

Planning and Optimal Control

C.3. Reinforcement Learning for Playing Games

Lars Schmidt-Thieme

Information Systems and Machine Learning Lab (ISMLL)
Institute for Computer Science
University of Hildesheim, Germany

Syllabus

A. Models for Sequential Data

- Tue. 22.10. (1) 1. Markov Models
- Tue. 29.10. (2) 2. Hidden Markov Models
- Tue. 5.11. (3) 3. State Space Models
- Tue. 12.11. (4) 3b. (ctd.)

B. Models for Sequential Decisions

- Tue. 19.11. (5) 1. Markov Decision Processes
- Tue. 26.11. (6) 1b. (ctd.)
- Tue. 3.12. (7) 1c. (ctd.)
- Tue. 10.12. (8) 2. Monte Carlo and Temporal Difference Methods
- Tue. 17.12. (9) 3. Q Learning
- Tue. 24.12. — — *Christmas Break* —
- Tue. 7.1. (10) 4. Policy Gradient Methods

C. Applications for Sequential Decisions

- Tue. 14.1. (11) 1. Cooperative Reinforcement Learning
- Tue. 21.1. (12) 2. Learning to Optimize
- Tue. 28.1. (13) 3. Reinforcement Learning for Games
- Tue. 4.2. (14) Q&A

Outline

1. Introduction
2. Improving a Policy via Monte Carlo Tree Search
3. A Policy Model That Learns to Improve
4. Evaluation

Outline

1. Introduction
2. Improving a Policy via Monte Carlo Tree Search
3. A Policy Model That Learns to Improve
4. Evaluation

Why Games?

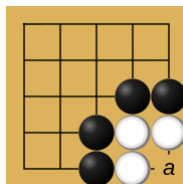
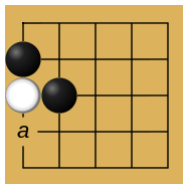
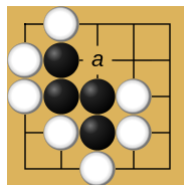
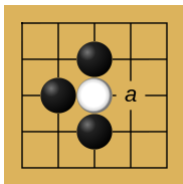
- ▶ ground truth mechanics
 - ▶ known and
 - ▶ simple (rules).
- ▶ consequences of actions taken in games can be assessed purely computationally
 - ▶ no costly/slow interactions required
 - ▶ with humans or
 - ▶ with the physical world.
- ▶ unlimited data (self play)
- ▶ difficult to win, requires intelligence.
- ▶ transition and reward model of the MDP are known, but its state space is too large to compute the optimal policy with the methods seen so far.
- ▶ unknown strategies of opponents render transitions stochastic.

Go / Rules

- ▶ start: board with 19×19 empty places.
- ▶ moves: first black, then white play alternately:
 - ▶ put an own stone on an empty place and
 - ▶ remove all sets of connected enemy stones that are not connected to an empty place (capture).
 - ▶ remove all sets of connected own stones that are not connected to an empty place (self capture).
 - ▶ moves leading back to an earlier position are forbidden.
- ▶ end: once both players consecutively passed.
- ▶ win: the player owning more places wins, counting
 - ▶ places occupied with own stones plus
 - ▶ connected empty places surrounded by own stones (territory).

Go / Rules / Capture

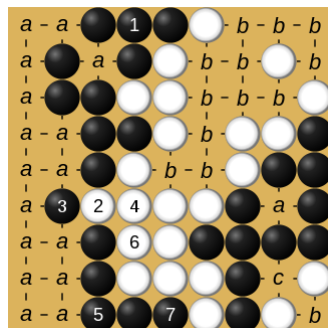
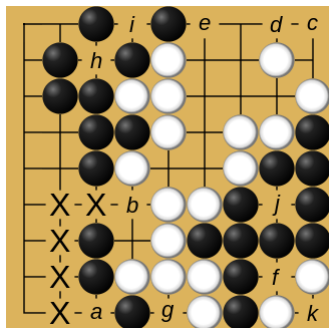
- ▶ put an own stone on an empty place and
 - ▶ remove all sets of connected enemy stones that are not connected to an empty place (capture).



[source: Wikipedia, Rules of Go, https://en.wikipedia.org/wiki/Rules_of_Go.]

Go / Rules / Capture

- ▶ win: the player owning more places wins, counting
 - ▶ places occupied with own stones plus
 - ▶ connected empty places surrounded by own stones (territory).



[source: Wikipedia, Rules of Go, https://en.wikipedia.org/wiki/Rules_of_Go.]

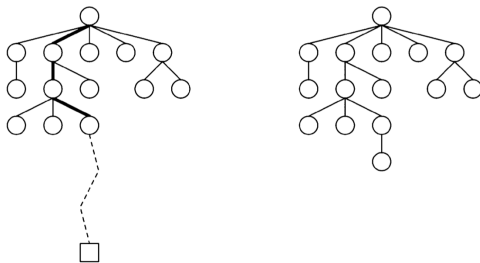
Outline

1. Introduction
2. Improving a Policy via Monte Carlo Tree Search
3. A Policy Model That Learns to Improve
4. Evaluation

AI Search Problem

- ▶ problem: find the optimal policy given an MDP model
- ▶ search approach: find the optimal action sequence by
 1. looking ahead into the state/action space
 - ▶ usually represented by a search tree
 2. computing values at the terminal nodes
 3. propagating the values back to the root node
 - ▶ finally: make decision based on expected values at the root node.
- ▶ adversarial search:
 - ▶ player chooses best action for himself. vs. opponent will choose worst action for the player.
 - ▶ minimax algorithm: in every node, alternatively
 - ▶ choose action with highest value (for the player; player's move)
 - ▶ choose action with the lowest value (for the player; opponent's move)
- ▶ most state spaces are too large for such a complete enumeration.

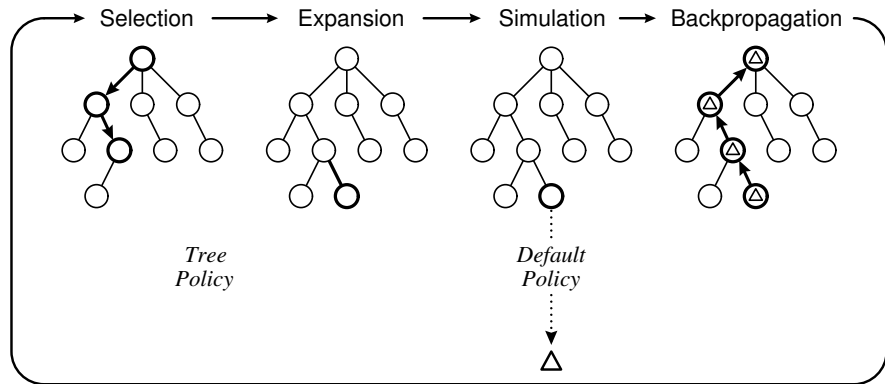
Monte Carlo Tree Search (MCTS) / Idea



[source: Browne et al. 2012]

- ▶ nodes = states
- ▶ edges = actions

Monte Carlo Tree Search (MCTS) / Steps



[source: Browne et al. 2012]

Monte Carlo Tree Search (MCTS)

- ▶ apply MCTS when the MDP is **not** known, but we have just estimates
 - ▶ for the value function, for the optimal policy
 - ▶ does not find the optimal policy, but hopefully will **improve** the actual estimates.
- ▶ given an initial policy π and estimated value function v find an improved action value function v' by
 - ▶ averaging over the values of follow-up states of an action a in a state s (any time step ahead)and an improved policy π'
 - ▶ by using the argmax policy of v' (“max child”)
 - ▶ here: by relative frequencies of the actions (“robust child”).

Monte Carlo Tree Search (MCTS)

- ▶ given an initial policy π and estimated value function v
 - simulate sequences from a generating policy (called **tree policy**)
 - ▶ that initially is π ,
 - ▶ later on shifts towards actions leading to high (estimated) value.
- ▶ to accomplish this, measure:
 - ▶ $N(s, a)$: how often an action a has been taken in state s during simulation so far.
 - ▶ $V(s, a)$: the total value seen in follow-up states of action a taken in state s (any time step ahead)
- ▶ a useful such generating policy (“PUCT algorithm”):

$$\tilde{\pi}(s, a; N, V, \pi) := \frac{V(s, a)}{N(s, a)} + c \frac{\sqrt{\sum_{a' \in A} N(s, a')}}{1 + N(s, a)} \pi(s, a)$$
 - ▶ where c is an exploitation/exploration trade-off weight.

Monte Carlo Tree Search (MCTS) / PUCT

$$\tilde{\pi}(s, a; N, V, \pi) = \left(\frac{1 + N_a}{1 + N_a + c\sqrt{N}} \hat{V}(a) + \frac{c\sqrt{N}}{1 + N_a + c\sqrt{N}} \pi(a) \right) \cdot \frac{1 + c\sqrt{N} \frac{1}{1+N_a}}{|A| + c\sqrt{N} \left(\sum_{a' \in A} \frac{1}{1+N_{a'}} \right)}$$

where $N_a := N(s, a)$, $N := \sum_{a' \in A} N(s, a')$, $\hat{V}(a) := \frac{V(s, a)}{N(s, a)}$, $\pi(a) := \pi(s, a)$

- ▶ **convex combination** of \hat{V} and π
 - ▶ $\tilde{\pi} \rightarrow \hat{V}$ with more samples ($N_a \rightarrow \infty$)
- ▶ 3rd **policy favoring actions taken rarely so far**
 - ▶ weight $1/(1 + N_a)$
 - ▶ doing nothing if all actions have been taken equally likely

Monte Carlo Tree Search (MCTS)

```

1 policy-mcts-improved( $s_0, \pi, v, K, c, \tau$ ):
2   create tree with root  $s_0$ , child indices  $A$  and edge attributes  $p, N, V$ :
3      $p(s_0, a) := \pi(s_0, a), N(s_0, a) := 0, V(s_0, a) := 0$  for all  $a \in A$ 
4   for  $k := 1, \dots, K$ :
5      $t := 0$ 
6     do
7        $a_t := \arg \max_{a \in A} \tilde{\pi}(s_t, a) := \frac{V(s_t, a)}{N(s_t, a)} + c \frac{\sqrt{\sum_{a' \in A} N(s_t, a')}}{1 + N(s_t, a)} p(s_t, a)$ 
8        $s_{t+1} := \text{execute-action}(s_t, a_t)$ 
9        $t := t + 1$ 
10    while a child node  $s_t$  of  $s_{t-1}$  for action  $a_{t-1}$  exists already in the tree
11    create child node  $s_t$  of  $s_{t-1}$  for action  $a_{t-1}$ 
12       $p(s_t, a) := \pi(s_t, a), N(s_t, a) := 0, V(s_t, a) := 0$  for all  $a \in A$ 
13       $v_t := v(s_t)$ 
14      for  $t' := 0, \dots, t - 1$ :
15         $N(s_{t'}, a_{t'}) := N(s_{t'}, a_{t'}) + 1$ 
16         $V(s_{t'}, a_{t'}) := V(s_{t'}, a_{t'}) + v_t$ 
17       $\pi'(a) := N(s_0, a)^{1/\tau} / \sum_{a' \in A} N(s_0, a')^{1/\tau}$ 
18    return  $\pi'$ 
  
```

where

- ▶ $s_0 \in S$ state to improve policy for.
- ▶ π original policy.
- ▶ $v \in \mathbb{R}^S$ original value function.
- ▶ $K \in \mathbb{N}$ number of simulation samples.
- ▶ τ annealing rate.
- ▶ c exploitation/exploration trade-off.
- ▶ returns π' an improved policy for state s_0 .

Outline

1. Introduction
2. Improving a Policy via Monte Carlo Tree Search
3. A Policy Model That Learns to Improve
4. Evaluation

Policy Model $\hat{\pi}$

- ▶ input: suitable representation of state s_t
 - ▶ richer: last t_0 states $s_{t-t_0+1}, \dots, s_{t-1}, s_t$
- ▶ output:
 - ▶ **value function**: estimated value \hat{v} of the input state.

$$v := \begin{cases} +1, & \text{if player wins} \\ -1, & \text{else} \end{cases}$$

- ▶ **improved policy**: estimated selection probabilities $\hat{\pi} \in \mathbb{R}^A$ of an look-ahead improved policy of $\hat{\pi}$ for this state, e.g.,

$$\pi := \text{policy-MCTS-improved}(\hat{\pi})$$

Policy Model $\hat{\pi}$ / Loss

multitask loss/objective function:

$$\ell(\hat{\pi}, \hat{v}; \pi, v) := (v - \hat{v})^2 + \pi^T \log^\circ \hat{\pi} + \lambda \|\theta\|^2$$

$$v := \begin{cases} +1, & \text{if player wins} \\ -1, & \text{else} \end{cases}$$

$$\pi := \text{policy-MCTS-improved}(\hat{\pi}), \in \mathbb{R}^A$$

$\hat{\pi}, \hat{v} :=$ current estimations of optimal policy and value
(=output of DNN)

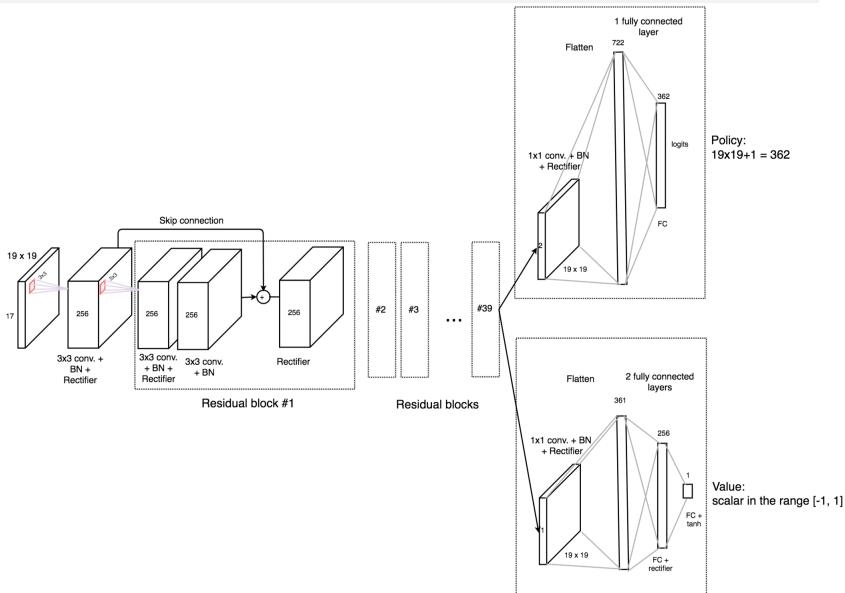
$\theta :=$ parameters of DNN

Policy Model $\hat{\pi}$ / AlphaGo Zero for Go (1/2)

- ▶ input:
 - ▶ board as 19×19 image with 17 binary channels
 - ▶ 8 for locations of white stones in last 8 positions
 - ▶ 8 for locations of black stones in last 8 positions
 - ▶ 1 for who's turn (same value for all pixels)
- ▶ 1 initial convolutional block:
 - ▶ convolution (3×3 , stride 1, 256 filters), batch normalization, rectifier.
- ▶ 19 residual blocks:
 - ▶ convolution (3×3 , stride 1, 256 filters), batch normalization, rectifier,
 - ▶ convolution (3×3 , stride 1, 256 filters), batch normalization,
 - ▶ skip connection adding block input, rectifier.

Policy Model $\hat{\pi}$ / AlphaGo Zero for Go (2/2)

- ▶ two heads for two outputs:
 - ▶ scalar \hat{v} :
 - ▶ convolution (1×1 , stride 1, 1 filter), batch normalization, rectifier, fully connected layer (size 256), rectifier, fully connected layer (size 1), tanh.
 - ▶ vectorized $(\hat{p}(a))_{a \in A}$, where A are all locations and action “pass”.
 - ▶ convolution (1×1 , stride 1, 2 filters), batch normalization, rectifier, fully connected layer (size $19^2 + 1 = 362$), logistic.
- ▶ 22.8M parameters
- ▶ with 19 residual blocks,
the receptive field of the last block is the whole board.



[source: https://medium.com/@jonathan_hui/alphago-zero-a-game-changer-14ef6e45eba5]

```
1 learn-policy( $N$ ):  
2    $\theta$  := randomly initialization  
3   do until convergence:  
4      $\mathcal{D} := \{\text{sample-policy}(\hat{\pi}(\theta)) \mid n = 1, \dots, N\}$   
5      $\theta := \text{update-model}(\theta, \mathcal{D})$   
6   return  $\theta$ 
```

where

- ▶ $N \in \mathbb{N}$ sample size per model update
- ▶ returns θ parameters of the policy model $\hat{\pi}$ (and value model \hat{v})

Required Computational Resources

data foundation:

- ▶ none
 - ▶ model learns only from generated data, not from human games.

computational resources:

- ▶ 3 days version:
 - ▶ 192 GPU days (= 3 days on 64 GPUs)
 - ▶ 4.9M games, 1,600 MCTS simulations/move
 - ▶ 700,000 minibatches a 2,048 positions each (possibly overlapping).
 - ▶ network with 20 residual blocks (40+ convolutional layers)
- ▶ 40 days version:
 - ▶ 2560 GPU days
 - ▶ network with 40 residual blocks (80+ convolutional layers)

Outline

1. Introduction
2. Improving a Policy via Monte Carlo Tree Search
3. A Policy Model That Learns to Improve
- 4. Evaluation**

Evaluation Criteria



Arpad Elo
(1903-1992)

- a. Elo rating:
 - ▶ a score that predicts how likely one player wins over another one
- b. accuracy of predicting the next move of a human expert player.
- c. accuracy of predicting the outcome of a match between professional players given a position.

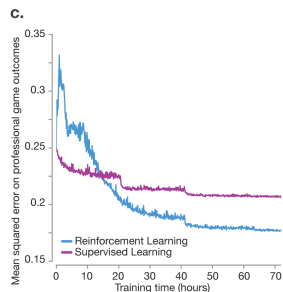
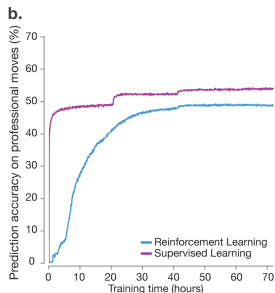
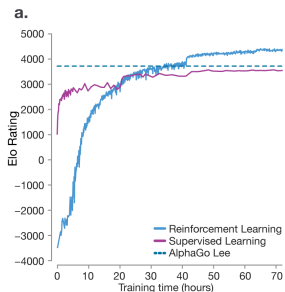
Elo Ratings

$$\begin{aligned}
 p(a \text{ wins against } b) &= \text{logistic}\left(\frac{1}{400}(r_a - r_b)\right) = \frac{1}{1 + e^{-\frac{1}{400}(r_a - r_b)}} \\
 &= \text{logistic}\left(\frac{1}{400}(r^T(e_a - e_b))\right)
 \end{aligned}$$

- ▶ Elo ratings are the weights of a logistic regression model for the player ID predictor.
- ▶ To evaluate alphaGoZero, the Elo ratings of alphaGo are used for reference.
 - ▶ and the alphaGo Elo ratings had been computed from its games against human Go masters.

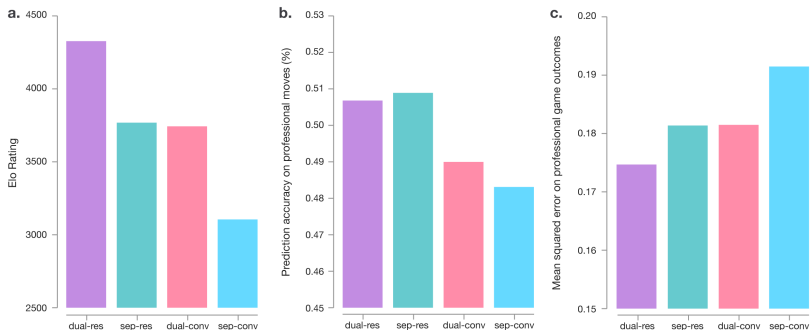
Note: e_a denotes the a -th unit vector in \mathbb{R}^N with N the number of players, i.e.,
 $(e_a)_n := \mathbb{I}(a = n), n = 1, \dots, N$.

Self-Play RL vs. Supervised Learning



[source: Silver et al. 2017]

Value/Policy Model Architectures



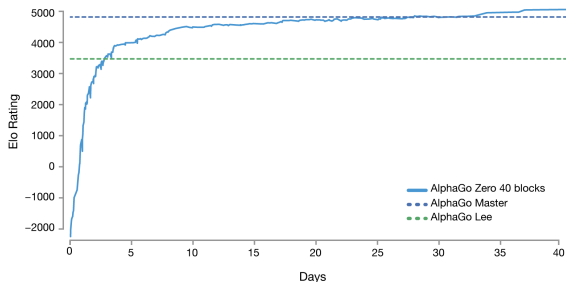
different architectures:

- ▶ multitask value/policy model (“**dual**”)
vs. two separate value and policy models (“**sep**”)
- ▶ residual deep neural network (“**res**”)
vs. purely convolutional deep neural network (“**conv**”)

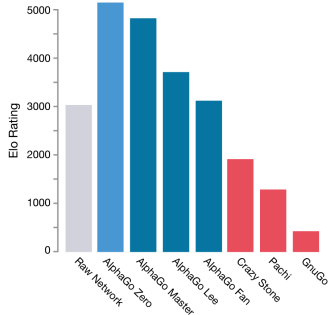
[source: Silver et al. 2017]

Self-Play RL (AlphaGo Zero) vs. Other Systems

a.



b.



[source: Silver et al. 2017]

Go as Testbed for Reinforcement Learning

- ▶ simple representation of the state
- ▶ simple representation of moves
- ▶ easily scalable problem sizes: vary board sizes
- ▶ hard task for humans
- ▶ easily scalable policy/value model
 - ▶ deep neural network, complexity scaled by number of layers

Further Readings

- ▶ About the policy model and the Go application:
 - ▶ Silver et al. [2017]
 - ▶ Some details, esp. about the MCTS used, are described more detailed in a precursor paper:
 - ▶ Silver et al. [2016]
- ▶ Monte Carlo Tree Search (MCTS):
 - ▶ a brief summary:
https://en.wikipedia.org/wiki/Monte_Carlo_tree_search
 - ▶ a survey: Browne et al. [2012]
 - ▶ The PUCT algorithm is an idiosyncratic combination of
 - ▶ UCT — *UCB applied to trees* [Kocsis and Szepesvári, 2006]
 - ▶ PUCB — *Predictors + UCB* [Rosin, 2011]
 - ▶ both extending UCB — *Upper Confidence Bounds* [Auer et al., 2002]
- ▶ Generally about adversarial search:
 - ▶ Russell et al. 2009, ch. 5

PUCT

UCB:

$$\pi^{(t)}(a) := \hat{Q}^{(t)}(a) + \sqrt{\frac{2 \log N^{(t)}}{N_a^{(t)}}}, \quad \text{with time } t = N$$

UCT:

$$\pi^{(t)}(s, a) := \hat{Q}^{(t)}(s, a) + 2c \sqrt{\frac{2 \log N_s^{(t)}}{N_{s,a}^{(t)}}}, \quad \text{with a constant } c$$

PUCB:

$$\pi^{(t)}(a) := \hat{Q}^{(t)}(a) + \sqrt{\frac{3 \log N^{(t)}}{2N_a^{(t)}}} - \frac{2}{M_a} \sqrt{\frac{\log N^{(t)}}{N^{(t)}}},$$

with action specific, time-invariant scalar predictors M_a

References

- Peter Auer, Nicola Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2-3):235–256, 2002.
- Cameron B. Browne, Edward Powley, Daniel Whitehouse, Simon M. Lucas, Peter I. Cowling, Philipp Rohlfshagen, Stephen Tavener, Diego Perez, Spyridon Samothrakis, and Simon Colton. A survey of monte carlo tree search methods. *IEEE Transactions on Computational Intelligence and AI in games*, 4(1):1–43, 2012.
- Levente Kocsis and Csaba Szepesvári. Bandit based monte-carlo planning. In *European Conference on Machine Learning*, pages 282–293. Springer, 2006.
- Christopher D. Rosin. Multi-armed bandits with episode context. *Annals of Mathematics and Artificial Intelligence*, 61(3):203–230, 2011.
- Stuart Jonathan Russell, Peter Norvig, John F. Canny, Jitendra M. Malik, and Douglas D. Edwards. *Artificial Intelligence: A Modern Approach*. Prentice hall Upper Saddle River, 3rd edition, 2009.
- David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, and Marc Lanctot. Mastering the game of Go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, and Adrian Bolton. Mastering the game of Go without human knowledge. *Nature*, 550(7676):354–359, 2017.