# AlphaGo Zero

Shahzeb Aamir

# Content:

1. What is Go?

2. What is AlphaGo?

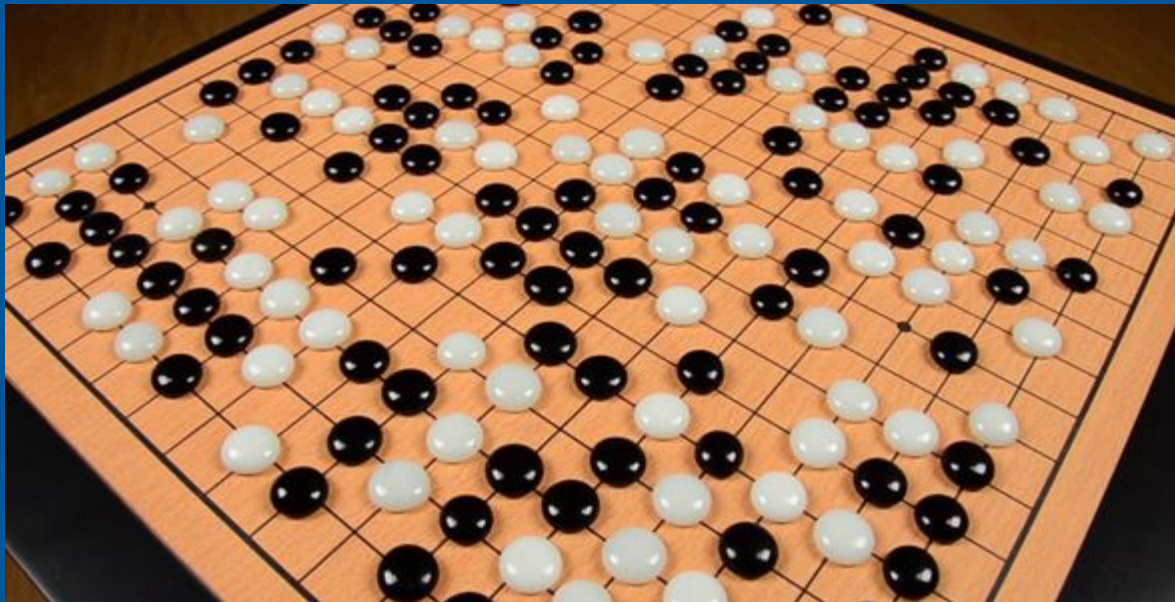3. Monte Carlo Tree Search
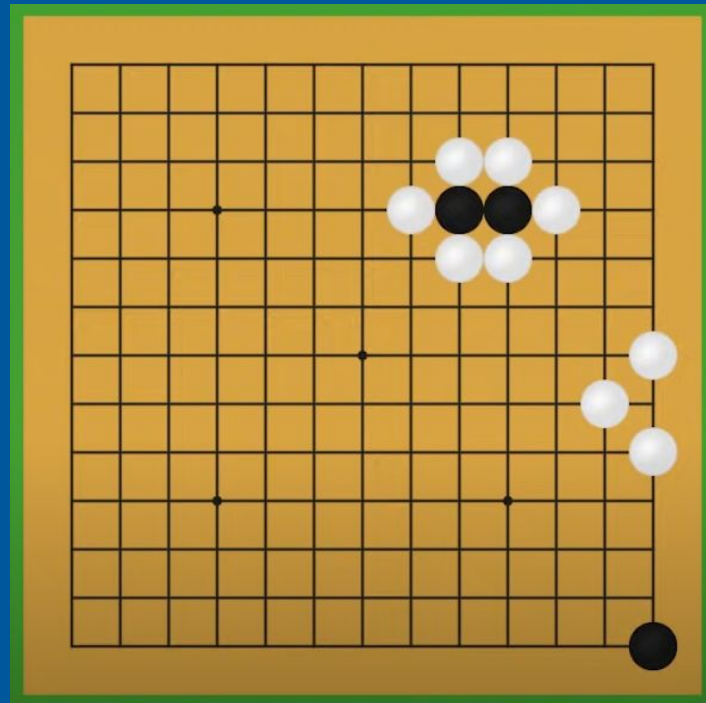
4. NN Architecture

## Go

19x19 grid

Turn-Based, Two players game

Goal:

Surround and capture opponents

stones, or strategically create

spaces of territory.

# Go

19x19 grid

Turn-Based, Two players game

Goal:

Surround and capture opponents

stones, or strategically create

spaces of territory.

# Go

19x19 grid

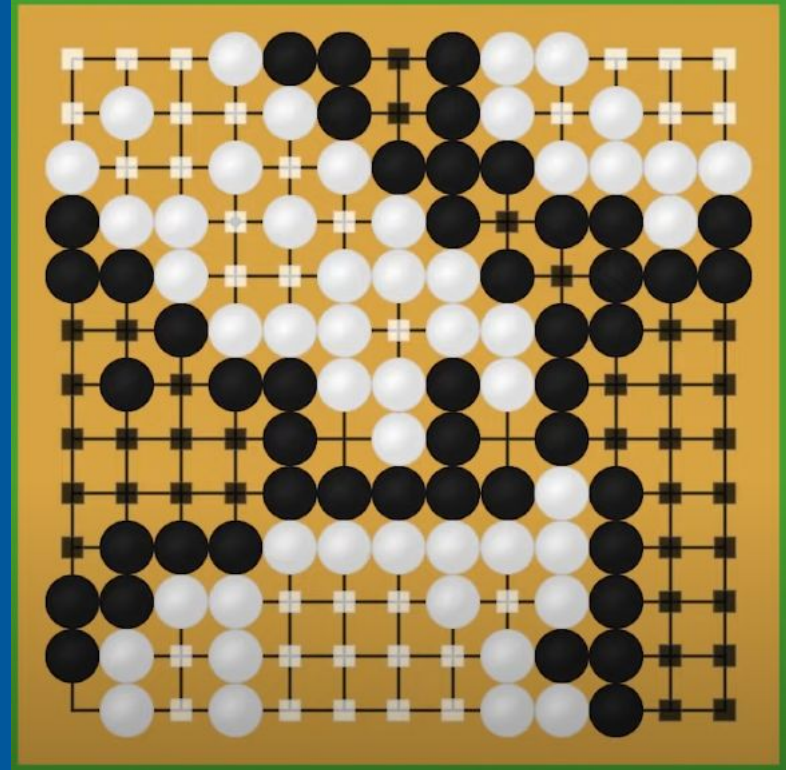Turn-Based, Two players game

Goal:

Surround and capture opponents

stones, or strategically create

spaces of territory.

Highest points of "empty spaces" win
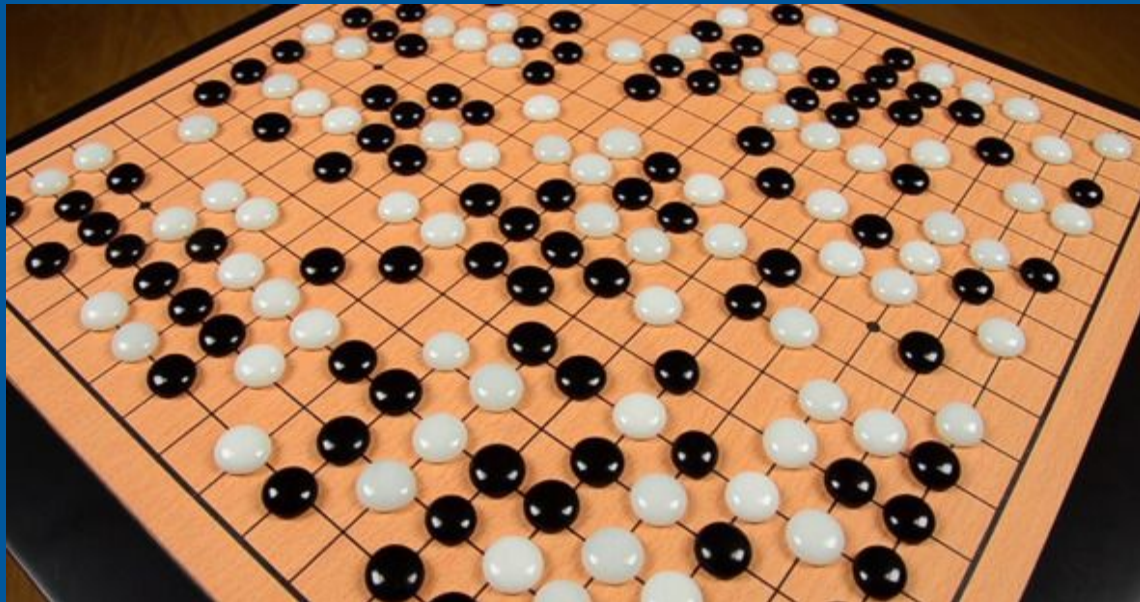
**Go**

Simple rules, right?

However, it has $10^{170}$ possible board

Configurations

More than the atoms in this universe

YES, you read right. More than the

atoms in this known universe

It's a googol times more complex than

chess

**What is AlphaGo**

AlphaZero, a single system that taught itself from scratch how to master the games of chess, shogi, and Go, beating a world-champion computer program in each case.

**What Problems does it solve?**

Turn-based, fully observable positions with definite sets of rules. The opponent goal is to prevent us from winning

**Not so intelligent Approach**

Brute force method. Search all possible moves and its subsequent branches to evaluate and select the best move

**A better Approach, and the idea behind AlphaGo**

A deep neural network estimates the most promising set of moves in a search tree

**How does it work**

Game State as Input → **DNN** → Estimate the value and proposed policy

Algorithm that performs intelligent search for possible moves based on the suggestion of DNN: Monte Carlo Tree Search



(a) Selection     (b) Expansion     (c) Simulation     (d) Backpropagation

# ALPHAGO ZERO CHEAT SHEET

The training pipeline for AlphaGo Zero consists of three stages, executed in parallel

## SELF PLAY
Create a 'training set'

The best current player plays 25,000 games against itself

See MCTS section to understand how AlphaGo Zero selects each move

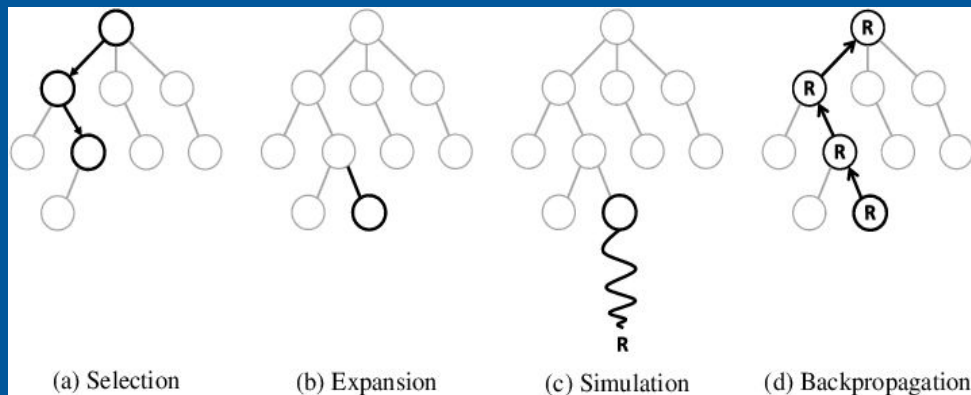At each move, the following information is stored

$\pi$

The game state
(see 'What is a Game State' section)

The search probabilities
(from the MCTS)

The winner
(+1 if this player won, -1 if the player lost - added once the game has finished)

## RETRAIN NETWORK
Optimise the network weights

A TRAINING LOOP

Sample a mini-batch of 2048 positions from the last 500,000 games

Retrain the current neural network on these positions
— The game states are the input (see 'Deep Neural Network Architecture')

Loss Function
Compares predictions from the neural network with the search probabilities and actual winner

PREDICTIONS

$P$   Cross-entropy
+
$V$   Mean-squared error
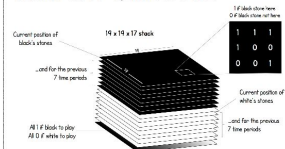+
Regularisation

ACTUAL

$\pi$

After every 1,000 training loops, evaluate the network

## EVALUATE NETWORK
Test to see if the new network is stronger

Play 400 games between the latest neural network and the current best neural network

Both players use MCTS to select their moves, with their respective neural networks to evaluate leaf nodes

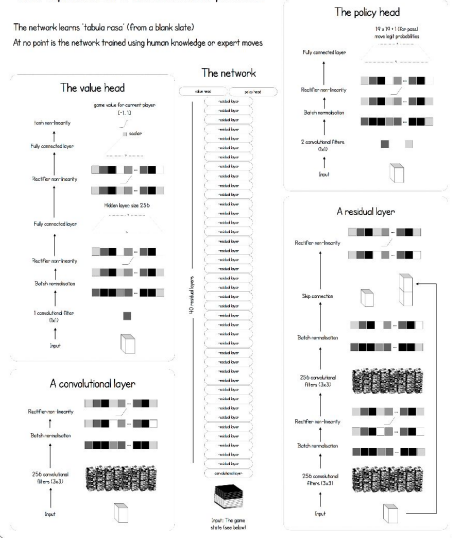Latest player must win 55% of games to be declared the new best player

I          II

## WHAT IS A 'GAME STATE'

1 if black stone here
0 if black stone not here

Current position of black's stones

19 x 19 x 17 stack

1 1 1
1 0 0
0 0 1

...and for the previous 7 time periods

Current position of white's stones

...and for the previous 7 time periods

All 1 if black to play
All 0 if white to play
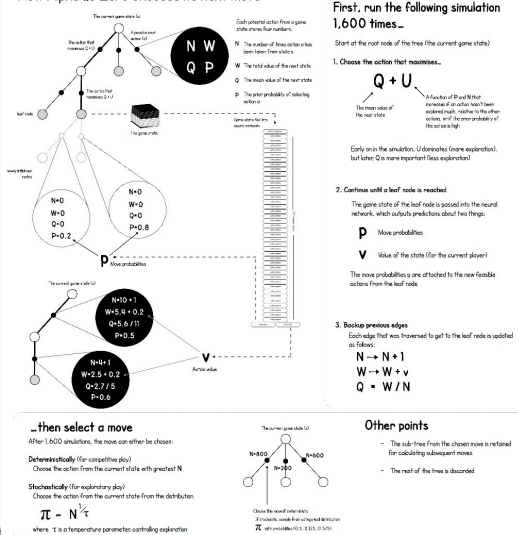
This stack is the input to the deep neural network

## THE DEEP NEURAL NETWORK ARCHITECTURE
How AlphaGo Zero assesses new positions

The network learns 'tabula rasa' (from a blank slate)

At no point is the network trained using human knowledge or expert moves

### The value head
game value for current player [-1, 1]

tanh non-linearity

Fully connected layer

Rectifier non-linearity

Fully connected layer

Hidden layer size 256

Rectifier non-linearity

Batch normalisation

1 convolutional filter (1x1)

Input

### The network
value head    policy head

### A convolutional layer

Rectifier non-linearity

Batch normalisation

256 convolutional filters (3x3)

Input

### The policy head
19 x 19 + 1 (for each legal move) probabilities

Fully connected layer

Rectifier non-linearity

Batch normalisation

2 convolutional filters (1x1)

Input

### A residual layer

Rectifier non-linearity

Skip connection

Batch normalisation

256 convolutional filters (3x3)

Rectifier non-linearity

Batch normalisation

256 convolutional filters (3x3)

Input

Input: The game state (see below)

## MONTE CARLO TREE SEARCH (MCTS)
How AlphaGo Zero chooses its next move

The current game state (s)

N W
Q P

Each potential action from a game state stores these four numbers:

$N$  The number of times action a has been taken from state s

$W$  The total value of the next state

$Q$  The mean value of the next state

$P$  The prior probability of selecting action a

leaf node

one game state

N=0
W=0
Q=0
P=0.2

N=0
W=0
Q=0
P=0.8

$P$  Move probabilities

N=10 + 1
W=5.4 + 0.2
Q=5.6 / 11
P=0.5

N=4+1
W=2.5 + 0.2
Q=2.7 / 5
P=0.6

$V$  Action value

### First, run the following simulation 1,600 times...

Start at the root node of the tree (the current game state)

**1. Choose the action that maximises...**

$$Q + U$$

The mean value of the next state

A function of P and N that increases if an action hasn't been explored much, relative to the other actions, or if the prior probability of the action is high

Early on in the simulation, U dominates (more exploration), but later, Q is more important (less exploration)

**2. Continue until a leaf node is reached**

The game state of the leaf node is passed into the neural network, which outputs predictions about two things:

$P$  Move probabilities

$V$  Value of the state (for the current player)

The move probabilities p are attached to the new feasible actions from the leaf node.

**3. Backup previous edges**

Each edge that was traversed to get to the leaf node is updated as follows:

$$N \rightarrow N + 1$$
$$W \rightarrow W + v$$
$$Q = W / N$$

## ...then select a move
After 1,600 simulations, the move can either be chosen:

The current game state (s)

N=800    N=600

N=200

**Deterministically** (for competitive play)
Choose the action from the current state with greatest N

**Stochastically** (for exploratory play)
Choose the action from the current state from the distribution:

$$\pi \sim N^{1/\tau}$$

where $\tau$ is a temperature parameter controlling exploration

Choose the move if deterministic
if stochastic, sample from categorical distribution
$\pi$ with probabilities (0.5, 0.125, 0.375)

## Other points

— The sub-tree from the chosen move is retained for calculating subsequent moves

— The rest of the tree is discarded

# How does it work

Three stages executed in parallel

## SELF PLAY

Create a 'training set'

The best current player plays 25,000 games against itself
See MCTS section to understand how AlphaGo Zero selects each move

At each move, the following information is stored

$\pi$

The game state
(see 'What is a Game State' section)

The search probabilities
(from the MCTS)

The winner
(+1 if this player won, -1 if this player lost - added once the game has finished)

## RETRAIN NETWORK

Optimise the network weights

A TRAINING LOOP

Sample a mini-batch of 2048 positions from the last 500,000 games

Retrain the current neural network on these positions
– The game states are the input (see 'Deep Neural Network Architecture')

Loss function
Compares predictions from the neural network with the search probabilities and actual winner

PREDICTIONS **P** **V**   Cross-entropy + Mean-squared error + Regularisation   $\pi$ ACTUAL

After every 1,000 training loops, evaluate the network

## EVALUATE NETWORK

Test to see if the new network is stronger

Play 400 games between the latest neural network and the current best neural network

Both players use MCTS to select their moves, with their respective neural networks to evaluate leaf nodes

Latest player must win 55% of games to be declared the new best player

I   II

**Monte Carlo Tree Search:**
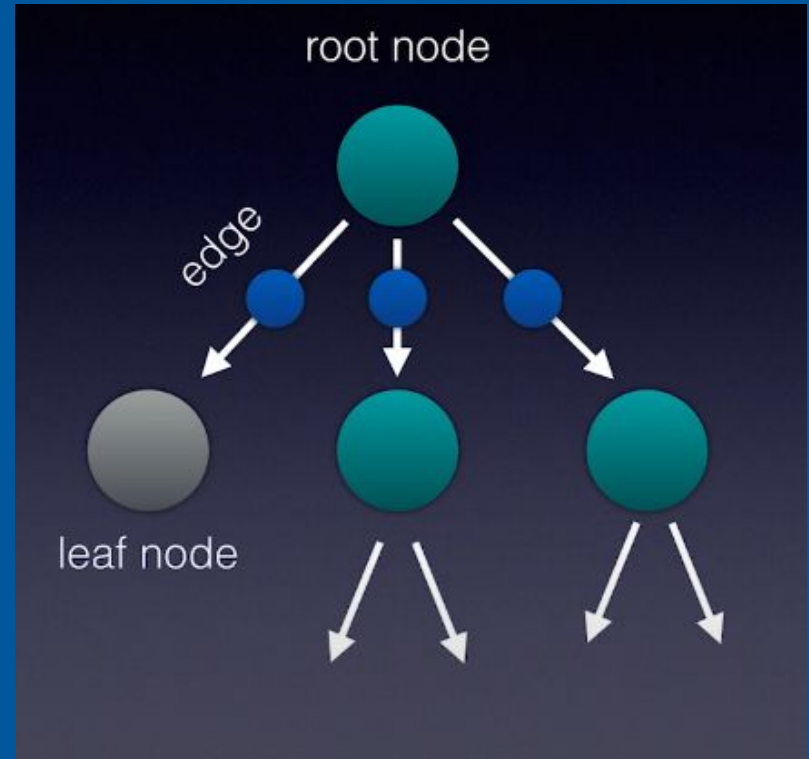
**The Phases of MCTS:**

SELECT (Root to Leaf that is most promising)

EXPAND (by using one more move)

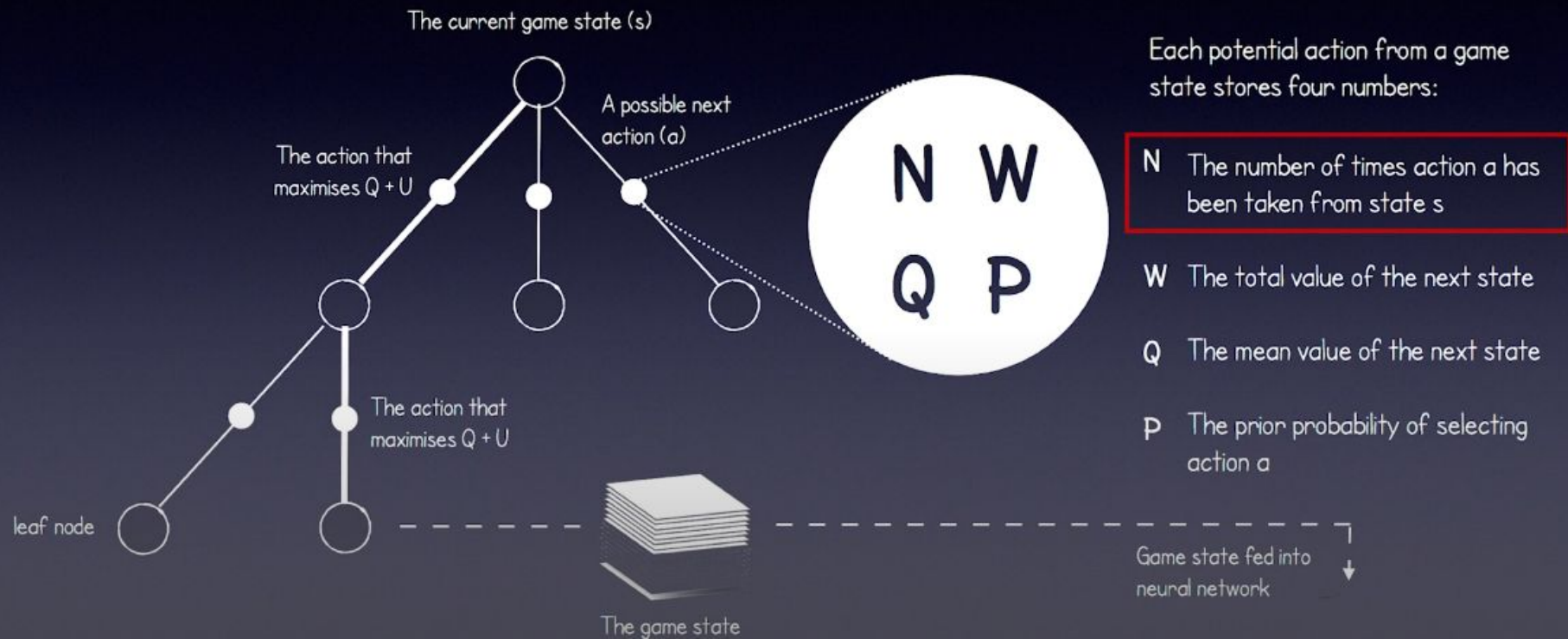BACK UP (and update all edges traversed using statistics)

X 1600 times

PLAY

**Monte Carlo Tree Search:**

**Statistics:**



The current game state (s)

A possible next action (a)

The action that maximises Q + U

The action that maximises Q + U

leaf node

The game state

N W
Q P

Each potential action from a game state stores four numbers:

N    The number of times action a has been taken from state s

W    The total value of the next state

Q    The mean value of the next state

P    The prior probability of selecting action a

Game state fed into neural network

**Monte Carlo Tree Search:**

**Details:**

1. Choose the action that maximises...

$$Q + U$$

The mean value of the next state

A function of **P** and **N** that increases if an action hasn't been explored much, relative to the other actions, or if the prior probability of the action is high

Early on in the simulation, U dominates (more exploration), but later, Q is more important (less exploration)

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

**Monte Carlo Tree Search:**

**Details:**

## 2. Continue until a leaf node is reached

The game state of the leaf node is passed into the neural network, which outputs predictions about two things:

**p**     Move probabilities

**v**     Value of the state (for the current player)

The move probabilities p are attached to the new feasible actions from the leaf node

**Monte Carlo Tree Search:**

**Details:**

3. Backup previous edges

Each edge that was traversed to get to the leaf node is updated as follows:

$$N \rightarrow N + 1$$
$$W \rightarrow W + v$$
$$Q = W / N$$

**Monte Carlo Tree Search:**

**Details: For Playing (Selecting the move)**

- **Deterministically** (for competitive play) choose the action with the greatest N

- **Stochastically** (for training) sample randomly from the probability distribution…

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

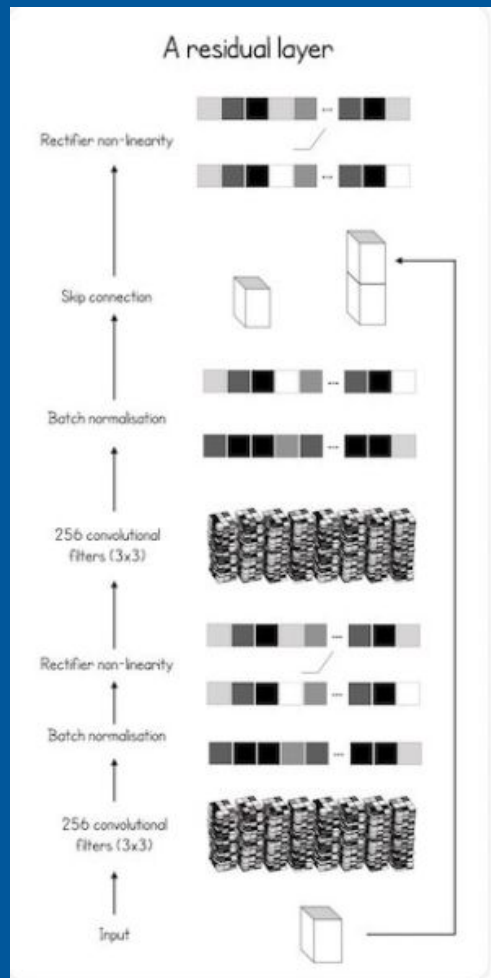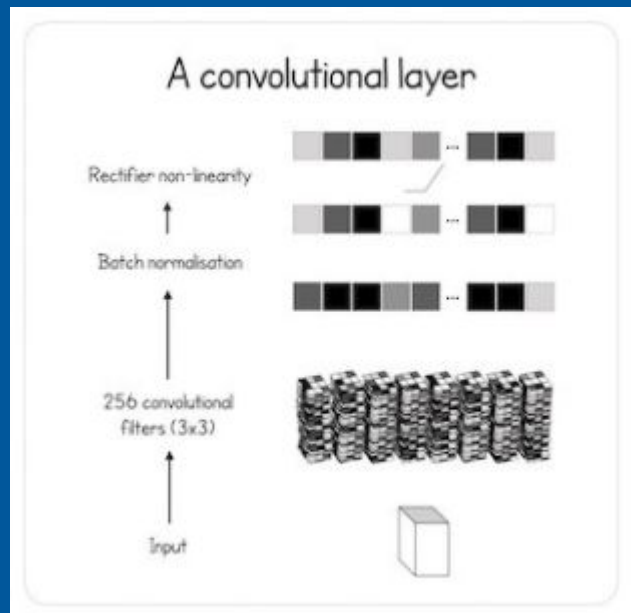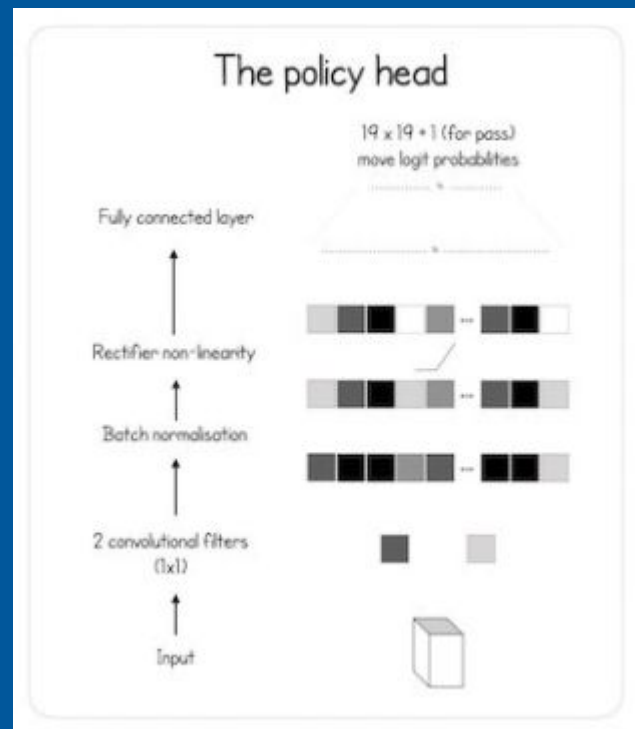Where $\tau$ is a temperature parameter, controlling exploration.

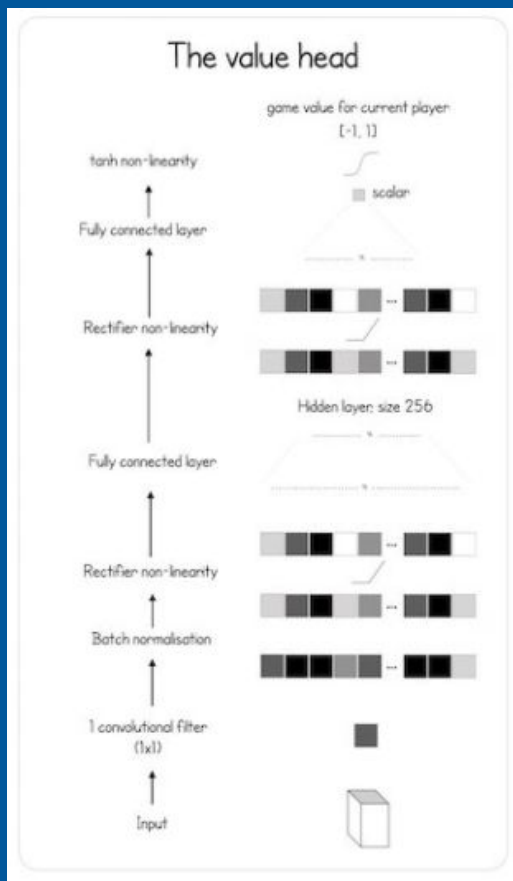The current game state (s)

N=800

N=600

N=200

Choose this move if deterministic

If stochastic, sample from categorical distribution

$\pi$ with probabilities (0.5, 0.125, 0.375)

**NN Architecture:**

**Game State:**

**NN Architecture:**



A convolutional layer

Rectifier non-linearity

Batch normalisation

256 convolutional filters (3x3)

Input

A residual layer

Rectifier non-linearity

Skip connection

Batch normalisation

256 convolutional filters (3x3)

Rectifier non-linearity

Batch normalisation

256 convolutional filters (3x3)

Input

**NN Architecture:**



$$L = (z - v)^2 - \pi \cdot \log(p) + c \cdot \|\theta\|^2$$

**AlphaGo beats Lee Sedol:**

**Thank You!**