

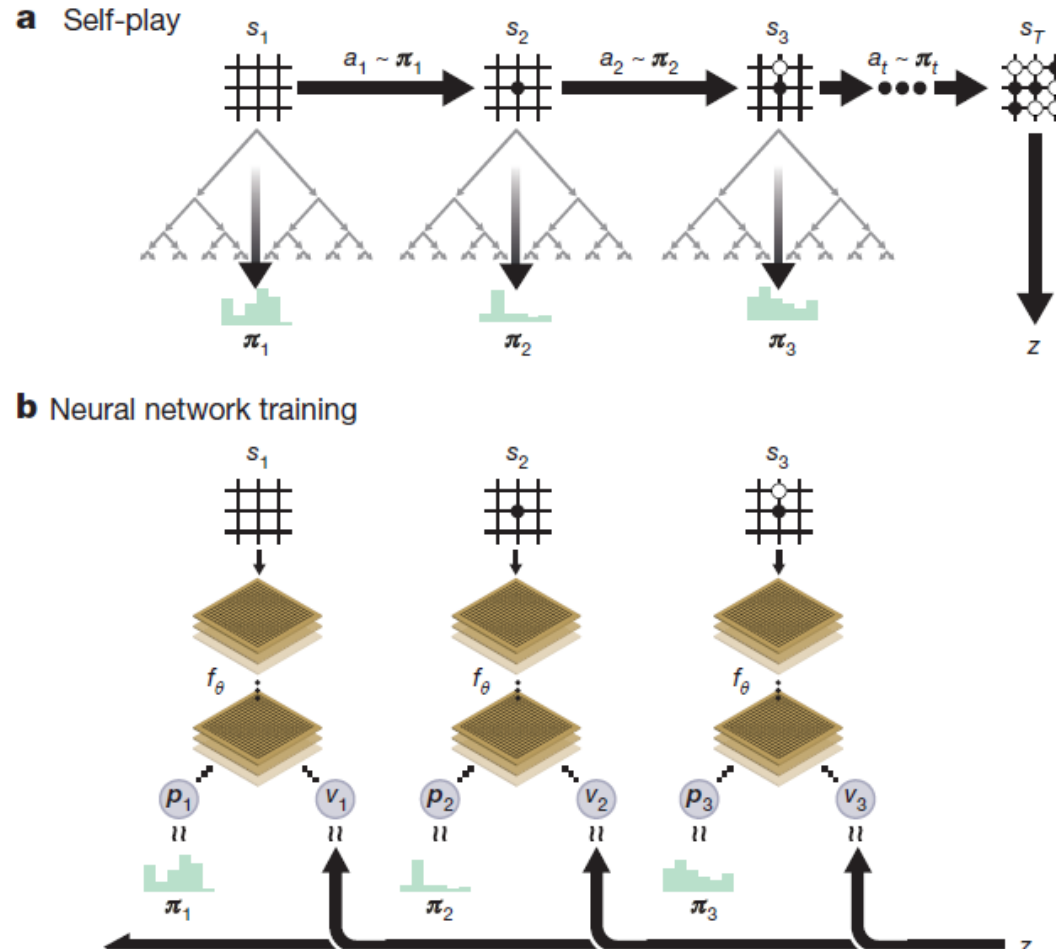
AlphaZero in Gomoku

Hongming Zhang

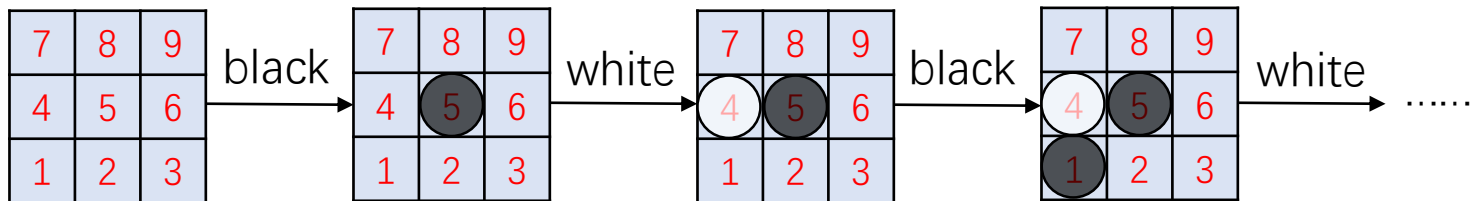
Bo Xu's Workgroup@CASIA

2018.11.30

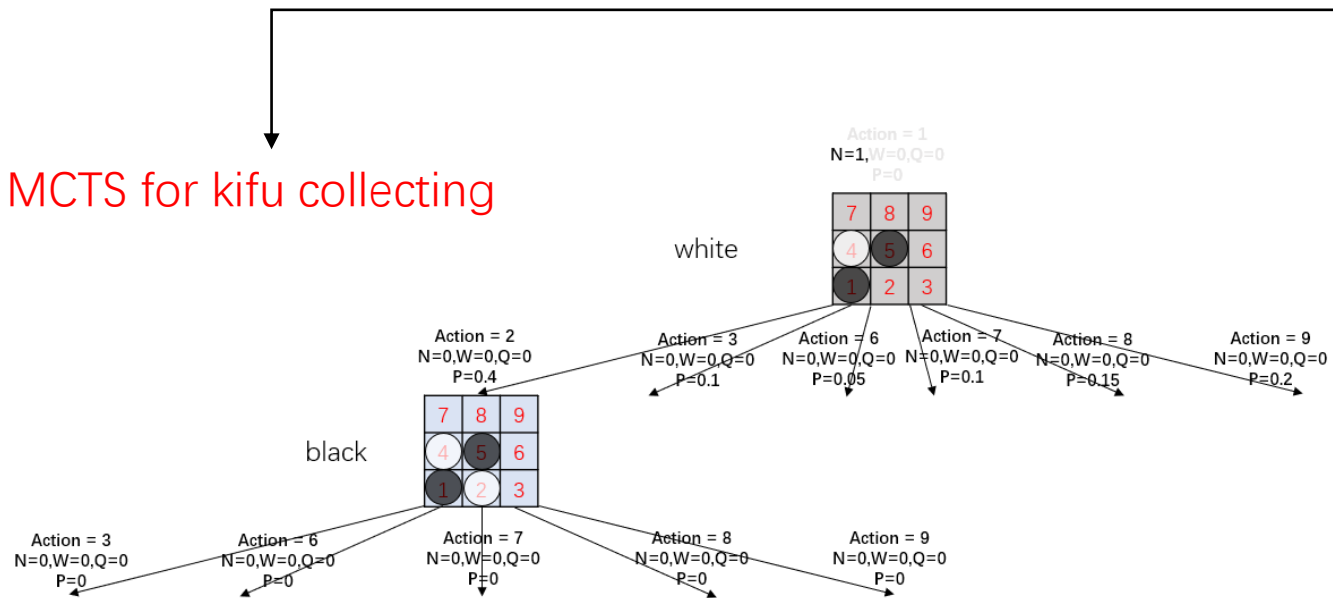
- Self-play reinforcement learning in AlphaGo Zero



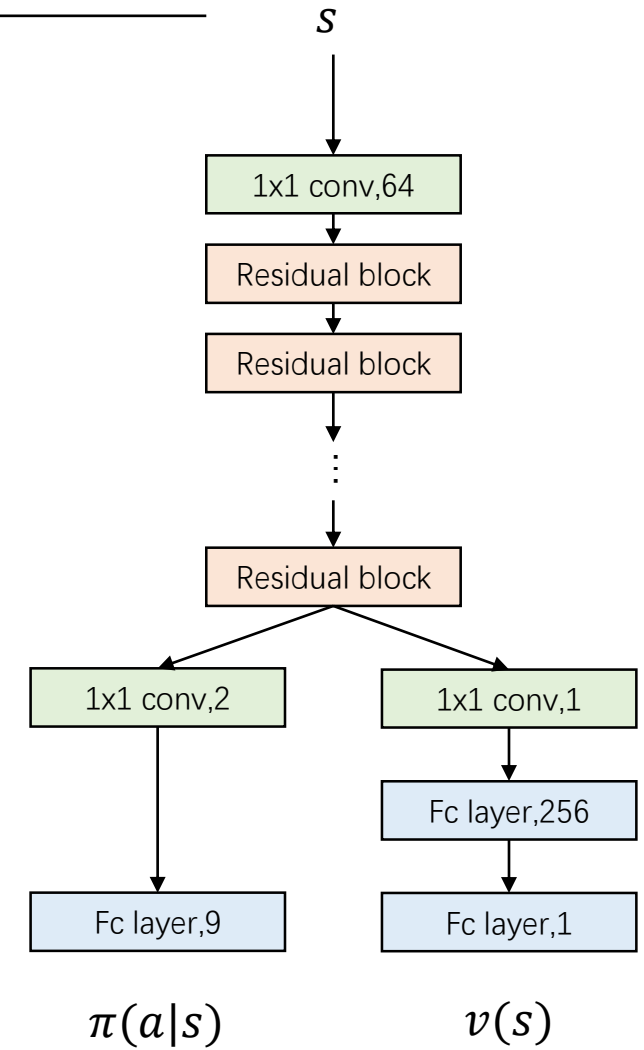
- To simplify, we take the 3x3 Gomoku for instance.
- The board here is only 3x3, the one that get 3 in a row will win the game.
- Numbers in the board denote actions, black and white denote two players.



MCTS for kifu collecting



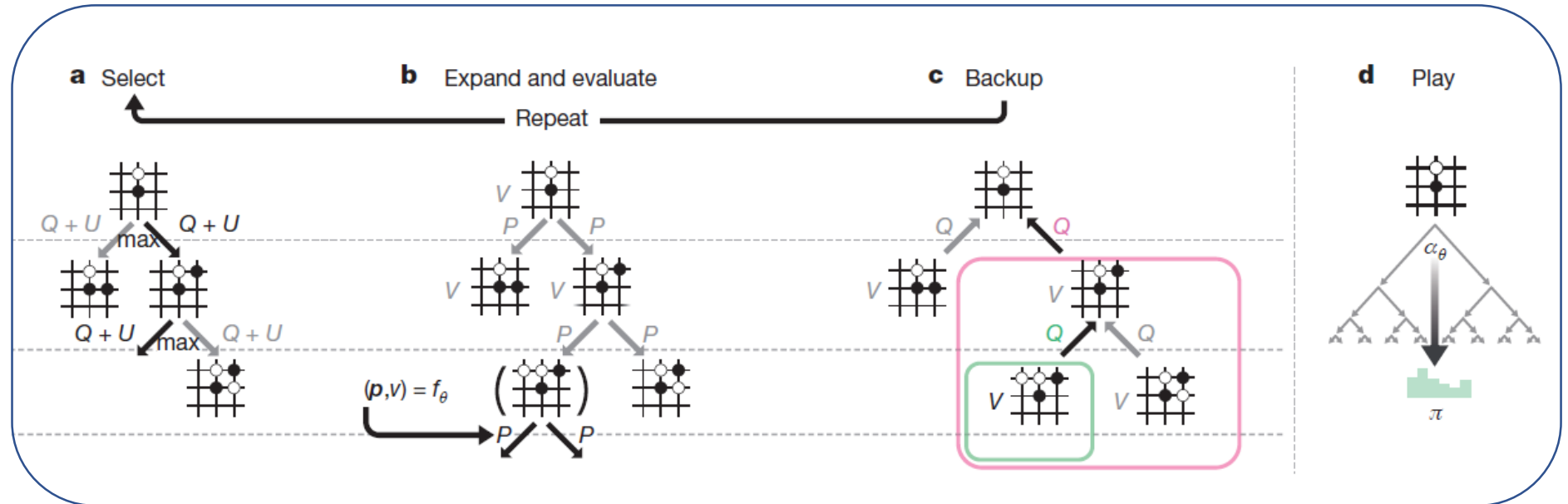
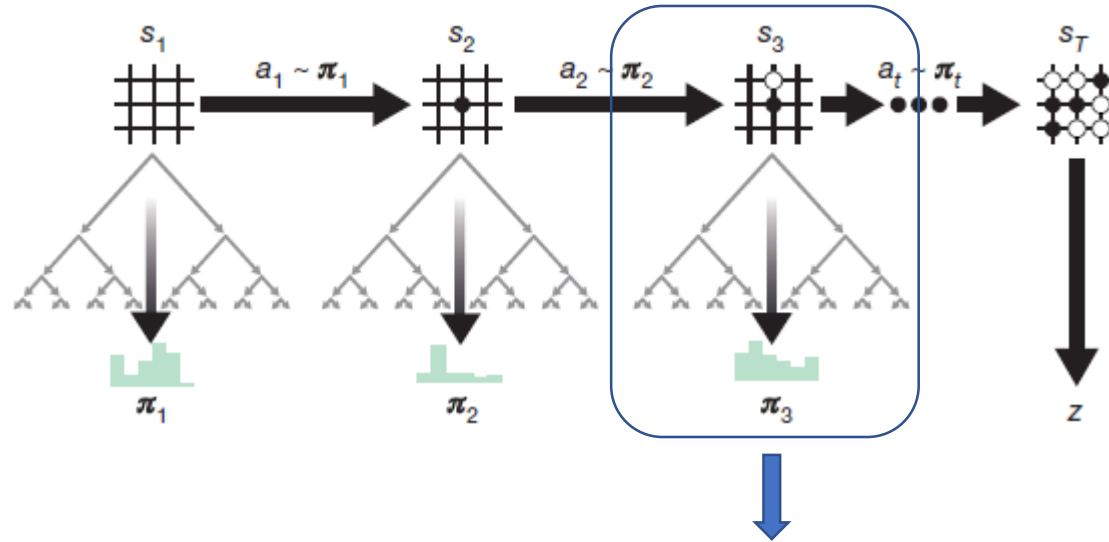
Deep learning for kifu fitting



➤ Algorithm schematic

➤ MCTS in AlphaGo Zero

- Select
- Expand and evaluate
- Backup
- Play



➤ Monte Carlo Tree Search in Gomoku

7	8	9
4	5	6
1	2	3

Real board state

- We assume the game start from this state.
- So the root node in the tree is from here!

Action = 1
N=0,W=0,Q=0
P=0

white

7	8	9
4	5	6
1	2	3

Tree from this node.

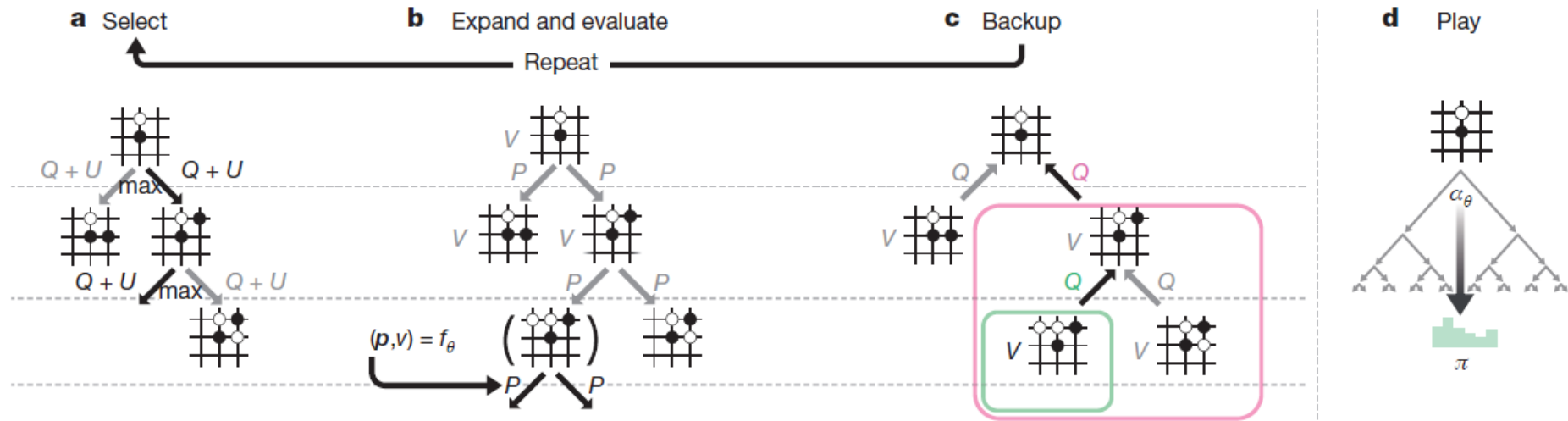
The information here is

- Action : action to arrive this state/node
- N : this state's/node's visit time
- W : this state's/node's total value
- Q : this state's/node's mean value

7	8	9
4	5	6
1	2	3

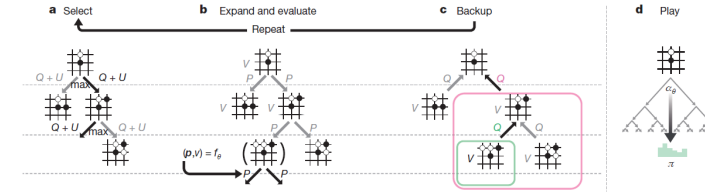
Real board state

Begin!



Select

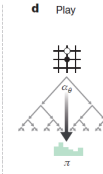
- If the node is in tree, we can select an action to go to the next node/state.
- If the node is leaf node, there is no node to go and we should expand the children nodes.



Action = 1
N=0,W=0,Q=0
P=0

7	8	9
4	5	6
1	2	3

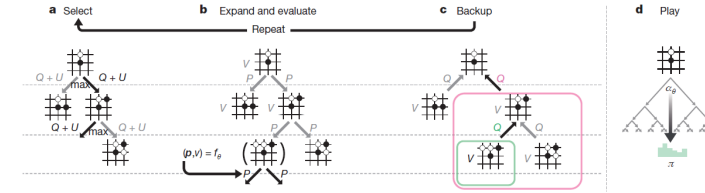
white



7	8	9
4	5	6
1	2	3

Real board state

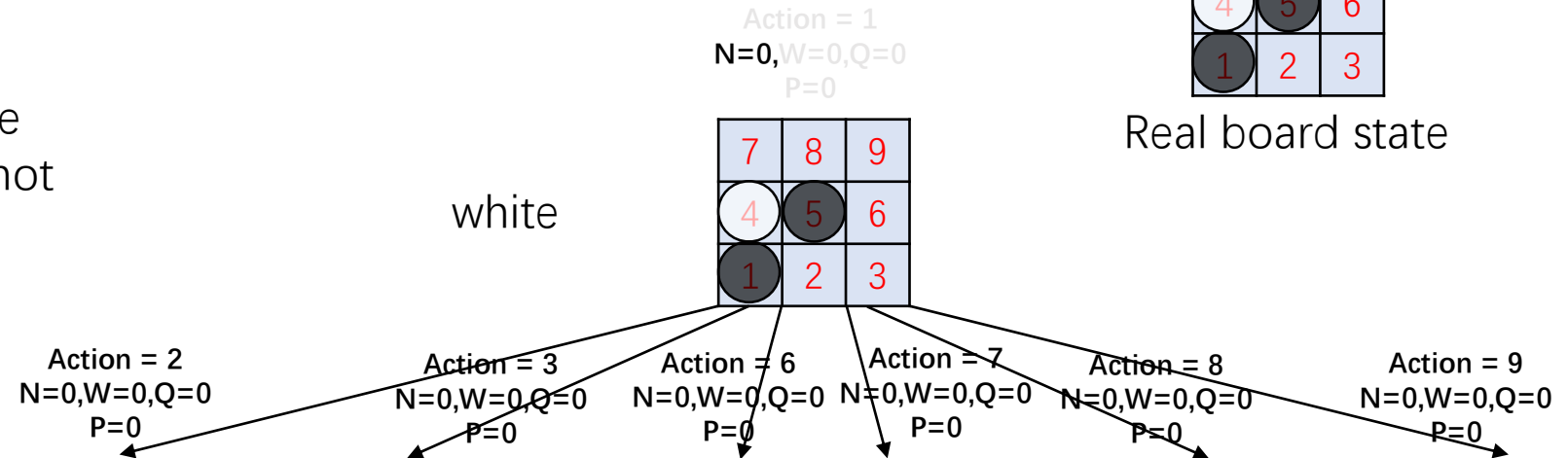
Expand and evaluate



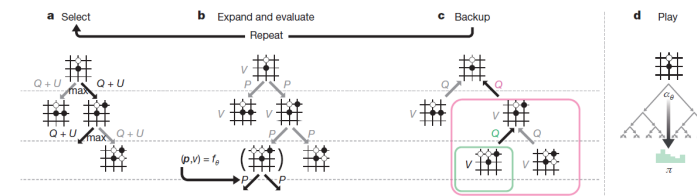
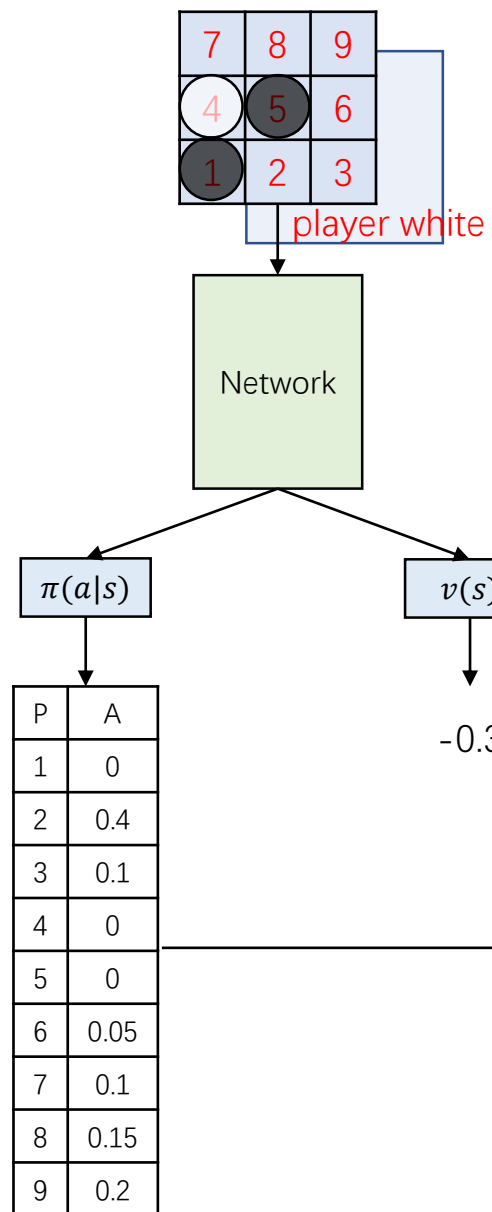
7	8	9
4	5	6
1	2	3

Real board state

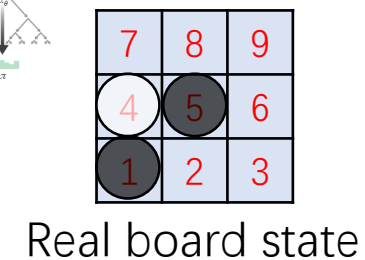
- We assumed this was root node
So the parameters in grey will not be used!



Expand and evaluate



Action = 1
 $N=0, W=0, Q=0$
 $P=0$



white

Action = 2
 $N=0, W=0, Q=0$
 $P=0$

Action = 3
 $N=0, W=0, Q=0$
 $P=0$

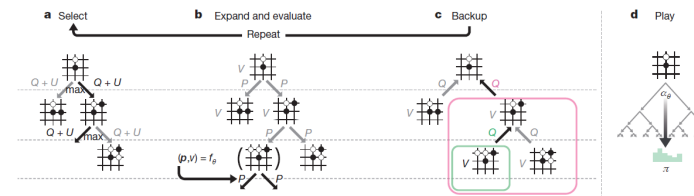
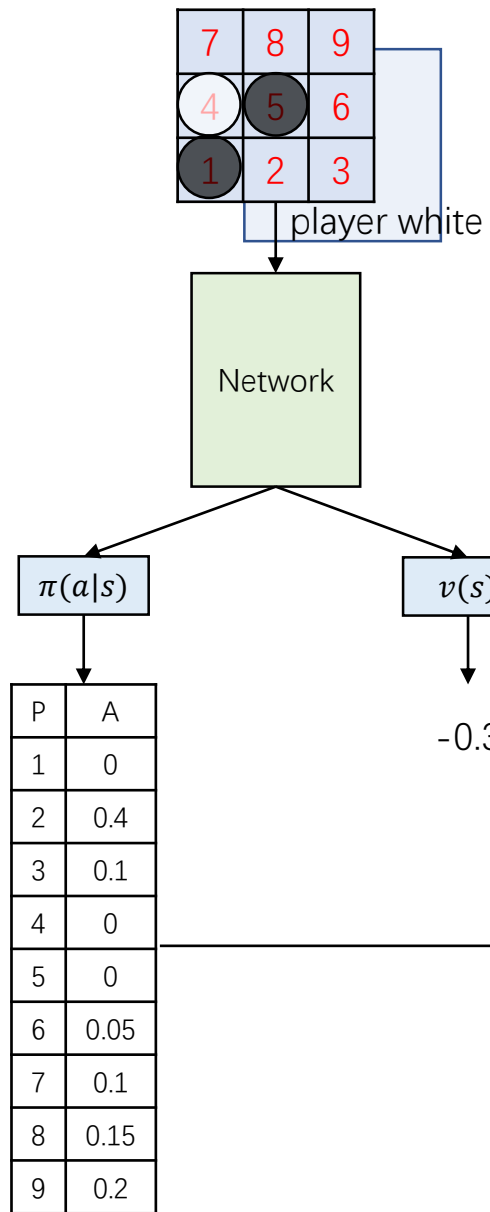
Action = 6
 $N=0, W=0, Q=0$
 $P=0$

Action = 7
 $N=0, W=0, Q=0$
 $P=0$

Action = 8
 $N=0, W=0, Q=0$
 $P=0$

Action = 9
 $N=0, W=0, Q=0$
 $P=0$

Expand and evaluate



Action = 1
 $N=0, W=0, Q=0$
 $P=0$

Real board state

white

Action = 2
 $N=0, W=0, Q=0$
 $P=0.4$

Action = 3
 $N=0, W=0, Q=0$
 $P=0.1$

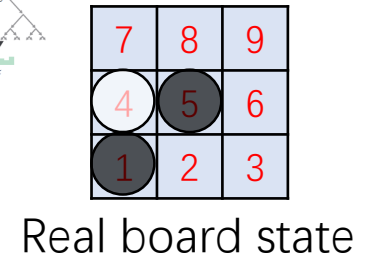
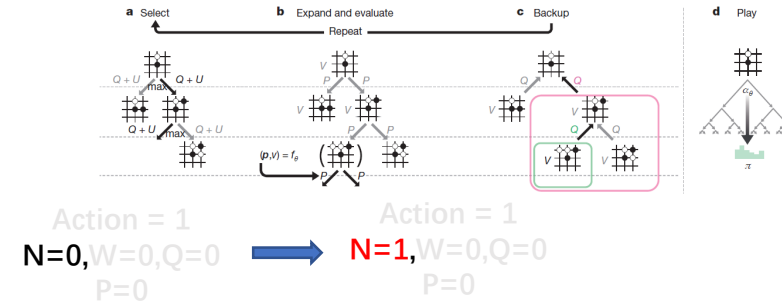
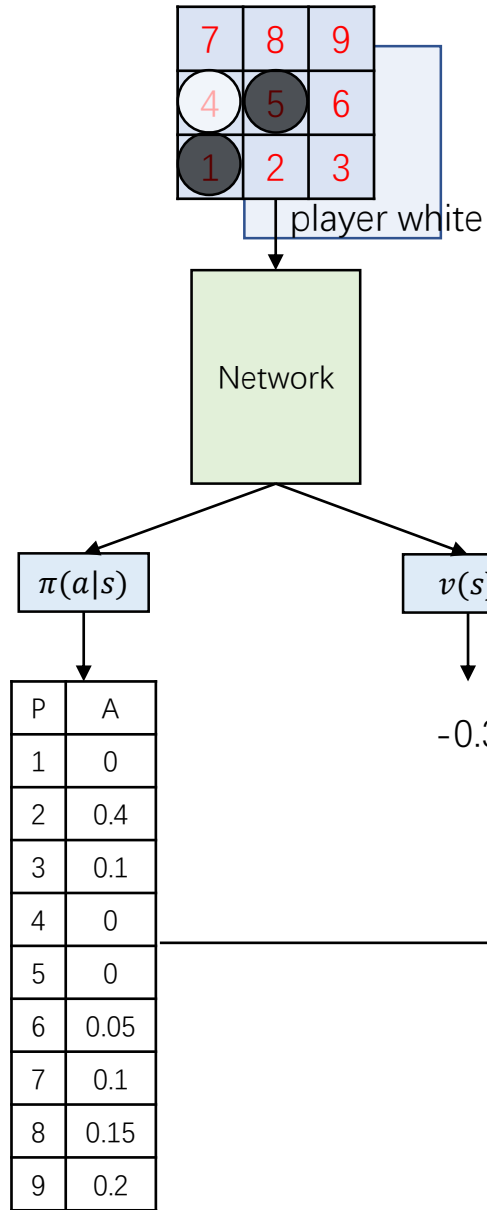
Action = 6
 $N=0, W=0, Q=0$
 $P=0.05$

Action = 7
 $N=0, W=0, Q=0$
 $P=0.1$

Action = 8
 $N=0, W=0, Q=0$
 $P=0.15$

Action = 9
 $N=0, W=0, Q=0$
 $P=0.2$

Backup



Action = 2
N=0, W=0, Q=0
P=0.4

Action = 3
N=0, W=0, Q=0
P=0.1

Action = 6
N=0, W=0, Q=0
P=0.05

Action = 7
N=0, W=0, Q=0
P=0.1

Action = 8
N=0, W=0, Q=0
P=0.15

Action = 9
N=0, W=0, Q=0
P=0.2

$$N(s, a) = N(s, a) + 1$$

$$W(s, a) = W(s, a) + v(s)$$

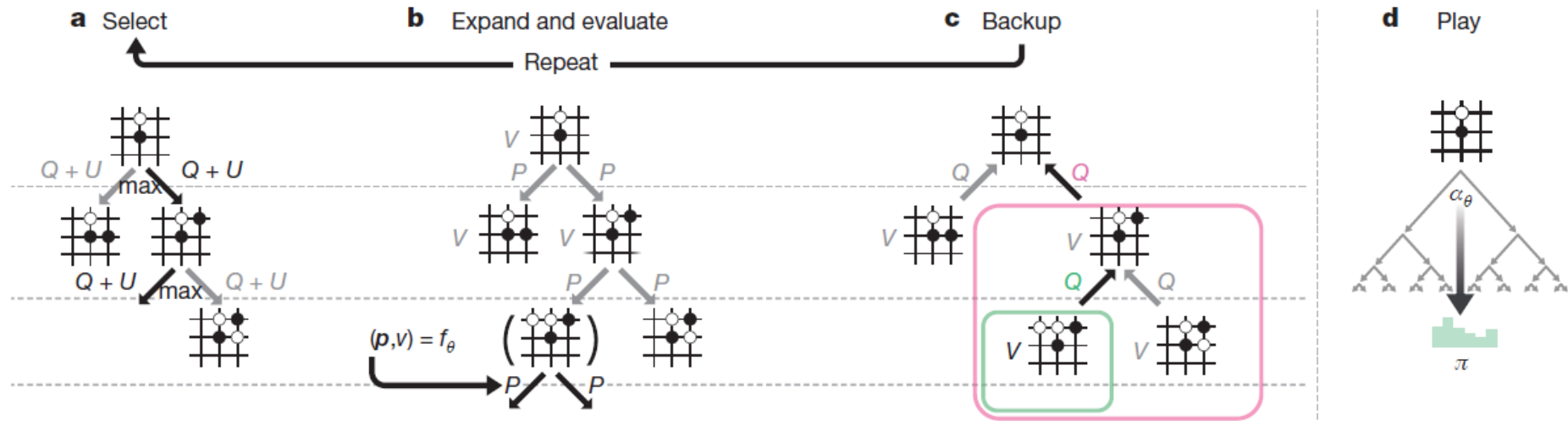
$$Q(s, a) = \frac{W(s, a)}{N(s, a)}$$

- We assumed this was root node before. So no need to backup value here!
- But we need to change the visit time from N to N+1 for later use.

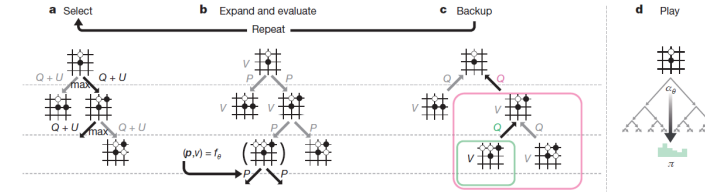
7	8	9
4	5	6
1	2	3

Real board state

Done !
And one more !



Select



7	8	9
4	5	6
1	2	3

Real board state

$$a = \operatorname{argmax}_a (Q(s, a) + U(s, a))$$

$$\text{Exploration : } U(s, a) = c_{puct} P(s, a) \frac{\sqrt{\sum_b N(s, b)}}{1 + N(s, a)}, c_{puct} = 5$$

$$\text{Exploitation : } Q(s, a) = \frac{W}{N}$$

white

Action = 2
 N=0, W=0, Q=0
 P=0.4

Action = 3
 N=0, W=0, Q=0
 P=0.1

Action = 6
 N=0, W=0, Q=0
 P=0.05

Action = 7
 N=0, W=0, Q=0
 P=0.1

Action = 8
 N=0, W=0, Q=0
 P=0.15

Action = 9
 N=0, W=0, Q=0
 P=0.2

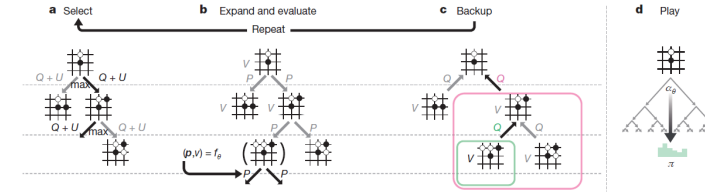
7	8	9
4	5	6
1	2	3

Select

$$a = \operatorname{argmax}_a (Q(s, a) + U(s, a))$$

$$\text{Exploration : } U(s, a) = c_{puct} P(s, a) \frac{\sqrt{\sum_b N(s, b)}}{1 + N(s, a)}, c_{puct} = 5$$

$$\text{Exploitation : } Q(s, a) = \frac{W}{N}$$



7	8	9
4	5	6
1	2	3

Real board state

white

7	8	9
4	5	6
1	2	3

Action = 2
 N=0, W=0, Q=0
 P=0.4

7	8	9
4	5	6
1	2	3

black

Action = 3
 N=0, W=0, Q=0
 P=0.1

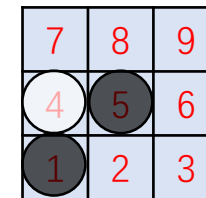
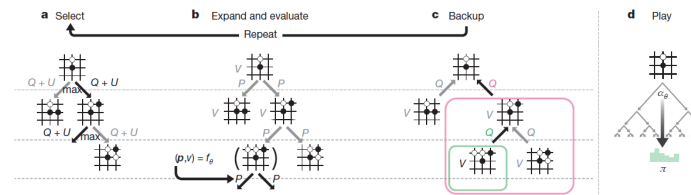
Action = 6
 N=0, W=0, Q=0
 P=0.05

Action = 7
 N=0, W=0, Q=0
 P=0.1

Action = 8
 N=0, W=0, Q=0
 P=0.15

Action = 9
 N=0, W=0, Q=0
 P=0.2

Expand and evaluate



Real board state



white

Action = 2
 $N=0, W=0, Q=0$
 $P=0.4$

Action = 3
 $N=0, W=0, Q=0$
 $P=0.1$

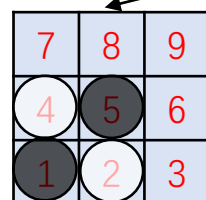
Action = 6
 $N=0, W=0, Q=0$
 $P=0.05$

Action = 7
 $N=0, W=0, Q=0$
 $P=0.1$

Action = 8
 $N=0, W=0, Q=0$
 $P=0.15$

Action = 9
 $N=0, W=0, Q=0$
 $P=0.2$

black



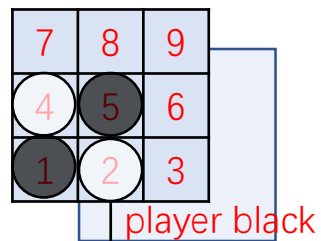
Action = 3
 $N=0, W=0, Q=0$
 $P=0$

Action = 6
 $N=0, W=0, Q=0$
 $P=0$

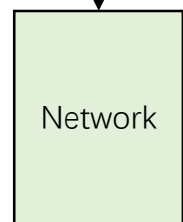
Action = 7
 $N=0, W=0, Q=0$
 $P=0$

Action = 8
 $N=0, W=0, Q=0$
 $P=0$

Action = 9
 $N=0, W=0, Q=0$
 $P=0$



player black



Network

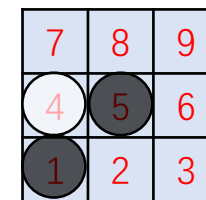
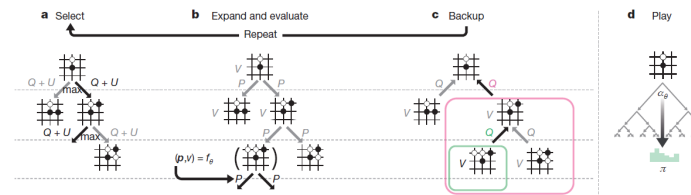
$v(s)$

-0.1

$\pi(a|s)$

P	A
1	0
2	0
3	0.2
4	0
5	0
6	0.15
7	0.1
8	0.25
9	0.3

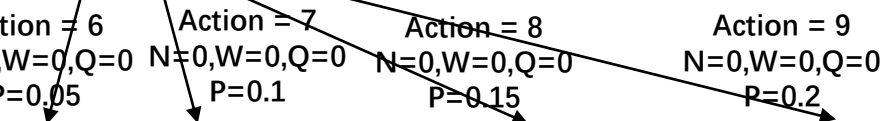
Expand and evaluate



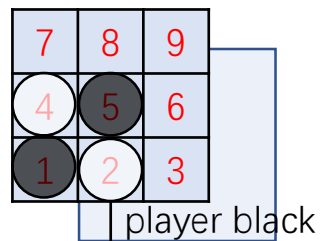
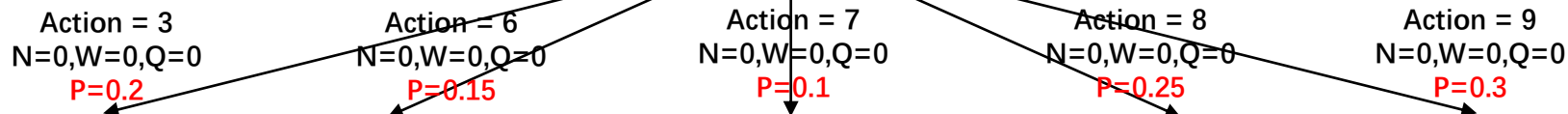
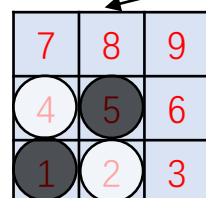
Real board state



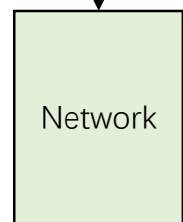
white



black



player black

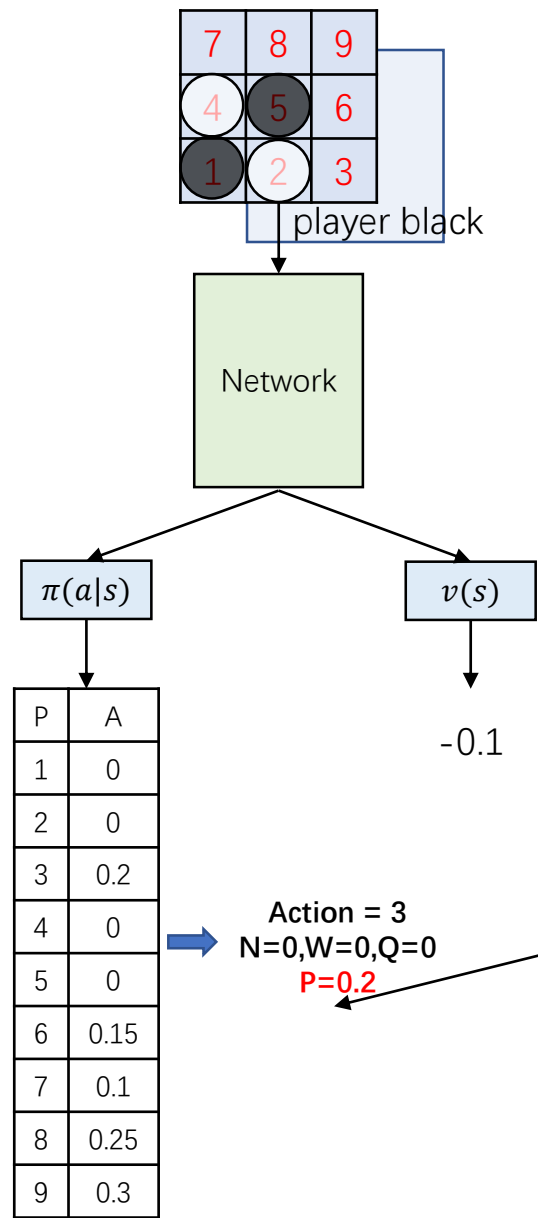


Network

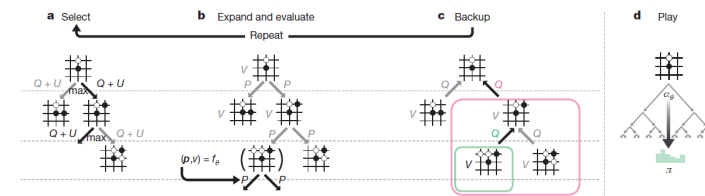


-0.1

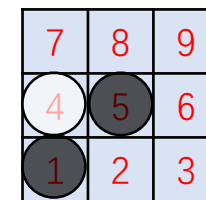
P	A
1	0
2	0
3	0.2
4	0
5	0
6	0.15
7	0.1
8	0.25
9	0.3



Backup



Action = 1
 $N=1, W=0, Q=0$
 $P=0$



Real board state

white



Action = 2
 $N=0, W=0, Q=0$
 $P=0.4$

Action = 2
 $N=1, W=0.1, Q=0.1$
 $P=0.4$

Action = 3
 $N=0, W=0, Q=0$
 $P=0.1$

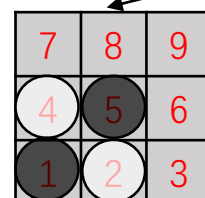
Action = 6
 $N=0, W=0, Q=0$
 $P=0.05$

Action = 7
 $N=0, W=0, Q=0$
 $P=0.1$

Action = 8
 $N=0, W=0, Q=0$
 $P=0.15$

Action = 9
 $N=0, W=0, Q=0$
 $P=0.2$

black



Action = 3
 $N=0, W=0, Q=0$
 $P=0.2$

Action = 6
 $N=0, W=0, Q=0$
 $P=0.15$

Action = 7
 $N=0, W=0, Q=0$
 $P=0.1$

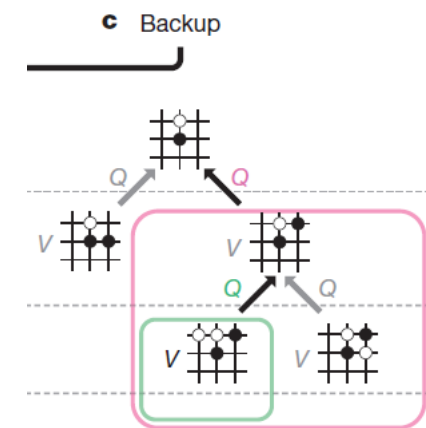
Action = 8
 $N=0, W=0, Q=0$
 $P=0.25$

Action = 9
 $N=0, W=0, Q=0$
 $P=0.3$

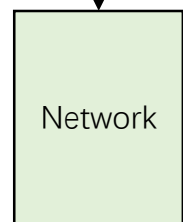
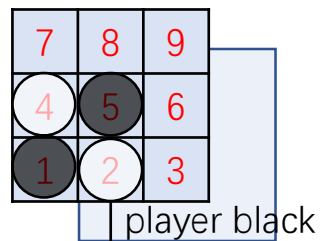
$$N(s, a) = N(s, a) + 1$$

$$W(s, a) = W(s, a) + v(s')$$

$$Q(s, a) = \frac{W(s, a)}{N(s, a)}$$



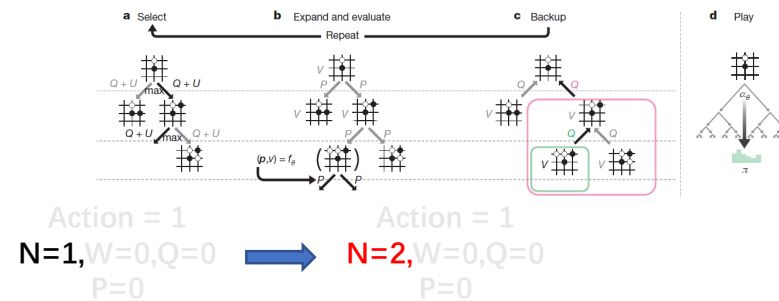
switching perspective
 $v_{white} = -v_{black} = 0.1$


 $\pi(a|s)$
 $v(s)$

-0.1

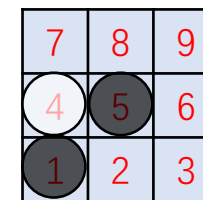
P	A
1	0
2	0
3	0.2
4	0
5	0
6	0.15
7	0.1
8	0.25
9	0.3

Backup



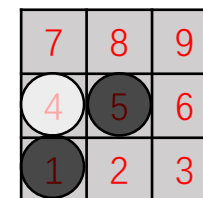
Action = 1
 $N=1, W=0, Q=0$
 $P=0$

Action = 1
 $N=2, W=0, Q=0$
 $P=0$



Real board state

white



Action = 2
 $N=1, W=0.1, Q=0.1$
 $P=0.4$

Action = 3
 $N=0, W=0, Q=0$
 $P=0.1$

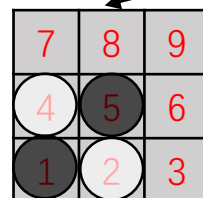
Action = 6
 $N=0, W=0, Q=0$
 $P=0.05$

Action = 7
 $N=0, W=0, Q=0$
 $P=0.1$

Action = 8
 $N=0, W=0, Q=0$
 $P=0.15$

Action = 9
 $N=0, W=0, Q=0$
 $P=0.2$

black



Action = 3
 $N=0, W=0, Q=0$
 $P=0.2$

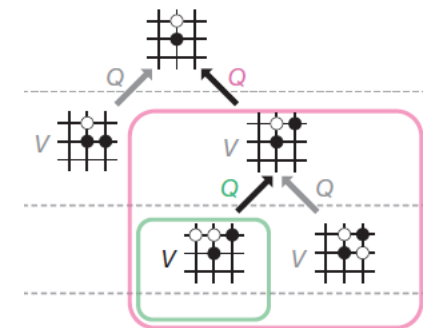
Action = 6
 $N=0, W=0, Q=0$
 $P=0.15$

Action = 7
 $N=0, W=0, Q=0$
 $P=0.1$

Action = 8
 $N=0, W=0, Q=0$
 $P=0.25$

Action = 9
 $N=0, W=0, Q=0$
 $P=0.3$

c Backup



switching perspective

 $v_{white} = -v_{black} = 0.1$

$$N(s, a) = N(s, a) + 1$$

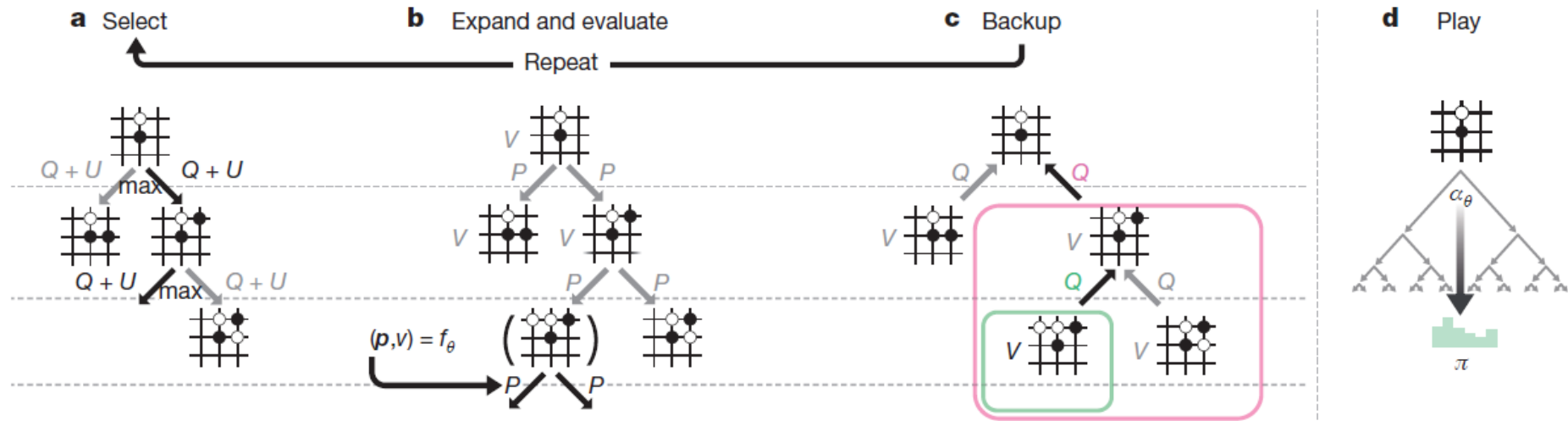
$$W(s, a) = W(s, a) + v(s')$$

$$Q(s, a) = \frac{W(s, a)}{N(s, a)}$$

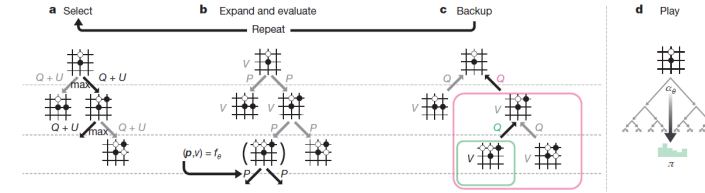
7	8	9
4	5	6
1	2	3

Real board state

Done !
And one more !



Select



Action = 1
N=2, W=0, Q=0
P=0

7	8	9
4	5	6
1	2	3

Real board state

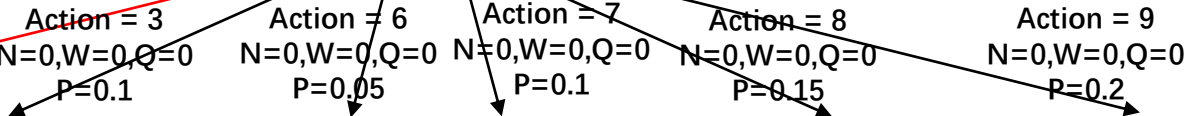
$$a = \operatorname{argmax}_a (Q(s, a) + U(s, a))$$

Exploration : $U(s, a) = c_{puct} P(s, a) \frac{\sqrt{\sum_b N(s, b)}}{1 + N(s, a)}$, $c_{puct} = 5$

Exploitation : $Q(s, a) = \frac{W}{N}$

white

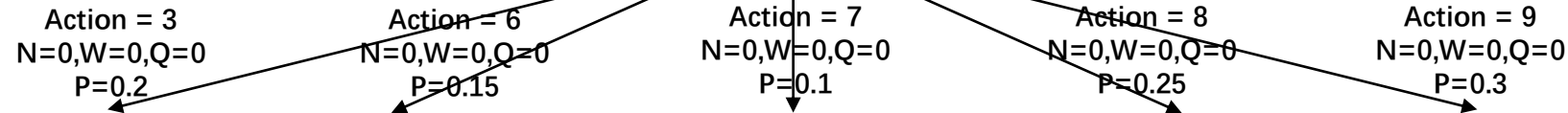
7	8	9
4	5	6
1	2	3



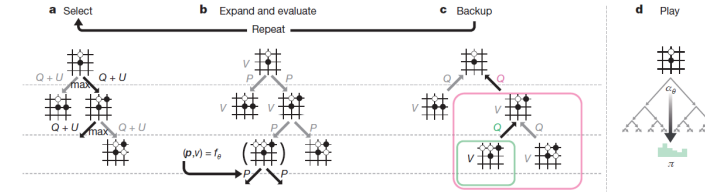
black

Action = 2
N=1, W=0.1, Q=0.1
P=0.4

7	8	9
4	5	6
1	2	3



Select



Action = 1
N=2, W=0, Q=0
P=0

7	8	9
4	5	6
1	2	3

Real board state

$$a = \operatorname{argmax}_a (Q(s, a) + U(s, a))$$

Exploration : $U(s, a) = c_{puct} P(s, a) \frac{\sqrt{\sum_b N(s, b)}}{1 + N(s, a)}$, $c_{puct} = 5$

Exploitation : $Q(s, a) = \frac{W}{N}$

white

7	8	9
4	5	6
1	2	3

Action = 2
N=1, W=0.1, Q=0.1
P=0.4

Action = 3
N=0, W=0, Q=0
P=0.1

Action = 6
N=0, W=0, Q=0
P=0.05

Action = 7
N=0, W=0, Q=0
P=0.1

Action = 8
N=0, W=0, Q=0
P=0.15

Action = 9
N=0, W=0, Q=0
P=0.2

black

7	8	9
4	5	6
1	2	3

Action = 3
N=0, W=0, Q=0
P=0.2

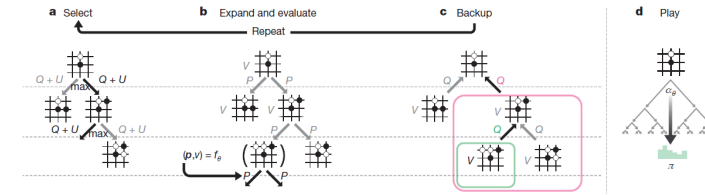
Action = 6
N=0, W=0, Q=0
P=0.15

Action = 7
N=0, W=0, Q=0
P=0.1

Action = 8
N=0, W=0, Q=0
P=0.25

Action = 9
N=0, W=0, Q=0
P=0.3

Expand and evaluate



7	8	9
4	5	6
1	2	3

Real board state

white

7	8	9
4	5	6
1	2	3

Action = 2
 $N=1, W=0.1, Q=0.1$
 $P=0.4$

Action = 3
 $N=0, W=0, Q=0$
 $P=0.1$

Action = 6
 $N=0, W=0, Q=0$
 $P=0.05$

Action = 7
 $N=0, W=0, Q=0$
 $P=0.1$

Action = 8
 $N=0, W=0, Q=0$
 $P=0.15$

Action = 9
 $N=0, W=0, Q=0$
 $P=0.2$

black

7	8	9
4	5	6
1	2	3

Action = 3
 $N=0, W=0, Q=0$
 $P=0.2$

Action = 6
 $N=0, W=0, Q=0$
 $P=0.15$

Action = 7
 $N=0, W=0, Q=0$
 $P=0.1$

Action = 8
 $N=0, W=0, Q=0$
 $P=0.25$

Action = 9
 $N=0, W=0, Q=0$
 $P=0.3$

white

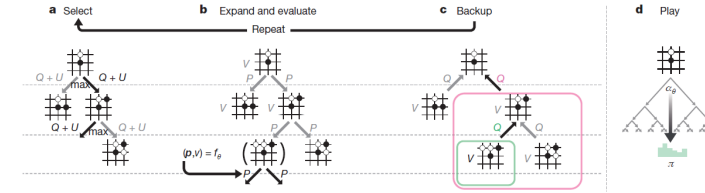
7	8	9
4	5	6
1	2	3

Terminal and
 white lose !

Reward = -1

- Note that we always get the reward from the current perspective.

Backup



7	8	9
4	5	6
1	2	3

Real board state

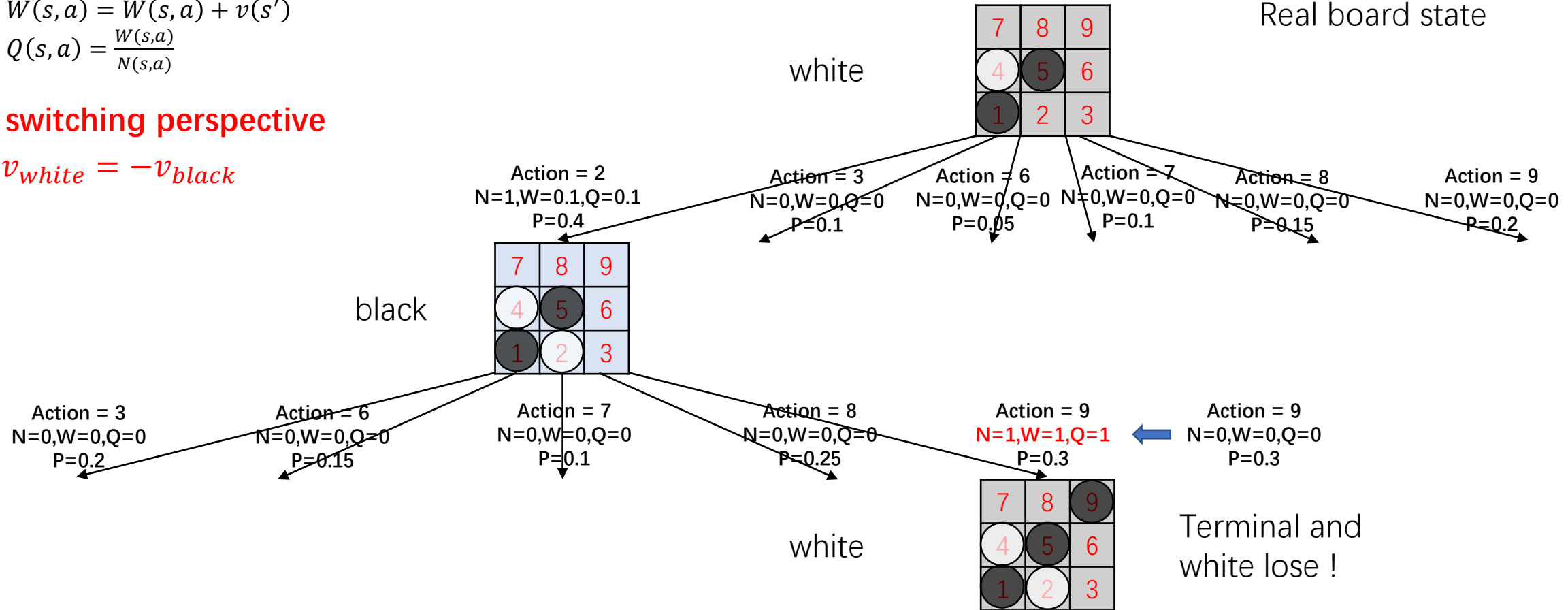
$$N(s, a) = N(s, a) + 1$$

$$W(s, a) = W(s, a) + v(s')$$

$$Q(s, a) = \frac{W(s, a)}{N(s, a)}$$

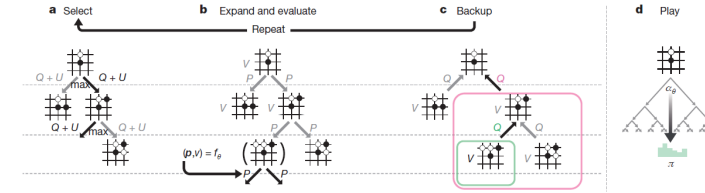
switching perspective

$$v_{white} = -v_{black}$$



- Now we have the real value so don't need a predictive value from network. \rightarrow Reward = -1

Backup



7	8	9
4	5	6
1	2	3

Real board state

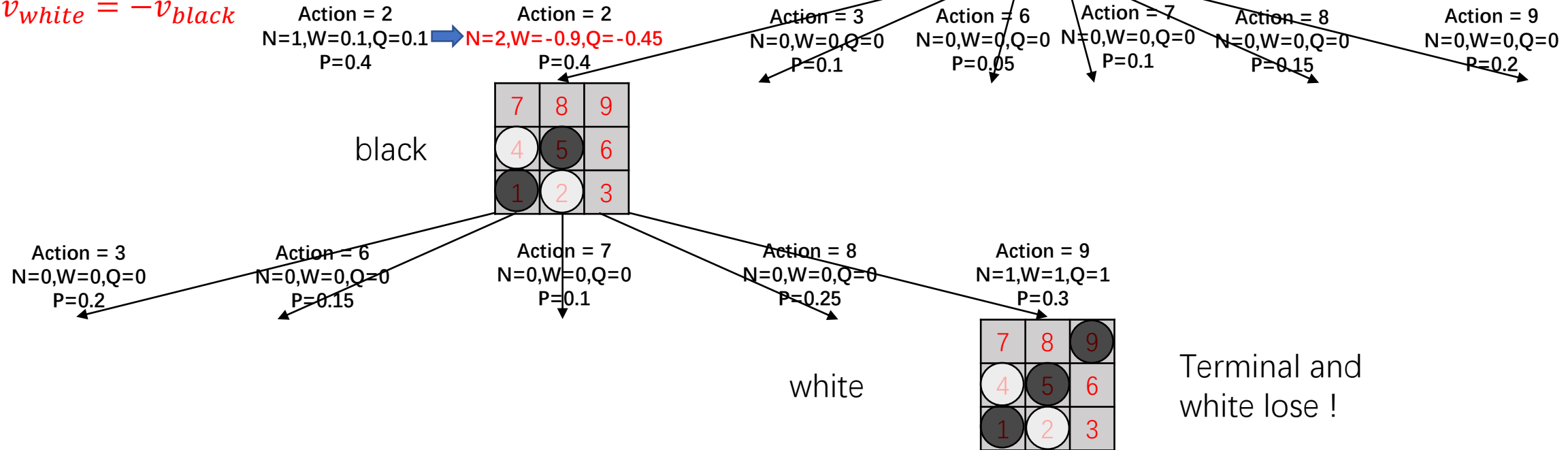
$$N(s, a) = N(s, a) + 1$$

$$W(s, a) = W(s, a) + v(s')$$

$$Q(s, a) = \frac{W(s, a)}{N(s, a)}$$

switching perspective

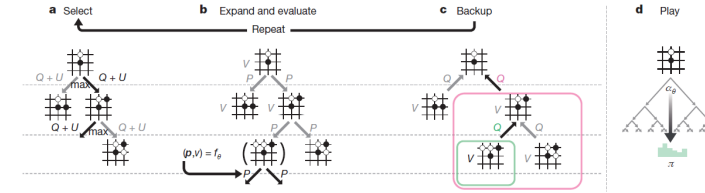
$$v_{white} = -v_{black}$$



- Now we have the real value so don't need a predictive value from network.

Reward = -1

Backup



7	8	9
4	5	6
1	2	3

Real board state

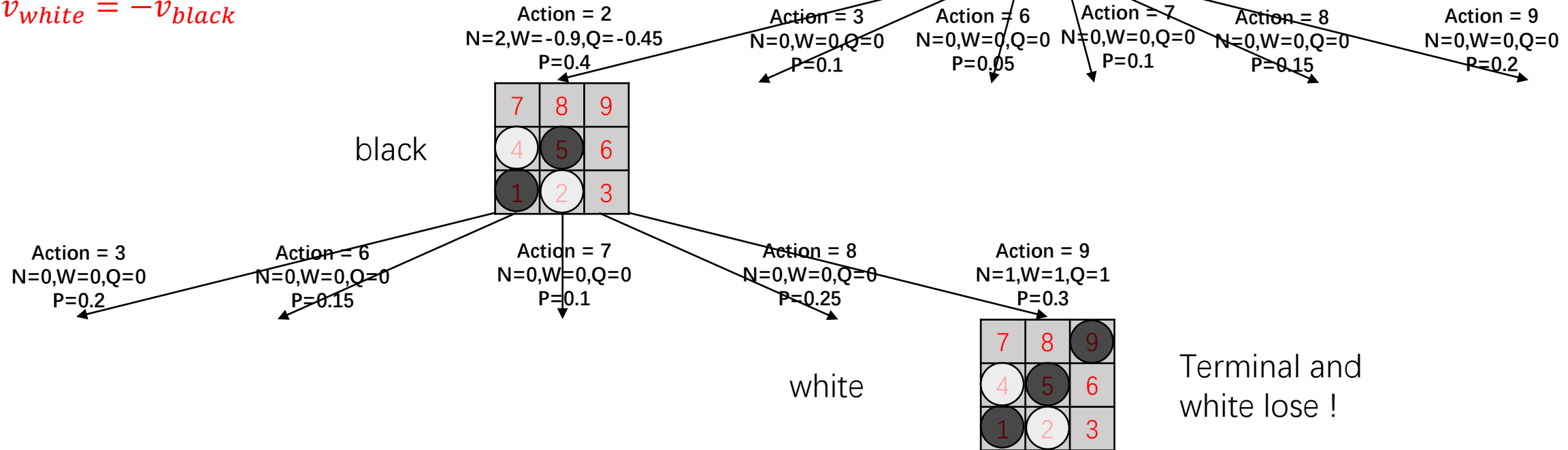
$$N(s, a) = N(s, a) + 1$$

$$W(s, a) = W(s, a) + v(s')$$

$$Q(s, a) = \frac{W(s, a)}{N(s, a)}$$

switching perspective

$$v_{white} = -v_{black}$$

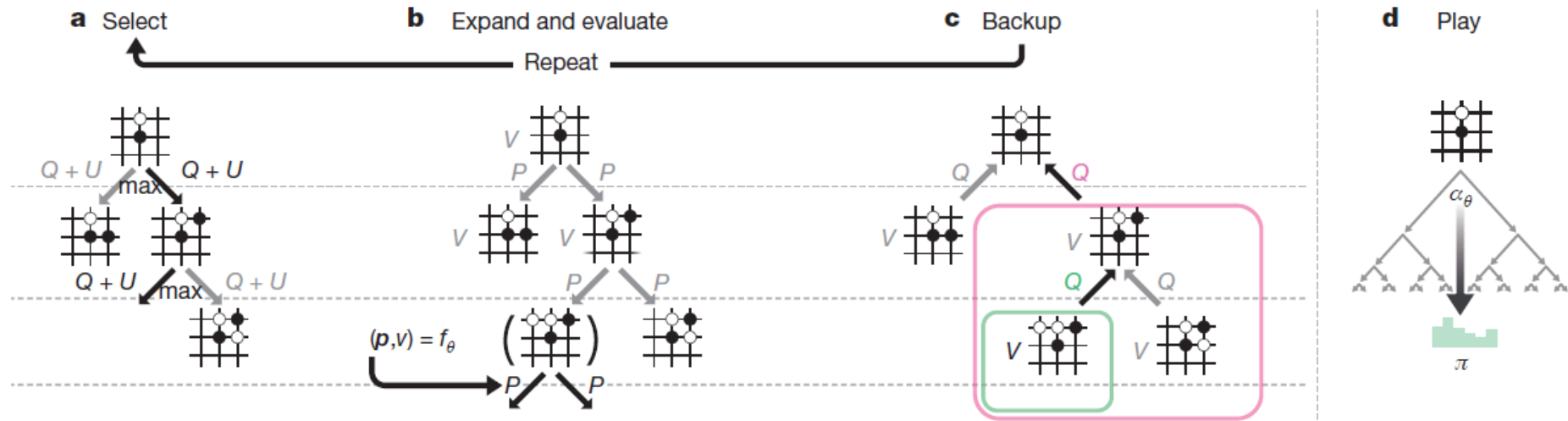


- Now we have the real value so don't need a predictive value from network.

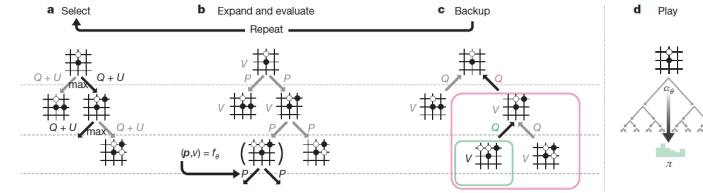
7	8	9
4	5	6
1	2	3

Real board state

Done !
And one more !



Select



Action = 1
N=3, W=0, Q=0
P=0

7	8	9
4	5	6
1	2	3

Real board state

$$a = \operatorname{argmax}_a (Q(s, a) + U(s, a))$$

Exploration : $U(s, a) = c_{puct} P(s, a) \frac{\sqrt{\sum_b N(s, b)}}{1 + N(s, a)}$, $c_{puct} = 5$

Exploitation : $Q(s, a) = \frac{W}{N}$

white

7	8	9
4	5	6
1	2	3

Action = 2
N=2, W=-0.9, Q=-0.45
P=0.4

Action = 3
N=0, W=0, Q=0
P=0.1

Action = 6
N=0, W=0, Q=0
P=0.05

Action = 7
N=0, W=0, Q=0
P=0.1

Action = 8
N=0, W=0, Q=0
P=0.15

Action = 9
N=0, W=0, Q=0
P=0.2

black

7	8	9
4	5	6
1	2	3

Action = 3
N=0, W=0, Q=0
P=0.2

Action = 6
N=0, W=0, Q=0
P=0.15

Action = 7
N=0, W=0, Q=0
P=0.1

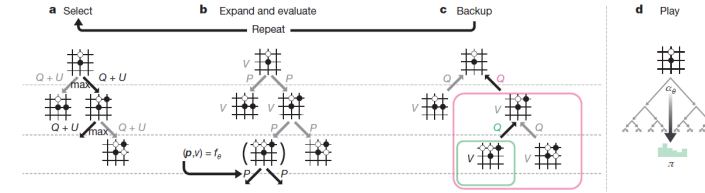
Action = 8
N=0, W=0, Q=0
P=0.25

Action = 9
N=1, W=1, Q=1
P=0.3

white

7	8	9
4	5	6
1	2	3

Select



$$a = \operatorname{argmax}_a (Q(s, a) + U(s, a))$$

$$\text{Exploration : } U(s, a) = c_{puct} P(s, a) \frac{\sqrt{\sum_b N(s, b)}}{1 + N(s, a)}, c_{puct} = 5$$

$$\text{Exploitation : } Q(s, a) = \frac{W}{N}$$

Real board state

7	8	9
4	5	6
1	2	3

white

7	8	9
4	5	6
1	2	3

Action = 2
N=2, W=-0.9, Q=-0.45
P=0.4

7	8	9
4	5	6
1	2	3

black

Action = 3
N=0, W=0, Q=0
P=0.1

Action = 6
N=0, W=0, Q=0
P=0.05

Action = 7
N=0, W=0, Q=0
P=0.1

Action = 8
N=0, W=0, Q=0
P=0.15

Action = 9
N=0, W=0, Q=0
P=0.2

7	8	9
4	5	6
1	2	3

Action = 3
N=0, W=0, Q=0
P=0.2

Action = 6
N=0, W=0, Q=0
P=0.15

Action = 7
N=0, W=0, Q=0
P=0.1

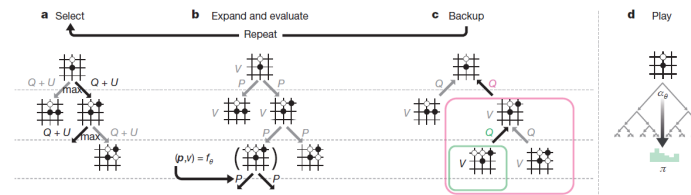
Action = 8
N=0, W=0, Q=0
P=0.25

Action = 9
N=1, W=1, Q=1
P=0.3

7	8	9
4	5	6
1	2	3

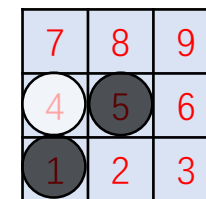
white

Expand and evaluate



Action = 1
 $N=3, W=0, Q=0$
 $P=0$

Real board state



white



Action = 2
 $N=2, W=-0.9, Q=-0.45$
 $P=0.4$

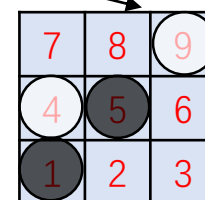
Action = 3
 $N=0, W=0, Q=0$
 $P=0.1$

Action = 6
 $N=0, W=0, Q=0$
 $P=0.05$

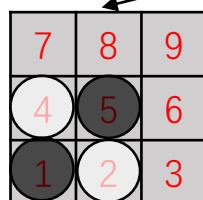
Action = 7
 $N=0, W=0, Q=0$
 $P=0.1$

Action = 8
 $N=0, W=0, Q=0$
 $P=0.15$

Action = 9
 $N=0, W=0, Q=0$
 $P=0.2$



black



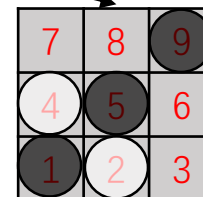
Action = 3
 $N=0, W=0, Q=0$
 $P=0.2$

Action = 6
 $N=0, W=0, Q=0$
 $P=0.15$

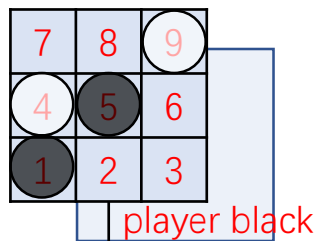
Action = 7
 $N=0, W=0, Q=0$
 $P=0.1$

Action = 8
 $N=0, W=0, Q=0$
 $P=0.25$

Action = 9
 $N=1, W=1, Q=1$
 $P=0.3$



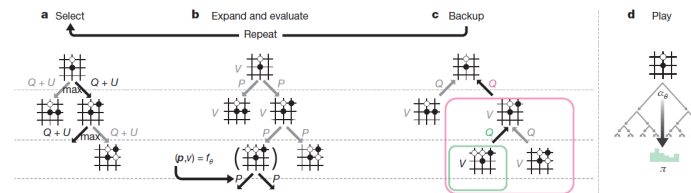
white


 $\pi(a|s)$
 $v(s)$

-0.2

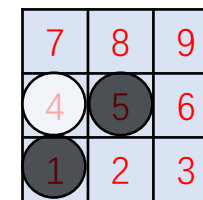
P	A
1	0
2	0.25
3	0.2
4	0
5	0
6	0.1
7	0.2
8	0.25
9	0

Expand and evaluate

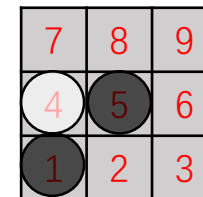


Action = 1
 $N=3, W=0, Q=0$
 $P=0$

Real board state



white



Action = 2
 $N=2, W=-0.9, Q=-0.45$
 $P=0.4$

Action = 3
 $N=0, W=0, Q=0$
 $P=0.1$

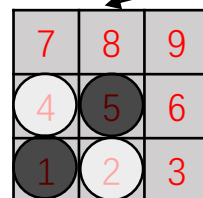
Action = 6
 $N=0, W=0, Q=0$
 $P=0.05$

Action = 7
 $N=0, W=0, Q=0$
 $P=0.1$

Action = 8
 $N=0, W=0, Q=0$
 $P=0.15$

Action = 9
 $N=0, W=0, Q=0$
 $P=0.2$

black



Action = 3
 $N=0, W=0, Q=0$
 $P=0.2$

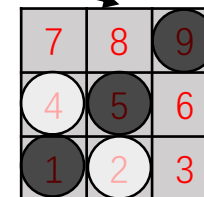
Action = 6
 $N=0, W=0, Q=0$
 $P=0.15$

Action = 7
 $N=0, W=0, Q=0$
 $P=0.1$

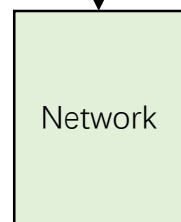
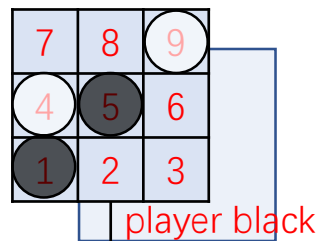
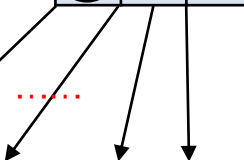
Action = 8
 $N=0, W=0, Q=0$
 $P=0.25$

Action = 9
 $N=1, W=1, Q=1$
 $P=0.3$

white



Action = 2
 $N=0, W=0, Q=0$
 $P=0.25$



Network

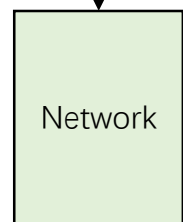
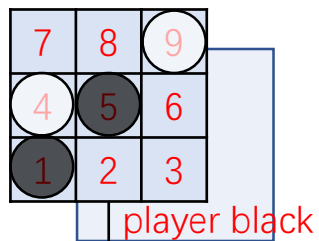
player black

 $v(s)$

-0.2

 $\pi(a|s)$

P	A
1	0
2	0.25
3	0.2
4	0
5	0
6	0.1
7	0.2
8	0.25
9	0



Network

 $\pi(a|s)$ $v(s)$

-0.2

P	A
1	0
2	0.25
3	0.2
4	0
5	0
6	0.1
7	0.2
8	0.25
9	0

Backup

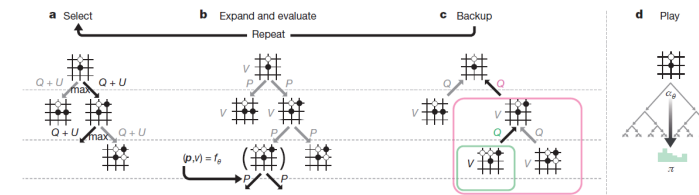
$$N(s, a) = N(s, a) + 1$$

$$W(s, a) = W(s, a) + v(s')$$

$$Q(s, a) = \frac{W(s, a)}{N(s, a)}$$

switching perspective

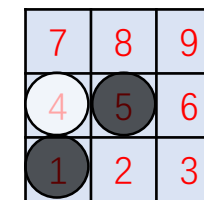
$$v_{white} = -v_{black}$$



Action = 1
 $N=3, W=0, Q=0$
 $P=0$

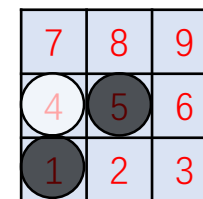
→

Action = 1
 $N=4, W=0, Q=0$
 $P=0$



Real board state

white



Action = 3

 $N=0, W=0, Q=0$
 $P=0.1$

Action = 6

 $N=0, W=0, Q=0$
 $P=0.05$

Action = 7

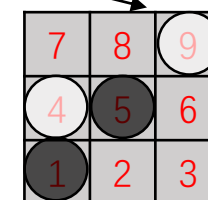
 $N=0, W=0, Q=0$
 $P=0.1$

Action = 8

 $N=0, W=0, Q=0$
 $P=0.15$

Action = 9
 $N=0, W=0, Q=0$
 $P=0.2$

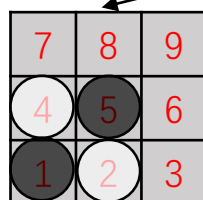
Action = 9

 $N=1, W=0.2, Q=0.2$
 $P=0.2$


Action = 2
 $N=0, W=0, Q=0$
 $P=0.25$

...

black



Action = 2

 $N=2, W=-0.9, Q=-0.45$
 $P=0.4$

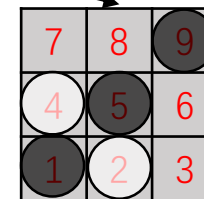
Action = 3
 $N=0, W=0, Q=0$
 $P=0.2$

Action = 6
 $N=0, W=0, Q=0$
 $P=0.15$

Action = 7
 $N=0, W=0, Q=0$
 $P=0.1$

Action = 8
 $N=0, W=0, Q=0$
 $P=0.25$

Action = 9
 $N=1, W=1, Q=1$
 $P=0.3$

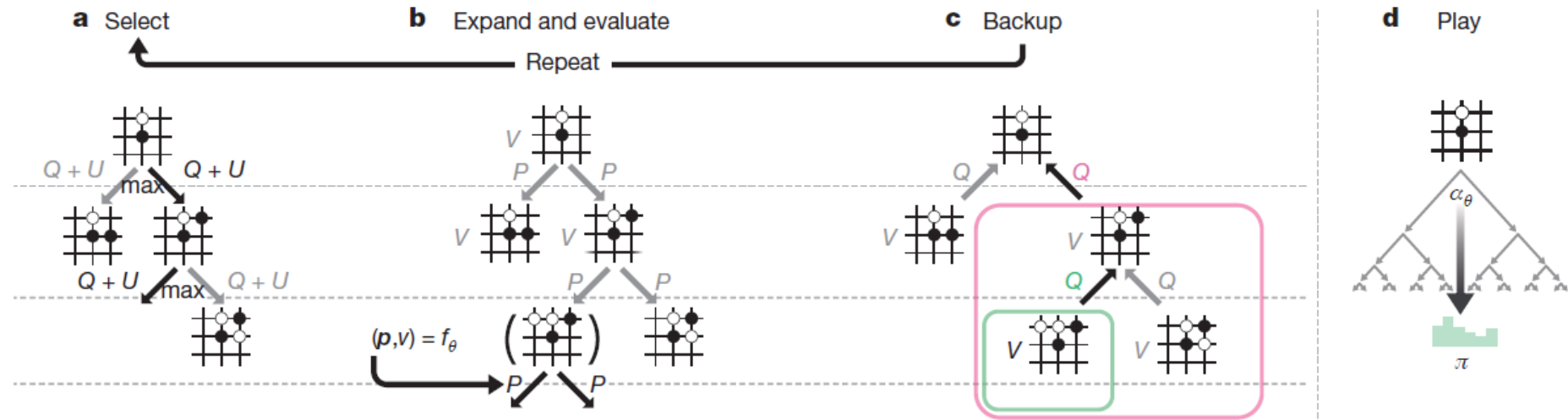


white

7	8	9
4	5	6
1	2	3

Real board state

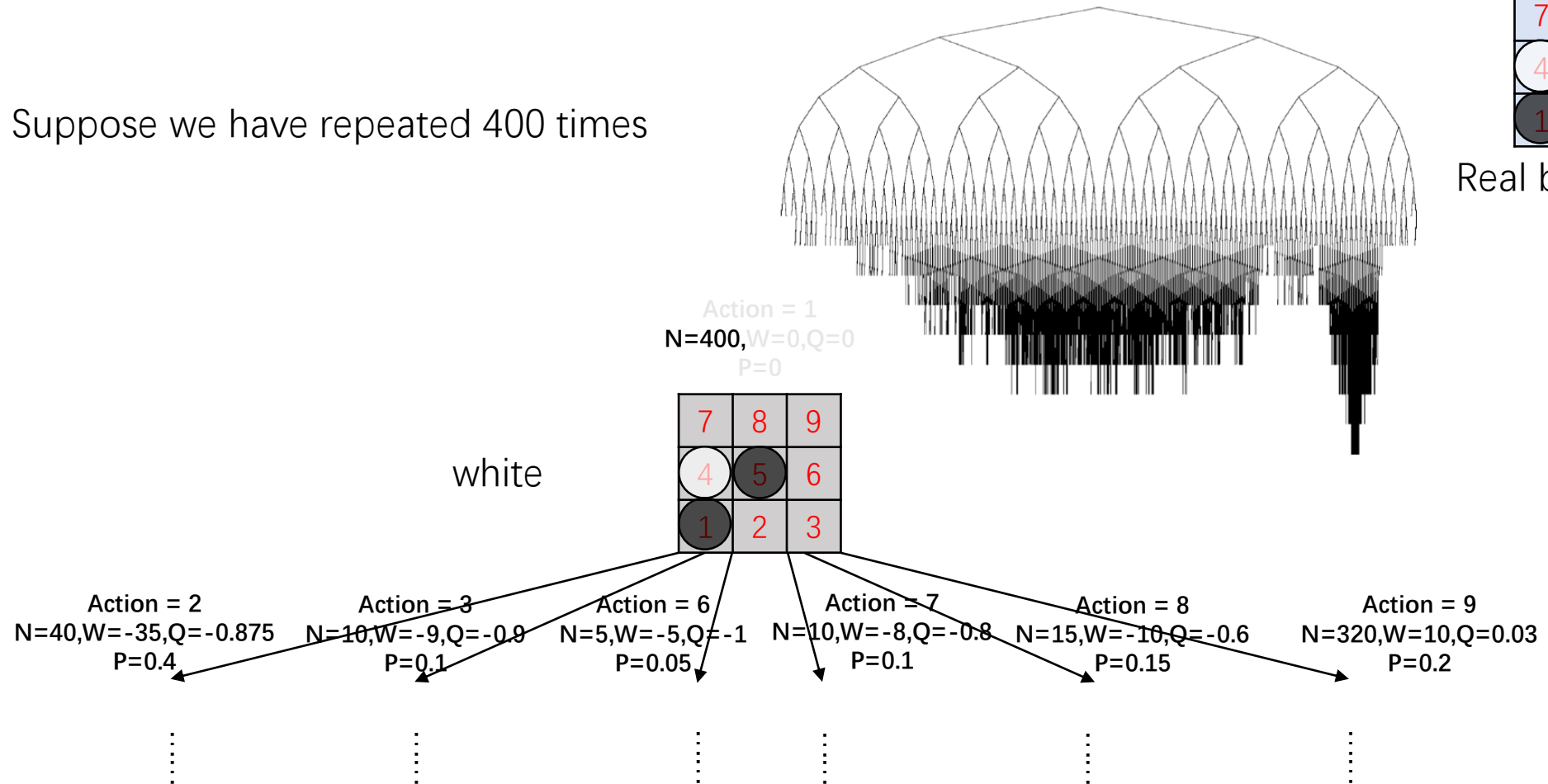
Done !



- Suppose we have repeated 400 times

7	8	9
4	5	6
1	2	3

Real board state



d Play

- ## Real board state

The diagram illustrates a policy network. At the top, a 3x3 grid represents a state space, with a black dot in the center and a white dot in the top-right corner. Below the grid is a triangular network structure. The central path of the network is labeled α_θ . The network branches out to a green bar chart labeled π , which represents the policy distribution over actions.

Action = 1
N=400, W=0, Q=0
P=0

7	8	9
4	5	6
1	2	3

Diagram illustrating a Markov Decision Process (MDP) with 6 states and 9 actions. The states are represented by vertical ellipses. The actions and their associated parameters are shown above the transitions:

- Action = 2: $N=40, W=-35, Q=-0.875, P=0.4$
- Action = 3: $N=10, W=-9, Q=-0.9, P=0.1$
- Action = 6: $N=5, W=-5, Q=-1, P=0.05$
- Action = 7: $N=10, W=-8, Q=-0.8, P=0.1$
- Action = 8: $N=15, W=-10, Q=-0.6, P=0.15$
- Action = 9: $N=320, W=10, Q=0.03, P=0.2$

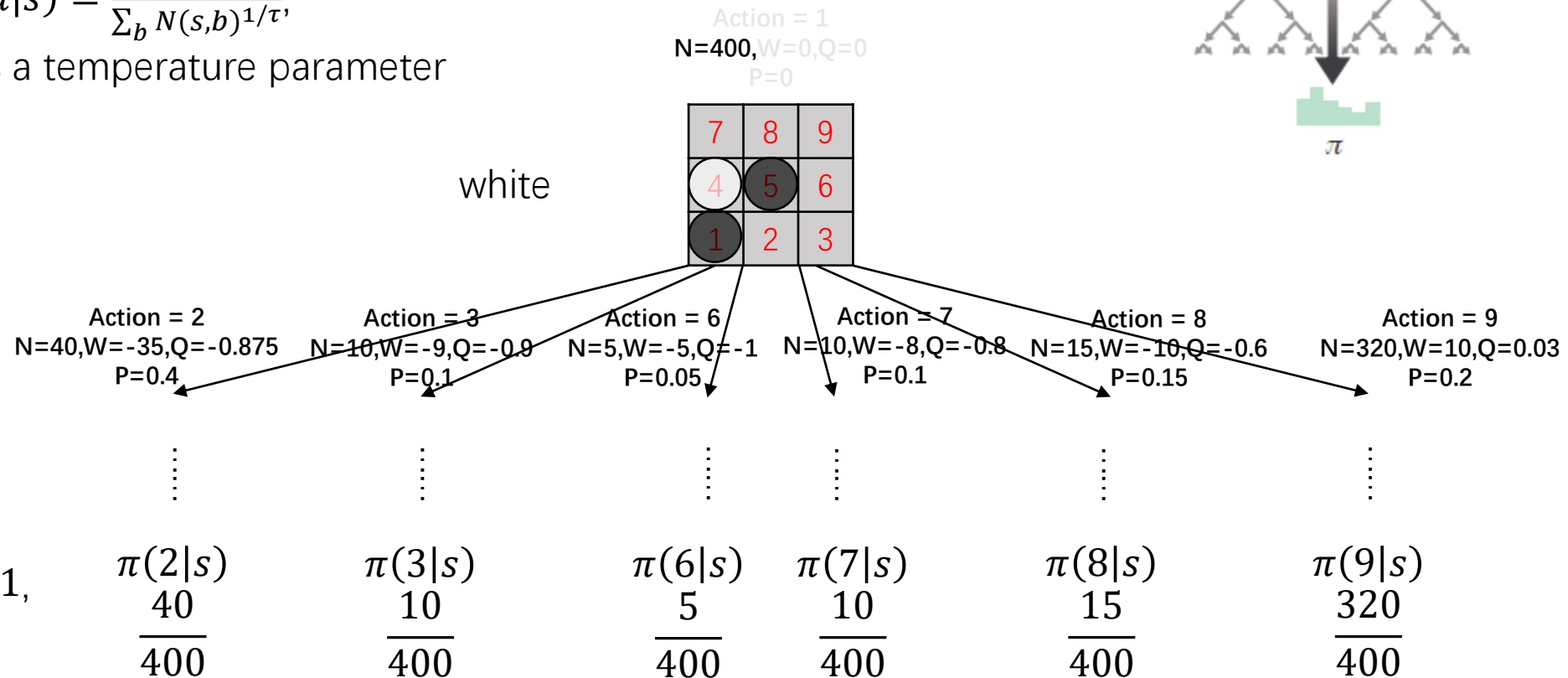
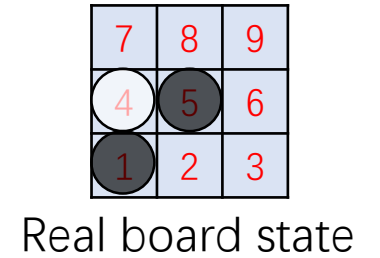
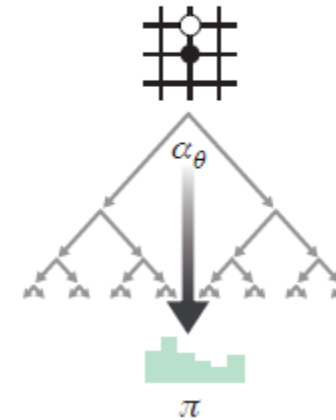
Play

- The move is chosen by calculating the probability

$$\pi(a|s) = \frac{N(s,a)^{1/\tau}}{\sum_b N(s,b)^{1/\tau}}$$

τ is a temperature parameter

d Play



If we set $\tau = 1$,
then

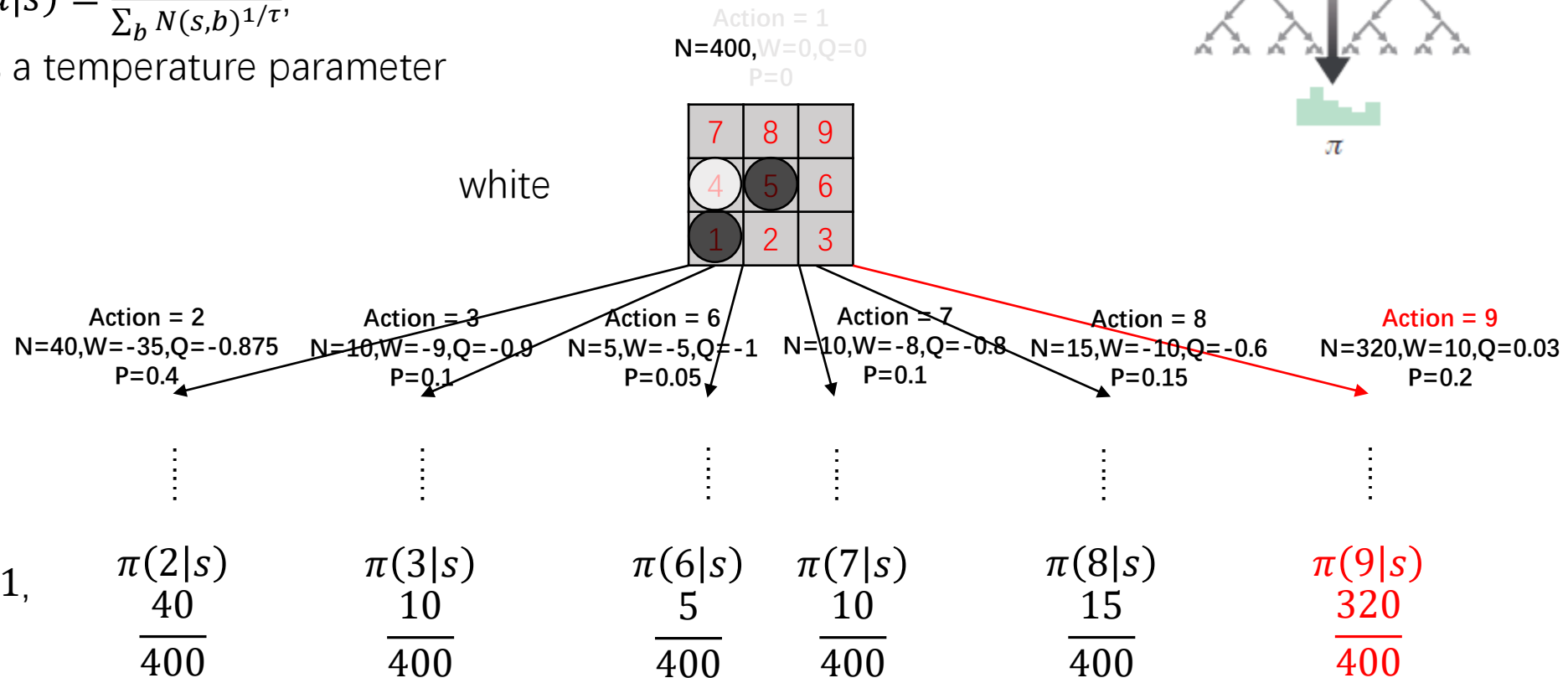
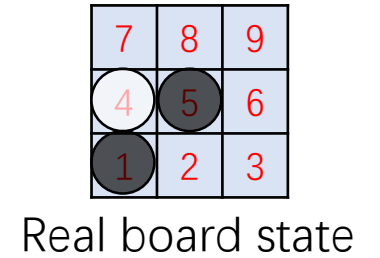
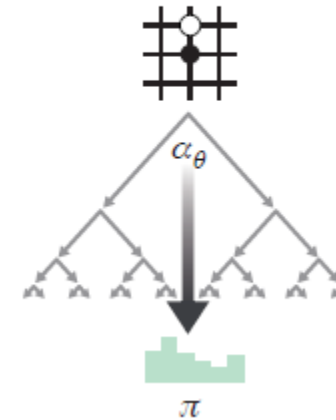
Play

- The move is chosen by calculating the probability

$$\pi(a|s) = \frac{N(s,a)^{1/\tau}}{\sum_b N(s,b)^{1/\tau}}$$

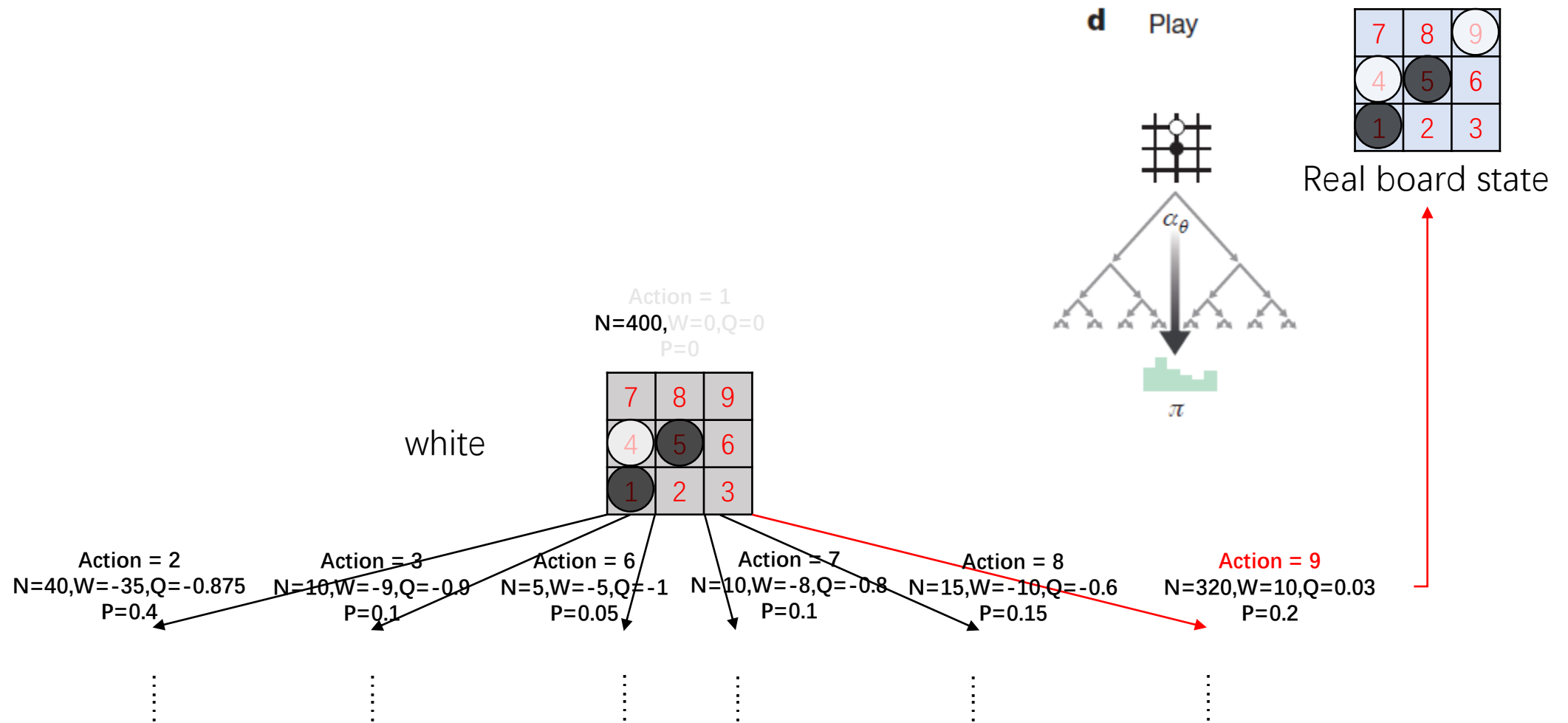
τ is a temperature parameter

d Play

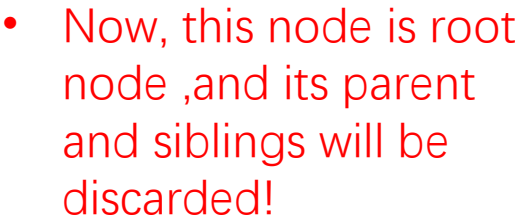


If we set $\tau = 1$,
then

Play



Do move in real board!



7	8	9
4	5	6
1	2	3

Real board state

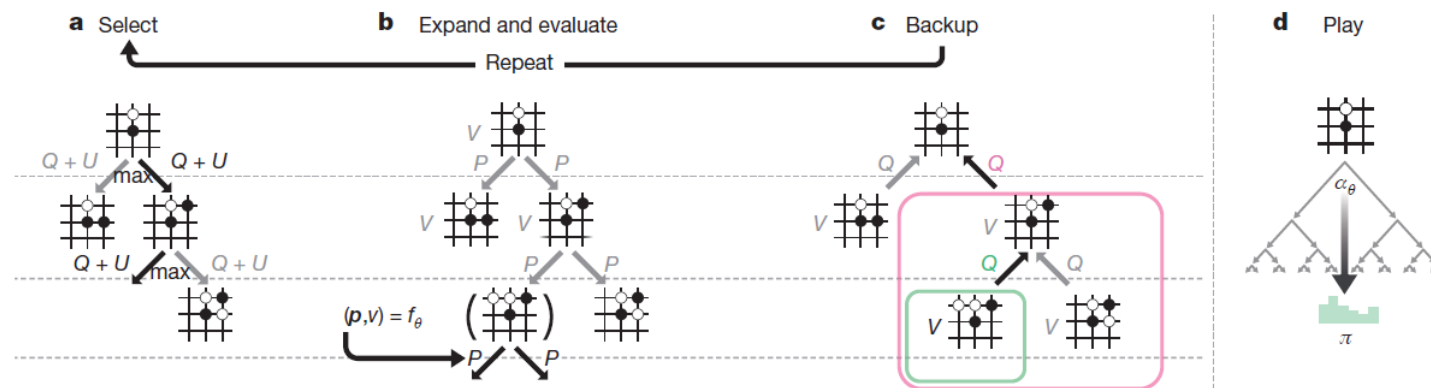
➤ Monte Carlo Tree Search in Gomoku

- The tree search algorithm will start from here again!
- And after every real move, the search algorithm will be repeated until game is over!

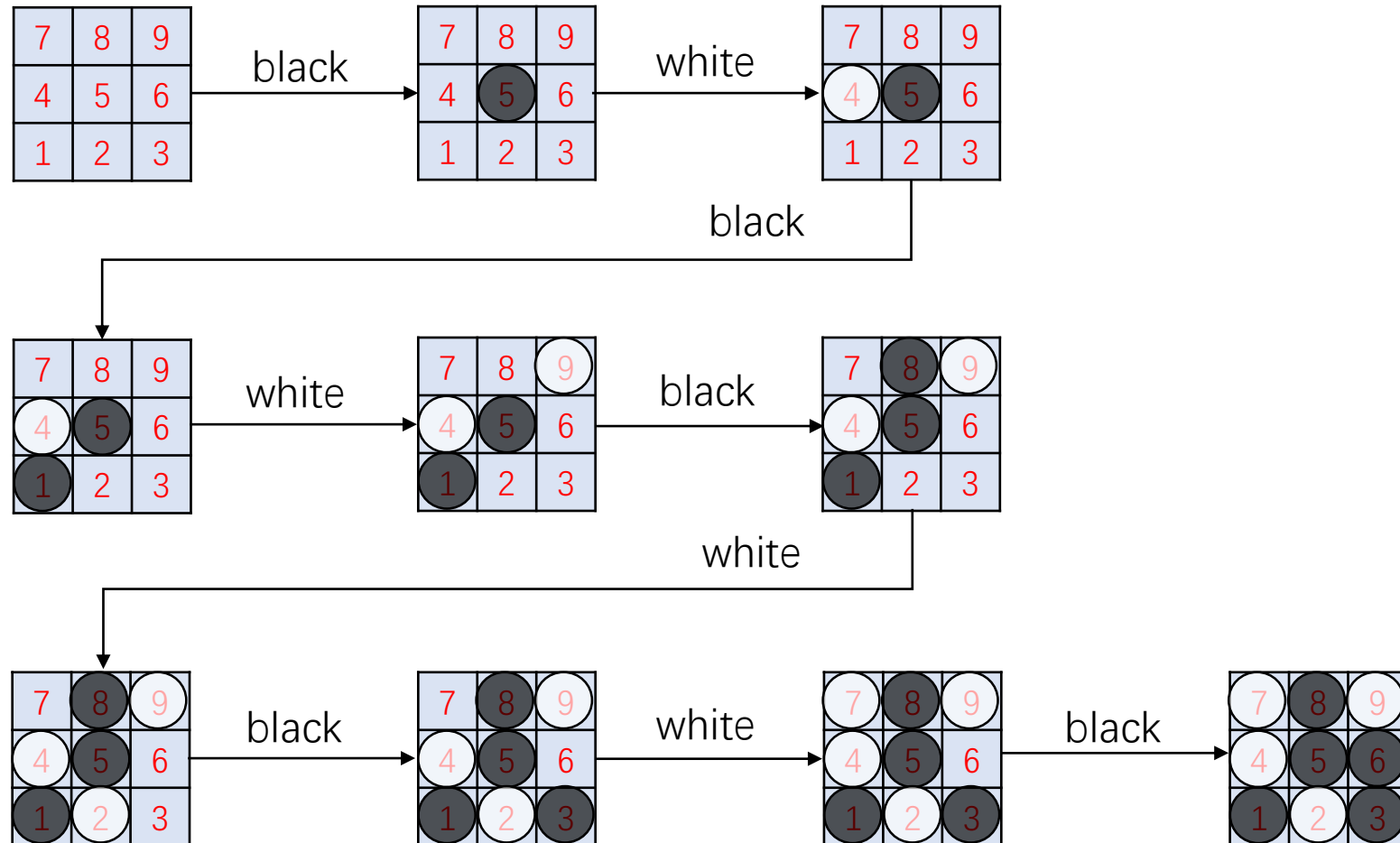
black

Action = 9
 $N=320, W=10, Q=0.03$
 $P=0.2$

7	8	9
4	5	6
1	2	3



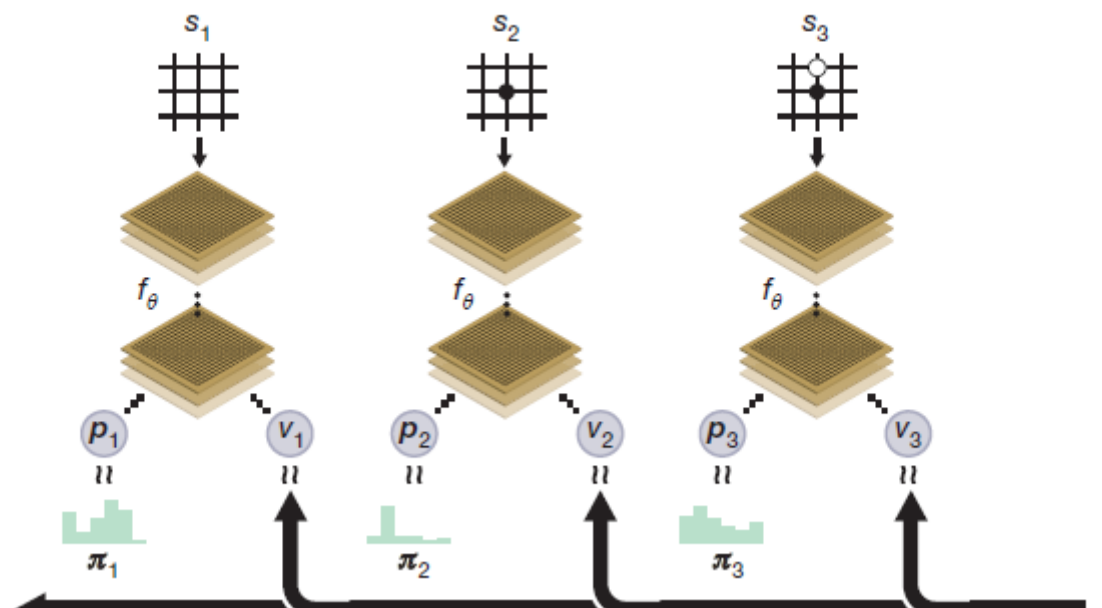
➤ One game is done!



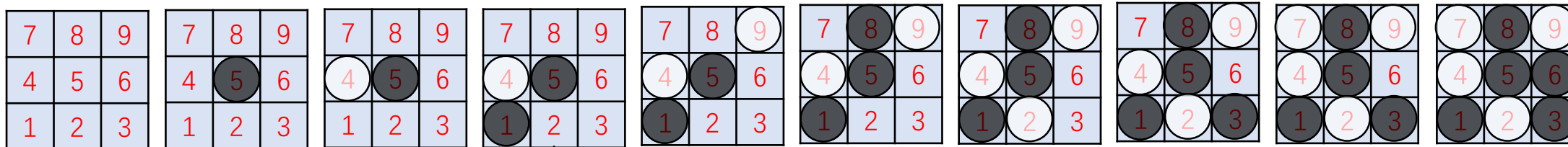
Game tie!
Reward = 0

- After self-play for $\#C$ games, Neural network training starts.

b Neural network training



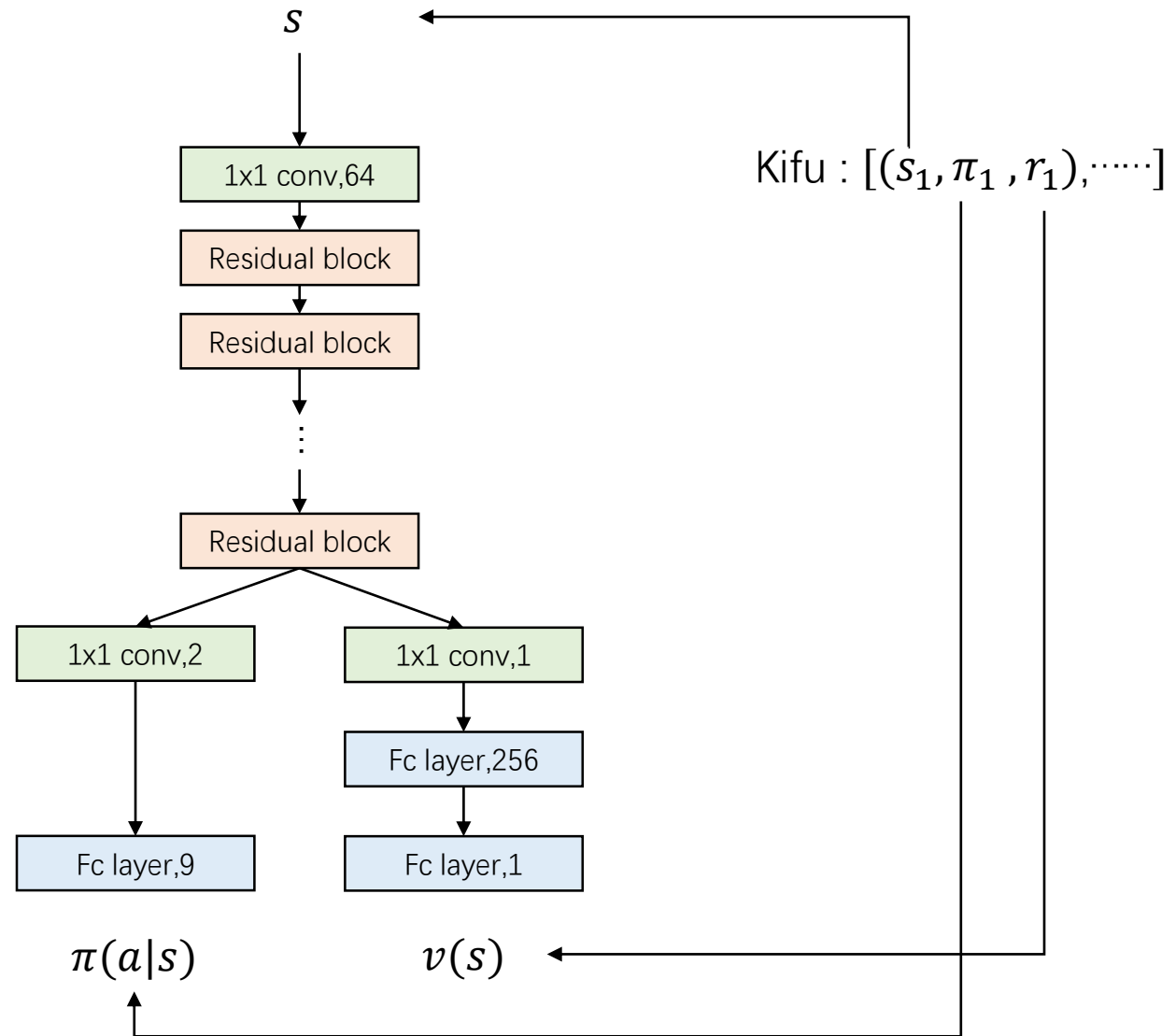
- We get the data : kifu



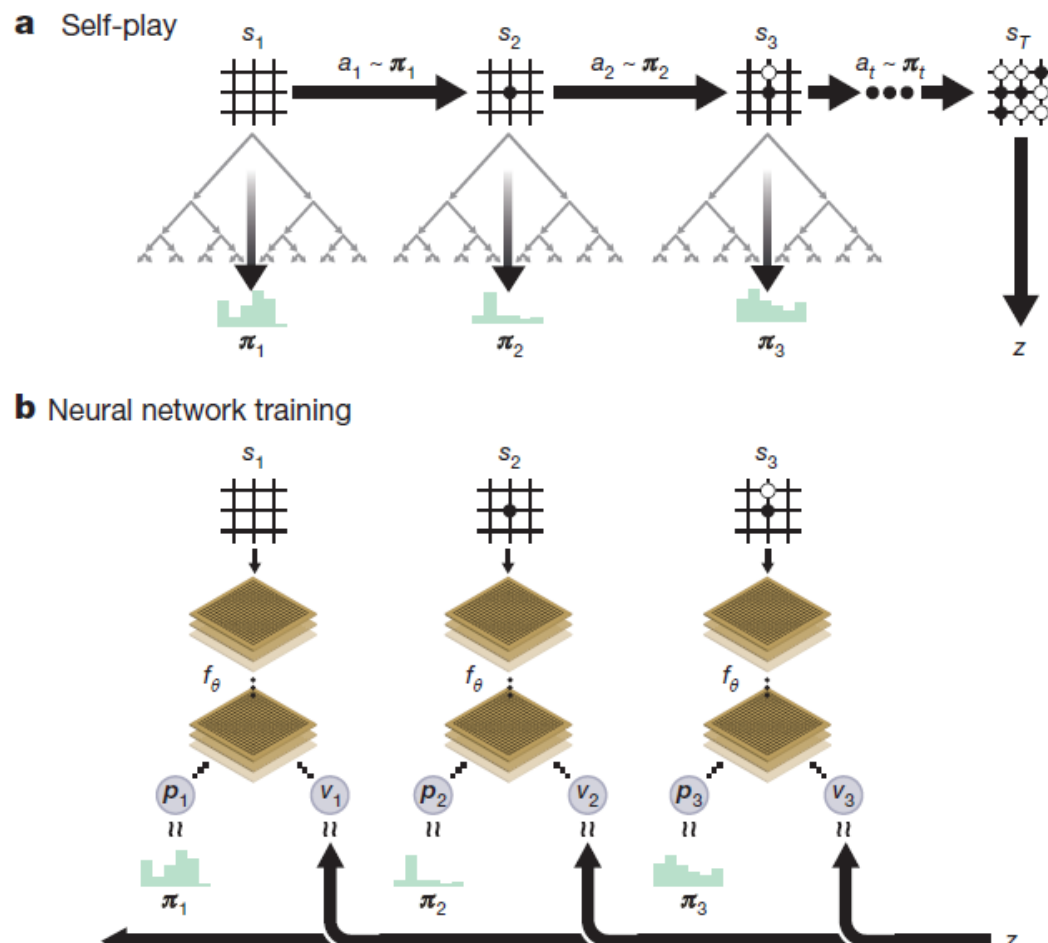
$$[(s_1, \pi(a|s_1) = \frac{N(s_1, a)^{\frac{1}{\tau}}}{\sum_b N(s_1, b)^{\frac{1}{\tau}}}, reward = 0), \\ \dots, \\ (s_{10}, \pi(a|s_{10}) = \frac{N(s_{10}, a)^{\frac{1}{\tau}}}{\sum_b N(s_{10}, b)^{\frac{1}{\tau}}}, reward = 0)]$$

$$s = \begin{array}{|c|c|c|} \hline 7 & 8 & 9 \\ \hline 4 & 5 & 6 \\ \hline 1 & 2 & 3 \\ \hline \end{array} \quad \pi(a|s) = \begin{array}{ccccccccc} \pi(1|s) & \pi(2|s) & \pi(3|s) & \pi(4|s) & \pi(5|s) & \pi(6|s) & \pi(7|s) & \pi(8|s) & \pi(9|s) \\ \hline 0 & 40 & 10 & 0 & 0 & 5 & 10 & 15 & 320 \\ \hline \frac{}{400} & \frac{}{400} & \frac{}{400} & \frac{}{400} & \frac{}{400} & \frac{}{400} & \frac{}{400} & \frac{}{400} & \frac{}{400} \\ \hline \end{array} \quad reward = 0$$

- Use the kifu to train the network!



- Repeat the whole process over and over
- The network can be more accurate and its strength will improve gradually!





Thanks !