

Stochastic Planning in Games

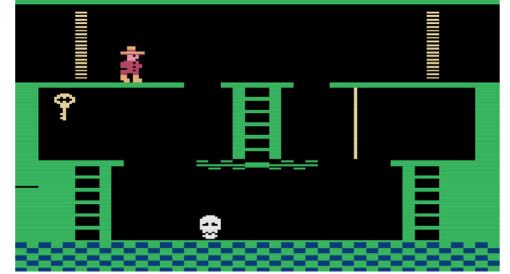
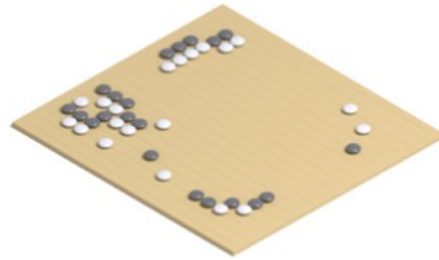
Peter Müller



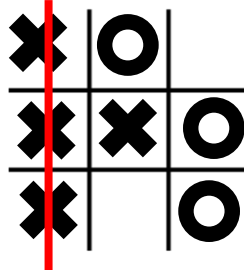
Why Games?



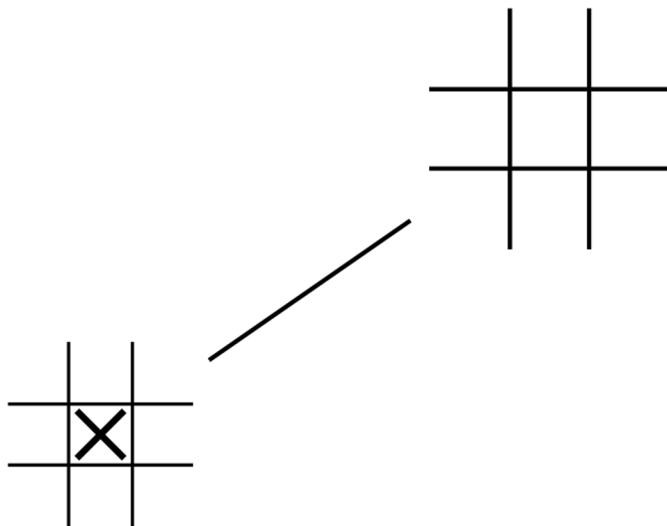
O	O	O
	O	X
	X	X



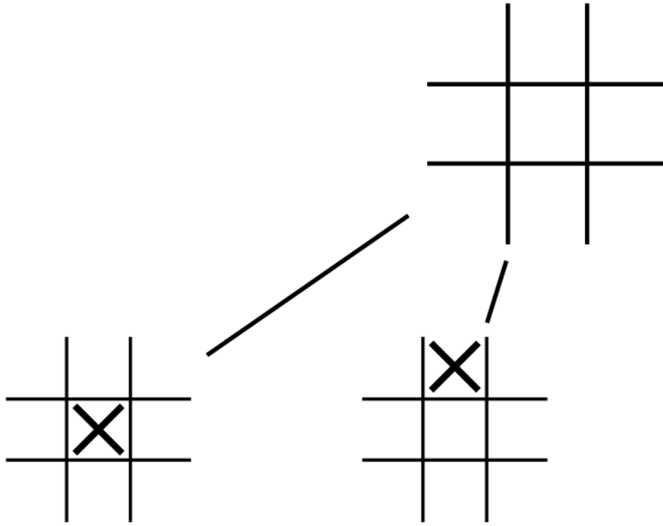
Tic-Tac-Toe



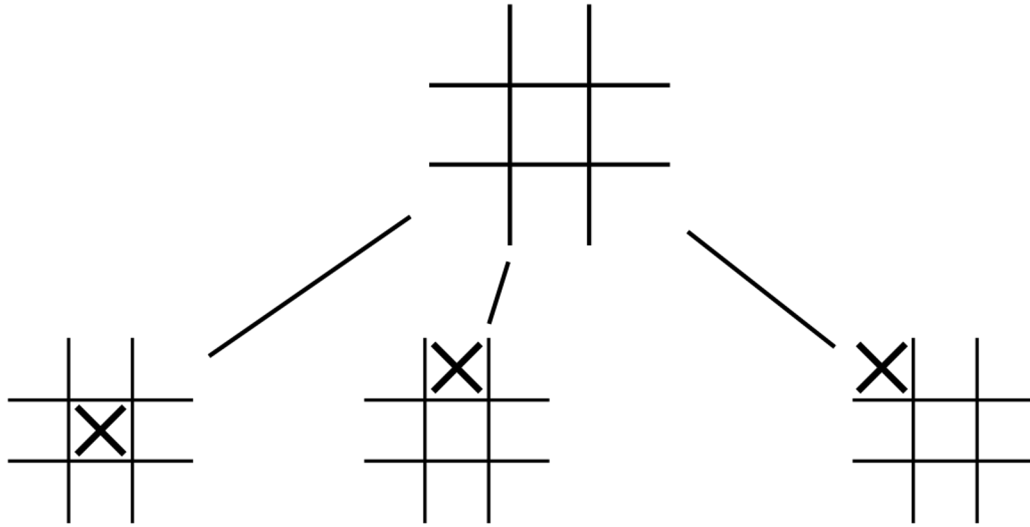
Tic-Tac-Toe



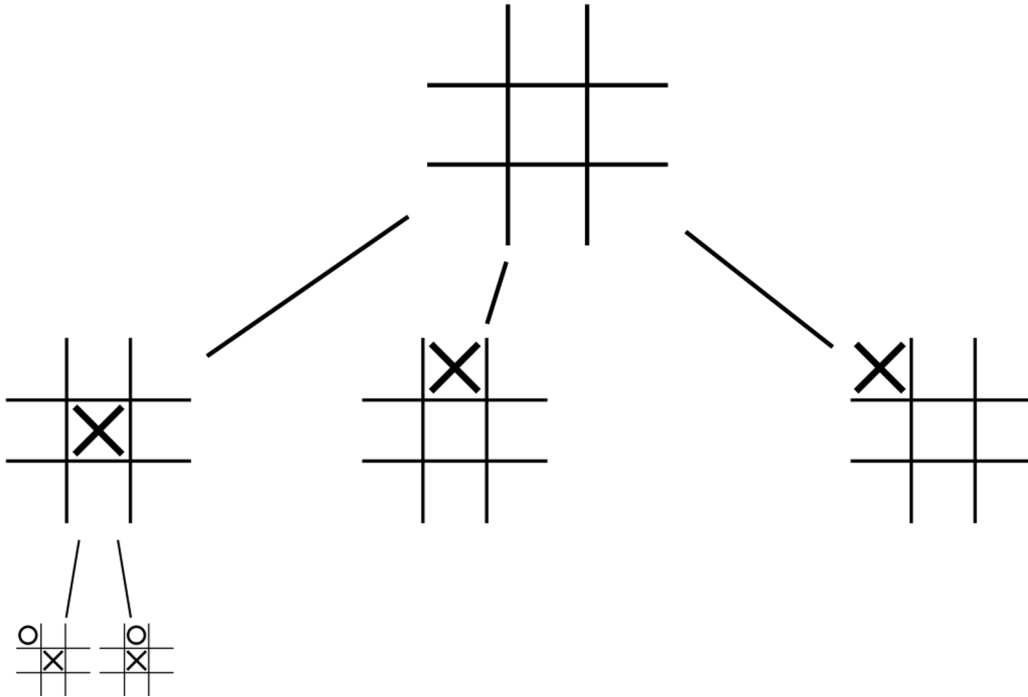
Tic-Tac-Toe



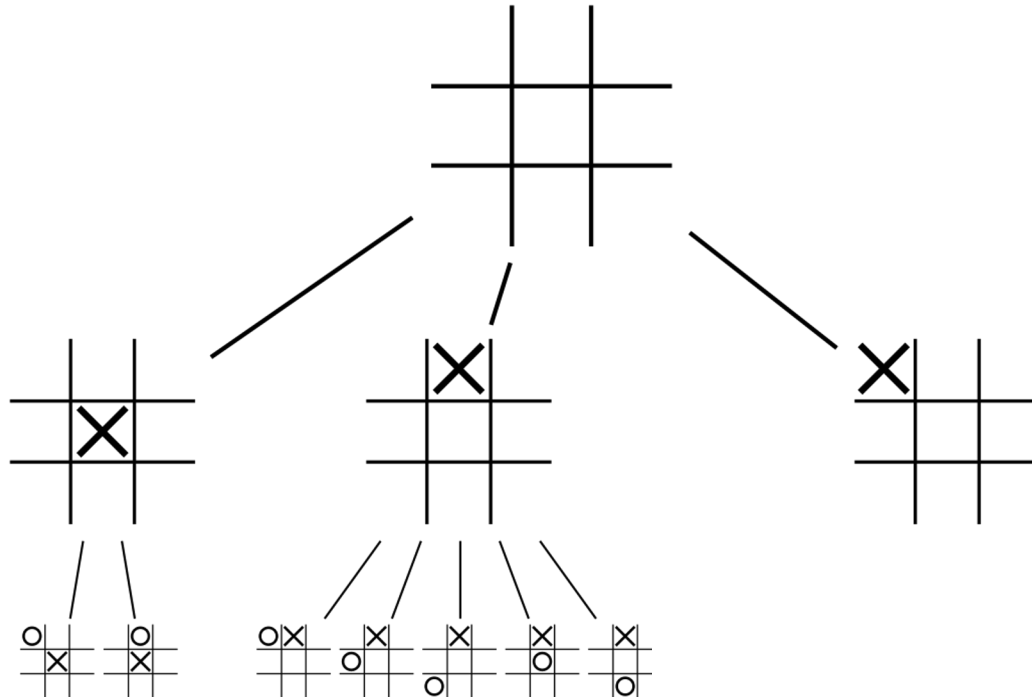
Tic-Tac-Toe



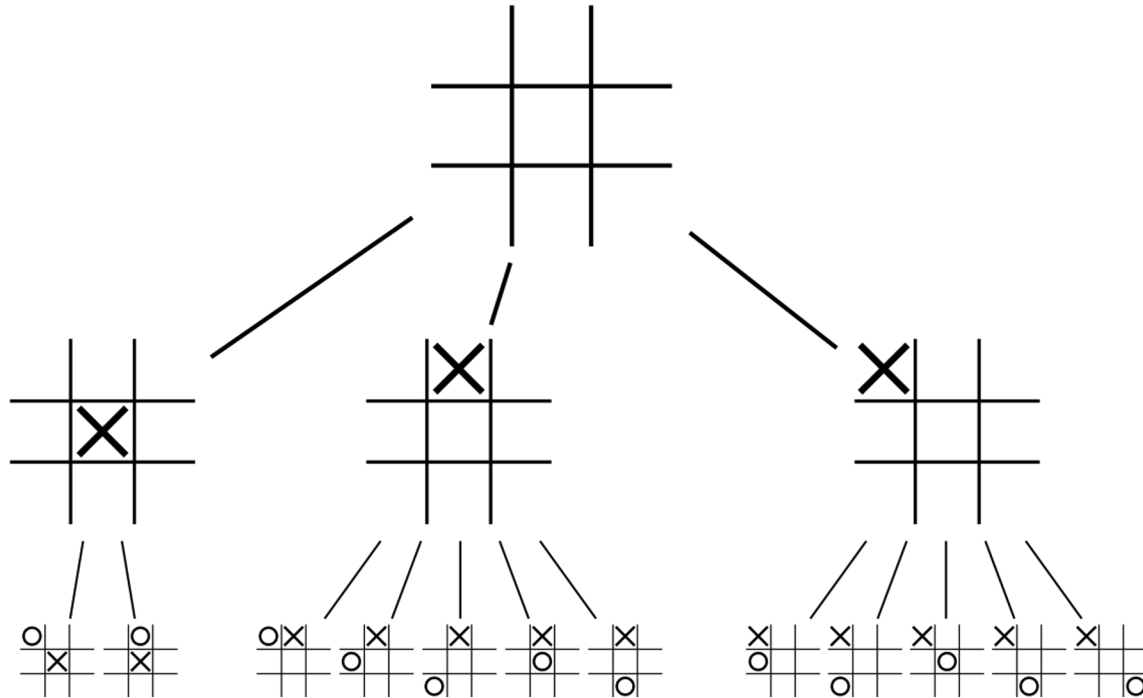
Tic-Tac-Toe



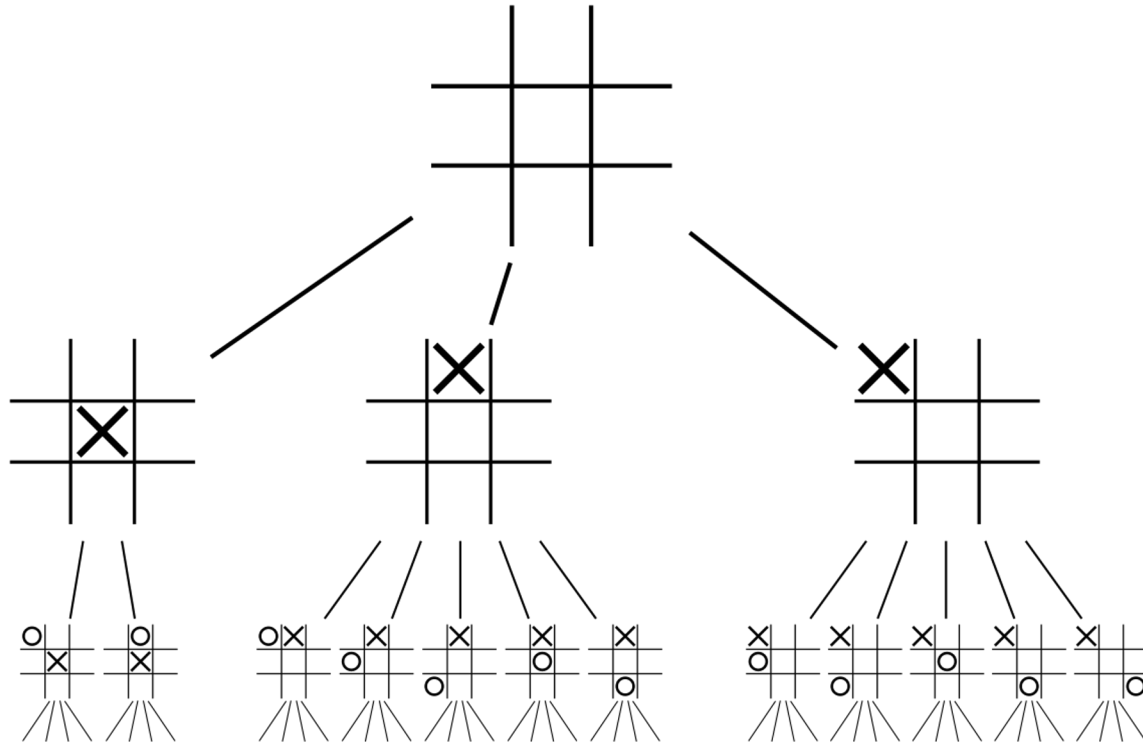
Tic-Tac-Toe



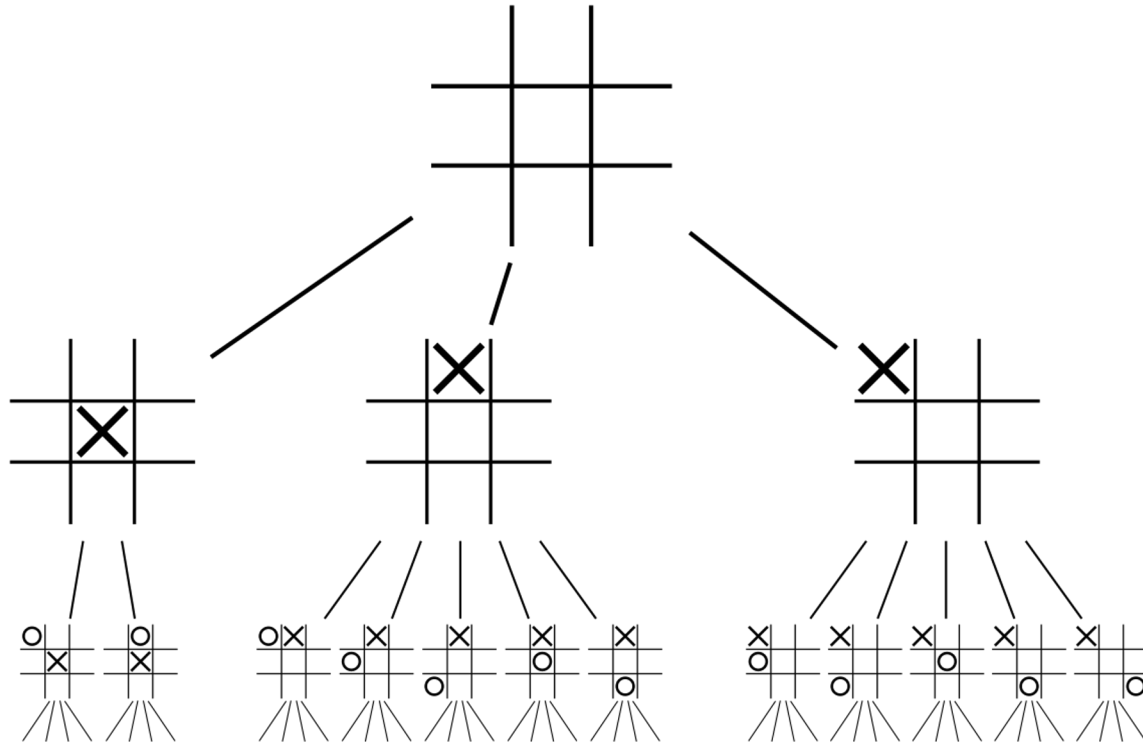
Tic-Tac-Toe



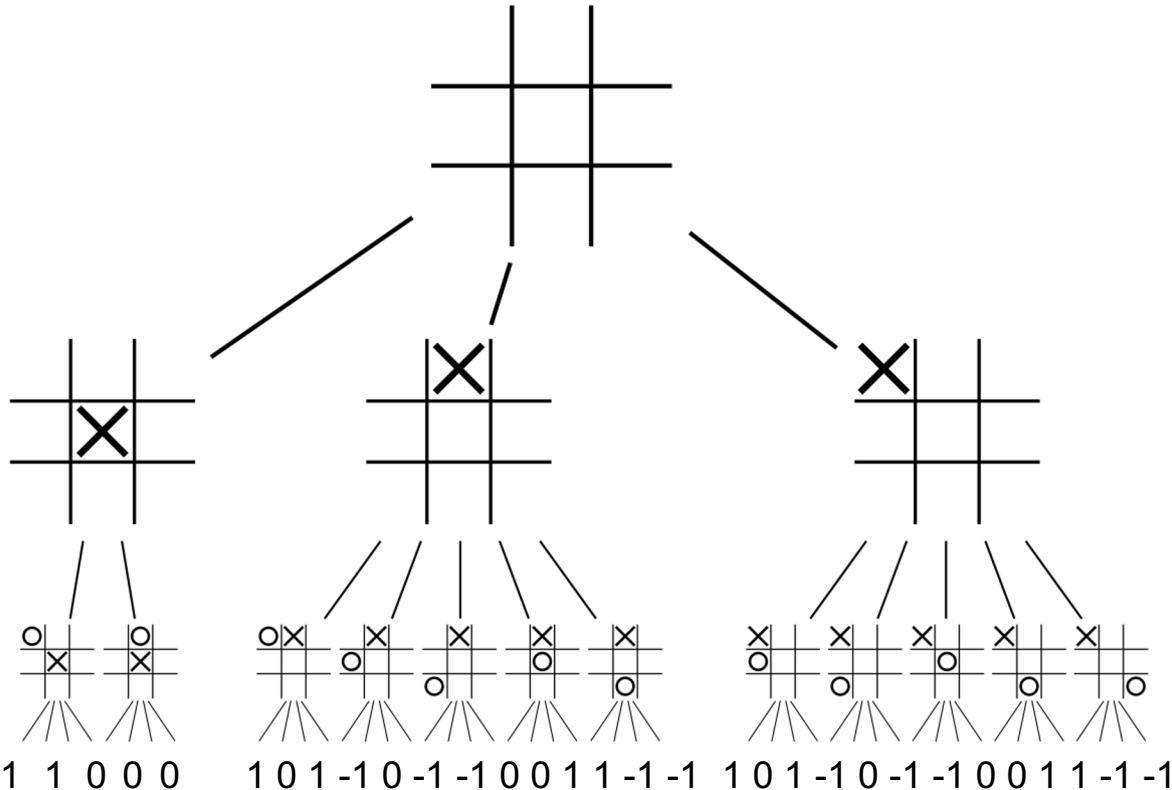
Tic-Tac-Toe



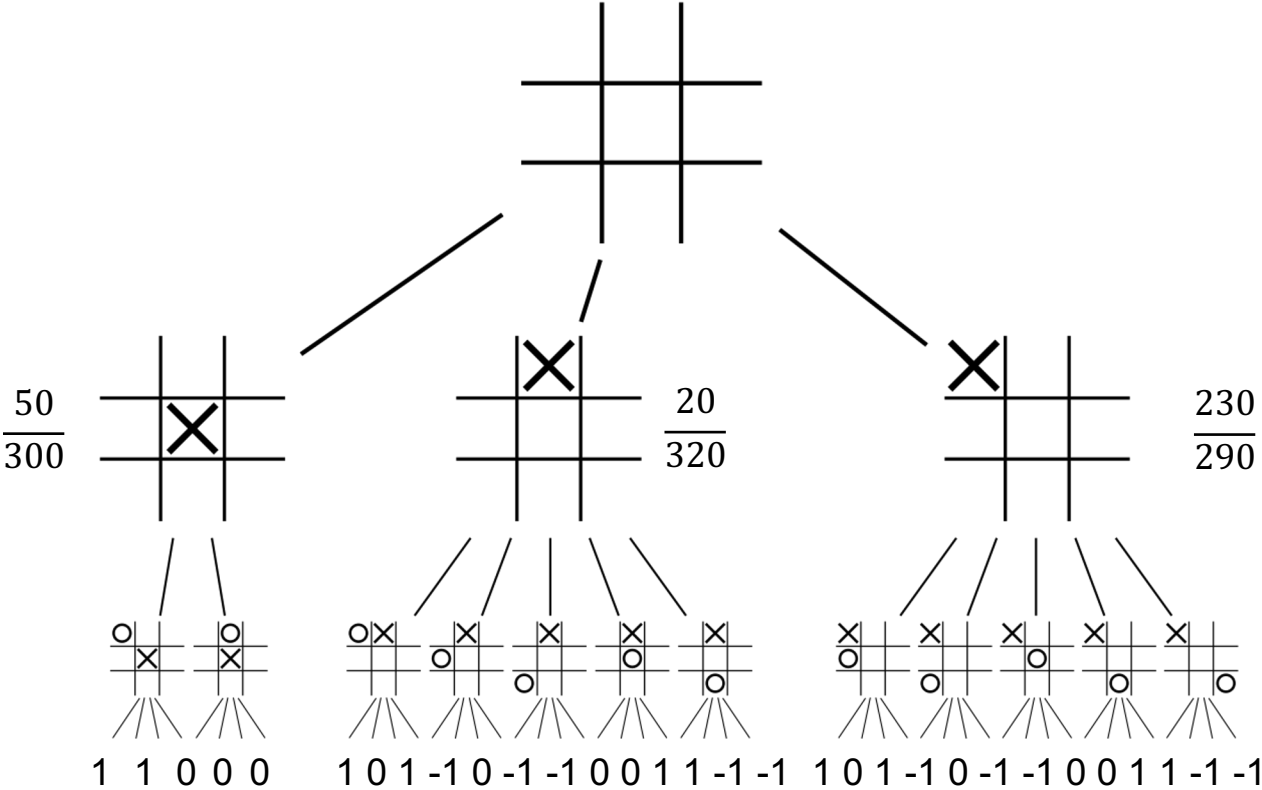
Tic-Tac-Toe



Tic-Tac-Toe



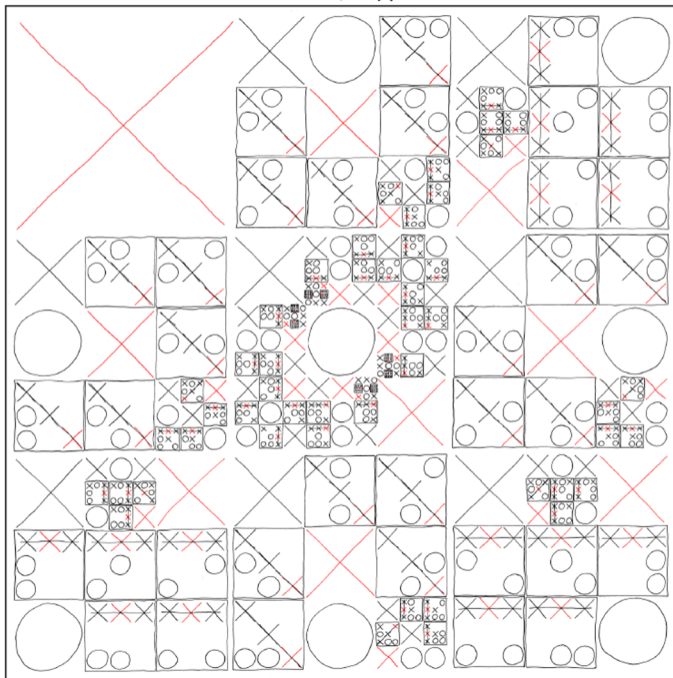
Tic-Tac-Toe



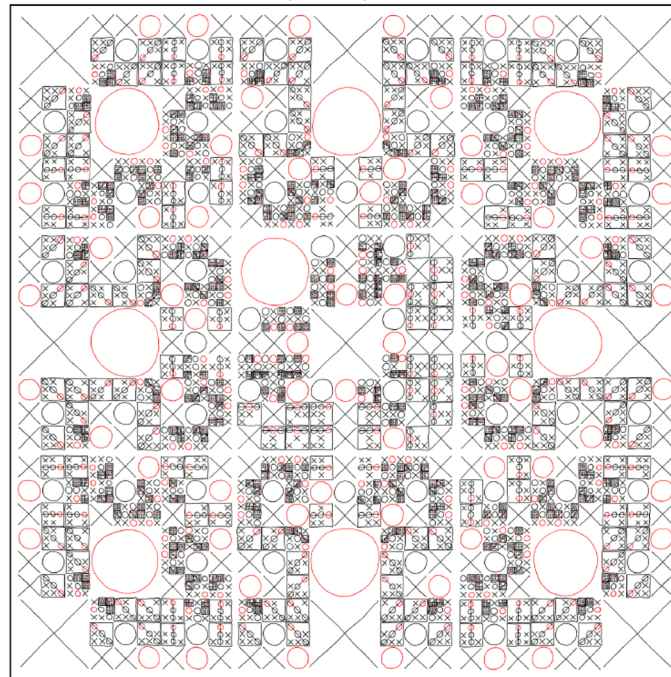
COMPLETE MAP OF OPTIMAL TIC-TAC-TOE MOVES

YOUR MOVE IS GIVEN BY THE POSITION OF THE LARGEST RED SYMBOL ON THE GRID. WHEN YOUR OPPONENT PICKS A MOVE, ZOOM IN ON THE REGION OF THE GRID WHERE THEY WENT. REPEAT.

MAP FOR X:



MAP FOR O:

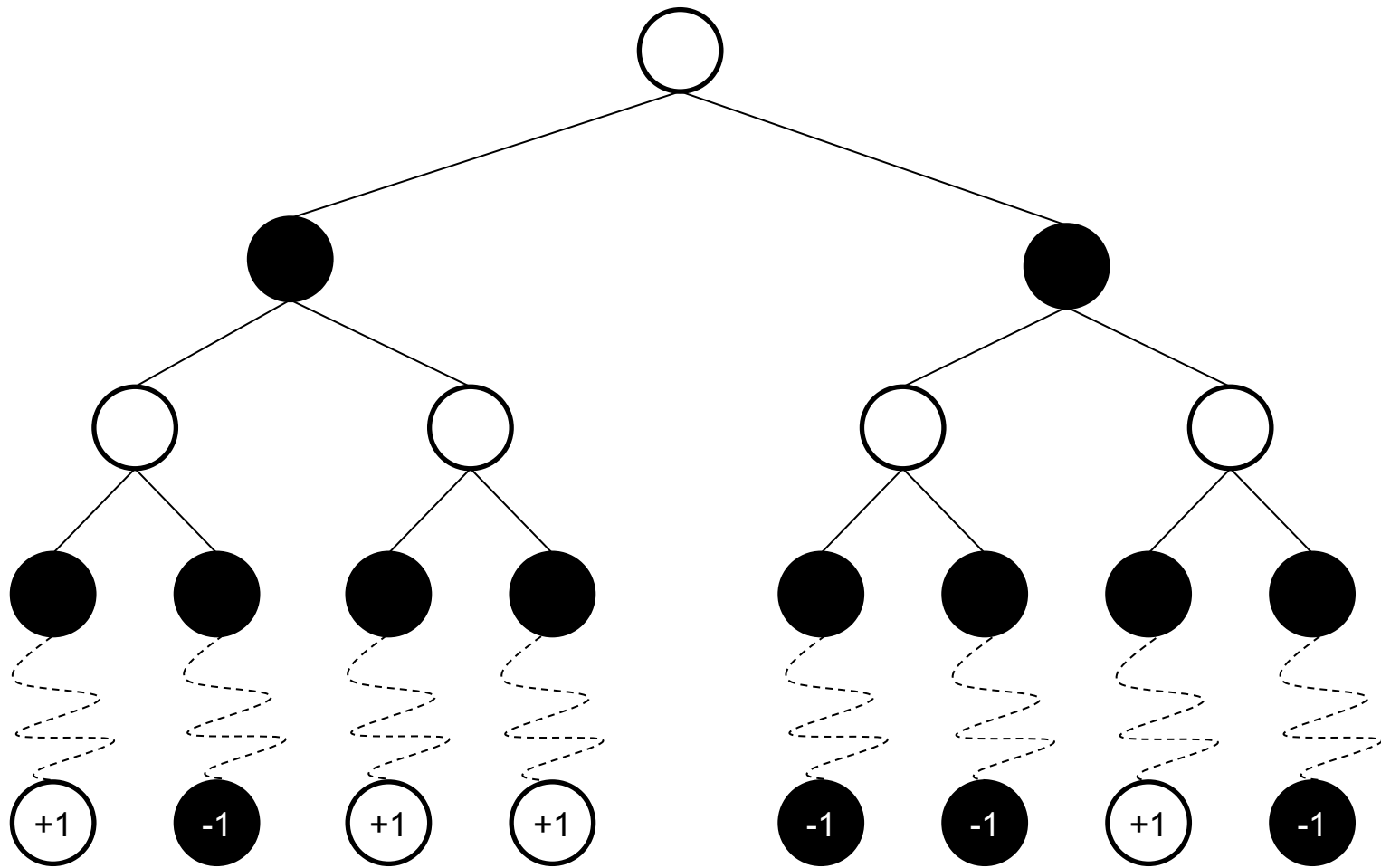


Chess



- Branching Factor: 35
- Game Length: 80





Evaluation Function



White: 2x Bishop (3 points) + 5x Pawn (1 point) = 11 points



Evaluation Function

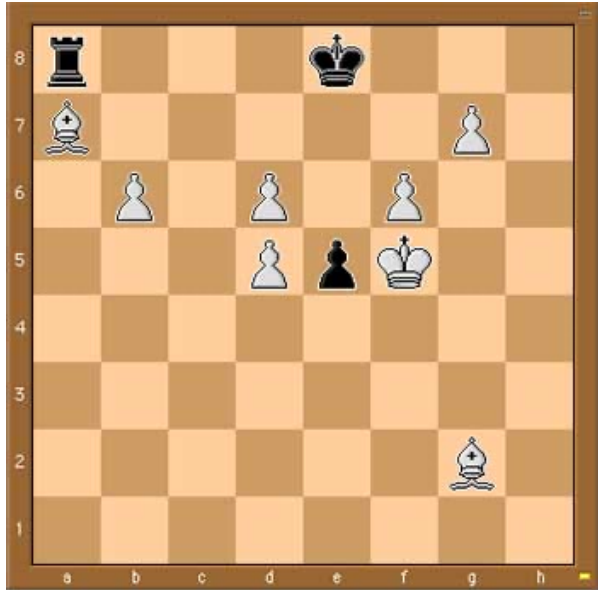


White: 2x Bishop (3 points) + 5x Pawn (1 point) = 11 points

Black: 1x Rook (5 points) + 1x Pawn (1 point) = 6 points



Evaluation Function



White: 2x Bishop (3 points) + 5x Pawn (1 point) = 11 points

Black: 1x Rook (5 points) + 1x Pawn (1 point) = 6 points

Eval: $11 - 6 = 5$ points



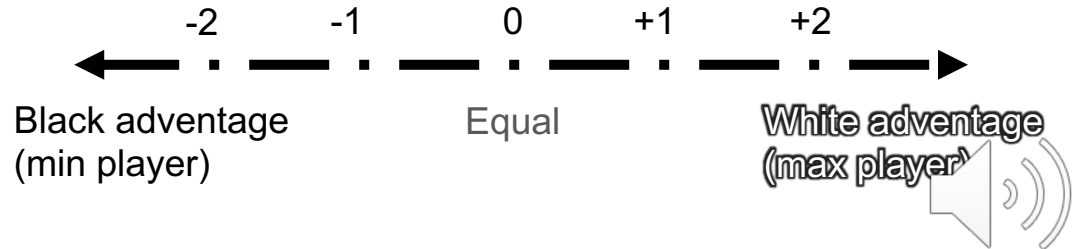
Evaluation Function



White: 2x Bishop (3 points) + 5x Pawn (1 point) = 11 points

Black: 1x Rook (5 points) + 1x Pawn (1 point) = 6 points

Eval: $11 - 6 = 5$ points



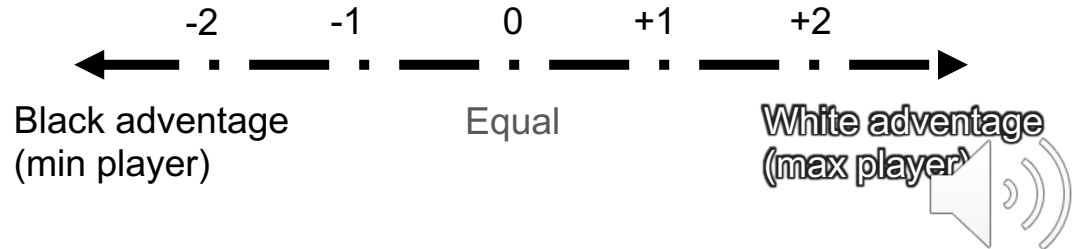
Evaluation Function

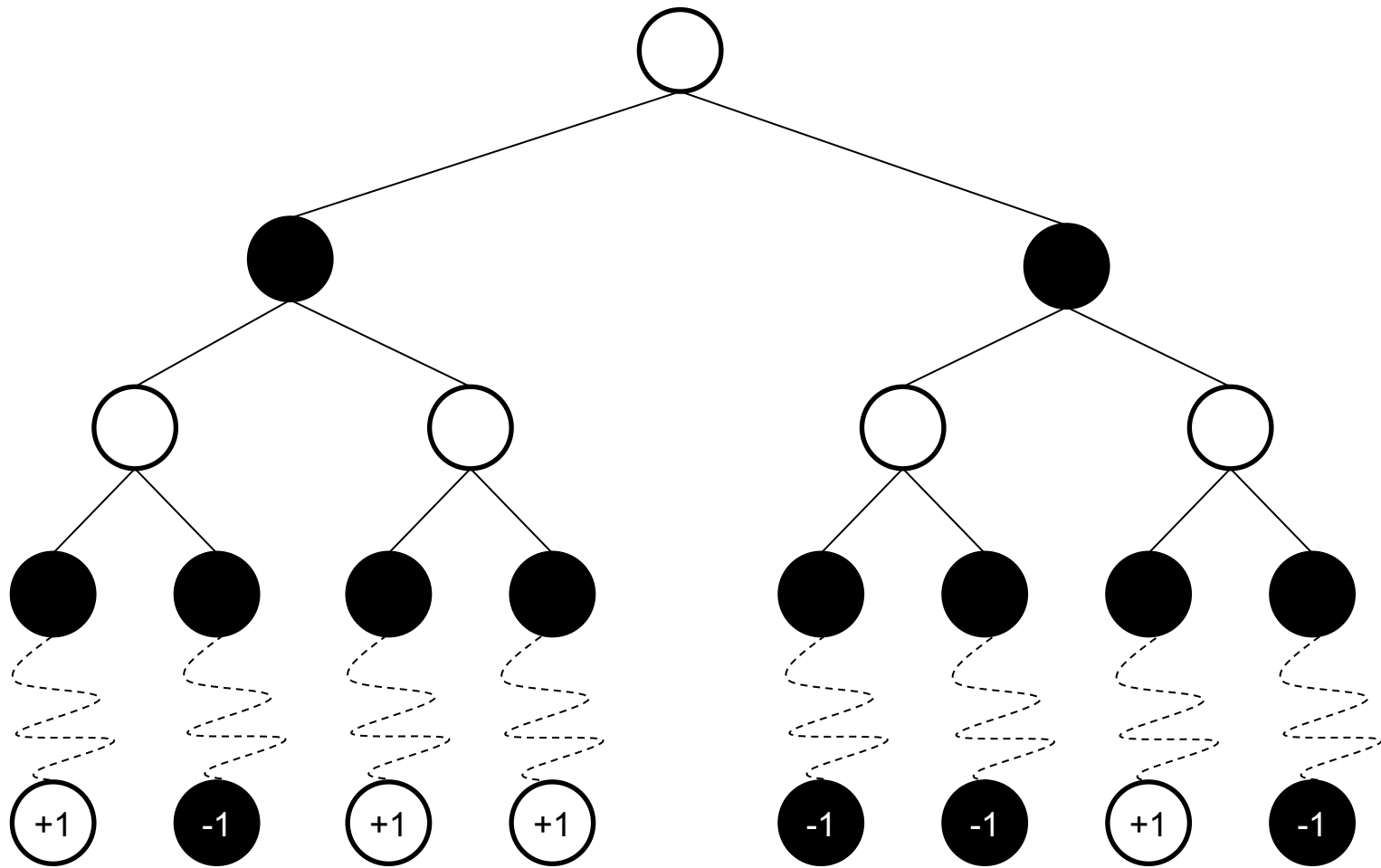


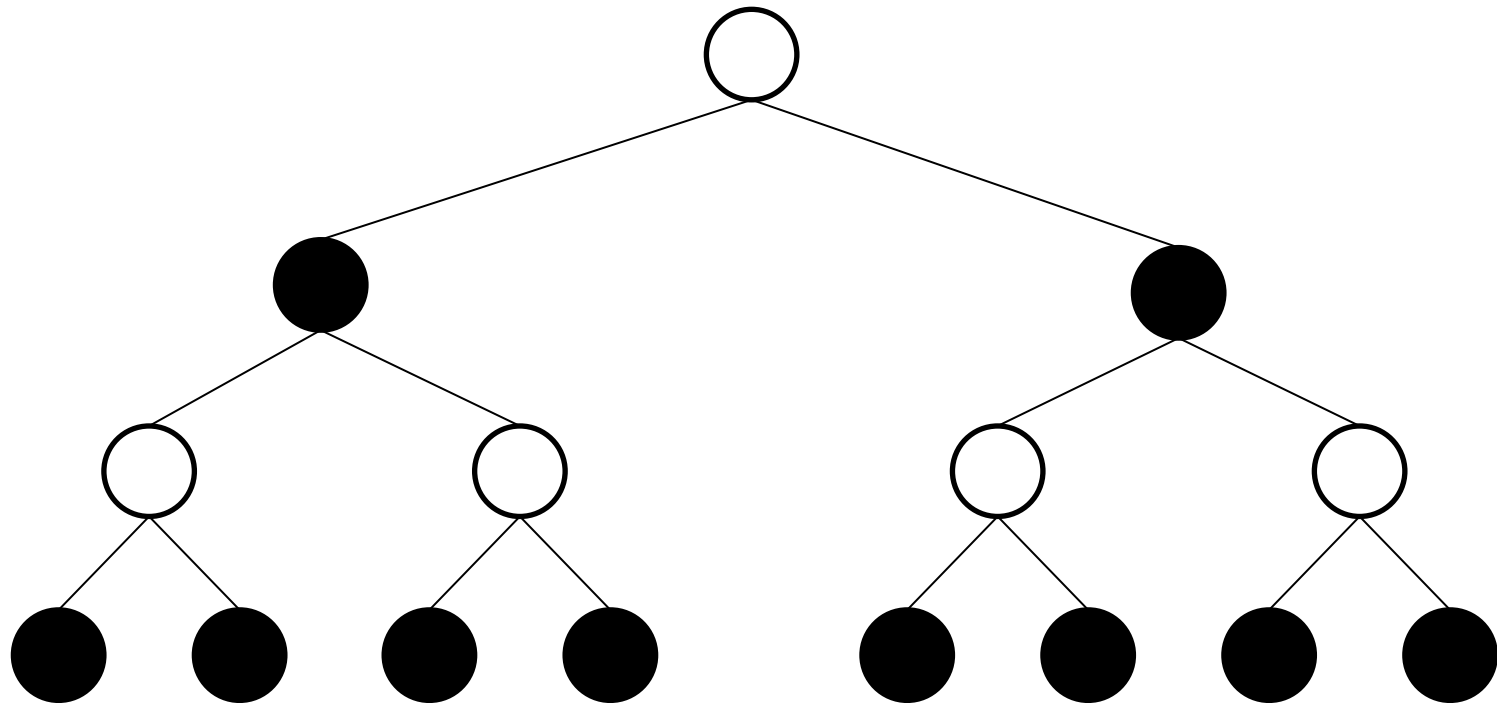
White: 2x Bishop (3 points) + 5x Pawn (1 point) = 11 points

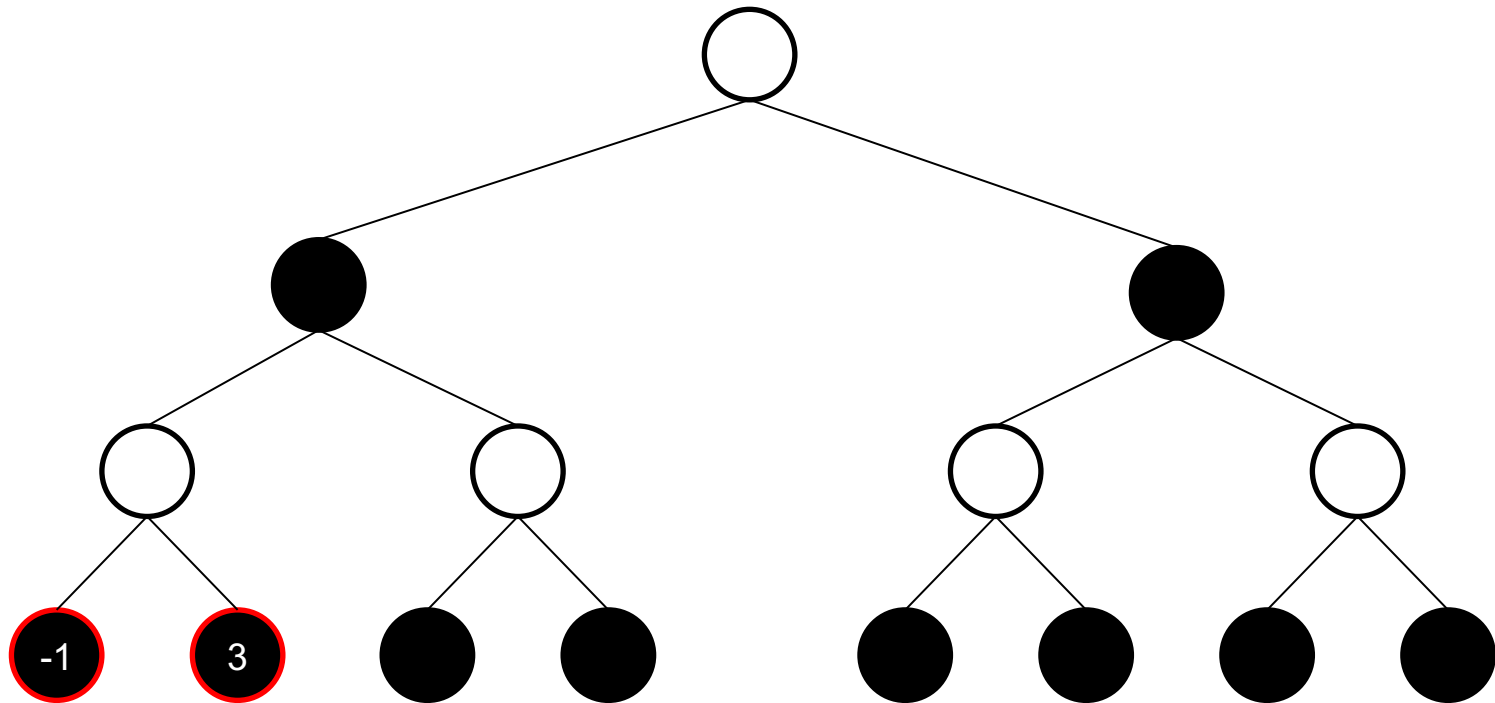
Black: 1x Rook (5 points) + 1x Pawn (1 point) = 6 points

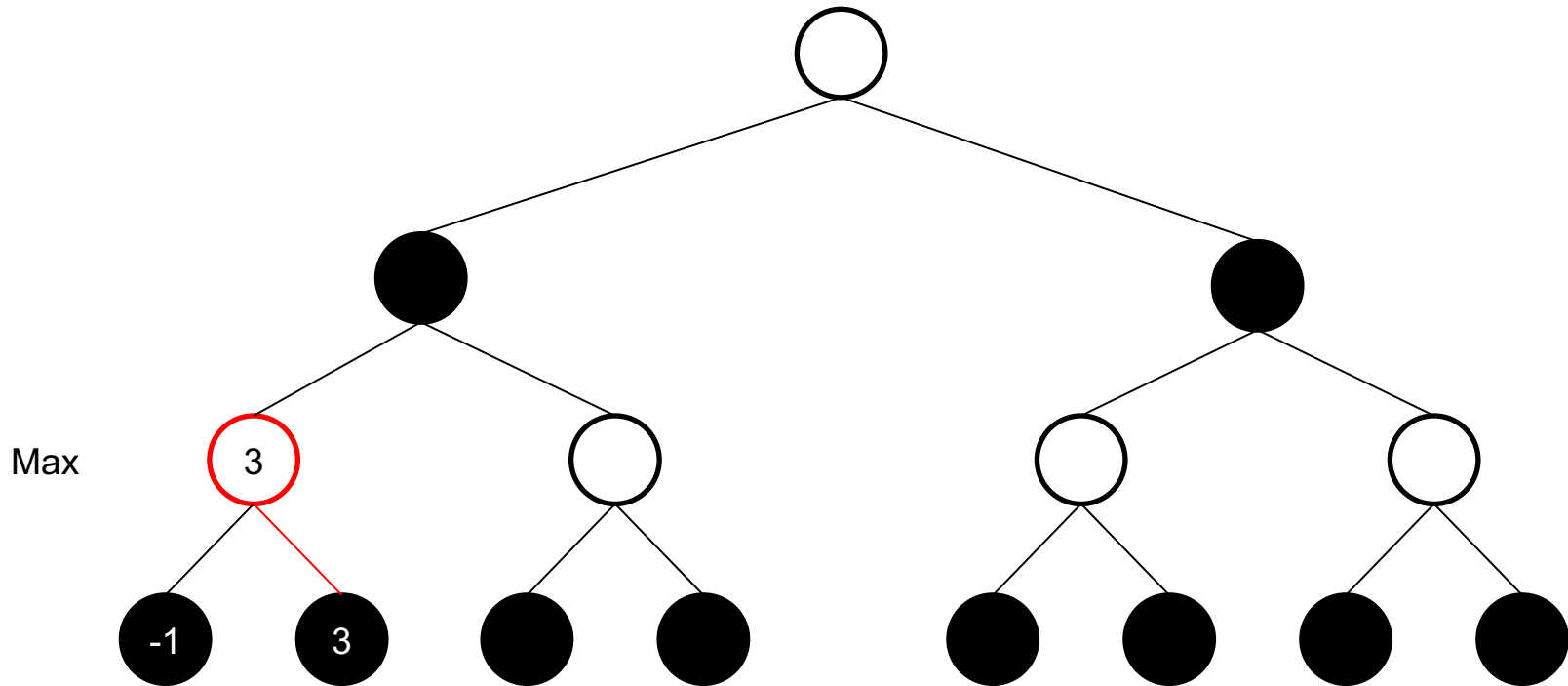
Eval: $11 - 6 = 5$ points

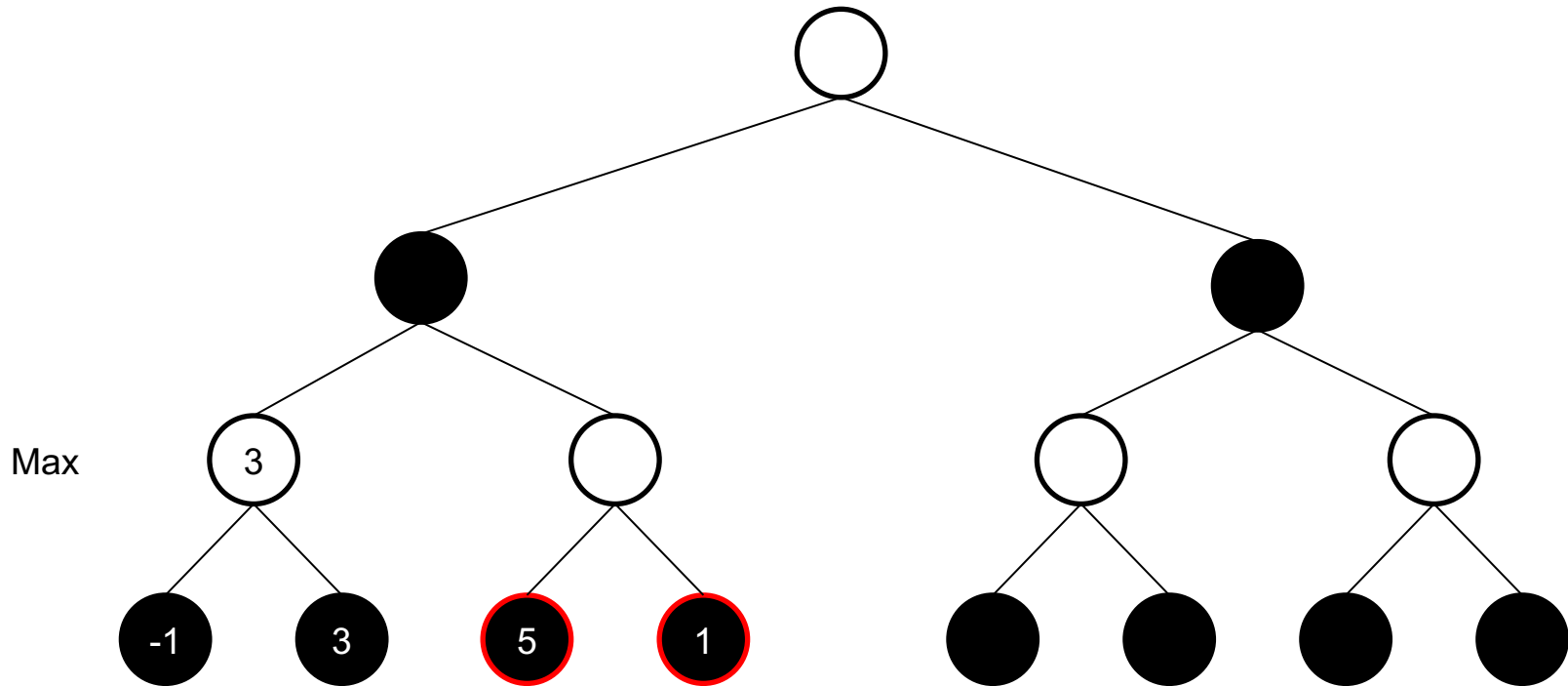






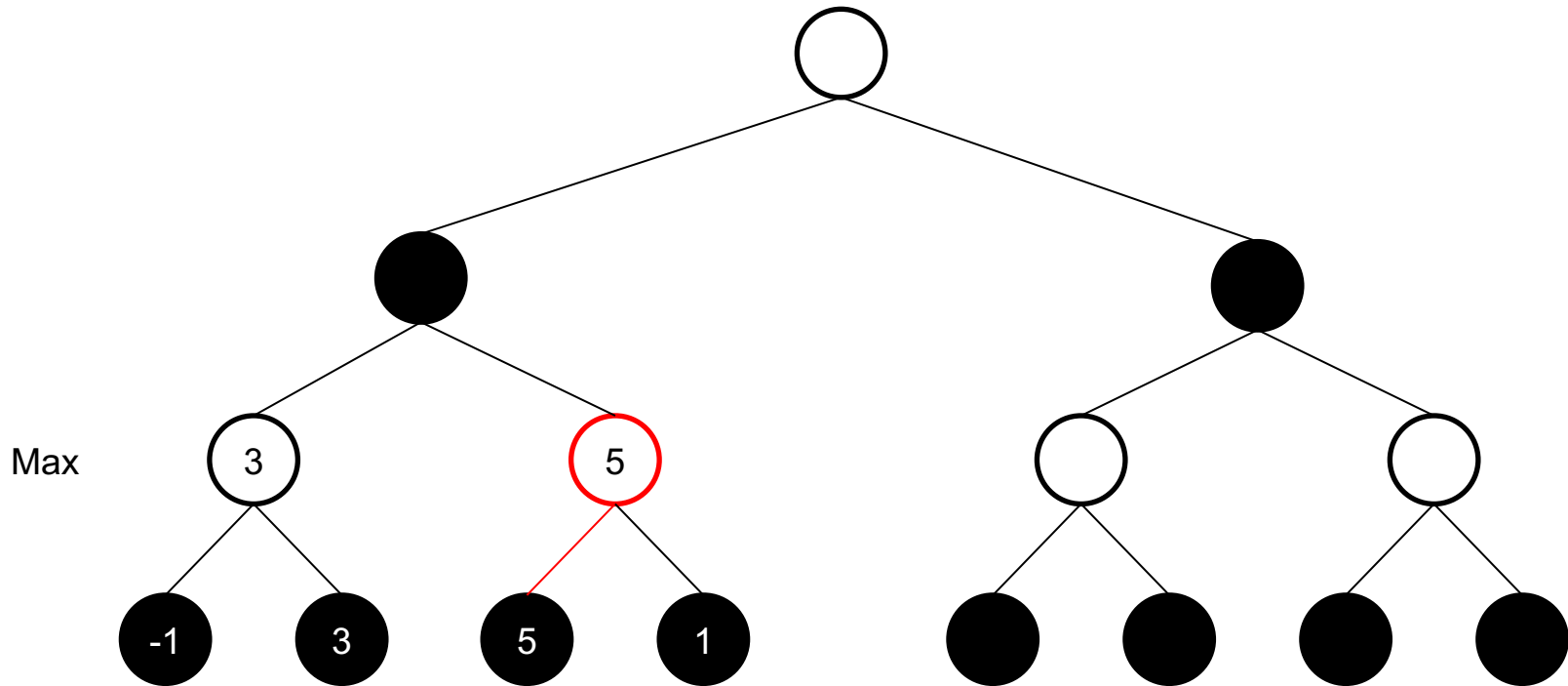


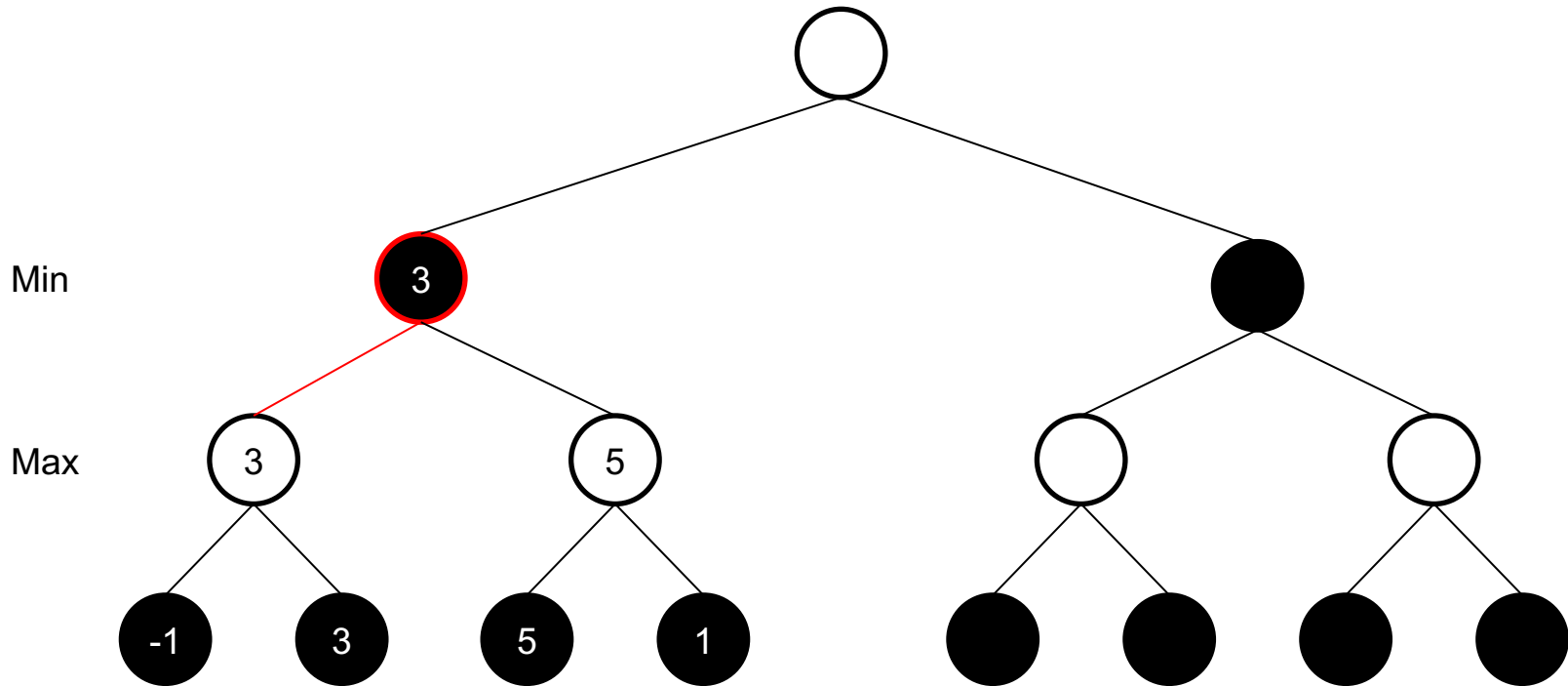


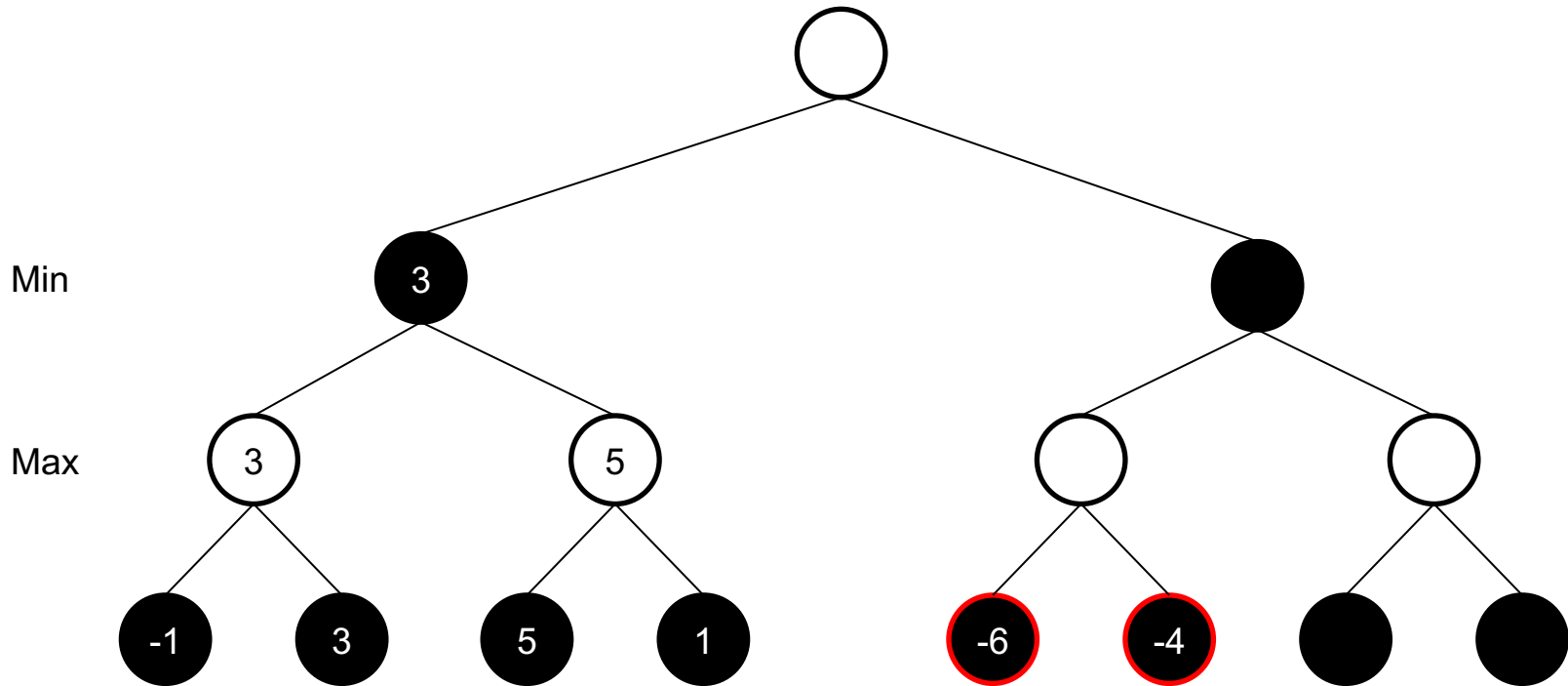


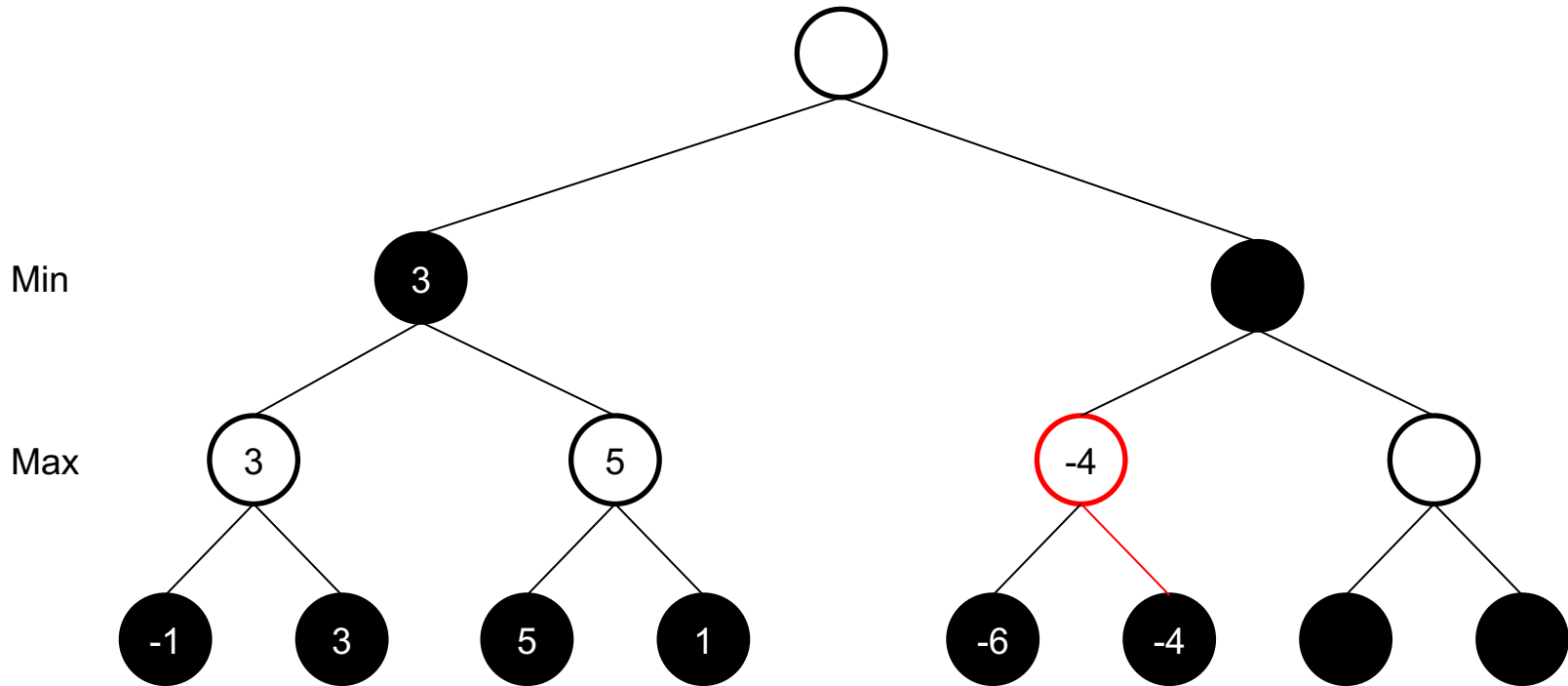
Max

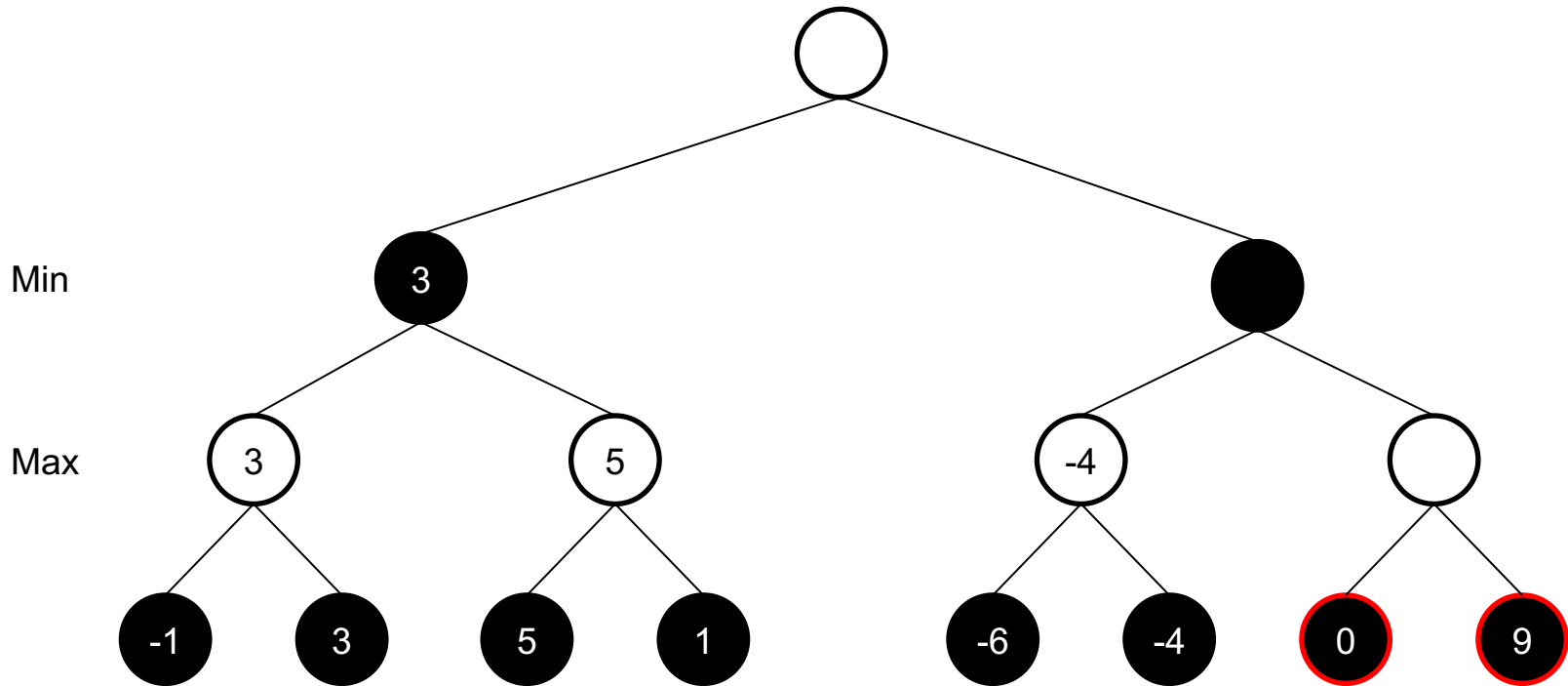


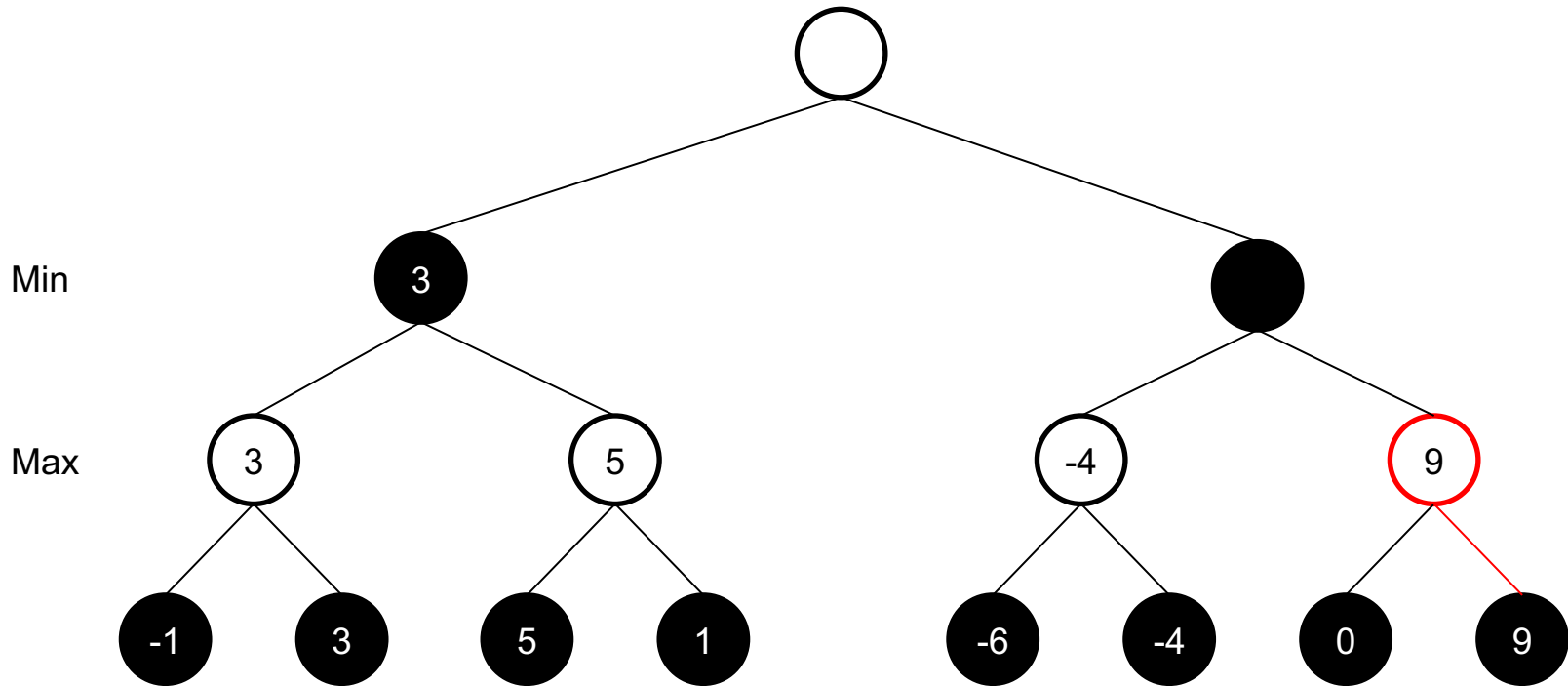


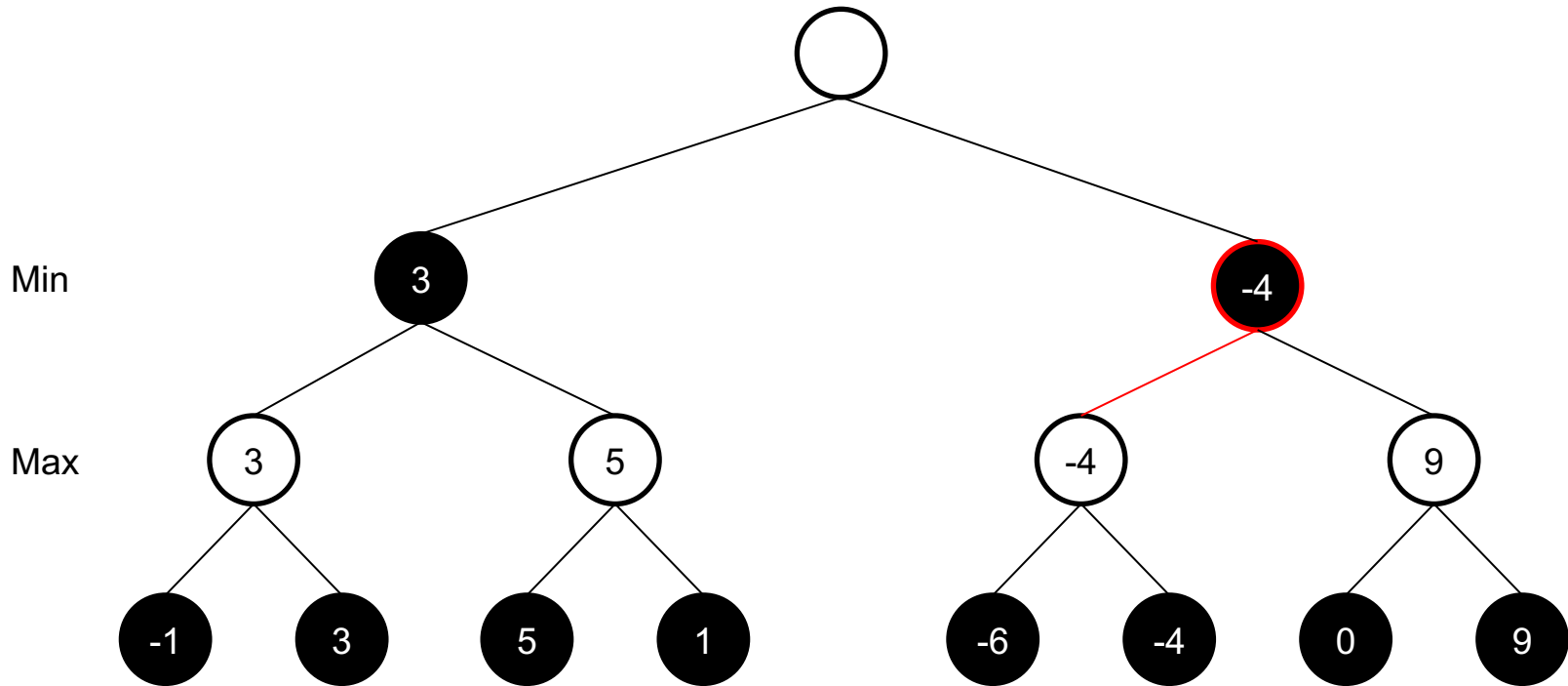












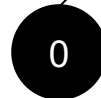
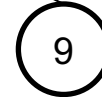
Max



Min



Max



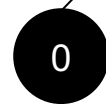
Max

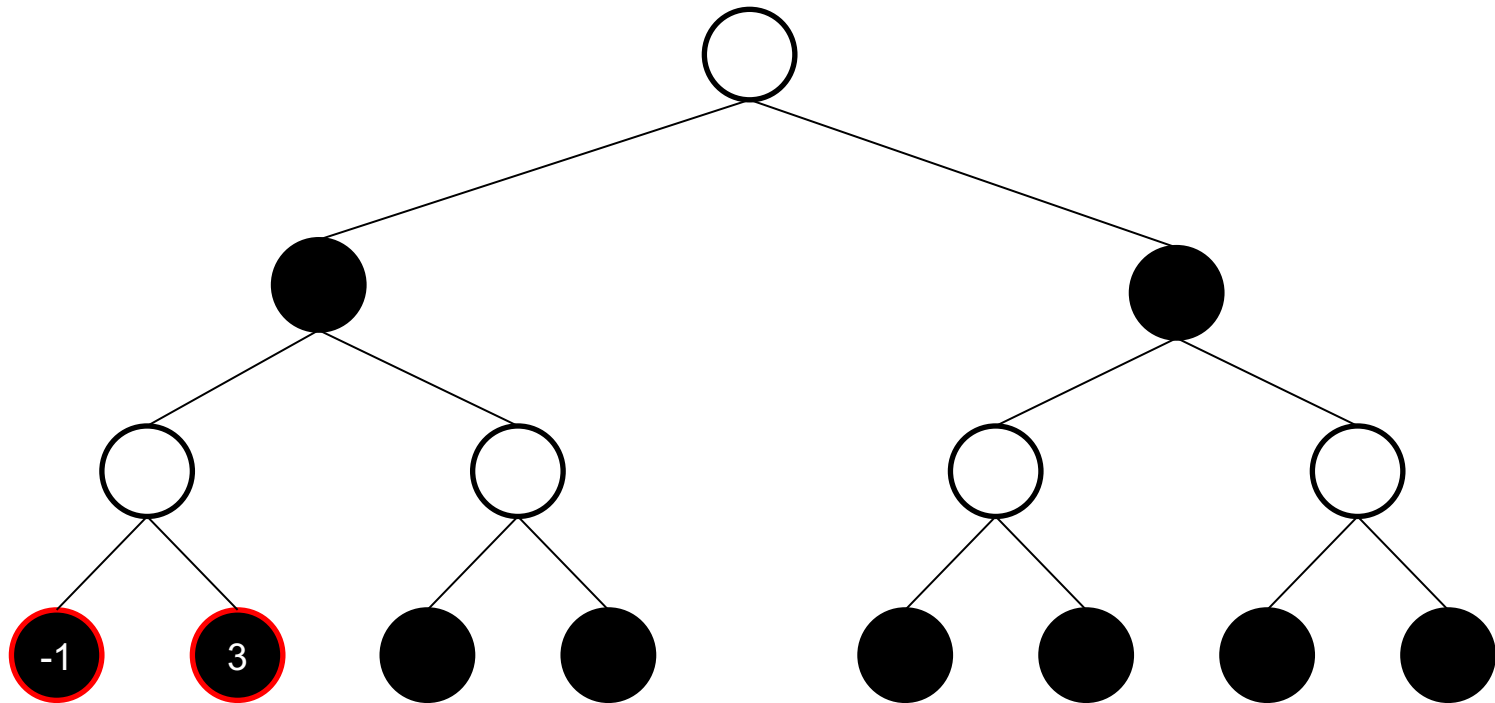


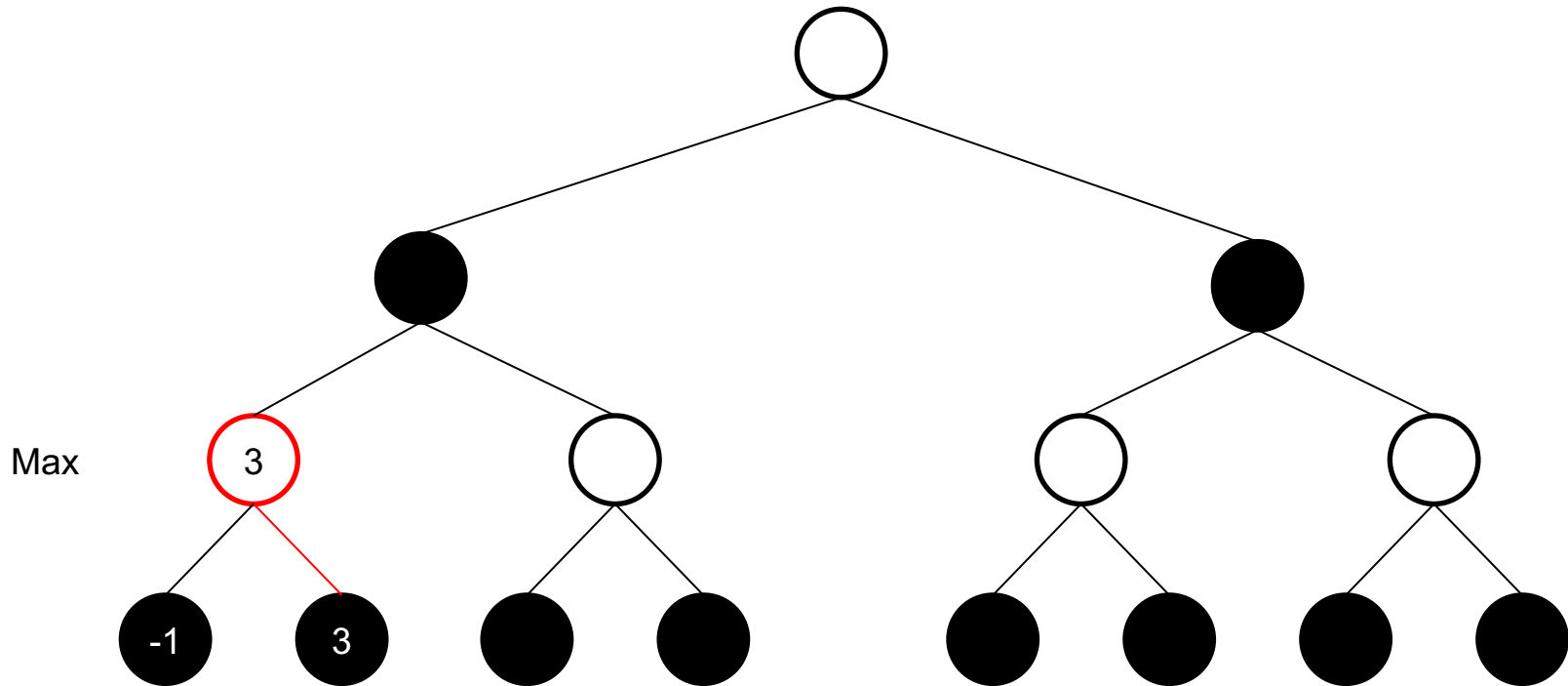
Min

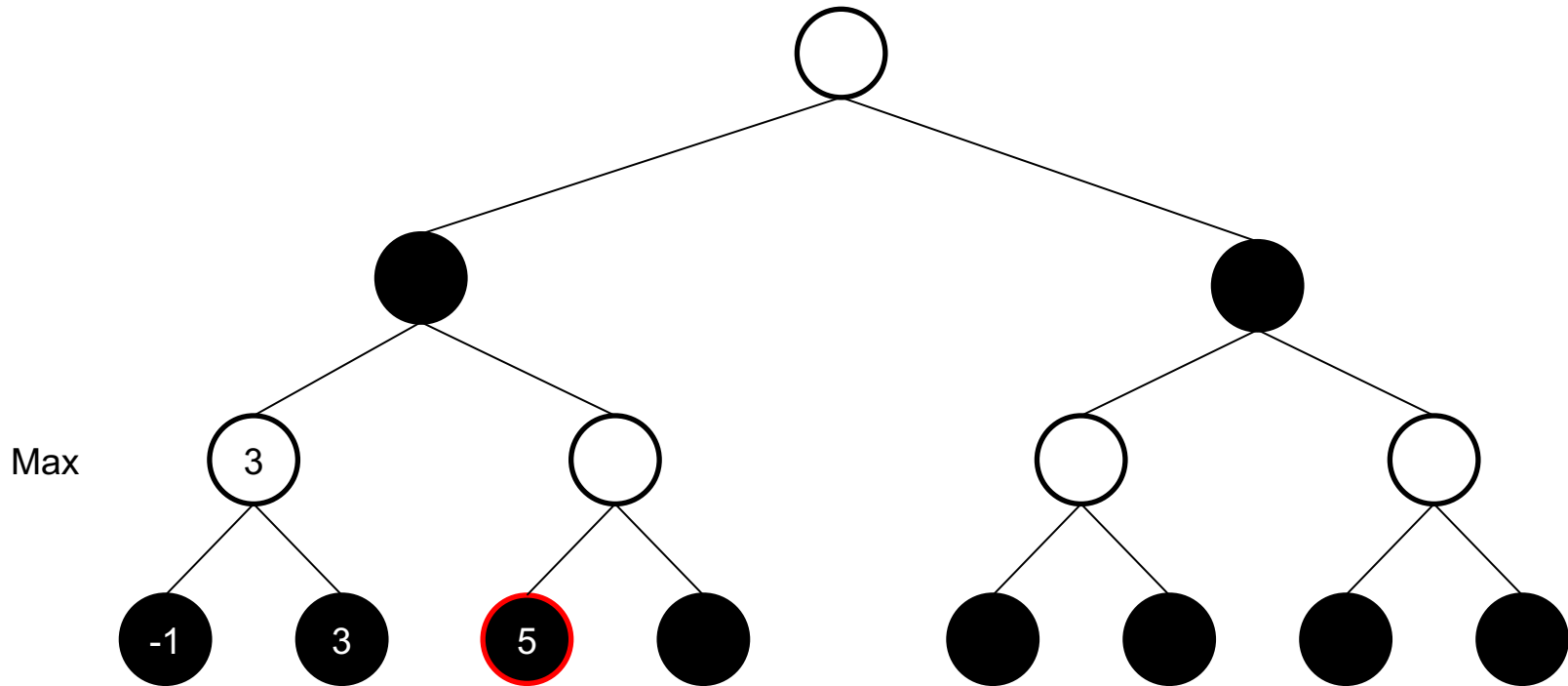


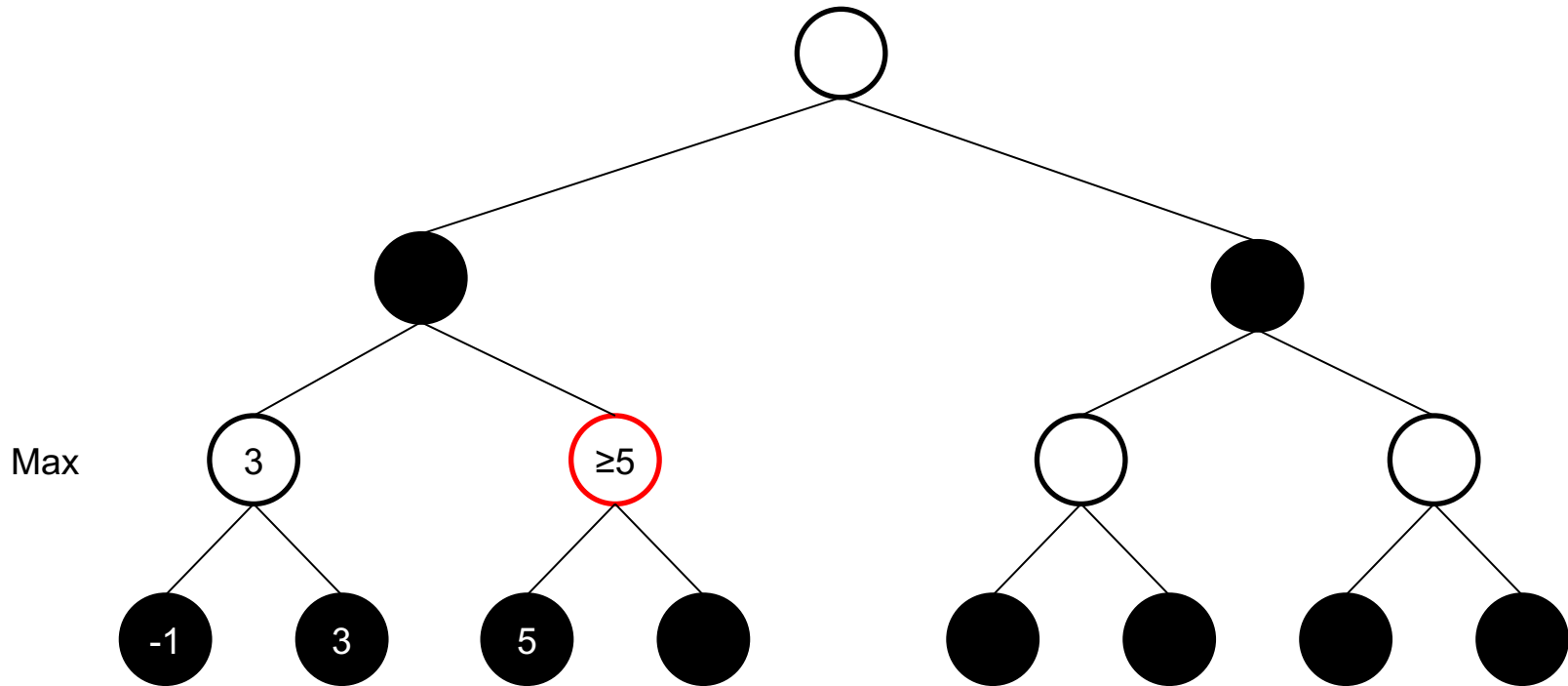
Max

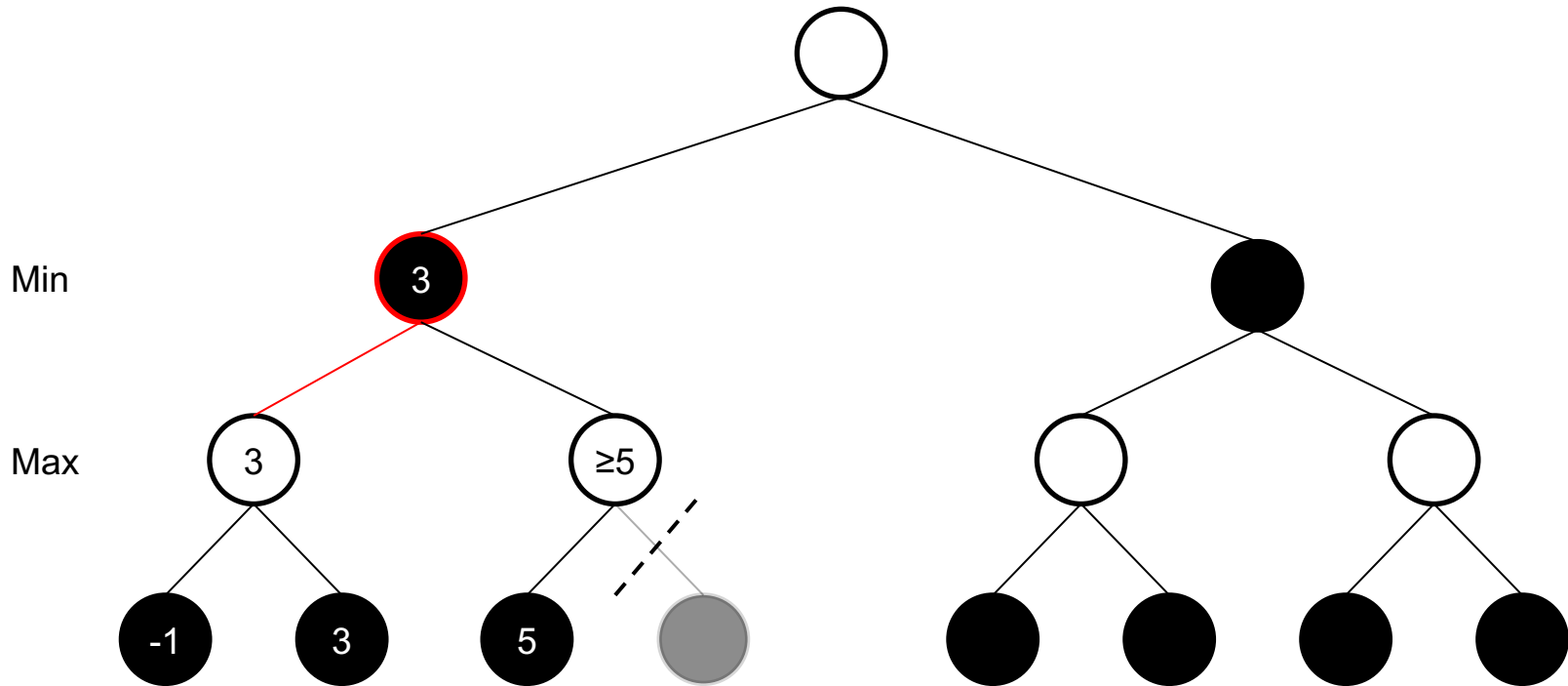








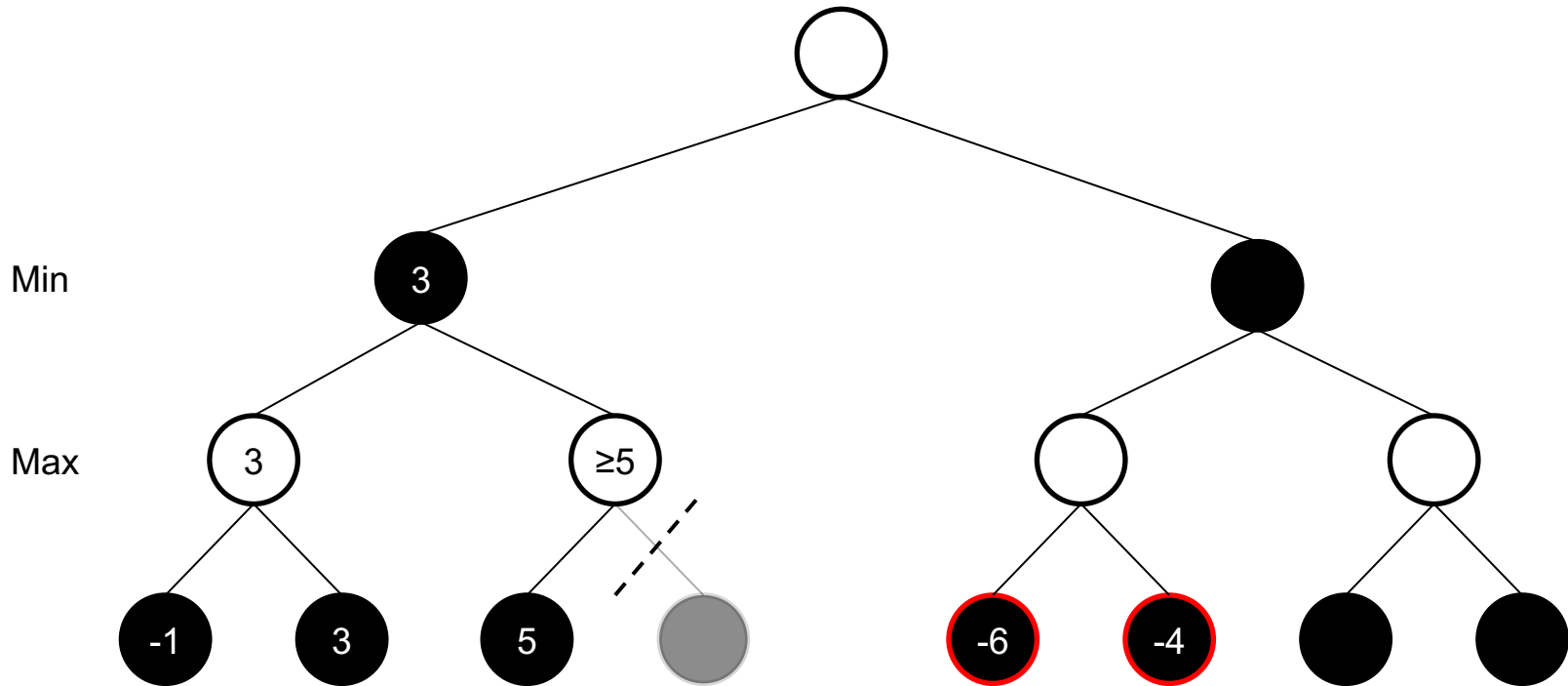


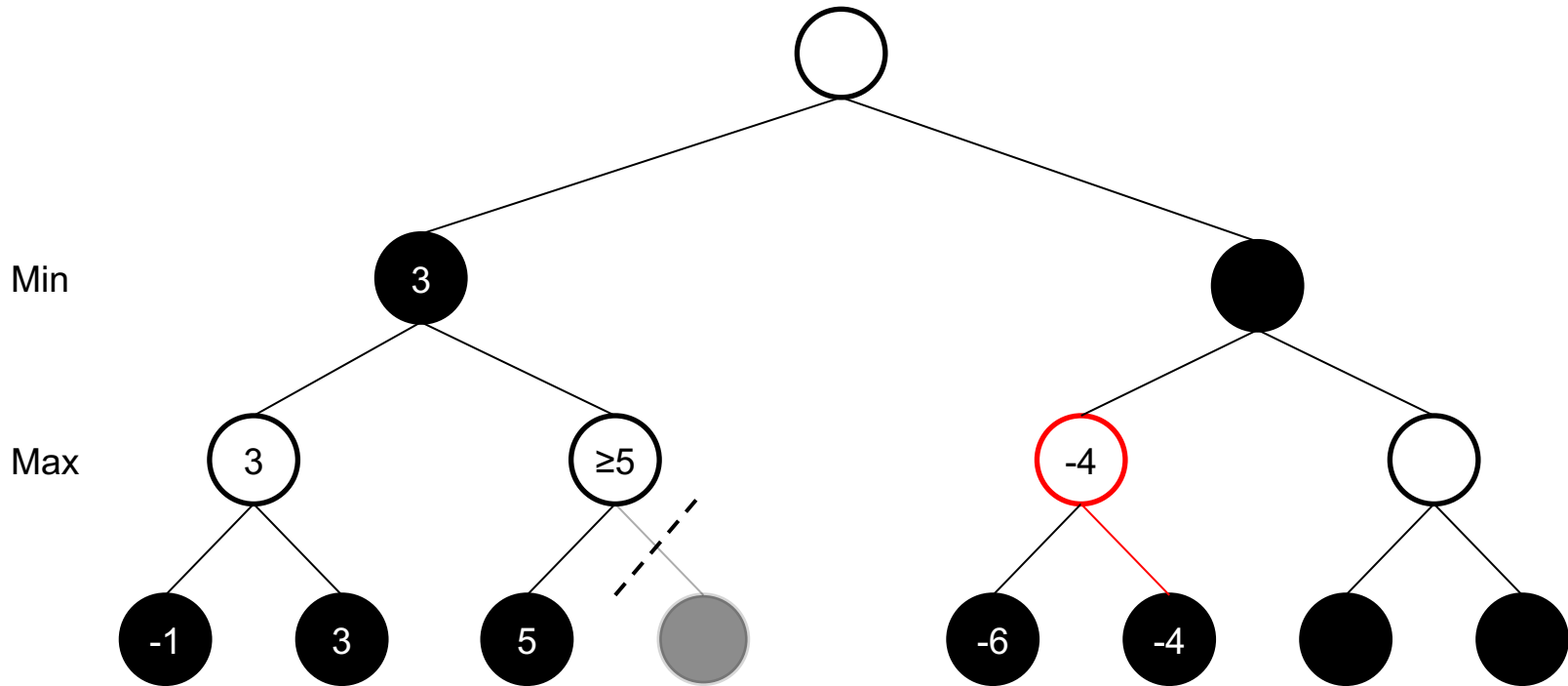


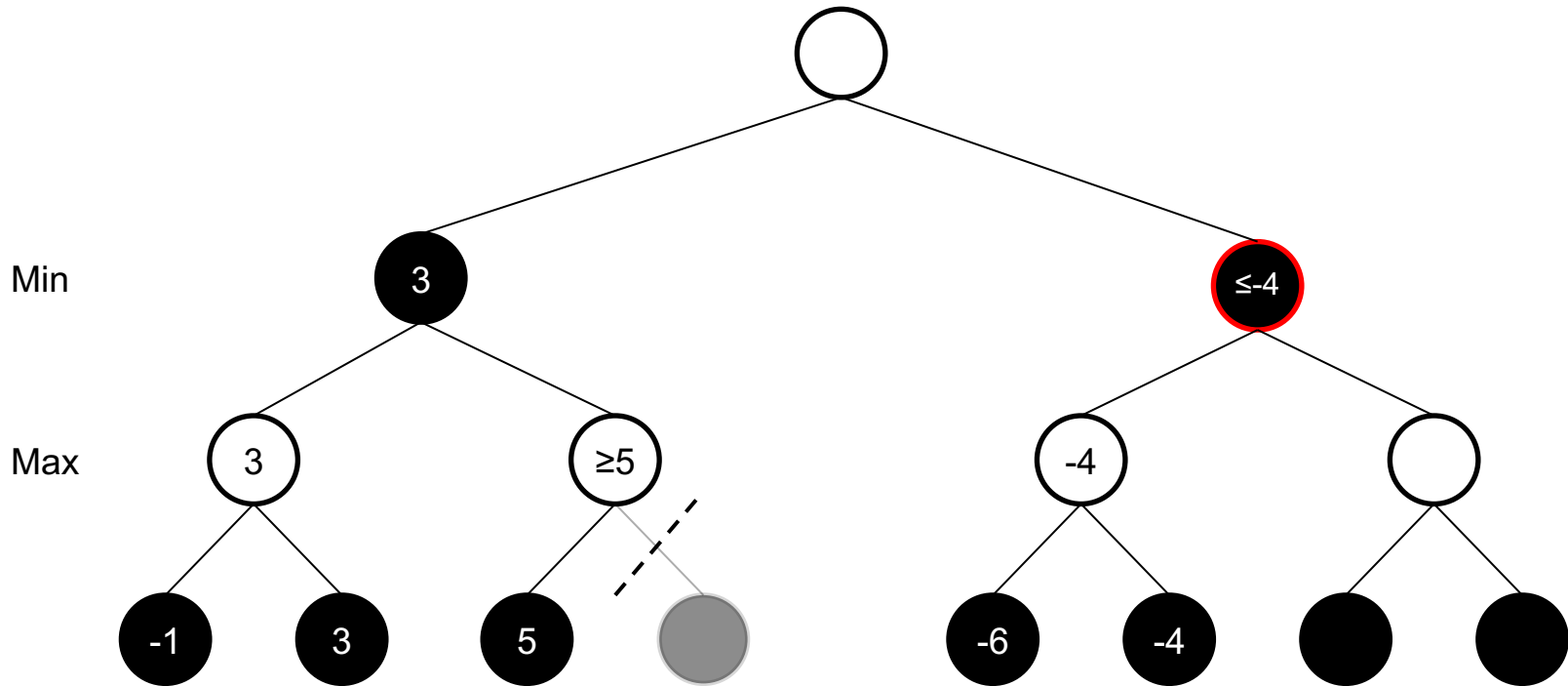
Min

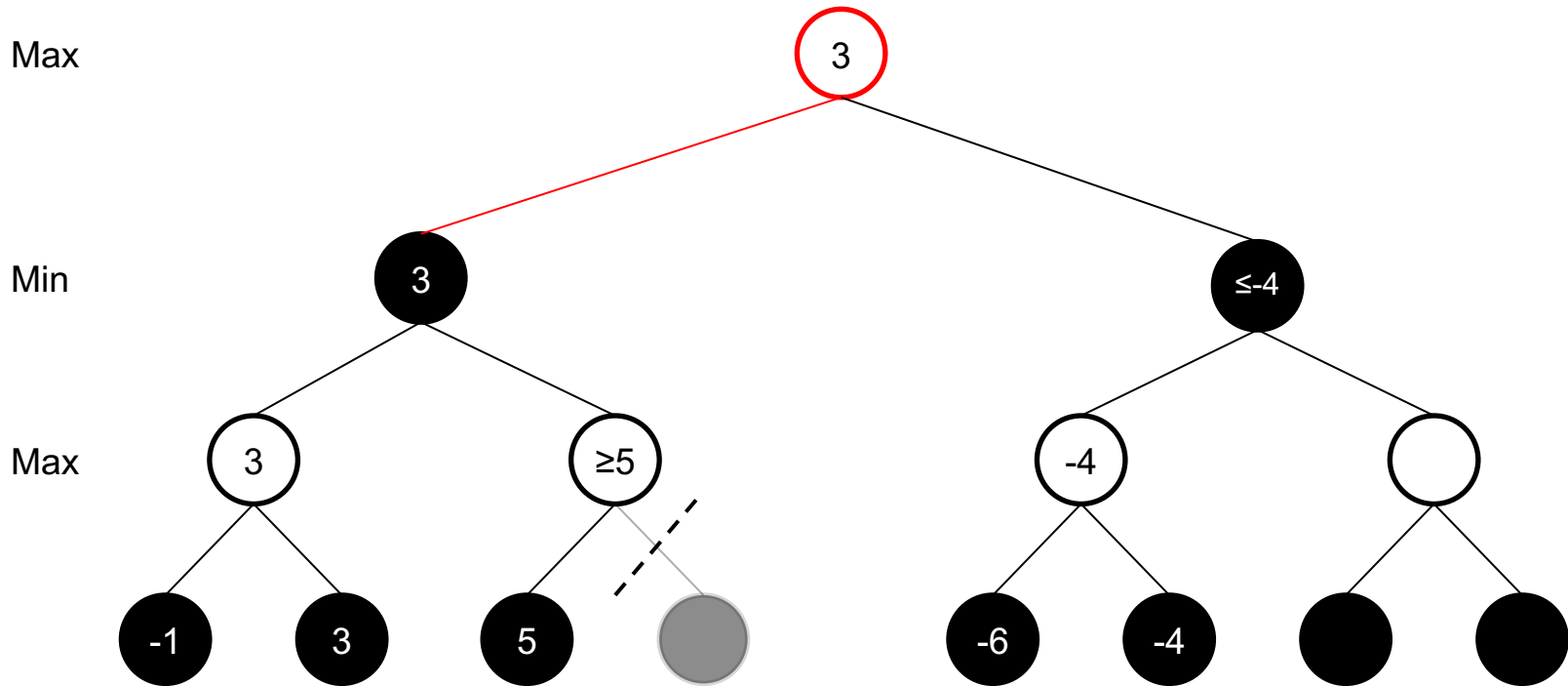
Max











Max



Min



Max



Max



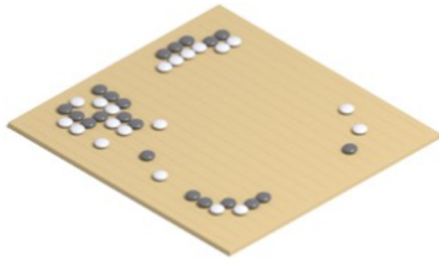
Min



Max



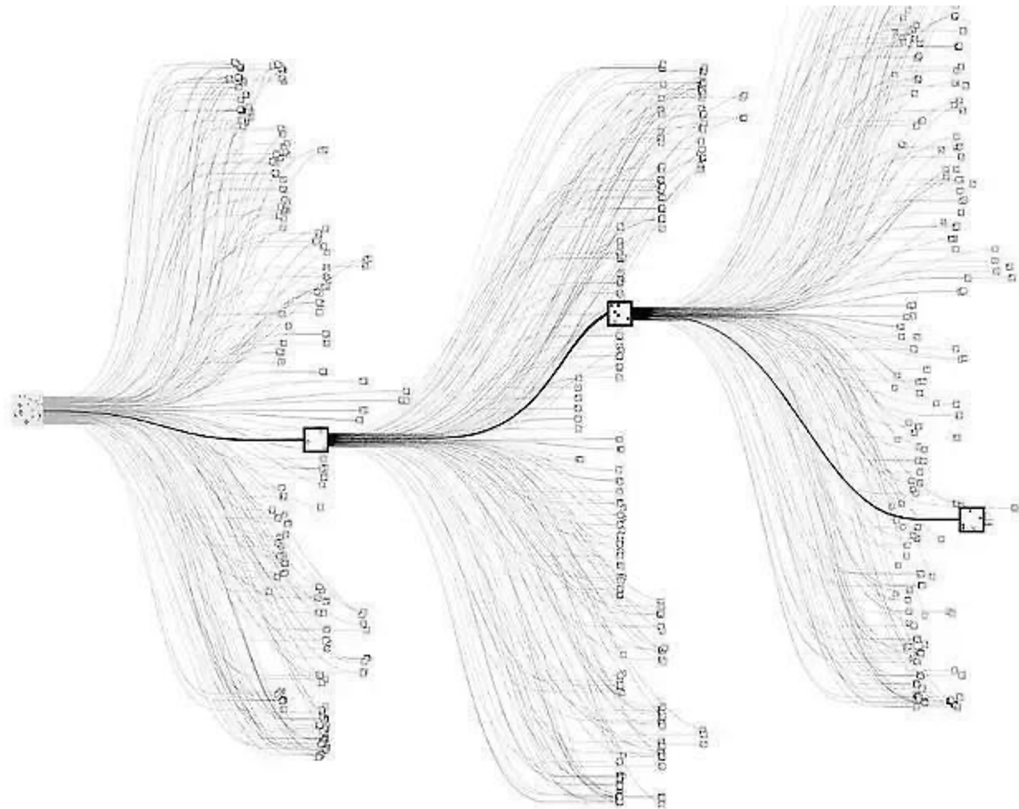
Go



- Branching Factor: 250
- Game Length: 150



Monte Carlo Tree Search



Monte Carlo Tree Search

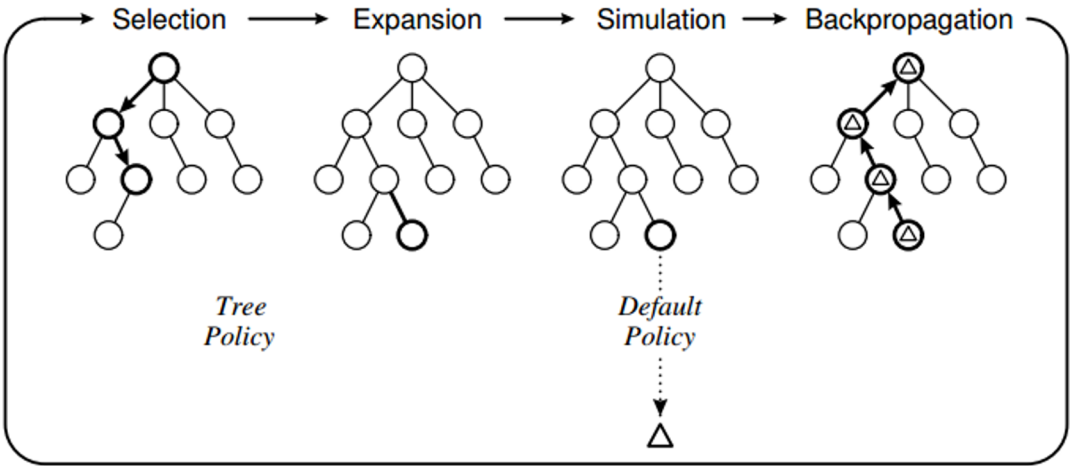
1. Selection (Tree Traversal)

$$UCB1(s_i) = \frac{w_i}{n_i} + C \sqrt{\frac{\ln N}{n_i}}$$

2. Expansion

3. Simulation (Rollout)

4. Backpropagation



Monte Carlo Tree Search

Exploitation vs. Exploration

- Exploit promising actions
- Explore little known actions



Monte Carlo Tree Search

Exploitation vs. Exploration

- Exploit promising actions
- Explore little known actions

$$UCB1(s_i) = \frac{w_i}{n_i} + C \sqrt{\frac{\ln N}{n_i}}$$



Monte Carlo Tree Search

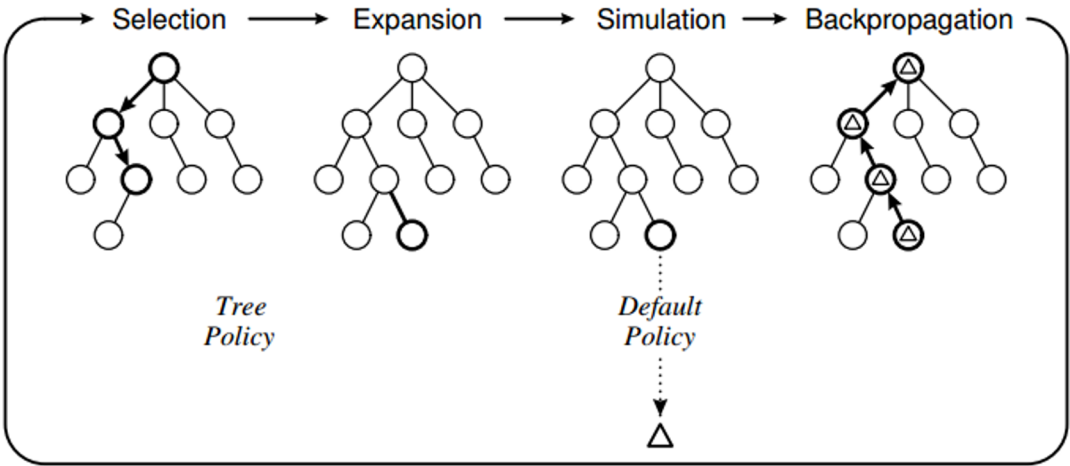
1. Selection (Tree Traversal)

$$UCB1(s_i) = \frac{w_i}{n_i} + C \sqrt{\frac{\ln N}{n_i}}$$

2. Expansion

3. Simulation (Rollout)

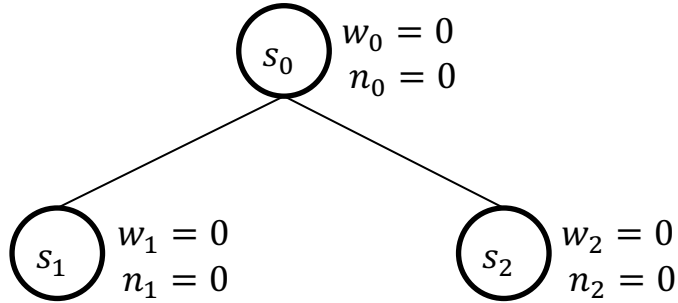
4. Backpropagation



MCTS

$$UCB1(s_i) = \frac{w_i}{n_i} + 2\sqrt{\frac{\ln N}{n_i}}$$

1. Selection
2. Expansion
3. Simulation
4. Backpropagation



MCTS

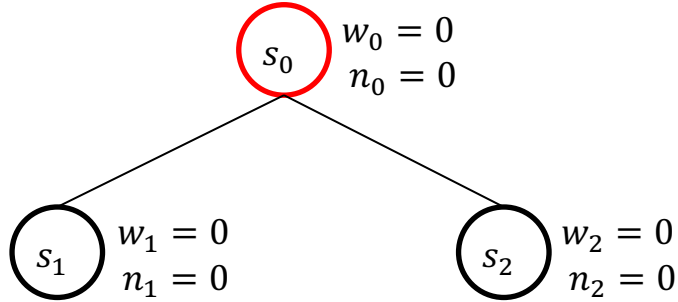
$$UCB1(s_i) = \frac{w_i}{n_i} + 2\sqrt{\frac{\ln N}{n_i}}$$

1. Selection

2. Expansion

3. Simulation

4. Backpropagation



MCTS

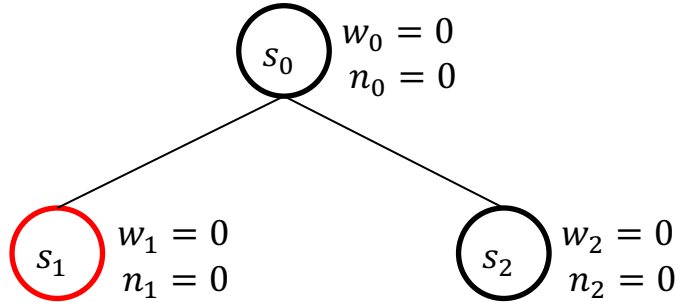
$$UCB1(s_i) = \frac{w_i}{n_i} + 2\sqrt{\frac{\ln N}{n_i}}$$

1. Selection

2. Expansion

3. Simulation

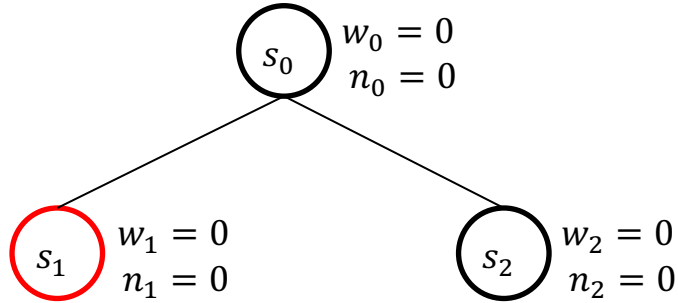
4. Backpropagation



MCTS

$$UCB1(s_i) = \frac{w_i}{n_i} + 2\sqrt{\frac{\ln N}{n_i}}$$

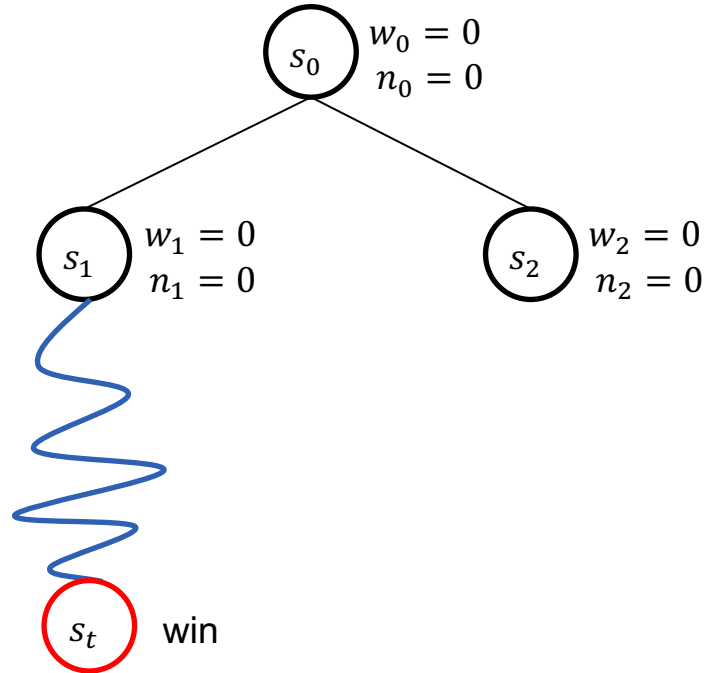
1. Selection
2. Expansion
3. Simulation
4. Backpropagation



MCTS

$$UCB1(s_i) = \frac{w_i}{n_i} + 2\sqrt{\frac{\ln N}{n_i}}$$

1. Selection
2. Expansion
3. Simulation
4. Backpropagation

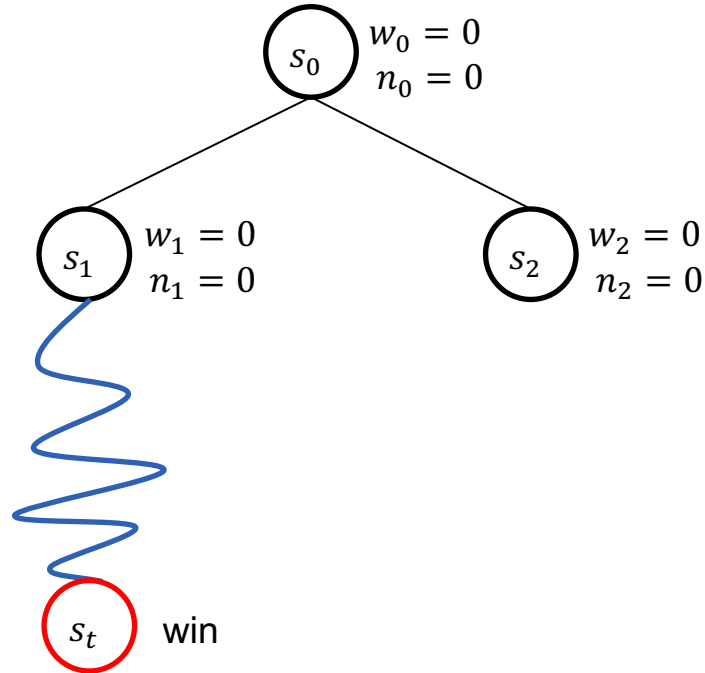


MCTS

$$UCB1(s_i) = \frac{w_i}{n_i} + 2\sqrt{\frac{\ln N}{n_i}}$$

1. Selection
2. Expansion
3. Simulation

4. Backpropagation

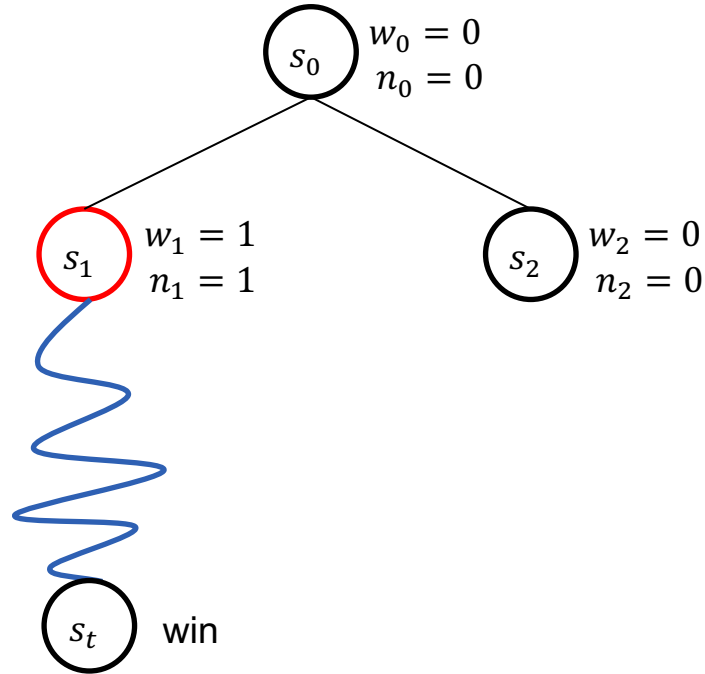


MCTS

$$UCB1(s_i) = \frac{w_i}{n_i} + 2 \sqrt{\frac{\ln N}{n_i}}$$

1. Selection
2. Expansion
3. Simulation

4. Backpropagation

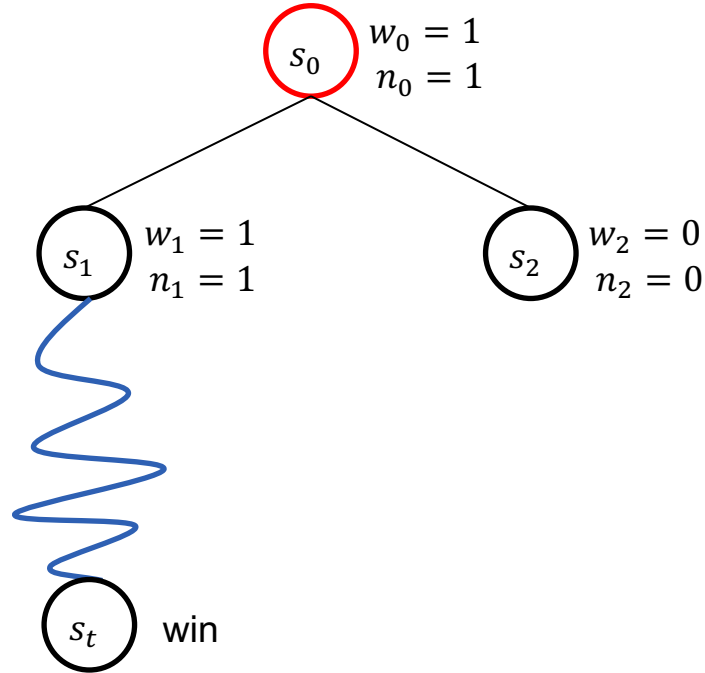


MCTS

$$UCB1(s_i) = \frac{w_i}{n_i} + 2\sqrt{\frac{\ln N}{n_i}}$$

1. Selection
2. Expansion
3. Simulation

4. Backpropagation



MCTS

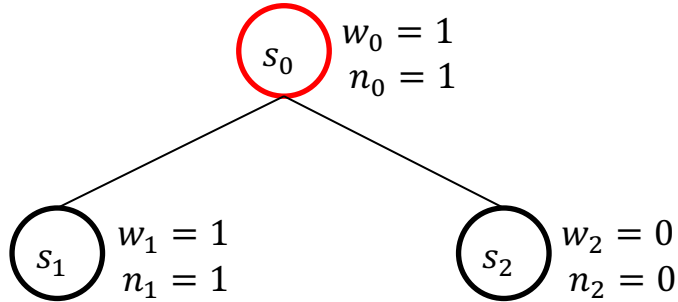
$$UCB1(s_i) = \frac{w_i}{n_i} + 2\sqrt{\frac{\ln N}{n_i}}$$

1. Selection

2. Expansion

3. Simulation

4. Backpropagation



MCTS

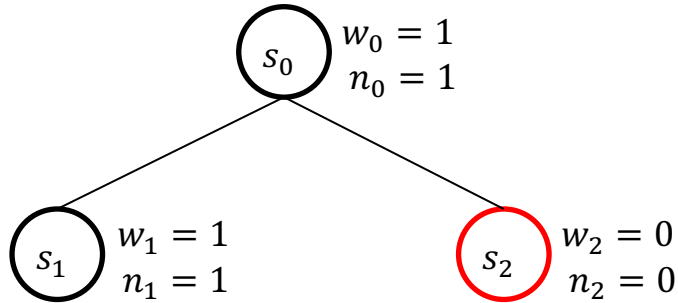
$$UCB1(s_i) = \frac{w_i}{n_i} + 2\sqrt{\frac{\ln N}{n_i}}$$

1. Selection

2. Expansion

3. Simulation

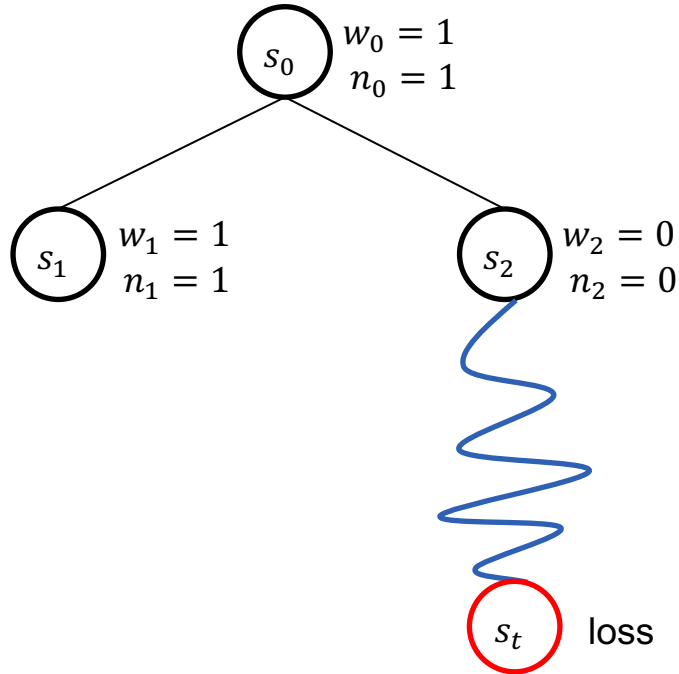
4. Backpropagation



MCTS

$$UCB1(s_i) = \frac{w_i}{n_i} + 2\sqrt{\frac{\ln N}{n_i}}$$

1. Selection
2. Expansion
3. Simulation
4. Backpropagation

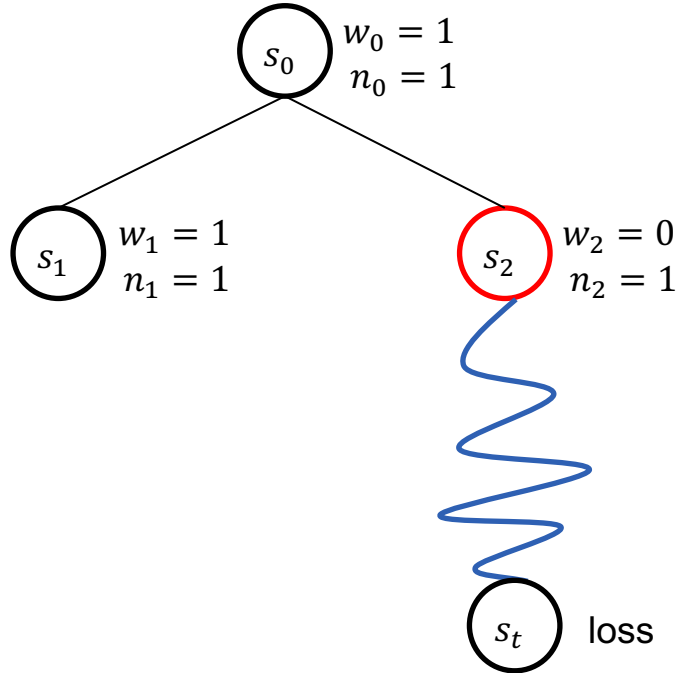


MCTS

$$UCB1(s_i) = \frac{w_i}{n_i} + 2 \sqrt{\frac{\ln N}{n_i}}$$

1. Selection
2. Expansion
3. Simulation

4. Backpropagation

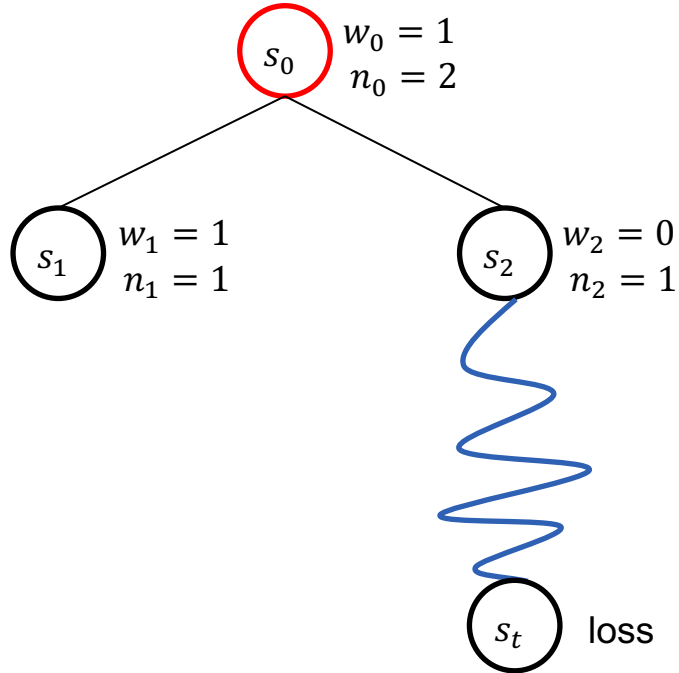


MCTS

$$UCB1(s_i) = \frac{w_i}{n_i} + 2 \sqrt{\frac{\ln N}{n_i}}$$

1. Selection
2. Expansion
3. Simulation

4. Backpropagation



MCTS

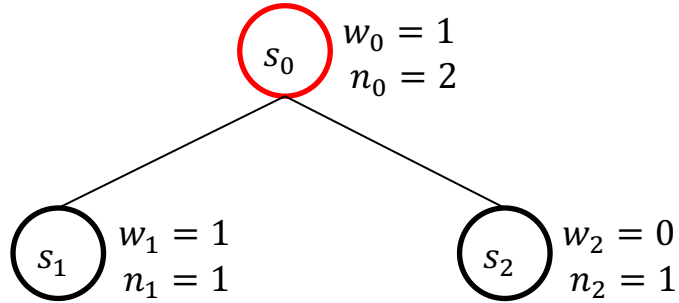
$$UCB1(s_i) = \frac{w_i}{n_i} + 2\sqrt{\frac{\ln N}{n_i}}$$

1. Selection

2. Expansion

3. Simulation

4. Backpropagation



MCTS

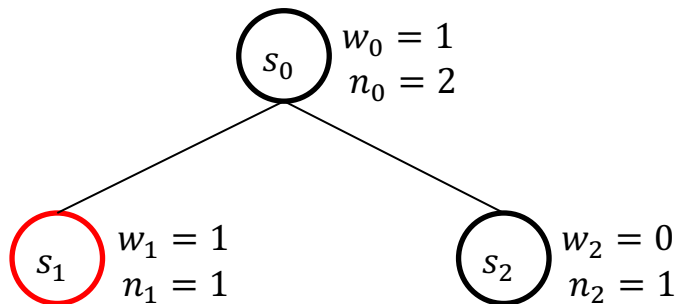
$$UCB1(s_i) = \frac{w_i}{n_i} + 2\sqrt{\frac{\ln N}{n_i}}$$

1. Selection

2. Expansion

3. Simulation

4. Backpropagation



MCTS

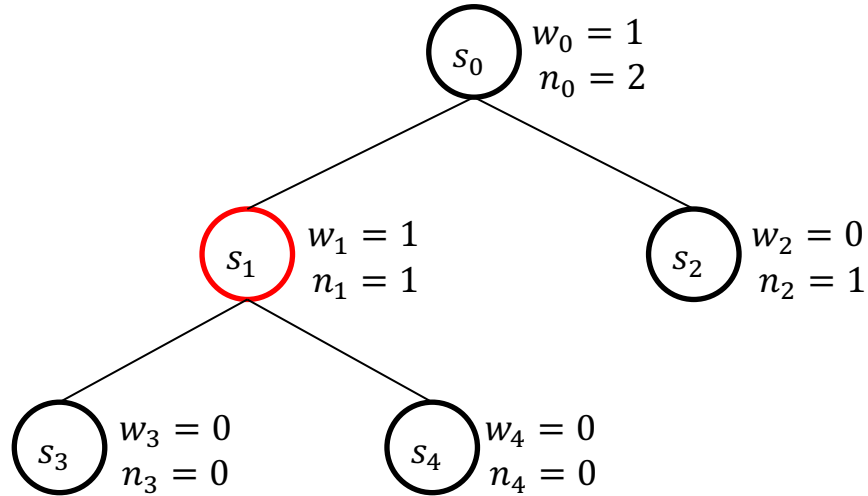
$$UCB1(s_i) = \frac{w_i}{n_i} + 2 \sqrt{\frac{\ln N}{n_i}}$$

1. Selection

2. Expansion

3. Simulation

4. Backpropagation



MCTS

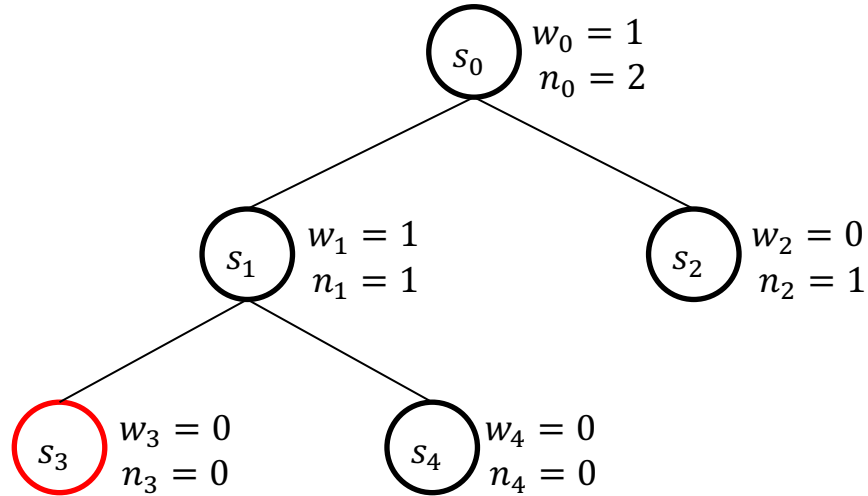
$$UCB1(s_i) = \frac{w_i}{n_i} + 2\sqrt{\frac{\ln N}{n_i}}$$

1. Selection

2. Expansion

3. Simulation

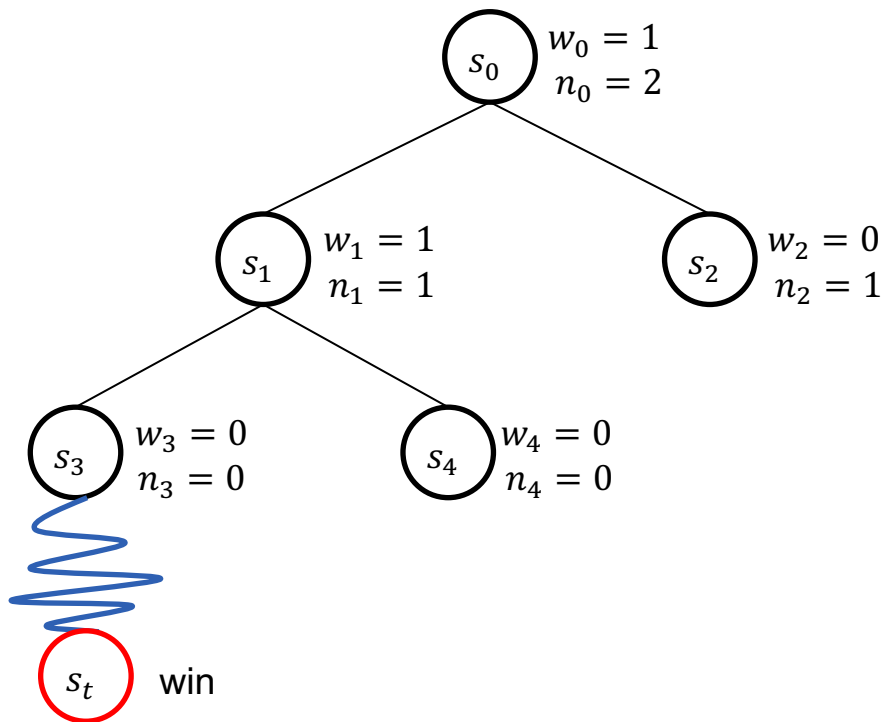
4. Backpropagation



MCTS

$$UCB1(s_i) = \frac{w_i}{n_i} + 2 \sqrt{\frac{\ln N}{n_i}}$$

1. Selection
2. Expansion
3. Simulation
4. Backpropagation

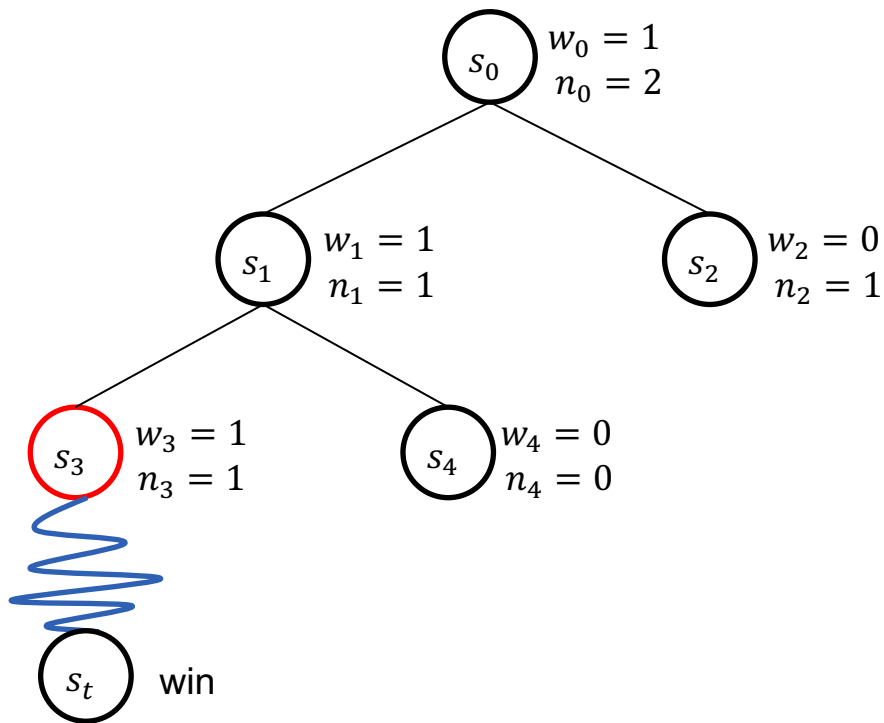


MCTS

$$UCB1(s_i) = \frac{w_i}{n_i} + 2 \sqrt{\frac{\ln N}{n_i}}$$

1. Selection
2. Expansion
3. Simulation

4. Backpropagation

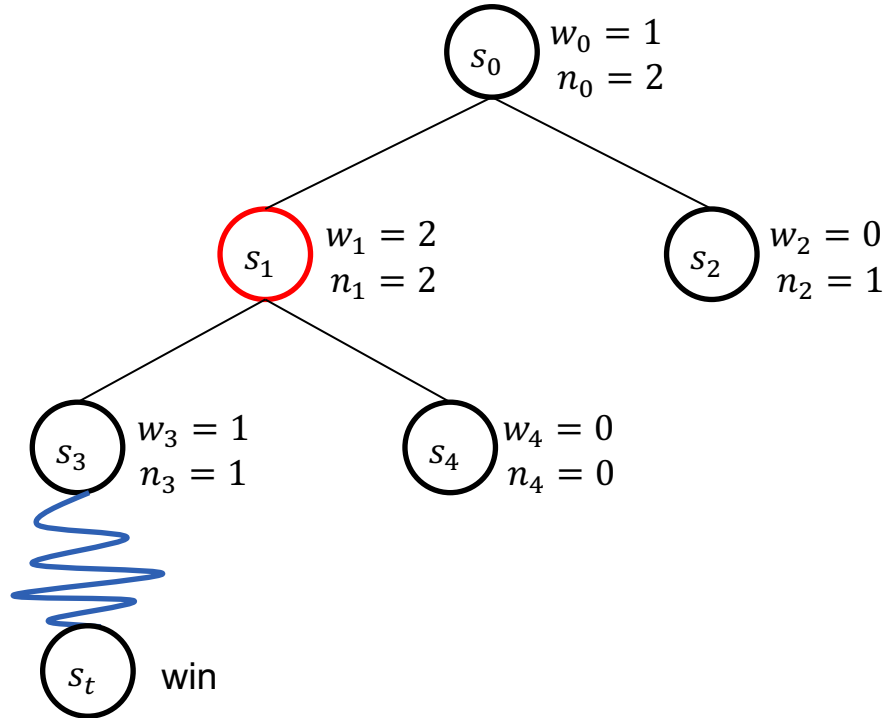


MCTS

$$UCB1(s_i) = \frac{w_i}{n_i} + 2 \sqrt{\frac{\ln N}{n_i}}$$

1. Selection
2. Expansion
3. Simulation

4. Backpropagation

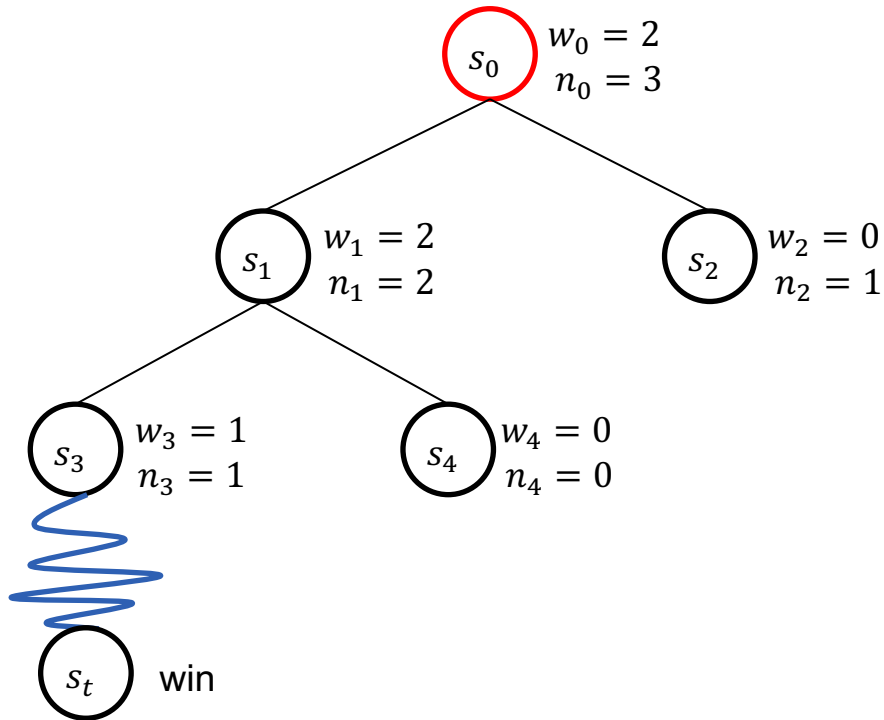


MCTS

$$UCB1(s_i) = \frac{w_i}{n_i} + 2 \sqrt{\frac{\ln N}{n_i}}$$

1. Selection
2. Expansion
3. Simulation

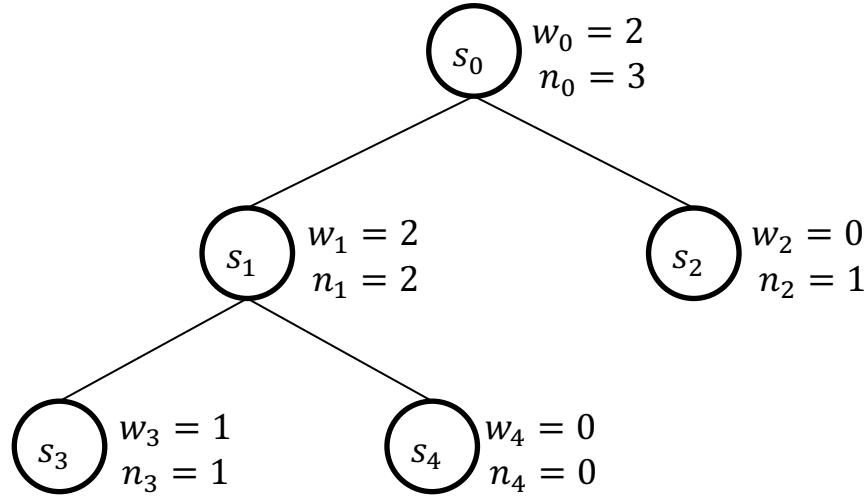
4. Backpropagation



MCTS

$$UCB1(s_i) = \frac{w_i}{n_i} + 2 \sqrt{\frac{\ln N}{n_i}}$$

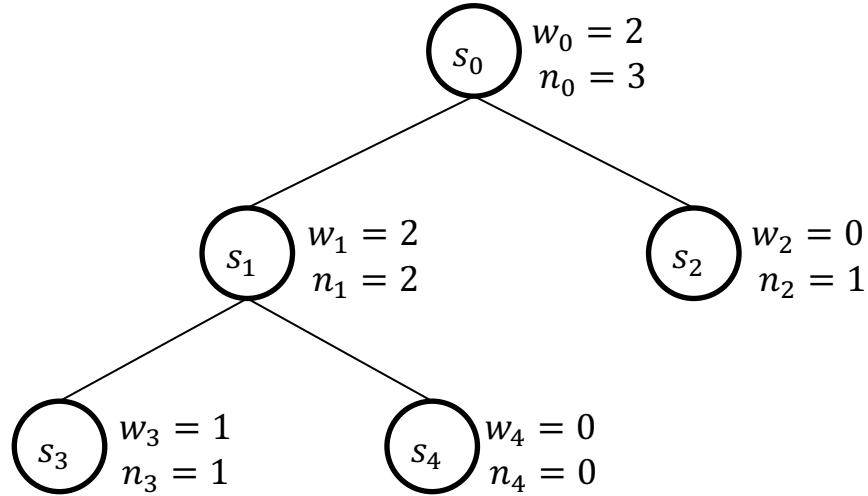
1. Selection
2. Expansion
3. Simulation
4. Backpropagation



MCTS

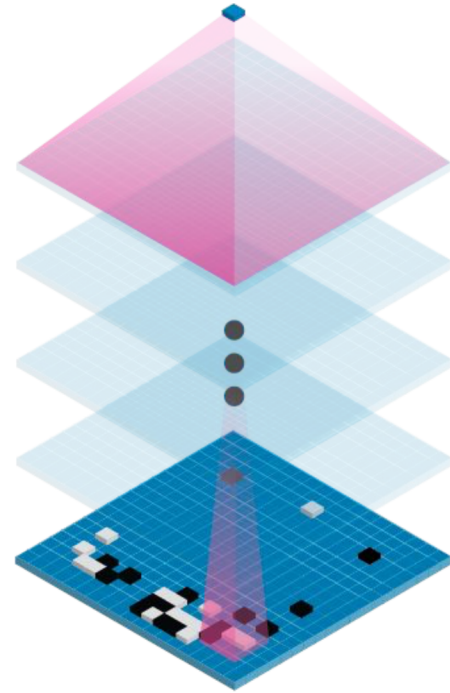
$$UCB1(s_i) = \frac{w_i}{n_i} + 2 \sqrt{\frac{\ln N}{n_i}}$$

1. Selection
2. Expansion
3. Simulation
4. Backpropagation



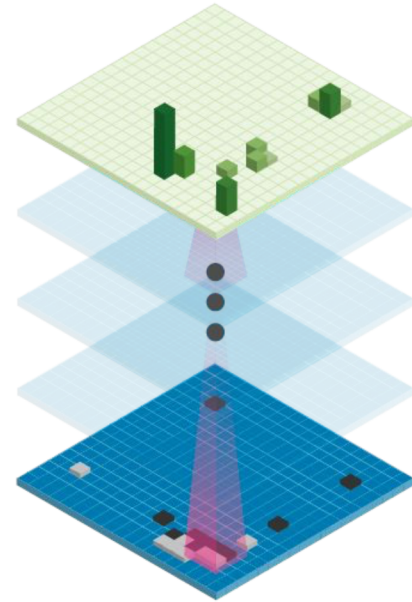
Value Network

How well are we doing?



Policy Network

What are the most likely actions?



AlphaGo

a

Rollout policy

SL policy network

RL policy network

Value network

p_π

p_σ

p_ρ

v_θ

Policy gradient

Classification

Classification

Self Play

Regression

Human expert positions

Self-play positions

Neural network

Data

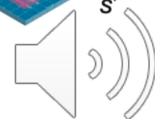
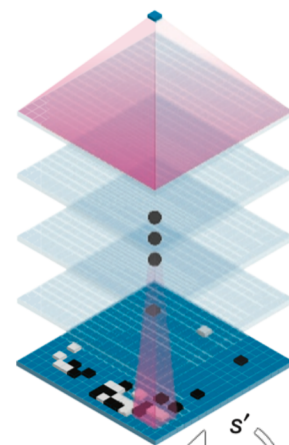
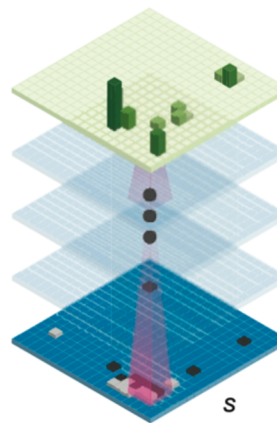
b

Policy network

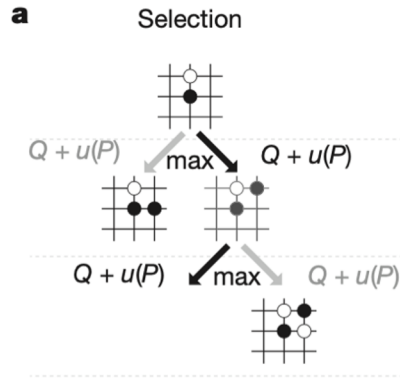
Value network

$p_{\sigma\rho}(a|s)$

$v_\theta(s')$



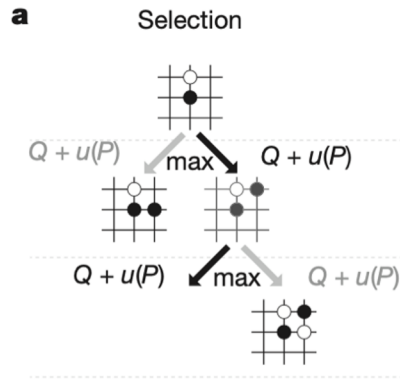
AlphaGo MCTS



$$u(a) = v(a) + p(a) \cdot pb_c$$



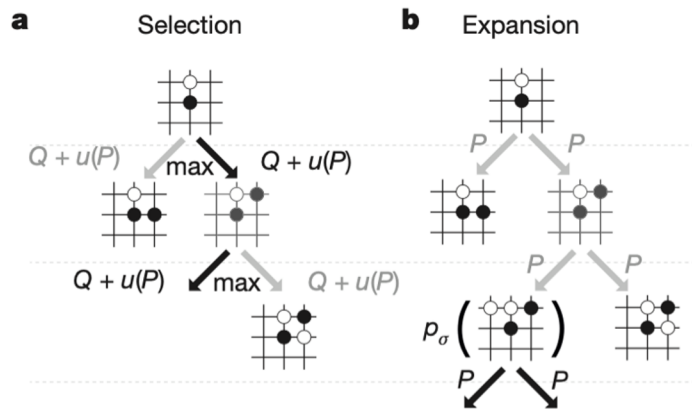
AlphaGo MCTS



$$u(a) = v(a) + p(a) \cdot pb_c$$



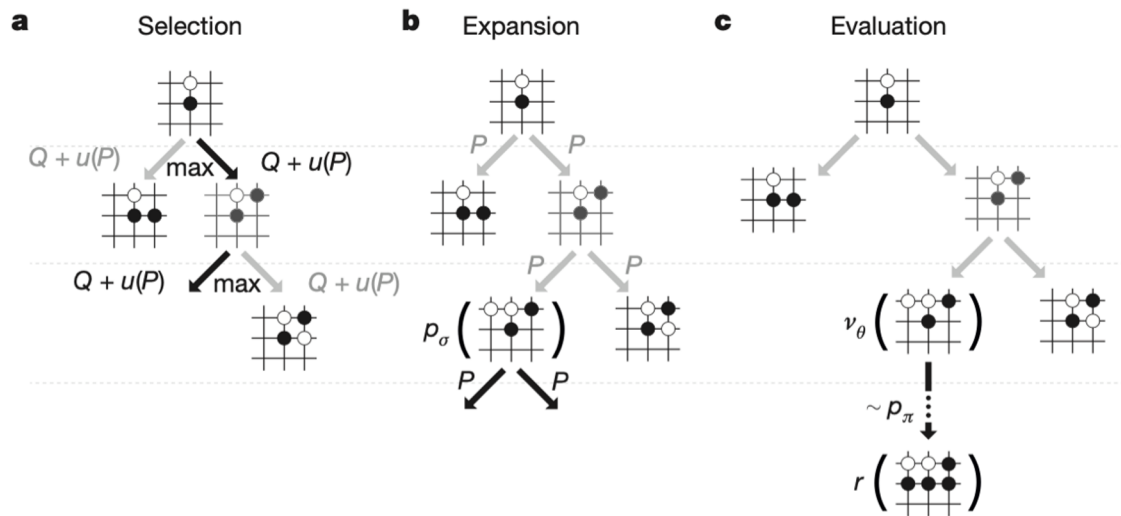
AlphaGo MCTS



$$u(a) = v(a) + p(a) \cdot pb_c$$



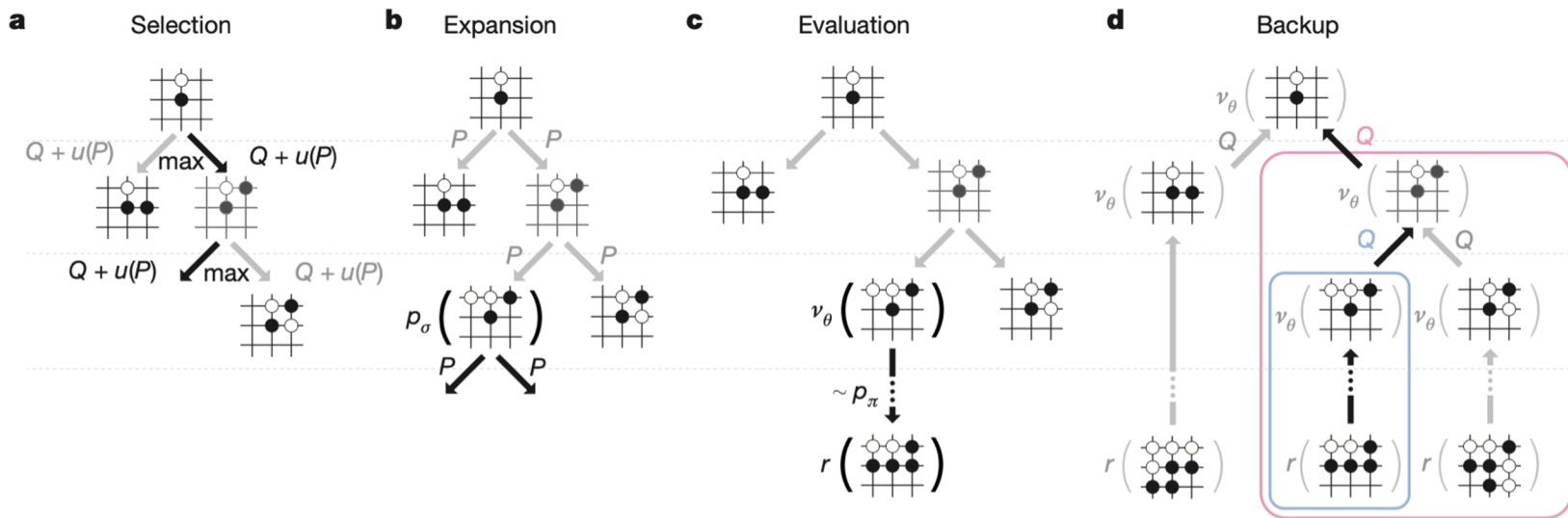
AlphaGo MCTS



$$u(a) = v(a) + p(a) \cdot pb_c$$



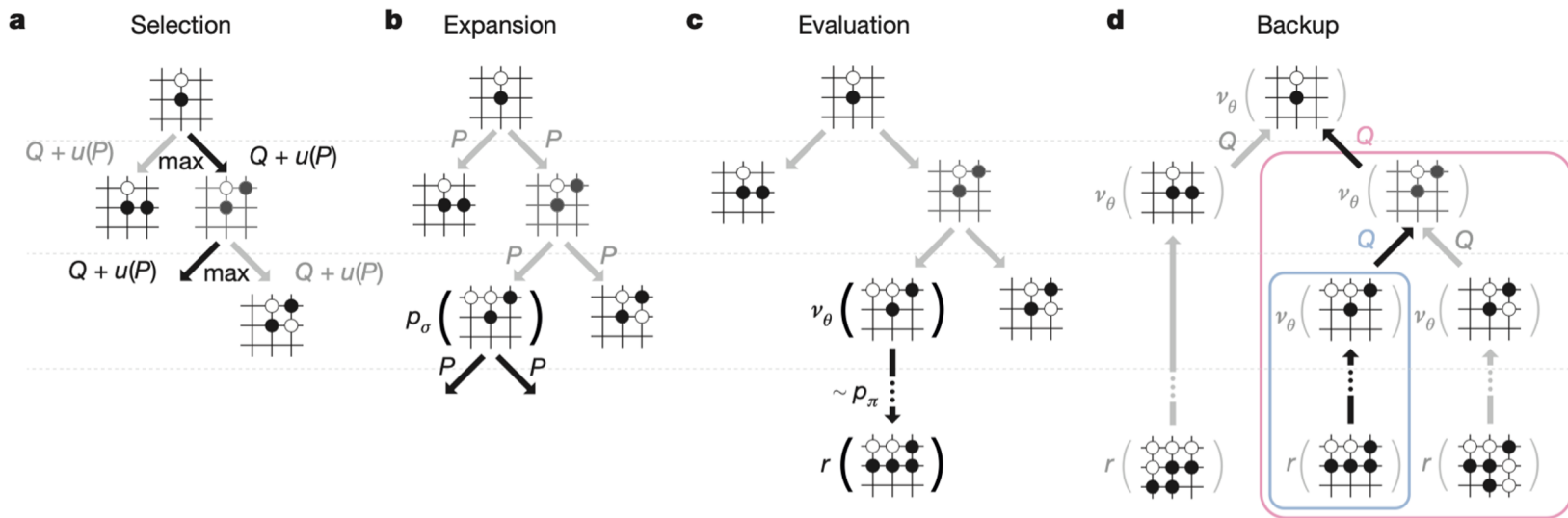
AlphaGo MCTS



$$u(a) = v(a) + p(a) \cdot pb_c$$



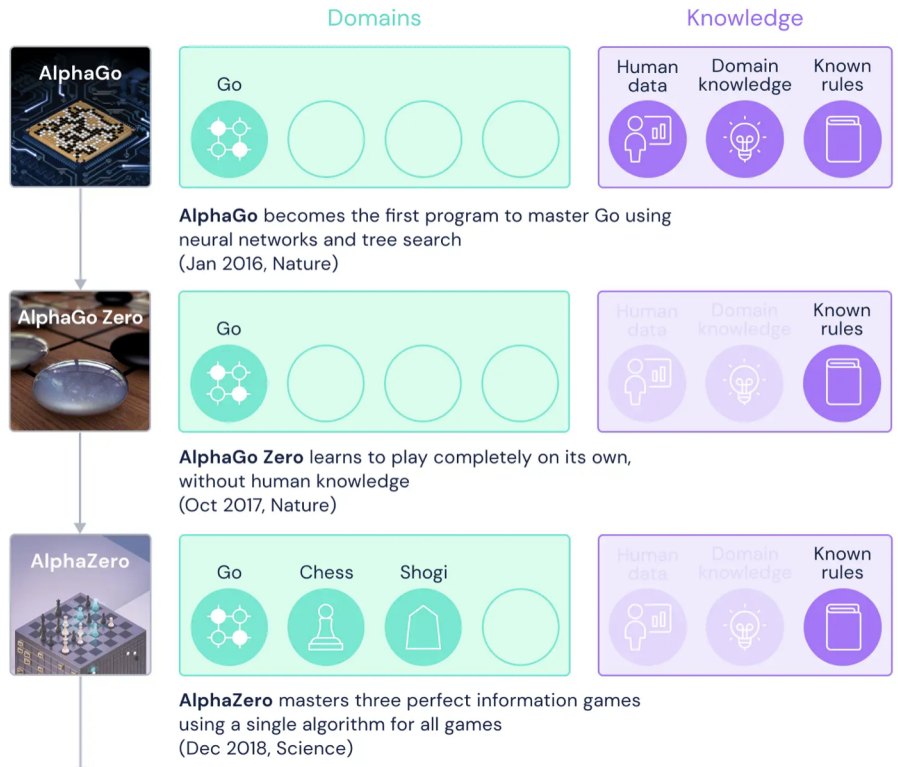
AlphaGo MCTS



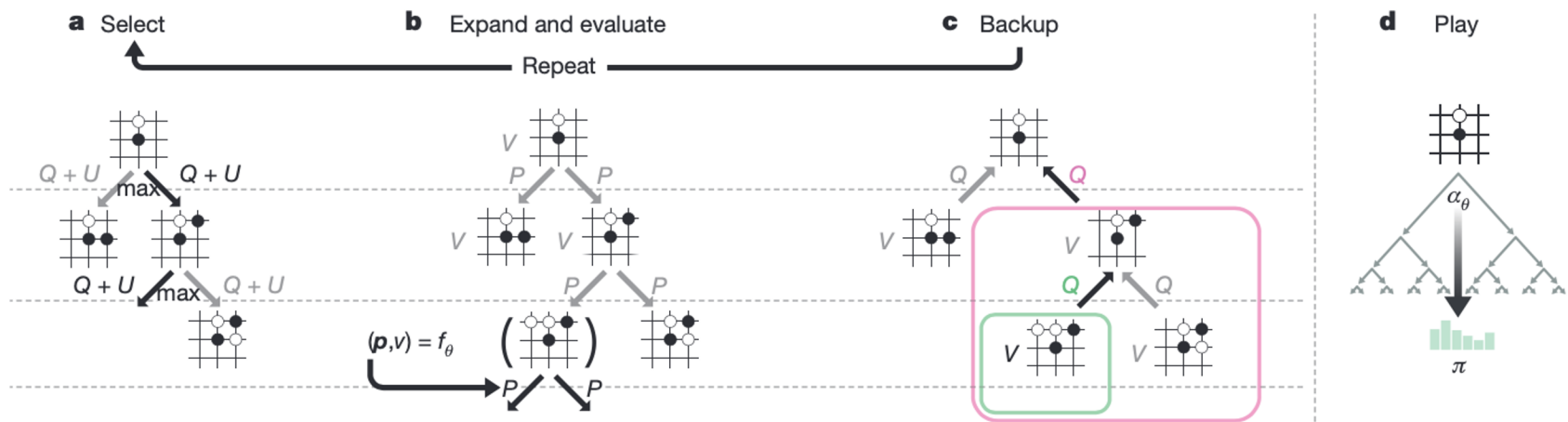
$$u(a) = v(a) + p(a) \cdot pb_c$$



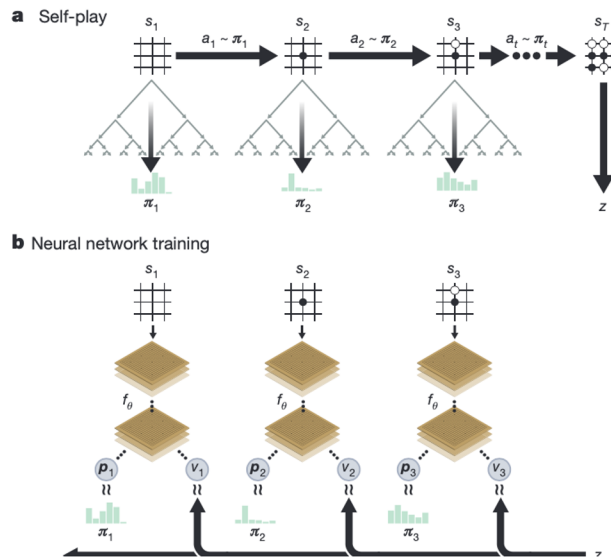




AlphaGo Zero MCTS



AlphaZero Training

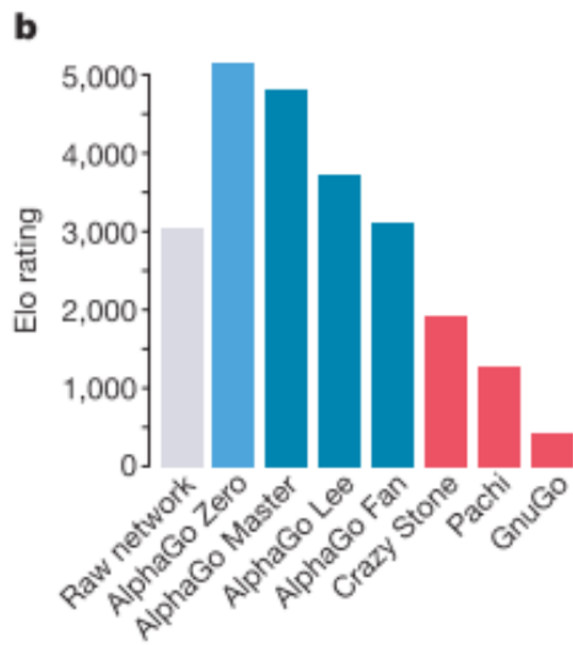


$$(\mathbf{p}, v) = f_{\theta}(s),$$

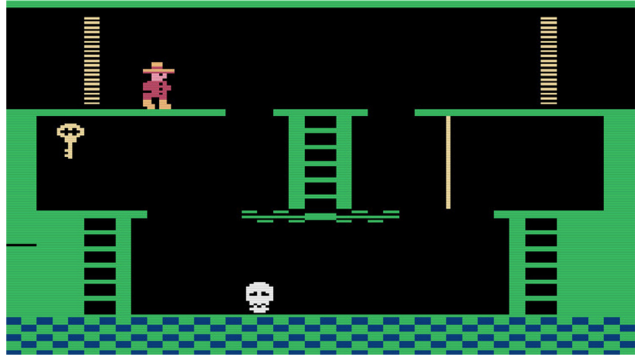
$$l = (z - v)^2 - \boldsymbol{\pi}^{\top} \log \mathbf{p} + c \|\boldsymbol{\theta}\|^2$$



AlphaGo Zero Results



Atari



- Image Input
- No Access to rules

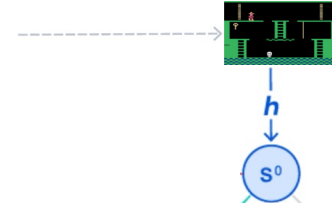


MuZero Planning



MuZero Planning

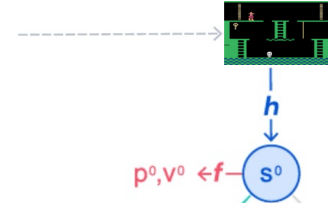
representation $h_{\theta}(o_1, \dots, o_t) = s^0$



MuZero Planning

representation $h_{\theta}(o_1, \dots, o_t) = s^0$

prediction $f_{\theta}(s^k) = \mathbf{p}^k, v^k$

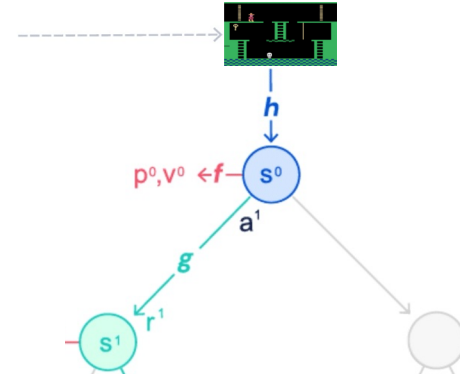


MuZero Planning

representation $h_{\theta}(o_1, \dots, o_t) = s^0$

prediction $f_{\theta}(s^k) = \mathbf{p}^k, v^k$

dynamics $g_{\theta}(s^{k-1}, a^k) = r^k, s^k$



MuZero Planning

representation

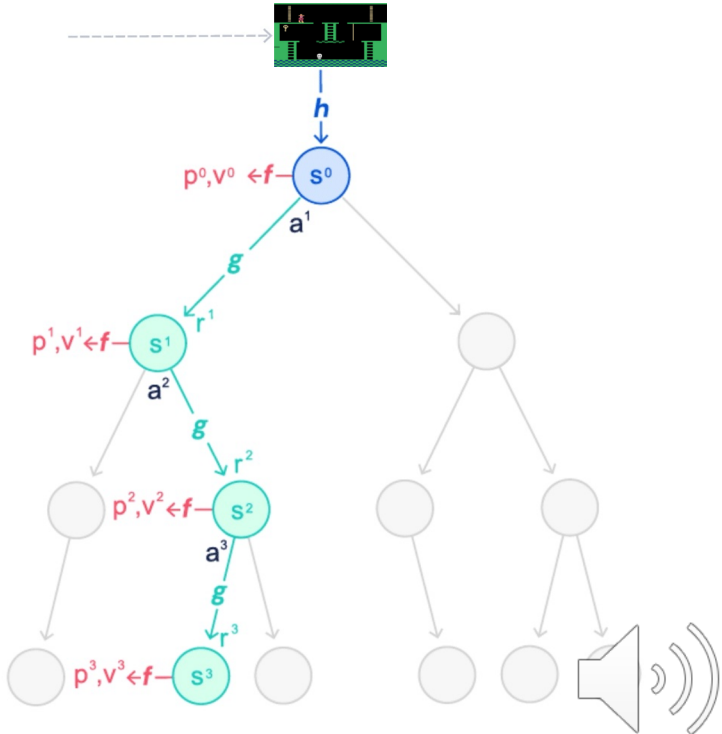
$$h_{\theta}(o_1, \dots, o_t) = s^0$$

prediction

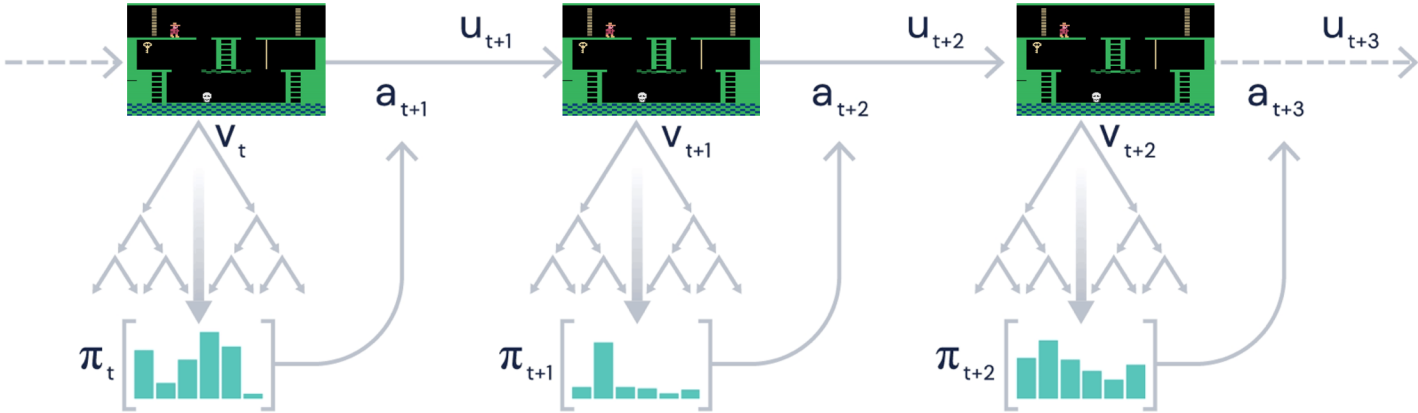
$$f_{\theta}(s^k) = p^k, v^k$$

dynamics

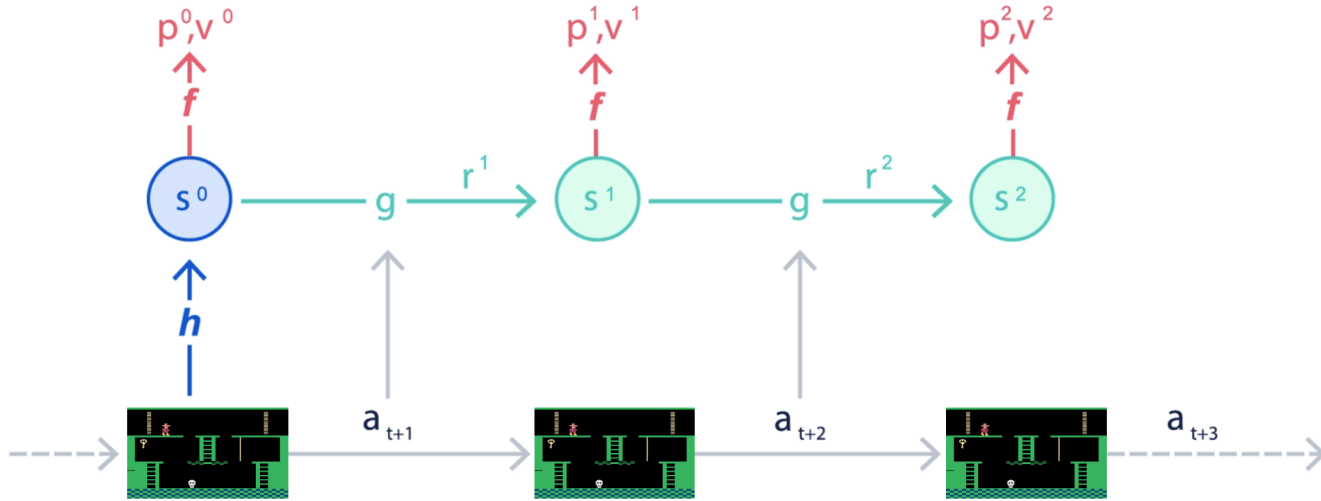
$$g_{\theta}(s^{k-1}, a^k) = r^k, s^k$$



MuZero Training Data Generation



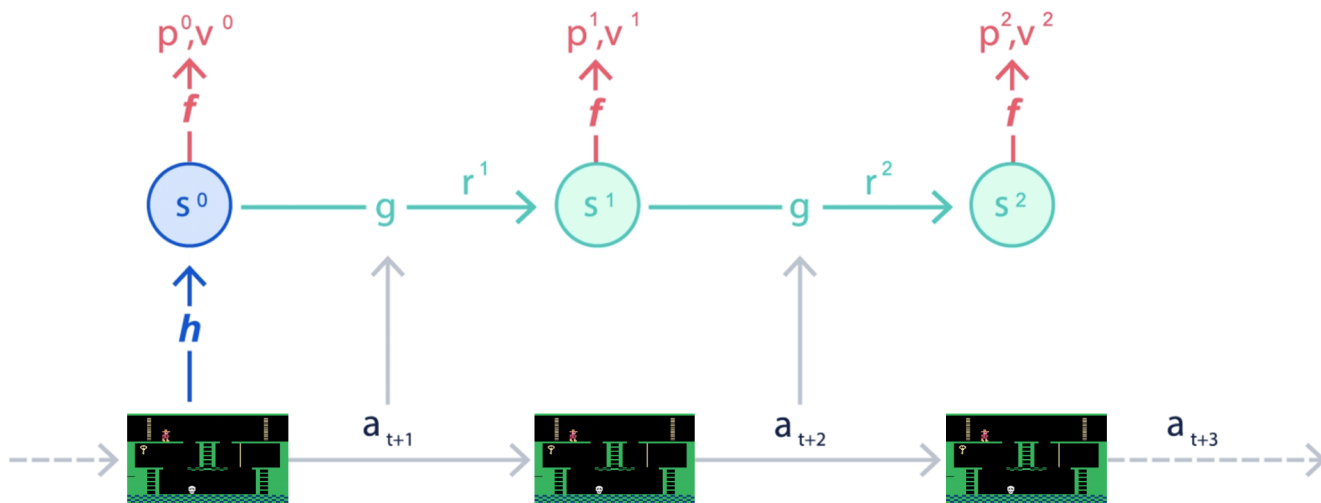
MuZero Training



$$l_t(\theta) = \sum_{k=0}^K l^p(\pi_{t+k}, p_t^k) + \sum_{k=0}^K l^v(z_{t+k}, v_t^k) + \sum_{k=1}^K l^r(u_{t+k}, r_t^k) + c\|\theta\|^2$$



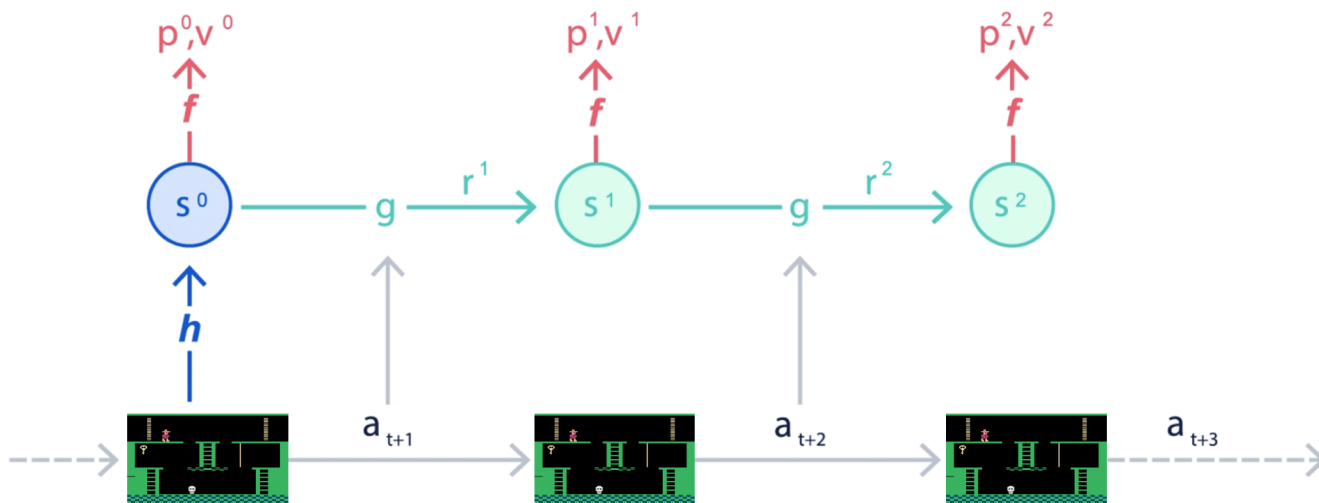
MuZero Training



$$l_t(\theta) = \sum_{k=0}^K l^p(\pi_{t+k}, p_t^k) + \sum_{k=0}^K l^v(z_{t+k}, v_t^k) + \sum_{k=1}^K l^r(u_{t+k}, r_t^k) + c\|\theta\|^2$$



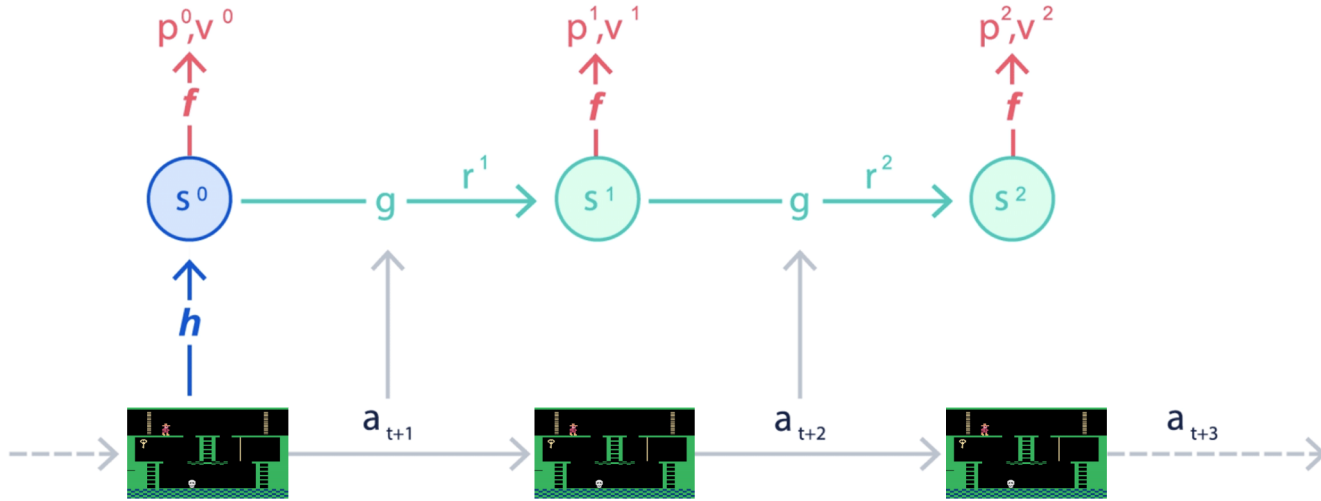
MuZero Training



$$l_t(\theta) = \sum_{k=0}^K l^p(\pi_{t+k}, p_t^k) + \sum_{k=0}^K l^v(z_{t+k}, v_t^k) + \sum_{k=1}^K l^r(u_{t+k}, r_t^k) + c\|\theta\|^2$$



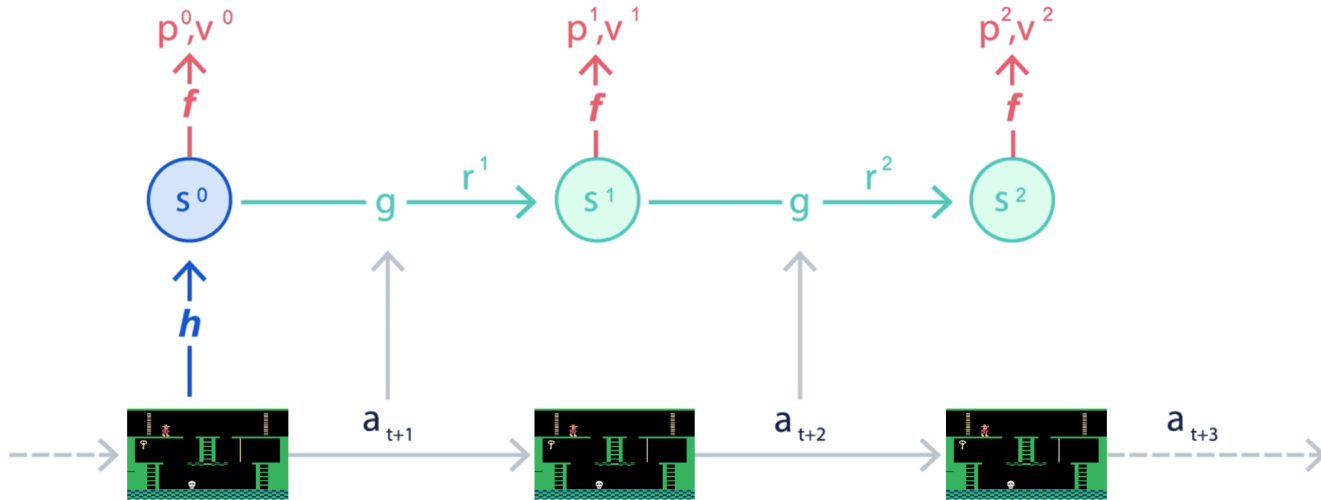
MuZero Training



$$l_t(\theta) = \sum_{k=0}^K l^p(\pi_{t+k}, p_t^k) + \sum_{k=0}^K l^v(z_{t+k}, v_t^k) + \sum_{k=1}^K l^r(u_{t+k}, r_t^k) + c\|\theta\|^2$$



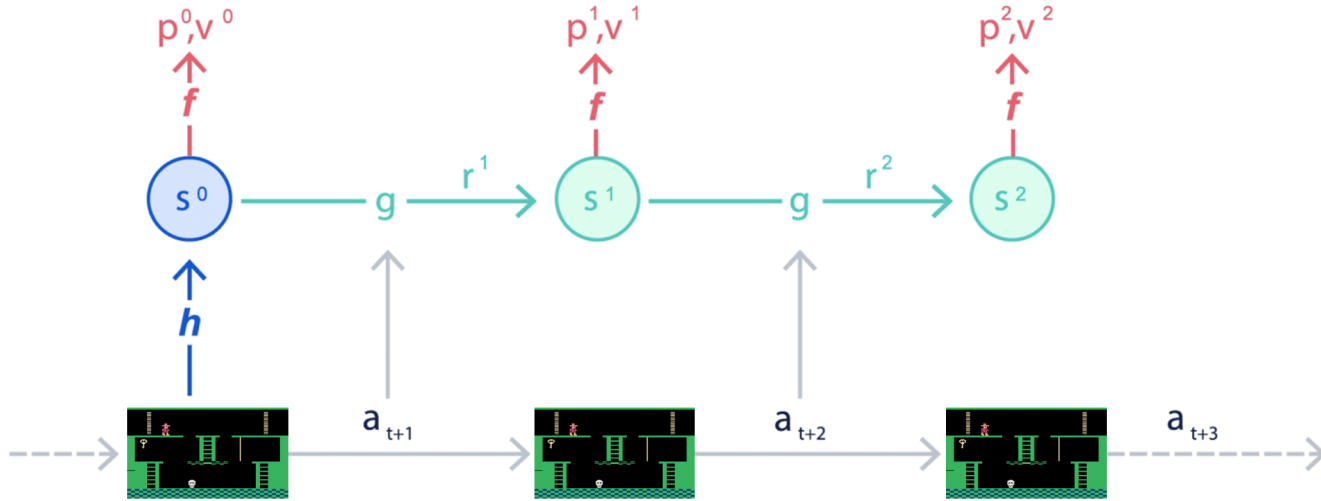
MuZero Training



$$l_t(\theta) = \sum_{k=0}^K l^p(\pi_{t+k}, p_t^k) + \sum_{k=0}^K l^v(z_{t+k}, v_t^k) + \sum_{k=1}^K l^r(u_{t+k}, r_t^k) + c\|\theta\|^2$$



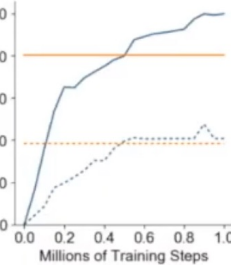
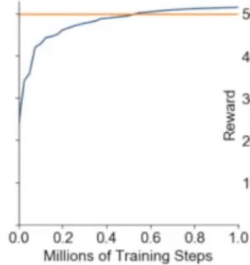
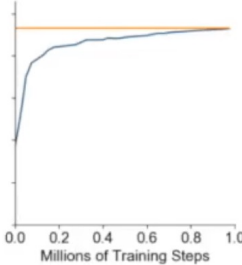
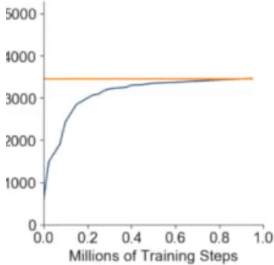
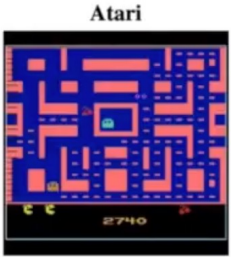
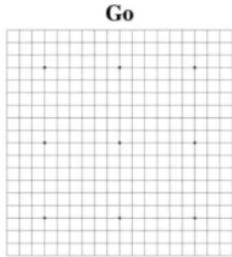
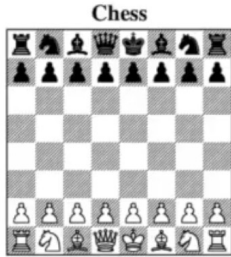
MuZero Training



$$l_t(\theta) = \sum_{k=0}^K l^p(\pi_{t+k}, p_t^k) + \sum_{k=0}^K l^v(z_{t+k}, v_t^k) + \sum_{k=1}^K l^r(u_{t+k}, r_t^k) + c\|\theta\|^2$$

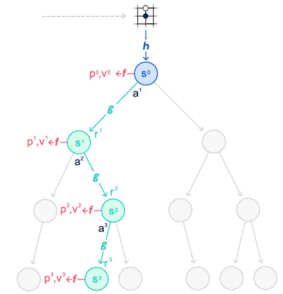
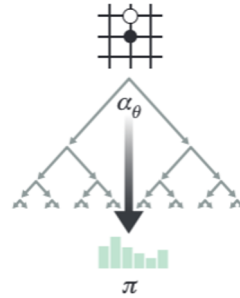
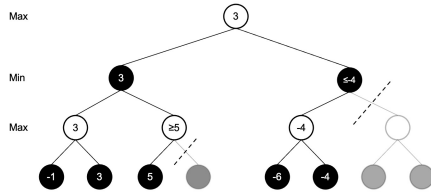
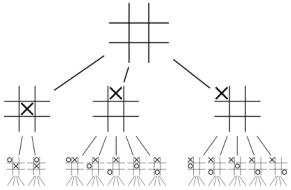
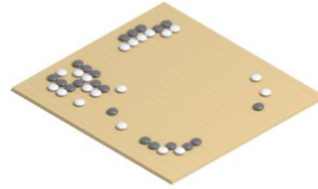


MuZero Results



Summary

0	0	0
	0	X
X	X	



References

- [2020-10-22 MuZero ICAPS talk.pdf](#)
- [AlphaZero: Shedding new light on the grand games of chess, shogi and Go](#)
- [MuZero: Mastering Go, chess, shogi and Atari without rules](#)
- https://en.wikipedia.org/wiki/Alpha-beta_pruning
- <https://www.nature.com/articles/nature16961.pdf>
- https://www.nature.com/articles/nature24270.epdf?author_access_token=VJXbVjaSHxFoctQ4p2k4tRgN0jAjWel9jnR3ZoTv0PVW4gB86EEpGqTRDtplz-2rmo8-KG06ggVobU5NSCFeHILHcVFUeMsbvwS-lxjqQGg98faovwixeTUgZAUMnRQ