



# Mastering the game of Go without human knowledge

Wei Wang 2019.12.10

East China Normal University  
Multi-Agent Artificial Intelligence Laboratory

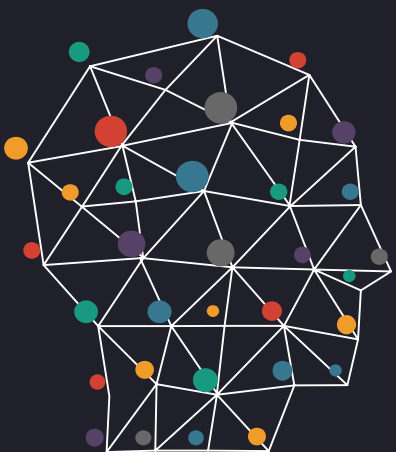
Silver D, Schrittwieser J, et al. Mastering the game of go without human knowledge[J]. Nature, 2017.



# CONTENTS

- 1. Introduction
- 2. Method
- 3. Experiment
- 4. Conclusion





# 1. Introduction

## 1.1 AlphaGo Fan & AlphaGo Zero

### AlphaGo Fan :

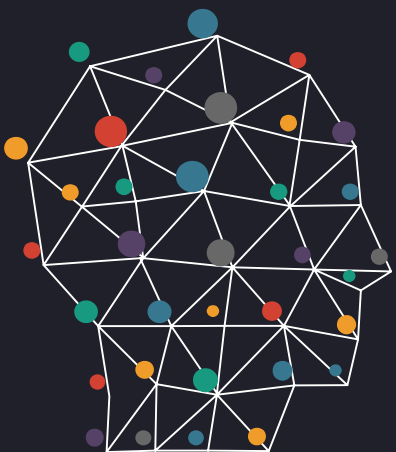
AlphaGo was the first program to achieve superhuman performance in Go.

Used two deep neural networks:

- A policy network that outputs move probabilities.
- A value network that outputs a position evaluation.

### AlphaGo Zero :

- It is trained by self-play reinforcement learning, starting from random play, without any supervision or use of human data.
- It uses only the black and white stones from the board as input features.
- It uses a single neural network, rather than separate policy and value networks.
- It uses a simpler tree search that relies upon this single neural network to evaluate positions and sample moves, without performing any Monte Carlo rollouts.



## 2. Method

## 2.1 Problem Description

### Go expression :

- $19 \times 19 = 361$  positions; white:1, black:-1, None:0
- Checkerboard state vector:

$$\vec{s} = (\underbrace{1, 0, -1, \dots}_{361})$$

- Next action vector:

$$\vec{a} = (\underbrace{1, 0, -1, \dots}_{361})$$

**Go problem :** Given a state arbitrarily, looking for the best policy, finally get the most site on the checkerboard.

## 2.2 Network Structure

### Deep neural network $f(\theta)$ :

- Parameters were used:  $\theta$
- Input: raw board representation  $s$  of the position and its history (Seven steps)
- Output: move probabilities and a value  $(p, v) = f\theta(s)$ 
  - The vector of move probabilities  $p$  represents the probability of selecting each move  $a$  (including pass),  $p_a = \Pr(a|s)$ .
  - The value  $v$  is a scalar evaluation, estimating the probability of the current player winning from position  $s$ .

The neural network consists of many residual blocks of convolutional layers with batch normalization and rectifier nonlinearities.

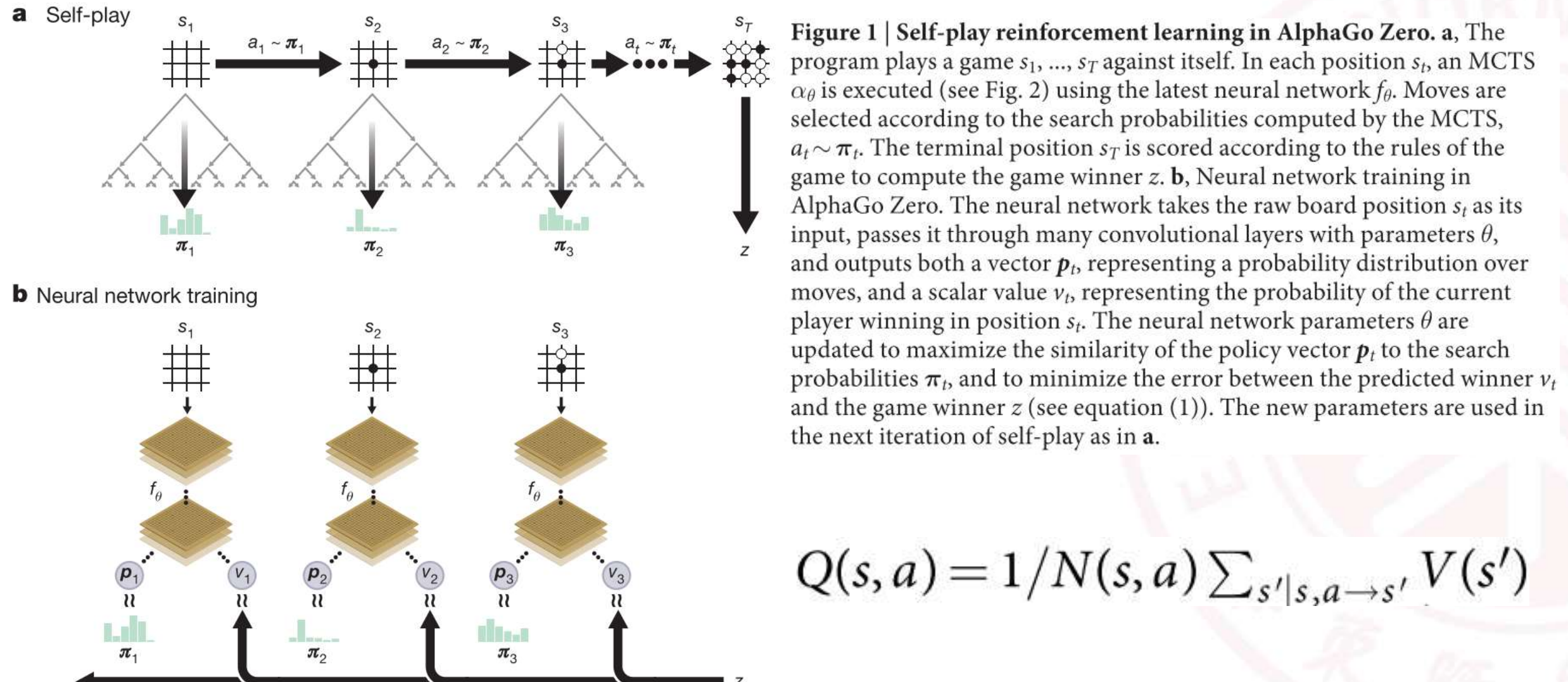
## 2.3 Reinforcement Learning

### Reinforcement Learning :

- Self - play
- In each position  $s$ , an MCTS search is executed, guided by the neural network  $f\theta$ .
- These search probabilities usually select much stronger moves than the raw move probabilities  $p$  of the neural network  $f\theta(s)$ .
- Self-play with search—using the improved MCTS based policy to select each move, then using the game winner  $z$  as a sample of the value—may be viewed as a powerful policy evaluation operator.
- Use these search operators repeatedly in a policy iteration procedure: the neural network's parameters are updated to make the move probabilities and value more closely match the improved search probabilities and self-play winner  $(\pi, z)$ ; these new parameters are used in the next iteration of self-play to make the search even stronger.

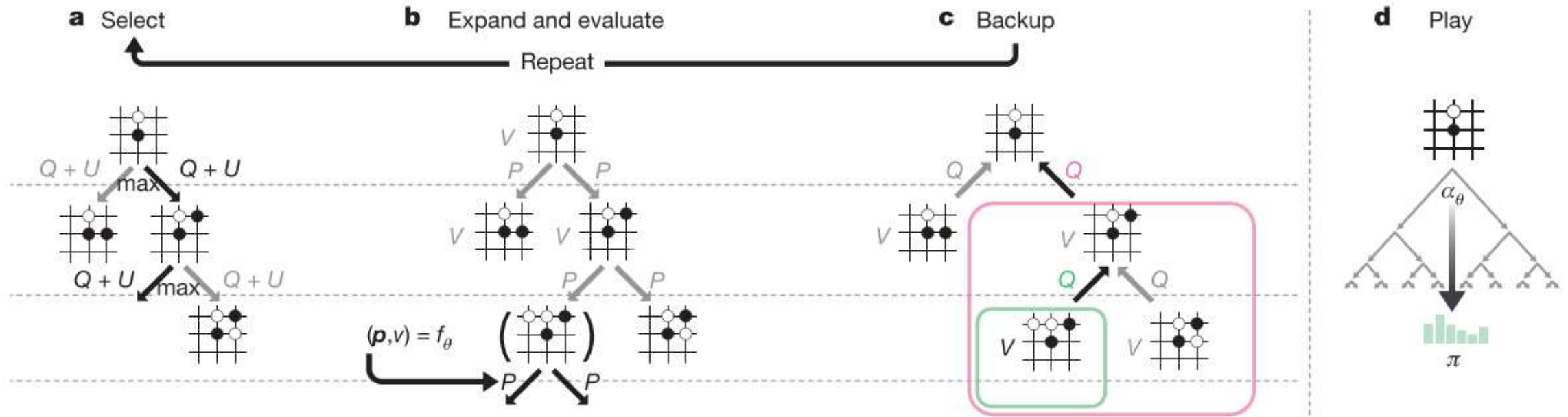


## 2.4 Self-play reinforcement learning



$$Q(s, a) = 1/N(s, a) \sum_{s'|s, a \rightarrow s'} V(s')$$

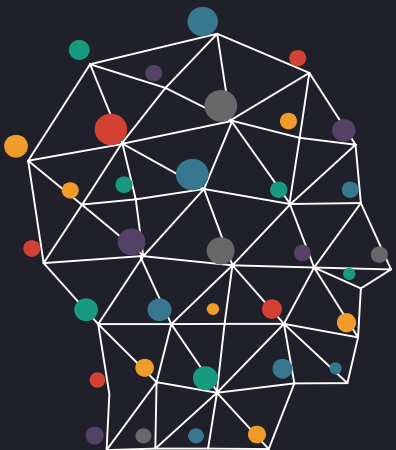
## 2.5 MCTS



**Figure 2 | MCTS in AlphaGo Zero.** **a**, Each simulation traverses the tree by selecting the edge with maximum action value  $Q$ , plus an upper confidence bound  $U$  that depends on a stored prior probability  $P$  and visit count  $N$  for that edge (which is incremented once traversed). **b**, The leaf node is expanded and the associated position  $s$  is evaluated by the neural network  $(P(s, \cdot), V(s)) = f_\theta(s)$ ; the vector of  $P$  values are stored in

the outgoing edges from  $s$ . **c**, Action value  $Q$  is updated to track the mean of all evaluations  $V$  in the subtree below that action. **d**, Once the search is complete, search probabilities  $\pi$  are returned, proportional to  $N^{1/\tau}$ , where  $N$  is the visit count of each move from the root state and  $\tau$  is a parameter controlling temperature.

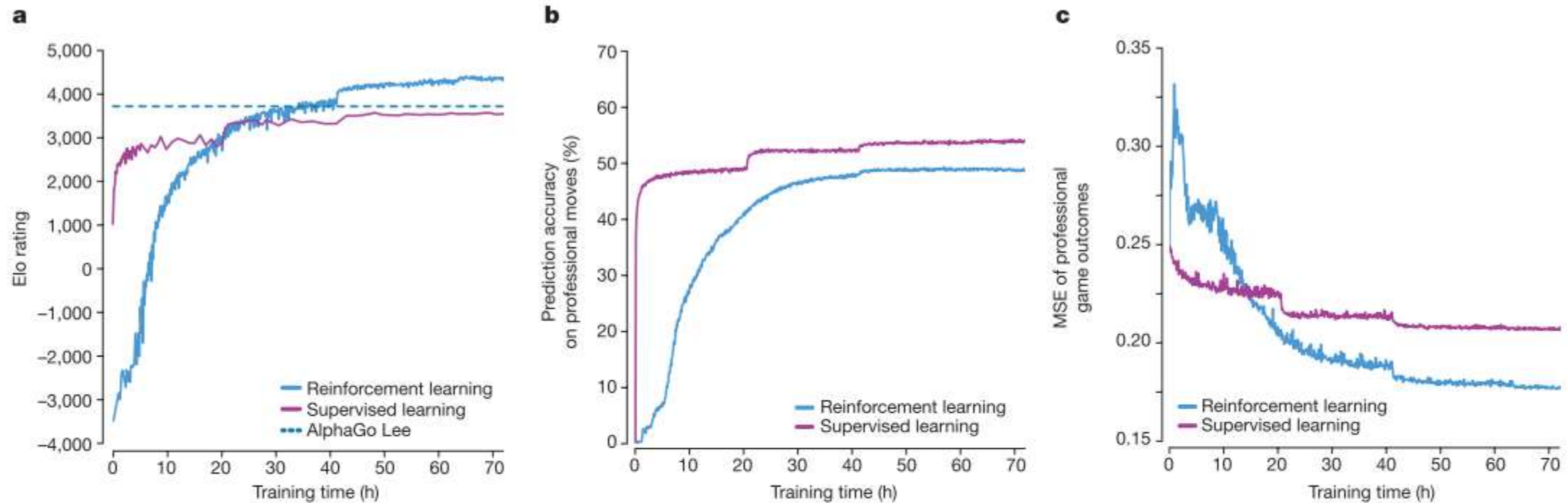
$$(p, v) = f_\theta(s) \quad \text{and} \quad l = (z - v)^2 - \pi^T \log p + c \|\theta\|^2$$



## 3. Experiment



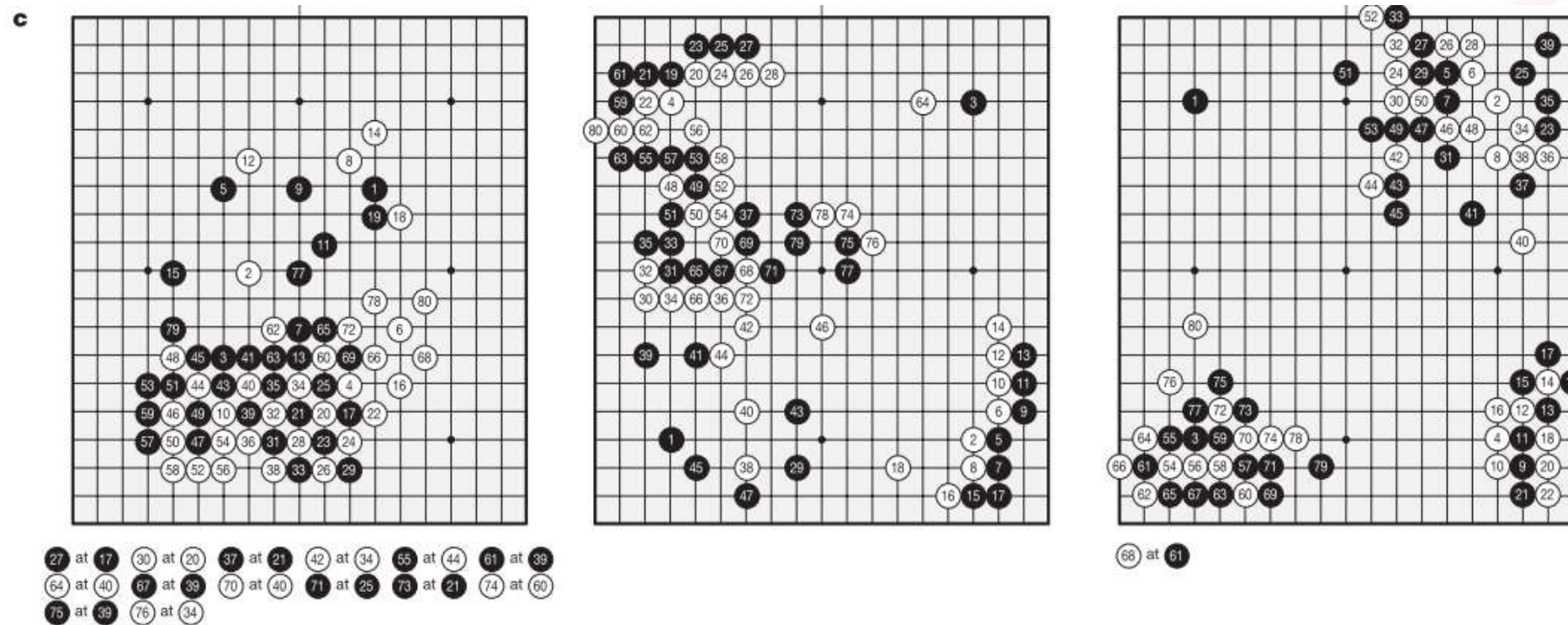
## 3.1 Empirical evaluation of AlphaGo Zero



**Figure 3 | Empirical evaluation of AlphaGo Zero.** **a**, Performance of self-play reinforcement learning. The plot shows the performance of each MCTS player  $\alpha_{\theta_i}$  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  $f_{\theta_i}$ , at each iteration of self-play  $i$ , in predicting human professional moves from the GoKifu dataset. The accuracy measures the

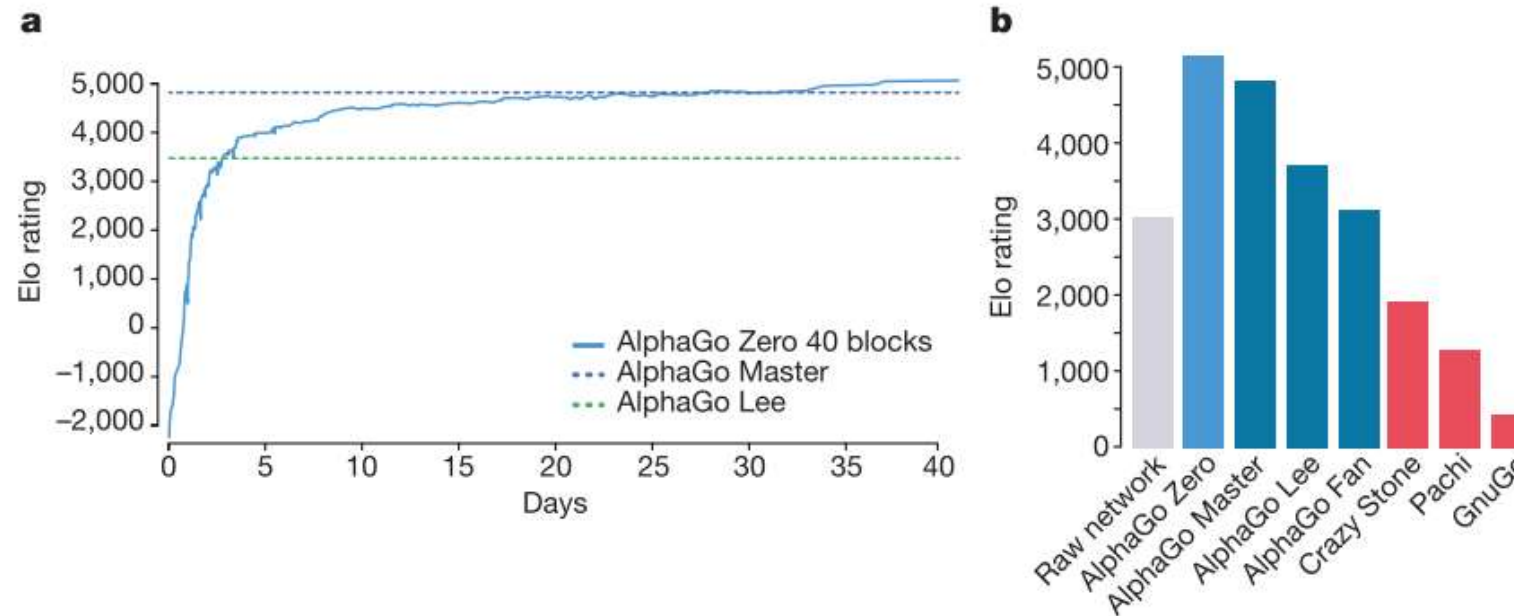
percentage of positions in which the neural network assigns the highest probability to the human move. The accuracy of a neural network trained by supervised learning is also shown. **c**, Mean-squared error (MSE) of human professional game outcomes. The plot shows the MSE of the neural network  $f_{\theta_i}$ , 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.

## 3.3 Knowledge learned by AlphaGo Zero



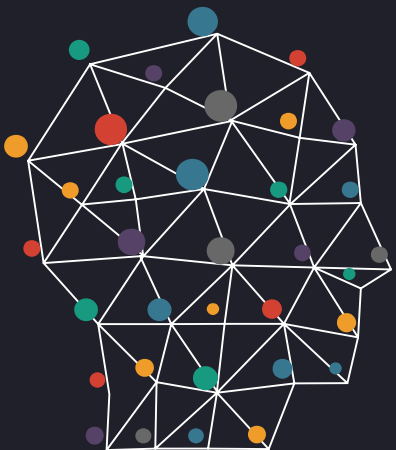


## 3.4 Final Performance of AlphaGo Zero



**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  $\alpha_{\theta_i}$  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). **b**, Final performance of AlphaGo Zero. AlphaGo Zero was trained for 40 days using a 40-block residual neural network. The plot shows the results of a tournament between: AlphaGo Zero, AlphaGo Master (defeated top human professionals 60–0 in online games), AlphaGo

Lee (defeated Lee Sedol), AlphaGo Fan (defeated Fan Hui), as well as previous Go programs Crazy Stone, Pachi and GnuGo. Each program was given 5 s of thinking time per move. AlphaGo Zero and AlphaGo Master played on a single machine on the Google Cloud; AlphaGo Fan and AlphaGo Lee were distributed over many machines. The raw neural network from AlphaGo Zero is also included, which directly selects the move  $a$  with maximum probability  $p_a$ , without using MCTS. Programs were evaluated on an Elo scale<sup>25</sup>: a 200-point gap corresponds to a 75% probability of winning.



## 4. Conclusion



## 4 Conclusion

Our results comprehensively demonstrate that a pure reinforcement learning approach is fully feasible, even in the most challenging of domains: it is possible to train to superhuman level, without human examples or guidance, given no knowledge of the domain beyond basic rules. Furthermore, a pure reinforcement learning approach requires just a few more hours to train, and achieves much better asymptotic performance, compared to training on human expert data. Using this approach, AlphaGo Zero defeated the strongest previous versions of AlphaGo, which were trained from human data using handcrafted features, by a large margin.

Humankind has accumulated Go knowledge from millions of games played over thousands of years, collectively distilled into patterns, proverbs and books. In the space of a few days, starting tabula rasa, AlphaGo Zero was able to rediscover much of this Go knowledge, as well as novel strategies that provide new insights into the oldest of games.





Mastering the game of Go without human knowledge

**Thanks!**