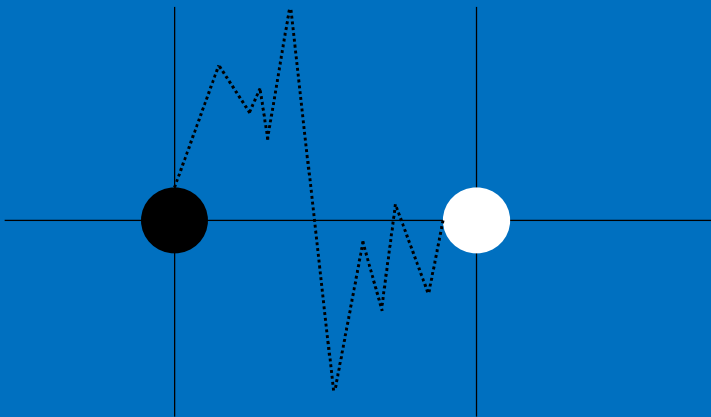# From AlphaGo to AlphaZero

1. Mastering the Game of Go with Deep Neural Networks and Tree Search (2016)
2. Mastering the Game of Go without Human Knowledge (2017)
3. Mastering Chess and Shogi by Self-Play with a General Reinforcement Learning Algorithm (2017)

- By Harvey
- 2022

# 1. Deep Reinforcement Learning in Different Games
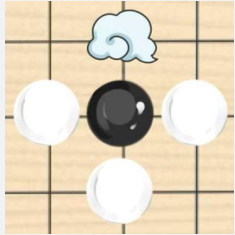
Assumption: AI can only see/do what humans can see/do



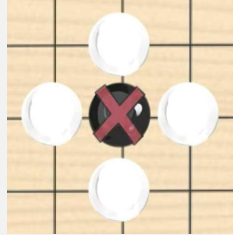| Games | Atari | Board Game | Dota / LOL | Starcraft II |
|---|---|---|---|---|
| Agents | 1 | 2 | $\geq 1$ (Generally 5 vs. 5) | $\geq 2$ (Generally 1 vs. 1) |
| Information | Perfect information | Perfect Information | Imperfect information (fog of war) | Imperfect information (fog of war) |
| State space | Limited | Limited but large | Unlimited | Unlimited |
| Action space | Limited | Limited but large | Limited but large | Unlimited |
| AI vs. Human | AI dominated | AI dominated | AI won in 1 vs. 1 Completely failed in 5 vs. 5 (OpenAI) | Completely failed in any competitive game but learn to do simple tasks (DeepMind) |

# 1. Go Game: Basic Rules
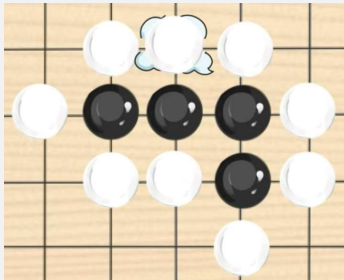
## How to play

### Liberty



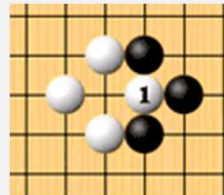- Liberty is when the area surrounding the go has no other go



- If a go has no liberty, it is consider "dead" and will be taken out

### Liberty (Multiple Go)



- In this case, black gos are "dead" and will be taken out

### Ko Rule



- One cannot place the go in "dead" area unless you can eat them
- Forbid sequential eat in the same position
- Unless you make a move that significantly changes the board position

### Life



- If go is surrounded and the area has more than two liberties (i.e. they cannot be taken out with one step), then this area of go is considered to be "alive" (life)
- In the case shown in the left, black gos cannot be placed inside the white gos. Hence, white gos are considered to be "alive"

## How to win



The winner is:
1. whoever occupies more areas of the board.
2. whoever has a larger number of "alive" go on the board .
3. rules are geo-dependent but with minor difference.

# Versions

| Version | Training hardware | Elo rating | When | |
|---|---|---|---|---|
| AlphaGo Fan | 176 GPUs | 3144 | Oct 2015 | 5:0 against Fan Hui |
| AlphaGo Lee | 48 TPUs distributed | 3738 | Mar 2016 | 4:1 against Lee Sedol |
| AlphaGo Master | 4 TPUs, single machine | 4858 | May 2017 | 60:0 against professional players |
| AlphaGo Zero | 4 TPUs, single machine | 5185 | Oct 2017 | 100:0 against AlphaGo Lee |
| AlphaZero | 4 TPUs, single machine | 5018 | Dec 2017 | 60:40 against AlphaGo Zero |

TPU: tensor processing unit

# AlphaGo

An "interesting" fact about AlphaGo:
The most powerful version consists of ~1900 CPUs and ~200 GPUs
Electricity bill is ~$3000 per game, training phase: 160,000 game $\approx$ $480 million

# Networks

- Policy network
  - Given state s, probability of action.
- Value network
  - Given state s, expected value of winning the game.
- Fast rollout
  - A much faster (smaller) network to make decisions.

# 2. Why Go and Prior Work

## Why Go Game?



- The most challenging classic board games from a computational perspective.
- 19 X 19 board positions.
- $250^{150}$ possible sequences of moves (chess: $35^{80}$).
- The objective is to occupy more territory.
  - Leads to highly sophisticated evaluation functions
  - No one before AlphaGo has successfully build an effective evaluation function.
- Exhaustive search is infeasible.

## Prior Work

**Reinforcement Learning**       **Self-play**       **Linear Value Function**



$$V_\theta(s) = \phi(s)^T \theta$$

**Tree Search**

### Minimax (Alpha-Beta) Tree Search

- Most game too large for Minimax Tree Search because it requires one to go all the way to the end of the game
- Truncate the tree by using approximated value function $V_\theta(s) \approx v^*(s)$
- Super-human performance in chess
- Not effective in Go

### Monte Carlo Tree Search

- Double approximation
- $V^n(s) \approx v^{p^n}(s) \approx v^*(s)$
- First, use $n$ Monte Carlo simulations to estimate the value function of a policy $p^n$
- Second, use the value function to approximate the optimal value function
- Why does it work? In the limit they are equivalent (mathematically).

**Why?**

Reinforcement Learning: Neurobiology foundations

Tree Search: Players tend to make forecasts (truncated) and predictions

# 2. AlphaGo: Structure Breakdown

## Training Pipeline

Human expert positions

Classification

Supervised Learning
policy network

Reinforcement Learning
policy network

Self-Play

Self-Play data

Minimize Mean Squared Errors

Regression

Value Network

(Silver et al. 2016)

## Testing/Tournament

Monte Carlo Tree Search (MCTS)

t
Current board position

t+1
If I choose an action

(Silver et al. 2016)

### Problems with a full-scale MCTS

- Tree depth: need to simulate all the way to the end of the game

- Tree breadth: enormous amounts of choices per time step

Tree Depth

Shrink the tree size

Tree Breadth

Use value network to reduce the tree depth.

Use policy network to reduce the tree breadth.

# 2. AlphaGo: Supervised Learning of Policy Networks

## Convolutional Neural Network

### CNN Training and Testing

- To predict human experts' move.

- 160,000 games (professional players).

- 29.4 million moves (28.4 million training and 1 million testing).

- Stochastic Gradient Descent update to maximize the action log likelihood

## Prediction accuracy: 55-57%

If I feed the neural network with a board status (state), the neural network can predict, on average 55-57% chance, the correct move that a human expert will do under the same board condition.

**Better accuracy or higher chances of winning?**
**After all, human (including experts) make mistakes/their strategy may not be optimal strategy**

**Supervised learning is not enough**

---

(Smiling) Puppy picture

Human expert positions

Pixels

48 feature planes
(each 19 X 19)
19 X 19 X 48

13 layers        1 fully connected layer        Softmax

+        +

In the picture, there is a
- (Smiling) Puppy (80%)
- (Angry) puppy (60%)
  - Kittie (10%)
    - Baby (1%)
      - ...

Human experts will take a move in
- Position 1 (56%)
- Position 2 (30%)
- Position 3 (20%)
- ...

# 2.1 AlphaGo: Policy and Policy Network

## Policy Based Learning

$P(a|s)$ **Formally, a policy maps states to actions**

- In policy gradient framework, it is a probability distribution over all states and actions (vs. value framework).
- Policy will tell the probability of I choose a particular action under current state.
- No value function, states and actions are allowed to be continuous and infinite.
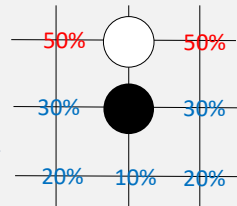- Often parameterized, a function $P_\theta(a|s)$.

### Policy Network in AlphaGo

Take board position **s** as input

Action probability distribution/matrix

$P_\theta(a|s)$

| 50% | n/a | 50% |
|---|---|---|
| 30% | n/a | 30% |
| 20% | 10% | 20% |

***Training:***
Stochastic gradient ascent to maximize the log likelihood

## Rollout Policy vs. Supervised Learning Policy

**Trade off between speed and accuracy**

| Choose an action | Rollout Policy | SL Policy |
|---|---|---|
| Speed | $2\mu s$ | $3ms$ |
| Accuracy | 24.4% | 57% |

**The key differences**

| Structure | Rollout Policy | SL Policy |
|---|---|---|
| activation function | Linear | Non-linear |
| Feature fed into CNN | Small feature size (A simple demo) | Full size (19 X 19 X 48) |
| Notation | $P_\pi(a|s)$ | $P_\sigma(a|s)$ |

# 2. AlphaGo: Structure Breakdown

## Training Pipeline

Human expert positions

**Classification** →

Supervised Learning
policy network

Reinforcement Learning
policy network    Self-Play

Self-play data

Minimize Mean Squared Errors

**Regression** →

Value Network

(Silver et al. 2016)

## Testing/Tournament

Monte Carlo Tree Search (MCTS)

t
Current board position

t+1
If I choose an action

⋮    ⋮

(Silver et al. 2016)

### Problems with a full-scale MCTS

- Tree depth: need to simulate all the way to the end of the game

- Tree breadth: enormous amount of choices per time step

Tree Depth

Use value network to reduce the tree depth

Shrink the tree size

Tree Breadth

Use policy network to reduce the tree breadth

# 2.2 Reinforcement Learning: Policy Network Iterations

**The aim here is to allow the policy to be improved by letting the AI play against itself**

$$P_\sigma(a|s) \Rightarrow P_\rho(a|s), \sigma \text{ and } \rho \text{ are weight parameters}$$



Opponents Pool

Old 1 (with parameters $\sigma_1$)
Old 2 (with parameters $\sigma_2$)
Old 3 (with parameters $\sigma_3$)
…

Randomly sampled

$\sigma_0$ vs $\sigma_i$

10,000 iterations / mini-batches

128 games per iteration/mini-batch

Per 500 iterations

Current AlphaGo (with parameters $\sigma_0$)

128 results (scores) Then replayed to determine policy gradient update

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

$$z_t = \pm 1$$

# 2. AlphaGo: Structure Breakdown

## Training Pipeline

Human expert positions

Classification →

Supervised Learning
policy network

Reinforcement Learning
policy network

Self-Play

Self-play data

Minimize Mean Squared Errors

Regression →

Value Network

(Silver et al. 2016)

## Testing/Tournament

Monte Carlo Tree Search (MCTS)

t
Current board position

t+1
If I choose an action

⋮          ⋮

(Silver et al. 2016)

### Problems with pure MCTS

- Tree depth: need to simulate all the way to the end of the game

- Tree breadth: enormous amount of choices per time step

Tree Depth

Shrink the tree size

Tree Breadth

Use value network to reduce the tree depth

Use policy network to reduce the tree breadth

# 2.3 Reinforcement Learning: Value Network and Position Evaluation

**Now we have a better policy (strategy), we want to predict the outcome of game when we make a move**
**i.e., given the current board position, how likely will I win the game.**

## Value Function

$$v^p(s) = \mathbb{E}[z_t | s_t = s, \, a_{t...T} \sim p]$$

- **Interpretation:** the value (chances of win) of the current state under the policy p.
- **Problem:** we do not know the optimal value function under perfect play $v^*(s)$
- **Solution:** use approximation
  - $v_\theta(s) \approx v^{p_\rho}(s) \approx v^*(s)$

### Training

- Parameterized value function $v_\theta(s)$ with weights $\theta$
- We know the game results for a given state (board position) under the strongest policy $p_\rho$, i.e. $v_\theta(s)$.
- **Optimization problem:** optimize $\theta$ so that $v_\theta(s)$ is close to $v^{p_\rho}(s)$
- **Approach:** standard regression (projection) method
- Minimize Mean Squared Errors:
  - $\left(v_\theta(s) - z^k\right)^2$ where $z^k = \pm 1$

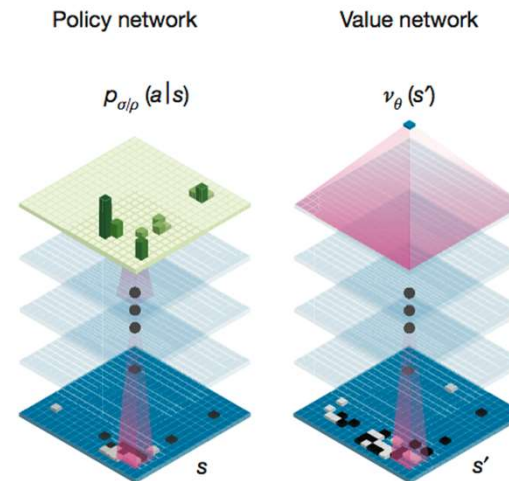## Overfitting

### Naïve approach leads to overfitting

- **Training:** Mean Squared Errors of 0.19
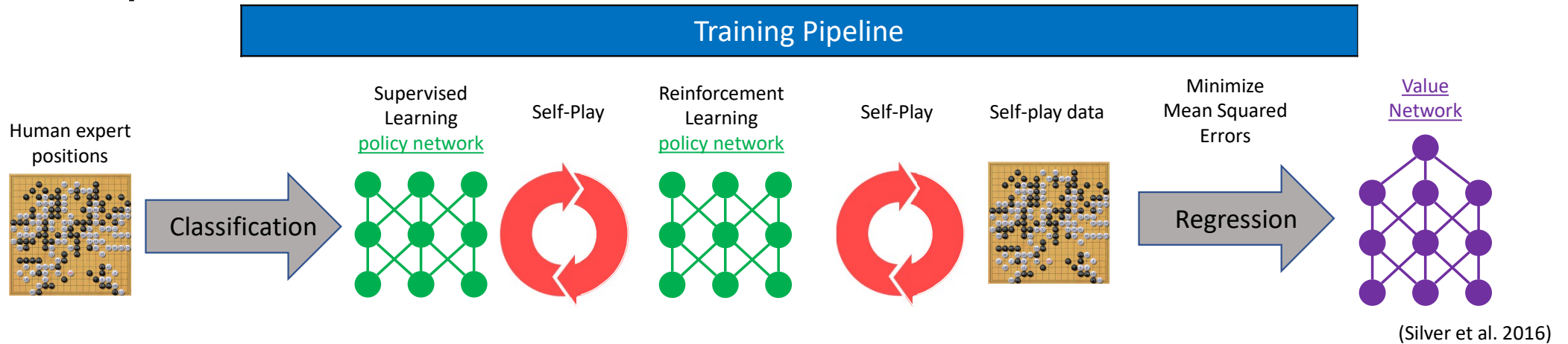- **Testing:** Mean Squared Errors of 0.37

### Solution and approach

- A new self-play dataset consisting of 30 million distinct positions, each sampled from a separate game
- Each game between the RL policy network and itself
- MSEs: 0.226 (Training) and 0.234 (Testing)



(Silver et al. 2016)

# 2. AlphaGo: Structure Breakdown

## Training Pipeline

Human expert positions

**Classification** →

Supervised Learning
policy network

Self-Play

Reinforcement Learning
policy network

Self-Play

Self-play data

Minimize Mean Squared Errors

**Regression** →

Value Network

(Silver et al. 2016)

## Testing/Tournament

Monte Carlo Tree Search (MCTS)

t
Current board position

t+1
If I choose an action

⋮ ⋮
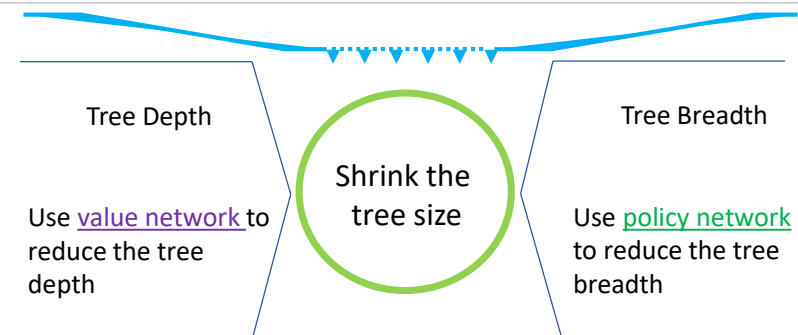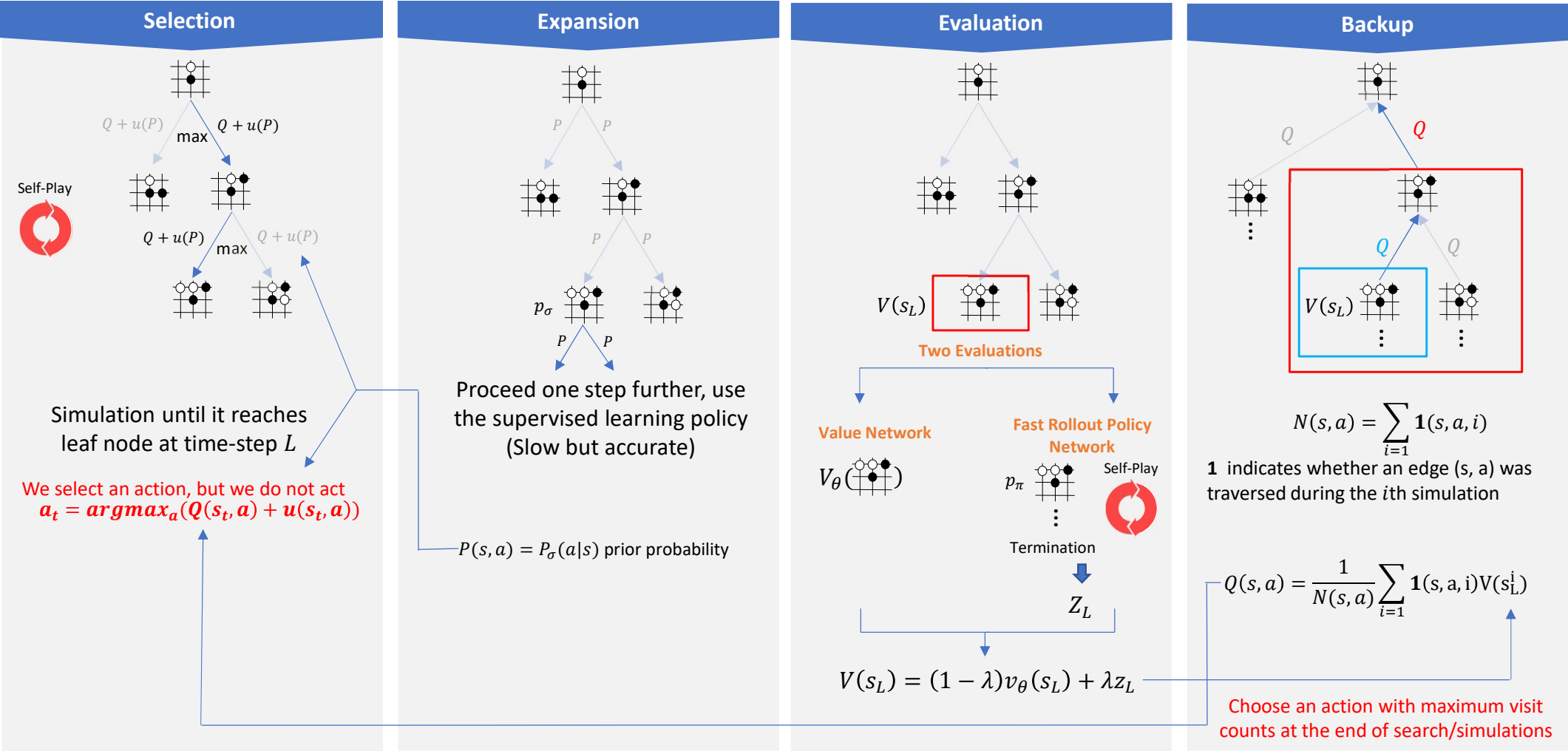
(Silver et al. 2016)

### Problems with pure MCTS

- Tree depth: need to simulate all the way to the end of the game

- Tree breadth: enormous amount of choices per time step

Tree Depth

Shrink the tree size

Tree Breadth

Use value network to reduce the tree depth

Use policy network to reduce the tree breadth

# 2.4 Monte Carlo Tree Search: Search with Policy and Value Networks



Asynchronous Policy and Value Monte Carlo Tree Search (APV-MCTS)

**Selection**

Self-Play

$Q + u(P)$   max   $Q + u(P)$

$Q + u(P)$   max   $Q + u(P)$

Simulation until it reaches
leaf node at time-step $L$

We select an action, but we do not act
$a_t = argmax_a(Q(s_t, a) + u(s_t, a))$

**Expansion**

$P$   $P$

$P$   $P$

$p_\sigma$

$P$   $P$

Proceed one step further, use
the supervised learning policy
(Slow but accurate)

$P(s, a) = P_\sigma(a|s)$ prior probability

**Evaluation**

$P$   $P$

$P$   $P$

$V(s_L)$

**Two Evaluations**

**Value Network**

$V_\theta(\quad)$

**Fast Rollout Policy
Network**

$p_\pi$

Self-Play

Termination

$Z_L$

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

**Backup**

$Q$   $Q$

$Q$   $Q$

$V(s_L)$

$$N(s, a) = \sum_{i=1} \mathbf{1}(s, a, i)$$

$\mathbf{1}$ indicates whether an edge (s, a) was
traversed during the $i$th simulation

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

Choose an action with maximum visit
counts at the end of search/simulations

# AlphaGo Zero

Simple Neural Network Structure
Completely tablua Rasa


An "interesting" fact about AlphaGo Zero:
The system consists of only 4 TPUs,
each of which is 15-30 times more efficient than GPUs in performance per-watt
And 64 GPUs, 19 CPUs.

# Improvements from AlphaGo

- No supervised learning
- Simplified input features (from 19x19x48 -> 19x19x17)
- Change network structure from CNN to ResNet.

# 3.1 Deep Neural Network- Structure Breakdown

**Single Neural Network Head**

**Policy Network**

$p_a = \Pr(a|s)$ probability of selecting each move from the current position

**But two outputs**

**Value Network**

$(p, v) = f_\theta(s)$

$v$ scalar evaluation, estimating the probability of current players winning from the current position

**Core RL Concepts**

Policy Evaluation → Value Function

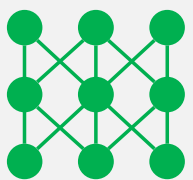**Policy Iteration**

Policy Improvement
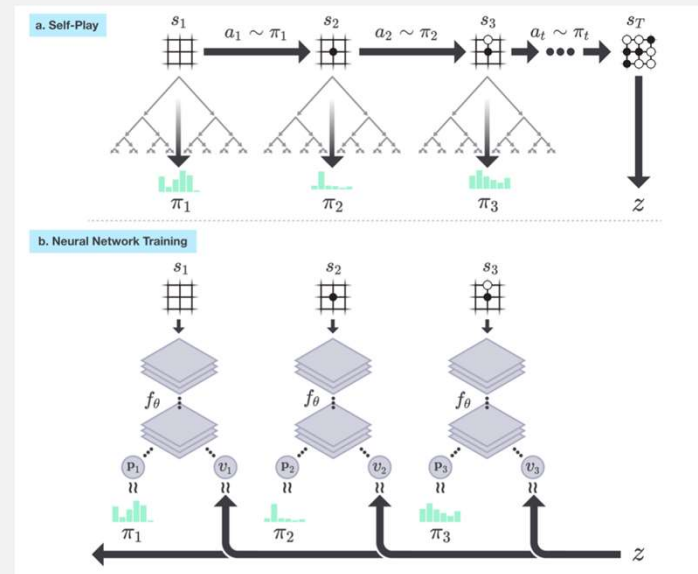
**Training Pipeline**

Reinforcement Learning

**Neural Network**

Self-Play

**Policy Evaluation**

MCTS $\Rightarrow \pi_t$, a move is selected according to **search policy** $a_t \sim \pi_t$
Terminal state $s_T$ will give us a game winner $Z$

a. Self-Play

$s_1 \quad a_1 \sim \pi_1 \quad s_2 \quad a_2 \sim \pi_2 \quad s_3 \quad a_t \sim \pi_t \quad s_T$

$\pi_1 \qquad \pi_2 \qquad \pi_3 \qquad z$

b. Neural Network Training

$s_1 \qquad s_2 \qquad s_3$

$f_\theta \qquad f_\theta \qquad f_\theta$

$p_1 \quad v_1 \quad p_2 \quad v_2 \quad p_3 \quad v_3$

$\pi_1 \qquad \pi_2 \qquad \pi_3 \qquad z$

At time t, feed the board position $s_t$ into policy $f_\theta(s) = f_\theta(s_t)$, which gives us $p_t$ and $v_t$.

**Policy Improvement**

Optimization problem: optimize weight parameters $\theta$ to minimize prediction errors $\boldsymbol{v_t} - \boldsymbol{Z_t}$ and maximize the similarity between neural network probability $\boldsymbol{p_t}$ and search probability $\boldsymbol{\pi_t}$

$\Rightarrow$ **Minimize Loss Function:** $l = (z - v)^2 - \pi^T \log \mathbf{p} + \mathbf{c}||\theta||^2$

# 3.2 Monte Carlo Tree Search: Search with a Single Neural Networ

## Asynchronous Policy and Value Monte Carlo Tree Search (APV-MCTS)

| Selection | Expansion and Evaluation | Backup | Play |
|---|---|---|---|

**Repeat**

**Selection:**

$Q + U(P)$  max  $Q + U(P)$

**Self-Play**

$Q + U(P)$  max  $Q + U(P)$

Simulation until it reaches leaf node at time-step $L$

We select but we do not act
$$a_t = argmax_a(Q(s_t, a) + U(s_t, a))$$

**Expansion and Evaluation:**

$V$

$P$  $P$

$V$  $V$

$P$  $P$

max

$$(P, V) = f_\theta(\quad)$$

$P$  $P$

$P(s, a) = P_\sigma(a|s)$ prior probability

$V$ is the probability value that I win from the current board position $S_t$ at time t (Value function)

**Backup:**

$Q$  $Q$

$Q$  $Q$

$$V(s_L)$$

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

$$Q(s, a) = \frac{1}{N(s, a)} \sum_{s'|s, a \to s'} V(s')$$

**Play:**

50%   50%

30%   30%

20%  10%  20%

$$\pi(a|s_0) = \frac{N(s_0, a)^{\frac{1}{\tau}}}{\sum_b N(s_0, b)^{\frac{1}{\tau}}}$$

Once an action is chosen, that sub-node becomes root, children remains and the rest of the old tree is discarded
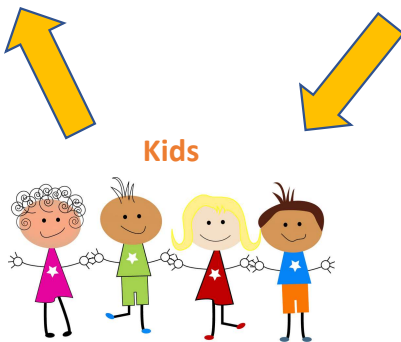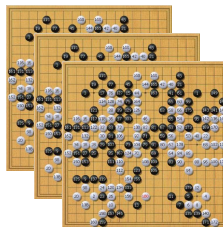
# 4. Humans, AlphaGo and AlphaGo Zero

## Humans

They develop new policies only when they become masters in Go

**Ancient Master / Experts**

**Dozens of manuals/books**
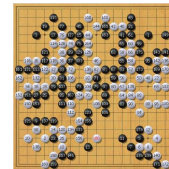Heritage from ancient knowledge in thousands of years

**Kids**

The problem is that there is limited evaluation process children will not (often are not allowed) to question ancient knowledge (mostly because of the Asian culture)

## AlphaGo

**Supervised Learning**

**Self-play**

Policy Evaluation

Humans are right

Humans are wrong

Policy Improvement

## AlphaGo Zero

1. Only basic rules
2. Start from completely random play

**Self-play**

Policy Evaluation

I am right

I am wrong

Policy Improvement

**Self-developed manuals**

- Concepts of shapes, territory, influence
- Joseki
- Fuseki
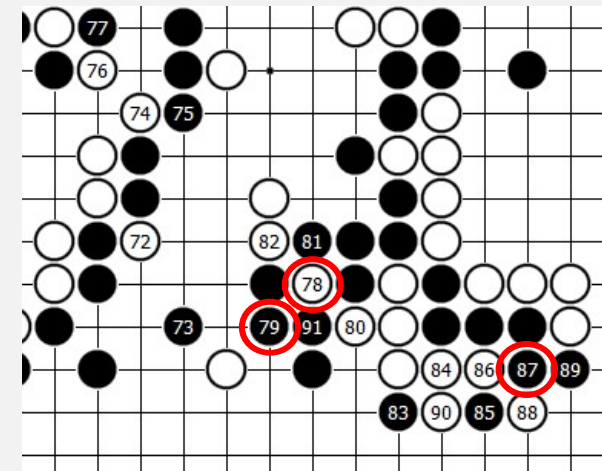- Tesuji
- Sente
- Shicho

# 5. Discussions

## Humans in Go game

- Humans are often wrong.
- Master manuals are more like guidance rather than mandatory knowledge.
- Self-evaluation and improvement are key to learning.
- The concept of tabula rasa is crucial in learning theories.

## AI in Go game

- Algorithms are much more important than data (model-based).
- People tend to assume that machine learning is all about big data and massive computations, which is wrong.
- AlphaGo Zero does not use any human data whatsoever, it always has the opponent at just right level (self-play), and improves itself from self-learning.
- Yet it performs much better, and significantly less computational requirement.

## The only game that AlphaGo resigned



- AlphaGo (Black) calculated that the probability of Lee Sedol making move 78 is 0.01%.
- Value network informed AlphaGo that move 79 had a winning probability of 70%.
- It realised (re-calculated) that the winning probability was 55% after move 87.

# 6. AlphaZero: General Version

## Generalize AlphaGo Zero in board games

### AlphaGo Zero vs. AlphaZero

| | AlphaGo Zero | AlphaZero |
|---|---|---|
| Value Function | Binary win/loss Probability of Winning | Expected outcome (Chess has win, loss, draw) |
| Rotation and Reflection | Invariance (faster training) | Variance (Asymmetric rules) |
| Policy Improvement | Current best player vs. New player after each iteration  If win by 55% then update the policy | Updated continuously, even during the iterations |

### Neural network architecture in different games

| | Chess | Shogi | Go |
|---|---|---|---|
| Input Feature | 119 | 362 | 17 |
| Policy Plane | 8 X 8 X 73 | 9 X 9 X 139 | 19 X 19 + 1 |
| Training Time | 9h | 12h | 34h |
| Training Games | 44 million | 24 million | 21 million |