

# Man vs. Machine In Weiqi

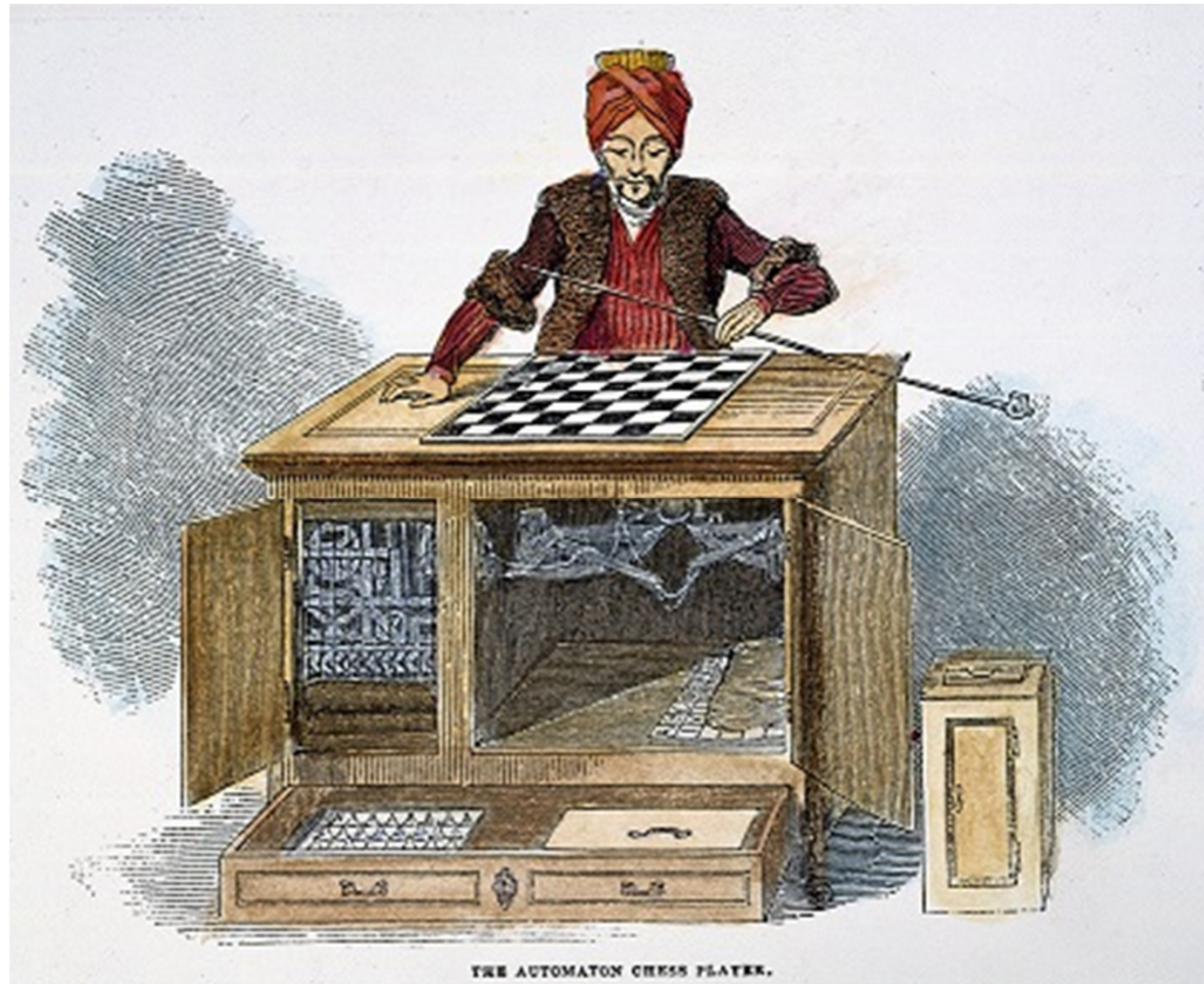
Recent development in Weiqi Machines (AlphaGo)

Li CHENG

31Mar16

# A Tale about the Chess Machine

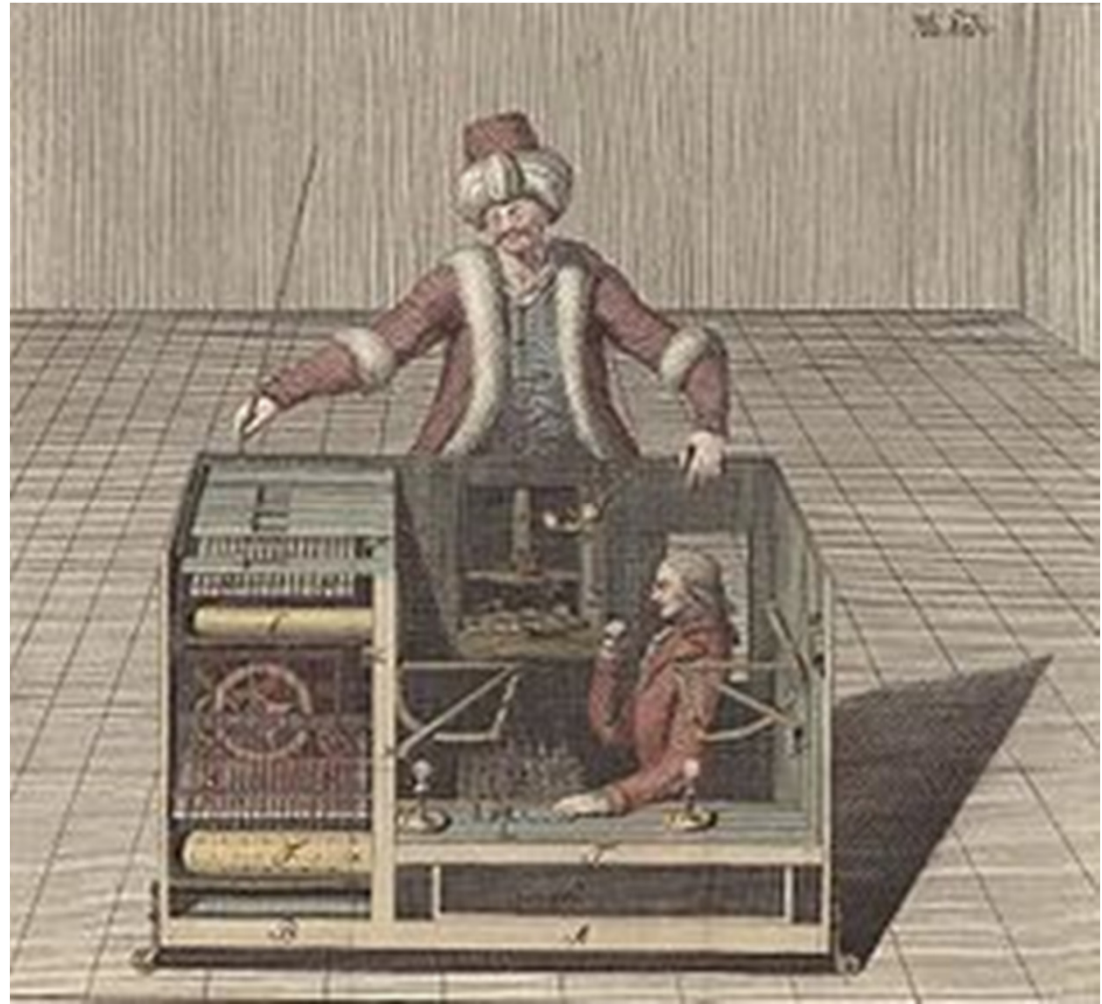
- Chess 1770



Wolfgang von Kempelen, "the Turk"

# A Tale about the Chess Machine

- Chess 1770



Wolfgang von Kempelen, "the Turk"

# The REAL Chess Machine

- Deep Blue vs. Garry Kasparov 1997, 3.5:2.5

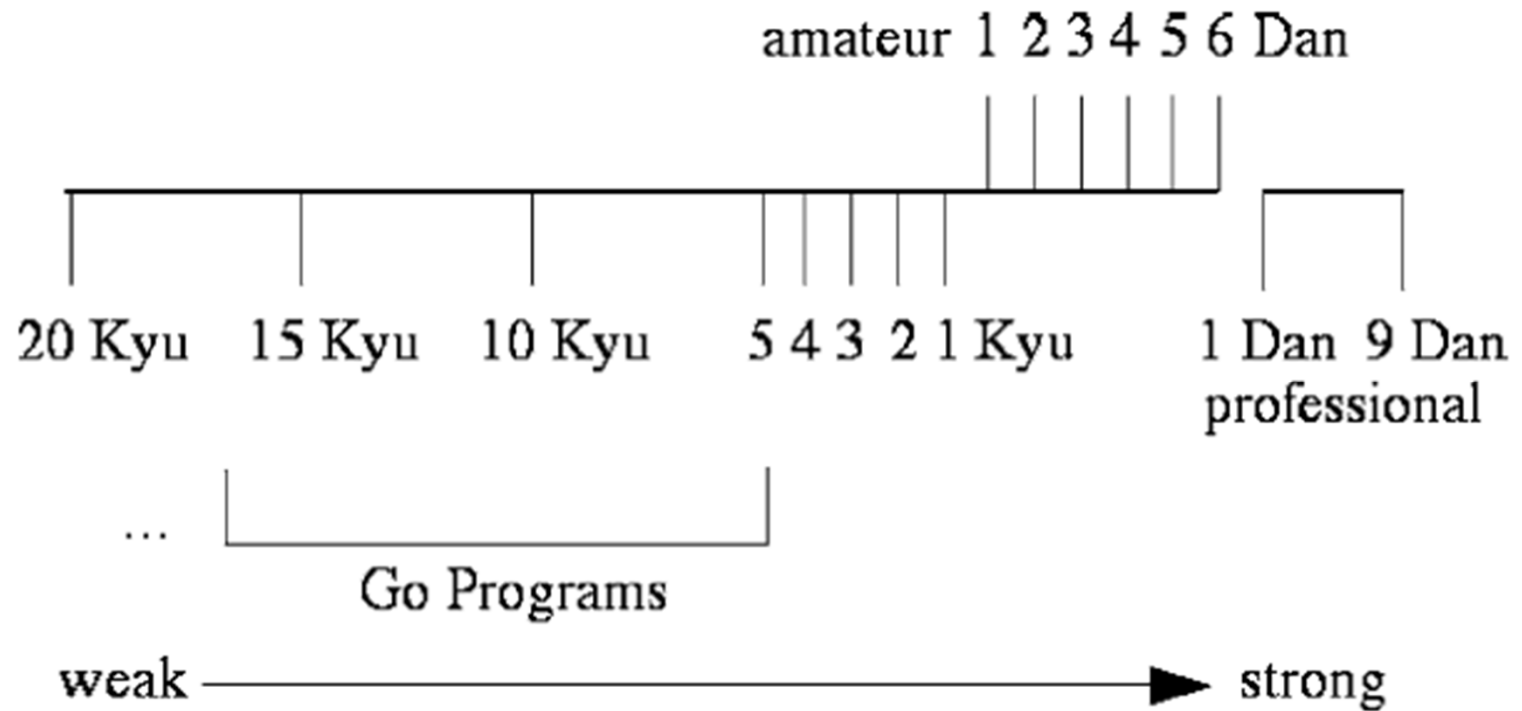




# Weiqi



# Weiqi



# Weiqi

Rank	Name	♂♀	Flag	Elo
1	<a href="#">Ke Jie</a>	♂		3615
2	<a href="#">Google AlphaGo</a>			3584
3	<a href="#">Park Jungwan</a>	♂		3551
4	<a href="#">Iyama Yuta</a>	♂		3534
5	<a href="#">Lee Sedol</a>	♂		3518
6	<a href="#">Shi Yue</a>	♂		3512
7	<a href="#">Kang Dongyun</a>	♂		3506
8	<a href="#">Park Yeonghun</a>	♂		3506
9	<a href="#">Mi Yuting</a>	♂		3504
10	<a href="#">Kim Jiseok</a>	♂		3498

Elo-like rating From (<http://www.goratings.org/>)

# Challenges of Weiqi Machines

- Weiqi has  $3^{361}$  ( about  $10^{170}$  ) states vs. only about  $10^{80}$  atoms in the entire universe
- Up to 361 legal moves
- Long-term effect of a move may only be revealed after hundreds of additional moves
- Last mind-sports game to be conquered



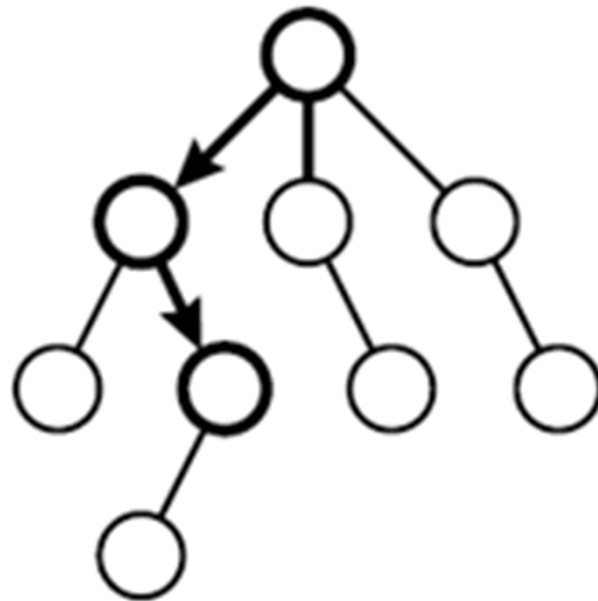


# Game Theory of Weiqi Machines

- A two-player, perfect information, zero-sum game (weiqi, chess, checker, ...)
- From MDP-> symmetric alternating Markov (SAM) game
- Two policies each for agent (black) and for opponent (white)
- Self-play policy: a policy used by both agent & opponent
- Minimax optimal policy for SAM game: a self play policy both maximizes agent's value function and minimized opponent's value function



# Game Tree Search of Weiqi Machines



# AlphaGo

## Mastering the game of Go with deep neural networks and tree search

David Silver<sup>1\*</sup>, Aja Huang<sup>1\*</sup>, Chris J. Maddison<sup>1</sup>, Arthur Guez<sup>1</sup>, Laurent Sifre<sup>1</sup>, George van den Driessche<sup>1</sup>, Julian Schrittwieser<sup>1</sup>, Ioannis Antonoglou<sup>1</sup>, Veda Panneershelvam<sup>1</sup>, Marc Lanctot<sup>1</sup>, Sander Dieleman<sup>1</sup>, Dominik Grewe<sup>1</sup>, John Nham<sup>2</sup>, Nal Kalchbrenner<sup>1</sup>, Ilya Sutskever<sup>2</sup>, Timothy Lillicrap<sup>1</sup>, Madeleine Leach<sup>1</sup>, Koray Kavukcuoglu<sup>1</sup>, Thore Graepel<sup>1</sup> & Demis Hassabis<sup>1</sup>

The game of Go has long been viewed as the most challenging of classic games for artificial intelligence owing to its enormous search space and the difficulty of evaluating board positions and moves. Here we introduce a new approach to computer Go that uses ‘value networks’ to evaluate board positions and ‘policy networks’ to select moves. These deep neural networks are trained by a novel combination of supervised learning from human expert games, and reinforcement learning from games of self-play. Without any lookahead search, the neural networks play Go at the level of state-

---

<sup>1</sup>Google DeepMind, 5 New Street Square, London EC4A 3TW, UK. <sup>2</sup>Google, 1600 Amphitheatre Parkway, Mountain View, California 94043, USA.

\*These authors contributed equally to this work.

# AlphaGo

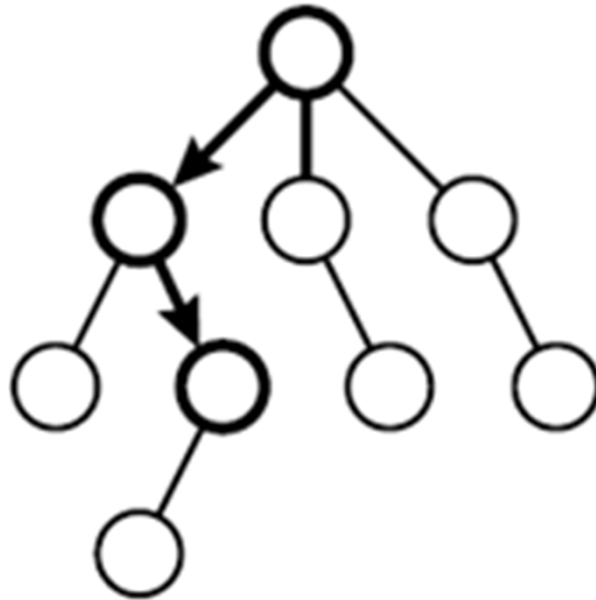
- An engineering feat (about 2-year time)
- In fact, many underground efforts over many years
- AlphaGo vs. Lee Seedal (9-15 Mar 2016): 4-1, Huge success, a real earthquake in Weiqi community
- 17Mar16: South Korean govt. will invest 1 trillion won (US\$0.86 Billion) in AI over next 5-year

# Three Pillars of AlphaGo

- (Monte Carlo) Tree Search or MCTS
- Reinforcement Learning (learning by self-play)
- Deep Learning (aka Deep Convolutional Neural Network, or DCNN)

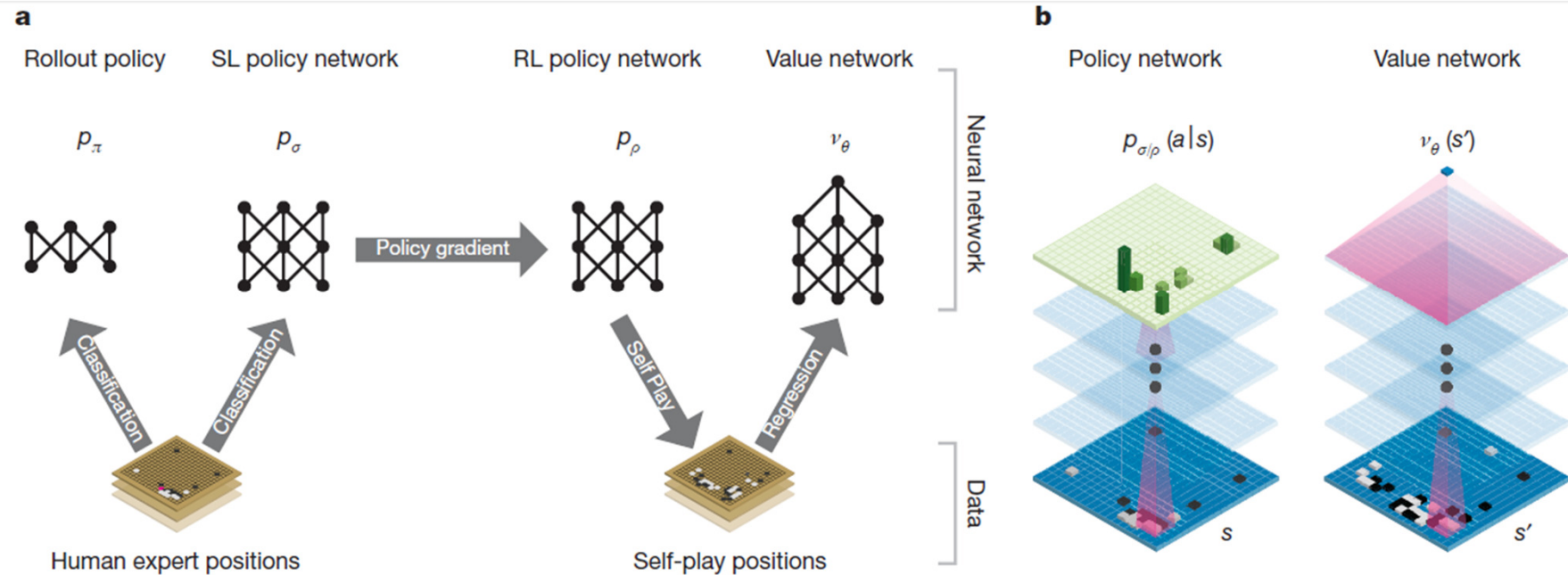


# Sketch the Main Ideas



# How AlphaGo Works

overview



**Figure 1 | Neural network training pipeline and architecture.** **a**, A fast rollout policy  $p_\pi$  and supervised learning (SL) policy network  $p_\sigma$  are trained to predict human expert moves in a data set of positions. A reinforcement learning (RL) policy network  $p_\rho$  is initialized to the SL policy network, and is then improved by policy gradient learning to maximize the outcome (that is, winning more games) against previous versions of the policy network. A new data set is generated by playing games of self-play with the RL policy network. Finally, a value network  $v_\theta$  is trained by regression to predict the expected outcome (that is, whether

the current player wins) in positions from the self-play data set. **b**, Schematic representation of the neural network architecture used in AlphaGo. The policy network takes a representation of the board position  $s$  as its input, passes it through many convolutional layers with parameters  $\sigma$  (SL policy network) or  $\rho$  (RL policy network), and outputs a probability distribution  $p_\sigma(a|s)$  or  $p_\rho(a|s)$  over legal moves  $a$ , represented by a probability map over the board. The value network similarly uses many convolutional layers with parameters  $\theta$ , but outputs a scalar value  $v_\theta(s')$  that predicts the expected outcome in position  $s'$ .

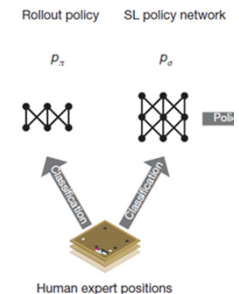
# How AlphaGo Works

Supervised learning of policy networks

- The SL network is a 13-layer CNN trained from 30m positions from KGS server, to predict expert's move
- Final layer is a softmax layer outputs probability over all legal moves
- Trained on randomly sampled (s,a) pairs, using SG to max likelihood of selecting move a at state s:

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

- 57% prediction accuracy with all input features, and 55.7% with only board position and move history, vs. 44.4% from state-of-the-arts



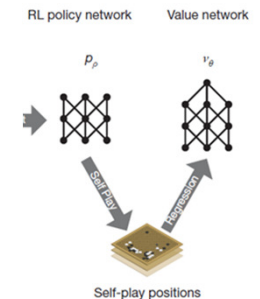
# How AlphaGo Works

Reinforcement learning of policy networks

- The RL policy network  $P_\rho$  starts with SL policy network  $P_\sigma$
- It learn by self-playing between itself and a randomly chosen previous version of itself
- The reward of curr state  $s$   $r(s)$  is obtained by playing multiple games to end (roll-outs) and obtaining each time either +1 (black win) or -1
- Params of the policy at curr state is updated as

$$\Delta\rho \propto \frac{\partial \log p_\rho(a_t|s_t)}{\partial \rho} z_t$$

- Evaluation of RL policy network is by sampling each move as  $a_t \sim P_\rho(\cdot|s_t)$  at curr state  $s_t$ .
- It won over SL policy network 80% of the games
- It won over Pachi (KGS 2D) 85% over 100,000 games



# How AlphaGo Works

Reinforcement learning of value networks

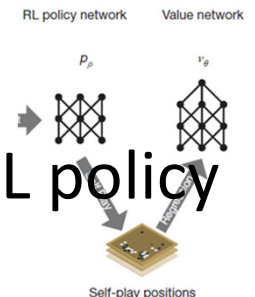
- Predict the outcome of curr state  $s$  played by policy  $p$  for both players  $v^p(s) = \mathbb{E}[z_t | s_t = s, a_{t...T} \sim p]$
- Ideally it is the optimal value func under perfect play,  $v^*(s)$ . In practice, it is  $v^{P_\rho}$ , the value function of our strongest play  $P_\rho$ . This is approximated by a  $\theta$ -parametric CNN, as

$$v_\theta(s) \approx v^{P_\rho}(s) \approx v^*(s)$$

- Same architecture as RL policy network. Only it now predicts a single value instead of a distribution.
- It is trained by regression on  $(s,a)$ , with SG to minimize the MSE between the predicted value  $v_\theta(s)$  and ideal outcome  $z$ , as

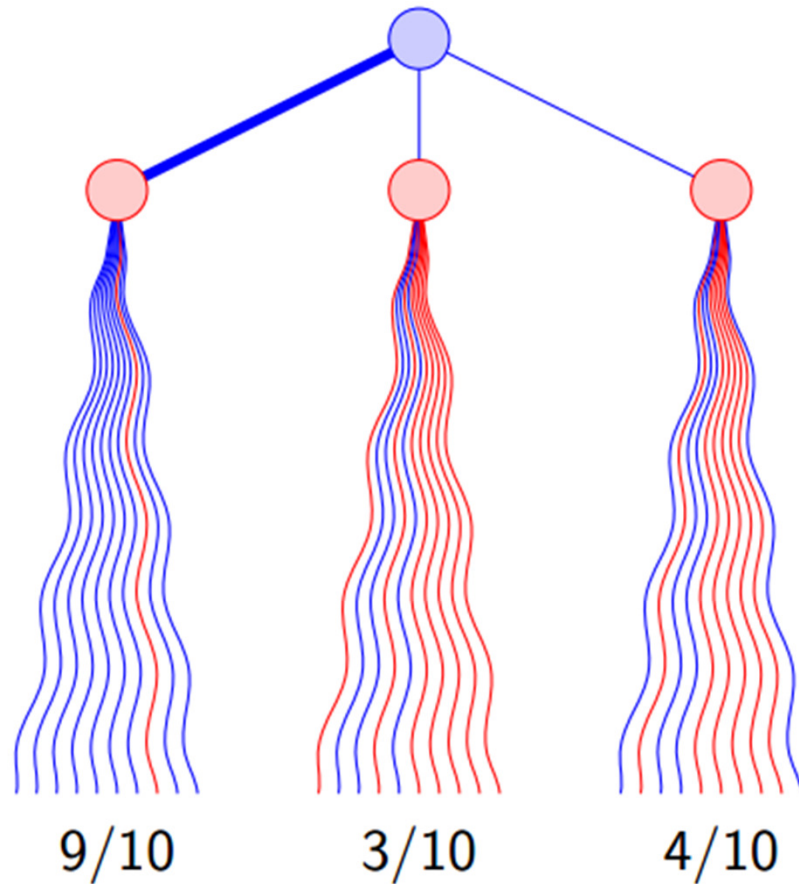
$$\Delta\theta \propto \frac{\partial v_\theta(s)}{\partial \theta} (z - v_\theta(s))$$

- Comparable with Monte Carlo rollouts using the RL policy network, but with 15K times less computation.





# Monte Carlo Tree Search



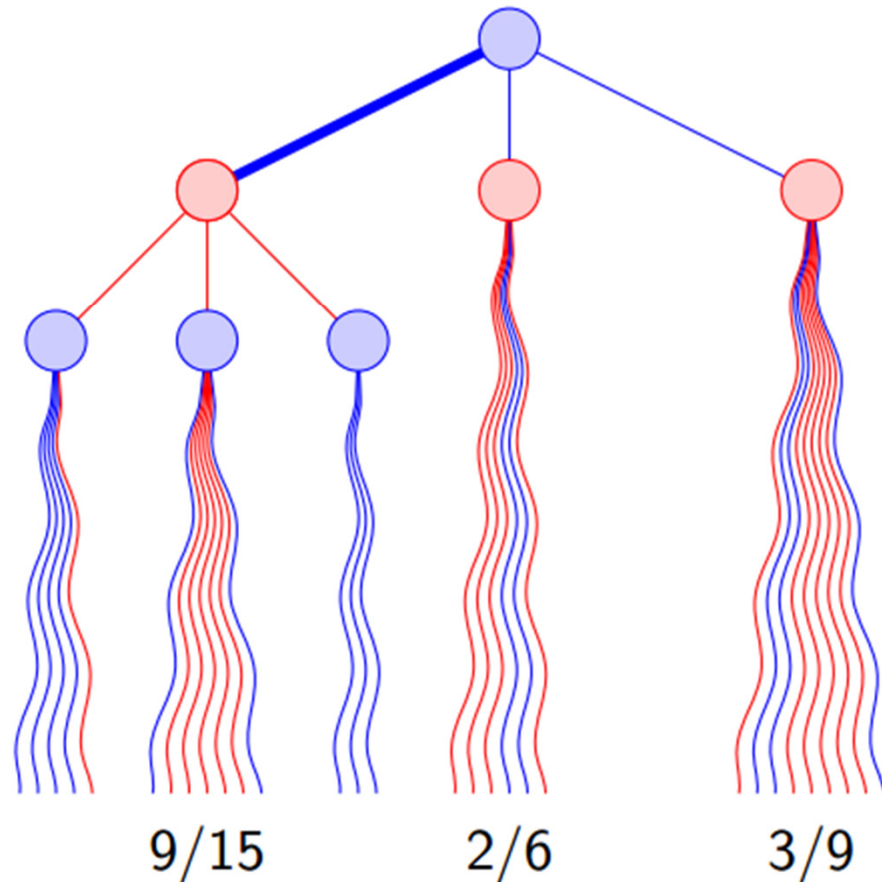
## Algorithm

- $N$  playouts for every move
- Pick the best winning rate
- 5,000 playouts/s on 19x19

## Problems

- Evaluation may be wrong
- For instance, if all moves lose immediately, except one that wins immediately.

# Monte Carlo Tree Search



## Principle

- More playouts to best moves
- Apply recursively
- Under some simple conditions: proven convergence to optimal move when  $\# \text{playouts} \rightarrow \infty$

From Remi Coulom, the author of Crazy Stone (<http://www.remi-coulom.fr/JFFoS/JFFoS.pdf>)

# How AlphaGo Works

———— MCTS

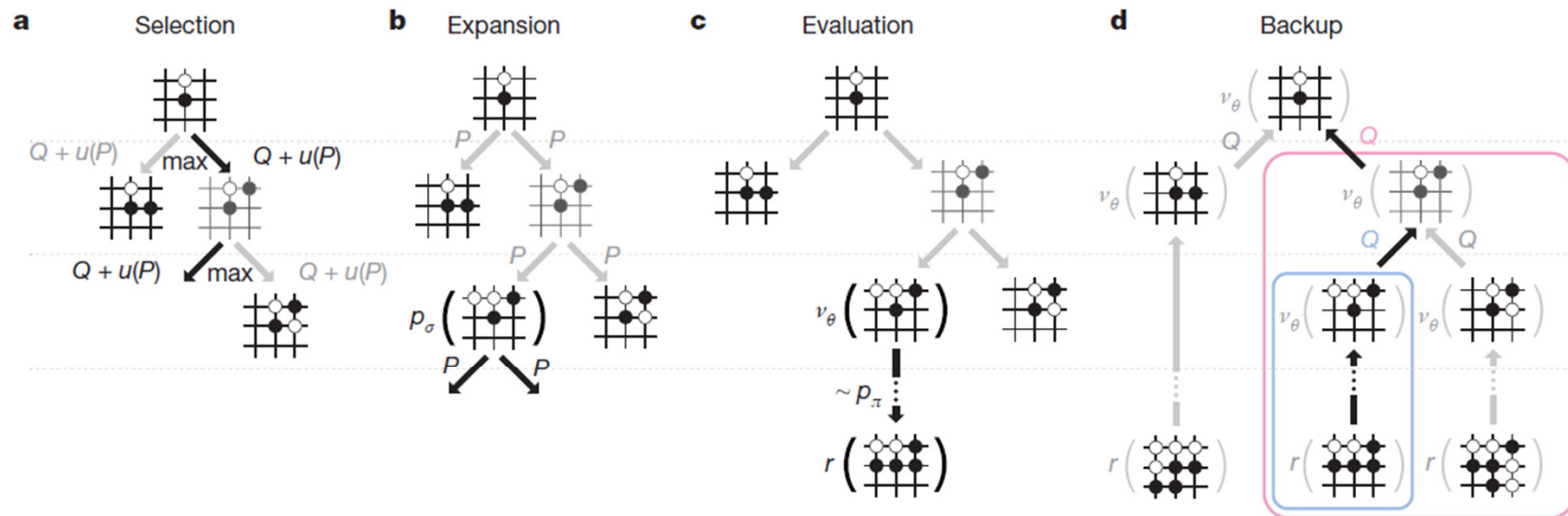


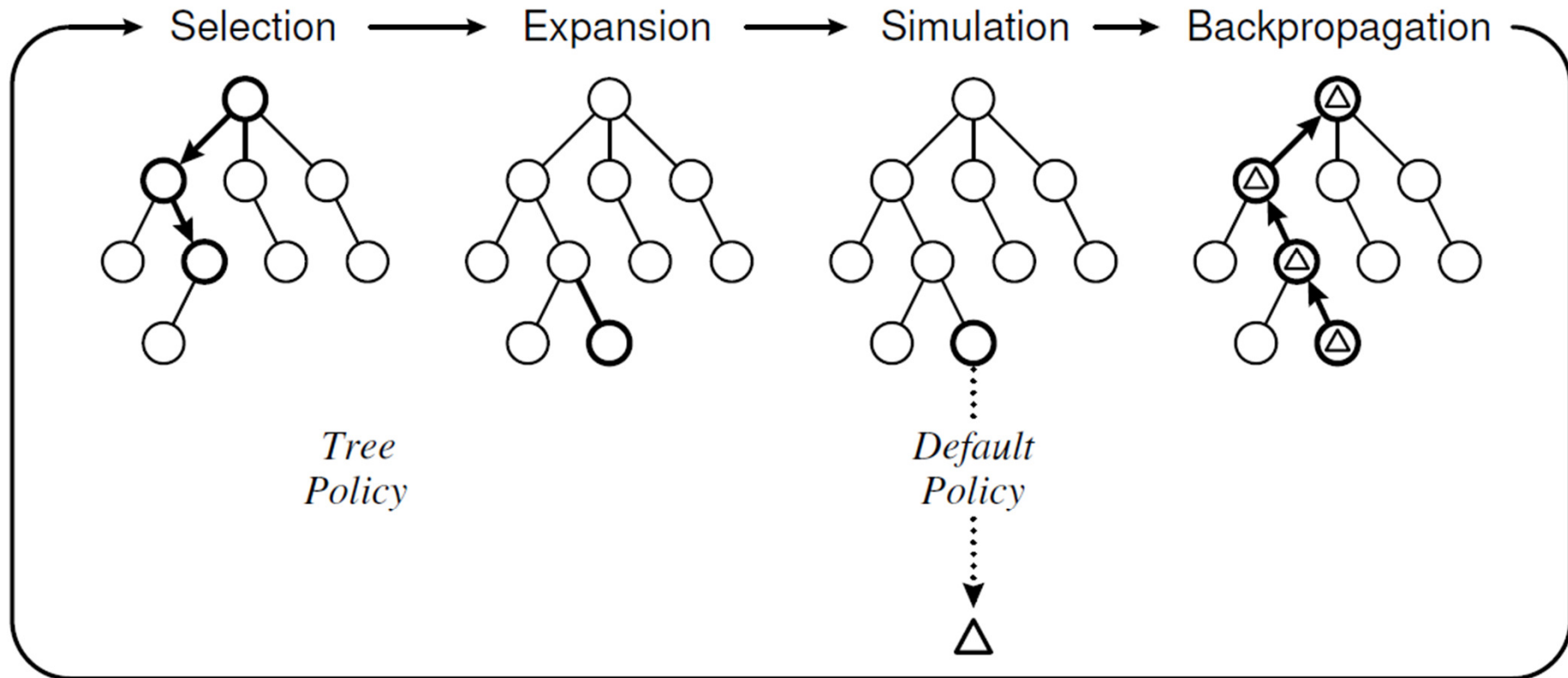
Figure 3 | Monte Carlo tree search in AlphaGo. **a**, Each simulation traverses the tree by selecting the edge with maximum action value  $Q$ , plus a bonus  $u(P)$  that depends on a stored prior probability  $P$  for that edge. **b**, The leaf node may be expanded; the new node is processed once by the policy network  $p_\sigma$  and the output probabilities are stored as prior probabilities  $P$  for each action. **c**, At the end of a simulation, the leaf node

is evaluated in two ways: using the value network  $v_\theta$ ; and by running a rollout to the end of the game with the fast rollout policy  $p_\pi$ , then computing the winner with function  $r$ . **d**, Action values  $Q$  are updated to track the mean value of all evaluations  $r(\cdot)$  and  $v_\theta(\cdot)$  in the subtree below that action.

From AlphaGo's 28Jan16 nature paper @ (<http://www.nature.com/nature/journal/v529/n7587/full/nature16961.html>)

# How AlphaGo Works

———— Nothing new in MCTS, see e.g.



One iteration of the general MCTS approach.

# How AlphaGo Works

————— Search with policy and value networks

- The RL policy network and value network are combined in a MCST to select actions by look-ahead search.
- The tree is traversed by simulation. An edge of the MCST contains an state-action value  $Q(s, a)$ . At each time step  $t$  of each simulation, an action is selected by

$$a_t = \underset{a}{\operatorname{argmax}}(Q(s_t, a) + u(s_t, a))$$

- Here  $u$  is a UCB value, usually in the form of  $\sqrt{\frac{2 \ln n}{n_j}}$ , where  $n_j$  is the number of time arm  $j$  is played and  $n$  is the total number of plays so far

- Value of the leaf node,  $V(s_L)$ , is a combination of two sources,

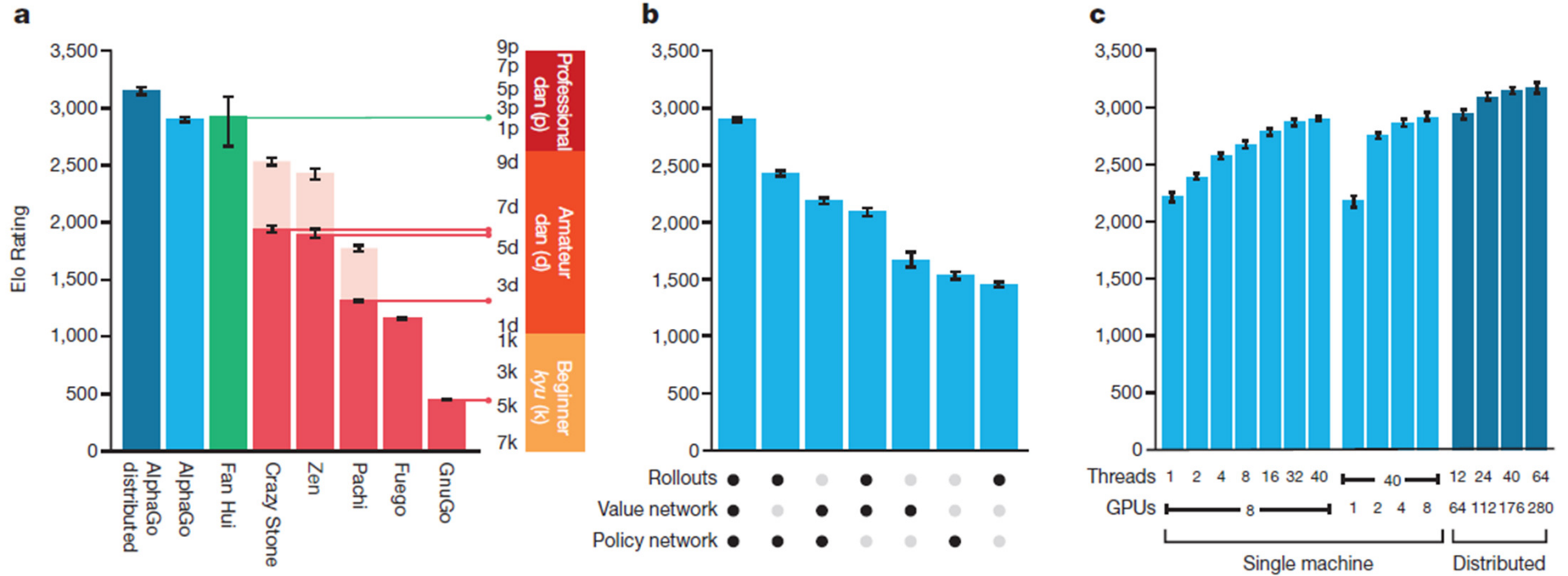
$$V(s_L) = (1 - \lambda)v_{\theta}(s_L) + \lambda z_L$$

Where  $v_{\theta}(s_L)$  is from the value network, and  $z_L$  is based on the fast roll-outs of policy  $P_{\pi}$ .



# AlphaGo

— evaluations



# AlphaGo

—— hardware

## Configuration and performance

Configuration ▼	Search threads ◆	No. of CPU ◆	No. of GPU ◆	Elo rating ◆
Single <sup>[6]</sup> p.10-11	40	48	1	2,151
Single	40	48	2	2,738
Single	40	48	4	2,850
Single	40	48	8	2,890
Distributed	12	428	64	2,937
Distributed	24	764	112	3,079
Distributed	40	1,202	176	3,140
Distributed	64	1,920	280	3,168

?Against H.  
FAN Oct15



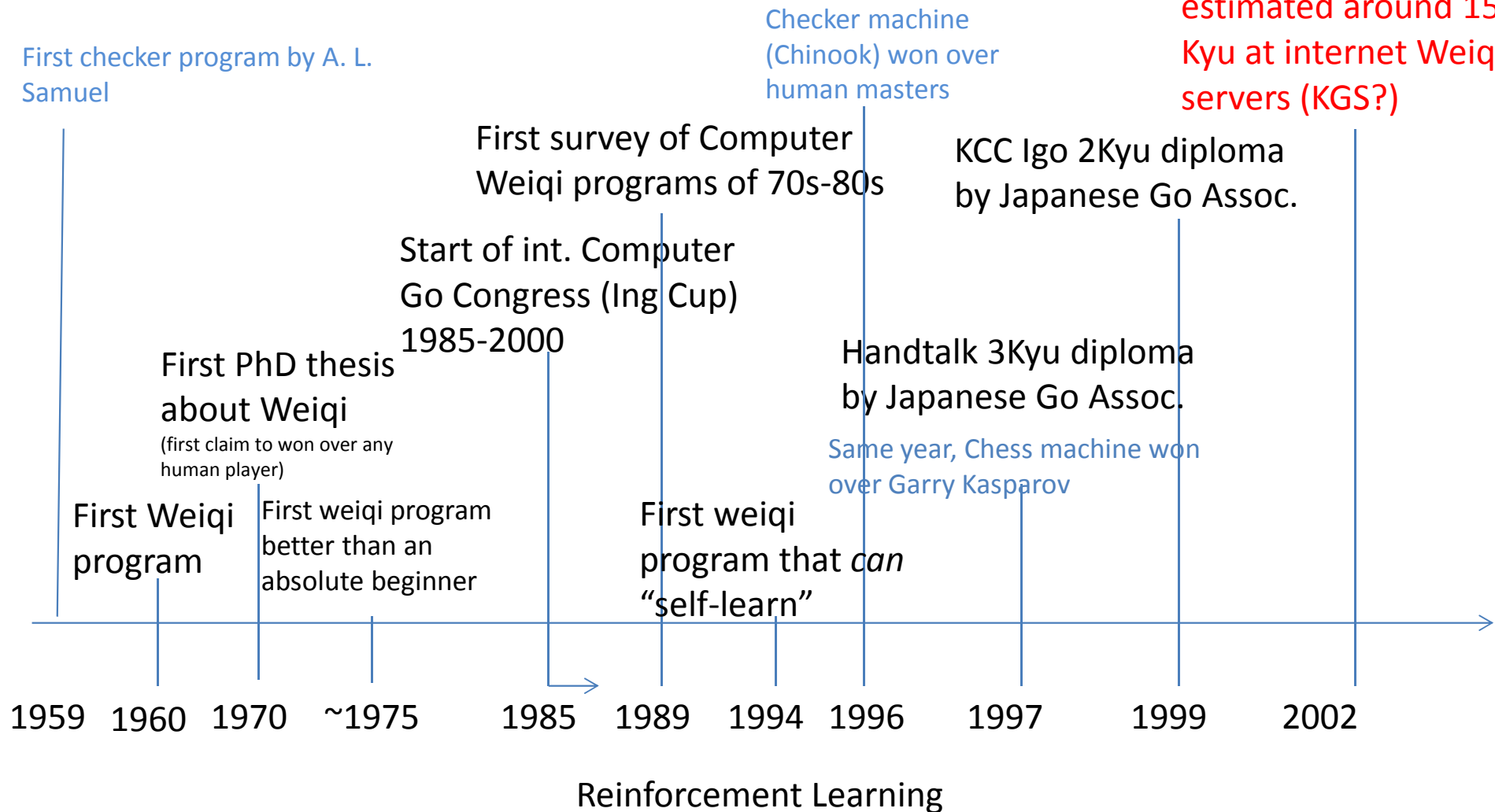
?Against S.  
Lee Mar16



From wikipedia, (<https://en.wikipedia.org/wiki/AlphaGo>)

# A bit History of Weiqi Machines

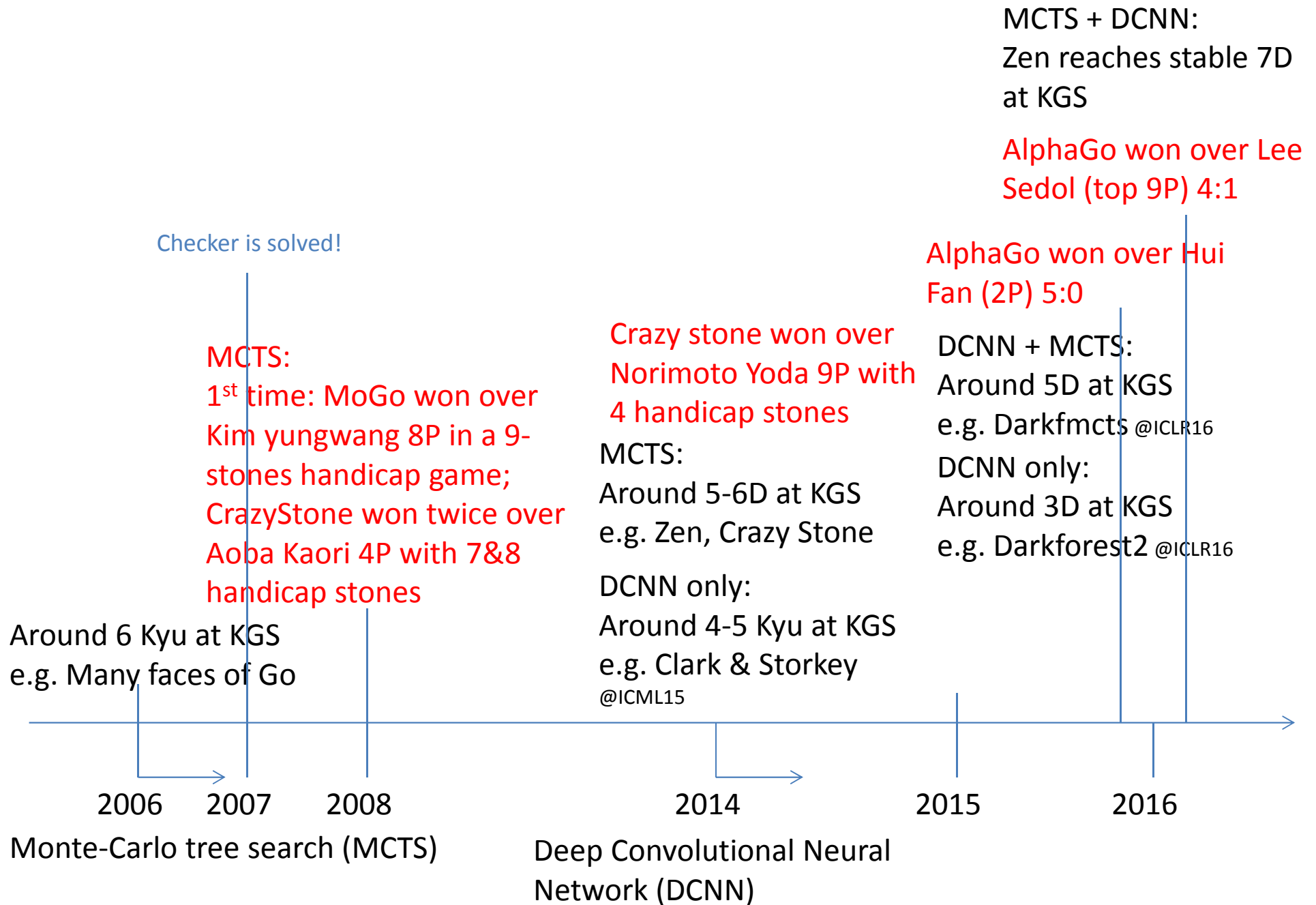
Best program:  
estimated around 15  
Kyu at internet Weiqi  
servers (KGS?)



# A bit **History of Weiqi Machines**

The research of Go programs is still in its infancy, but we shall see that to bring Go programs to a level comparable with current Chess programs, investigations of a totally different kind than used in computer chess are needed.

-- John McCarthy 1990





# Future Look of Weiqi Machines

— three versions of “Solved” Problems

- Ultra-weakly solved: From start of game, know the optimal outcome (of black)
- Weakly solved: From start of game, best play (of black) to end regardless of the opponent’s play
- Strongly solved: From any state of a game, best play (of black) to end

# How AlphaGo Works

————— Here search with policy and value networks

- The RL policy network and value network are combined in a MCST to select actions by look-ahead search.
- The tree is traversed by simulation. An edge of the MCST contains an state-action value  $Q(s, a)$ . At each time step  $t$  of each simulation, an action is selected by

$$a_t = \operatorname{argmax}_a (Q(s_t, a) + u(s_t, a))$$

- Here  $u$  is a bonus value  $u(s, a) \propto \frac{P(s, a)}{1 + N(s, a)}$ ,  $N(s, a)$  is visit count,

and  $P(s, a) = p_\sigma(a|s)$  is the prior prob for an action  $a$ . we also have

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

- Value of the leaf node,  $V(s_L)$ , is a combination of the two sources,

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

Where  $v_\theta(s_L)$  is from the value network, and  $z_L$  is based on the fast roll-outs of policy  $P_\pi$ .