



# AlphaGo

How it works

PPT provider: Shane (Seungwhan) Moon

PhD student, Carnegie Mellon University

# AlphaGo vs European Champion (Fan Hui 2-Dan<sup>\*</sup>)<sup>rank</sup>



**October 5 – 9, 2015**

**<Official match>**

- Time limit: 1 hour
- AlphaGo Wins (5:0)

# AlphaGo vs World Champion (Lee Sedol 9–Dan)



**March 9 – 15, 2016**

**<Official match>**

- Time limit: 2 hours
- reward: 1 million USD

**Venue: Seoul, Four Seasons Hotel**

**alphaGo Won**  
**4:1**

# Lee Sedol 9-dan vs AlphaGo



Born in south Korean

**12**, became a professional player

**20**, became the world champion

# Computer Go AI?

DeepMind in London

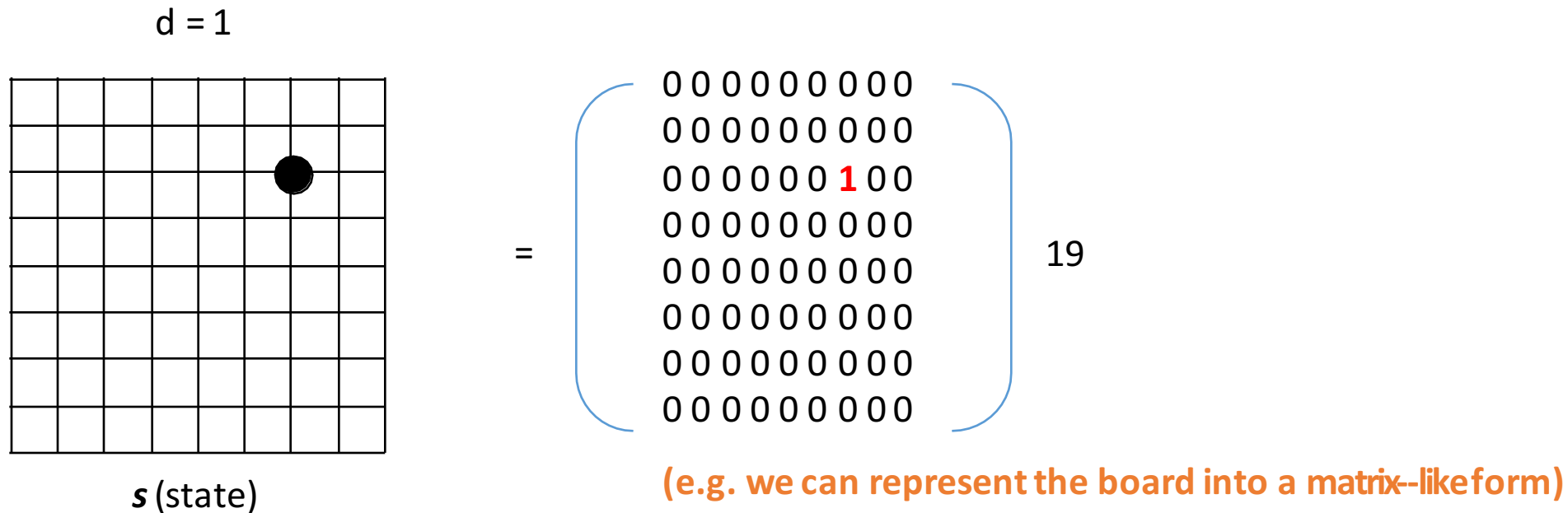
**2010**, start

2014, Google

**2015**, alphaGO



# Computer Go AI – Definition

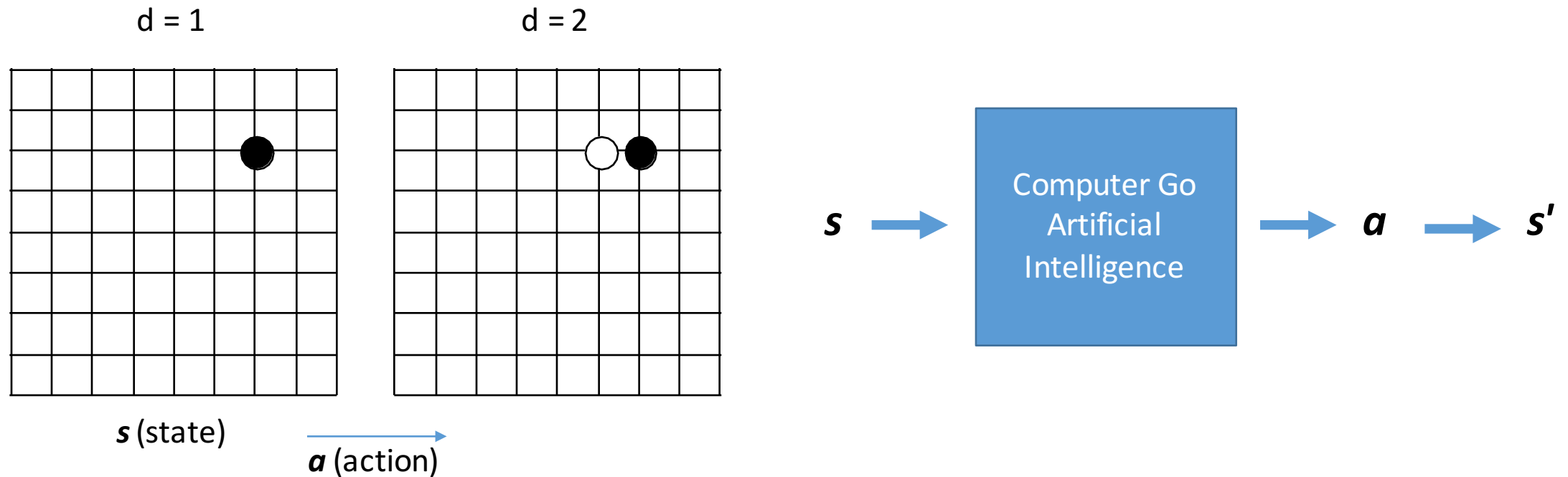


**\* The actual model uses other features than board positions as well, (Extended table 4 in paper)**

Extended Data Table 4 | Input features for rollout and tree policy

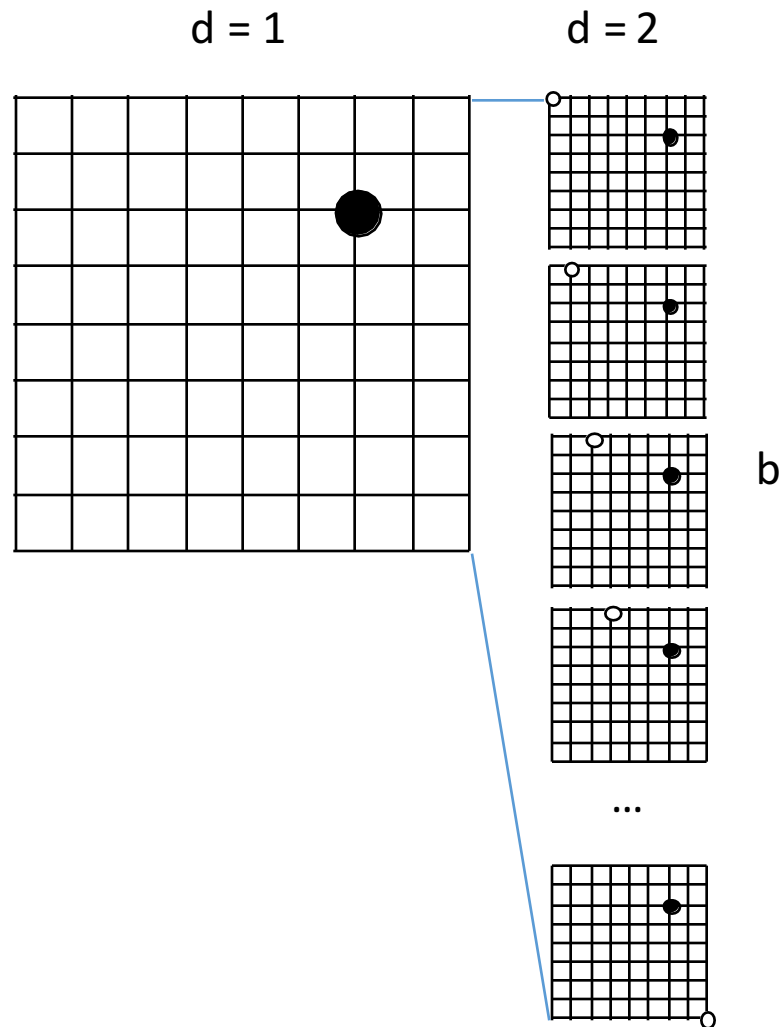
Feature	# of patterns	Description
Response	1	Whether move matches one or more response pattern features
Save atari	1	Move saves stone(s) from capture
Neighbour	8	Move is 8-connected to previous move
Nakade	8192	Move matches a <i>nakade</i> pattern at captured stone
Response pattern	32207	Move matches 12-point diamond pattern near previous move
Non-response pattern	69338	Move matches $3 \times 3$ pattern around move
Self-atari	1	Move allows stones to be captured
Last move distance	34	Manhattan distance to previous two moves
Non-response pattern	32207	Move matches 12-point diamond pattern centred around move

# Computer Go AI – Definition



Given  $s$ , pick the best  $a$

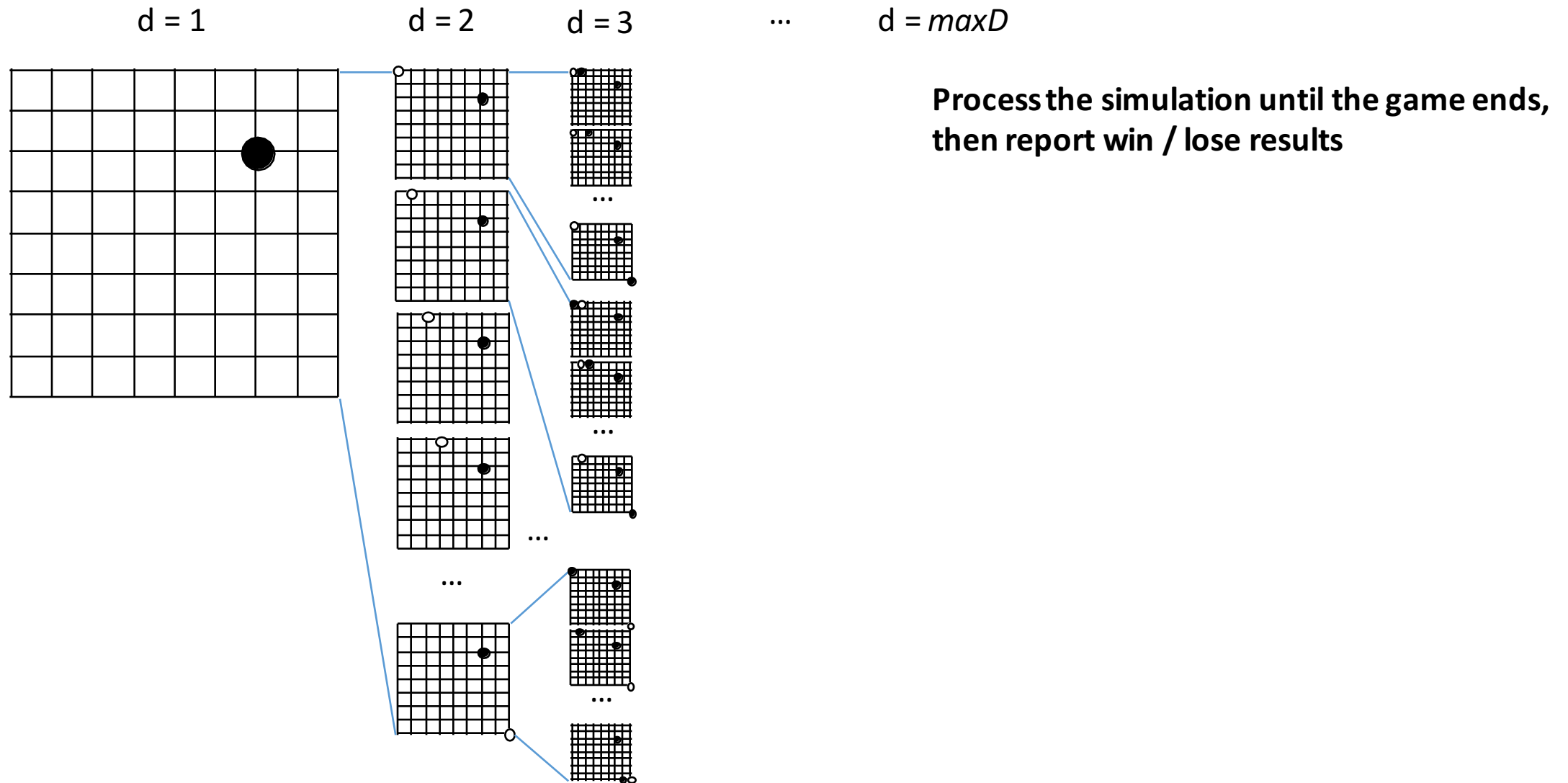
# Computer Go AI – An Implementation Idea?



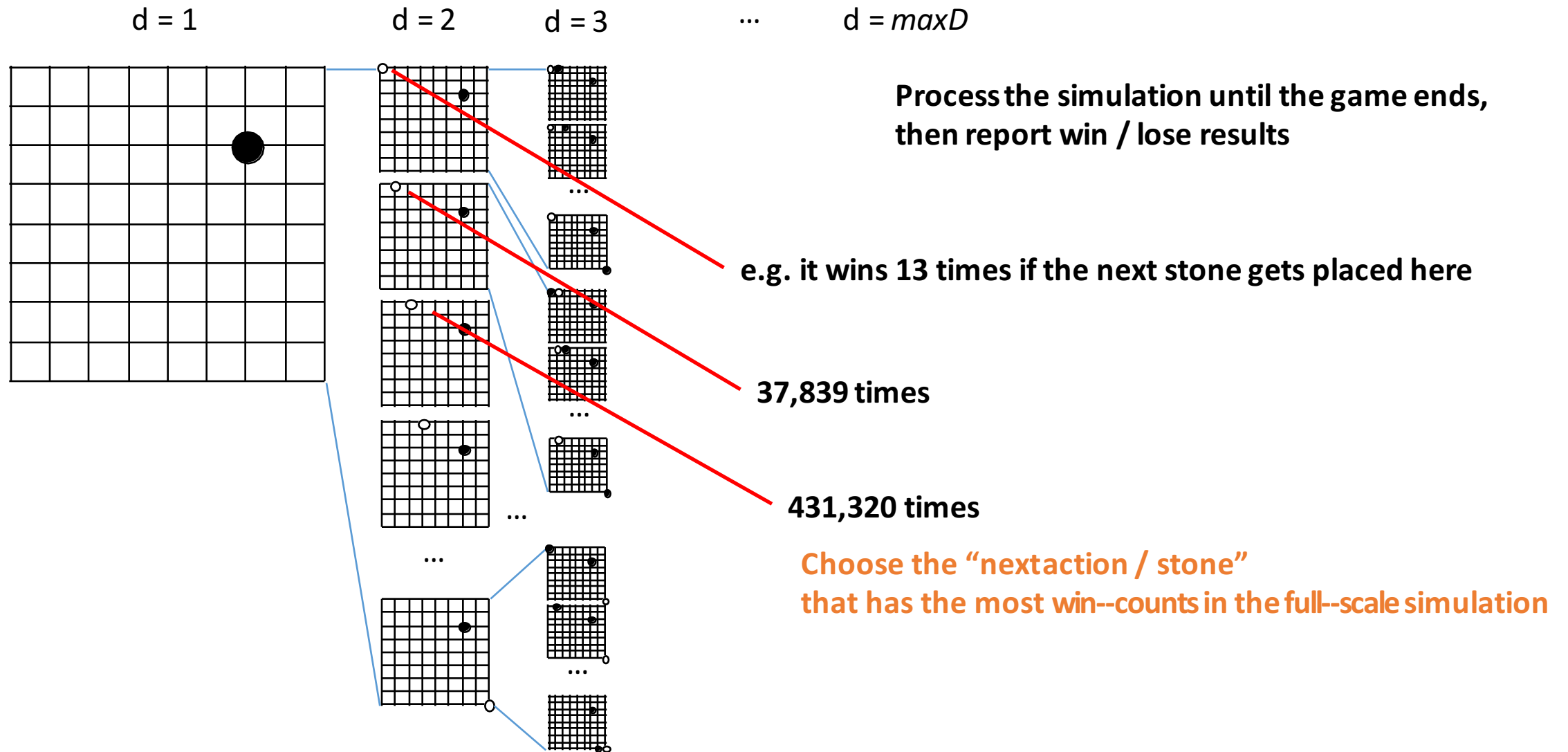
How about simulating all possible board positions?



# Computer Go AI – An Implementation Idea?



# Computer Go AI – An Implementation Idea?

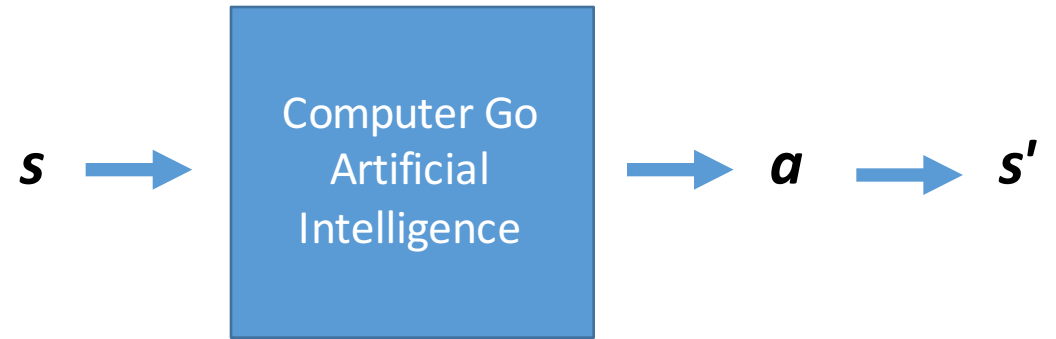
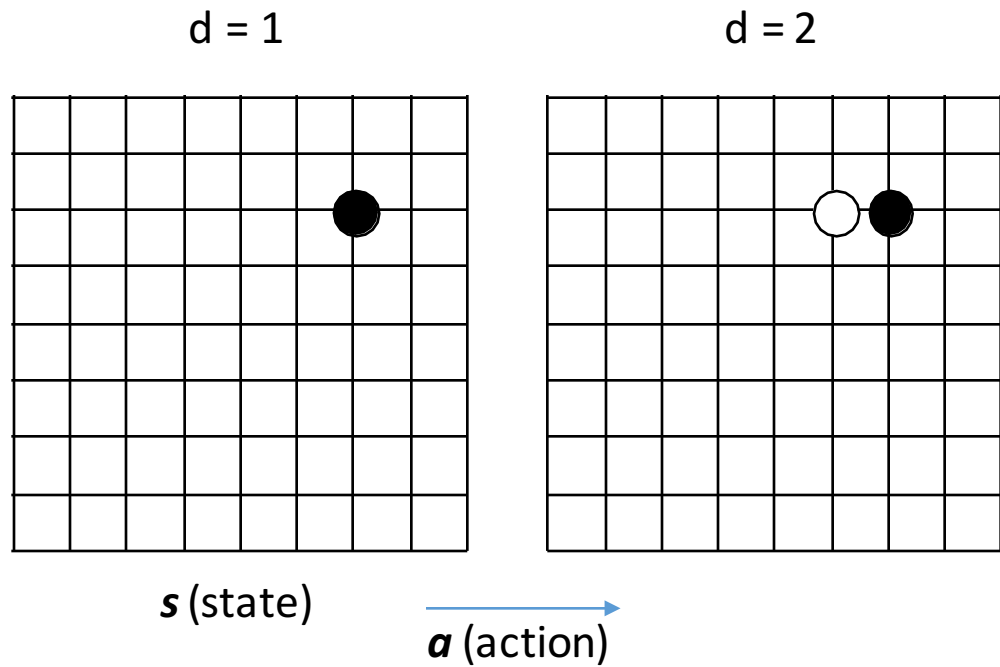


This is NOT possible; it is said the possible configurations of the board exceeds the number of atoms in the universe

**b ~ 250**  
**d ~ 150**



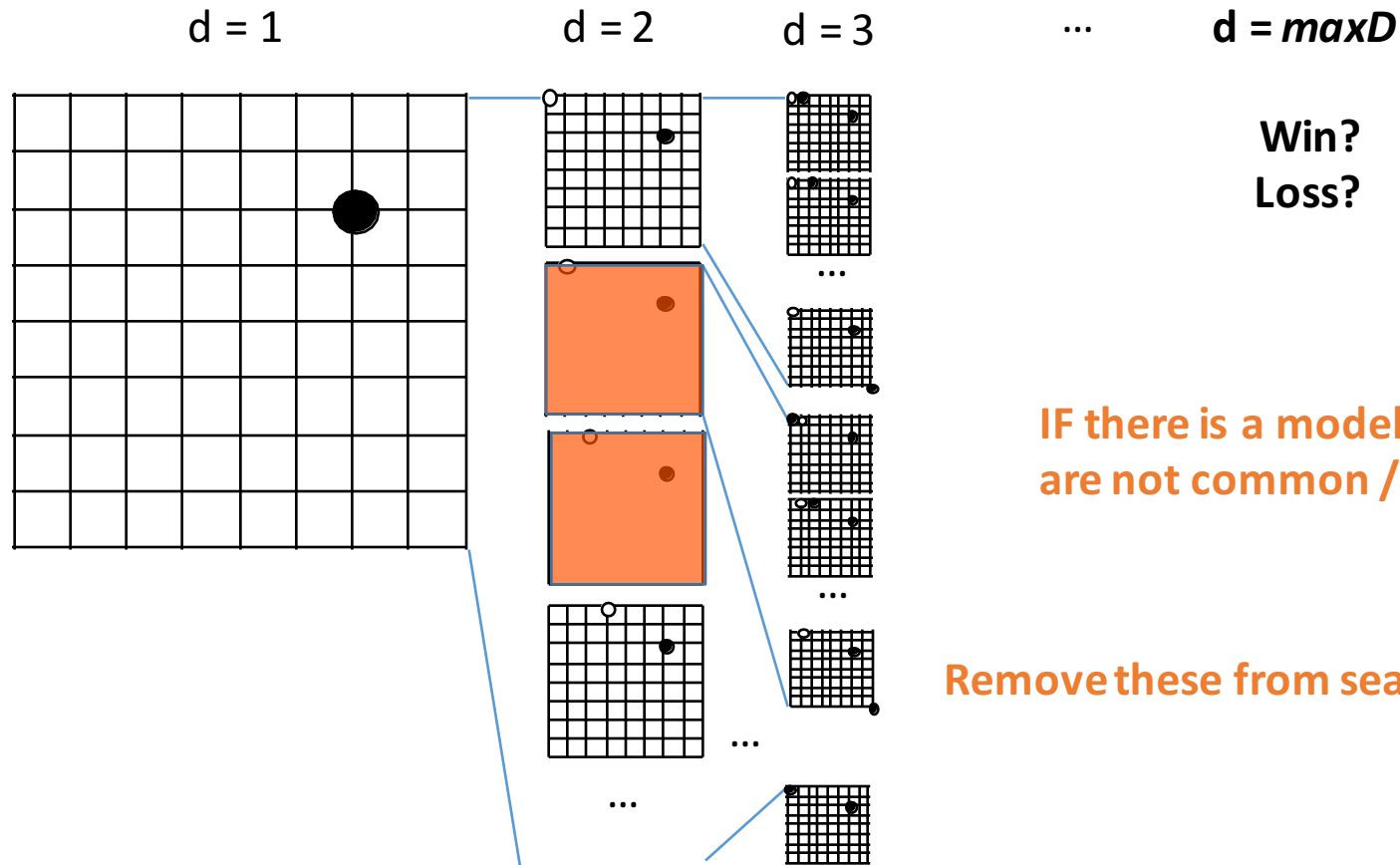
**5659799424266695229693199556804869862926581998836961**  
**3684891343062096883241205281967122071574755988121587**  
**2027456054864353492670804579578458515243255533277988**  
**43383789062500000000000000000000000000000000000000**  
**000**  
**00**



**Key: To Reduce Search Space !!!**

# Reducing Search Space

## 1. Reducing “action candidates” (Breadth Reduction)



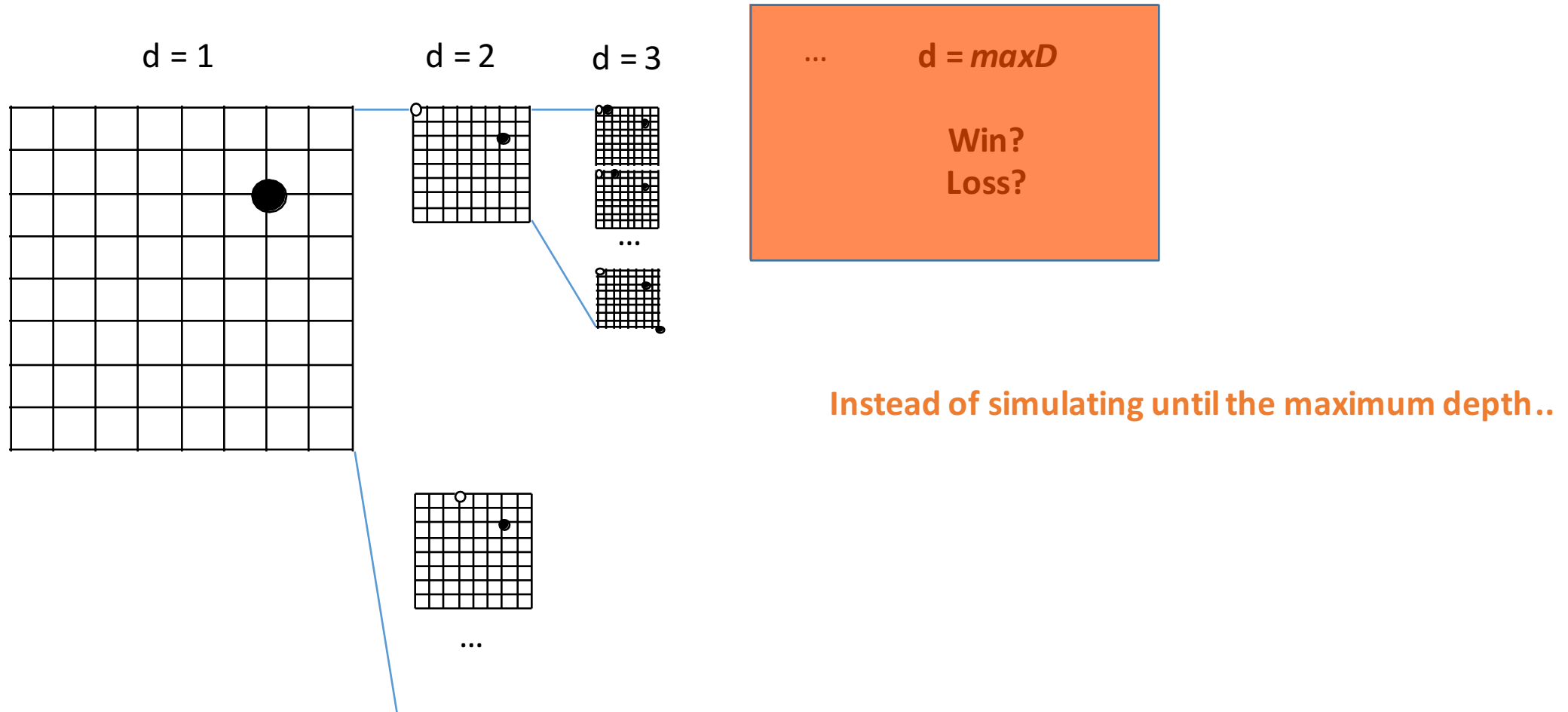
Win?  
Loss?

IF there is a model that can tell you that these moves are not common / probable (e.g. by experts, etc.) ...

Remove these from search candidates in advance (breadth reduction)

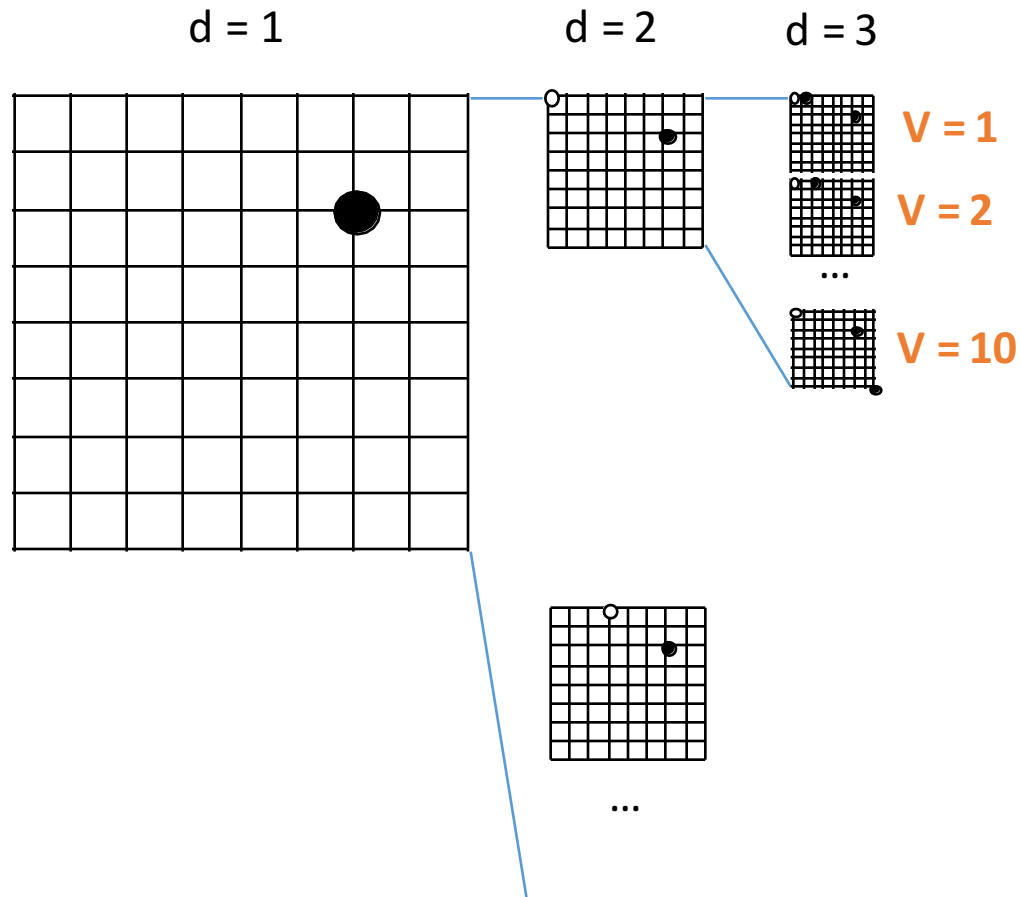
# Reducing Search Space

## 2. Position evaluation ahead of time (Depth Reduction)



# Reducing Search Space

## 2. Position evaluation ahead of time (Depth Reduction)



IF there is a function that can measure:  
 $V(s)$ : "board evaluation of state  $s$ "

# Reducing Search Space

- 1. Reducing “action candidates” (Breadth Reduction)**
- 2. Position evaluation ahead of time (Depth Reduction)**



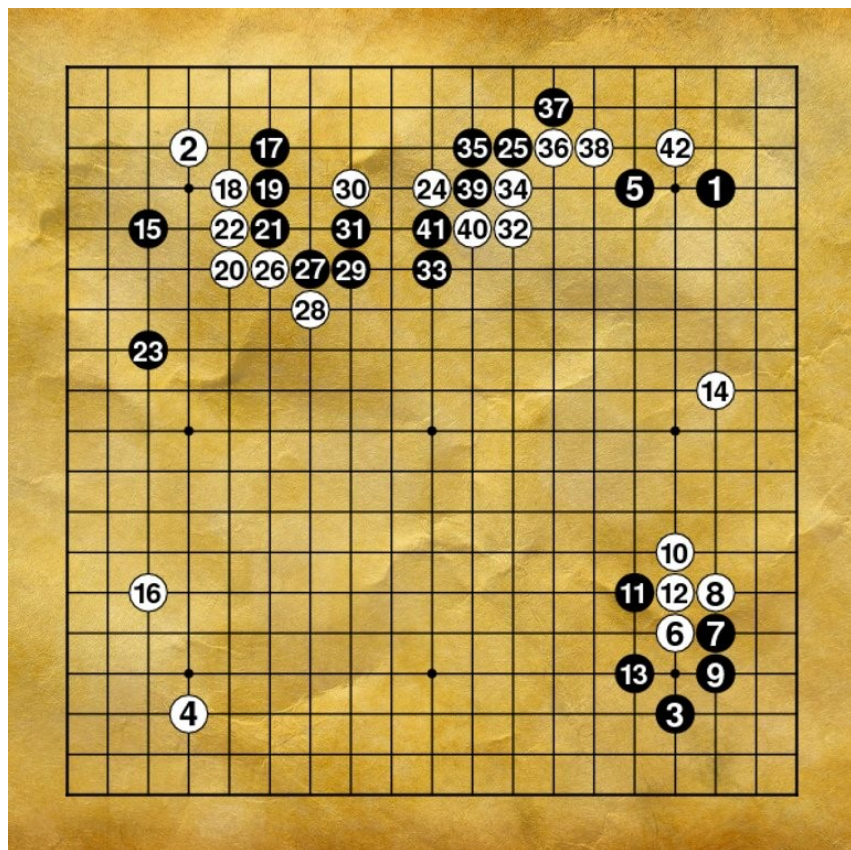
1. Reducing “action candidates”

Learning:  **$P(\text{next action} \mid \text{current state})$**

**$= P(a \mid s)$**

# 1. Reducing “action candidates”

## (1) Imitating expert moves (supervised learning)



Current State

Next Action

s1

a1

s2

a2

s3

a3

Prediction  
Model

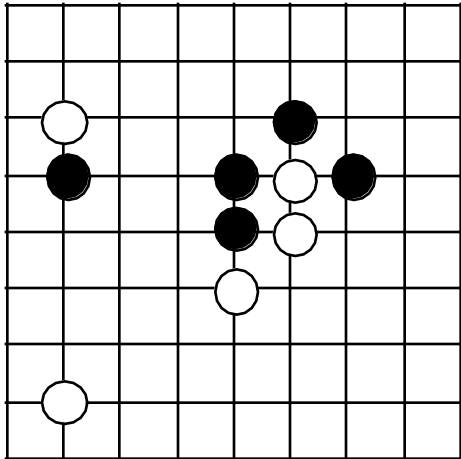
**Data:** Online Go experts (5~9 dan) (KGS)

160K games, 30M board positions

# 1. Reducing “action candidates”

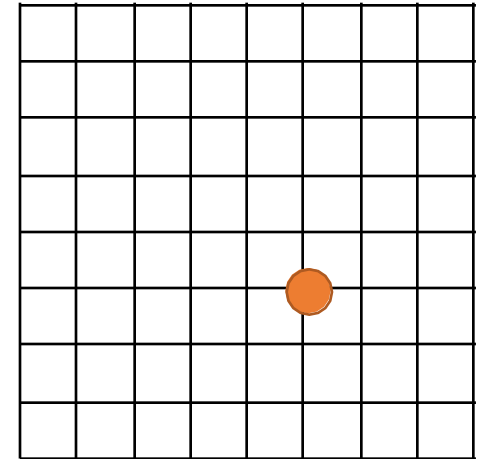
(1) Imitating expert moves (supervised learning)

Current Board



**Prediction Model**

Next Action



There are  $19 \times 19 = 361$   
possible actions  
(with different probabilities)

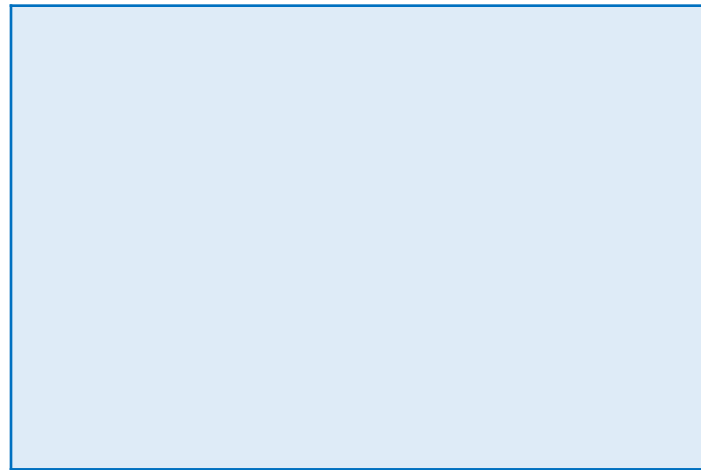
# 1. Reducing “action candidates”

(1) Imitating expert moves (supervised learning)

Current Board

```
00 000 0000  
00 000 1000  
0 -100 1 -1100  
0 1 00 1 -1000  
00 00 -10000  
00 000 0000  
0 -10000 0000  
00 000 0000
```

$s$



$g: s \rightarrow p(a|s)$

```
00 0000 000  
00 0000 000  
00 0000 000  
00 0000.2 0.100  
00 0000 0.4 0.200  
00 0000.1 000  
00 0000 000  
00 0000 000
```

$p(a|s)$

argmax

$a$

Next Action

```
00 00000000  
00 00000000  
00 00000000  
00 00000000  
00 00000000  
00 0000 1 000  
00 00000000  
00 00000000  
00 00000000
```

# 1. Reducing “action candidates”

(1) Imitating expert moves (supervised learning)

Current Board

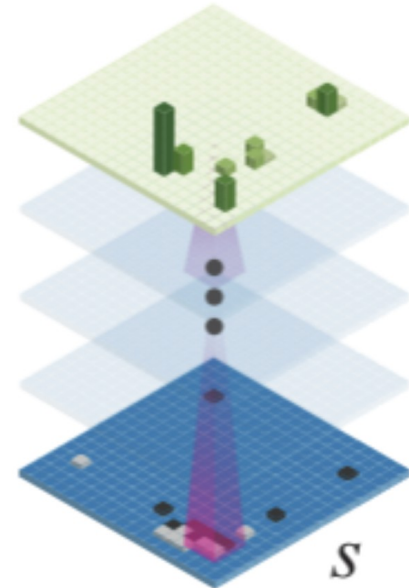
```
00 000 0000  
00 000 1000  
0 -100 1 -11 00  
0 1 00 1 -1 000  
00 00 -10000  
00 000 0000  
0 -1000 0000  
00 000 0000
```

$s$

Supervised Learning Policy

**Prediction  
Model**

$g: s \rightarrow p(a|s)$



$p(a|s)$

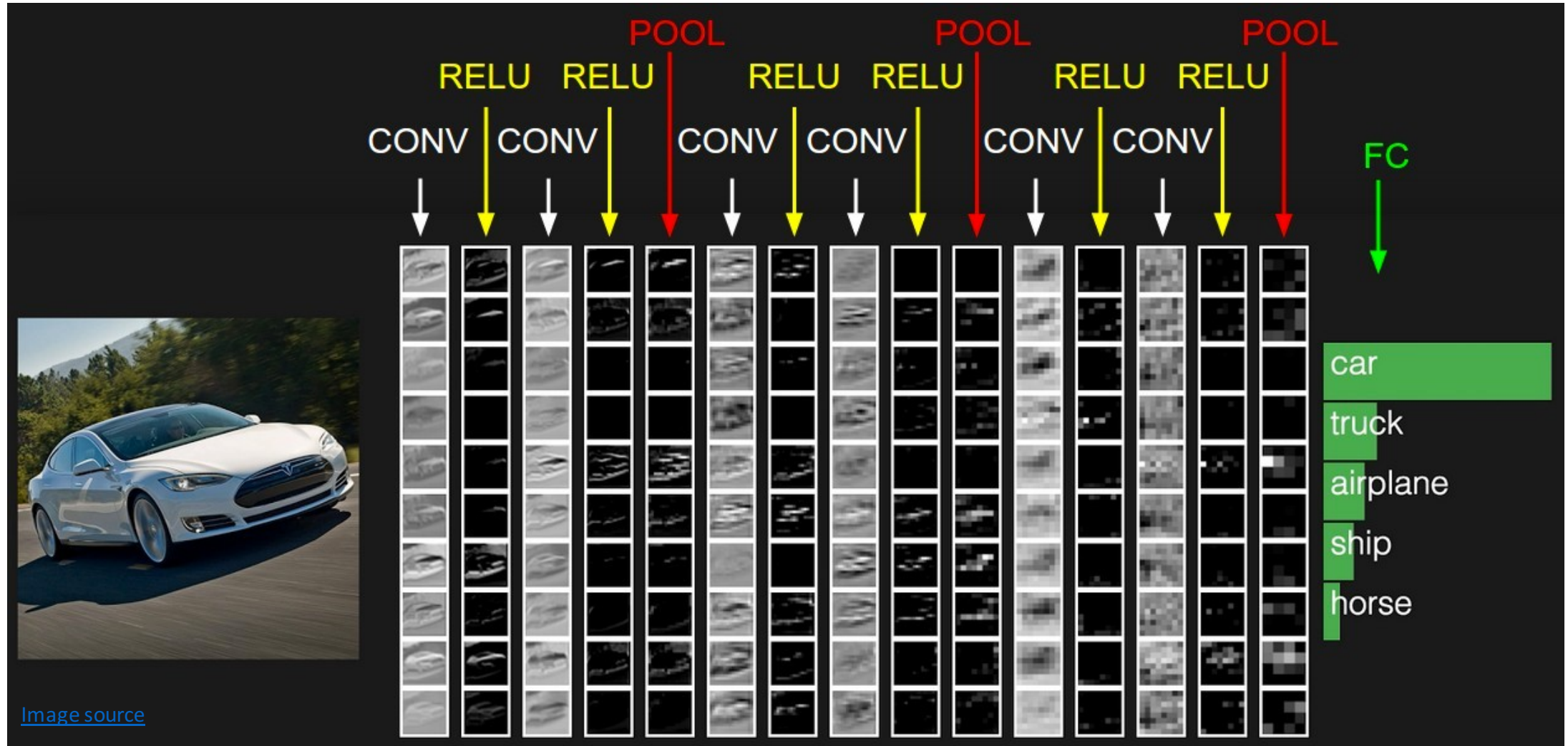
argmax

$a$

Next Action

```
000000000  
000000000  
000000000  
000000000  
000000000  
000001000  
000000000  
000000000  
000000000
```

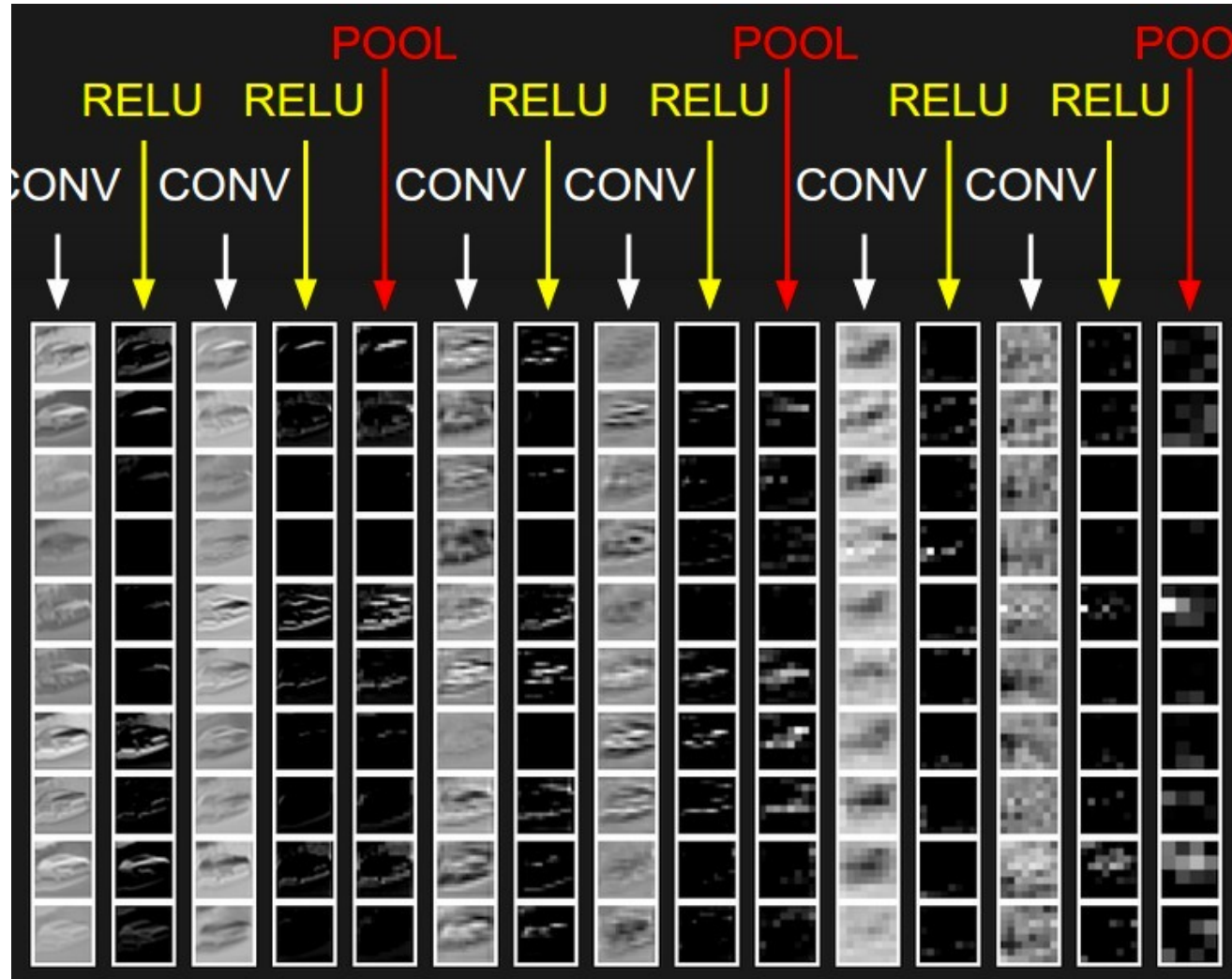
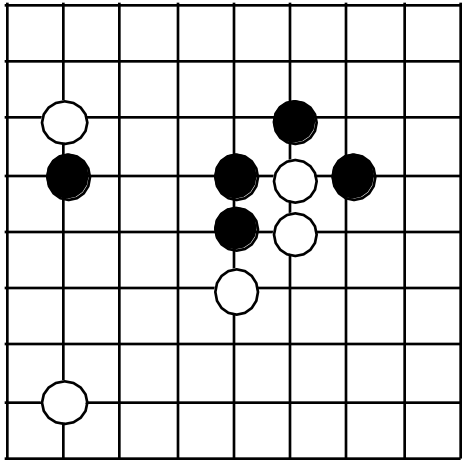
# Convolutional Neural Network (CNN)



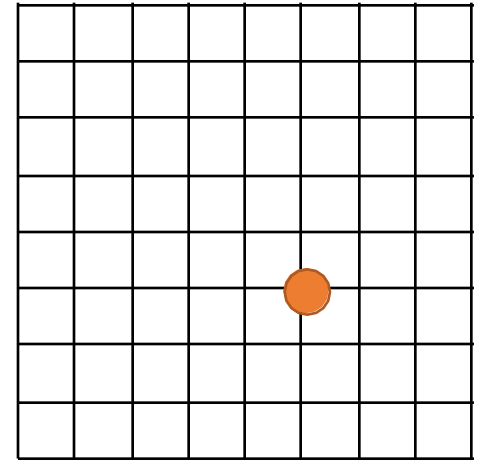
CNN is a powerful model for image recognition tasks; it abstracts out the input image through convolution layers

# Convolutional Neural Network (CNN)

Current Board



Next Action

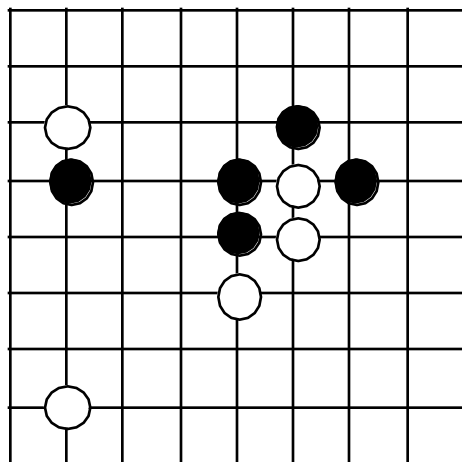


And they use this CNN model (similar architecture) to evaluate the board position; which learns "some" spatial invariance

# 1. Reducing “action candidates”

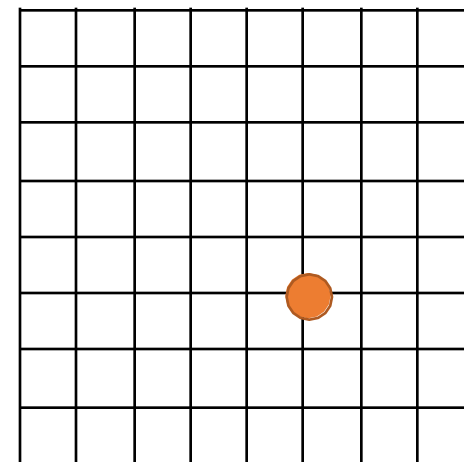
(1) Imitating expert moves (supervised learning)

Current Board



**Expert Moves Imitator Model  
(w/ CNN)**

Next Action



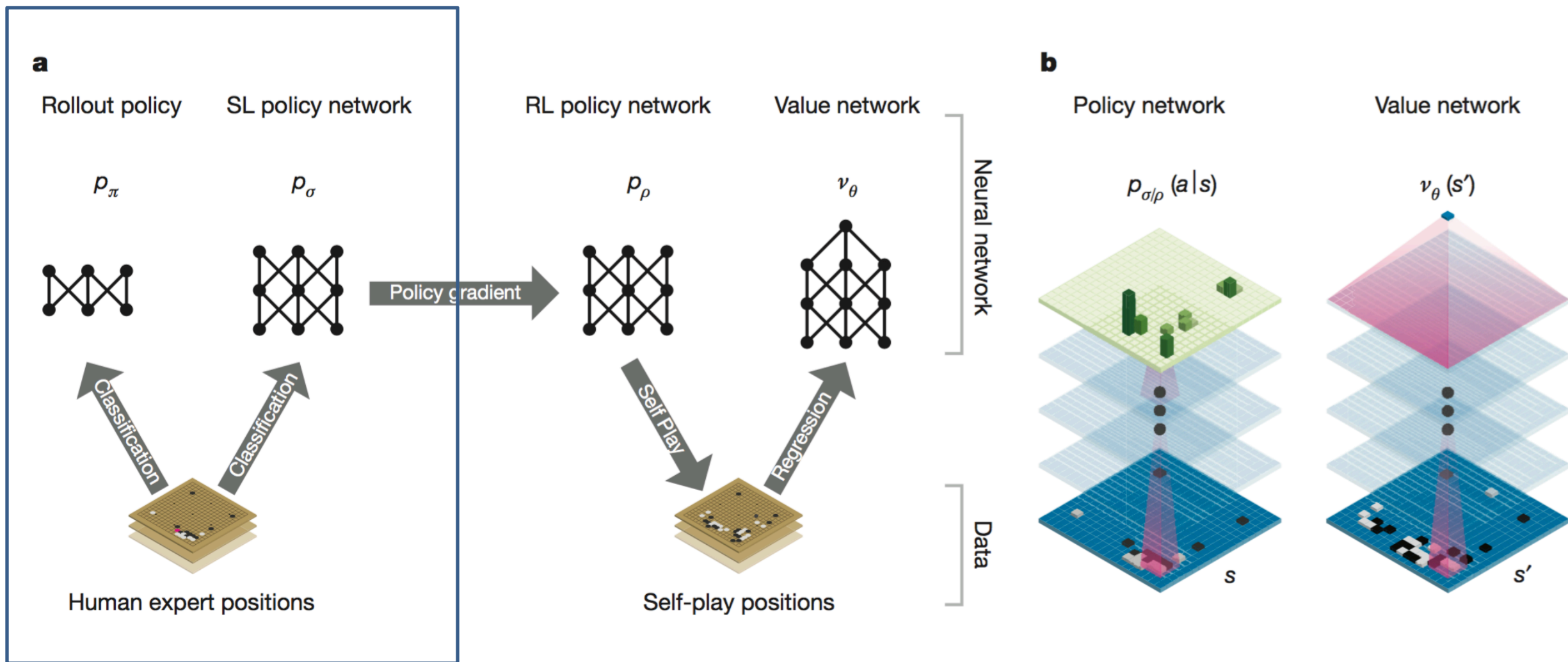
**50** GPUs, 3 week

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

**ACC: 57%**



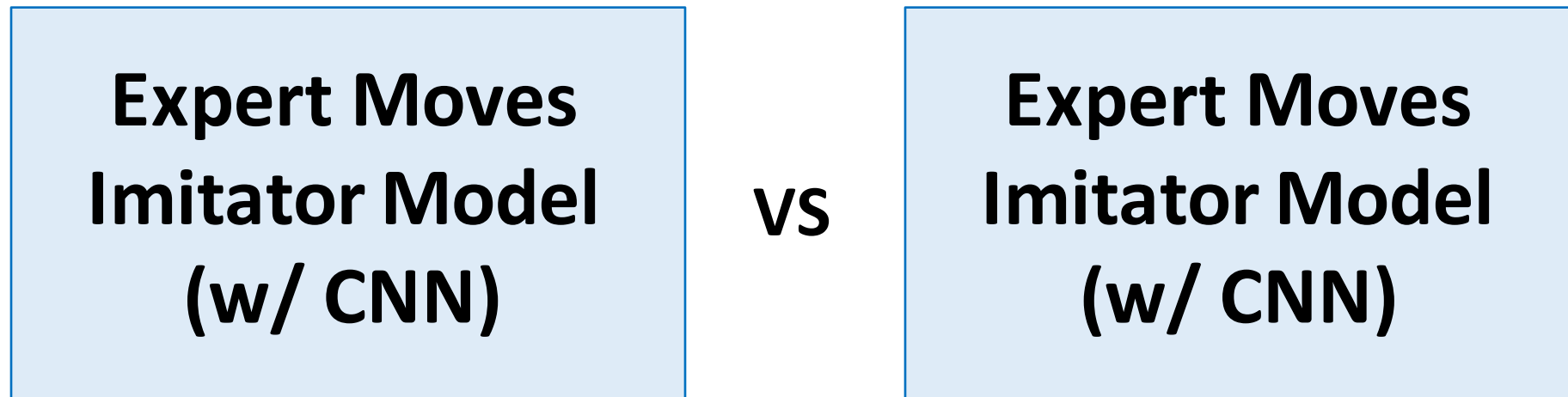
# AlphaGo



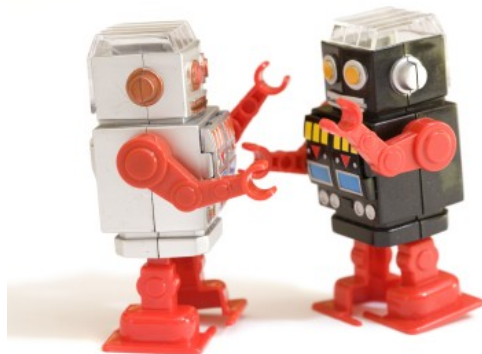
# 1. Reducing “action candidates”

(2) Improving through self-plays (reinforcement learning)

Improving by playing against itself



Goal shift:

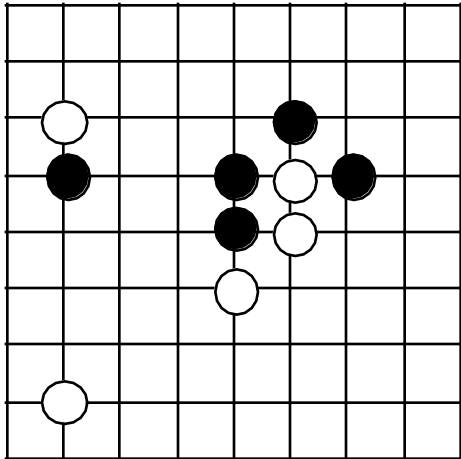


**Imitation -> winning !!**

# 1. Reducing “action candidates”

(2) Improving through self-plays (reinforcement learning)

Board position



**Expert Moves Imitator Model  
(w/ CNN)**

win/loss

**Loss**

$z = -1$

**Win**

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

$z = +1$

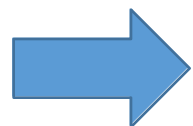
# 1. Reducing “action candidates”

(2) Improving through self-plays (reinforcement learning)

Older models vs. newer models



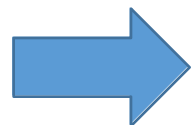
It uses the same topology as the expert moves imitator model, and just uses the updated parameters



**Return:** board positions, win/lose info

# 1. Reducing “action candidates”

(2) Improving through self-plays (reinforcement learning)



**Return:** board positions, win/lose info

1. Reducing “action candidates”

(2) Improving through self-plays (reinforcement learning)

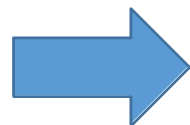
**Expert Moves  
Imitator Model**

VS

**Most Updated  
Model**

Supervised Learning Policy

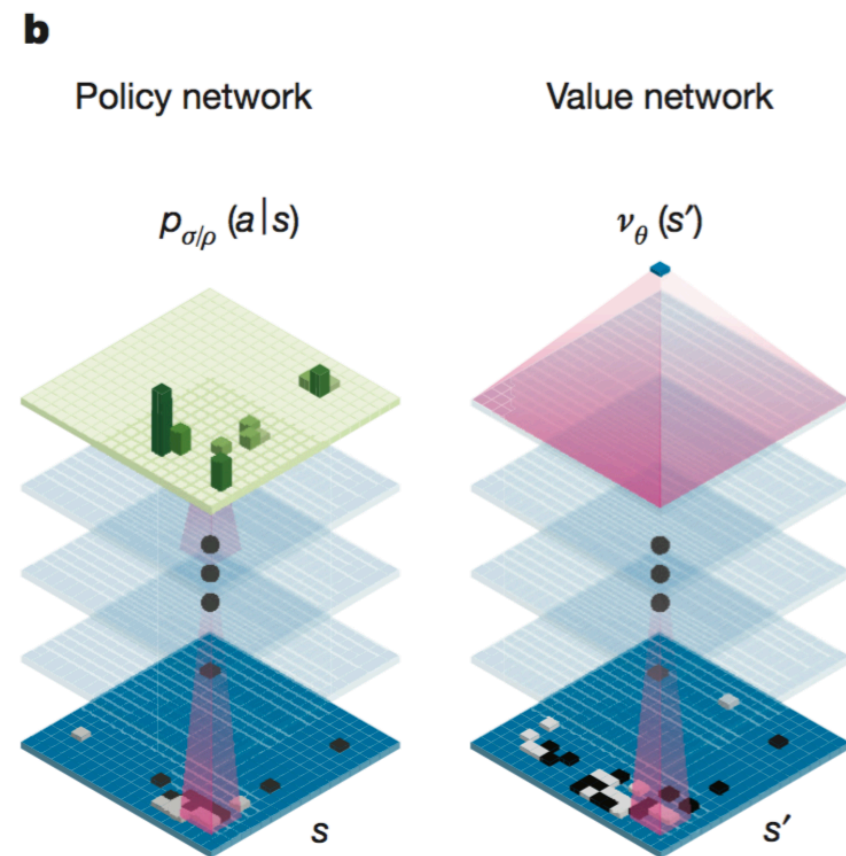
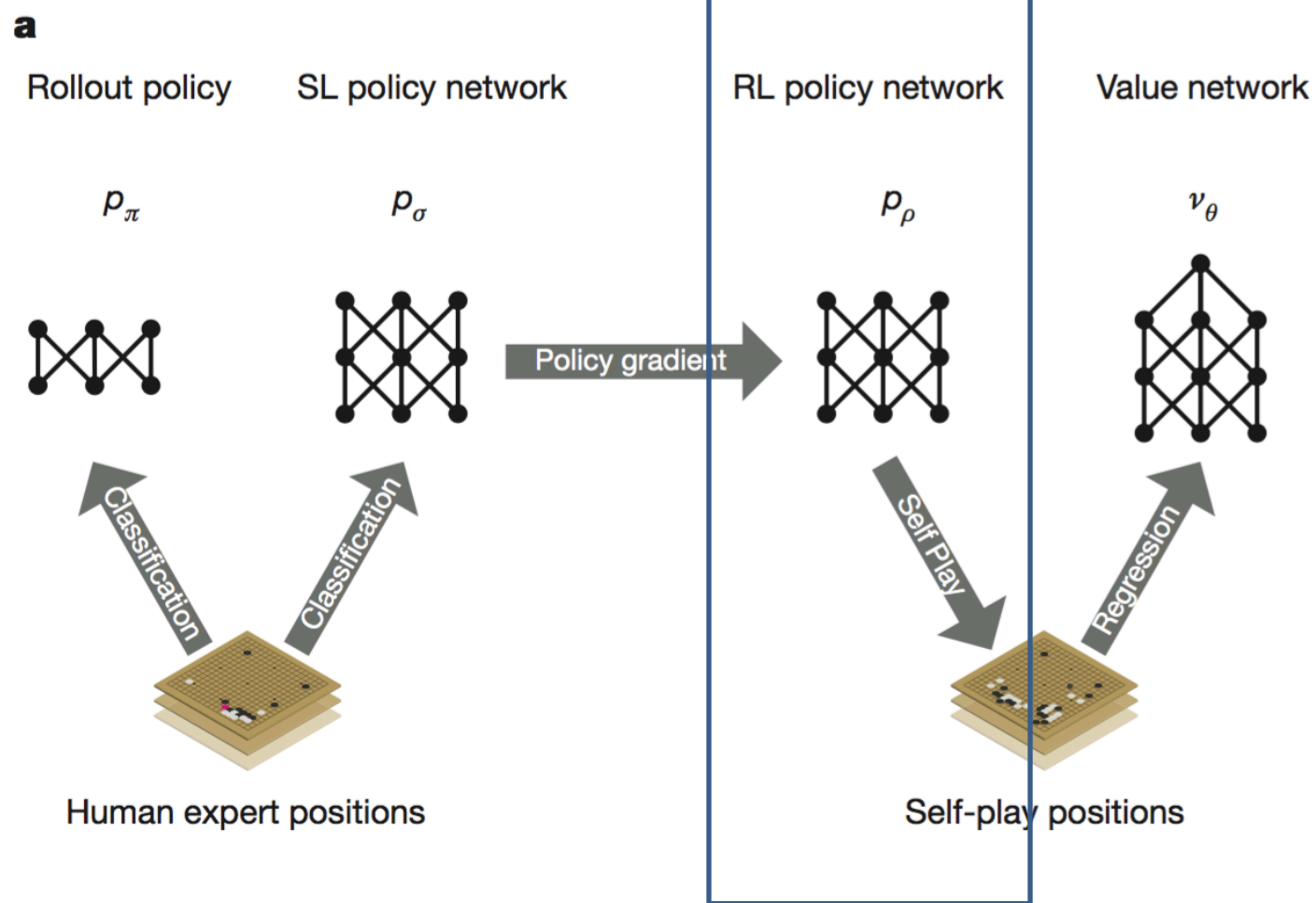
Reinforcement Learning Policy



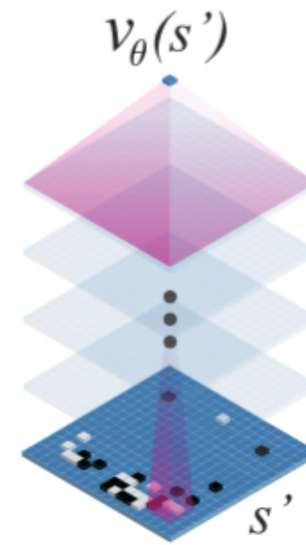
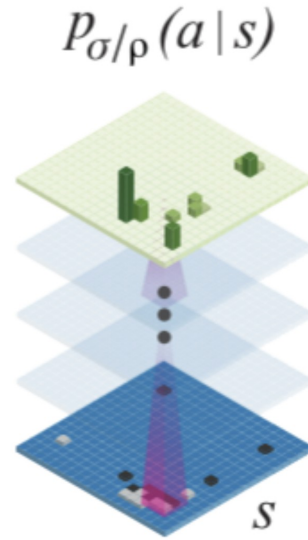
The final model wins **80%** of the time  
when playing against the first model

**50** GPUs, 1 day

# AlphaGo

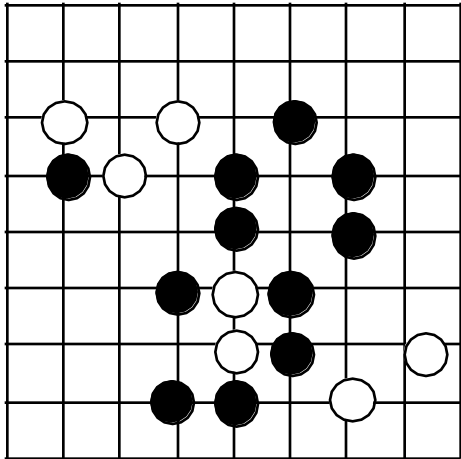


# 2. Board Evaluation



Adds a regression layer to the model  
Predicts values between 0~1  
Close to 1: a good board position  
Close to 0: a bad board position

Board Position



**Updated Model  
ver 1,000,000**

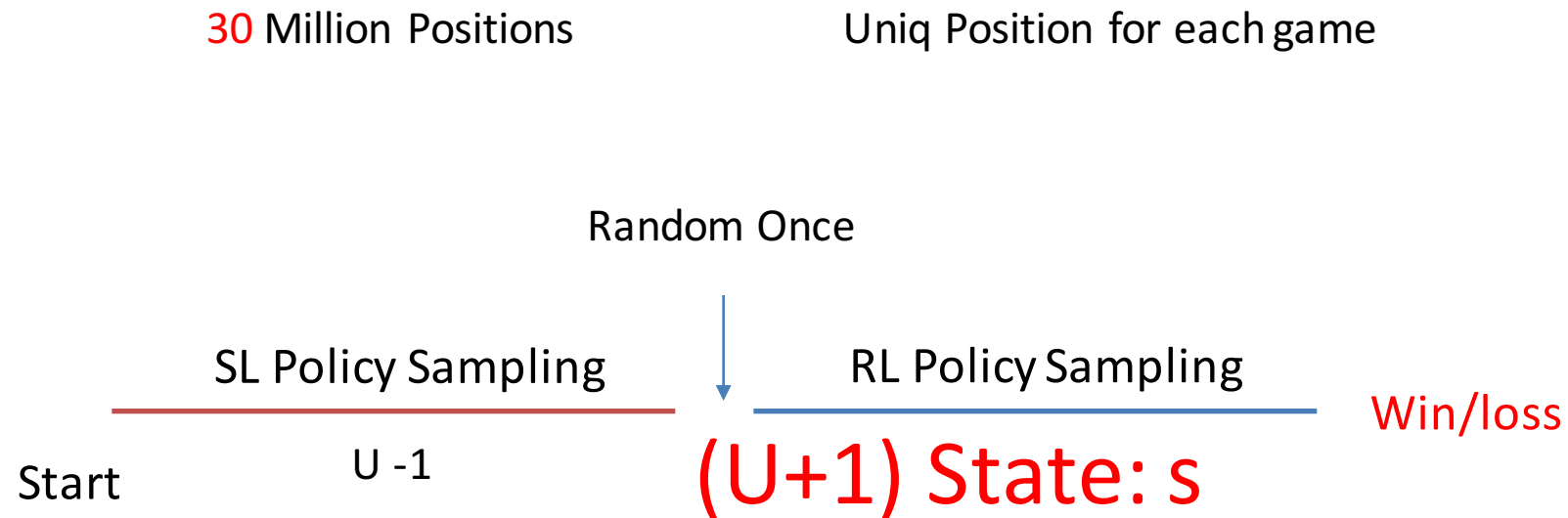
**Value  
Prediction  
Model  
(Regression)**

Win / Loss

**Win  
(0~1)**

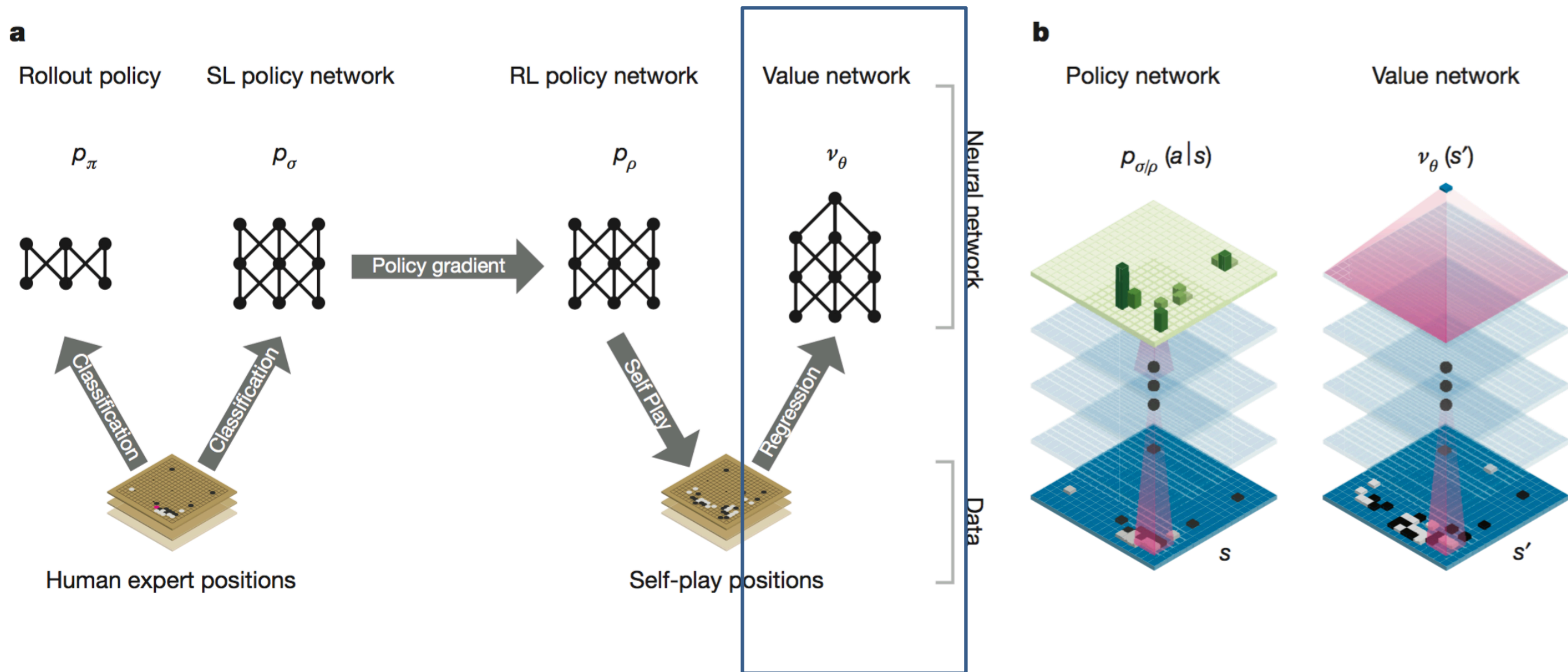


## 2. Board Evaluation



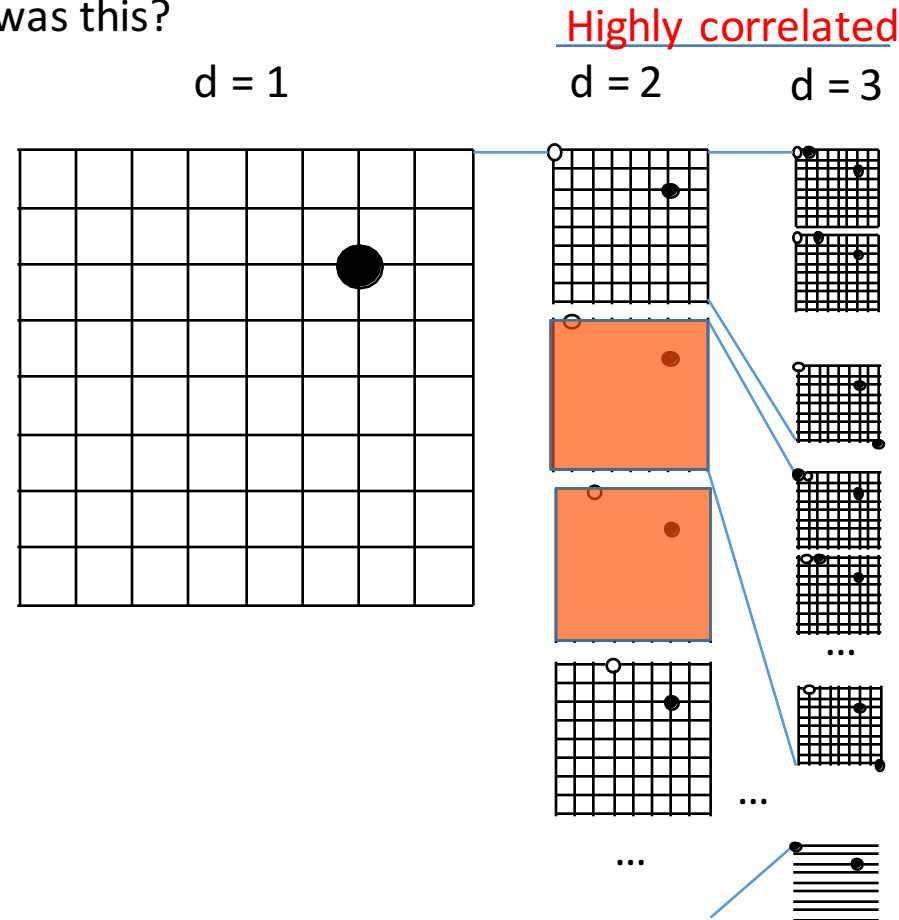
50 GPUs, one week **Training:**  $\Delta\theta \propto \frac{\partial v_{\theta}(s)}{\partial \theta} (z - v_{\theta}(s))$

# AlphaGo



# In-class Question

Q(2). To train the value network, a set of distinct positions, each from a different game, was constructed. The policy network is then played against itself from that position to get an estimate of the value of that state. The learning update is then only applied to the initial state from which the play started, as no updates are made to states in the rest of the game. Why was this?



Over-represents

# Reducing Search Space

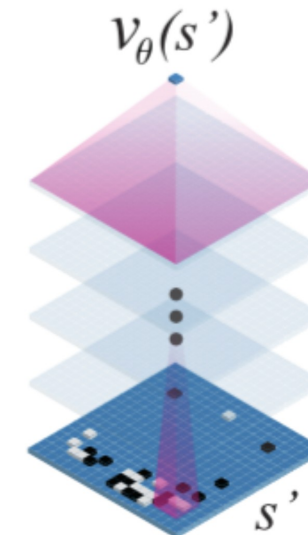
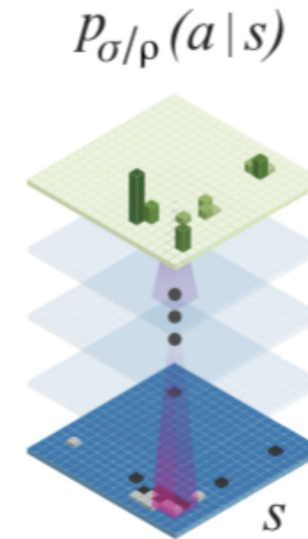
## 1. Reducing “action candidates” (Breadth Reduction)

### Policy Network

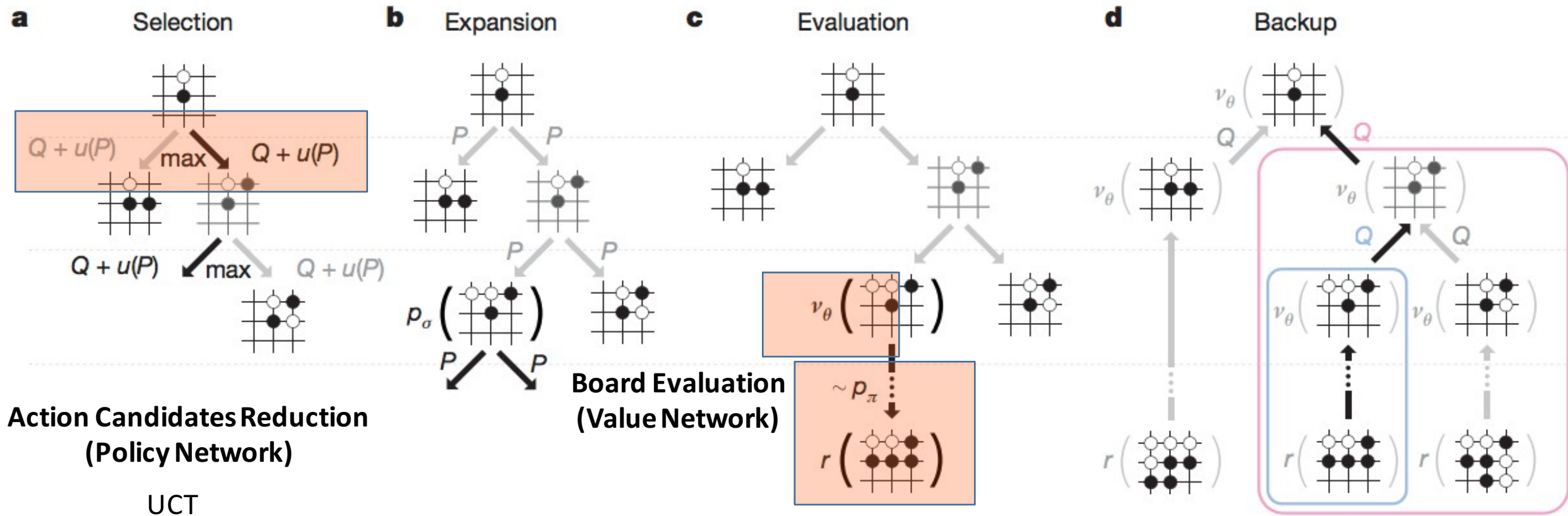
Reinforcement learning policy +  
(rollout policy, simple model to  
speed up the sampling)

## 2. Board Evaluation (Depth Reduction)

### Value Network



# Looking ahead (w/ Monte Carlo Search Tree) + UCT



Action Candidates Reduction  
(Policy Network)

Board Evaluation  
(Value Network)

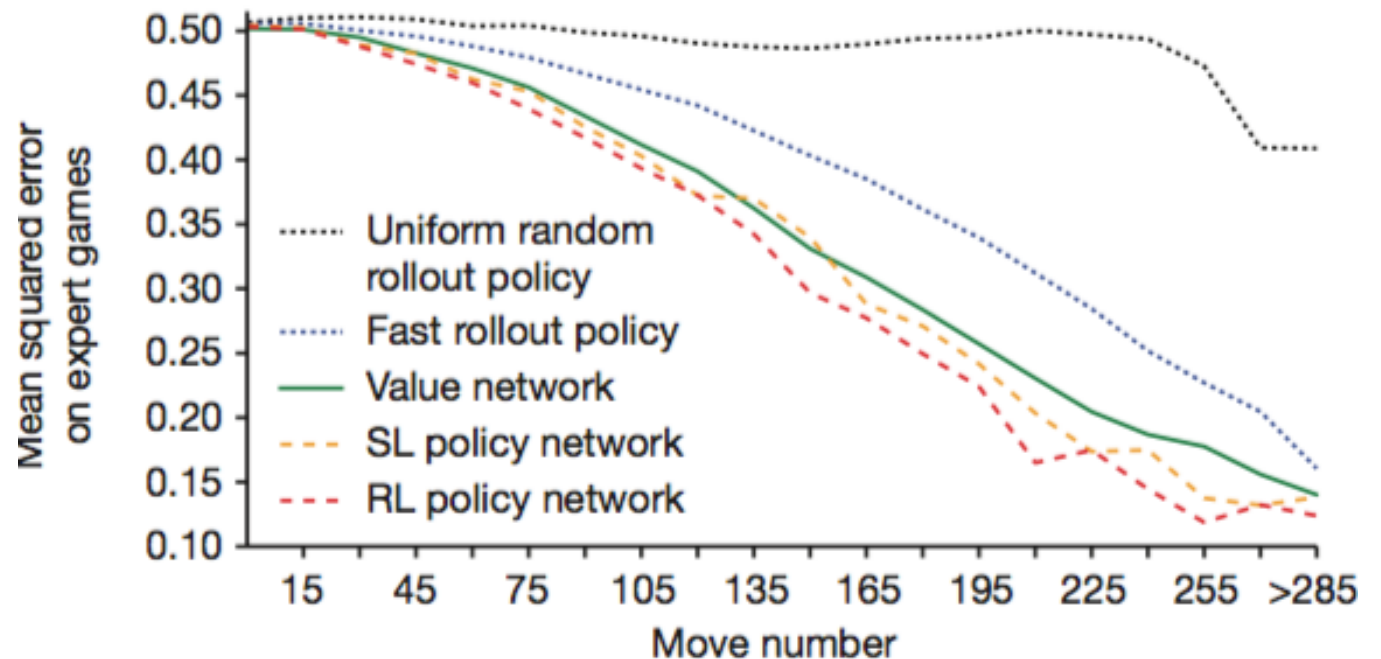
(Rollout): Faster version of estimating  $p(a | s)$   
 $\rightarrow$  uses shallow networks (3 ms  $\rightarrow$  2 $\mu$ s)

UCT

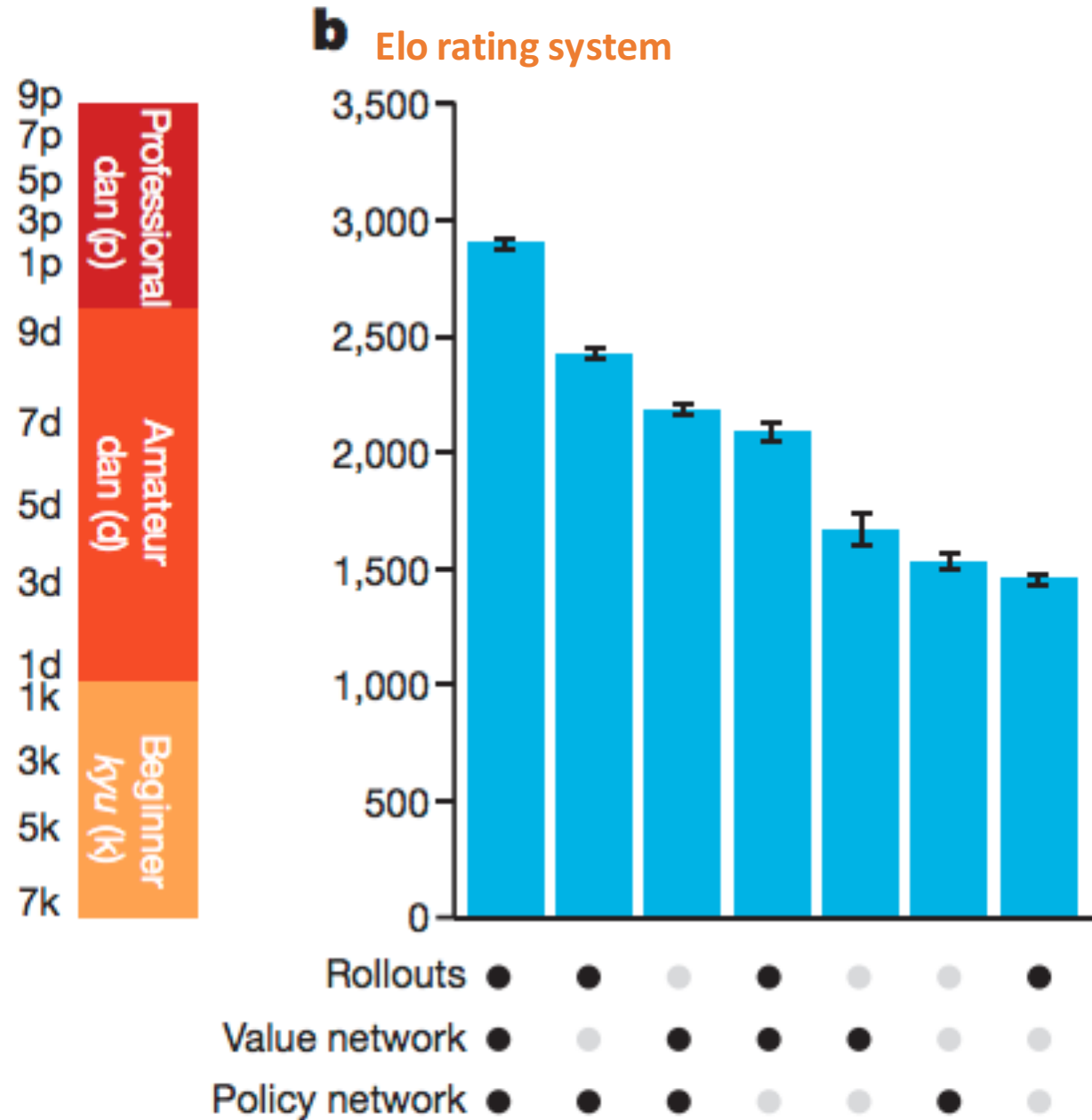
# In-class Question

Q(1). For AlphaGo, two policies are learned directly from expert games, one using a deep network, and the other a linear approximator. The linear approximator, which was used to generate moves during the playouts, is much faster to compute but is less accurate at predicting good moves. Explain why they might have made the decision to use the less accurate policy? In addition, note that they could have uniformly selected moves from all possible moves, which would have been even faster to computer. Discuss what this suggests about how we should develop policies used to generate playouts in Monte Carlo Tree Search.

2 us VS 3 ms



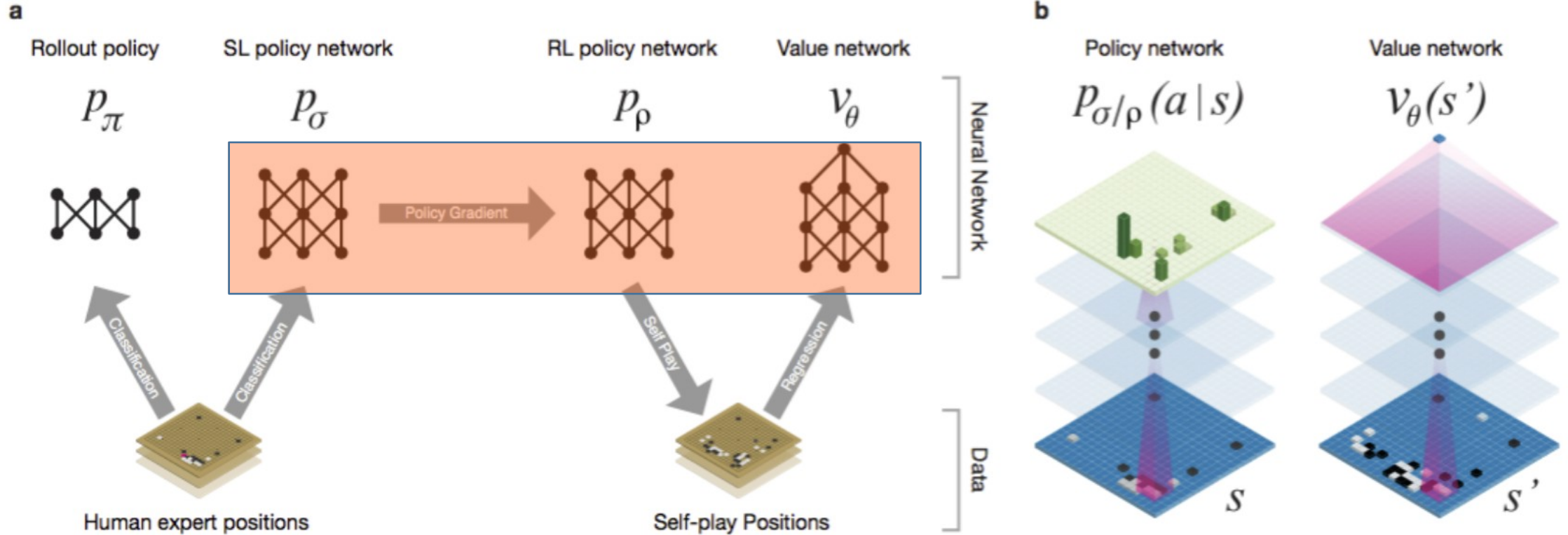
# Results



Performance with different combinations of AlphaGo components

# Takeaways

Use the networks trained for a certain task (with different loss objectives) for several other tasks





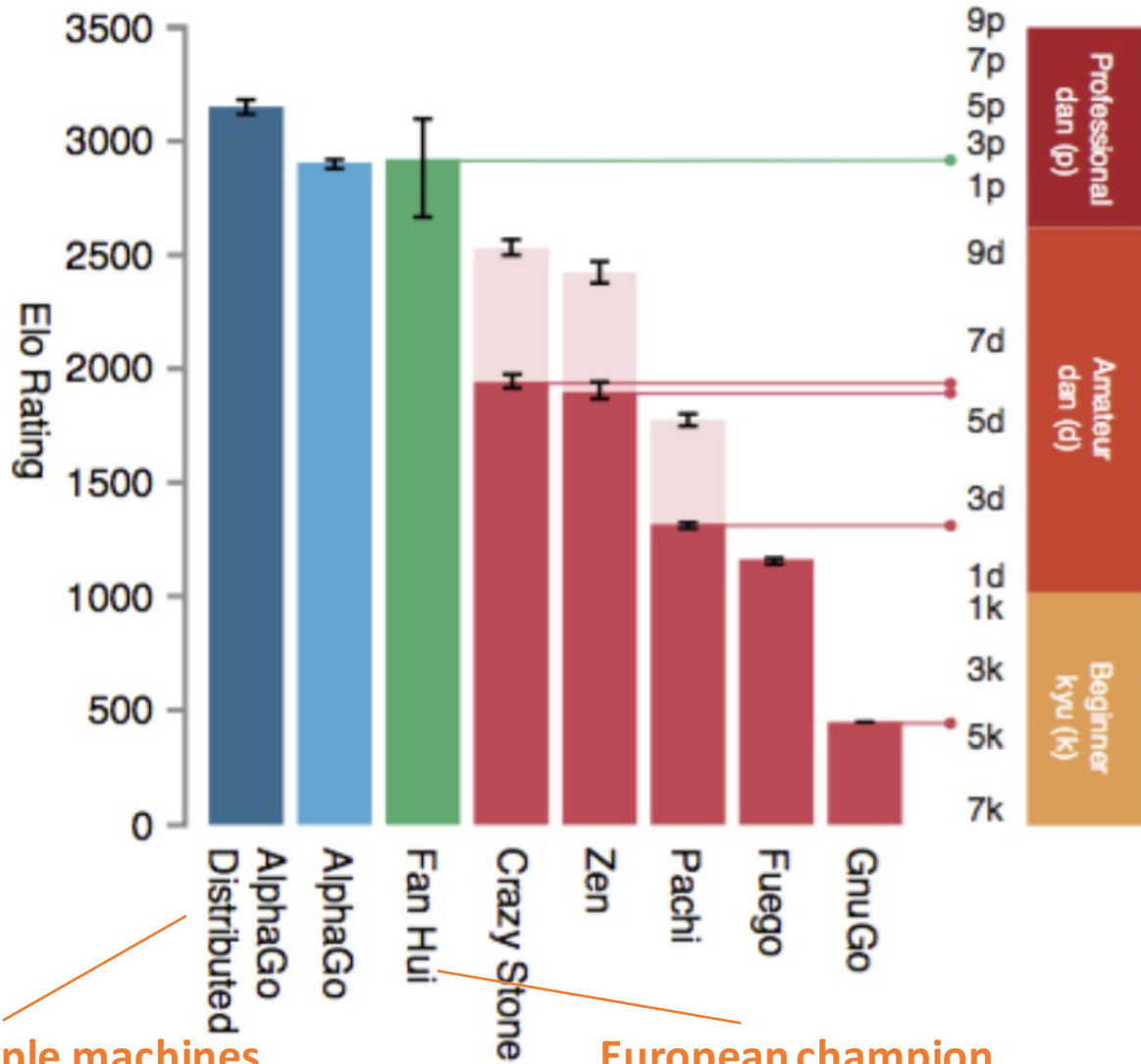
# Lee Sedol vs AlphaGo

## Energy Consumption

Lee Sedol	AlphaGo
<ul style="list-style-type: none"><li>- Recommended calories for a man per day : ~2,500 cal</li><li>- Assumption: Lee consumes the entire amount of per-day calories in this one game <math>2,500 \text{ cal} * 4.184 \text{ J/cal}</math></li></ul> <p><b><math>\approx 10\text{k [J]}</math></b></p>	<ul style="list-style-type: none"><li>- Assumption: CPU: ~100 W, GPU: ~300 W</li><li>- <b>1,202 CPUs, 176 GPUs</b></li></ul> <p><math>170,000 \text{ J/sec} * 5 \text{ hr} * 3,600 \text{ sec/hr}</math></p> <p><b><math>\approx 3,000\text{M [J]}</math></b></p> <p><b>= 300 k Lee</b></p>

A very, very rough calculation ;)

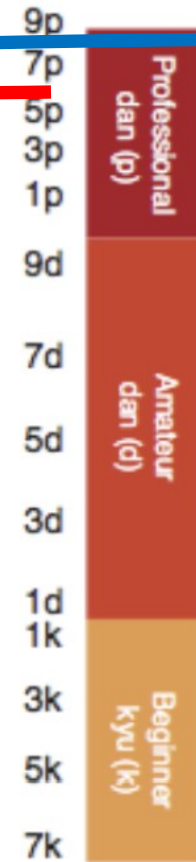
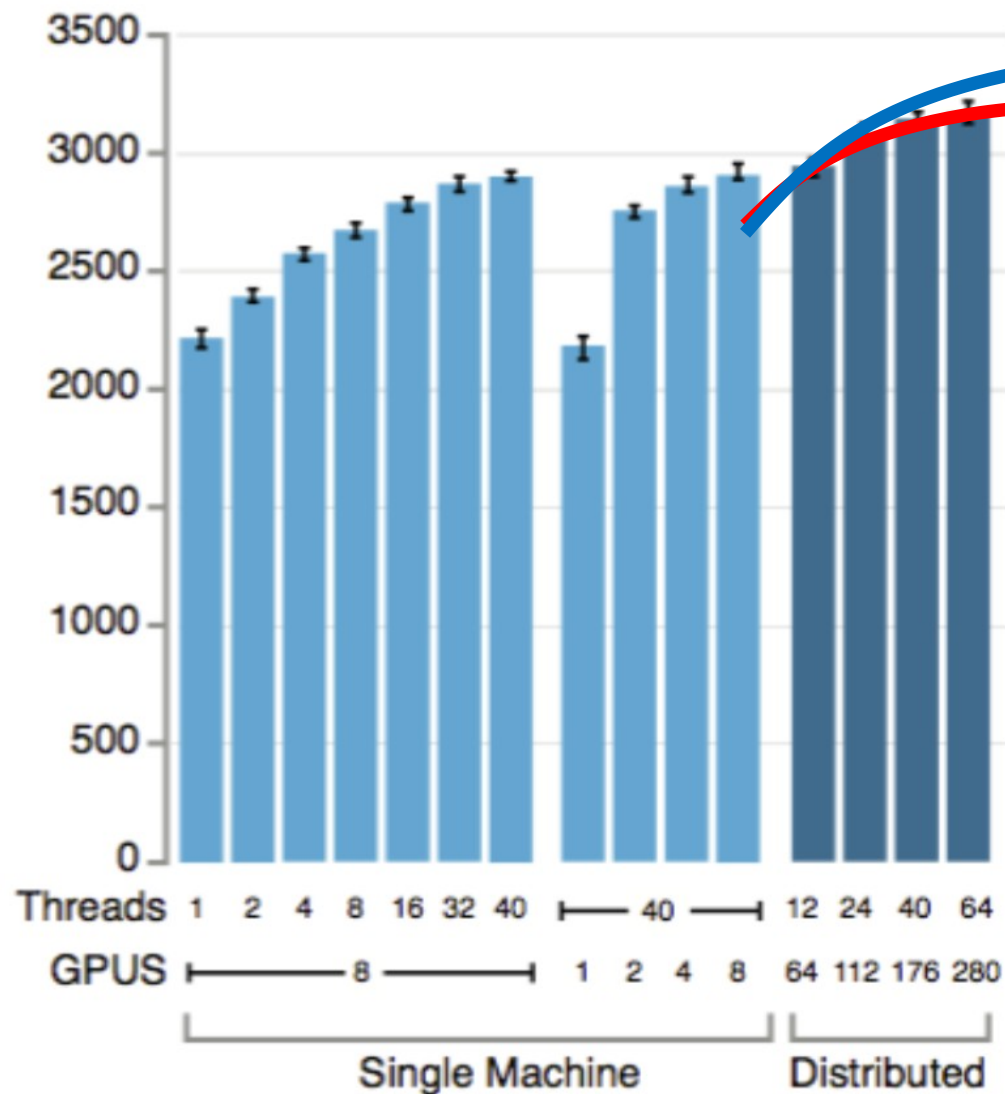
# AlphaGo is estimated to be around ~5-dan



= multiple machines

European champion

# Taking CPU / GPU resources to virtually infinity?



No one knows  
how it will converge

But Google has promised not to use more CPU/GPUs than they used for Fan Hui for the game with Lee

# AlphaGo learns millions of Go games everyday

AlphaGo will presumably converge to some point eventually.

However, in the Nature paper they don't report how AlphaGo's performance improves as a function of times AlphaGo plays against itself (self--plays).

# Conclusion

First Time

Computer can defeat professional player in Go

# Reference

- Silver, David, et al. "Mastering the game of Go with deep neural networks and tree search." *Nature* 529.7587 (2016): 484–489.