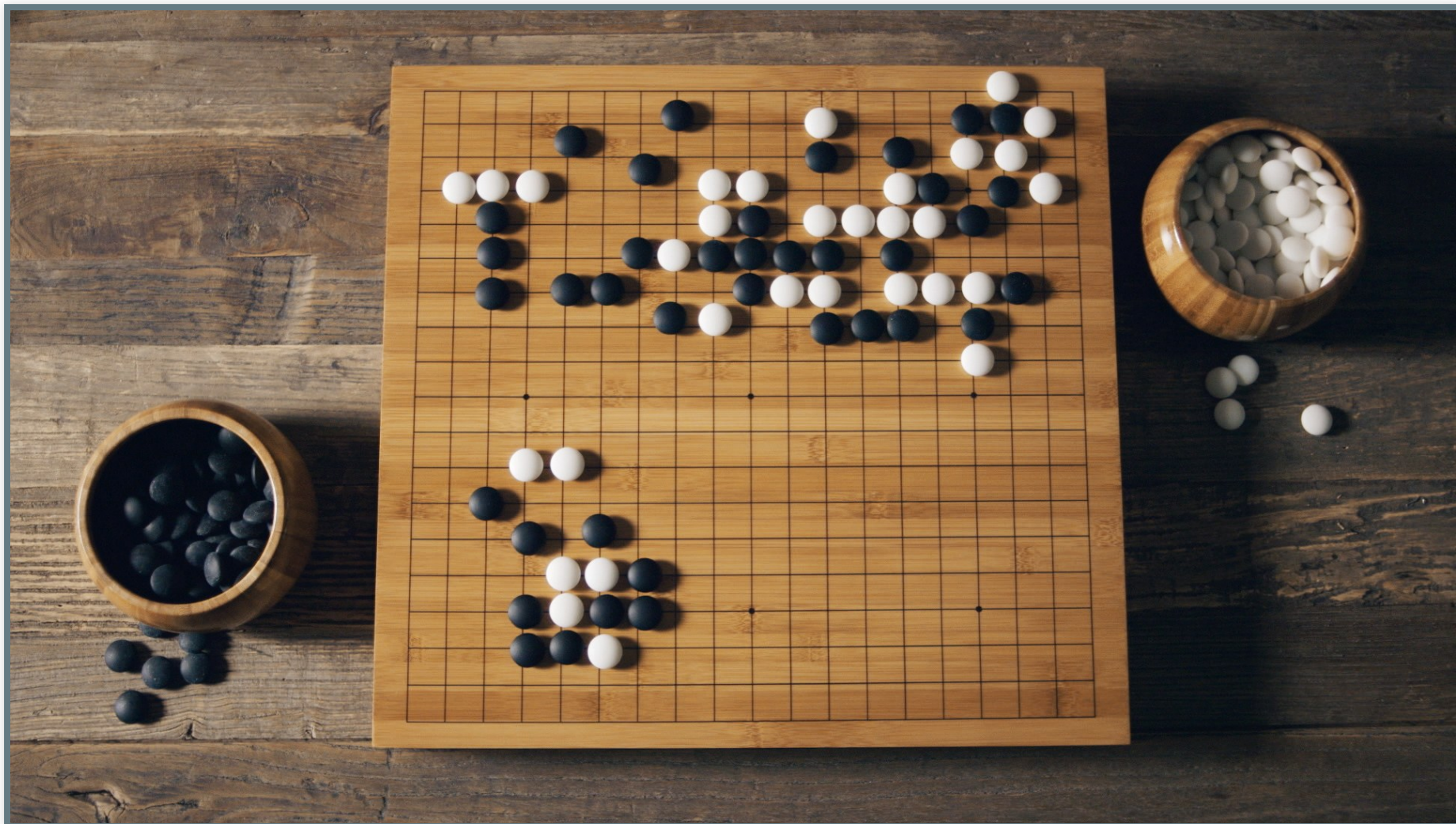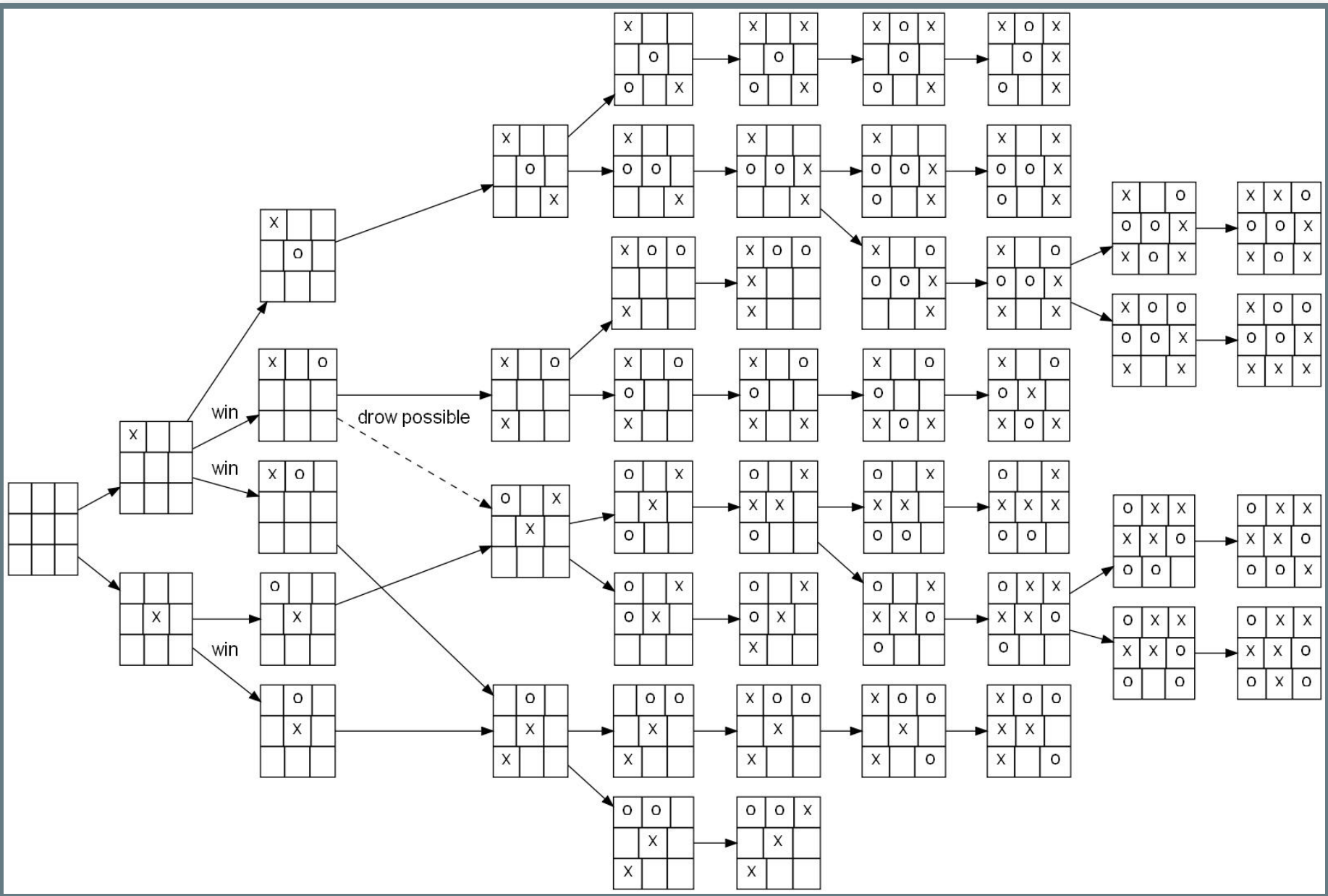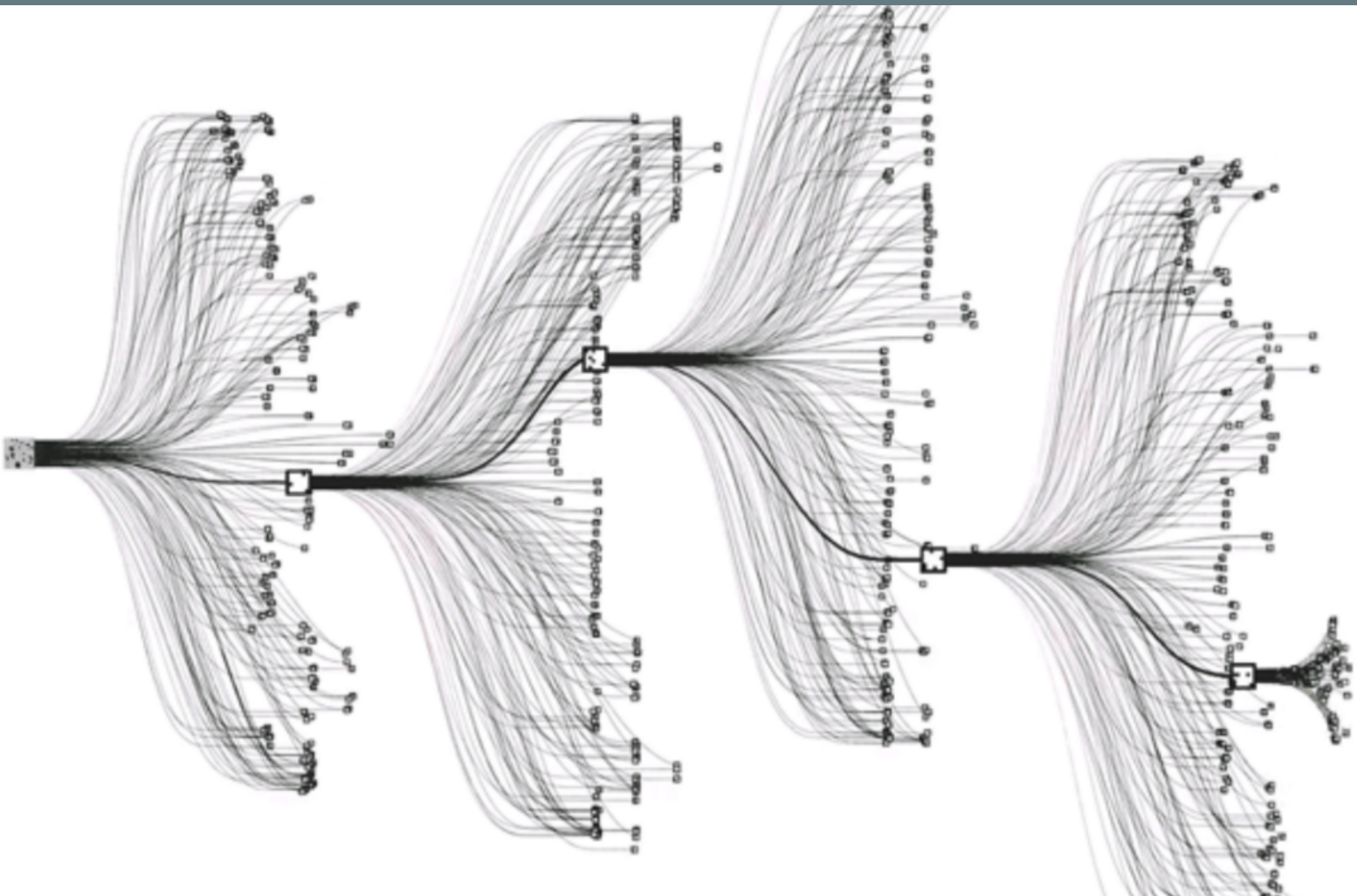$[\alpha_0]$ t0 3l0

With are only 3 equations ...

# Outline

1. The math behind Go
2. From Crazy Stone -> AlphaGO
3. AlphaGo vs AlphaZero
4. Policy Iteration
5. Policy Improvement (Math alert!)
6. Policy Evaluation
7. Code and demo

- $10^{170}$ possible states
- $10^{360}$ possible games for each starting state
- 250 legal moves from each state
- 150 moves for each match

# AI in Go

*"The mystery of Go, the ancient game that computers still can't win" - Wired 2014*

- Go is constructive
- Difficult to build an evaluation functional
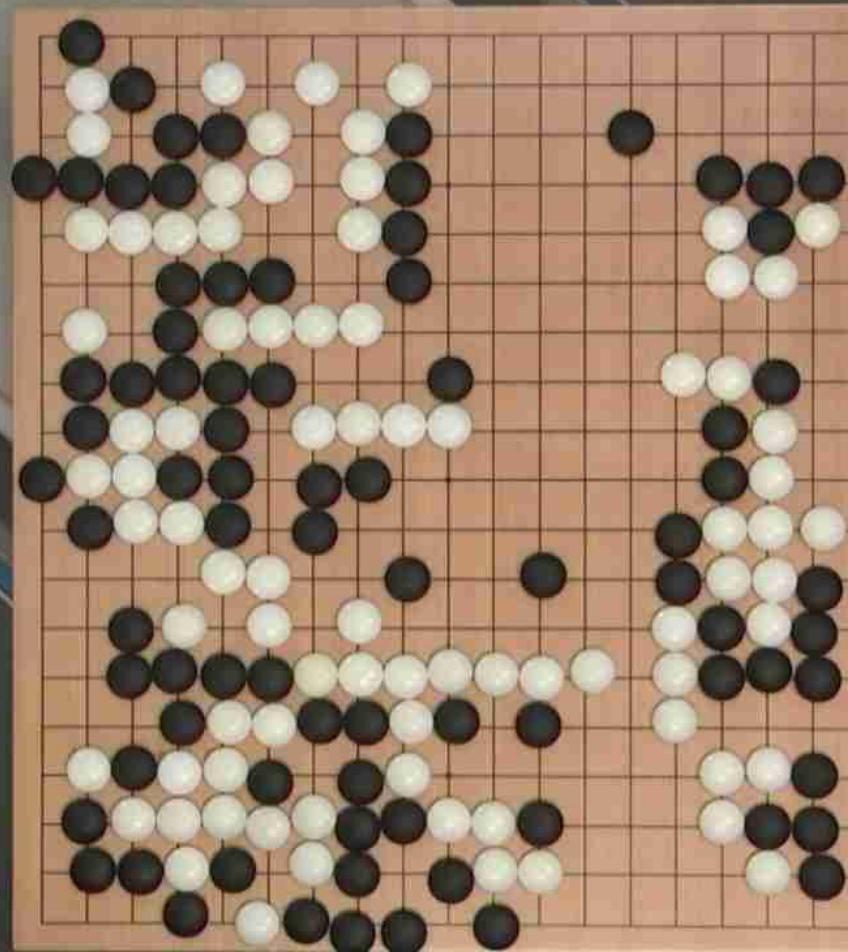- Humans describe more as intuititive game

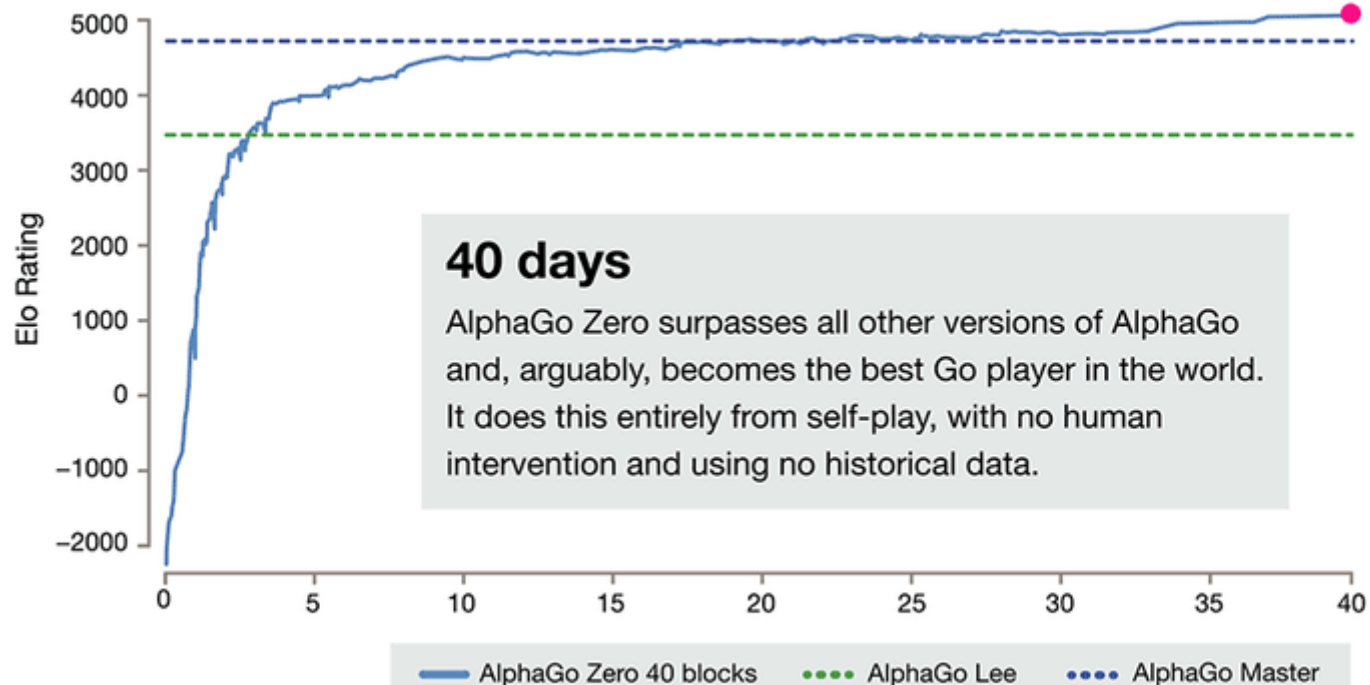- Adversarial
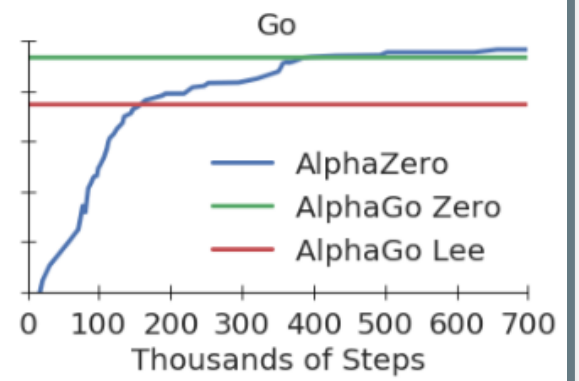- Deterministic
- Fully observable

# CrazyStone

AlphaGo

**40 days**

AlphaGo Zero surpasses all other versions of AlphaGo and, arguably, becomes the best Go player in the world. It does this entirely from self-play, with no human intervention and using no historical data.

Legend:
- AlphaGo Zero 40 blocks
- AlphaGo Lee
- AlphaGo Master
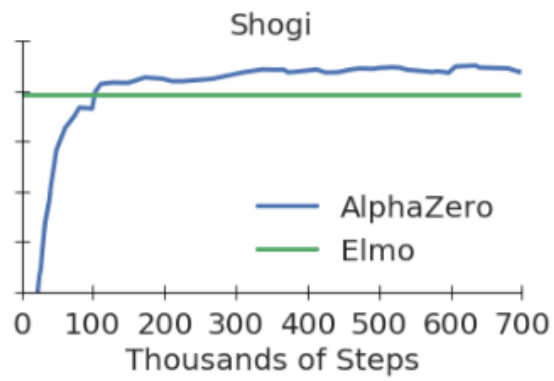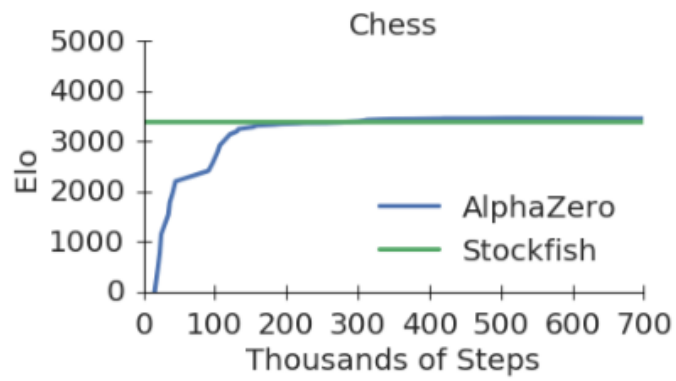
# AlphaGo Zero vs AlphaZero

- No data augmentation
- No threshold update
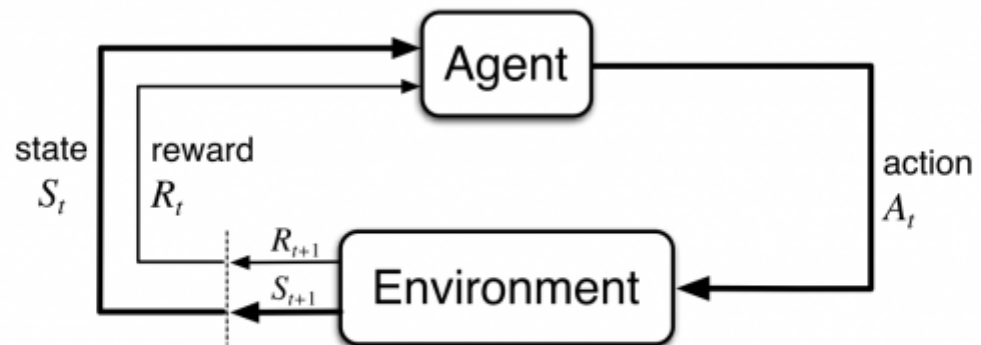- Diff. exploration noise for each game

# Reinforcement Learning

Figure 3.1: The agent–environment interaction in a Markov decision process.

$$v_\pi(s) = E_\pi \left[ \sum_t \gamma^t R_t \mid S_t \right]$$

where:

$$\pi(a \mid s) = P(a \mid s) \ \forall s \in S$$
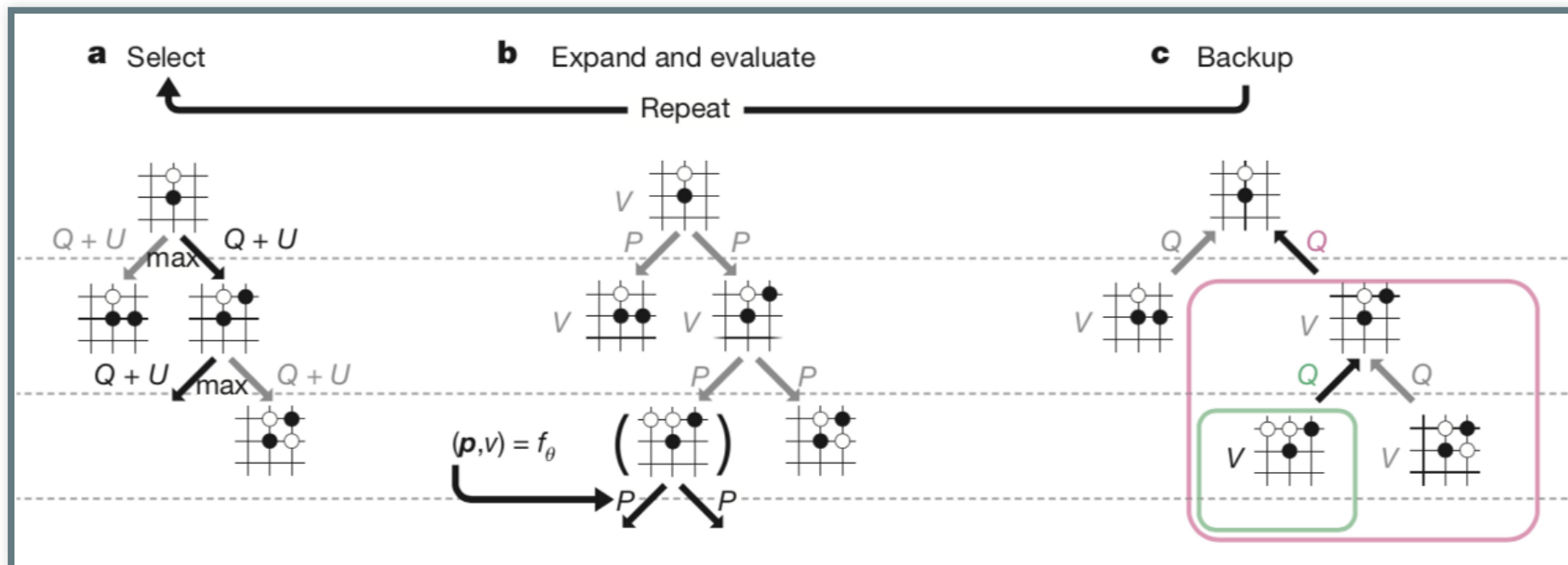
# Policy Iteration

- Add sudo code policy iterat

# Policy Improvement

# Monte-Carlo Tree Search

MCTS is an algorithm to perform sampling based lookahead search.
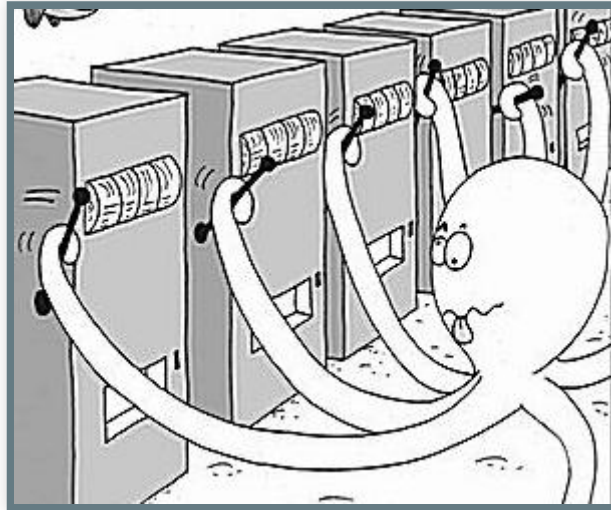
- Add sudo code MCTS

**a** Select

Q + U ⟶max⟶ Q + U

Q + U ⟶max⟶ Q + U

**b** Expand and evaluate

Repeat

V

P    P

V    V

P    P

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

P    P

**c** Backup

Q    Q

V

Q    Q

V    V

With the backup operation we keep track of:

- N(s,a) visit count
- W(s,a) total action value
- Q(s,a) mean action value
- P(s,a) prior probability

# Exploration

- $\epsilon - greedy$
- Bandits

$$cP(s, a) \frac{\sqrt{\sum_b N(s,b)}}{1+N(s,a)}$$

# Policy Evaluation

# Training

# Architecture

# Demo

Thank you!

github/mosc