# AlphaGo

# Go in numbers


**3,000**
Years Old


**40M**
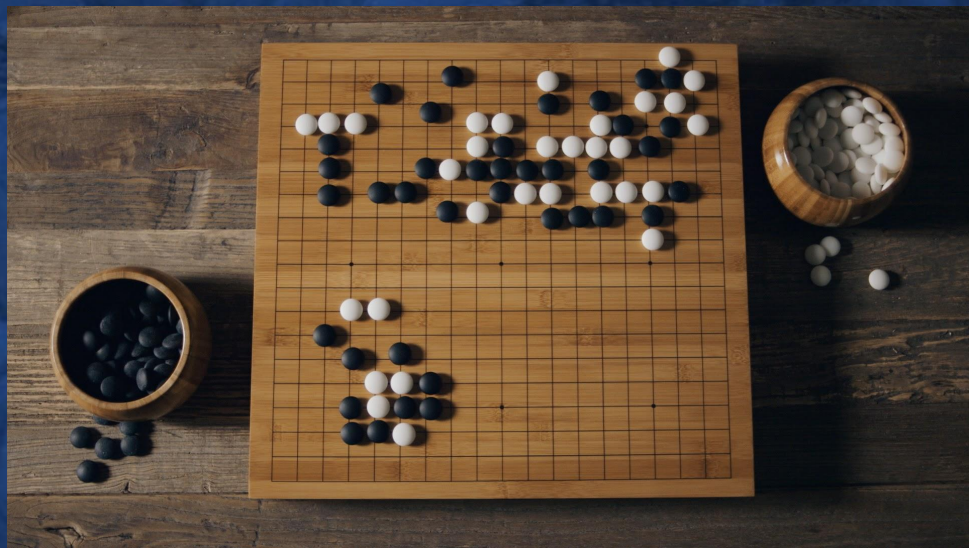Players


**10^170**
Positions

Google DeepMind

# Why is Go hard for computers to play?

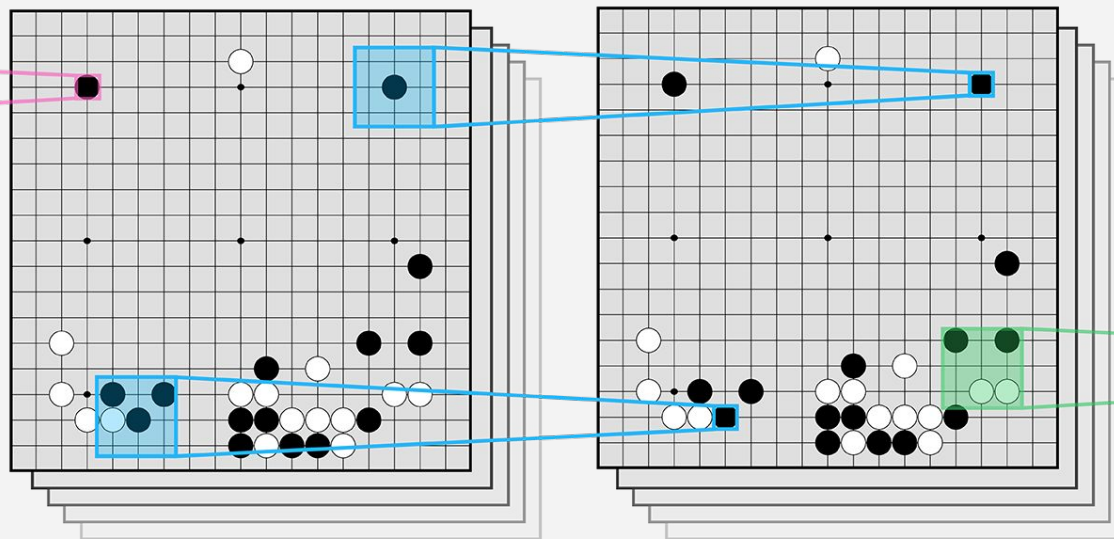Game tree complexity = $b^d$

Brute force search intractable:

1.  Search space is huge

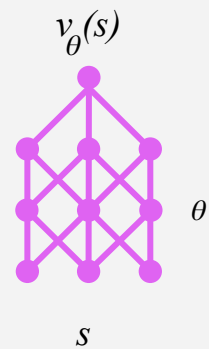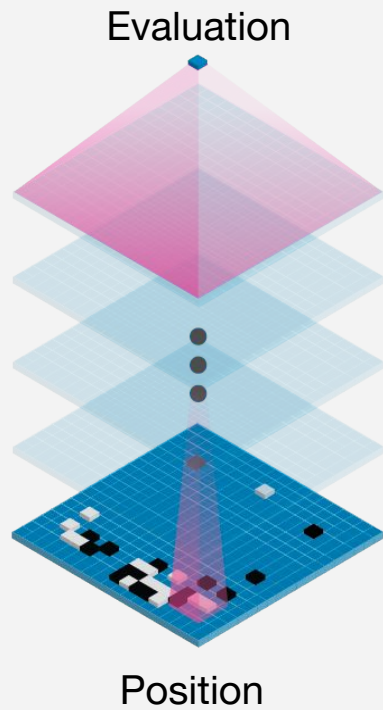2.  "Impossible" for computers to evaluate who is winning
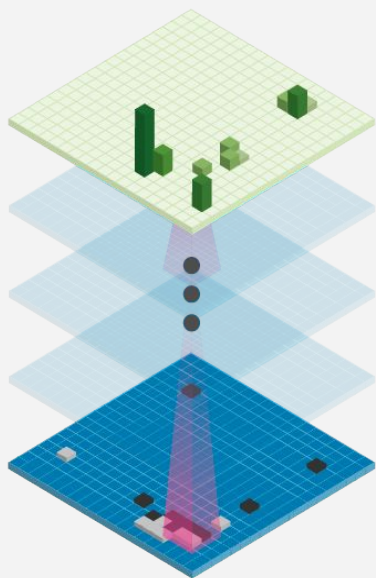
# Convolutional neural network

# Value network



Evaluation

Position

$v_\theta(s)$

$\theta$

$s$

Google DeepMind

# Policy network

Move probabilities



Position

$p_\sigma(a|s)$

$\sigma$

$s$

Google DeepMind
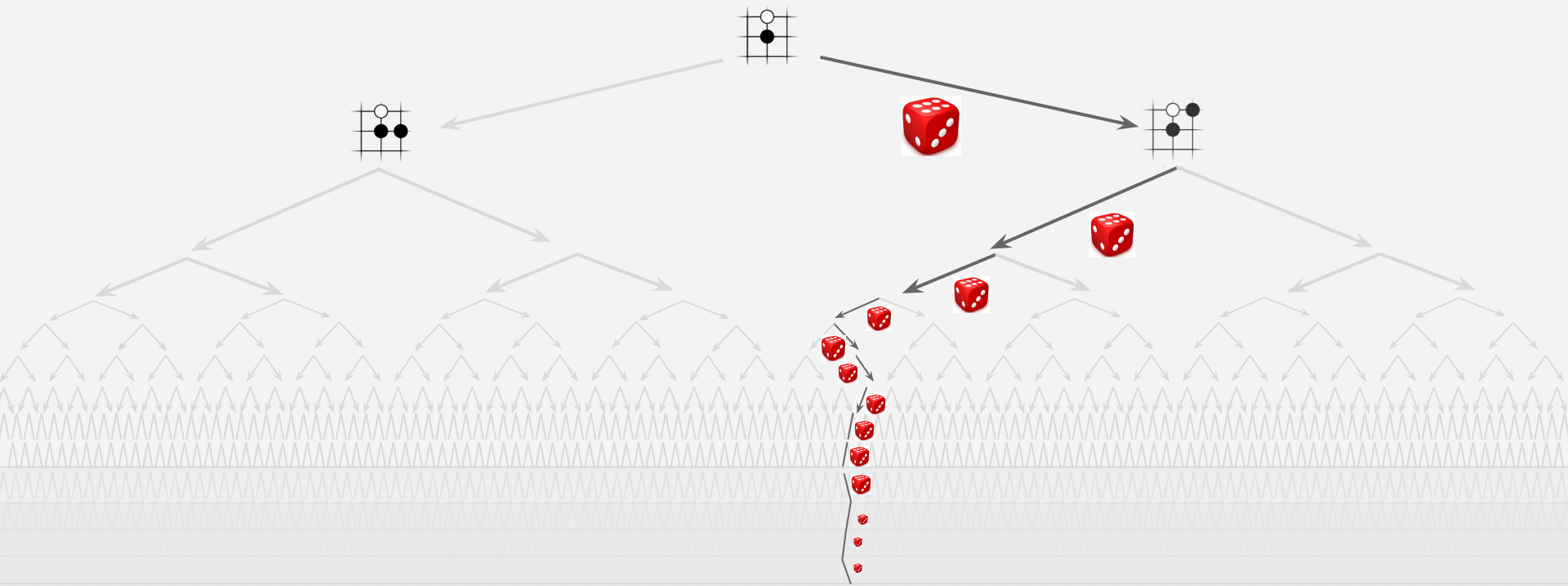
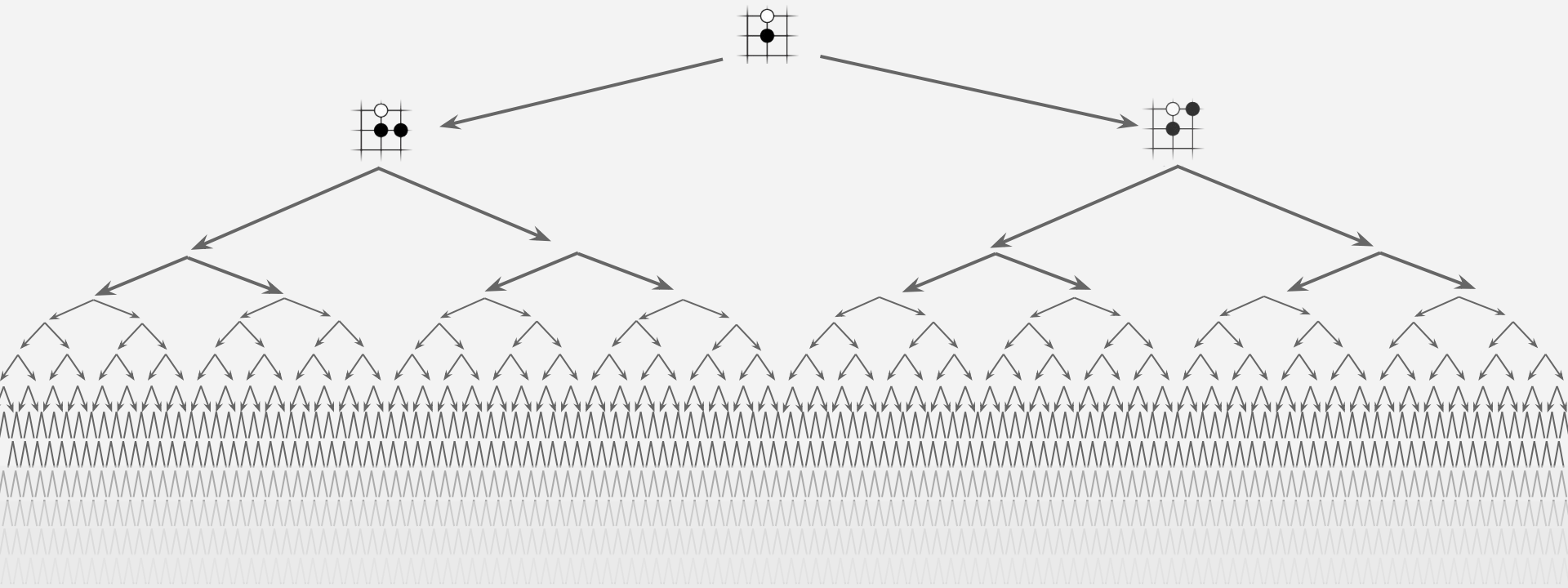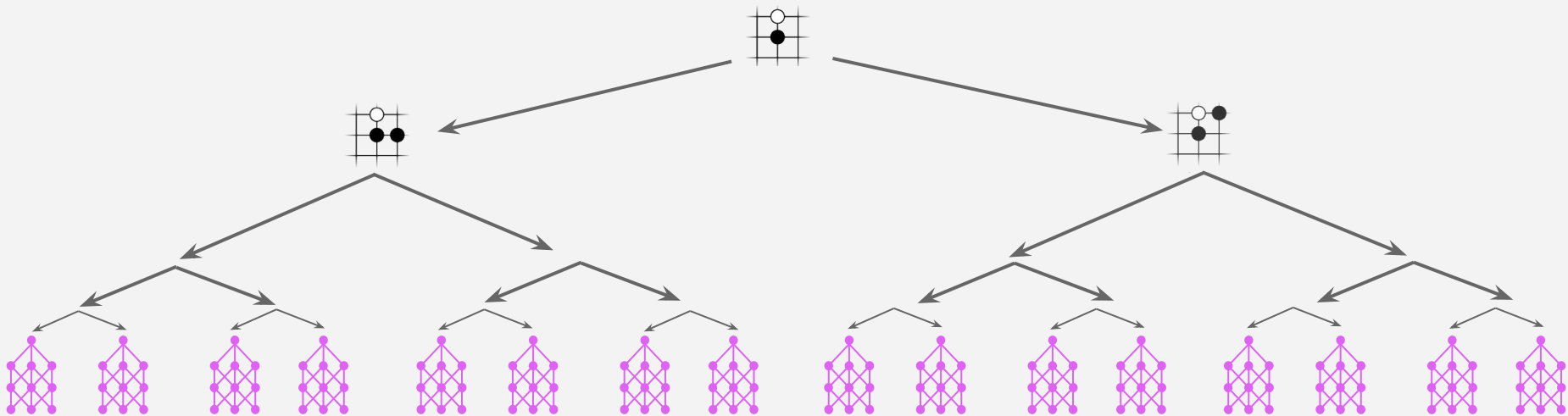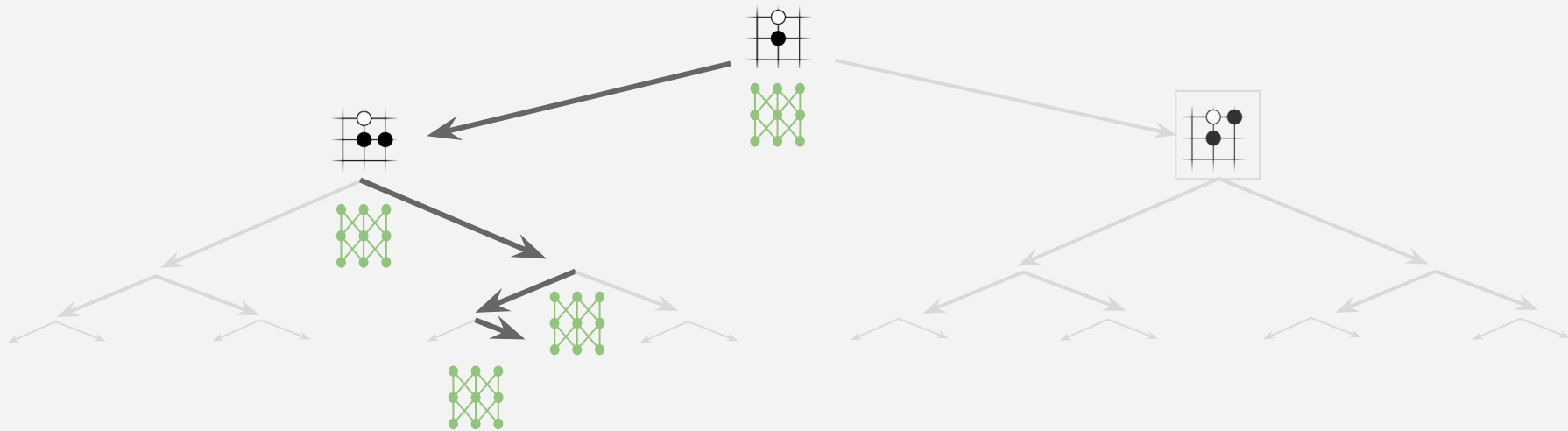# Exhaustive search

# Monte-Carlo rollouts

# Reducing depth with value network

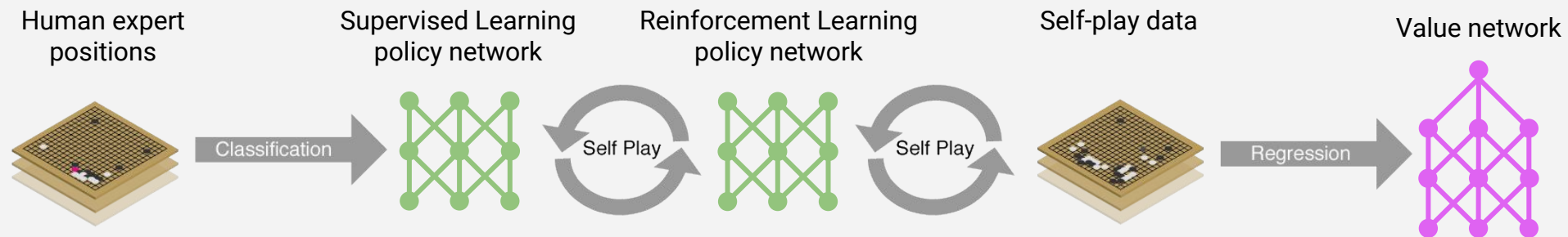# Reducing depth with value network

# Reducing breadth with policy network

# Neural network training pipeline



Human expert positions → Classification → Supervised Learning policy network → Self Play → Reinforcement Learning policy network → Self Play → Self-play data → Regression → Value network

Google DeepMind

# Supervised learning of policy networks

**Policy network:** 12 layer convolutional neural network

**Training data:** 30M positions from human expert games (KGS 5+ dan)

**Training algorithm:** maximise likelihood by stochastic gradient descent

$$\Delta\sigma \propto \frac{\partial \log p_\sigma(a|s)}{\partial \sigma}$$

**Training time:** 4 weeks on 50 GPUs using Google Cloud

**Results:** 57% accuracy on held out test data (state-of-the art was 44%)

# Reinforcement learning of policy networks

**Policy network:** 12 layer convolutional neural network

**Training data:** games of self-play between policy network

**Training algorithm:** maximise wins $z$ by policy gradient reinforcement learning

$$\Delta\sigma \propto \frac{\partial \log p_\sigma(a|s)}{\partial \sigma} z$$

**Training time:** 1 week on 50 GPUs using Google Cloud

**Results:** 80% vs supervised learning. Raw network ~3 amateur dan.

Google DeepMind

# Reinforcement learning of value networks

**Value network:** 12 layer convolutional neural network
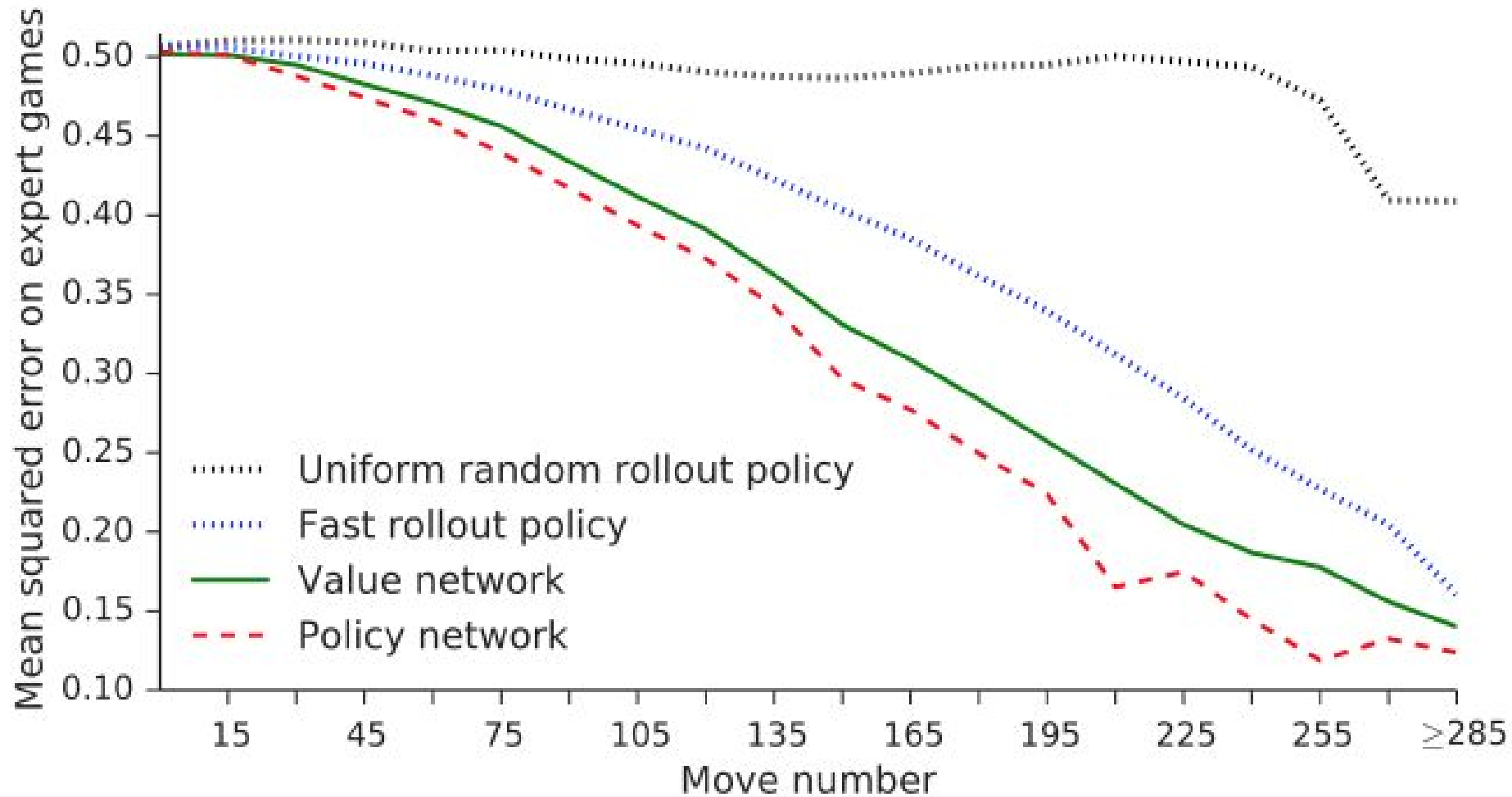
**Training data:** 30 million games of self-play

**Training algorithm:** minimise MSE by stochastic gradient descent

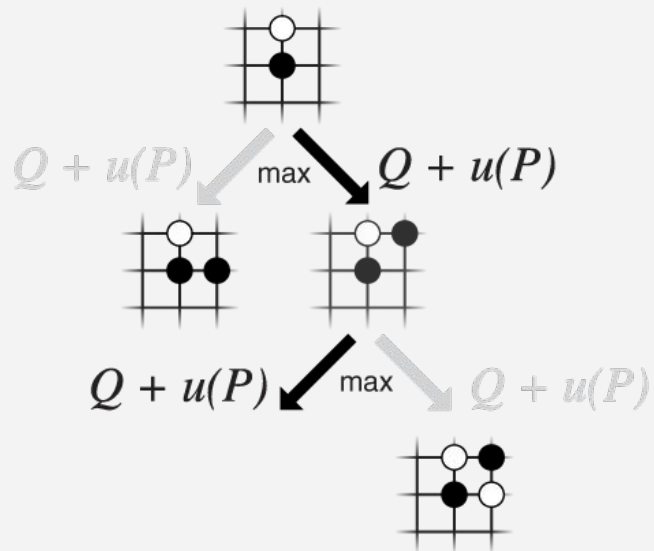$$\Delta\theta \propto \frac{\partial v_\theta(s)}{\partial \theta}(z - v_\theta(s))$$

**Training time:** 1 week on 50 GPUs using Google Cloud

**Results:** First strong position evaluation function - previously thought impossible
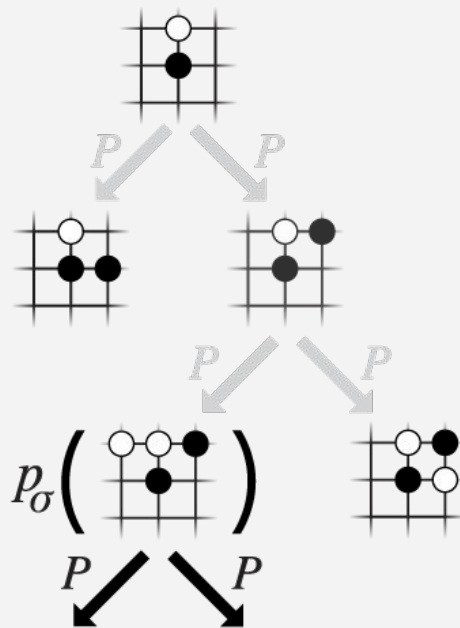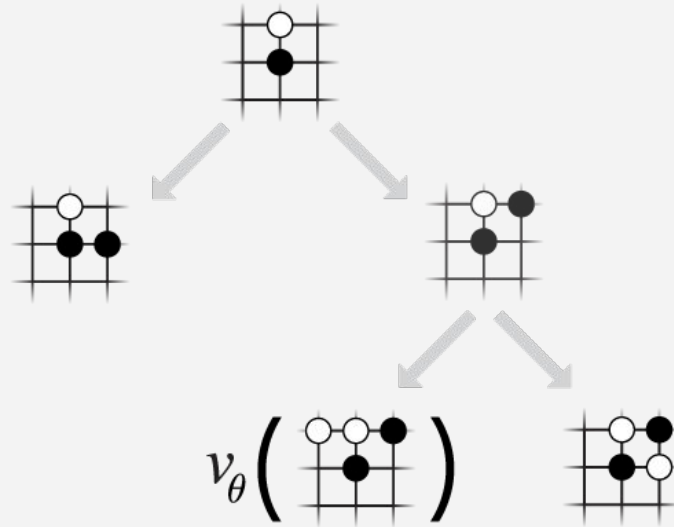
# Monte-Carlo tree search in AlphaGo: **selection**



$Q + u(P)$    max    $Q + u(P)$

$Q + u(P)$    max    $Q + u(P)$

$P$    prior probability
$Q$    action value

Google DeepMind

# Monte-Carlo tree search in AlphaGo: **expansion**



$p_\sigma$   Policy network
$P$   prior probability

Google DeepMind

# Monte-Carlo tree search in AlphaGo: **evaluation**



$v_\theta$    Value network

# Monte-Carlo tree search in AlphaGo: **rollout**
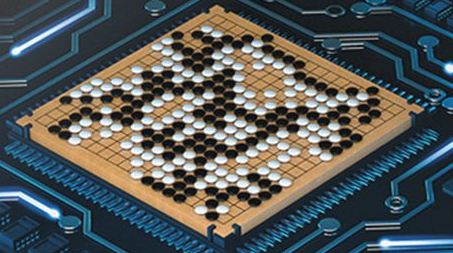


$v_\theta$  Value network
$r$  Game scorer

Google DeepMind

# Monte-Carlo tree search in AlphaGo: **backup**



$Q$   Action value
$v_\theta$   Value network
$r$   Game scorer

# nature

*At last — a computer program that can beat a champion Go player* **PAGE 484**

# ALL SYSTEMS GO

**CONSERVATION**
## SONGBIRDS À LA CARTE
*Illegal harvest of millions of Mediterranean birds*
PAGE 452

**RESEARCH ETHICS**
## SAFEGUARD TRANSPARENCY
*Don't let openness backfire on individuals*
PAGE 459

**POPULAR SCIENCE**
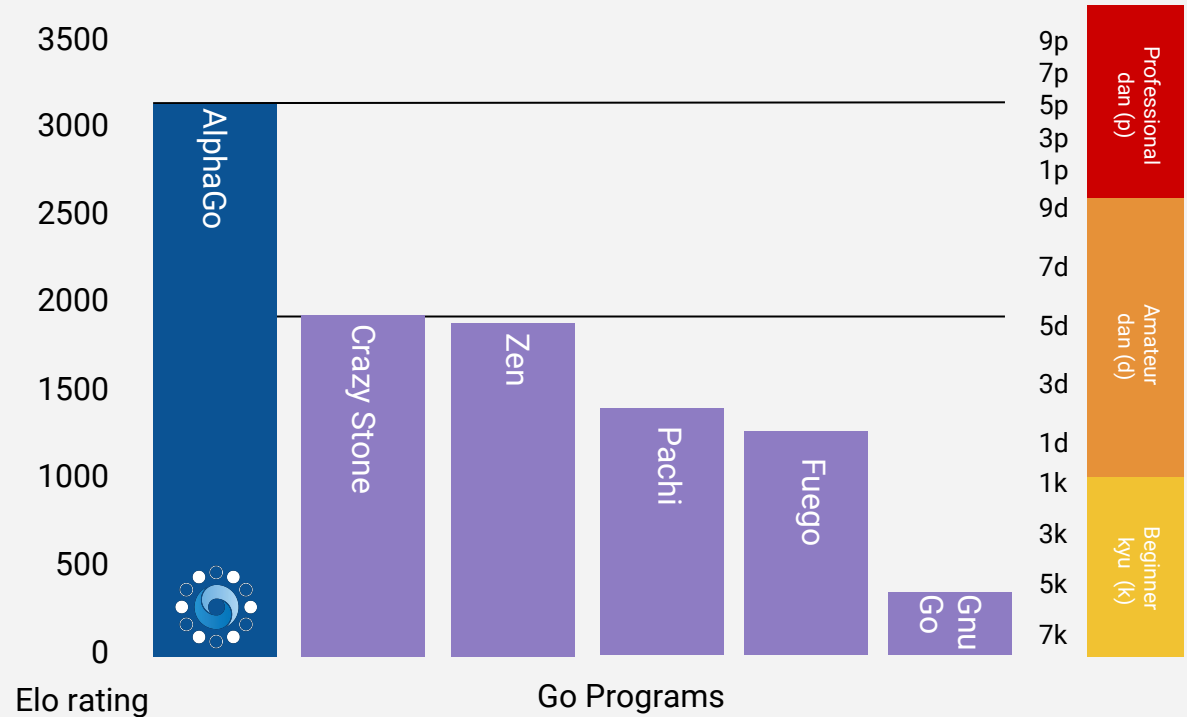## WHEN GENES GOT 'SELFISH'
*Dawkins's calling card 40 years on*
PAGE 462

# Evaluating Nature AlphaGo against computers

494/495 against computer opponents

>75% winning rate with 4 stone handicap

Even stronger using distributed machines



Elo rating

Go Programs

Bar chart showing Elo ratings for Go programs: AlphaGo (~3100), Crazy Stone (~1900), Zen (~1850), Pachi (~1400), Fuego (~1250), Gnu Go (~350). Right side scale shows ranks: Professional dan (p) 9p, 7p, 5p, 3p, 1p; Amateur dan (d) 9d, 7d, 5d, 3d, 1d; Beginner kyu (k) 1k, 3k, 5k, 7k.

# Evaluating Nature AlphaGo against humans

**Fan Hui** (2p): European Champion 2013 - 2016

Match was played in October 2015

AlphaGo won the match 5-0

**First program ever to beat a professional**

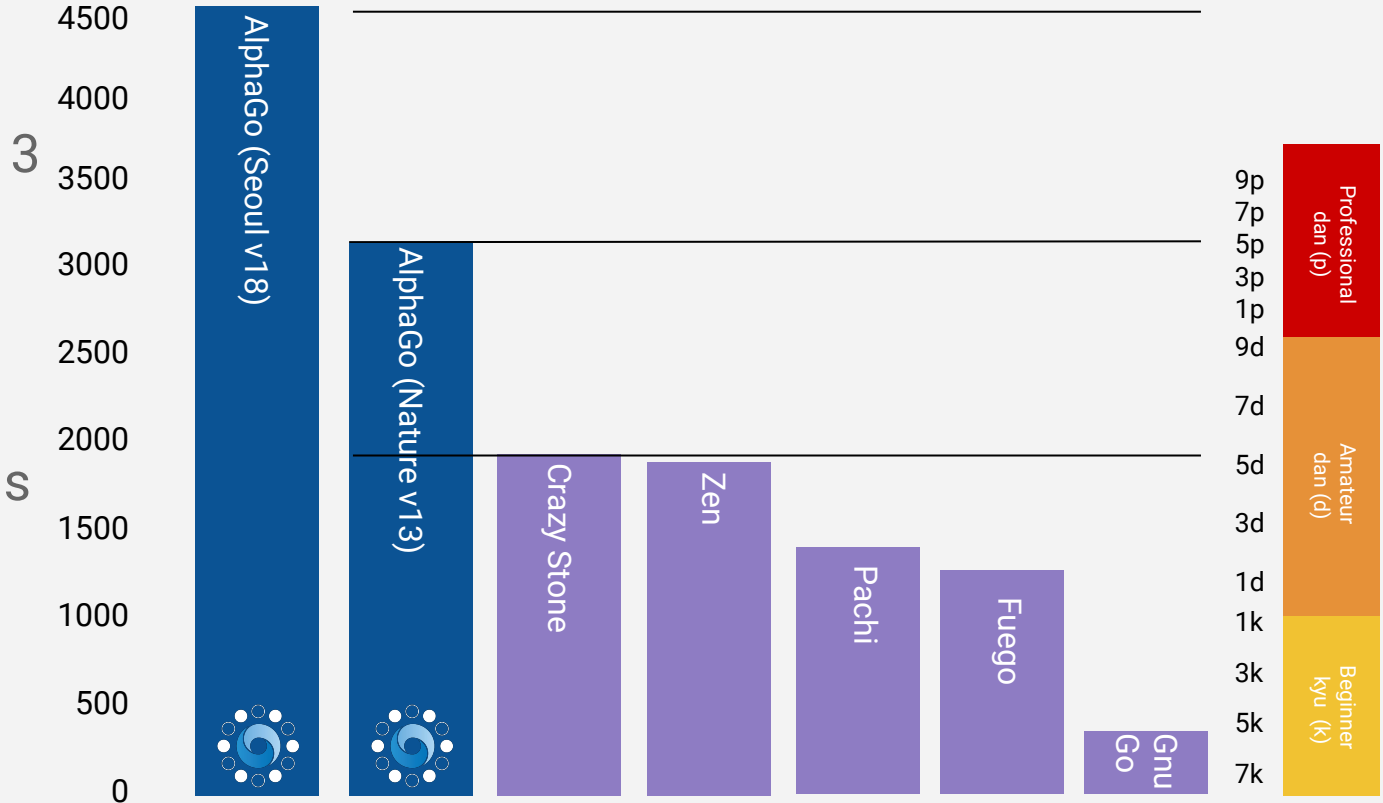on a full size 19x19 in an even game

# Seoul AlphaGo: Improvements

- Improved value network

- Improved policy network

- Improved search

- Improved hardware (TPU vs GPU)

Google DeepMind

# Evaluating Seoul AlphaGo against computers

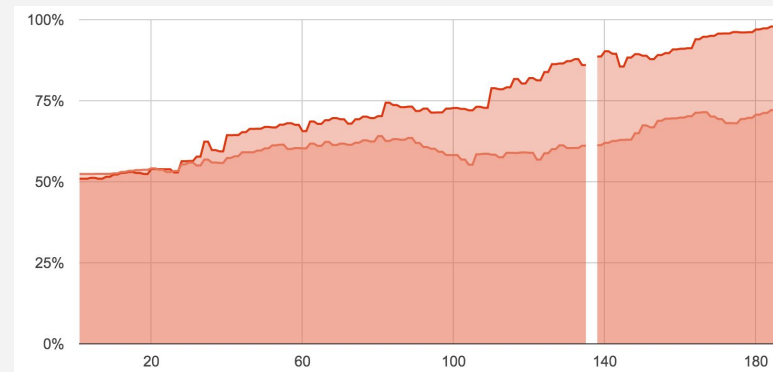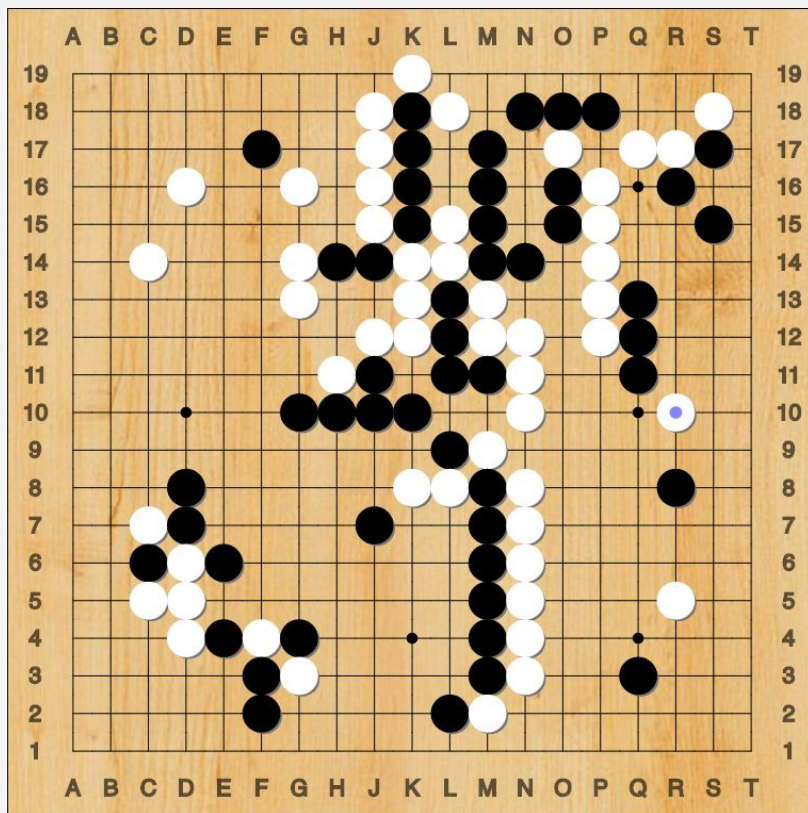Beats Nature AlphaGo with 3 to 4 stones handicap

CAUTION: ratings based on self-play results

| Rating | Program |
|---|---|
| | AlphaGo (Seoul v18) |
| | AlphaGo (Nature v13) |
| | Crazy Stone |
| | Zen |
| | Pachi |
| | Fuego |
| | Gnu Go |

4500
4000
3500
3000
2500
2000
1500
1000
500
0

9p
7p
5p
3p
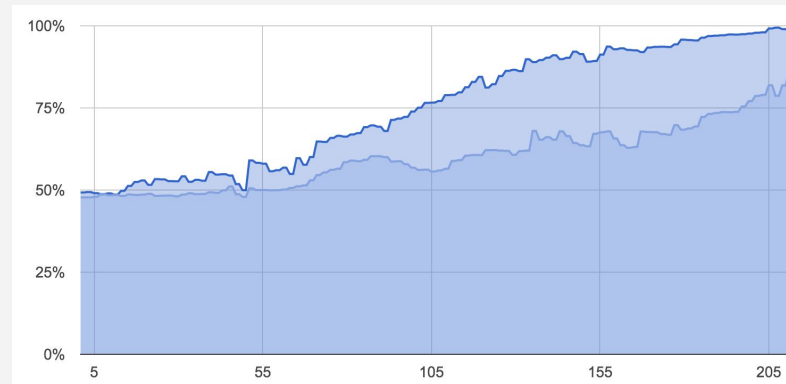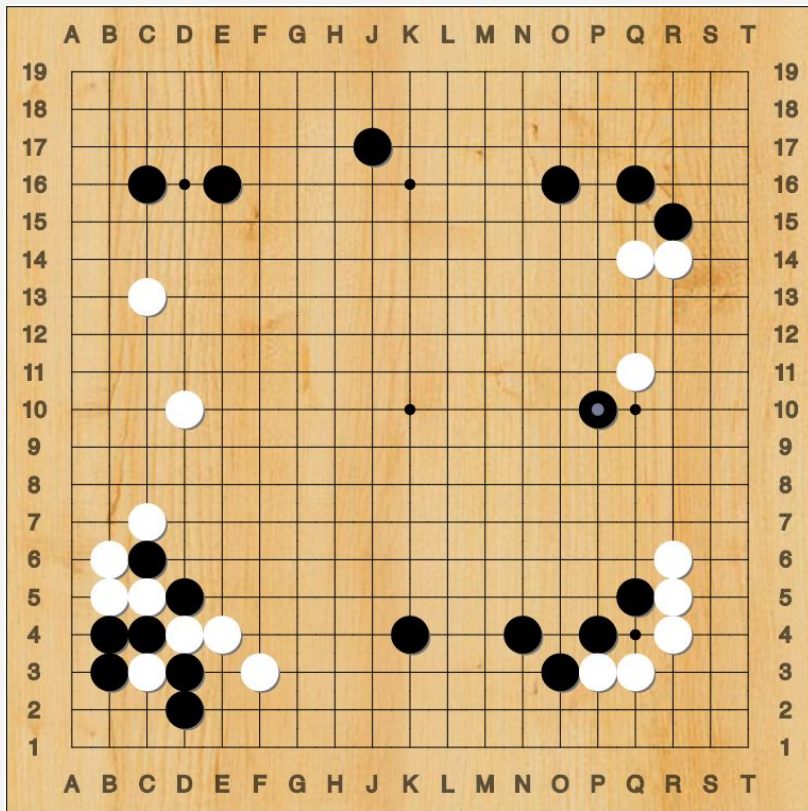1p
Professional dan (p)

9d
7d
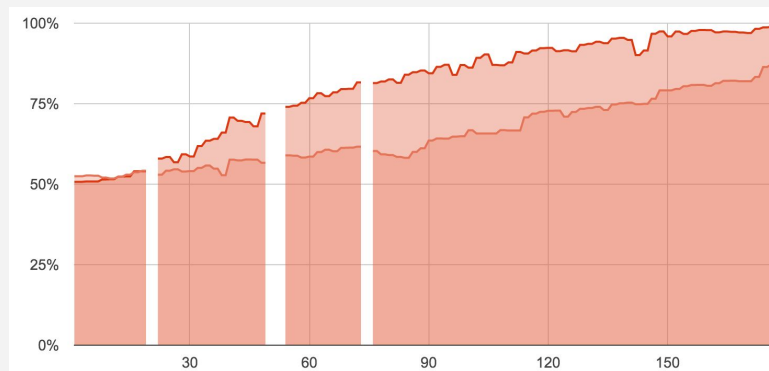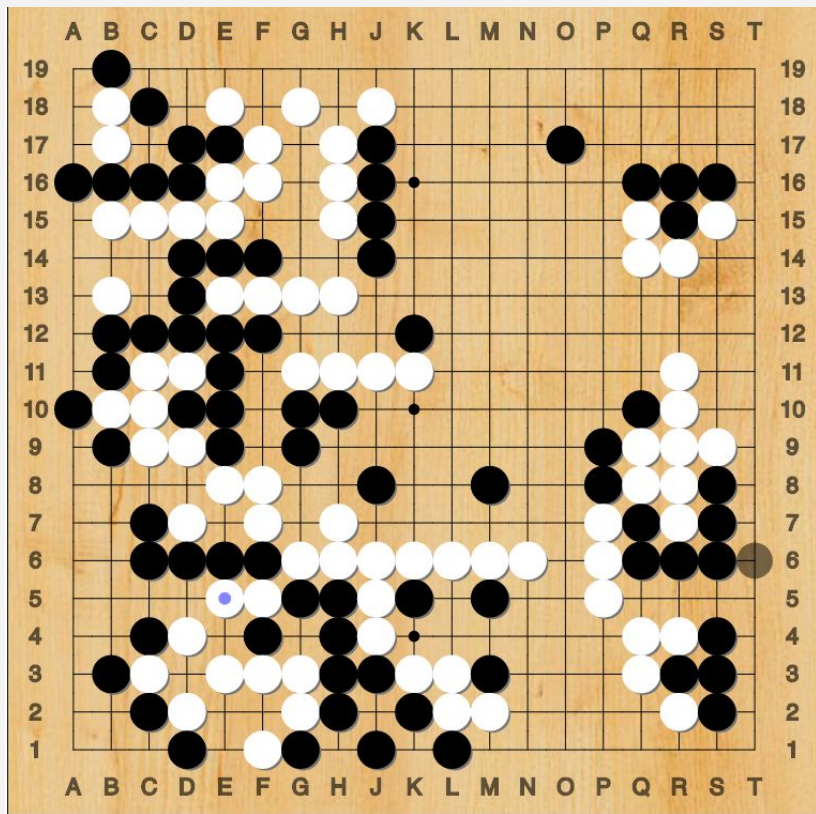5d
3d
1d
Amateur dan (d)

1k
3k
5k
7k
Beginner kyu (k)
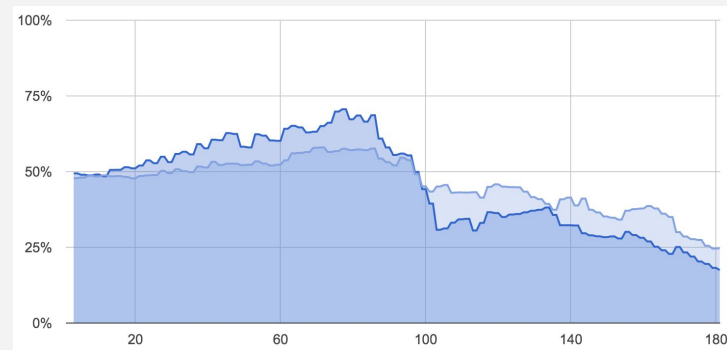
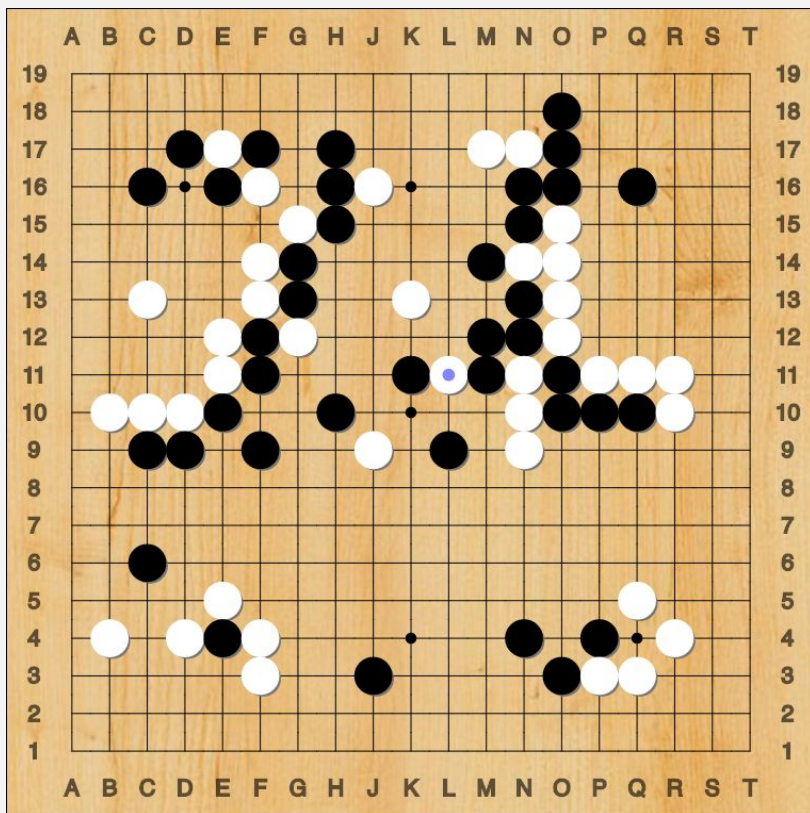# AlphaGo vs Lee Sedol: Game 1
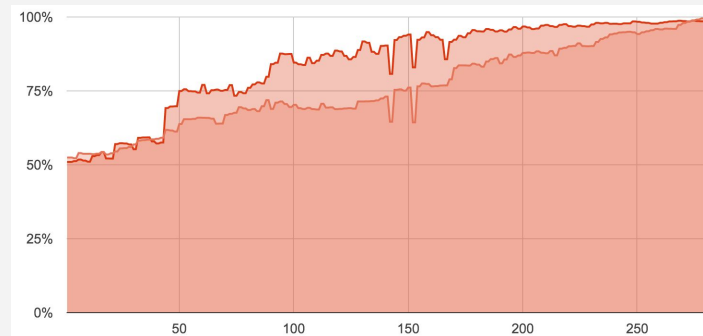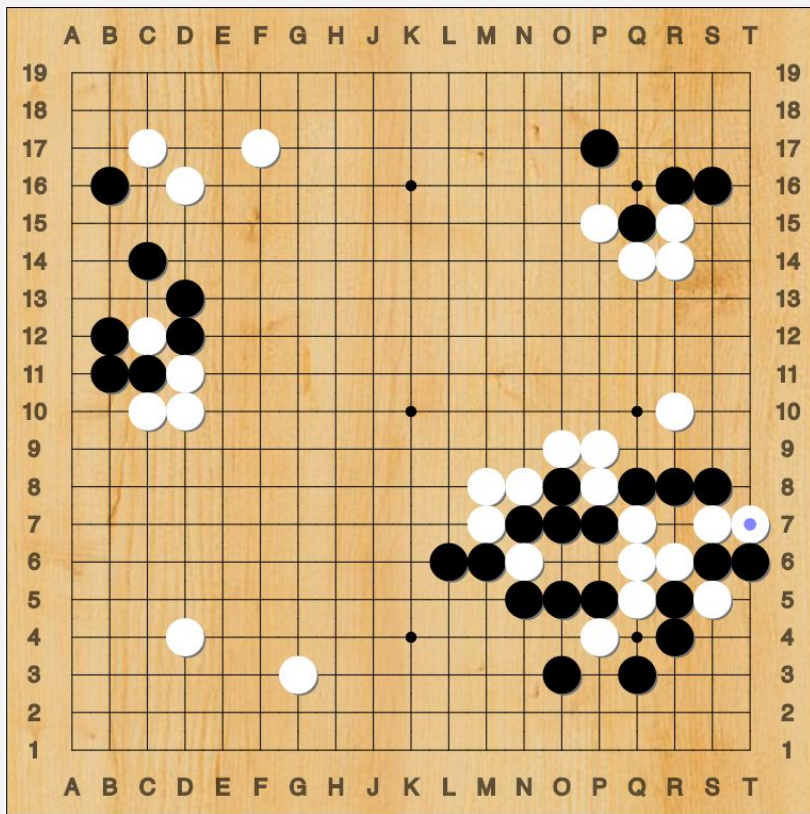
# AlphaGo vs Lee Sedol: Game 2

# AlphaGo vs Lee Sedol: Game 3

# AlphaGo vs Lee Sedol: Game 4

# AlphaGo vs Lee Sedol: Game 5

# Deep Blue

Handcrafted chess knowledge

Alpha-beta search guided by heuristic evaluation function
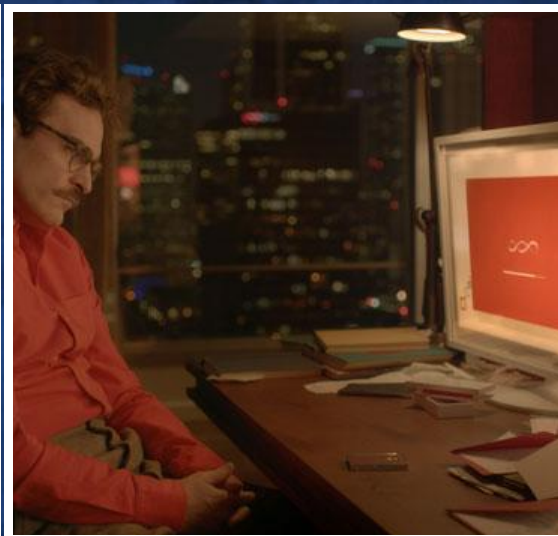
200 million positions / second

# AlphaGo

Knowledge learned from expert games and self-play

Monte-Carlo search guided by policy and value networks

60,000 positions / second

Google DeepMind

# What's Next?

# AlphaGo Team

Dave Silver • Aja Huang • Chris Maddison • Arthur Guez • Laurent Sifre • George Van Den Driessche • Julian Schrittwieser • Ioannis Antonoglou • Veda Panneershelvam • Yutian Chen

Marc Lanctot • Sander Dieleman • Dominik Grewe • John Nham • Nal Kalchbrenner • Tim Lillicrap • Maddy Leach • Koray Kavukcuoglu • Thore Graepel • Demis Hassabis

With thanks to: Lucas Baker, David Szepesvari, Malcolm Reynolds, Ziyu Wang, Nando De Freitas, Mike Johnson, Ilya Sutskever, Jeff Dean, Mike Marty, Sanjay Ghemawat.