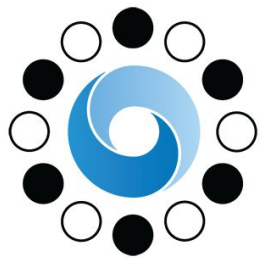


Mastering the Game of Go Without Human Knowledge

Lead: Liam Hinzman

Facilitators: Tahseen Shabab and Susan Cheng

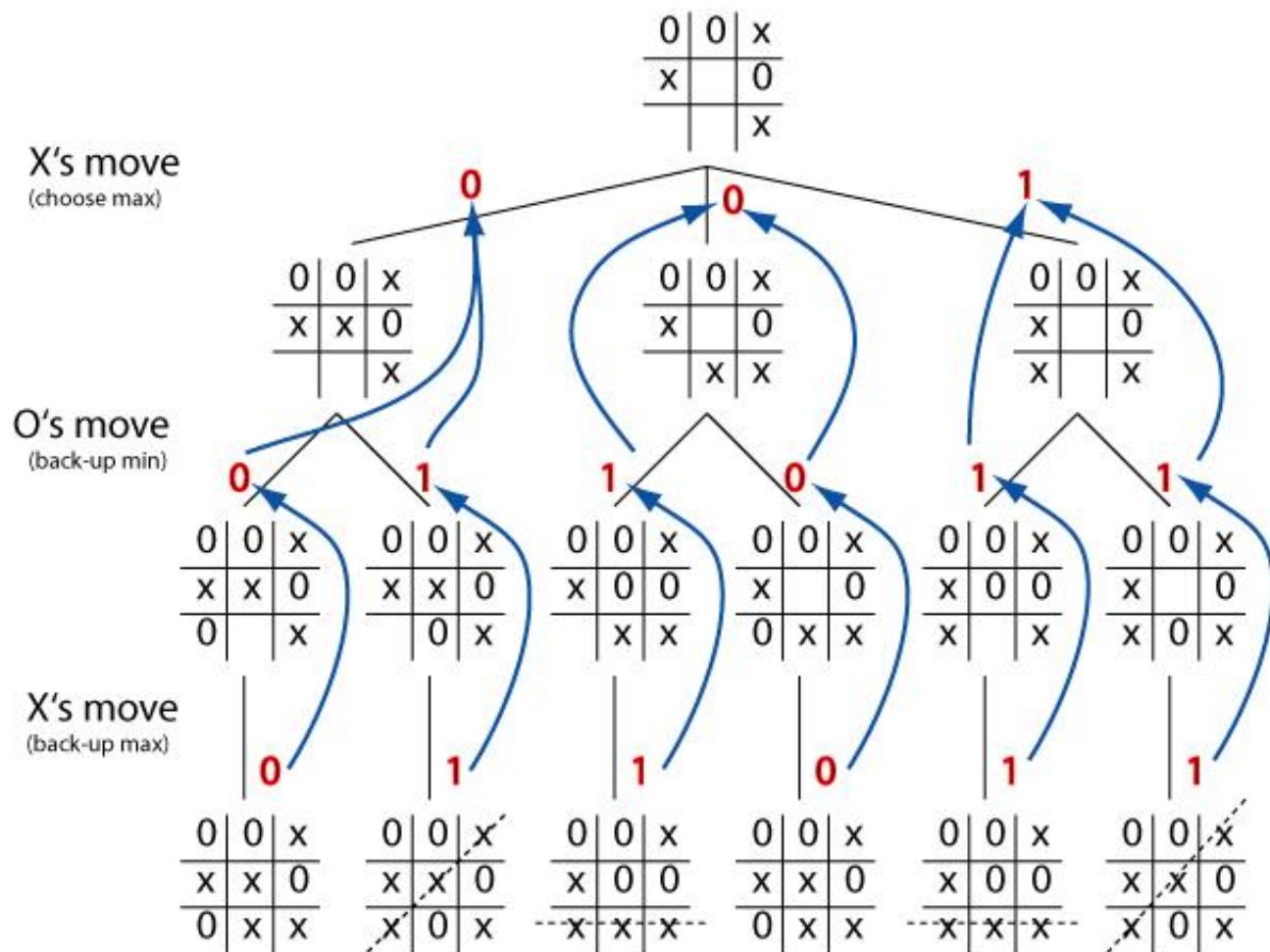


AlphaGo Zero

Overview

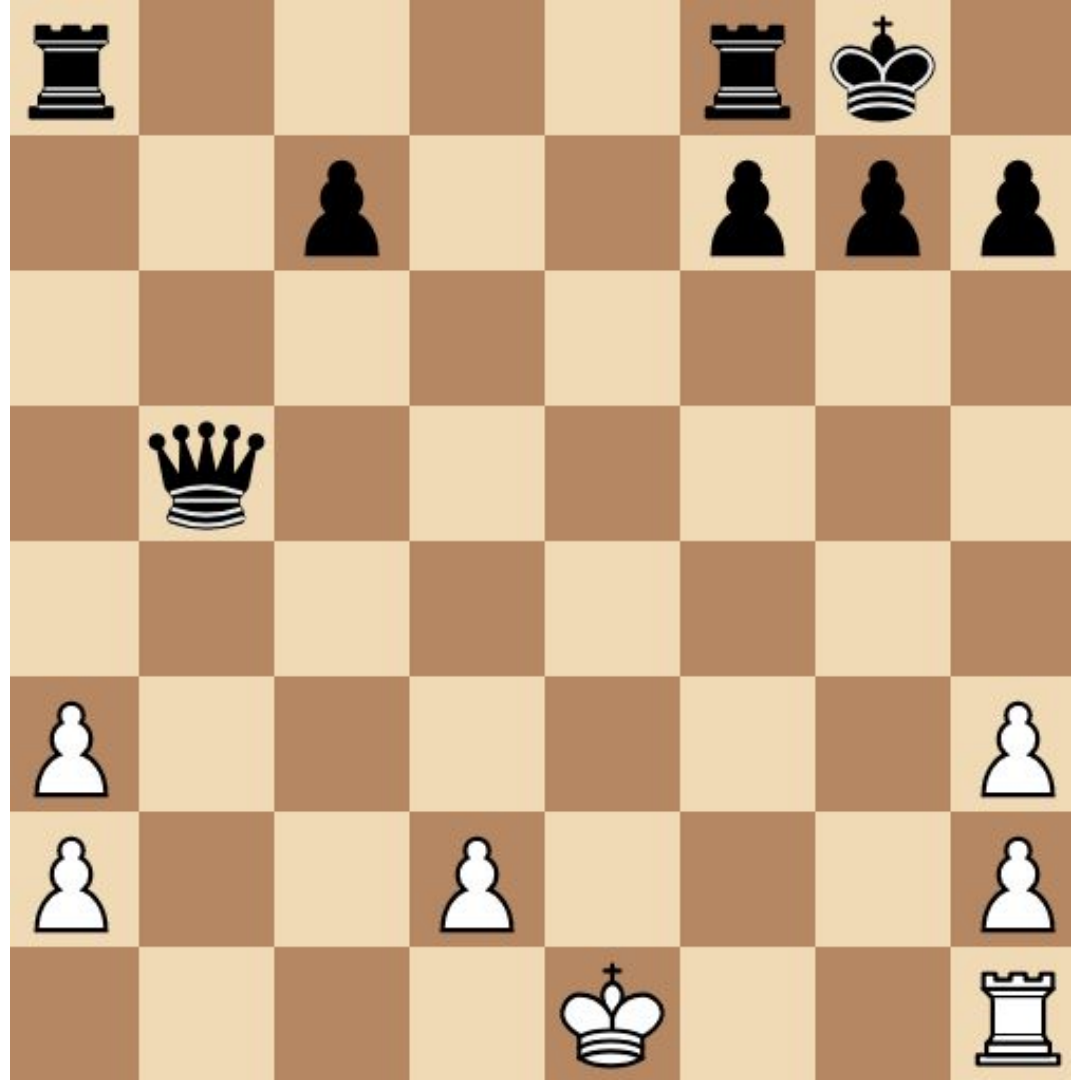
- Brief History of AI in Games
- What is Go and Why Should You Care?
- How AlphaGo Zero Works
- Results
- Discussion

Minimax



Heuristics

Reduces search depth

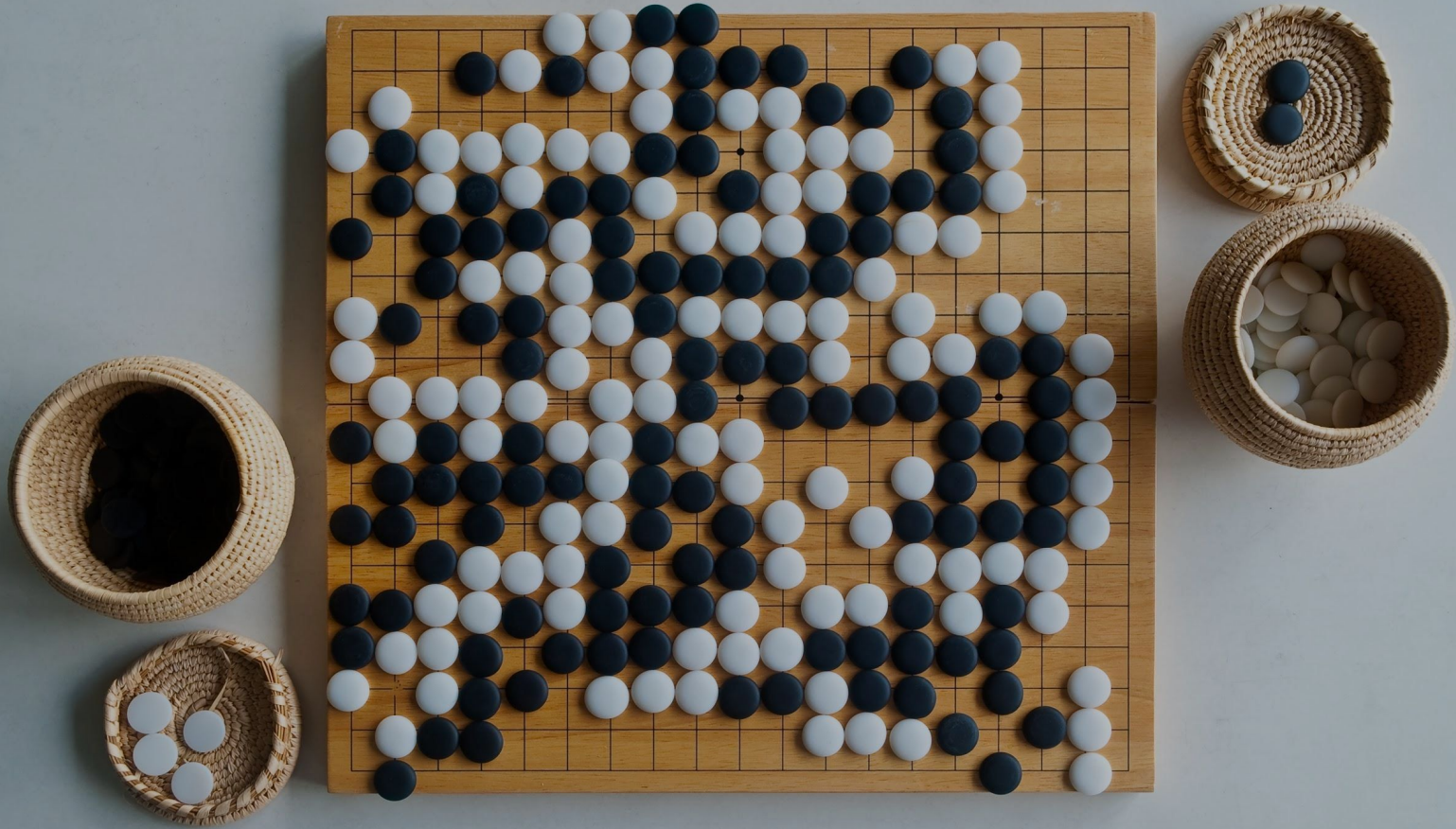


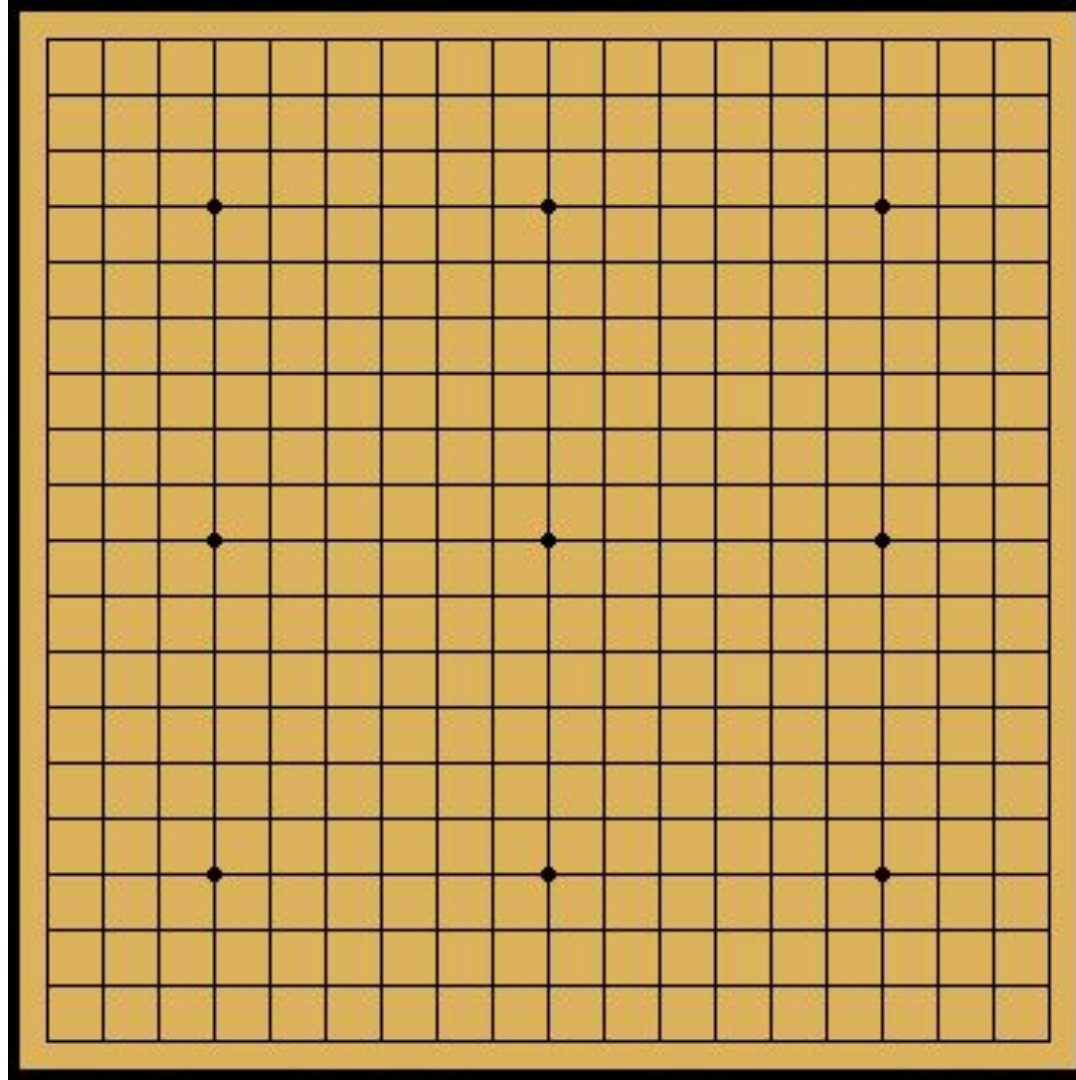
Deep Blue

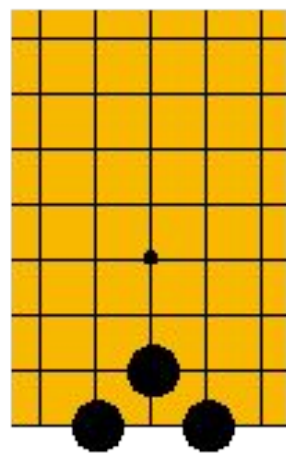
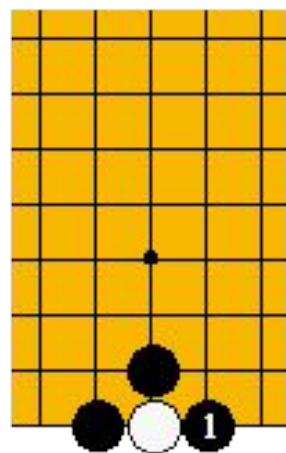
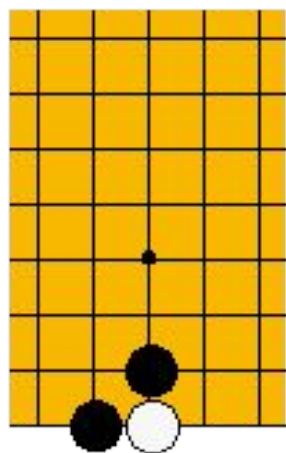
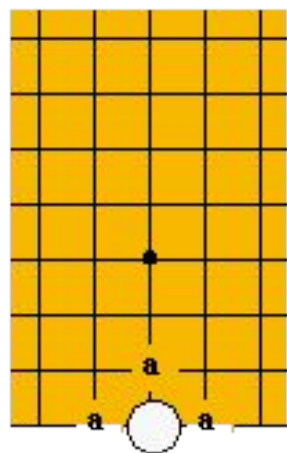
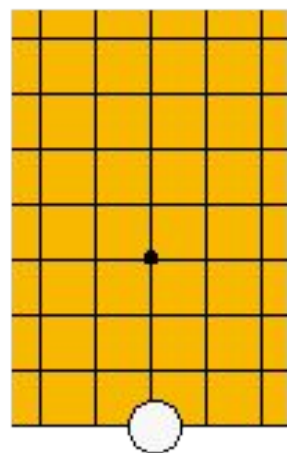
- 126 million positions per second
- Hand-designed Heuristics



The Game of Go

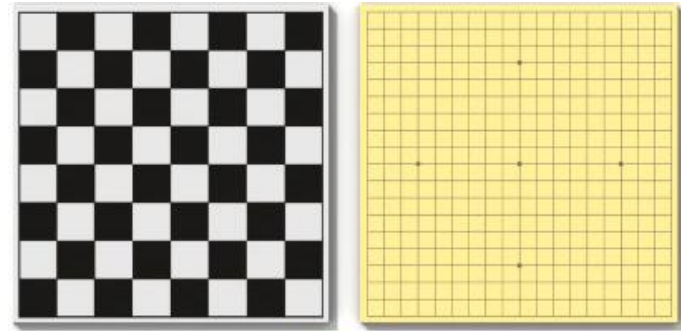






Go is Incredibly Complex

Go is Hard for Computers



GRID SIZE

8 x 8

19 x 19

AVERAGE NUMBER OF MOVE CHOICES PER TURN

35

200-300

LENGTH OF TYPICAL GAME

60 moves

200 moves

NUMBER OF POSSIBLE GAME POSITIONS

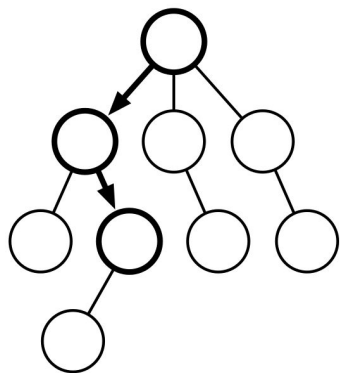
10^{44}

10^{170}

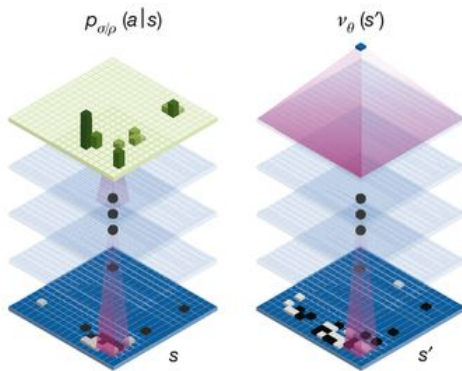
EXPLOSION OF CHOICES (starting from average game position)

35	Move 1	200
1225	Move 2	40 000
42 875	Move 3	8 000 000
1 500 625	Move 4	1 600 000 000

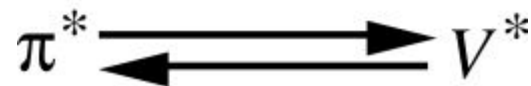
How AlphaGo Zero Works



Monte-Carlo Tree Search



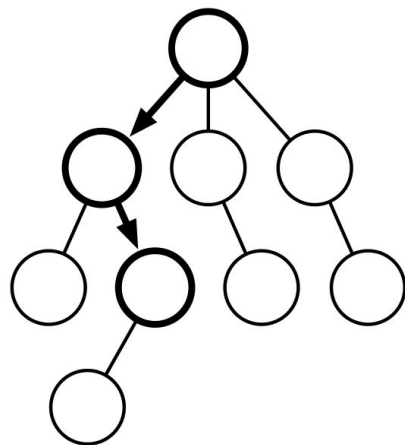
Residual Network



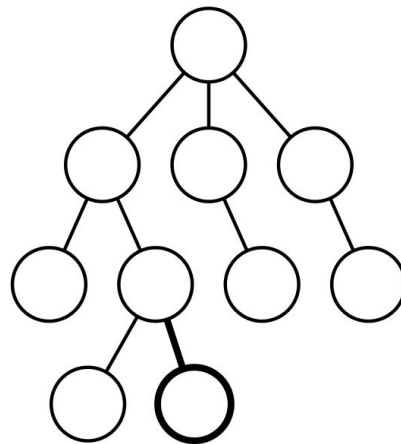
Policy Iteration

Monte-Carlo Tree Search (MCTS)

Selection



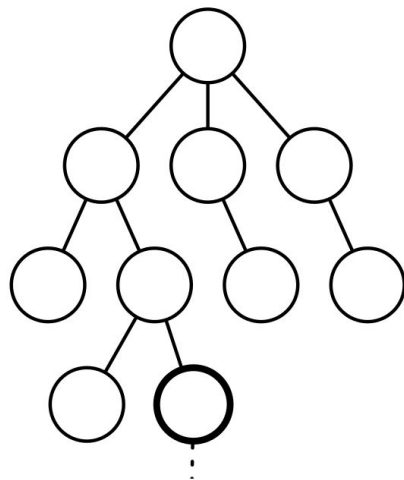
Expansion



Tree Policy

Monte-Carlo Tree Search (MCTS)

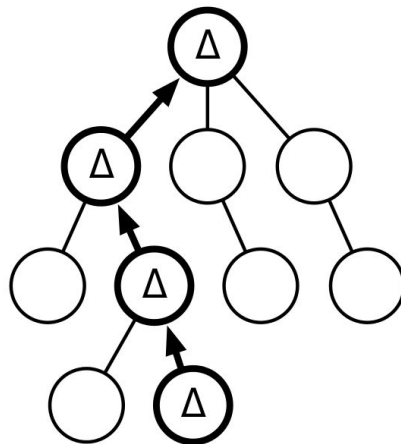
Sampling



Default Policy



Backpropagation



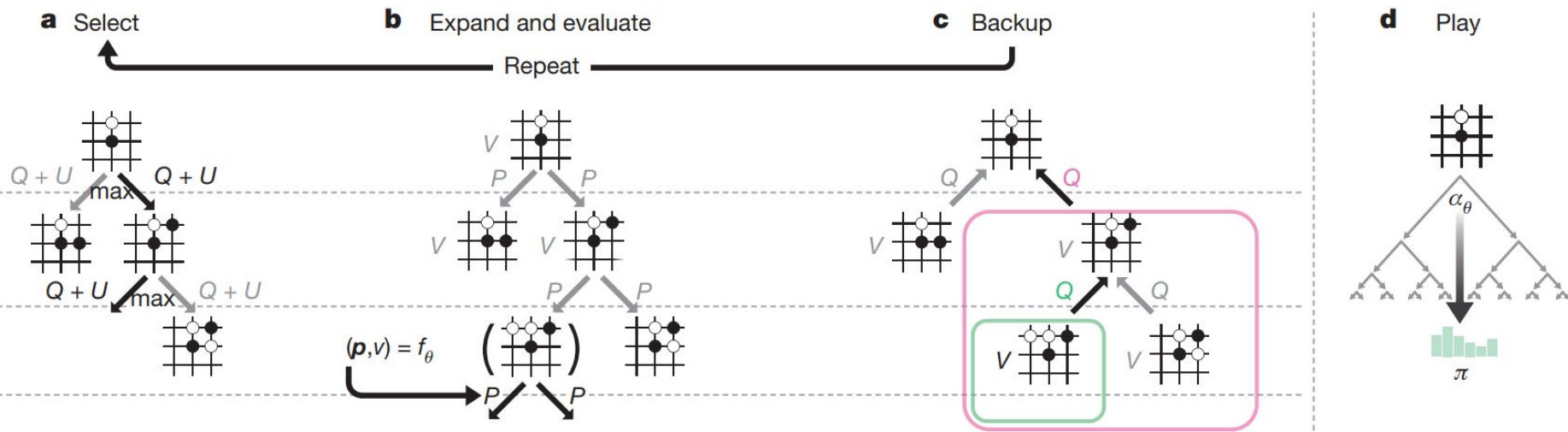
MCTS: Advantages

- Aheuristic
- Online-search
- Works well on large trees

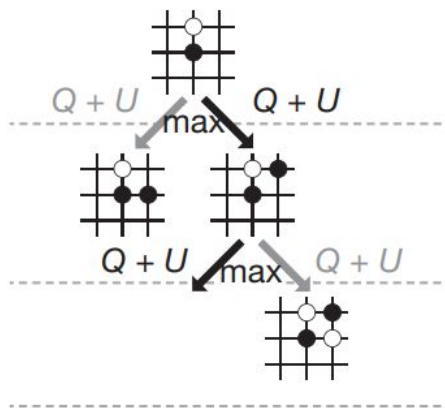
MCTS: Disadvantages

- Many simulation are required
- No generalization between similar states
- Performance is dependent on “rollout” policy

MCTS in AlphaGo Zero

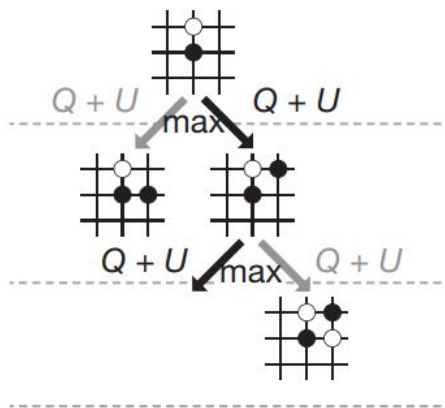


Upper Confidence Bound for Trees (UCT)



$$Q(s, a) + c_{puct} \cdot P(s, a) \cdot \frac{\sqrt{\sum_b N(s, b)}}{1 + N(s, a)}$$

Upper Confidence Bound for Trees (UCT)



$$\underbrace{Q(s, a)}_{\text{Exploitation}} + \underbrace{c_{puct} \cdot P(s, a) \cdot \frac{\sqrt{\sum_b N(s, b)}}{1 + N(s, a)}}_{\text{Exploration}}$$

Exploitation

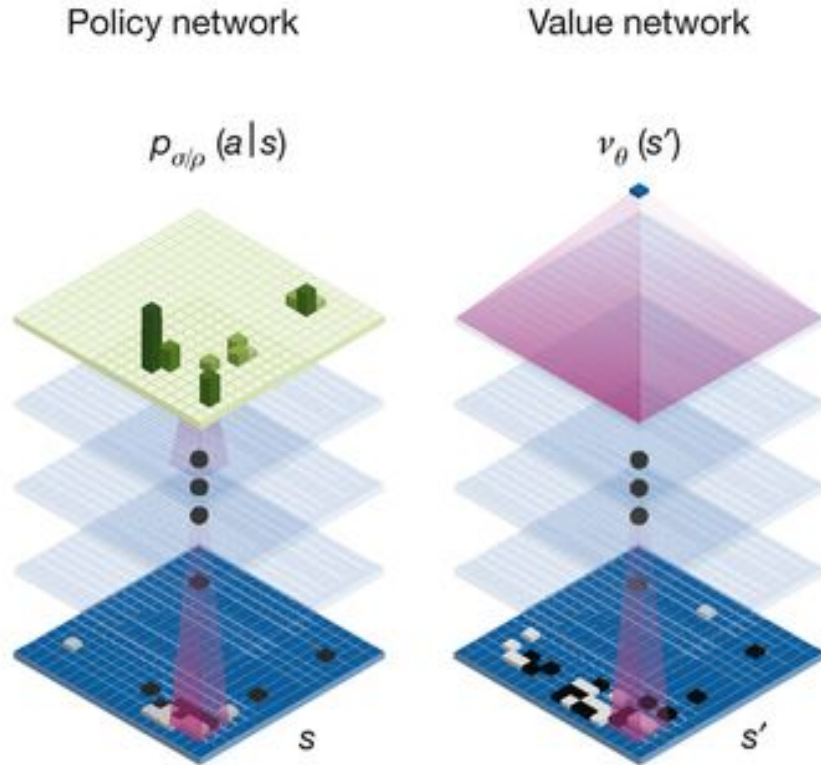
Exploration

Upper Confidence Bound for Trees (UCT)

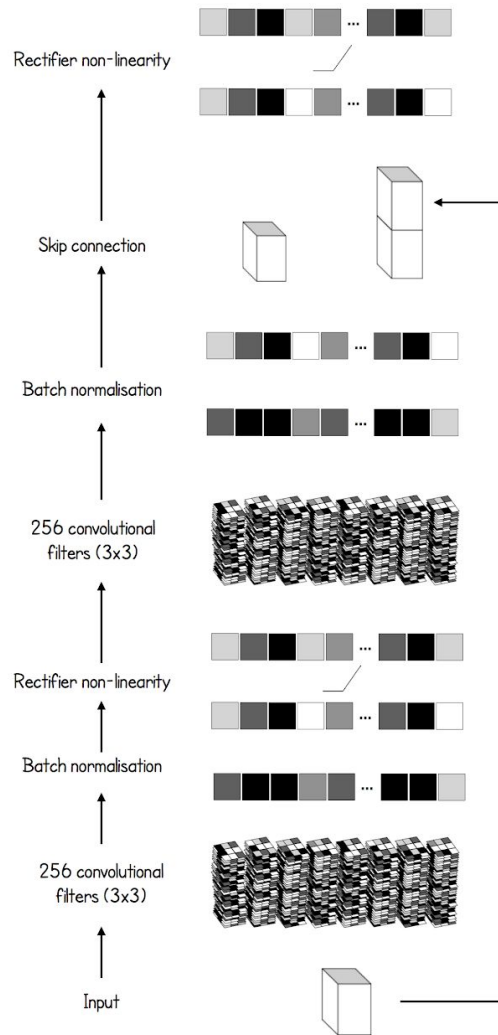
s	State
a	Action
Q(s, a)	Expected Reward
P(s, a)	Policy
N(s, a)	# of state visits
c _{puct}	Hyperparameter

$$\underbrace{Q(s, a)}_{\text{Exploitation}} + \underbrace{c_{puct} \cdot P(s, a) \cdot \frac{\sqrt{\sum_b N(s, b)}}{1 + N(s, a)}}_{\text{Exploration}}$$

AlphaGo Zero's Network Architecture

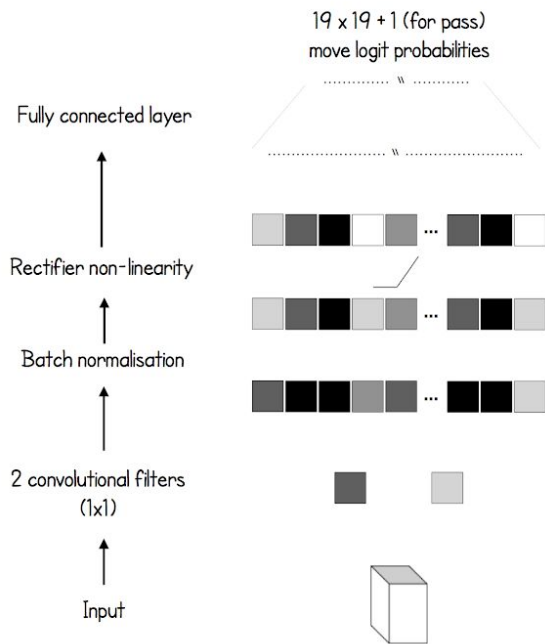


Residual Layer

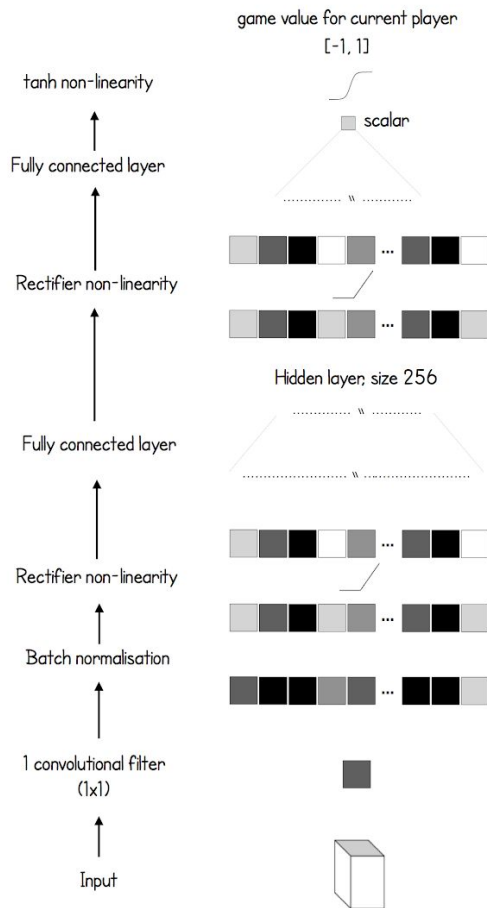


Dual Heads

The policy head



The value head



Training

Self-play Worker

 π 

Training Worker

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

Evaluator

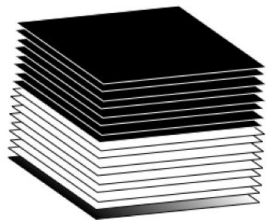
$$\pi' > \pi$$

How AlphaGo Zero Chooses a Move

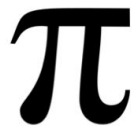
1600
Simulations

$$\pi \sim N^{1/\tau}$$

Self-Play Workers



The game state



The search probabilities



The winner

Training Worker

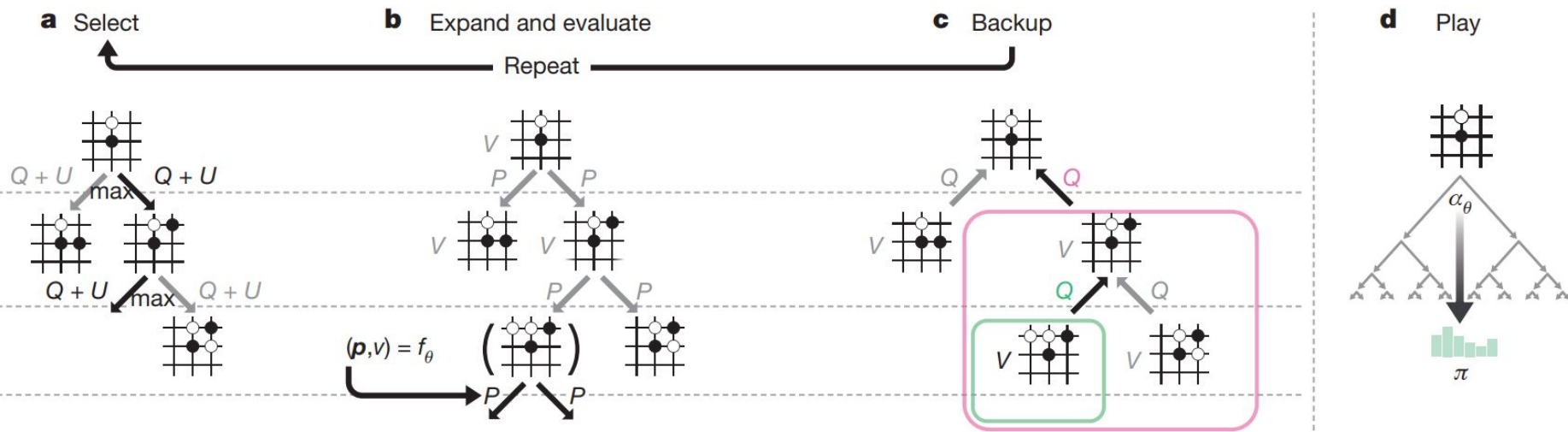
$$l = (z - v)^2 - \boldsymbol{\pi}^T \log \boldsymbol{p} + c \|\boldsymbol{\theta}\|^2$$

Evaluator

400 Games

55% Win Rate

MCTS in AlphaGo Zero



5 Minute Break

AlphaGo Zero

Entirely self-play

Input is game board

Single network

No rollouts

vs

AlphaGo

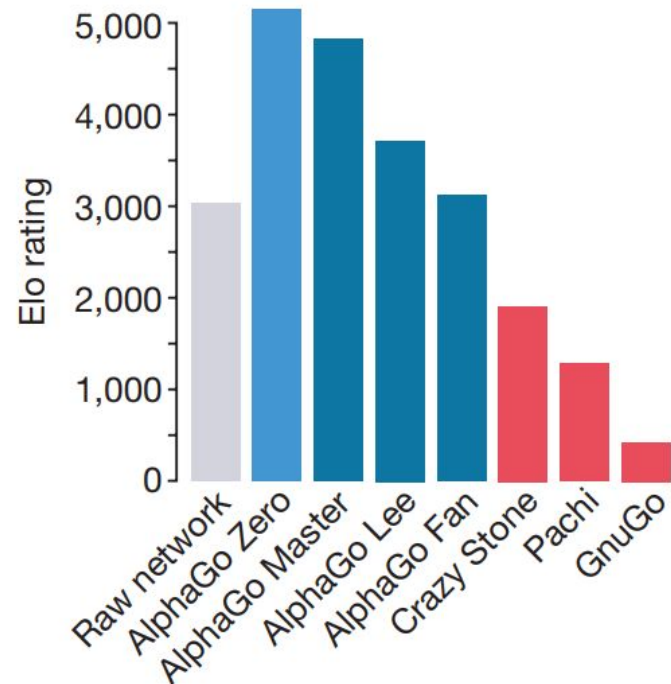
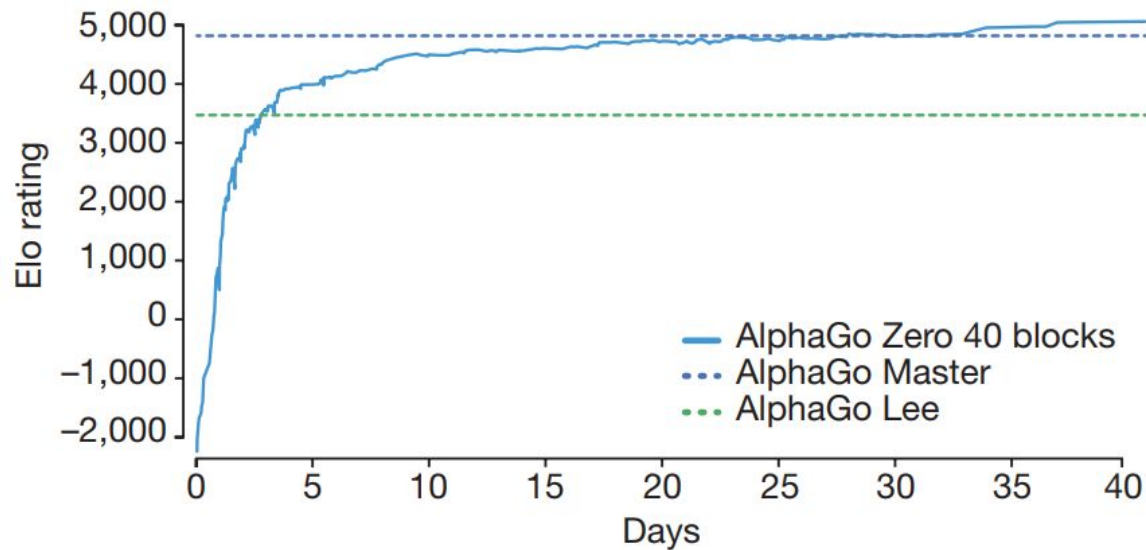
Supervised learning + self-play

Input is hand-crafted features

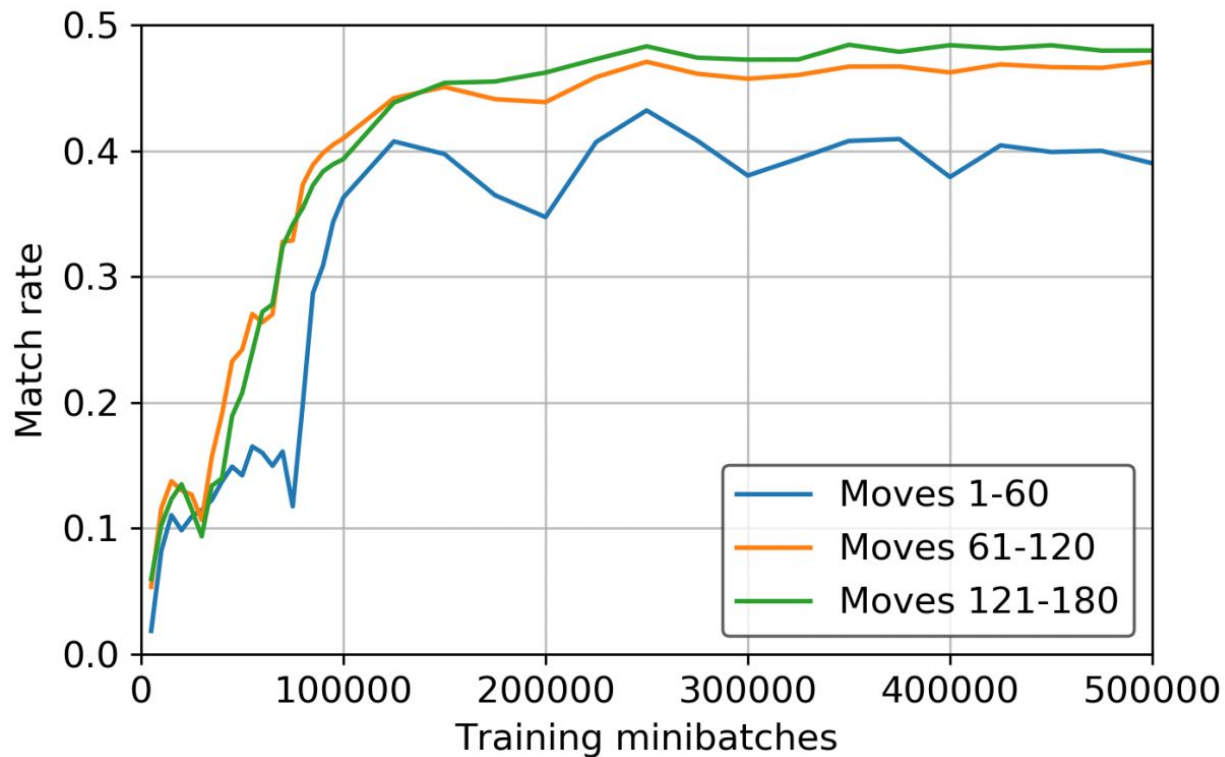
Two networks

Rollouts were used

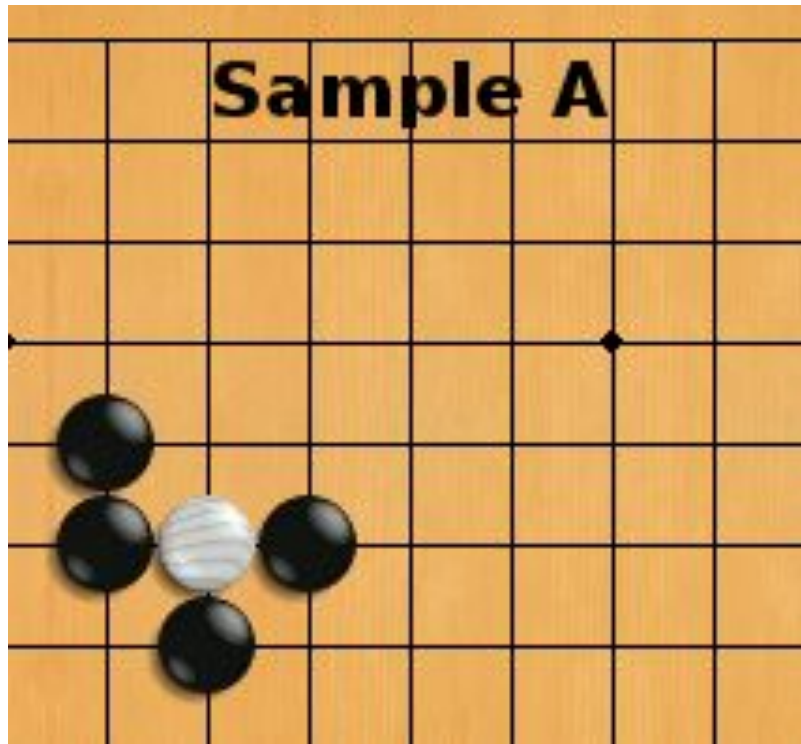
Results



Learning Stages



Ladders

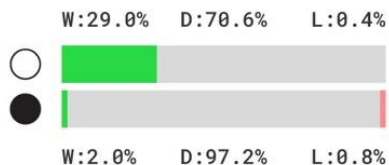


AlphaZero

Chess



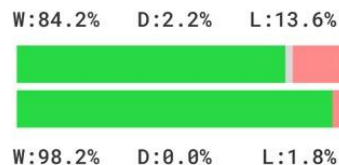
AlphaZero vs. Stockfish



Shogi



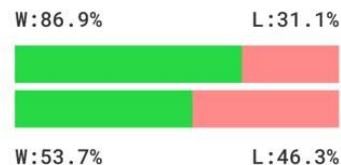
AlphaZero vs. Elmo



Go



AlphaZero vs. AGO



AZ wins ■ AZ draws ■ AZ loses ■ AZ white ○ AZ black ●

AlphaGo Zero's Gift



Discussion



Discussion

How can the AlphaGo Zero algorithm be extended to different games?

How can the sample efficiency of AlphaGo Zero be improved?

A very stable training environment is need for the algorithm.

Can this be alleviated to let AlphaZero applied to real-world problems?

Resources

Mastering the Game of Go without Human Knowledge

David Silver 2017 NIPS Talk

ELF OpenGo: An Analysis and Open Reimplementation of AlphaZero

David Silver's PhD Thesis: Reinforcement Learning and
Simulation-Based Search in Computer Go

A Brief History of Game AI Up To AlphaGo - Andrey Kurenkov

AlphaGo Zero Demystified - Dylan Djian