

Six - AlphaGo

A landmark achievement in artificial intelligence.

IBM Deep Blue



Figure 1: First defeat of a world chess champion by a machine in 1997

Deep Blue was handcrafted by programmers & chess grandmasters

AlphaGo *learnt* from human moves & self play

AlphaGo evaluated fewer positions

- **width** - policy network select states more intelligently
- **depth** - value function evaluate states more precisely

Why Go?

Long held as the most challenging classic game for artificial intelligence

- massive search space
- more legal positions than atoms in universe
- difficult to evaluate positions & moves
- sparse & delayed reward

Difficult to evaluate positions

- chess you can evaluate positions by summing the value of all the pieces
- go - it's just stones on the board, equal numbers each side

Components of the AlphaGo agent

Three policy networks $\pi(s)$

- fast rollout policy network – linear function
- supervised learning policy – 13 layer convolutional NN
- reinforcement learning policy – 13 layer convolutional NN

One value function $V(s)$ - convolutional neural network

Combined together using Monte Carlo tree search

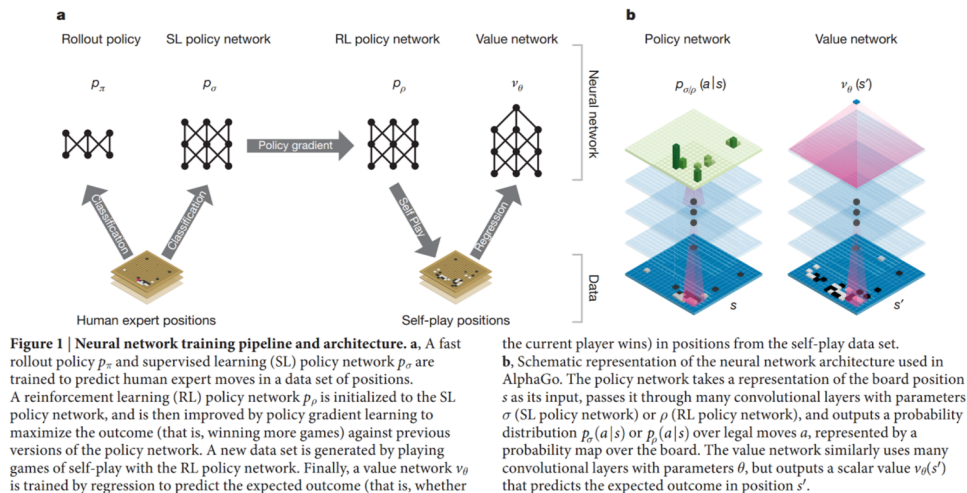


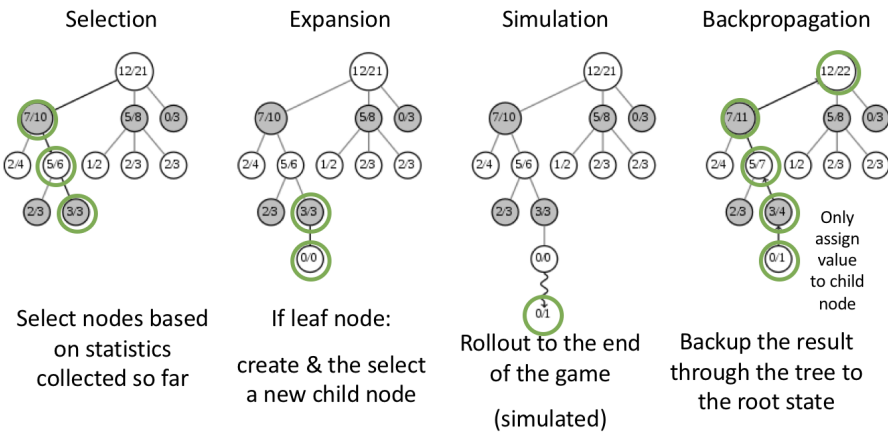
Figure 2: Learning of AlphaGo

Monte Carlo Tree Search

Value & policy networks combined using MCTS

Basic idea = analyse most promising next moves

Planning algorithm - simulated (not actual experience) - roll out to end of game (a simulated Monte Carlo return)



https://en.wikipedia.org/wiki/Monte_Carlo_tree_search

Figure 3: fig

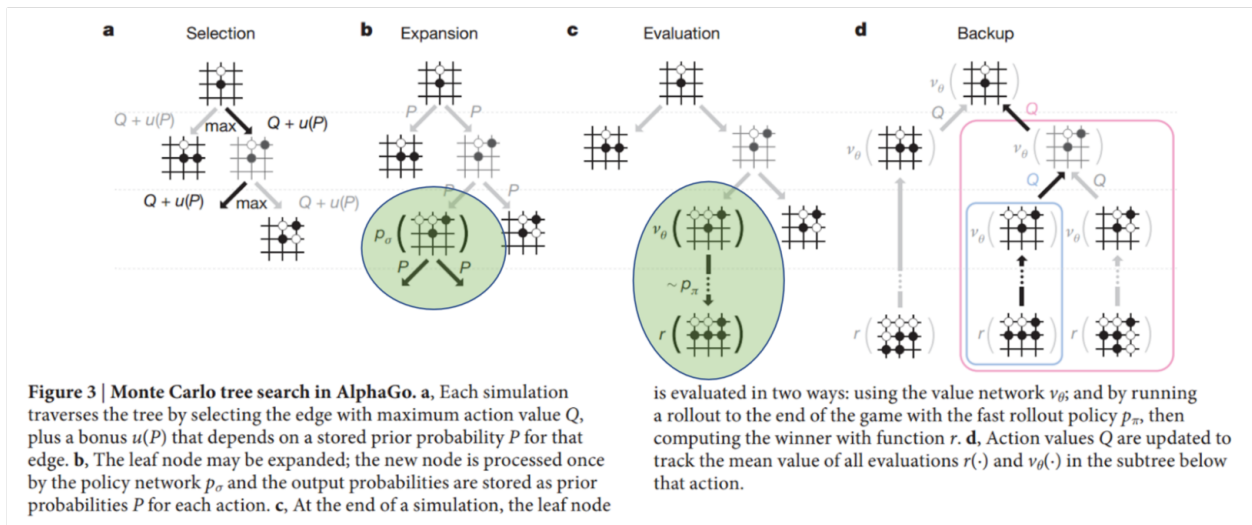
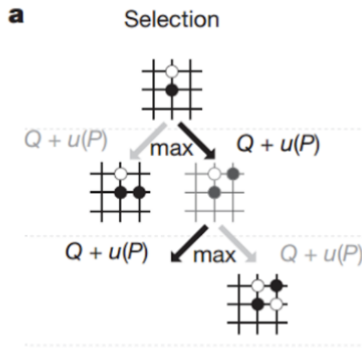


Figure 4: MCTS in AlphaGo

AlphaGo, in context – Andrej Karpathy

Convenient properties of Go

- fully deterministic
- fully observed



Bonus $u(s, a)$ penalizes
more visits
to encourages
exploration

Each edge (s, a) keeps track of statistics

action value $Q(s, a)$

visit count $N(s, a)$

prior probability $P(s, a)$ SL policy
network

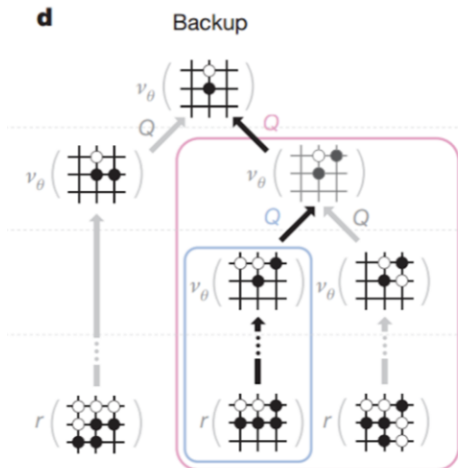
Action selected according too

$$a = \underset{a}{\operatorname{argmax}} [Q(s, a) - u(s, a)]$$

where

$$u(s, a) \propto P(s, a) / [1 - N(s, a)]$$

Figure 5: fig



After we finish our rollout – we calculate
a state value for our leaf node s_L

$$V(s_L) = (1 - \lambda)v_\theta(s) + \lambda z$$

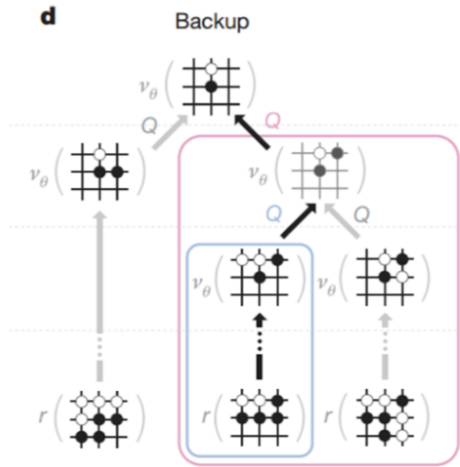
λ mixing parameter

v_θ value network estimate

z simulated result of rollout

We are combining the value network
with the MCTS rollout

Figure 6: fig



Then use our combined estimate $V(s_L)$ to update

action value $Q(s, a)$

visit count $N(s, a)$

$$Q(s, a) = \frac{1}{N(s, a)} \sum_{i=1}^n V(s_L^i)$$

We are averaging across all visits in the simulation

After all simulations finished - **select the most visited action from the root state**

Figure 7: fig

- discrete action space
- access to perfect simulator
- relatively short episodes
- evaluation is clear
- huge datasets of human play

DeepMind take advantage of properties of Go that will not be available in real world applications of reinforcement learning.

AlphaGo Zero

Key ideas

- simpler
- search
- adversarial
- machine knowledge only

Training time & performance

- AG Lee trained over several months
- AG Zero beat AG Lee 100-0 after 72 hours of training

Computational efficiency

- AG Lee = distributed w/ 48 TPU
- AG Zero = single machine w/ 4 TPU

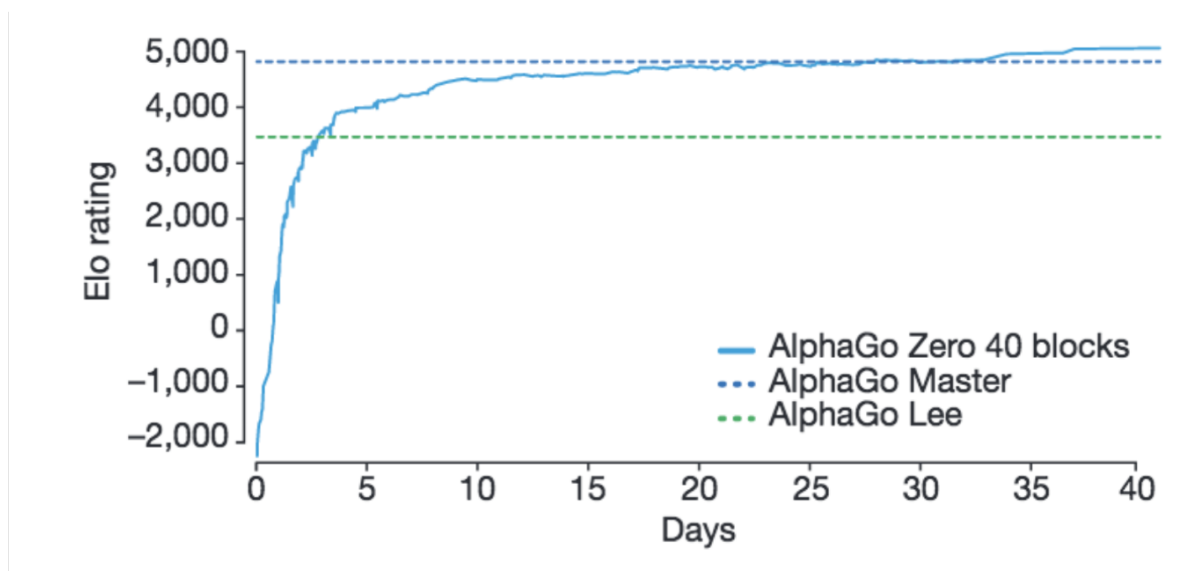


Figure 6 | Performance of AlphaGo Zero. a, Learning curve for AlphaGo Zero using a larger 40-block residual network over 40 days. The plot shows the performance of each player α_{θ_i} from each iteration i of our reinforcement learning algorithm. Elo ratings were computed from evaluation games between different players, using 0.4 s per search (see Methods)

Figure 8: fig

AlphaGo Zero innovations

Learns using only self play

- no learning from human expert games
- no feature engineering
- learn purely from board positions

Single neural network - combine the policy & value networks

RL surpasses SL after around 1 day

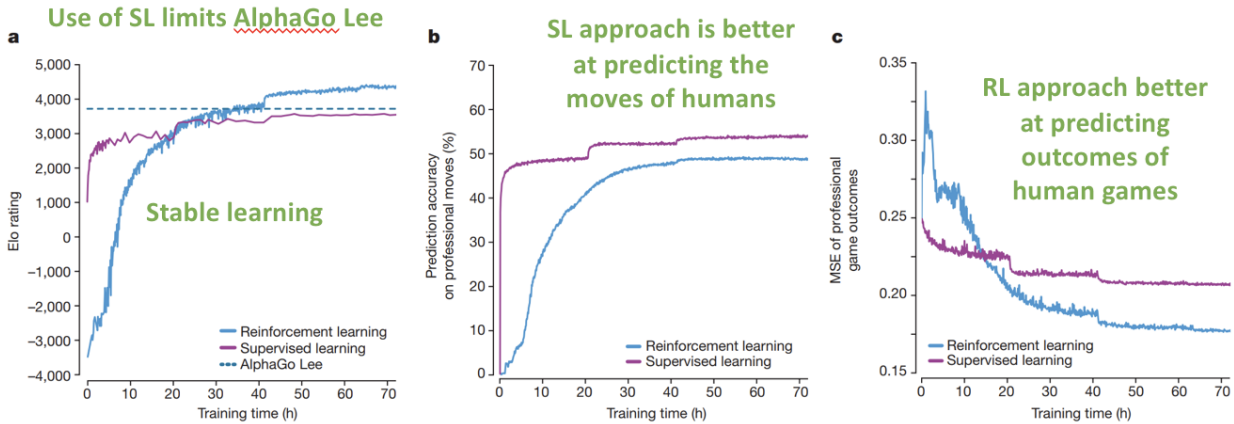


Figure 3 | Empirical evaluation of AlphaGo Zero. **a**, Performance of self-play reinforcement learning. The plot shows the performance of each MCTS player α_{θ_i} from each iteration i of reinforcement learning in AlphaGo Zero. Elo ratings were computed from evaluation games between different players, using 0.4 s of thinking time per move (see Methods). For comparison, a similar player trained by supervised learning from human data, using the KGS dataset, is also shown. **b**, Prediction accuracy on human professional moves. The plot shows the accuracy of the neural network f_{θ_i} at each iteration of self-play i , in predicting human professional moves from the GoKifu dataset. The accuracy measures the

percentage of positions in which the neural network assigns the highest probability to the human move. The accuracy of a neural network trained by supervised learning is also shown. **c**, Mean-squared error (MSE) of human professional game outcomes. The plot shows the MSE of the neural network f_{θ_i} at each iteration of self-play i , in predicting the outcome of human professional games from the GoKifu dataset. The MSE is between the actual outcome $z \in \{-1, +1\}$ and the neural network value v , scaled by a factor of $\frac{1}{4}$ to the range of 0–1. The MSE of a neural network trained by supervised learning is also shown.

Figure 9: fig

MCTS only during acting (not during learning)

Use of residual networks

Search in AlphaGo Zero

Policy evaluation

Policy is evaluated through self play

This creates high quality training signals - the game result

Policy improvement

MCTS is used during acting to create the improved policy

The improved policy generated during acting becomes the target policy during training

[Keynote David Silver NIPS 2017 Deep Reinforcement Learning Symposium AlphaZero](#)

DeepMind AlphaGo AMA

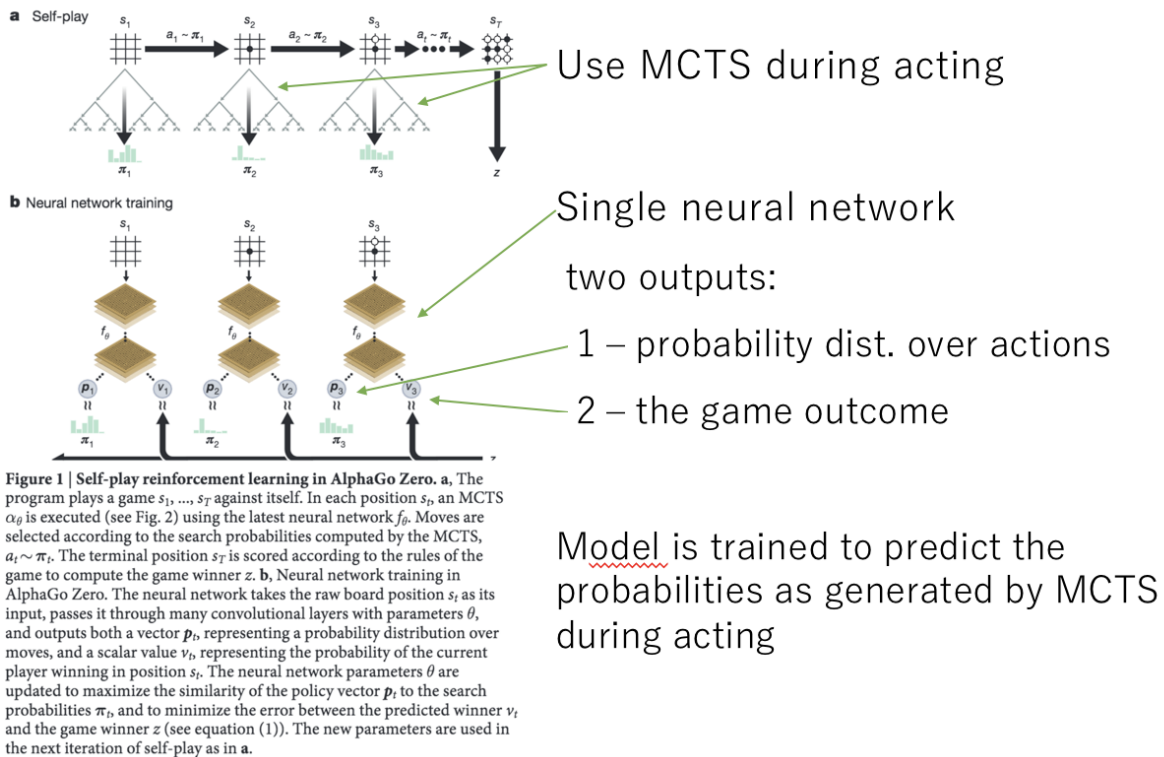


Figure 10: Acting and learning

AMA: We are David Silver and Julian Schrittwieser from DeepMind's AlphaGo team. Ask us anything. (self.MachineLearning)

submitted 3 days ago * (last edited 1 day ago) by David_Silver

DeepMind - announcement

this post was submitted on 17 Oct 2017

245 points (97% upvoted)

shortlink: <https://redd.it/76xjb5>

Figure 11: fig

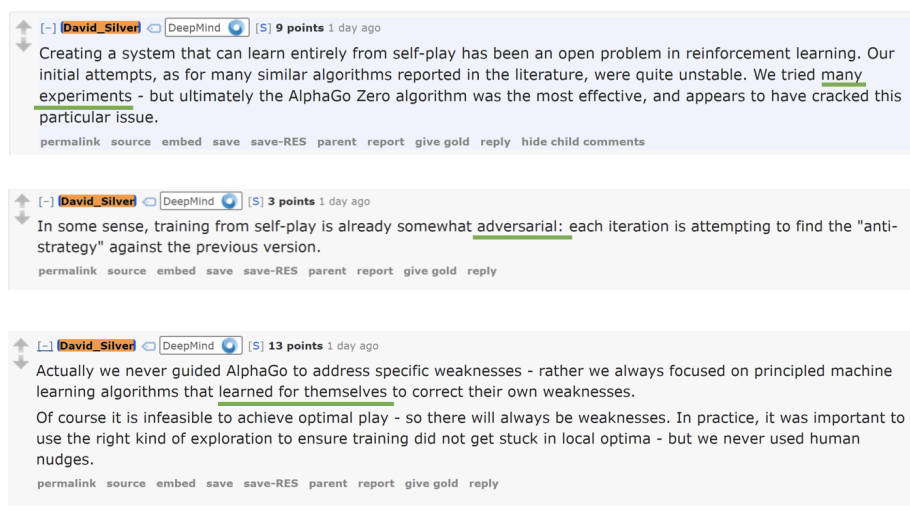


Figure 12: fig