

# Lecture 13: Monte Carlo Tree Search

Emma Brunskill

CS234 Reinforcement Learning.

Winter 2026

- With some slides from or derived from David Silver

# Class Structure

- Last time: Fast / sample efficient Reinforcement Learning
- **This Time: MCTS**
- Next time: Rewards in RL

# AlphaZero and Monte Carlo Tree Search

- Responsible in part for one of the greatest achievements in AI in the last decade— becoming a better Go player than any human
- Incorporates a number of interesting ideas

# Table of Contents

1 Simulation-Based Search

2 AlphaZero

# Computing Action for Current State Only

- So far in class, compute a policy for whole state space
- Key idea: Use additional local computation to make a better decision for right now

# Simple Monte-Carlo Search

- Given a model  $\mathcal{M}_v$  and a **simulation policy**  $\pi$
- For each action  $a \in \mathcal{A}$ 
  - Simulate  $K$  episodes from current (real) state  $s_t$

$$\{s_t, a, R_{t+1}^k, \dots, S_T^k\}_{k=1}^K \sim \mathcal{M}_v, \pi$$

- Evaluate actions by mean return (**Monte-Carlo evaluation**)

$$Q(s_t, a) = \frac{1}{K} \sum_{k=1}^K G_t \xrightarrow{P} q_\pi(s_t, a) \quad (1)$$

- Select current (real) action with maximum value

$$a_t = \operatorname{argmax}_{a \in \mathcal{A}} Q(s_t, a)$$

- This is essentially doing 1 step of policy improvement

# Expectimax Tree

- Can we do better than 1 step of policy improvement?
- If have a MDP model  $\mathcal{M}_v$
- Can compute optimal  $Q(s, a)$  values for current state by constructing an **expectimax** tree



# Expectimax Tree

- Can we do better than 1 step of policy improvement?
- If have a MDP model  $\mathcal{M}_v$
- Can compute optimal  $q(s, a)$  values for current state by constructing an expectimax tree
- Limitations: Size of tree scales as  $(|S||A|)^H$

# Monte-Carlo **T**ree Search (MCTS)

- Given a model  $\mathcal{M}_v$
- Build a **search tree** rooted at the current state  $s_t$
- Samples actions and next states
- Iteratively construct and update tree by performing  $K$  simulation episodes starting from the root state
- After search is finished, select current (real) action with maximum value in search tree

$$a_t = \operatorname{argmax}_{a \in A} Q(s_t, a)$$

# Upper Confidence Tree (UCT) Search

- How to select what action to take during a simulated episode?
- UCT: borrow idea from bandit literature and treat each node where can select actions as a multi-armed bandit (MAB) problem
- Maintain an upper confidence bound over reward of each arm

# Upper Confidence Tree (UCT) Search

- How to select what action to take during a simulated episode?
- UCT: borrow idea from bandit literature and treat each **node** where can select actions as a multi-armed bandit (MAB) problem
- Maintain an upper confidence bound over reward of each arm at a node

$$Q(s, a, i) = \frac{1}{N(i, a)} \sum_{k=1}^{N(i, a)} G_k(i, a) + c \sqrt{\frac{O(\log N(i))}{N(i, a)}}$$

- where  $N(i, a)$  is the number of times selected arm  $a$  at node  $i$ ,  $G_k(i, a)$  is the  $k$ -th return (discounted sum of rewards) from node  $i$  following action  $a$ , and
- For simulated episode  $k$  at node  $i$ , select action/arm with highest upper bound to simulate and expand (or evaluate) in the tree

$$a_{ik} = \arg \max Q(s, a, i)$$

- This implies that the policy used to simulate episodes with (and expand/update the tree) can change across each episode

# Advantages of MC Tree Search

- Highly selective best-first search
- Evaluates states dynamically
- Uses sampling to break curse of dimensionality
- Works for “black-box” models (only requires samples)
- Computationally efficient, anytime, parallelisable

# Table of Contents

1 Simulation-Based Search

2 AlphaZero

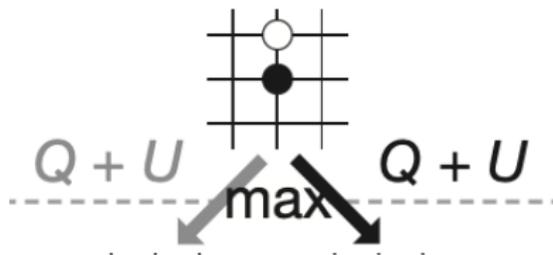
# Case Study: the Game of Go

- Go is 2500 years old
- Hardest classic board game
- Grand challenge task (John McCarthy)
- Traditional game-tree search has failed in Go
- Trailer: [▶ AlphaGo trailer link](#)
- Check your understanding: does playing Go involve learning to make decisions in a world where dynamics and reward model are unknown?



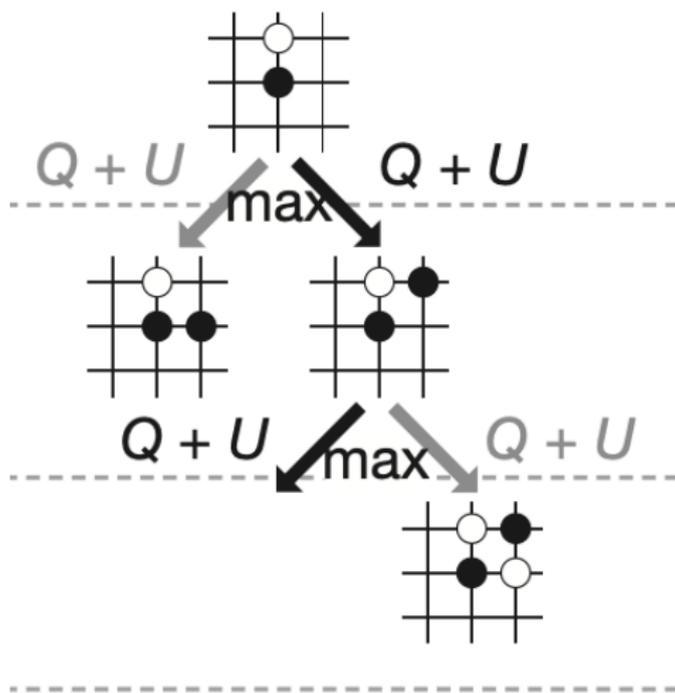
# Selecting a Move in a Single Game: Start at Root<sup>1</sup>

- Inspired by Upper Confidence Tree Search but many changes

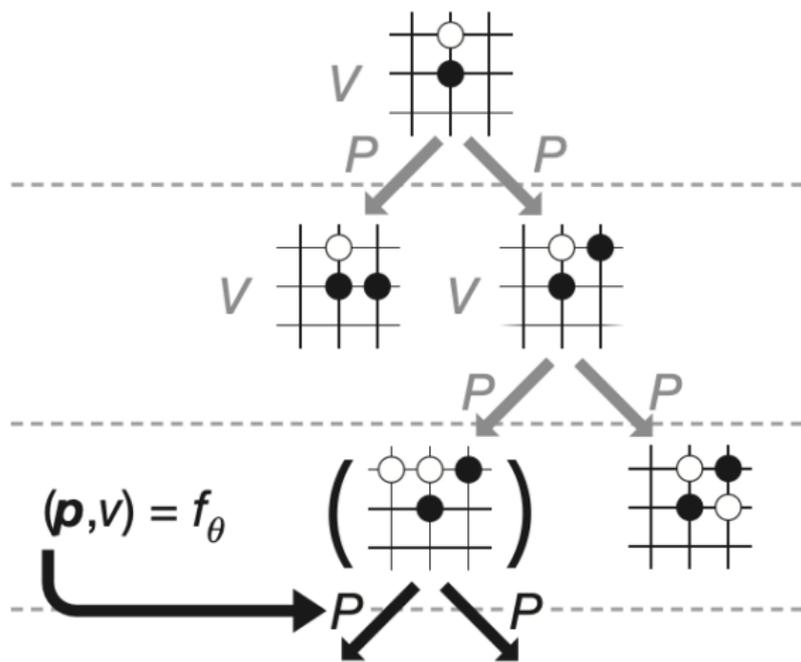


<sup>1</sup>Images from Silver et al. Nature 2017

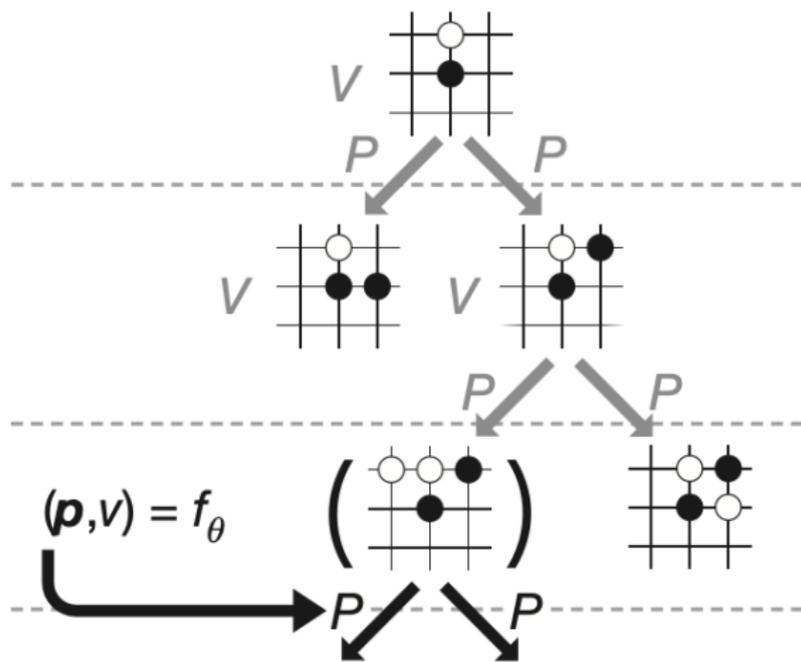
# Selecting a Move in a Single Game: Repeatedly Expand<sup>2</sup>



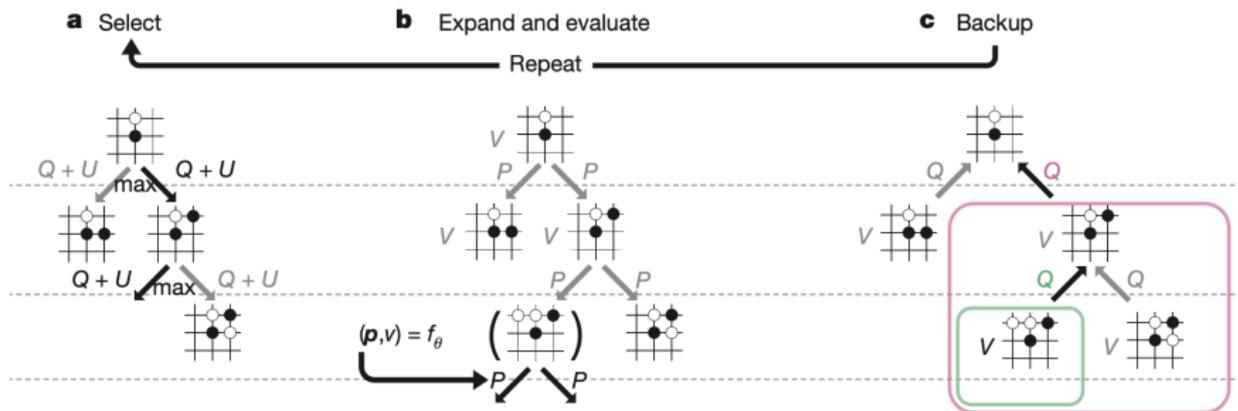
# Selecting a Move in a Single Game: Note Using Network Predictions for Action Probabilities<sup>3</sup>



# Selecting a Move in a Single Game: At Leaf, Plug in Network Predictions for Value<sup>4</sup>

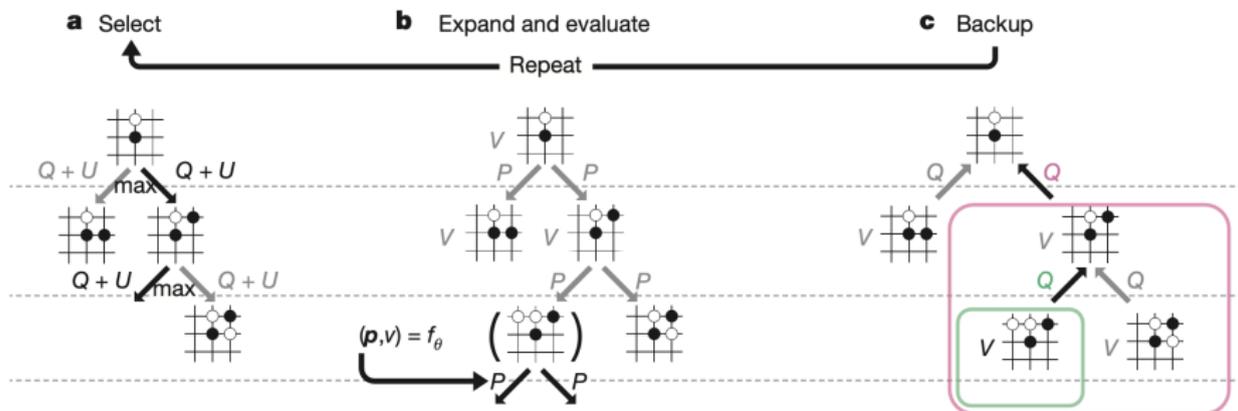


# Selecting a Move in a Single Game: Update Ancestors<sup>5</sup>



<sup>5</sup>Images from Silver et al. Nature 2017

# Selecting a Move in a Single Game: Repeat Many Times<sup>6</sup>

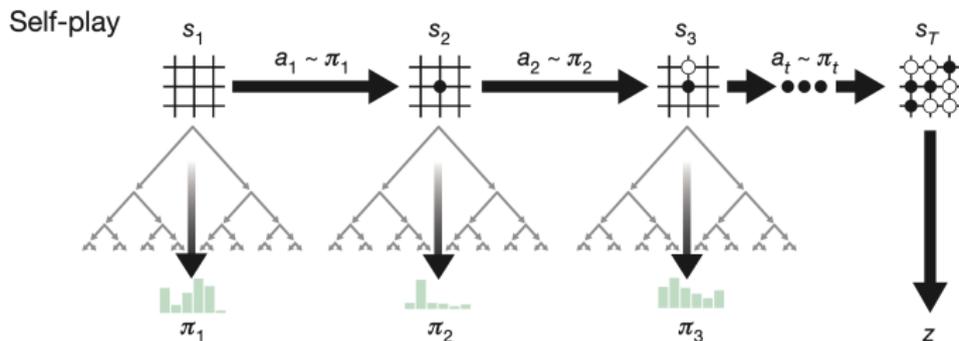


- Repeat roll out and backup process many times
- Note: inside the network alternating whether opponent or agent is "maximizing" its value. Therefore tree is mimicking a min-max tree
- At end, compute a policy for root node by

$$\pi(s) \propto N(s, a)^{\frac{1}{\tau}} \quad (2)$$

<sup>6</sup>Images from Silver et al. Nature 2017

# Self Play a Game<sup>7</sup>



- Select an action according to root policy, take action, and repeat whole process
- Repeat until game ends\* and observe a win or loss

<sup>7</sup>Images from Silver et al. Nature 2017

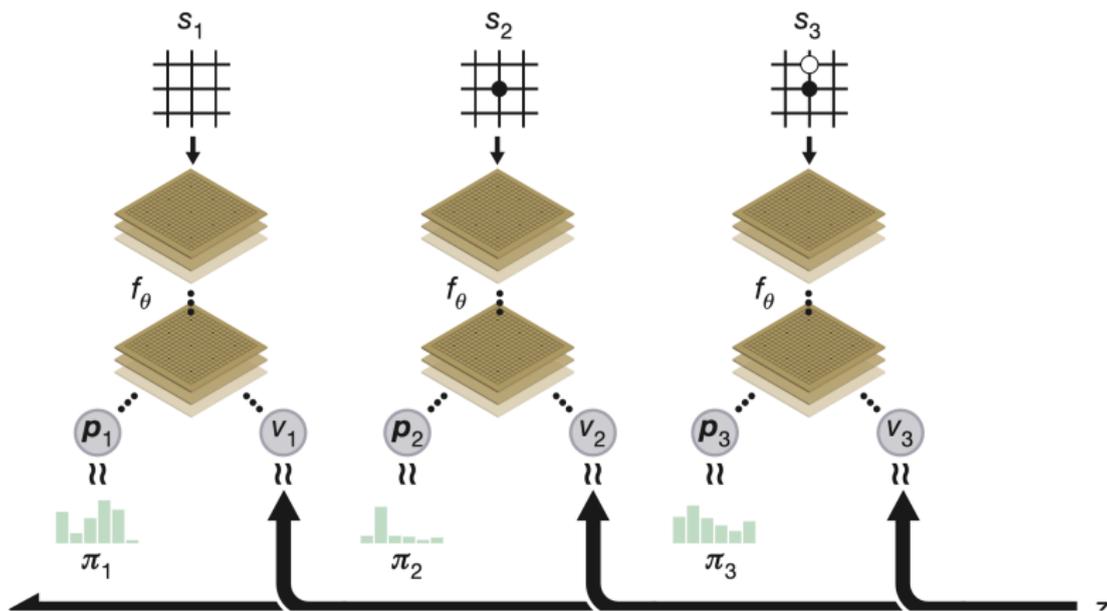
# Advantages of Self Play for Go

- Bottleneck is only computation, no humans needed
- Self-play also provides a well-matched player
- Optional check your understanding: how does this help with policy training? What is the reward density?

# Self Play for Go: Solution

- Bottleneck is only computation, no humans needed
- Self-play also provides a well-matched player
- Check your understanding: how does this help with policy training?  
What is the reward density?  
Rewards will be quite dense as both players are evenly matched. This provides a form of curriculum learning.

# Train Neural Network to Predict Policies and Values <sup>8</sup>



<sup>8</sup>Images from Silver et al. Nature 2017

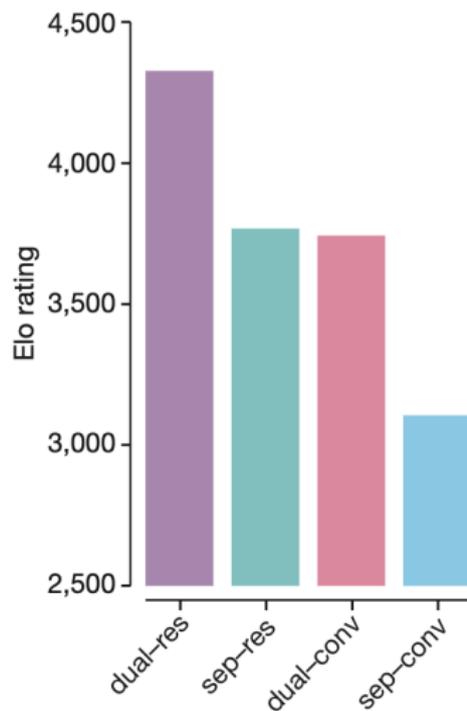
# AlphaGo and AlphaZero

- Self Play
- Strategic Computation
- Highly selective best-first search
- Power of Averaging
- Local Computation
- Learn and Update Heuristics

# AlphaGo and AlphaZero: Recap and Evaluation

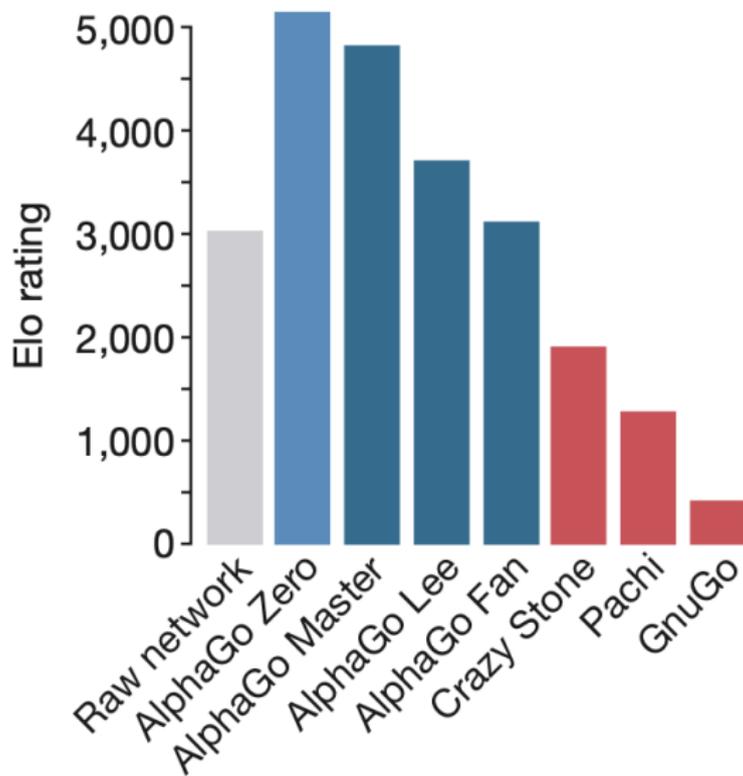
- Features:
  - Self Play
  - Strategic Computation
  - Highly selective best-first search
  - Power of Averaging
  - Local Computation
  - Learn and Update Heuristics
- Evaluation Questions
  - What is the influence of architecture?
  - What is the impact of using MCTS (on top of learning a policy / value function)?
  - How does it compare to human play or using human play?

# Impact of Architecture<sup>9</sup>



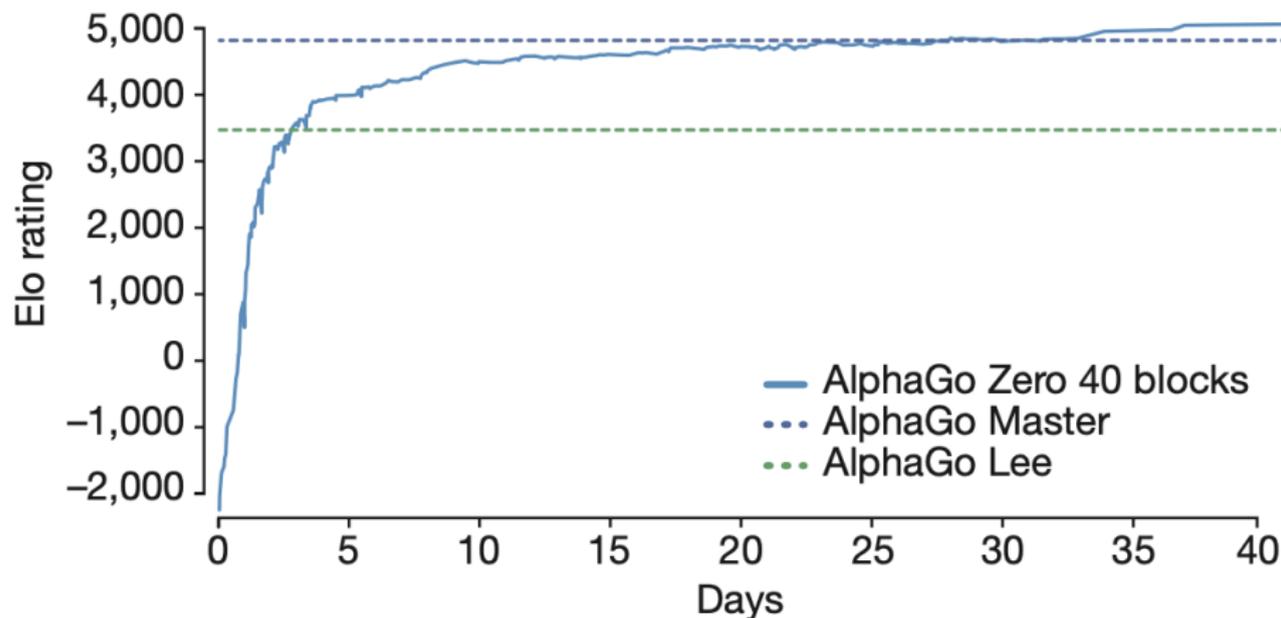
<sup>9</sup>Images from Silver et al. Nature 2017

# Impact of MCTS<sup>10</sup>



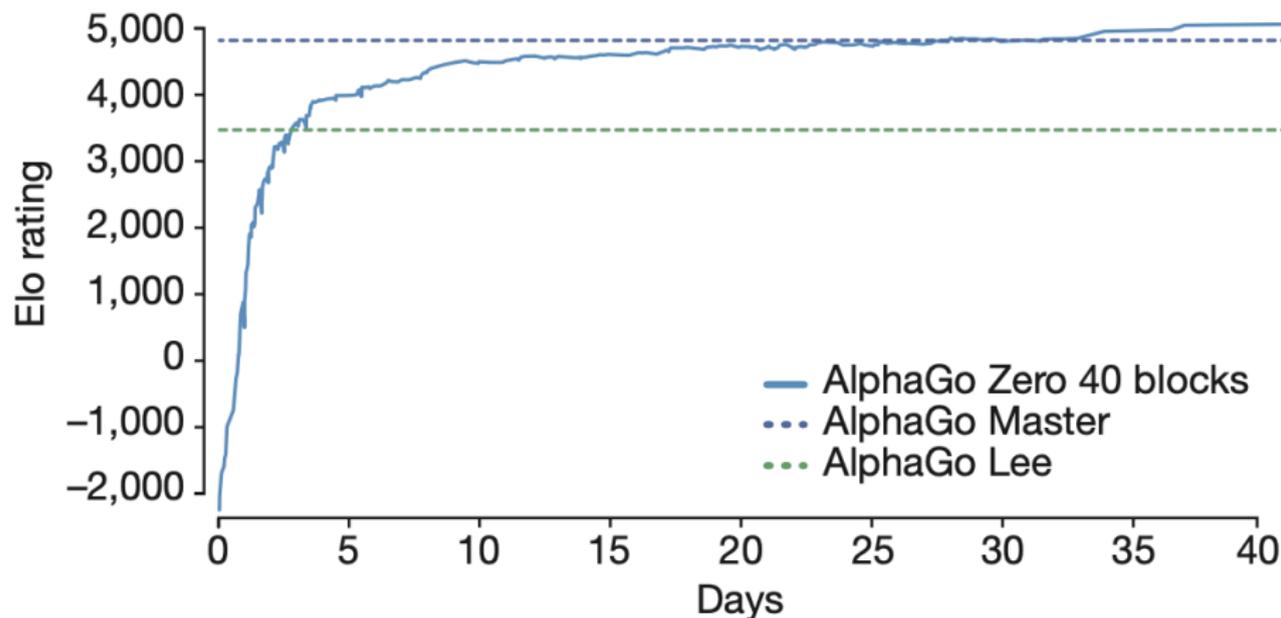
<sup>10</sup>Images from Silver et al. Nature 2017

# Overall performance<sup>11</sup>



<sup>11</sup>Images from Silver et al. Nature 2017

# Need for Human Data?<sup>12</sup>



<sup>12</sup>Images from Silver et al. Nature 2017

- These ideas are useful beyond Go. For chess, shogi, other games...
- For discovering faster matrix multiplication (AlphaTensor)
- For discovering faster sorting algorithms (AlphaDev)
- Beautiful insight: using RL to vastly speed up problems can represent as (extremely large) search problems

# Class Structure

- Last time: Guest lecture
- **This Time: MCTS**
- Next time: Quiz

# Optional Check Your Understanding: MCTS

- MCTS involves deciding on an action to take by doing tree search where it picks actions to maximize  $Q(S, A)$  and samples states.  
Select all
  - 1 Given a MDP, MCTS may be a good choice for short horizon problems with a small number of states and actions.
  - 2 Given a MDP, MCTS may be a good choice for long horizon problems with a large action space and a small state space
  - 3 Given a MDP, MCTS may be a good choice for long horizon problems with a large state space and small action space
  - 4 Not sure

F F T

# Optional Check Your Understanding: MCTS Solutions

- MCTS involves deciding on an action to take by doing tree search where it picks actions to maximize  $Q(S, A)$  and samples states.  
Select all
  - 1 Given a MDP, MCTS may be a good choice for short horizon problems with a small number of states and actions.
  - 2 Given a MDP, MCTS may be a good choice for long horizon problems with a large action space and a small state space
  - 3 Given a MDP, MCTS may be a good choice for long horizon problems with a large state space and small action space
  - 4 Not sure

F F T

# In more depth: Upper Confidence Tree (UCT) Search.

## Optional CYU

- UCT: borrow idea from bandit literature and treat each tree node where can select actions as a multi-armed bandit (MAB) problem
- Maintain an upper confidence bound over reward of each arm and select the best arm
- Check your understanding: Why is this slightly strange? Hint: why were upper confidence bounds a good idea for exploration/exploitation? Is there an exploration/exploitation problem during simulated episodes?<sup>13</sup>

---

<sup>13</sup>Relates to metalevel reasoning (for an example related to Go see "Selecting Computations: Theory and Applications", Hay, Russell, Tolpin and Shimony 2012)

# Optional Check Your Understanding: UCT Search

- In Upper Confidence Tree (UCT) search we treat each tree node as a multi-armed bandit (MAB) problem, and use an upper confidence bound over the future value of each action to help select actions for later rollouts. Select all that are true
  - 1 This may be useful since it will prioritize actions that lead to later good rewards
  - 2 UCB minimizes regret. UCT is minimizing regret within rollouts of the tree. (If this is true, think about if this a good idea?)
  - 3 Not sure

T. T

## In more depth: Upper Confidence Tree (UCT) Search

- UCT: borrow idea from bandit literature and treat each tree node where can select actions as a multi-armed bandit (MAB) problem
- Maintain an upper confidence bound over reward of each arm and select the best arm
- Hint: why were upper confidence bounds a good idea for exploration/exploitation? Is there an exploration/exploitation problem during simulated episodes?<sup>14</sup>

---

<sup>14</sup>Relates to metalevel reasoning (for an example related to Go see "Selecting Computations: Theory and Applications", Hay, Russell, Tolpin and Shimony 2012)

# Refresh Your Understanding

Select all that are true:

- Upper confidence bounds are used to balance exploration and leveraging the acquired information to achieve high reward
- These algorithms can be used in bandits and Markov decision processes
- If the reward model is known, there is no benefit to using an upper confidence bound algorithm

# Refresh Your Understanding

Select all that are true:

- Upper confidence bounds are used to balance exploration and leveraging the acquired information to achieve high reward
- These algorithms can be used in bandits and Markov decision processes
- If the reward model is known, there is no benefit to using an upper confidence bound algorithm

True. True. Depends on setting. In bandits, no additional gain. In RL, if the dynamics model is not known, there will be a gain.