# Master
# Reinforcement Learning 2022
# Lecture 6:
# Two-Agent Self-Play

Aske Plaat

liacs Leiden Institute of Advanced Computer Science

# Different Approaches

- Model-free

  - Value-based [2,3]

  - Policy-based [4]

- Model-based

  - Learned [5]

  - Perfect; Two-Agent [6]

- Multi-agent [7]

- Hierarchical Reinforcement Learning (Sub-goals) [8]

- Meta Learning [9]

# Motivation

# Overview

- MCTS: a well-known RL planner

- What if Internal Transition Function is Perfect
  & your Environment is yourself?
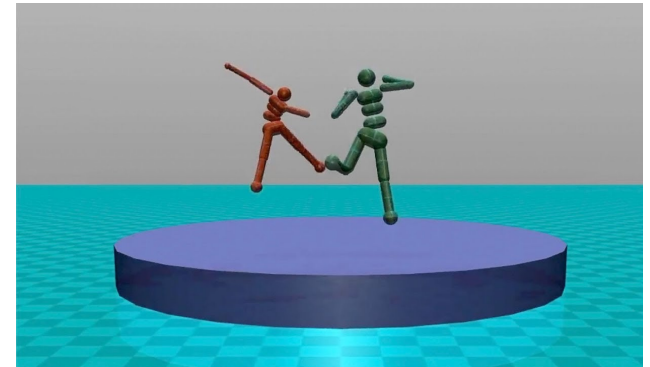
- AlphaZero

- Curriculum Learning

# What if Internal Transition function is Perfect?

- Previous chapter showed that accuracy of model is important.
  What if we have a perfect transition function?
  What if we can also use it to learn?

- Then World Champions get beaten:

    - Backgammon

    - Go

    - Chess

    - Shogi

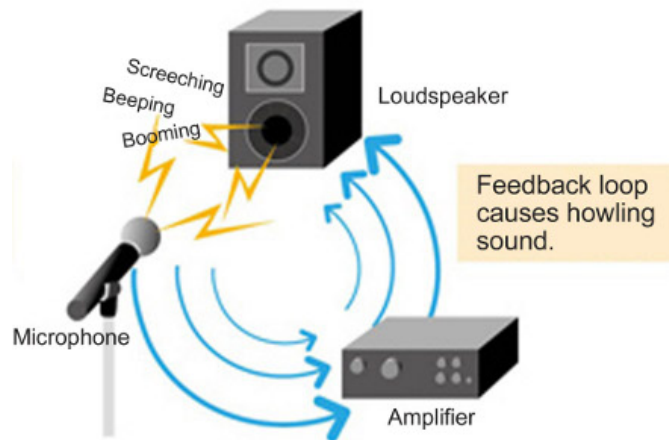- Today is about Best Case, when everything fits together and works

# Self-Play is old

- Most two-agent board game-playing programs choose (versions of) themselves as opponent for simulation or learning.

- Minimax (1949) is Self-Play

- Samuel's checkers players (1950-1960) used self-play outcomes for modifying evaluation weigths.

- TD-Gammon (1992) used Self-Play learning

- However, Self-Play is potentially unstable due to feedback and deadly triad

- It is overcome in AlphaGo in different ways

# Surprising Self-Play

- Find high quality examples to train RL on using RL. RL at three levels



- RL suffers from feedback, feedback creates instability

- What methods have been used to overcome feedback?

# AlphaZero:
# Three Levels of Self Play

1. **Move**-level self play (minimax, MCTS)

2. **Example**-level self play (learning, Actor Critic)

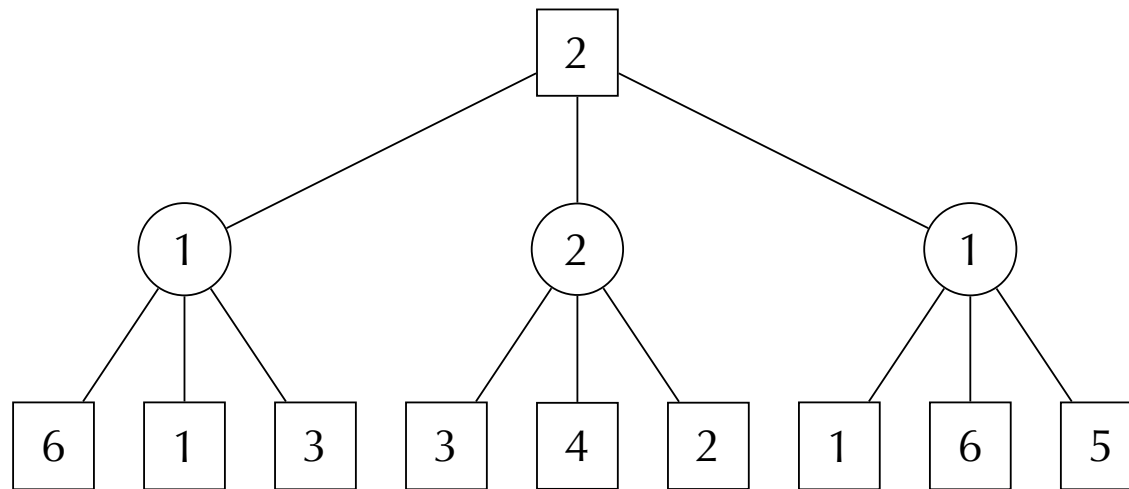3. **Game**-level self play (curriculum, self transcending player)

# 1. Move-level self play

# Minimax

- Assume you play best move, and opponent has your knowledge

# Minimax

```
                          ┌───┐
                          │ 2 │
                          └───┘
              ┌─────────────┼─────────────┐
            ( 1 )         ( 2 )         ( 1 )
           ┌──┼──┐       ┌──┼──┐       ┌──┼──┐
         ┌─┐ ┌─┐ ┌─┐   ┌─┐ ┌─┐ ┌─┐   ┌─┐ ┌─┐ ┌─┐
         │6│ │1│ │3│   │3│ │4│ │2│   │1│ │6│ │5│
         └─┘ └─┘ └─┘   └─┘ └─┘ └─┘   └─┘ └─┘ └─┘
```

- Two-agent zero sum: my win is your loss

- Max/min/max/min/max/min

- Max: Square

- Min: Circle

- b=branching factor, d=search depth

# Minimax

```python
INF = 99999

def eval(n):
    if n['type'] == 'LEAF':
        return n['value']
    else:
        error("Calling eval not on LEAF")

def minimax(n):
    if n['type'] == 'LEAF':
        return eval(n)
    elif n['type'] == 'MAX':
        g = -INF
        for c in n['children']:
            g = max(g, minimax(c))
    elif n['type'] == 'MIN':
        g = INF
        for c in n['children']:
            g = min(g, minimax(c))
    else:
        error("Wrong node type")
    return g

print("Minimax value: ", minimax(root))
```
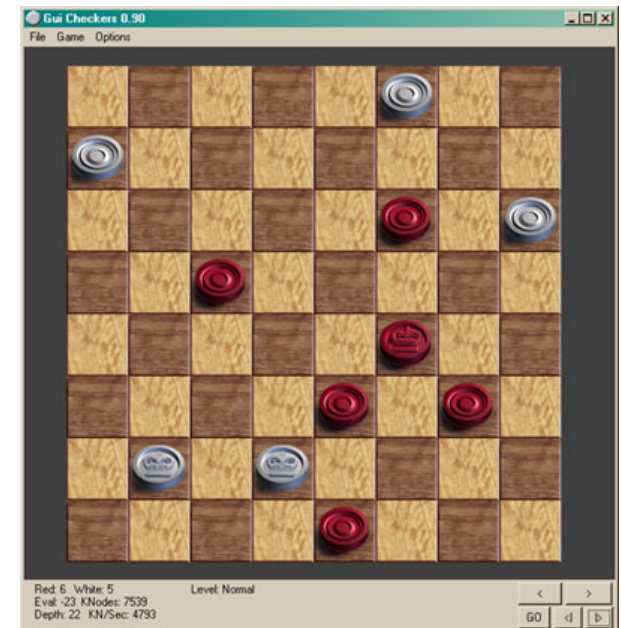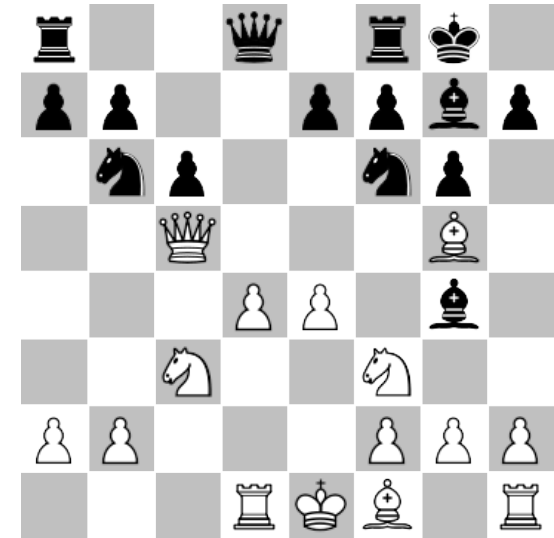
Listing 4.2: Minimax code

# State Space Complexity

| Name | board | state space | zero sum | information | turn |
|---|---|---|---|---|---|
| Chess | $8 \times 8$ | $10^{47}$ | zero sum | perfect | turn |
| Checkers | $8 \times 8$ | $10^{18}$ | zero sum | perfect | turn |
| Othello | $8 \times 8$ | $10^{28}$ | zero sum | perfect | turn |
| Backgammon | 24 | $10^{20}$ | zero sum | chance | turn |
| Go | $19 \times 19$ | $10^{170}$ | zero sum | perfect | turn |
| Shogi | $9 \times 9$ | $10^{71}$ | zero sum | perfect | turn |
| Poker | card | $10^{161}$ | non-zero | imperfect | turn |
| StarCraft | real time strategy | $10^{1685}$ | non-zero | imperfect | simultaneous |

- $10^{47}$… 1 ns/position -> $10^{38}$ s -> $10^{30}$ earth-year $10^{21}$ times the age of the known universe/position

# Heuristics



- Material (pawns, bishops, knights, …)

- Mobility (# actions)

- Center control

- King Safety

- …

# Alpha-Beta



- A cutoff is an action (e) that is so strong for my opponent that I will not play b (because a is better )  and hence we can stop searching b

# After Chess?

# Go!

# Go

| Name | board | state space | zero sum | information | turn |
|---|---|---|---|---|---|
| Chess | $8 \times 8$ | $10^{47}$ | zero sum | perfect | turn |
| Checkers | $8 \times 8$ | $10^{18}$ | zero sum | perfect | turn |
| Othello | $8 \times 8$ | $10^{28}$ | zero sum | perfect | turn |
| Backgammon | 24 | $10^{20}$ | zero sum | chance | turn |
| Go | $19 \times 19$ | $10^{170}$ | zero sum | perfect | turn |
| Shogi | $9 \times 9$ | $10^{71}$ | zero sum | perfect | turn |
| Poker | card | $10^{161}$ | non-zero | imperfect | turn |
| StarCraft | real time strategy | $10^{1685}$ | non-zero | imperfect | simultaneous |

- Worldwide most popular combinatorial game of strategy

- Much more complex than Chess

# Heuristic Planning

- Successful in games with tactical play where efficient heuristics can be found

  - Pieces move, and material is a good indicator of the score

- In Go, board is large, pieces do not move, material is typically balanced, and "influence" turned out to be difficult to program efficiently
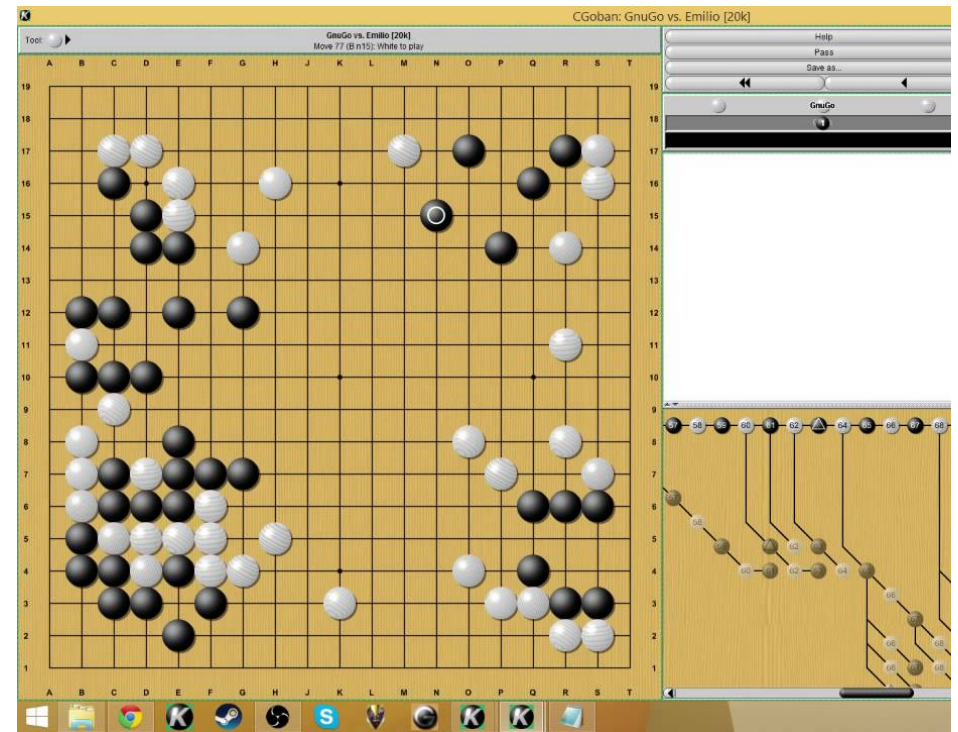


Now what?

# Traditional Go program

- Minimax

- Forward Pruning: only try "sensible" actions:
  connect, defend, territory jump
  Like a knowledge-based expert system

- Influence calculation for scoring

- Weak amateur level (10 kyu)

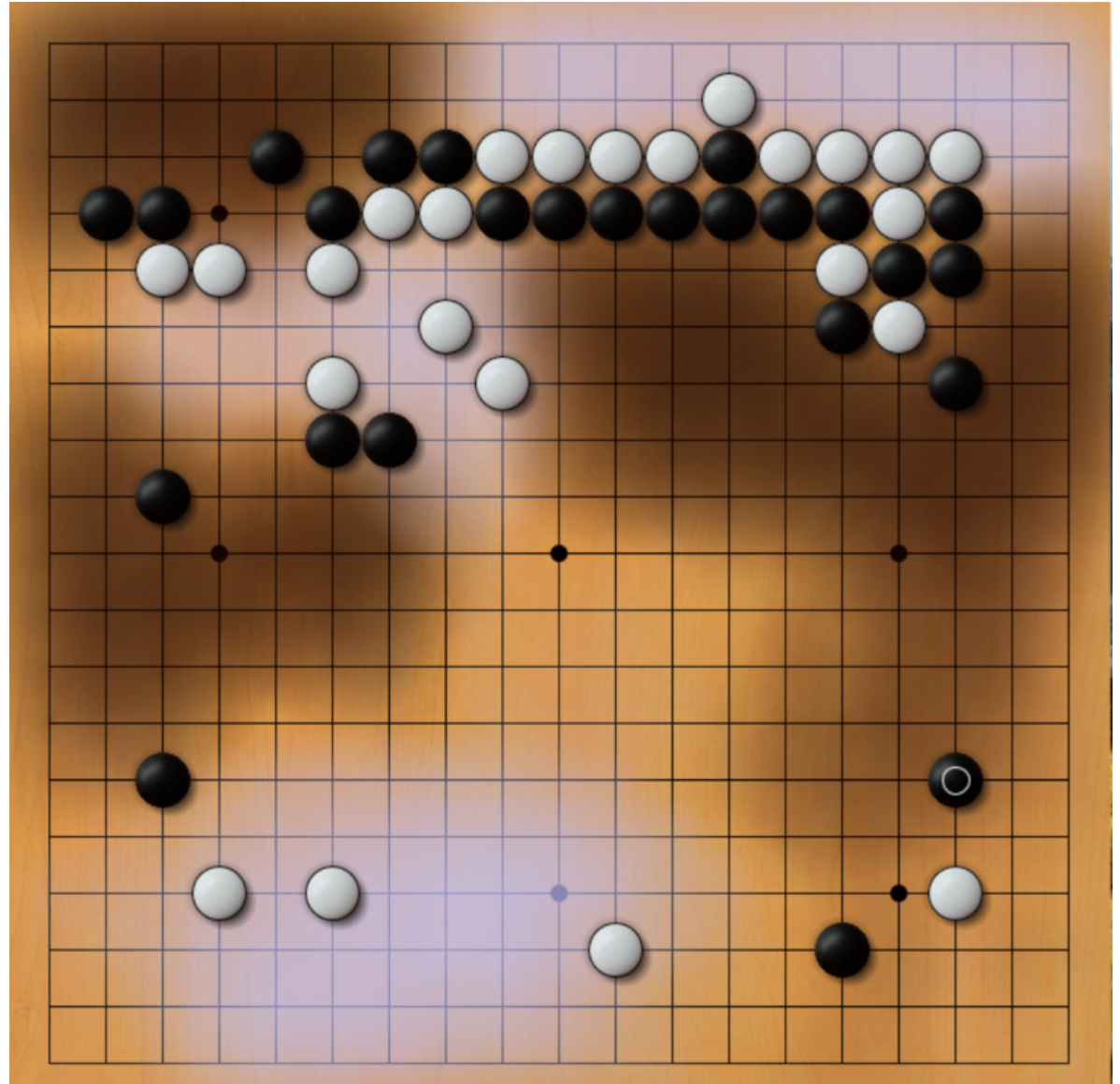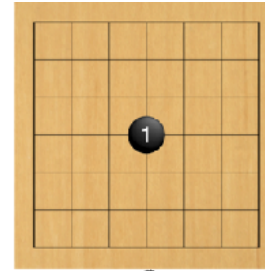- Years of no real progress

- Then: MCTS

# MCTS

# Adaptive Sampling

- No Efficient Heuristics for influence

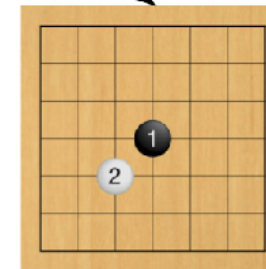- Rigid Search does not work well in large, flat state space
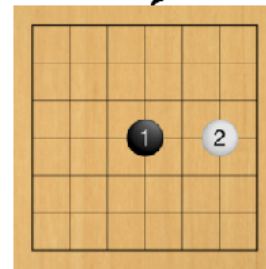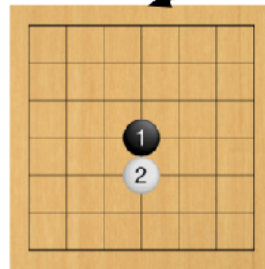
# Monte Carlo playouts



**Current Game**
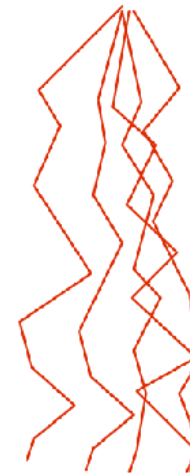
**Possible Moves**

**Monte Carlo Playouts**

**Wins**            53%            26%            37%            . . .

# Monte Carlo Playouts

- Chess: b=10. Full width

- Go: b=200. Forward pruning

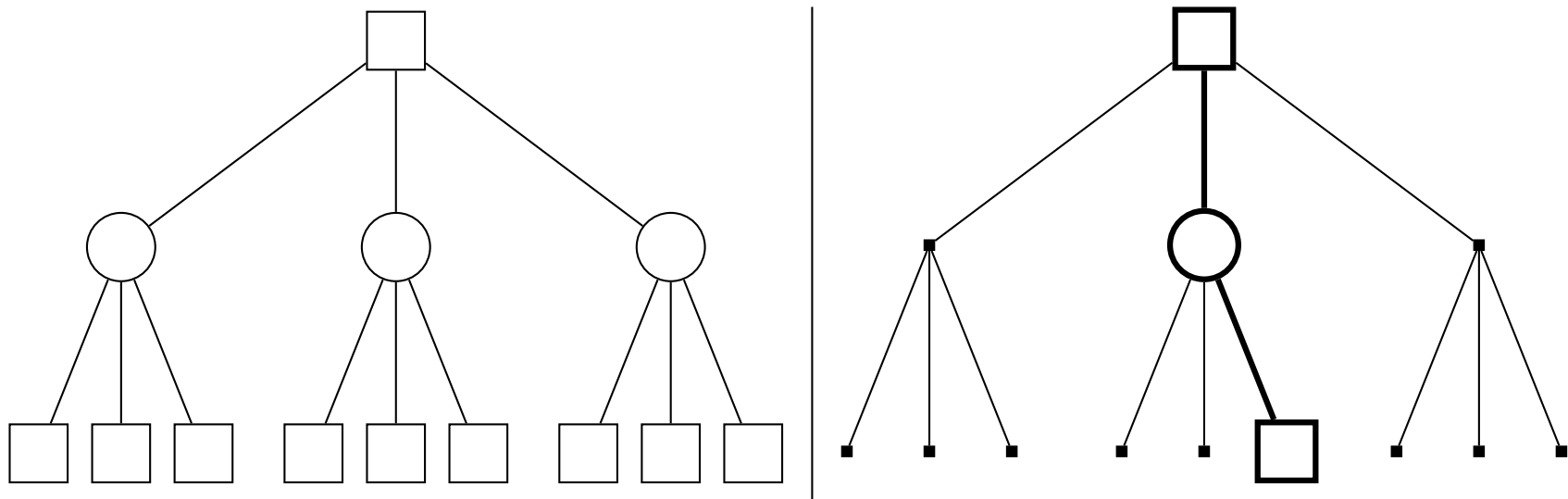- Playout: Not search **Tree b^d** but search **Path d**

Figure 4.1: Searching a Tree vs a Path

# Monte Carlo vs Minimax

- Minimax: Best of all actions

- Monte Carlo: Average of random playouts



Figure 4.1: Searching a Tree vs a Path

# Monte Carlo Tree Search

- Brügmann 1993 tried all playouts from the root (**flat**)
  Results were better than random, but **not great** -> MC

- Coulom 2006 (after work of others) tried it **recursively**, in a tree. This did give **good results** -> MCTS

- Kocsis & Szepesvari 2006 suggested the UCT selection rule to balance exploration and exploitation.
  (Based on extensive work in multi-armed bandit theory)
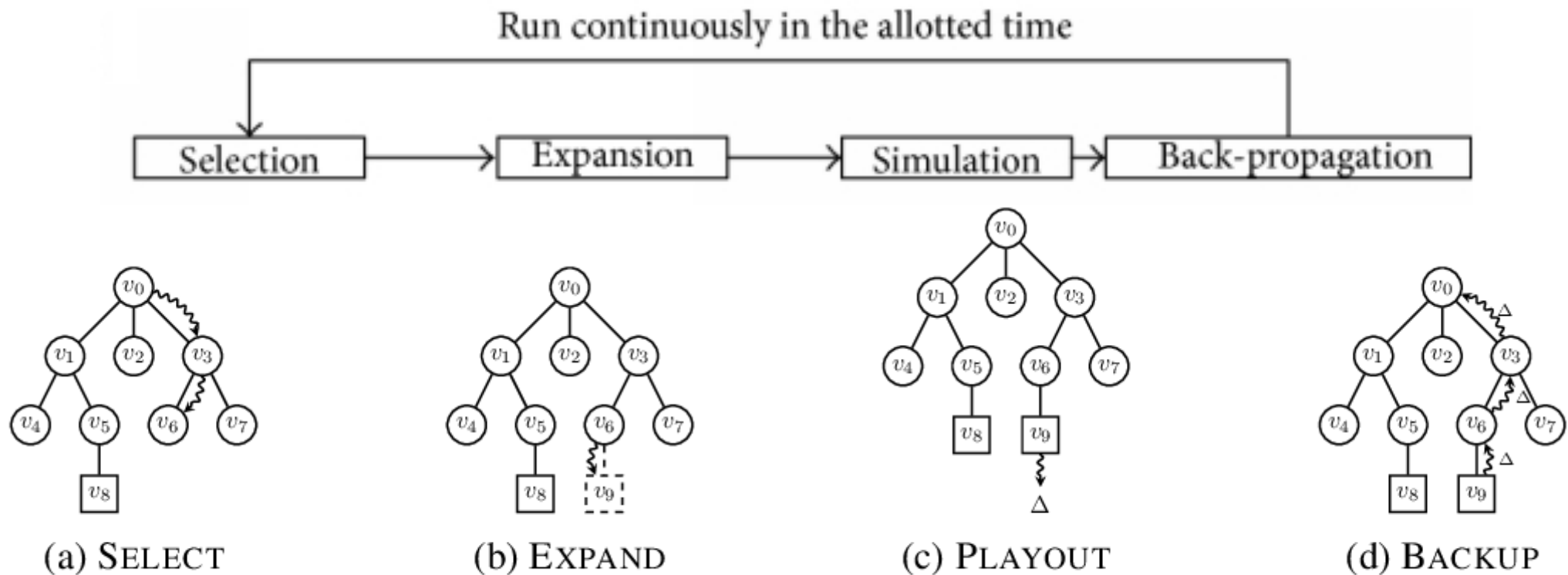
# Four Operations



Run continuously in the allotted time

Selection → Expansion → Simulation → Back-propagation

(a) SELECT  (b) EXPAND  (c) PLAYOUT  (d) BACKUP

Figure 1: One iteration of MCTS.

# UCT Selection formula

$$\text{UCT}(j) = \frac{w_j}{n_j} + C_p \sqrt{\frac{\ln n}{n_j}}$$

- $\text{UCT}_j = \text{winrate}_j + \text{Cp} * \text{newness}_j$

- Select the child j with the highest UCT value

- Winrate is for exploitation of what is known to be good

- Newness is for exploration of lesser-searched subtrees

- Larger Cp means more exploration

- **U**pper **C**onfidence bounds applied to **T**rees

# Example

# Example



- 1

- Select r

- expand a, playout a -> d=+, update a&r win++, visit++

- 2

- Select r

- expand b, playout b -> e=-, update b&r, visit++

- 3

- select UCT(a) 1/1+1*sqrt(ln 1/1) = 1, UCT(b) 0/1+1*sqrt(ln 1/1)=0

- Playout a -> c=+, update a&r, win++, visit++

- 4

- Select r

- select UCT(a) 2/2+1*sqrt(ln 2/2) = 1.588, UCT(b) 0/1+1*sqrt(ln 2/1)=0.832

- expand c, playout c=+, update c&a&r, win++, visit++

# Tree Shape

# Impact MCTS playing strength Go programs

# MCTS Go



Go AI Strength History

# Questions?

# 2. Example-level self play

# AlphaGo 1 2 3



1. AlphaGo: The Champion

2. AlphaGo Zero: Tabula Rasa
The Self-Learner



3. AlphaZero: Three games: Chess, Shogi, Go.
The Generalist

# AlphaGo Structure

- 4 nets

  - fast rollout policy

  - slow sl policy

  - slow rl policy

  - value net

- 3 learning methods

  - supervised small patterns fast rollout policy

  - supervised database grandmaster games

  - reinforcement from database de-correlated self-play games



Rollout policy $p_\pi$    SL policy network $p_\sigma$    RL policy network $p_\rho$    Value network $\nu_\theta$

Neural network

Policy gradient

Classification   Classification   Self Play   Regression

Traditional way to collect training data

Human expert positions

Self-play positions

Data

Generating a lot of additional training data by self-play.

# Policy & Value & Playout



Figure 1: One iteration of MCTS.

(a) SELECT

(b) EXPAND

(c) PLAYOUT

(d) BACKUP

# AlphaGo Structure

## AlphaGo Overview

**Input Board Position
as 19 x 19 Image
48 Feature Planes**

**Policy Network**

Expert Games

130 000 Games
30 M Positions

Supervised Learning

SL Policy
Position --> Next Move
Accuracy: 56%

**Fast Policy Network**

Expert Games

140 000 Patterns
130 000 Games
30 M Positions

Supervised Learning

Fast Policy
Pattern --> Next Move
Accuracy: 24%

**Reinforcement Learning
Policy Network**

Self-Play Games

1.3 M Games by
various versions
of RL Policy

Reinforcement Learning

RL Policy
Position --> Next Move
Wins 80% vs. SL Policy

**Value Network**

Self-Play Games

30 M Positions
by fixed version
of RL Policy

Reinforcement Learning

Position --> Win Probability
15 000 times faster than
MCTS Rollouts evaluations

**AlphaGo**

**Monte Carlo
Tree Search**

# AlphaGo Performance

# AlphaGo Matches



- 2015 Fan Hui London

- 2016 Lee Sedol Seoul

- 2017 Ke Jie Wuzhen

# nature

THE INTERNATIONAL WEEKLY JOURNAL OF SCIENCE

*At last* — a computer program that can beat a champion Go player PAGE 484

## ALL SYSTEMS GO

**CONSERVATION**
SONGBIRDS
À LA CARTE
*Illegal harvest of millions of Mediterranean birds*
PAGE 452

**RESEARCH ETHICS**
SAFEGUARD
TRANSPARENCY
*Don't let openness backfire on individuals*
PAGE 459

**POPULAR SCIENCE**
WHEN GENES
GOT 'SELFISH'
*Dawkins's calling card 40 years on*
PAGE 462

NATUREASIA.COM
28 January 2016
Vol 529 No. 7587

## ARTICLE

### Mastering the game of Go with deep neural networks and tree search

## ARTICLE

### Mastering the game of Go without human knowledge

# AlphaGo Zero

- Faster
  days, not weeks

- Better
  Higher Elo

- Elegant
  1 network

# Self-Play Loop

- Generate a sequence of own training examples

games: *self play*



tournaments: $(s, \pi, z)$ *examples*

```
1  for it in range (1, max_iterations): # do a curric. of self-play tourn.
2      for game in range(1, max_games): # play a tourn. of games; then train
3          trim(triples) # if buffer full: replace old entries
4          while not game_over(): # generate the moves of one game
5              game_pairs += mcts(eval(net)) # move is a (state, action) pair
6          triples += add(games_pairs, game_outcome()) # add win/lose to buf
7      net = train(net, triples) # retrain with (state, action, outc) triples
```

# AlphaGo Zero Overview

- Zero-knowledge

- One net (double-headed)

- One learning method: Self-Play

- Tabula Rasa: Only the rules & input/output layers, zero heuristics, zero grandmaster games

- Curriculum learning

# AlphaGo Zero Performance



**0 days**
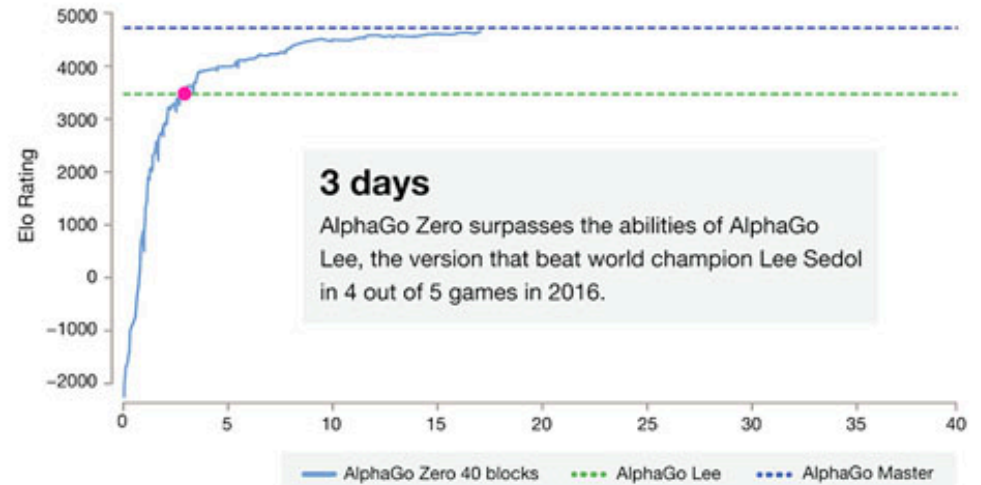AlphaGo Zero has no prior knowledge of the game and only the basic rules as an input.

AlphaGo Zero 40 blocks ···· AlphaGo Lee ···· AlphaGo Master

**3 days**
AlphaGo Zero surpasses the abilities of AlphaGo Lee, the version that beat world champion Lee Sedol in 4 out of 5 games in 2016.

AlphaGo Zero 40 blocks ···· AlphaGo Lee ···· AlphaGo Master

**21 days**
AlphaGo Zero reaches the level of AlphaGo Master, the version that defeated 60 top professionals online and world champion Ke Jie in 3 out of 3 games in 2017.

AlphaGo Zero 40 blocks ···· AlphaGo Lee ···· AlphaGo Master

**40 days**
AlphaGo Zero surpasses all other versions of AlphaGo and, arguably, becomes the best Go player in the world. It does this entirely from self-play, with no human intervention and using no historical data.

AlphaGo Zero 40 blocks ···· AlphaGo Lee ···· AlphaGo Master

# AlphaGo Zero Performance

# AlphaGo Zero Performance

# AlphaGo Zero Structure

- 1 net: ResNet with policy head and value head
  Combined loss-function

- 1 learning: RL Self-Play

- Tabula Rasa

# AlphaGo Zero Structure

Input: Board state (encoded)



$f_\theta$

$p_3$    $v_3$

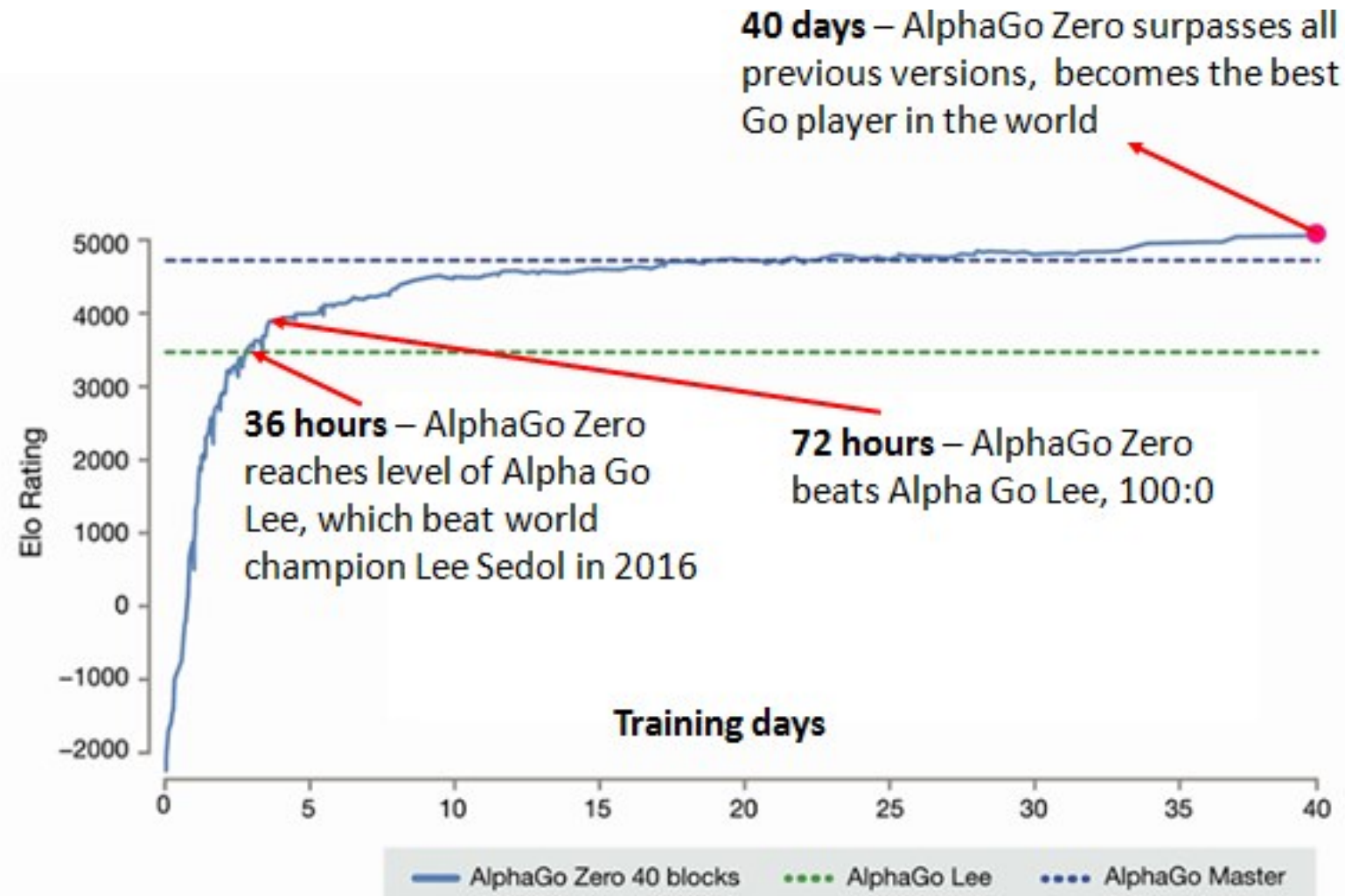Outputs: $P$ – Policy, $v$ - value

## One convolution block

128 filters (3X3 kernel, stride 1) + Batch norm + relu

## 19 Res blocks

Each block has 128 filters (3X3 kernel, stride 1) + Batch norm + relu + 128 filters (3X3 kernel, stride 1) + Batch norm + residual connection + relu

## One Output Block

Policy: convo of 32 filters (1X1 kernel, stride 1) + batch norm + relu + linear + softmax

Value: convo of 3 filters (1X1 kernel, stride 1) + batch norm + relu + linear + relu + linear + tanh

# AlphaGo Zero Networks



AG0: Comparison of Various Neural Network Architectures

[Silver et al. 2017b]

# AlphaGo Zero

- One net

- No Random Playout

- No Games database

# Two Questions

- Why Faster than AlphaGo?

- How can Self-Play (RL using RL for examples) ever be Stable? Deadly Triad Squared!

# AlphaGo Zero

- Stable

  - Extra Exploration

  - De-correlation

- How?

  - MCTS & Noise & Exploration & Replay Buffer & Many games

  - AlphaGo Zero's nets are not optimized against themselves, but against MCTS-improved versions of themselves

# AlphaGo Zero

- How Faster?

  - Curriculum learning

# 3. Game-level self play

# Curriculum Learning

- AlphaGo Zero learns better than AlphaGo

- AlphaGo Zero learns faster than AlphaGo. Why?

- Curriculum learning: start with easy examples

- Many small steps are faster than one large step

# Learning to Play



## 3 hours

AlphaGo Zero plays like a human beginner, forgoing long term strategy to focus on greedily capturing as many stones as possible.

Captured Stones

# Learning to Play



**19 hours**

AlphaGo Zero has learnt the fundamentals of more advanced Go strategies such as life-and-death, influence and territory.

# Learning to Play



## 70 hours

AlphaGo Zero plays at super-human level. The game is disciplined and involves multiple challenges across the board.

(68) at (61)

Captured Stones

# Curriculum

# AlphaZero
# General

# Science

AAAS

A DIGITAL
# PRODIGY
AlphaZero teaches
itself chess, shogi, and Go

# AlphaZero Overview

- Same net, same search, same tabula rasa self-play

- Different Input/Output layers

- Go, Chess, Shogi

# AlphaZero Structure

Input: Board state (encoded)

## One convolution block

128 filters (3X3 kernel, stride 1) + Batch norm + relu

## 19 Res blocks

Each block has 128 filters (3X3 kernel, stride 1) + Batch norm + relu + 128 filters (3X3 kernel, stride 1) + Batch norm + residual connection + relu

$f_\theta$

$p_3$ $v_3$

## One Output Block

Policy: convo of 32 filters (1X1 kernel, stride 1) + batch norm + relu + linear + softmax

Value: convo of 3 filters (1X1 kernel, stride 1) + batch norm + relu + linear + relu + linear + tanh

Outputs: $P$ – Policy, $v$ - value

# AlphaZero Performance

| Chess | Shogi | Go |
|---|---|---|



AlphaZero vs. Stockfish

AlphaZero vs. Elmo

AlphaZero vs. AG0

| | Chess | Shogi | Go |
|---|---|---|---|
| White (top bar) | W:29.0%  D:70.6%  L:0.4% | W:84.2%  D:2.2%  L:13.6% | W:68.9%  L:31.1% |
| Black (bottom bar) | W:2.0%  D:97.2%  L:0.8% | W:98.2%  D:0.0%  L:1.8% | W:53.7%  L:46.3% |

AZ wins  ■  AZ draws  ■  AZ loses  ■  AZ white ◯  AZ black ●

# AlphaZero Performance

# AlphaZero Conclusions

- First time learning, neural nets, and MCTS work in Chess

- Decades of heuristic planning research are surpassed

- Three Games share a general essence, since same architecture works (except I/O)

- Not same net. Net trained for Chess does not work for Shogi

- First architecture achieving very high performance in three games

# Curriculum and other learning

# Curriculum Learning & Friends

- **Learning** is Generalization from example to example

- **Curriculum** learning easy to hard concepts

- **Multi-task** learning two tasks at the same time

- **Transfer** learning from problem to problem

# Game State

# 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 this player lost - added once the game has finished)

## RETRAIN NETWORK
Optimise the network weights

### A TRAINING LOOP

Sample a mini-batch of 2,048 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 + $\pi$ — ACTUAL

v — Mean-squared error + 🏆

+ Regularisation

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

## WHAT IS A 'GAME STATE'

19 x 19 x 17 stack

Current position of black's stones

1 if black stone here
0 if black stone not here

...and for the previous 7 time periods
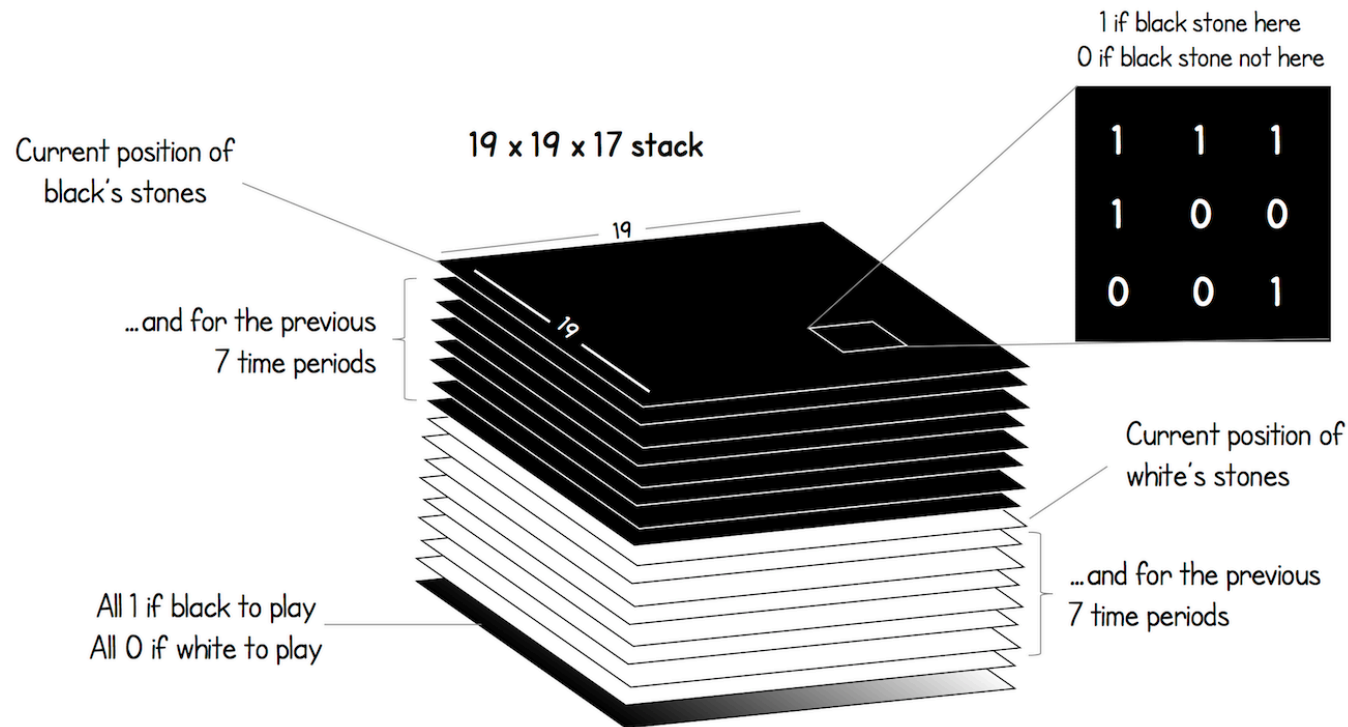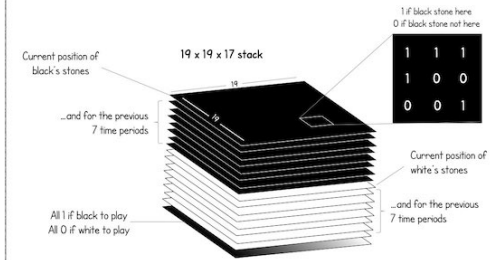
Current position of white's stones

...and for the previous 7 time periods

All 1 if black to play
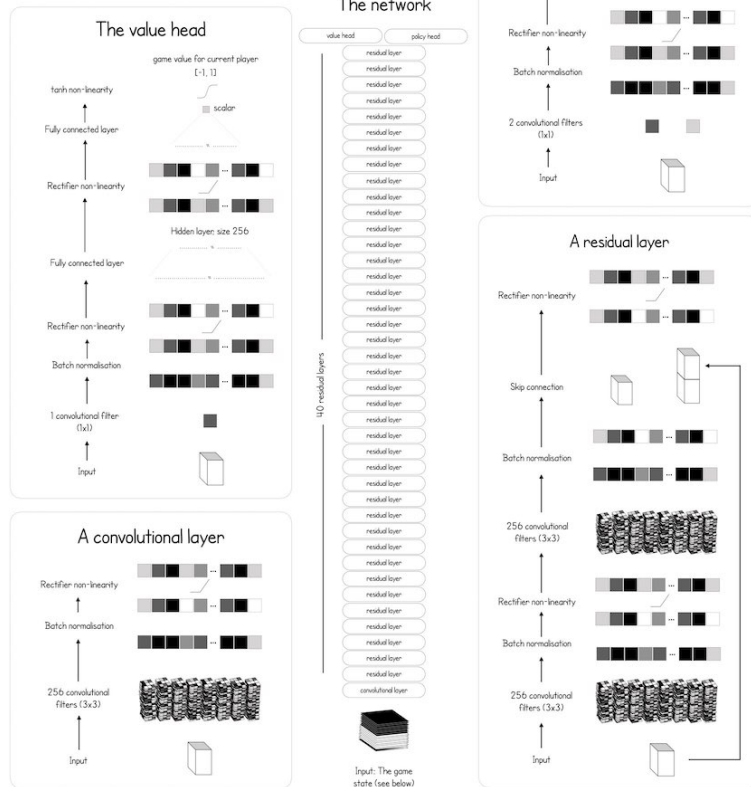All 0 if white to play

This stack is the input to the deep neural network

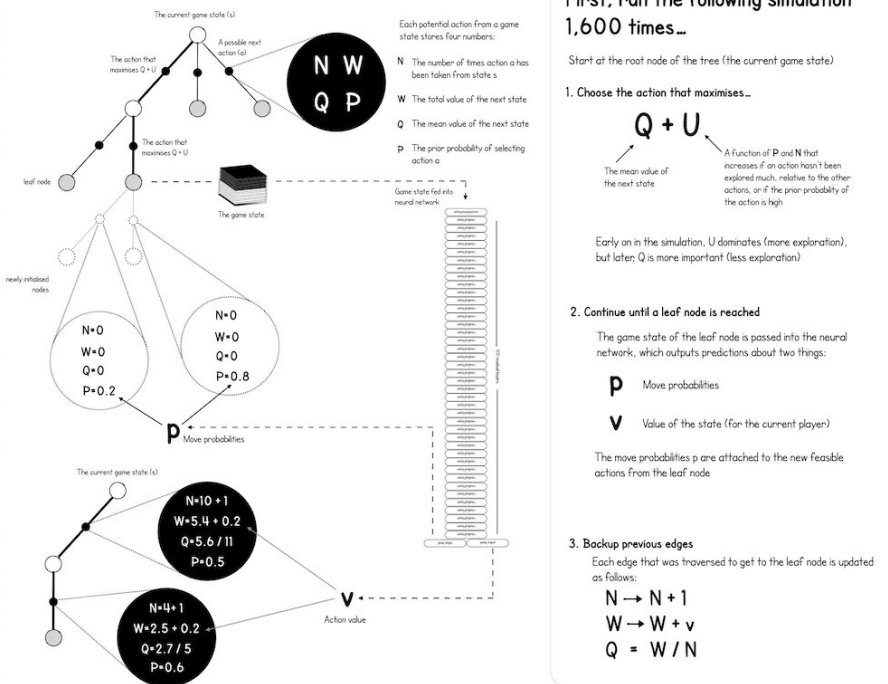## THE DEEP NEURAL NETWORK ARCHITECTURE
How AlphaGo Zero assesses new positions

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

scalar

Hidden layer size 256

Fully connected layer

Rectifier non-linearity

Batch normalisation

1 convolutional filter (1x1)

Input

### The network

value head    policy head

residual layer (×40)

Input: The game state (see below)

### A convolutional layer

Rectifier non-linearity

Batch normalisation

256 convolutional filters (3x3)

Input

### The policy head

19 x 19 + 1 (for pass) move logit 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

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

The current game state (s)

The action that maximises Q + U

A possible next action (a)

The action that maximises Q + U

leaf node

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

The game state

Game state fed into neural network

newly initialised nodes

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

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

p
Move probabilities

The current game state (s)

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:

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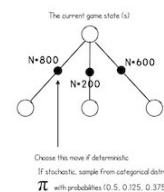$$Q = W / N$$

### ...then select a move

After 1,600 simulations, the move can either be chosen:

**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

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

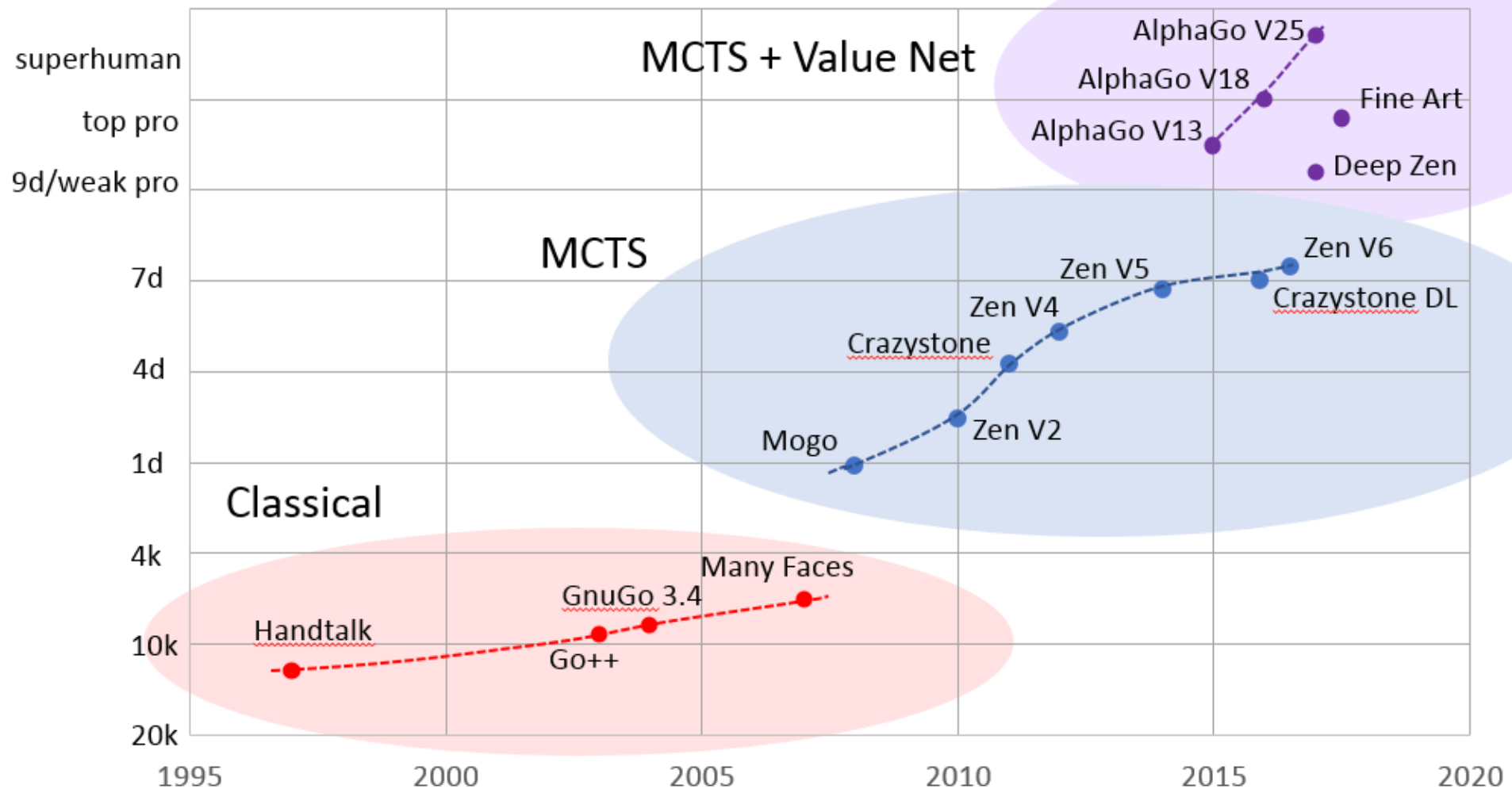where $\tau$ is a temperature parameter, controlling exploration

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)

### Other points

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

– The rest of the tree is discarded

# AlphaGo Performance



Go AI Strength History

# Open Source AlphaZero Reimplementations

- Leela

- ELF Facebook

- AlphaZero General Stanford

- PhoenixGo Tencent

- PolyGames Facebook