

# AlphaZero

Mastering Chess and Shogi  
by Self-Play

with a General Reinforcement Learning Algorithm

---

Karel Ha

article by Google DeepMind

AI Seminar, 19<sup>th</sup> December 2017



The Alpha\* Timeline

AlphaGo

AlphaGo Zero (AG0)

AlphaZero

Conclusion

## The Alpha\* Timeline

---







- professional 2 dan



- professional 2 dan
- European Go Champion in 2013, 2014 and 2015



- professional 2 dan
- European Go Champion in 2013, 2014 and 2015
- European Professional Go Champion in 2016

# AlphaGo (AlphaGo Fan) vs. Fan Hui

## AlphaGo (AlphaGo Fan) vs. Fan Hui



**AlphaGo won 5:0** in a formal match on October 2015.

# AlphaGo (AlphaGo Fan) vs. Fan Hui



**AlphaGo won 5:0** in a formal match on October 2015.

[AlphaGo] is very strong and stable, it seems like a wall. ... I know AlphaGo is a computer, but if no one told me, maybe I would think the player was a little strange, but a very strong player, a real person.

# Lee Sedol “The Strong Stone”



# Lee Sedol “The Strong Stone”



- professional 9 dan



# Lee Sedol “The Strong Stone”



- professional 9 dan
- the 2<sup>nd</sup> in international titles

# Lee Sedol “The Strong Stone”



- professional 9 dan
- the 2<sup>nd</sup> in international titles
- the 5<sup>th</sup> youngest to become a professional Go player in South Korean history

# Lee Sedol “The Strong Stone”




- professional 9 dan
- the 2<sup>nd</sup> in international titles
- the 5<sup>th</sup> youngest to become a professional Go player in South Korean history
- Lee Sedol would win 97 out of 100 games against Fan Hui.

# Lee Sedol “The Strong Stone”




- professional 9 dan
- the 2<sup>nd</sup> in international titles
- the 5<sup>th</sup> youngest to become a professional Go player in South Korean history
- Lee Sedol would win 97 out of 100 games against Fan Hui.
- “Roger Federer” of Go



I heard Google DeepMind's AI is surprisingly strong and getting stronger, but I am confident that I can win, at least this time.

Lee Sedol



I heard Google DeepMind's AI is surprisingly strong and getting stronger, but I am confident that I can win, at least this time.


---

**Lee Sedol**

...even beating AlphaGo by 4:1 may allow the Google DeepMind team to claim its de facto victory and the defeat of him [Lee Sedol], or even humankind.

---

**interview in JTBC  
Newsroom**



I heard Google DeepMind's AI is surprisingly strong and getting stronger, but I am confident that I can win, at least this time.

---

Lee Sedol

...even beating AlphaGo by 4:1 may allow the Google DeepMind team to claim its de facto victory and the defeat of him [Lee Sedol], or even humankind.

---

interview in JTBC  
Newsroom

# AlphaGo (AlphaGo Lee) vs. Lee Sedol





# AlphaGo (AlphaGo Lee) vs. Lee Sedol



In March 2016 **AlphaGo won 4:1** against the legendary Lee Sedol.

# AlphaGo (AlphaGo Lee) vs. Lee Sedol



In March 2016 **AlphaGo won 4:1** against the legendary Lee Sedol. AlphaGo won all but the 4<sup>th</sup> game; all games were won by resignation.

# AlphaGo (AlphaGo Lee) vs. Lee Sedol




In March 2016 **AlphaGo won 4:1** against the legendary Lee Sedol. AlphaGo won all but the 4<sup>th</sup> game; all games were won by resignation.

# AlphaGo Master



<https://deepmind.com/research/alphago/match-archive/master/>

# AlphaGo Master




In January 2017, DeepMind revealed that AlphaGo had played a series of unofficial online games against some of the strongest professional Go players under the pseudonyms "Master" and "Magister".

<https://deepmind.com/research/alphago/match-archive/master/>



# AlphaGo Master



In January 2017, DeepMind revealed that AlphaGo had played a series of unofficial online games against some of the strongest professional Go players under the pseudonyms "Master" and "Magister".

This AlphaGo was an improved version of the AlphaGo that played Lee Sedol in 2016.

# AlphaGo Master




In January 2017, DeepMind revealed that AlphaGo had played a series of unofficial online games against some of the strongest professional Go players under the pseudonyms "Master" and "Magister".

This AlphaGo was an improved version of the AlphaGo that played Lee Sedol in 2016.

Over one week, AlphaGo played 60 online fast time-control games.

# AlphaGo Master



In January 2017, DeepMind revealed that AlphaGo had played a series of unofficial online games against some of the strongest professional Go players under the pseudonyms "Master" and "Magister".

This AlphaGo was an improved version of the AlphaGo that played Lee Sedol in 2016.

Over one week, AlphaGo played 60 online fast time-control games.

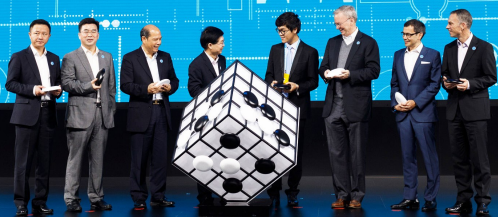
**AlphaGo won this series of games 60:0.**



中国围棋协会 Google 浙江省体育局

# 中国乌镇 围棋峰会 顶尖棋手 + DeepMind AlphaGo 共创棋妙未来

The Future of Go Summit in Wuzhen Legendary players and DeepMind's AlphaGo explore the mysteries of Go



# 中国乌镇 围棋峰会

The Future of Go Summit in Wuzhen



1:1 match vs Ke Jie



Team Go



Pair Go

# 中国乌镇 围棋峰会

The Future of Go Summit in Wuzhen



■ 23 May - 27 May 2017 in Wuzhen, China

# 中国乌镇 围棋峰会

The Future of Go Summit in Wuzhen



- 23 May - 27 May 2017 in Wuzhen, China
- Team Go vs. AlphaGo **0:1**

# 中国乌镇 围棋峰会

## The Future of Go Summit in Wuzhen



- 23 May - 27 May 2017 in Wuzhen, China
- Team Go vs. AlphaGo **0:1**
- AlphaGo vs. world champion Ke Jie **3:0**


# AlphaGo Zero

Starting from scratch

# AlphaGo Zero

Starting from scratch

defeated AlphaGo Lee by **100 games to 0**


A chessboard with various pieces including a king, pawns, and shogi stones. The board is illuminated with a warm, golden light, creating long shadows. The pieces are scattered across the board, with a white king in the center foreground and several pawns and shogi stones around it.

AI system that mastered  
chess, Shogi and Go to  
“superhuman levels” within  
a handful of hours

---

# AlphaZero





AI system that mastered  
chess, Shogi and Go to  
“superhuman levels” within  
a handful of hours

---

# AlphaZero

defeated AlphaGo Zero (version with 20 blocks trained for 3 days)  
by **60 games to 40**



# ① AlphaGo Fan



- 1 AlphaGo Fan
- 2 AlphaGo Lee



- 1 AlphaGo Fan
- 2 AlphaGo Lee
- 3 AlphaGo Master



- 1 AlphaGo Fan
- 2 AlphaGo Lee
- 3 AlphaGo Master
- 4 AlphaGo Zero



- 1 AlphaGo Fan
- 2 AlphaGo Lee
- 3 AlphaGo Master
- 4 AlphaGo Zero
- 5 AlphaZero



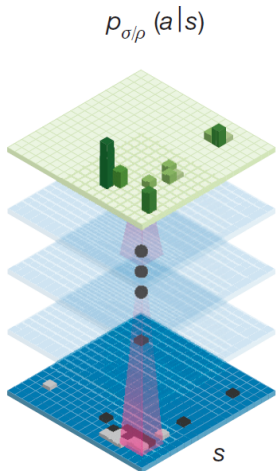
**AlphaGo**

---

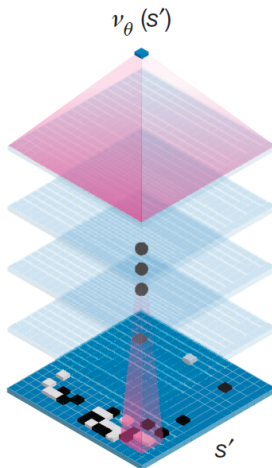


# Policy and Value Networks

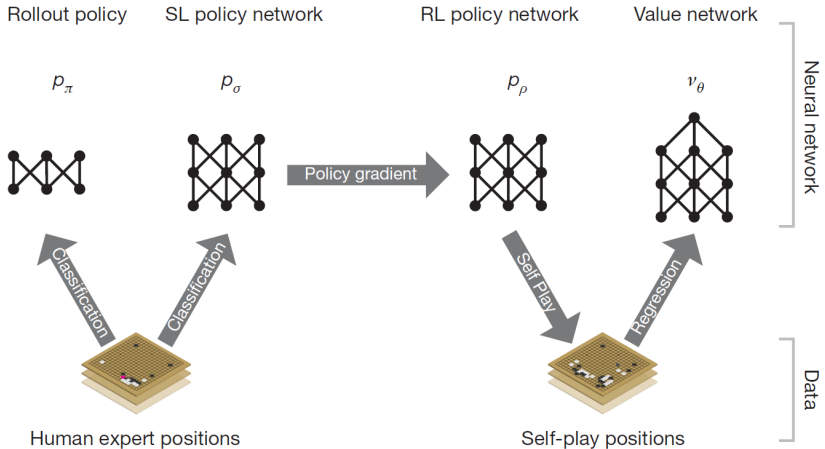
Policy network



Value network



# Training the (Deep Convolutional) Neural Networks



# AlphaGo Zero (AG0)

---

# AG0: Differences Compared to AlphaGo {Fan, Lee, Master}

AlphaGo {Fan, Lee, Master} × AlphaGo Zero:

## AG0: Differences Compared to AlphaGo {Fan, Lee, Master}

AlphaGo {Fan, Lee, Master} × AlphaGo Zero:

- supervised learning from human expert positions × from scratch by self-play reinforcement learning (“tabula rasa”)

# AG0: Differences Compared to AlphaGo {Fan, Lee, Master}

AlphaGo {Fan, Lee, Master} × AlphaGo Zero:

- supervised learning from human expert positions × from scratch by self-play reinforcement learning (“tabula rasa”)
- additional (auxiliary) input features × only the black and white stones from the board as input features

## AG0: Differences Compared to AlphaGo {Fan, Lee, Master}

AlphaGo {Fan, Lee, Master} × AlphaGo Zero:

- supervised learning from human expert positions × from scratch by self-play reinforcement learning (“tabula rasa”)
- additional (auxiliary) input features × only the black and white stones from the board as input features
- separate policy and value networks × single neural network

# AG0: Differences Compared to AlphaGo {Fan, Lee, Master}

AlphaGo {Fan, Lee, Master} × AlphaGo Zero:

- supervised learning from human expert positions × from scratch by self-play reinforcement learning (“tabula rasa”)
- additional (auxiliary) input features × only the black and white stones from the board as input features
- separate policy and value networks × single neural network
- tree search using also Monte Carlo rollouts × simpler tree search using only the single neural network to both evaluate positions and sample moves



# AG0: Differences Compared to AlphaGo {Fan, Lee, Master}

AlphaGo {Fan, Lee, Master} × AlphaGo Zero:

- supervised learning from human expert positions × from scratch by self-play reinforcement learning (“tabula rasa”)
- additional (auxiliary) input features × only the black and white stones from the board as input features
- separate policy and value networks × single neural network
- tree search using also Monte Carlo rollouts × simpler tree search using only the single neural network to both evaluate positions and sample moves
- (AlphaGo Lee) distributed machines + 48 tensor processing units (TPUs) × single machines + 4 TPUs

# AG0: Differences Compared to AlphaGo {Fan, Lee, Master}

AlphaGo {Fan, Lee, Master} × AlphaGo Zero:

- supervised learning from human expert positions × from scratch by self-play reinforcement learning (“tabula rasa”)
- additional (auxiliary) input features × only the black and white stones from the board as input features
- separate policy and value networks × single neural network
- tree search using also Monte Carlo rollouts × simpler tree search using only the single neural network to both evaluate positions and sample moves
- (AlphaGo Lee) distributed machines + 48 tensor processing units (TPUs) × single machines + 4 TPUs
- (AlphaGo Lee) several months of training time × 72 h of training time (outperforming AlphaGo Lee after 36 h)

AG0 achieves this via

AG0 achieves this via

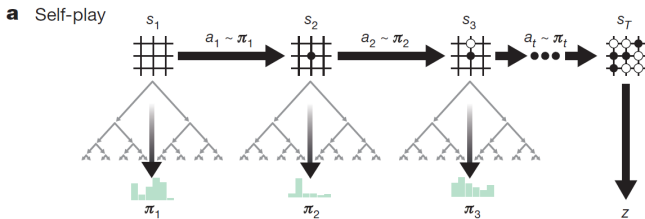
- a new reinforcement learning algorithm

AG0 achieves this via

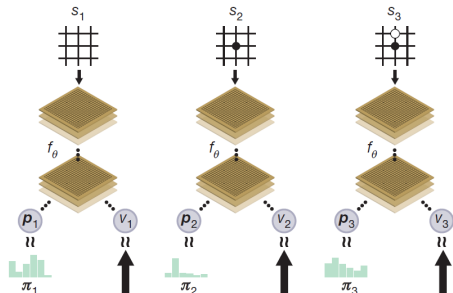
- a new reinforcement learning algorithm
- with lookahead search inside the training loop

# AG0: Self-Play Reinforcement Learning

# AG0: Self-Play Reinforcement Learning



**b** Neural network training



deep neural network  $f_\theta$  with parameters  $\theta$ :



deep neural network  $f_{\theta}$  with parameters  $\theta$ :

- input: raw board representation  $s$

deep neural network  $f_{\theta}$  with parameters  $\theta$ :

- input: raw board representation  $s$
- output:

# AG0: Self-Play Reinforcement Learning – Neural Network

deep neural network  $f_{\theta}$  with parameters  $\theta$ :

- input: raw board representation  $s$
- output:
  - move probabilities  $\mathbf{p}$

deep neural network  $f_{\theta}$  with parameters  $\theta$ :

- input: raw board representation  $s$
- output:
  - move probabilities  $\mathbf{p}$
  - value  $v$  of the board position

deep neural network  $f_\theta$  with parameters  $\theta$ :

- input: raw board representation  $s$
- output:
  - move probabilities  $\mathbf{p}$
  - value  $v$  of the board position
  - $f_\theta(s) = (\mathbf{p}, v)$

deep neural network  $f_\theta$  with parameters  $\theta$ :

- input: raw board representation  $s$
- output:
  - move probabilities  $\mathbf{p}$
  - value  $v$  of the board position
  - $f_\theta(s) = (\mathbf{p}, v)$
- specifics:

deep neural network  $f_\theta$  with parameters  $\theta$ :

- input: raw board representation  $s$
- output:
  - move probabilities  $\mathbf{p}$
  - value  $v$  of the board position
  - $f_\theta(s) = (\mathbf{p}, v)$
- specifics:
  - (20 or 40) residual blocks (of convolutional layers)

deep neural network  $f_\theta$  with parameters  $\theta$ :

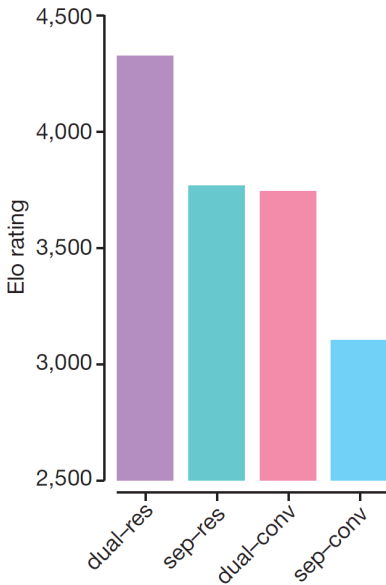
- input: raw board representation  $s$
- output:
  - move probabilities  $\mathbf{p}$
  - value  $v$  of the board position
  - $f_\theta(s) = (\mathbf{p}, v)$
- specifics:
  - (20 or 40) residual blocks (of convolutional layers)
  - **batch normalization**



deep neural network  $f_\theta$  with parameters  $\theta$ :

- input: raw board representation  $s$
- output:
  - move probabilities  $\mathbf{p}$
  - value  $v$  of the board position
  - $f_\theta(s) = (\mathbf{p}, v)$
- specifics:
  - (20 or 40) residual blocks (of convolutional layers)
  - batch normalization
  - rectifier non-linearities

# AG0: Comparison of Various Neural Network Architectures



# AG0: Self-Play Reinforcement Learning – Steps

0. random weights  $\theta_0$

## AG0: Self-Play Reinforcement Learning – Steps

0. random weights  $\theta_0$
1. at each iteration  $i > 0$ , self-play games are generated:

# AG0: Self-Play Reinforcement Learning – Steps

0. random weights  $\theta_0$
1. at each iteration  $i > 0$ , self-play games are generated:
  - i. MCTS samples search probabilities  $\pi_t$  based on the neural network from the previous iteration  $f_{\theta_{i-1}}$ :

$$\pi_t = \alpha_{\theta_{i-1}}(s_t)$$

for each time-step  $t = 1, 2, \dots, T$

# AG0: Self-Play Reinforcement Learning – Steps

0. random weights  $\theta_0$
1. at each iteration  $i > 0$ , self-play games are generated:
  - i. MCTS samples search probabilities  $\pi_t$  based on the neural network from the previous iteration  $f_{\theta_{i-1}}$ :

$$\pi_t = \alpha_{\theta_{i-1}}(s_t)$$

for each time-step  $t = 1, 2, \dots, T$

- ii. move is sampled from  $\pi_t$

# AG0: Self-Play Reinforcement Learning – Steps

0. random weights  $\theta_0$
1. at each iteration  $i > 0$ , self-play games are generated:
  - i. MCTS samples search probabilities  $\pi_t$  based on the neural network from the previous iteration  $f_{\theta_{i-1}}$ :

$$\pi_t = \alpha_{\theta_{i-1}}(s_t)$$

for each time-step  $t = 1, 2, \dots, T$

- ii. move is sampled from  $\pi_t$
- iii. data  $(s_t, \pi_t, z_t)$  for each  $t$  are stored for later training

# AG0: Self-Play Reinforcement Learning – Steps

0. random weights  $\theta_0$
1. at each iteration  $i > 0$ , self-play games are generated:
  - i. MCTS samples search probabilities  $\pi_t$  based on the neural network from the previous iteration  $f_{\theta_{i-1}}$ :

$$\pi_t = \alpha_{\theta_{i-1}}(s_t)$$

for each time-step  $t = 1, 2, \dots, T$

- ii. move is sampled from  $\pi_t$
- iii. data  $(s_t, \pi_t, z_t)$  for each  $t$  are stored for later training
- iv. new neural network  $f_{\theta_i}$  is trained in order to minimize the loss

$$l = (z - v)^2 - \pi^\top \log \mathbf{p} + c \|\theta\|^2$$



## AG0: Self-Play Reinforcement Learning – Steps

0. random weights  $\theta_0$
1. at each iteration  $i > 0$ , self-play games are generated:
  - i. MCTS samples search probabilities  $\pi_t$  based on the neural network from the previous iteration  $f_{\theta_{i-1}}$ :

$$\pi_t = \alpha_{\theta_{i-1}}(s_t)$$

for each time-step  $t = 1, 2, \dots, T$

- ii. move is sampled from  $\pi_t$
- iii. data  $(s_t, \pi_t, z_t)$  for each  $t$  are stored for later training
- iv. new neural network  $f_{\theta_i}$  is trained in order to minimize the loss

$$l = (z - v)^2 - \pi^\top \log \mathbf{p} + c \|\theta\|^2$$

## AG0: Self-Play Reinforcement Learning – Steps

0. random weights  $\theta_0$
1. at each iteration  $i > 0$ , self-play games are generated:
  - i. MCTS samples search probabilities  $\pi_t$  based on the neural network from the previous iteration  $f_{\theta_{i-1}}$ :

$$\pi_t = \alpha_{\theta_{i-1}}(s_t)$$

for each time-step  $t = 1, 2, \dots, T$

- ii. move is sampled from  $\pi_t$
- iii. data  $(s_t, \pi_t, z_t)$  for each  $t$  are stored for later training
- iv. new neural network  $f_{\theta_i}$  is trained in order to minimize the loss

$$l = (z - v)^2 - \pi^\top \log \mathbf{p} + c \|\theta\|^2$$

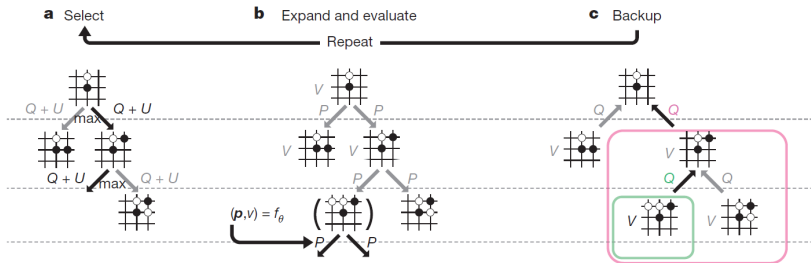
Loss  $l$  makes  $(\mathbf{p}, v) = f_\theta(s)$  more closely match the improved search probabilities and self-play winner  $(\pi, z)$ .

## AG0: Monte Carlo Tree Search (1/2)

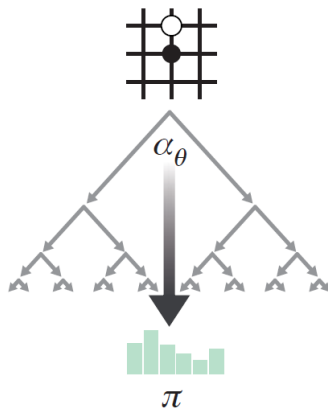
Monte Carlo Tree Search (MCTS) in AG0:

# AG0: Monte Carlo Tree Search (1/2)

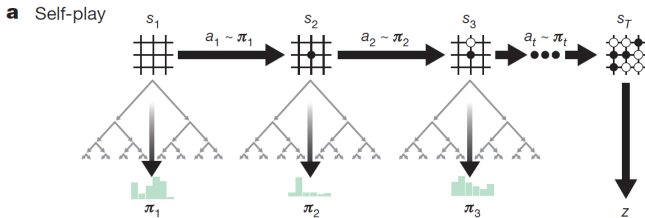
## Monte Carlo Tree Search (MCTS) in AG0:



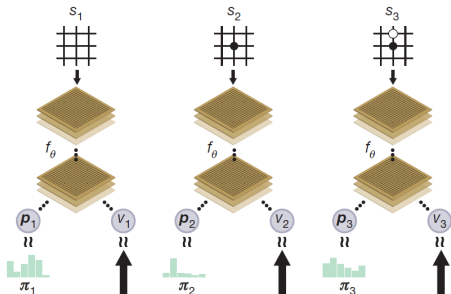
## **d** Play



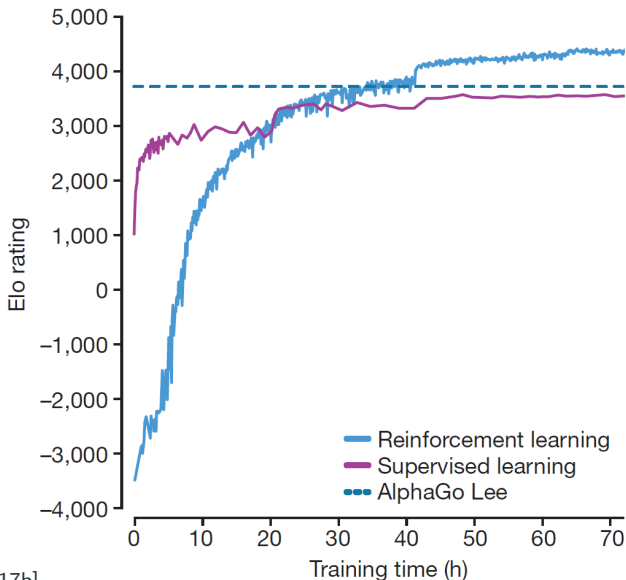
# AG0: Self-Play Reinforcement Learning – Review



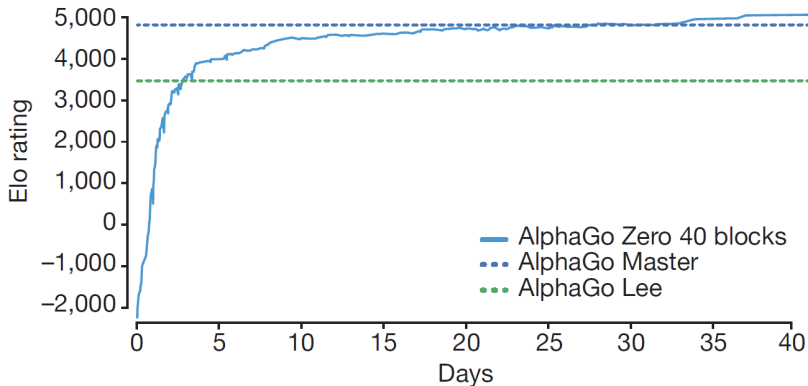
**b** Neural network training



# AG0: Elo Rating over Training Time (RL vs. SL)

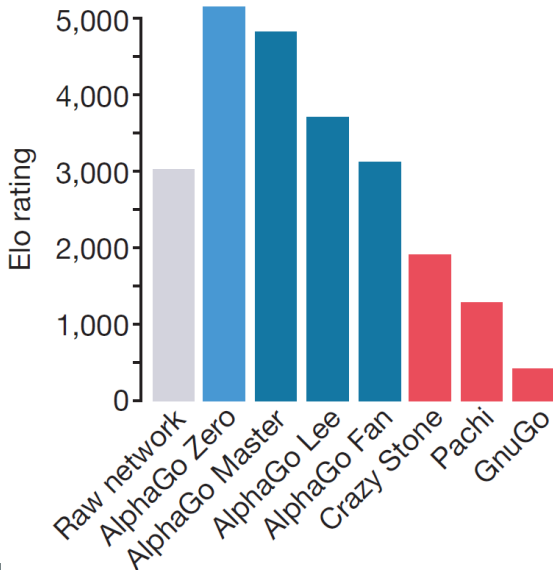


## AG0: Elo Rating over Training Time (AG0 with 40 blocks)

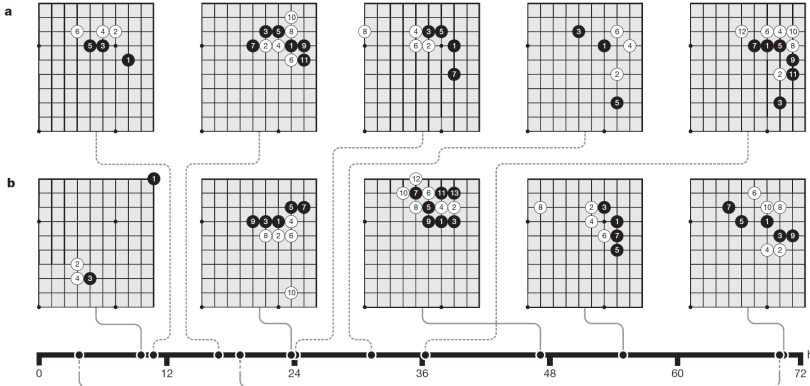




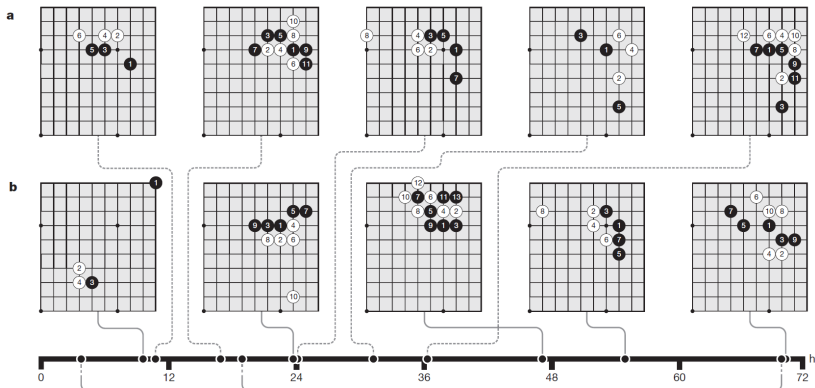
# AG0: Tournament between AI Go Programs



# AG0: Discovered Joseki (Corner Sequences)

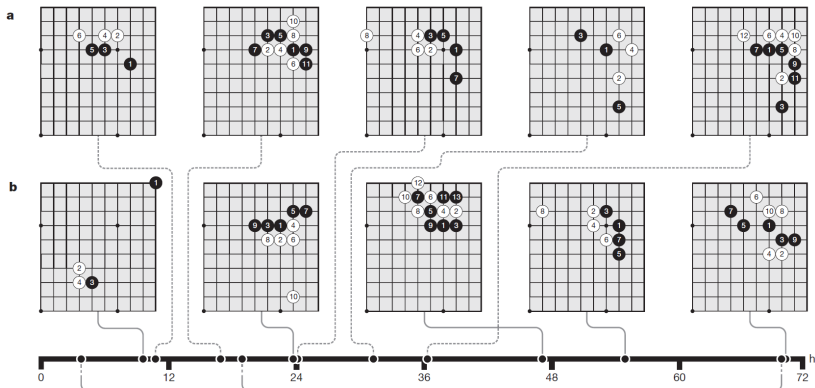


# AG0: Discovered Joseki (Corner Sequences)



a five human *joseki*

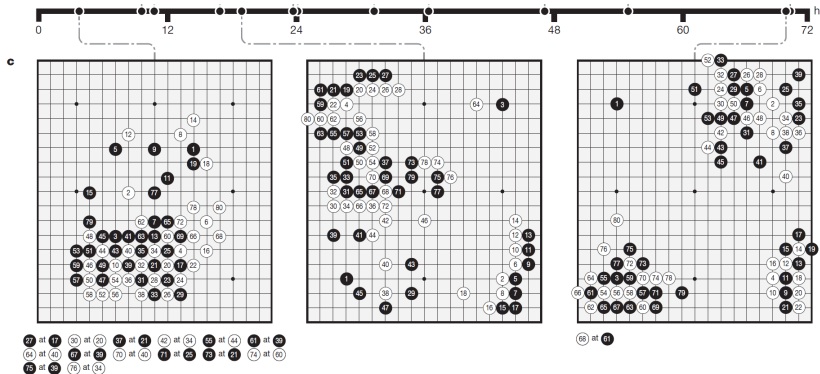
# AG0: Discovered Joseki (Corner Sequences)



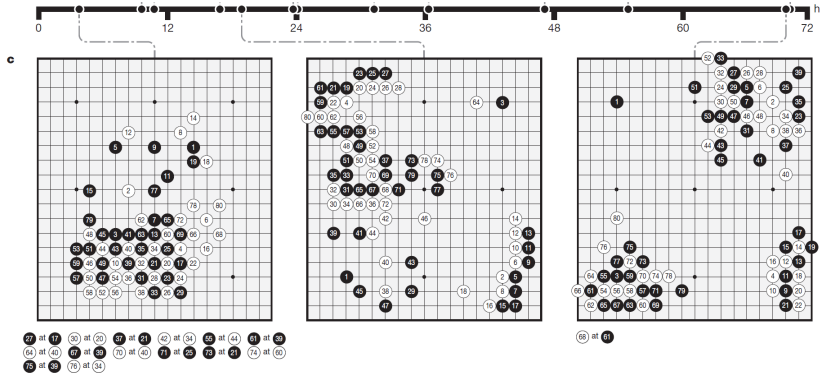
a five human *joseki*

b five novel *joseki* variants eventually preferred by AG0

# AG0: Discovered Playing Styles

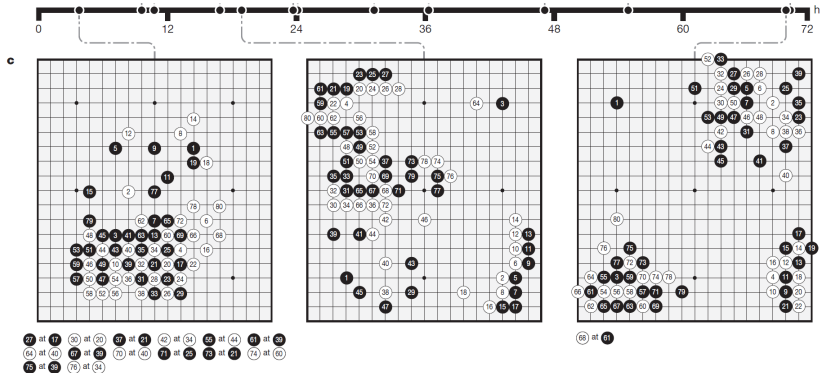


# AG0: Discovered Playing Styles

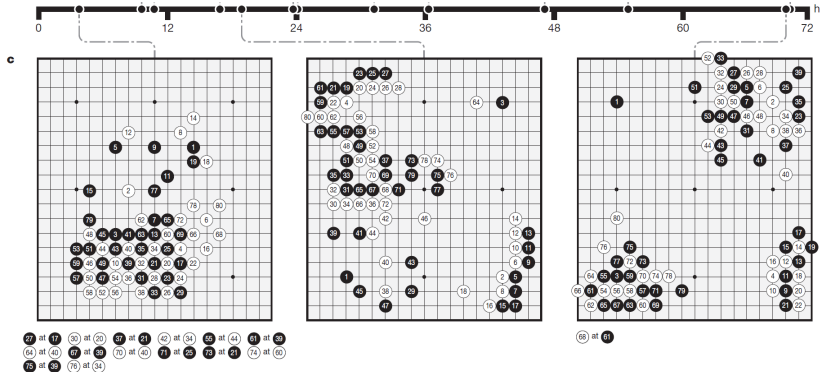


at 3 h greedy capture of stones

# AG0: Discovered Playing Styles



# AG0: Discovered Playing Styles



at 3 h greedy capture of stones

at 19 h the fundamentals of Go concepts (life-and-death, influence, territory...)

at 70 h remarkably balanced game (multiple battles, complicated *ko* fight, a half-point win for white...)



# AlphaZero

---





To watch such a strong programme like Stockfish, against whom most top players would be happy to win even one game out of a hundred, being completely taken apart is certainly definitive.

---

**Viswanathan Anand**



To watch such a strong programme like Stockfish, against whom most top players would be happy to win even one game out of a hundred, being completely taken apart is certainly definitive.

---

**Viswanathan Anand**

It's like chess from another dimension.

---

**Demis Hassabis**

# AlphaZero: Differences Compared to AlphaGo Zero

AlphaGo Zero × AlphaZero:

# AlphaZero: Differences Compared to AlphaGo Zero

AlphaGo Zero  $\times$  AlphaZero:

- binary outcome (win / loss)  $\times$  expected outcome (including draws or potentially other outcomes)

# AlphaZero: Differences Compared to AlphaGo Zero

AlphaGo Zero × AlphaZero:

- binary outcome (win / loss) × expected outcome (including draws or potentially other outcomes)
- board positions transformed before passing to neural networks (by randomly selected rotation or reflection) × no data augmentation

# AlphaZero: Differences Compared to AlphaGo Zero

AlphaGo Zero × AlphaZero:

- binary outcome (win / loss) × expected outcome (including draws or potentially other outcomes)
- board positions transformed before passing to neural networks (by randomly selected rotation or reflection) × no data augmentation
- games generated by the best player from previous iterations (margin of 55 %) × continual update using the latest parameters (without the evaluation and selection steps)



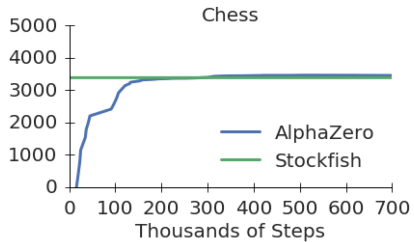
# AlphaZero: Differences Compared to AlphaGo Zero

AlphaGo Zero × AlphaZero:

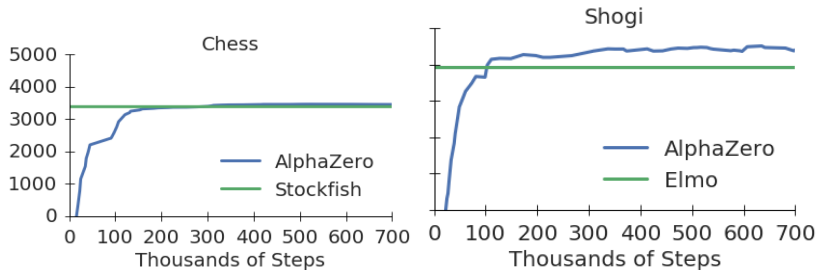
- binary outcome (win / loss) × expected outcome (including draws or potentially other outcomes)
- board positions transformed before passing to neural networks (by randomly selected rotation or reflection) × no data augmentation
- games generated by the best player from previous iterations (margin of 55 %) × continual update using the latest parameters (without the evaluation and selection steps)
- hyper-parameters tuned by Bayesian optimisation × reused the same hyper-parameters without game-specific tuning

# AlphaZero: Elo Rating over Training Time

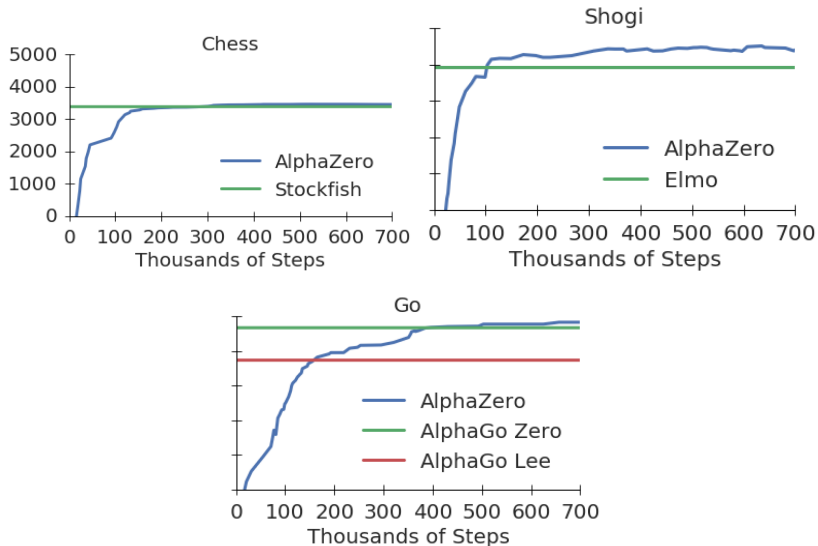
# AlphaZero: Elo Rating over Training Time



# AlphaZero: Elo Rating over Training Time



# AlphaZero: Elo Rating over Training Time



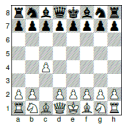
# AlphaZero: Tournament between AI Programs

Game	White	Black	Win	Draw	Loss
Chess	<i>AlphaZero</i>	<i>Stockfish</i>	25	25	0
	<i>Stockfish</i>	<i>AlphaZero</i>	3	47	0
Shogi	<i>AlphaZero</i>	<i>Elmo</i>	43	2	5
	<i>Elmo</i>	<i>AlphaZero</i>	47	0	3
Go	<i>AlphaZero</i>	<i>AG0 3-day</i>	31	–	19
	<i>AG0 3-day</i>	<i>AlphaZero</i>	29	–	21

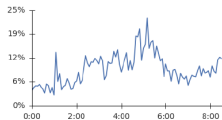
(Values are given from AlphaZero's point of view.)

# AlphaZero: Openings Discovered by the Self-Play (1/2)

A10: English Opening

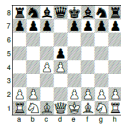


w 20/30/0, b 8/40/2

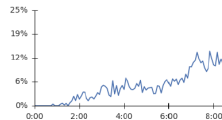


1...e5 g3 d5 cxd5 ♘f6 ♙g2 ♚xd5 ♜f3

D06: Queens Gambit

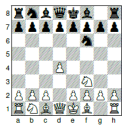


w 16/34/0, b 1/47/2

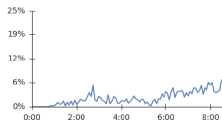


2...c6 ♘c3 ♚f6 ♚f3 a6 g3 c4 a4

A46: Queens Pawn Game

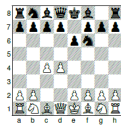


w 24/26/0, b 3/47/0

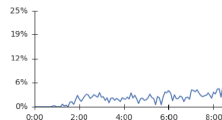


2...d5 c4 e6 ♘c3 ♙e7 ♙f4 O-O e3

E00: Queens Pawn Game



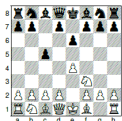
w 17/33/0, b 5/44/1



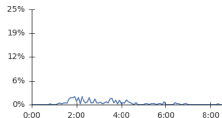
3.♘f3 d5 ♘c3 ♙b4 ♙g5 h6 ♚a4 ♜c6

# AlphaZero: Openings Discovered by the Self-Play (2/2)

B40: Sicilian Defence

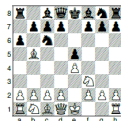


w 17/31/2, b 3/40/7

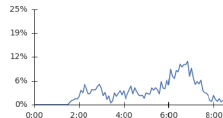


3.d4 cxd4 ♘xd4 ♘c6 ♘c3 ♖c7 ♙e3 a6

C60: Ruy Lopez (Spanish Opening)

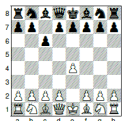


w 27/22/1, b 6/44/0

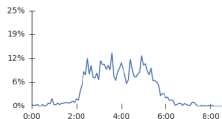


4.♙a4 ♙e7 O-O ♘f6 ♚e1 b5 ♙b3 O-O

B10: Caro-Kann Defence

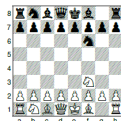


w 25/25/0, b 4/45/1

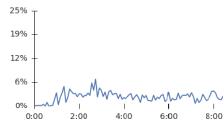


2.d4 d5 e5 ♙f5 ♘f3 e6 ♙e2 a6

A05: Reti Opening



w 13/36/1, b 7/43/0



2.c4 e6 d4 d5 ♘c3 ♙e7 ♙f4 O-O



## Conclusion

---

- challenging decision-making

# Difficulties of Go

- challenging decision-making
- intractable search space

# Difficulties of Go

- challenging decision-making
- intractable search space
- complex optimal solution

It appears infeasible to directly approximate using a policy or value function!

- Monte Carlo tree search

# AlphaZero: Summary

- Monte Carlo tree search
- effective move selection and position evaluation

# AlphaZero: Summary

- Monte Carlo tree search
- effective move selection and position evaluation
  - through deep convolutional neural networks

# AlphaZero: Summary

- Monte Carlo tree search
- effective move selection and position evaluation
  - through deep convolutional neural networks
  - trained by new self-play reinforcement learning algorithm



# AlphaZero: Summary

- Monte Carlo tree search
- effective move selection and position evaluation
  - through deep convolutional neural networks
  - trained by new self-play reinforcement learning algorithm
- new search algorithm combining

# AlphaZero: Summary

- Monte Carlo tree search
- effective move selection and position evaluation
  - through deep convolutional neural networks
  - trained by new self-play reinforcement learning algorithm
- new search algorithm combining
  - evaluation by a single neural network

# AlphaZero: Summary

- Monte Carlo tree search
- effective move selection and position evaluation
  - through deep convolutional neural networks
  - trained by new self-play reinforcement learning algorithm
- new search algorithm combining
  - evaluation by a single neural network
  - Monte Carlo tree search

# AlphaZero: Summary

- Monte Carlo tree search
- effective move selection and position evaluation
  - through deep convolutional neural networks
  - trained by new self-play reinforcement learning algorithm
- new search algorithm combining
  - evaluation by a single neural network
  - Monte Carlo tree search
- more efficient when compared to previous AlphaGo versions

# AlphaZero: Summary

- Monte Carlo tree search
- effective move selection and position evaluation
  - through deep convolutional neural networks
  - trained by new self-play reinforcement learning algorithm
- new search algorithm combining
  - evaluation by a single neural network
  - Monte Carlo tree search
- more efficient when compared to previous AlphaGo versions
  - single machine

# AlphaZero: Summary

- Monte Carlo tree search
- effective move selection and position evaluation
  - through deep convolutional neural networks
  - trained by new self-play reinforcement learning algorithm
- new search algorithm combining
  - evaluation by a single neural network
  - Monte Carlo tree search
- more efficient when compared to previous AlphaGo versions
  - single machine
  - 4 TPUs

# AlphaZero: Summary

- Monte Carlo tree search
- effective move selection and position evaluation
  - through deep convolutional neural networks
  - trained by new self-play reinforcement learning algorithm
- new search algorithm combining
  - evaluation by a single neural network
  - Monte Carlo tree search
- more efficient when compared to previous AlphaGo versions
  - single machine
  - 4 TPUs
  - hours rather than months of training time

# Novel approach



## Novel approach

During the matches (against Stockfish and Elmo), AlphaZero evaluated **thousands of times fewer** positions than Deep Blue against Kasparov.

## Novel approach

During the matches (against Stockfish and Elmo), AlphaZero evaluated **thousands of times fewer** positions than Deep Blue against Kasparov.

It compensated this by:

- selecting those positions **more intelligently** (the neural network)

## Novel approach

During the matches (against Stockfish and Elmo), AlphaZero evaluated **thousands of times fewer** positions than Deep Blue against Kasparov.

It compensated this by:

- selecting those positions **more intelligently** (the neural network)
- evaluating them **more precisely** (the same neural network)

## Novel approach

During the matches (against Stockfish and Elmo), AlphaZero evaluated **thousands of times fewer** positions than Deep Blue against Kasparov.

It compensated this by:

- selecting those positions **more intelligently** (the neural network)
- evaluating them **more precisely** (the same neural network)

## Novel approach

During the matches (against Stockfish and Elmo), AlphaZero evaluated **thousands of times fewer** positions than Deep Blue against Kasparov.

It compensated this by:

- selecting those positions **more intelligently** (the neural network)
- evaluating them **more precisely** (the same neural network)

Deep Blue relied on a handcrafted evaluation function.

## Novel approach

During the matches (against Stockfish and Elmo), AlphaZero evaluated **thousands of times fewer** positions than Deep Blue against Kasparov.

It compensated this by:

- selecting those positions **more intelligently** (the neural network)
- evaluating them **more precisely** (the same neural network)

Deep Blue relied on a handcrafted evaluation function.

AlphaZero was trained **tabula rasa** from self-play. It used **general-purpose** learning.

## Novel approach

During the matches (against Stockfish and Elmo), AlphaZero evaluated **thousands of times fewer** positions than Deep Blue against Kasparov.

It compensated this by:

- selecting those positions **more intelligently** (the neural network)
- evaluating them **more precisely** (the same neural network)

Deep Blue relied on a handcrafted evaluation function.

AlphaZero was trained **tabula rasa** from self-play. It used **general-purpose** learning.

This approach is not specific to the game of Go. The algorithm can be used **for much wider class** of AI problems!

**Thank you!**

**Questions?**



**Backup Slides**

# Input Features of AlphaZero's Neural Networks

Go		Chess		Shogi	
Feature	Planes	Feature	Planes	Feature	Planes
P1 stone	1	P1 piece	6	P1 piece	14
P2 stone	1	P2 piece	6	P2 piece	14
		Repetitions	2	Repetitions	3
				P1 prisoner count	7
				P2 prisoner count	7
Colour	1	Colour	1	Colour	1
		Total move count	1	Total move count	1
		P1 castling	2		
		P2 castling	2		
		No-progress count	1		
Total	17	Total	119	Total	362

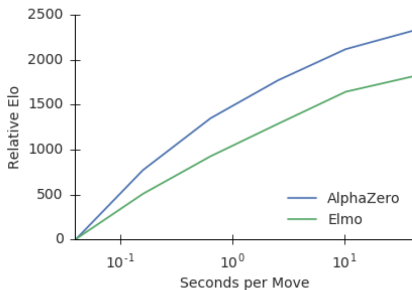
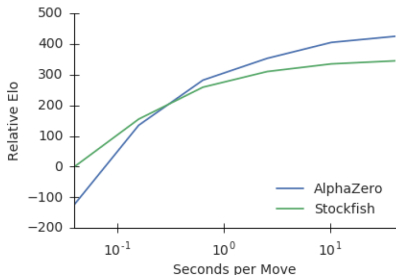
# AlphaZero: Statistics of Training

	Chess	Shogi	Go
Mini-batches	700k	700k	700k
Training Time	9h	12h	34h
Training Games	44 million	24 million	21 million
Thinking Time	800 sims 40 ms	800 sims 80 ms	800 sims 200 ms

# AlphaZero: Evaluation Speeds

Program	Chess	Shogi	Go
<i>AlphaZero</i>	80k	40k	16k
<i>Stockfish</i>	70,000k		
<i>Elmo</i>		35,000k	

# Scalability When Compared to Other Programs



[Silver et al. 2017a]

# Further Reading I

## AlphaGo:

- **Google Research Blog**

<http://googleresearch.blogspot.cz/2016/01/alphago-mastering-ancient-game-of-go.html>

- an article in **Nature**

<http://www.nature.com/news/google-ai-algorithm-masters-ancient-game-of-go-1.19234>

- a **reddit** article claiming that AlphaGo is even stronger than it appears to be:

"AlphaGo would rather win by less points, but with higher probability."

[https://www.reddit.com/r/baduk/comments/49y17z/the\\_true\\_strength\\_of\\_alphago/](https://www.reddit.com/r/baduk/comments/49y17z/the_true_strength_of_alphago/)

- a video of how AlphaGo works (put in layman's terms) <https://youtu.be/qWcfiPi9gUU>

## Articles by Google DeepMind:

- **Atari player:** a DeepRL system which combines Deep Neural Networks with Reinforcement Learning (Mnih et al. 2015)

- **Neural Turing Machines** (Graves, Wayne, and Danihelka 2014)

## Artificial Intelligence:

- **Artificial Intelligence course at MIT**

<http://ocw.mit.edu/courses/electrical-engineering-and-computer-science/6-034-artificial-intelligence-fall-2010/index.htm>

## Further Reading II

- **Introduction to Artificial Intelligence at Udacity**  
<https://www.udacity.com/course/intro-to-artificial-intelligence--cs271>
- **General Game Playing course** <https://www.coursera.org/course/ggp>
- **Singularity** <http://waitbutwhy.com/2015/01/artificial-intelligence-revolution-1.html> + Part 2
- **The Singularity Is Near** (Kurzweil 2005)

Combinatorial Game Theory (founded by John H. Conway to study endgames in Go):

- **Combinatorial Game Theory course** <https://www.coursera.org/learn/combinatorial-game-theory>
- **On Numbers and Games** (Conway 1976)
- **Computer Go as a sum of local games: an application of combinatorial game theory** (Müller 1995)

Chess:

- **Deep Blue beats G. Kasparov in 1997** <https://youtu.be/NJarxpYyoFI>

Machine Learning:

- **Machine Learning course**  
<https://youtu.be/hPKJBXkyTK://www.coursera.org/learn/machine-learning/>
- **Reinforcement Learning** <http://reinforcementlearning.ai-depot.com/>
- **Deep Learning** (LeCun, Bengio, and Hinton 2015)

# Further Reading III

- **Deep Learning course** <https://www.udacity.com/course/deep-learning--ud730>
- **Two Minute Papers** <https://www.youtube.com/user/keeroz>
- **Applications of Deep Learning** <https://youtu.be/hPKJBXkyTKM>

Neuroscience:

- <http://www.brainfacts.org/>



# References I



Conway, John Horton (1976). "On Numbers and Games". In: *London Mathematical Society Monographs* 6.



Graves, Alex, Greg Wayne, and Ivo Danihelka (2014). "Neural Turing Machines". In: *arXiv preprint arXiv:1410.5401*.



Kurzweil, Ray (2005). *The Singularity is Near: When Humans Transcend Biology*. Penguin.



LeCun, Yann, Yoshua Bengio, and Geoffrey Hinton (2015). "Deep Learning". In: *Nature* 521.7553, pp. 436–444.



Mnih, Volodymyr et al. (2015). "Human-Level Control through Deep Reinforcement Learning". In: *Nature* 518.7540, pp. 529–533. URL:  
<https://storage.googleapis.com/deepmind-data/assets/papers/DeepMindNature14236Paper.pdf>.



Müller, Martin (1995). "Computer Go as a Sum of Local Games: an Application of Combinatorial Game Theory". PhD thesis. TU Graz.



Silver, David et al. (2016). "Mastering the Game of Go with Deep Neural Networks and Tree Search". In: *Nature* 529.7587, pp. 484–489.



Silver, David et al. (2017a). "Mastering Chess and Shogi by Self-Play with a General Reinforcement Learning Algorithm". In: *arXiv preprint arXiv:1712.01815*.



Silver, David et al. (2017b). "Mastering the Game of Go without Human Knowledge". In: *Nature* 550.7676, pp. 354–359.