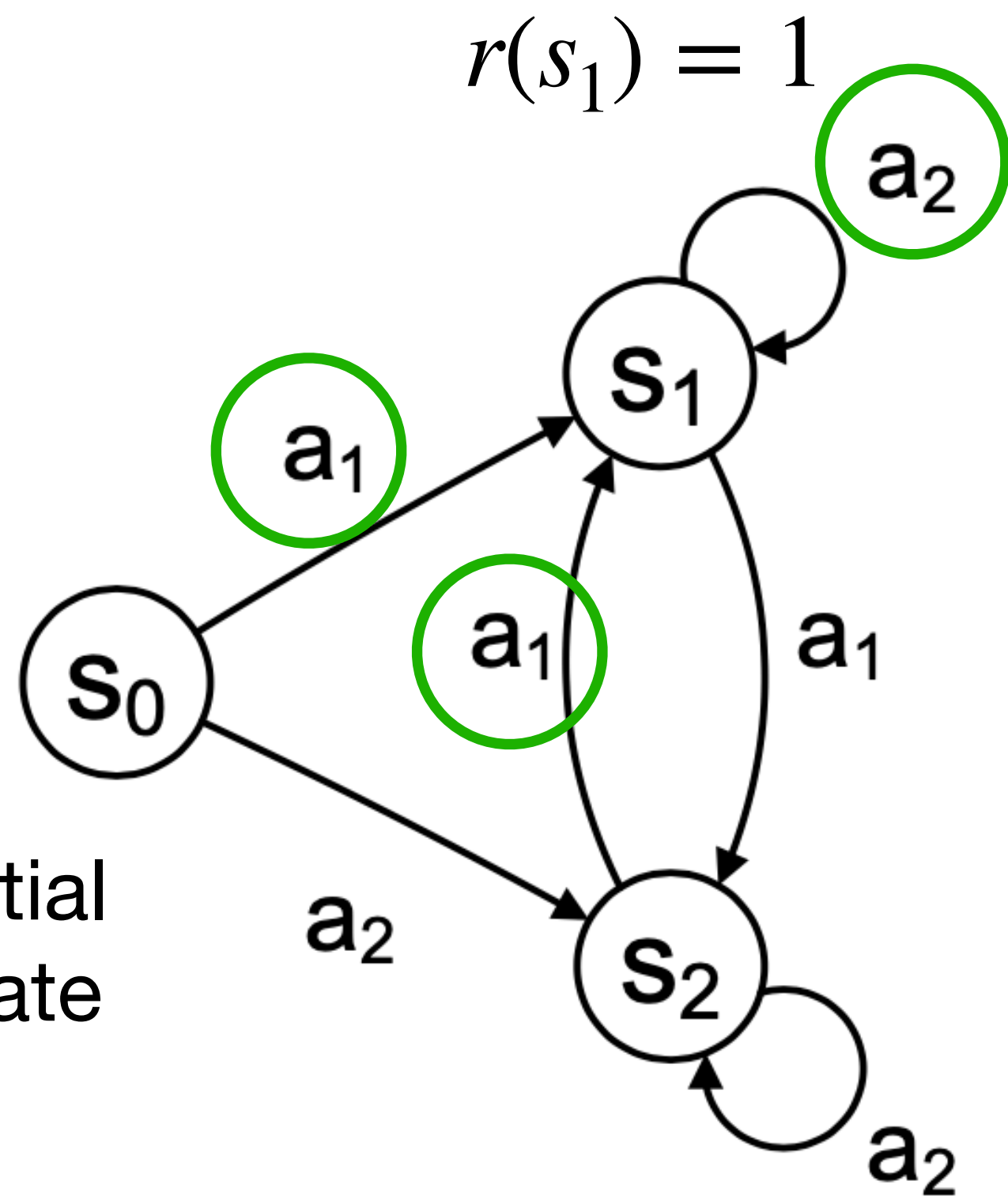


# Supervised Learning Approach: Behavior Cloning



*Mapping from state (image) to control (steering direction)*

# Distribution Shift Example ( $|V^{\pi^*} - V^{\hat{\pi}}| \leq H^2\epsilon$ )



Assume SL returns the policy  $\hat{\pi}$ :

$$\hat{\pi}(s_0) = \begin{cases} a_1 & \text{w/ prob } 1 - H\epsilon \\ a_2 & \text{w/ prob } H\epsilon \end{cases}, \quad \hat{\pi}(s_1) = a_2, \quad \hat{\pi}(s_2) = a_2$$

This policy has good supervised learning error:

$$\mathbb{E}_{\tau \sim \rho_{\pi^*}} \left[ \frac{1}{H} \sum_{h=0}^{H-1} \mathbf{1} [\hat{\pi}(s_h) \neq \pi^*(s_h)] \right] = \epsilon$$

note: while  $\hat{\pi}(s_2) \neq \pi^*(s_2)$ , state  $s_2$  is never visited under  $\pi^*$

We have **quadratic degradation** (in  $H$ ):

$$V_0^{\hat{\pi}}(s_0) = (1 - H\epsilon) \cdot V_0^{\pi^*}(s_0) + H\epsilon \cdot 0 = V_0^{\pi^*}(s_0) - \epsilon H(H - 1)$$

**Intuition:** once we make a mistake at  $s_0$ , we end up in  $s_2$  which is not in the training data!

Opt policy:  $\pi^*(s_0) = \pi^*(s_2) = a_1,$

$$\pi^*(s_1) = a_2$$

Under  $\rho_{\pi^*}$ , trajectory is  $s_0, s_1, s_1, \dots$

$$\rho_{\pi^*}(s_h = s_2) = 0$$

$$V_0^{\pi^*}(s_0) = H - 1$$

# The DAgger algorithm

Initialize  $\pi^0$ , and dataset  $\mathcal{D} = \emptyset$

For  $t = 0 \rightarrow T - 1$ :

1. W/  $\pi^t$ , generate dataset of trajectories  $\mathcal{D}^t = \{\tau_1, \tau_2, \dots\}$

where for all trajectories  $s_h \sim \rho_{\pi^t}$ ,  $a_h = \pi^*(s_h)$

2. **Data aggregation:**  $\mathcal{D} = \mathcal{D} \cup \mathcal{D}^t$

3. **Update policy via Supervised-Learning:**  $\pi^{t+1} = \text{SL}(\mathcal{D})$

In practice, the DAgger algorithm requires less human labeled data than BC.

[\[Informal Theorem\]](#) Under more assumptions + assuming  $\epsilon$  SL error is achievable, the DAgger algorithm has error:  $|V^{\pi^*} - V^{\hat{\pi}}| \leq H\epsilon$



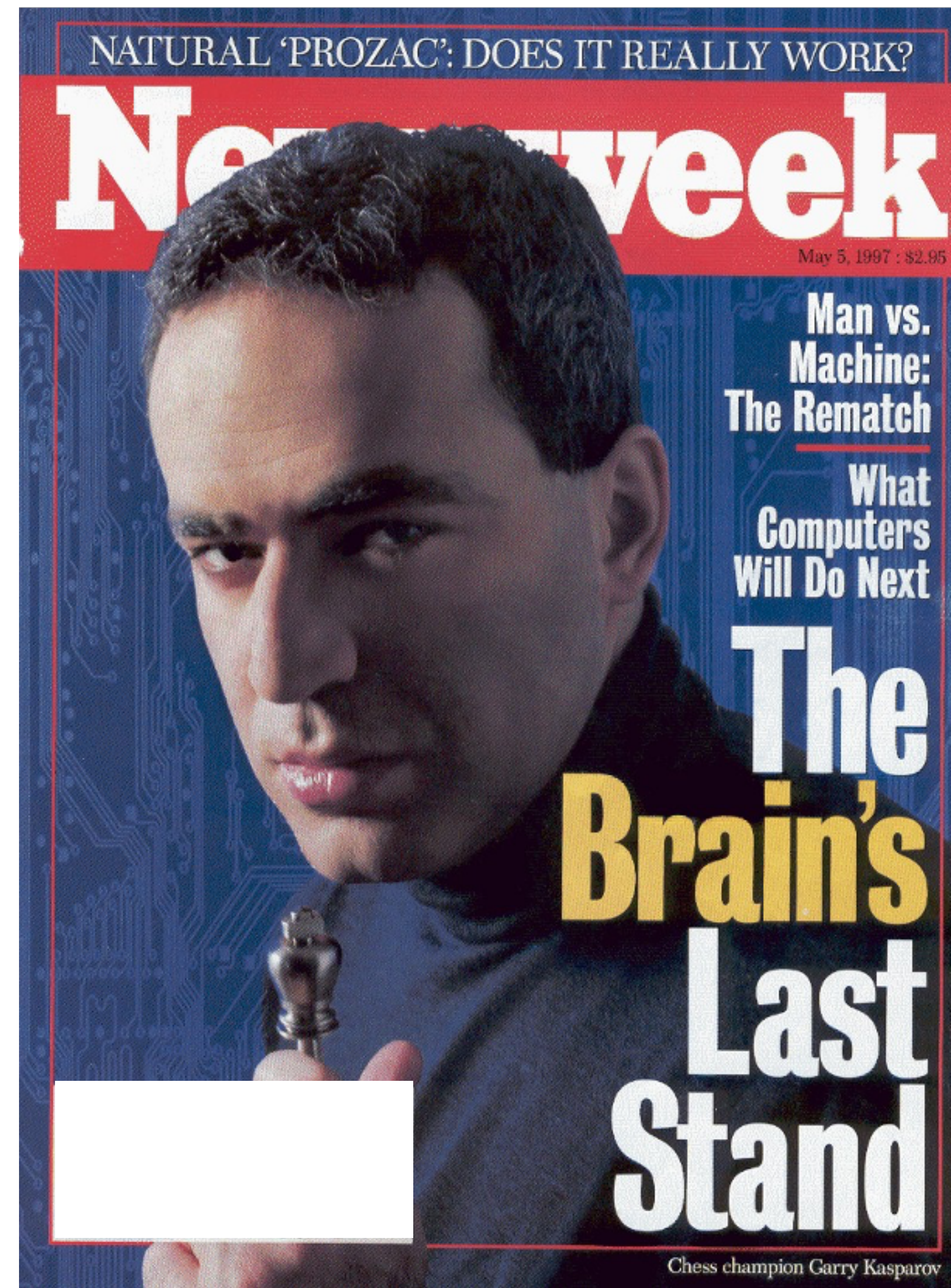
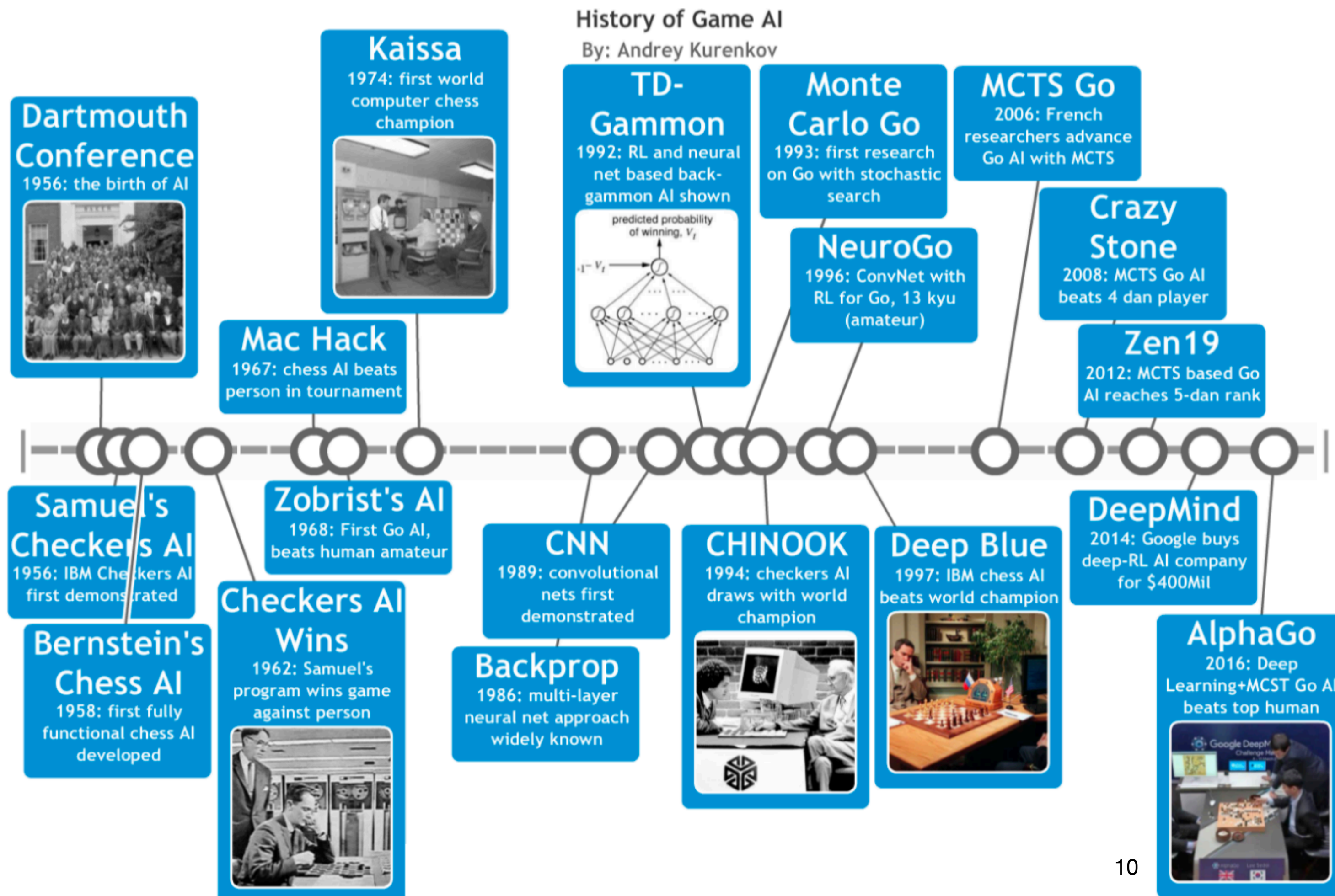
# Today

- ✓ • Feedback from last lecture
- ✓ • Recap
  - Game Playing: AlphaBeta Search/Rule Based Systems
  - MCTS
  - AlphaZero and Self-Play



# Fascination with AI and Games...

- **DeepBlue v. Kasparov (1997)**
  - winning in chess wasn't a good indicator of "progress in AI"





# Two-player, deterministic games

We will focus on games that are:

- deterministic
- two-player (alternating turns)
- zero sum (one player wins and the other loses)
- fully observable (by both players)
- stationary (only game state and whose turn it is matters)

E.g.,

- Tic-tac-toe
- Chess
- Go

Notation:

- Game states  $S$ , initial state  $s_0 \in S$
- Set of actions available in state  $s$ :  $A(s)$
- Dynamics  $P(s, a) \in S$
- Maximum game length  $H$
- Score at terminal state  $r(s)$  (sign determines winner)

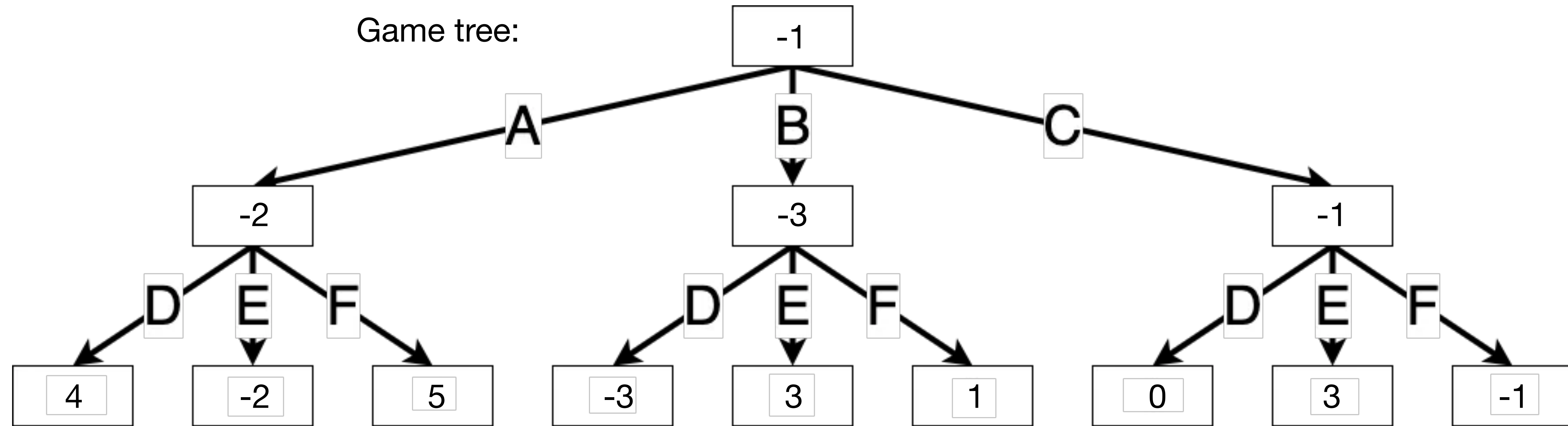
Still an MDP, but two competing players make it a bit different than earlier RL setup

# Simple example of min-max search

$H = 2$ , player 1 takes action A, B, or C  
then player 2 takes action D, E, F

Outcome $r(s)$	D	E	F
A	4	-2	5
B	-3	3	1
C	0	3	-1

Game tree:



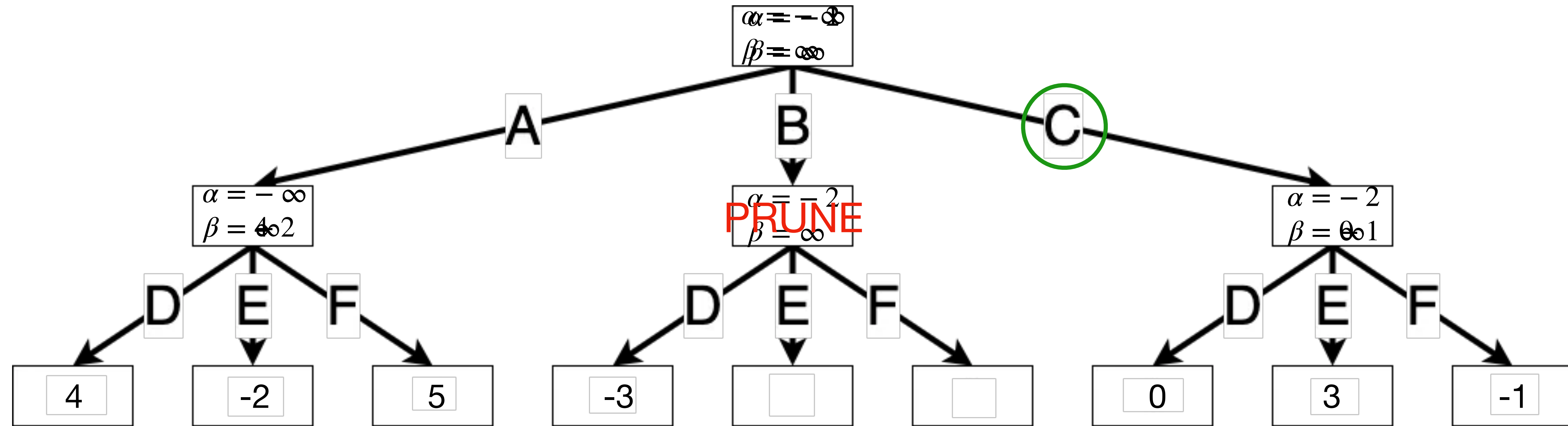
Basically dynamic programming! Numbers in boxes are value function  $V(s)$



# Alpha-beta search

Pruning can speed up search without losing exactness

- $\alpha(s)$  is lower-bound for  $V^*(s)$
- $\beta(s)$  is upper-bound for  $V^*(s)$
- Bounds sometimes allow pruning



The order that actions are considered can matter a lot

# Today

- ✓ • Feedback from last lecture
- ✓ • Recap
- ✓ • Game Playing: AlphaBeta Search/Rule Based Systems
  - MCTS
  - AlphaZero and Self-Play

# Monte Carlo Tree Search (MCTS)

For now, assume game outcome just win or lose:  $r(s) \in \{-1, 1\}$

- Alpha-beta search evaluates non-leaf nodes via a min-max approach
  - Even with pruning, requires searching **a LOT** of paths down tree
- Idea of MCTS is evaluate non-leaf nodes via **sampling** (Monte Carlo)
- High-level: at each iteration, **MCTS** does the following
  - Defines a game-playing strategy (policy for both players) that is a simple function of a set of statistics computed from existing samples
  - Plays the game to completion via this strategy and records outcome
  - Updates statistics used to define game-playing strategy
- Strategy gradually improves with more iterations/samples, so **can fit in any computational budget**
- Samples are concentrated around more promising strategies

# “Pure” MCTS Algorithm

**Input:** game state (“root node”  $R$ ), #iterations  $N$ , exploration constant  $C$

For iteration  $t = 1, \dots, N$

1. **Obtain the  $t$ -th sample trajectory:** Starting at  $R$ , while current state  $s \notin \{\text{win, lose}\}$

a. For player  $X \in \{0, 1\}$ , at current state  $s$ , let  $s' = P(s, a)$  and define:

$$\text{UCBscore}_t(s, a) = \frac{\text{\#wins for player } X \text{ from } s'}{\text{\#visits to } s'} + C \sqrt{\frac{\log(\text{\#visits to } s)}{\text{\#visits to } s'}}$$

b. “Take” action:

$$\hat{a} = \arg \max_a \text{UCBscore}(s, a)$$

2. **Update stats:** For all visited states  $s'$  in this trajectory,

c. update visit counts:

$$[\text{\#visits to } s'] = [\text{\#visits to } s'] + 1$$

d. for winner  $X$  and if  $s$  was visited by  $X$ :

$$[\text{\#wins for } X \text{ at } s'] = [\text{\#wins for } X \text{ at } s'] + 1$$

(data structure: only need to keep track of stats at visited states)

**Output:** return the action  $\hat{a} = \arg \max_a \text{UCBscore}_N(R, a)$



# Improving MCTS

- MCTS re-runs at every game step (root node gets updated to current state)
- “Pure” MCTS can work well for small games, but what can go wrong?
- For large games, most states never visited...  
so UCB basically just samples trajectories **randomly** after a certain point!
- **Solution:**
  - Fix a strategy  $\pi$  and a look-ahead horizon  $T$
  - Only use UCB strategy for choosing actions for  $T$  steps, use  $\pi$  after
  - Note since MCTS re-runs at every game step,  $\pi$ 's use gets later and later
- Need a good strategy  $\pi$ ... or, a good value function approximation  $\hat{V}(s)$ :  
After  $T$  steps, instead of using  $\pi$ , stop and record  $\hat{V}(s)$  as game outcome
- $\hat{V}(s)$  could be learned from offline/expert data and improved online

# Today

- ✓ • Feedback from last lecture
- ✓ • Recap
- ✓ • Game Playing: AlphaBeta Search/Rule Based Systems
- ✓ • MCTS
  - AlphaZero and Self-Play

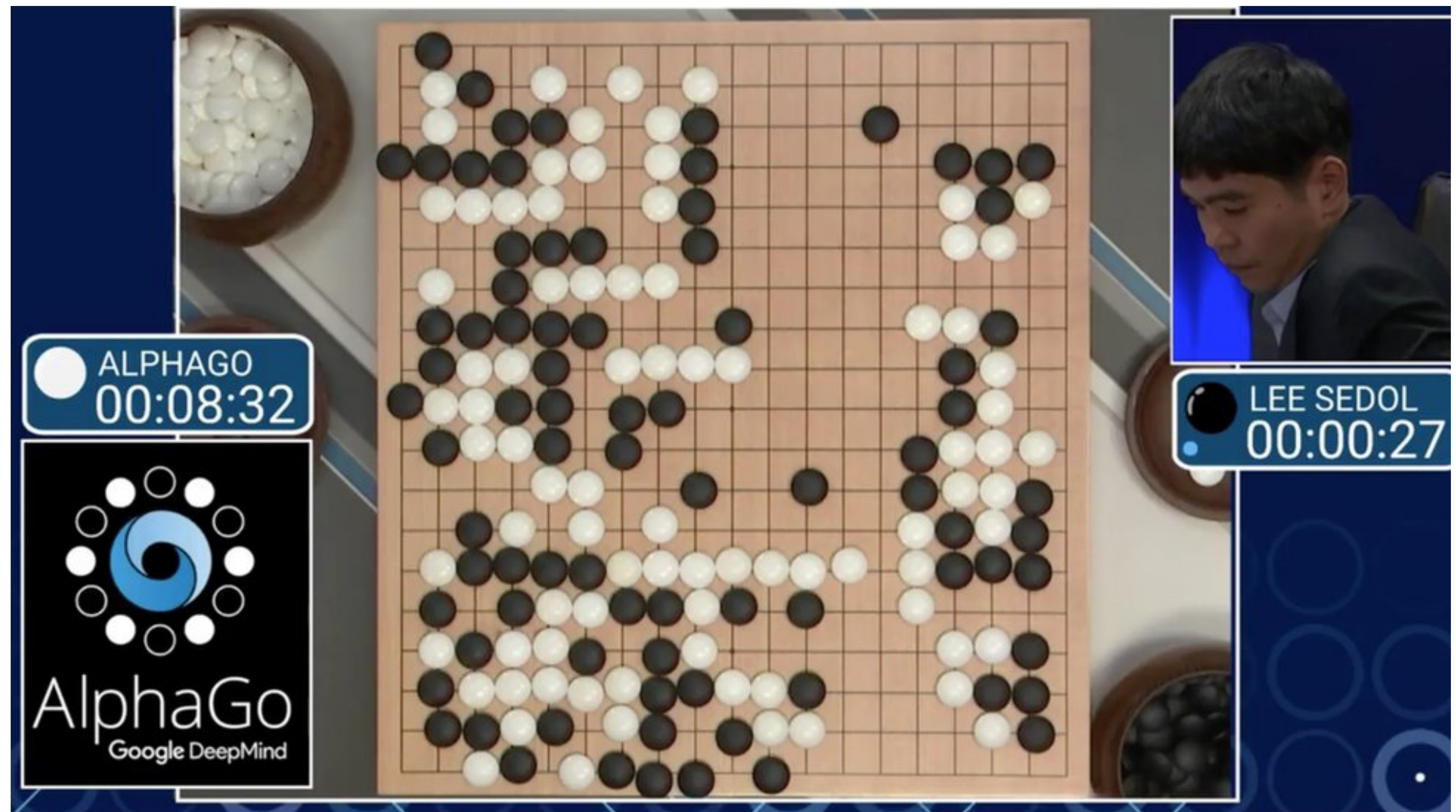


# AlphaGo

## AlphaGo versus Lee Sedol 4-1

Seoul, South Korea, 9-15 March 2016

Game one	AlphaGo W+R
Game two	AlphaGo B+R
Game three	AlphaGo W+R
Game four	Lee Sedol W+R
Game five	AlphaGo W+R



- Lots of moving parts:
  - **Imitation Learning:** first, the algo estimates the values from historical games.
  - It then uses an **MCTS-stye lookahead** with **learned value functions**.
- **AlphaZero** (2017) is a simpler more successful approach that uses self-play



# AlphaZero

- AlphaZero: MCTS + DeepLearning + self-play
  - MCTS subroutine has a value network and policy network
    - a **value network** estimating the value for the state of the board  $\hat{V}_\theta(s)$
    - A **policy network**  $\pi_\theta(a | s)$  that is a probability vector over all possible actions
  - Use these for MCTS, then play agent against self and use self-play data to learn better  $\theta$ ; iterate



# AlphaZero MCTS subroutine (without self-play)

**Input:** game state (“root node”  $R$ ), #iterations  $N$ , exploration constant  $C$ , look-ahead horizon  $T$ , **value network**  $\hat{V}_\theta(s)$ , **policy network**  $\pi_\theta(a | s)$

For iteration  $t = 1 : N$

1. **Obtain the  $t$ -th sample trajectory:** For  $T$  steps starting from  $R$ ,

a. For player  $X \in \{0,1\}$ , at current state  $s$ , define  $s' = P(s, a)$  and define:

$$\text{UCBscore}_t(s, a) = \bar{\hat{V}}(s') \cdot (-1)^X + C \cdot \pi_\theta(a | s) \cdot \sqrt{\frac{\log(\#\text{visits to } s)}{\#\text{visits to } s'}}$$

b. “Take” action:

$$\hat{a} = \arg \max_a \text{UCBscore}_t(s, a)$$

2. **Update stats:** For all visited states  $s'$  in this “roll-out”, letting  $s_T$  be the last sampled state

c. Update counts:  $[\#\text{visits to } s'] = [\#\text{visits to } s'] + 1$

d. Update average value estimate:  $\bar{\hat{V}}(s') \leftarrow \frac{[\#\text{visits to } s']}{[\#\text{visits to } s'] + 1} \bar{\hat{V}}(s') + \frac{1}{[\#\text{visits to } s'] + 1} \hat{V}_\theta(s_T)$

**Output:** return the action  $\hat{a} = \arg \max_a \text{UCBscore}_N(R, a)$

# Self-play

- Iterate the following:
  - **Self-play:** Play against self  $M$  times using current MCTS strategy
  - **Supervised Learning:** Use  $M$  self-play game trajectories to update:
    - $\hat{V}_\theta$  with squared error loss wrt game outcomes (similar to in fitted VI or baseline estimation)
    - $\pi_\theta$  with negative log likelihood loss wrt actions taken in game (similar to in BC)
    - In practice, combine loss functions into single SL problem with shared  $\theta$
- AlphaZero uses no historical data, only self-play
- Performance improvement was pretty astronomical!



## Chess [\[ edit \]](#)

In AlphaZero's chess match against Stockfish 8 (2016 [TCEC](#) world champion), each program was given one minute per move. Stockfish was allocated 64 threads and a [hash](#) size of 1 GB,<sup>[1]</sup> a setting that Stockfish's [Tord Romstad](#) later criticized as suboptimal.<sup>[7][note 1]</sup> AlphaZero was trained on chess for a total of nine hours before the match. During the match, AlphaZero ran on a single machine with four application-specific [TPUs](#). In 100 games from the normal starting position, AlphaZero won 25 games as White, won 3 as Black, and drew the remaining 72.<sup>[8]</sup> In a series of twelve, 100-game matches (of unspecified time or resource constraints) against Stockfish starting from the 12 most popular human openings, AlphaZero won 290, drew 886 and lost 24.<sup>[1]</sup>

## Shogi [\[ edit \]](#)

AlphaZero was trained on shogi for a total of two hours before the tournament. In 100 shogi games against elmo (World Computer Shogi Championship 27 summer 2017 tournament version with YaneuraOu 4.73 search), AlphaZero won 90 times, lost 8 times and drew twice.<sup>[8]</sup> As in the chess games, each program got one minute per move, and elmo was given 64 threads and a hash size of 1 GB.<sup>[1]</sup>

## Go [\[ edit \]](#)

After 34 hours of self-learning of Go and against AlphaGo Zero, AlphaZero won 60 games and lost 40.<sup>[1][8]</sup>

Cup

Event	Year	Time Controls	Result	Ref
Cup 1	2018	30+10	<b>1st</b>	<a href="#">[63]</a>
Cup 2	2019	30+5	<b>2nd</b> <sup>[note 1]</sup>	<a href="#">[64]</a>
Cup 3	2019	30+5	<b>2nd</b>	<a href="#">[65]</a>
Cup 4	2019	30+5	<b>1st</b>	<a href="#">[66]</a>
Cup 5	2020	30+5	<b>1st</b>	<a href="#">[67]</a>
Cup 6	2020	30+5	<b>3rd</b>	<a href="#">[68]</a>
Cup 7	2020	30+5	<b>1st</b>	<a href="#">[69]</a>
Cup 8	2021	30+5	<b>1st</b>	<a href="#">[70]</a>
Cup 9	2021	30+5	<b>1st</b>	<a href="#">[71]</a>
Cup 10	2022	30+3	<b>1st</b>	<a href="#">[72]</a>
Cup 11	2023	30+3	<b>2nd</b>	<a href="#">[73]</a>



# Today

- ✓ • Feedback from last lecture
- ✓ • Recap
- ✓ • Game Playing: AlphaBeta Search/Rule Based Systems
- ✓ • MCTS
- ✓ • AlphaZero and Self-Play

# Summary:

1. Search is powerful: MCTS
2. Search + learning is better: AlphaZero

Attendance:

[bit.ly/3RcTC9T](https://bit.ly/3RcTC9T)



Feedback:

[bit.ly/3RHtlxy](https://bit.ly/3RHtlxy)

