

Research review

Paper: **Mastering the game of Go with deep neural networks and tree search**

(<https://storage.googleapis.com/deepmind-media/alphago/AlphaGoNaturePaper.pdf>)

The paper introduces a novel approach for Go gaming by combining MonteCarlo Tree Search (MCTS) together with Deep Learning. Strongest Go programs prior AlphaGo were based only on MCTS and therefore have been outperformed in 494 out of 495 games by AlphaGo. Moreover, AlphaGo could win games with handicap opponent stones with minimum winrate of 77%.

AlphaGo uses pipeline of policy functions P and value function v . Policy functions are used to predict next move based on current board state. Value function is used as heuristic and gives a score of current board state.

The pipeline is described as following:

- Supervised Deep Learning classifier P_{sigma} based on Convolutional Neural Network with ReLU activations and Softmax probabilities is trained for move likelihood prediction.
- Another shallow (and fast) classifier P_{pi} trained using just softmax over small pattern features. This one is used for rapid action sampling during MonteCarlo rollouts.
- Reinforcement learning (RL) network P_{ro} with same structure and initial weights as P_{sigma} is then trained by playing against previous iterations of itself.
- Train RL-based value network v is then trained based on P_{ro} policy.
- Combine policy and value network with MCTS -- this is AlphaGo itself.

Most notable result is, of course, outstanding winning rate in comparison with other programs and, especially, first time that a Go program has defeated a professional human player.

But in addition, several intermediate results are also very interesting and promising:

- Small improvements of supervised deep learning training of P_{sigma} can significantly improve strength of further Go-playing program.
- Only P_{ro} policy network can defeat many other programs and is ranked at 2 amateur dan.
- State evaluation function (value network v) approaches the accuracy of MonteCarlo rollouts with **15000 less** computations.
- Having that Go game has significantly more possible states because of its breadth and branching factors, AlphaGo evaluated thousands of times fewer positions than Deep Blue did in its famous chess match.

Deep Neural Networks allow to create much more intelligent AI software without handcrafting

of heuristics. The cost of this is computational resources needed for runtime execution and, most important, for training. AlphaGo used an asynchronous multi-thread search that used 48 CPUs and 8 GPUs. Distributed version of AlphaGo exploited multiple machines with total 1202 CPUs and 176 GPUs.