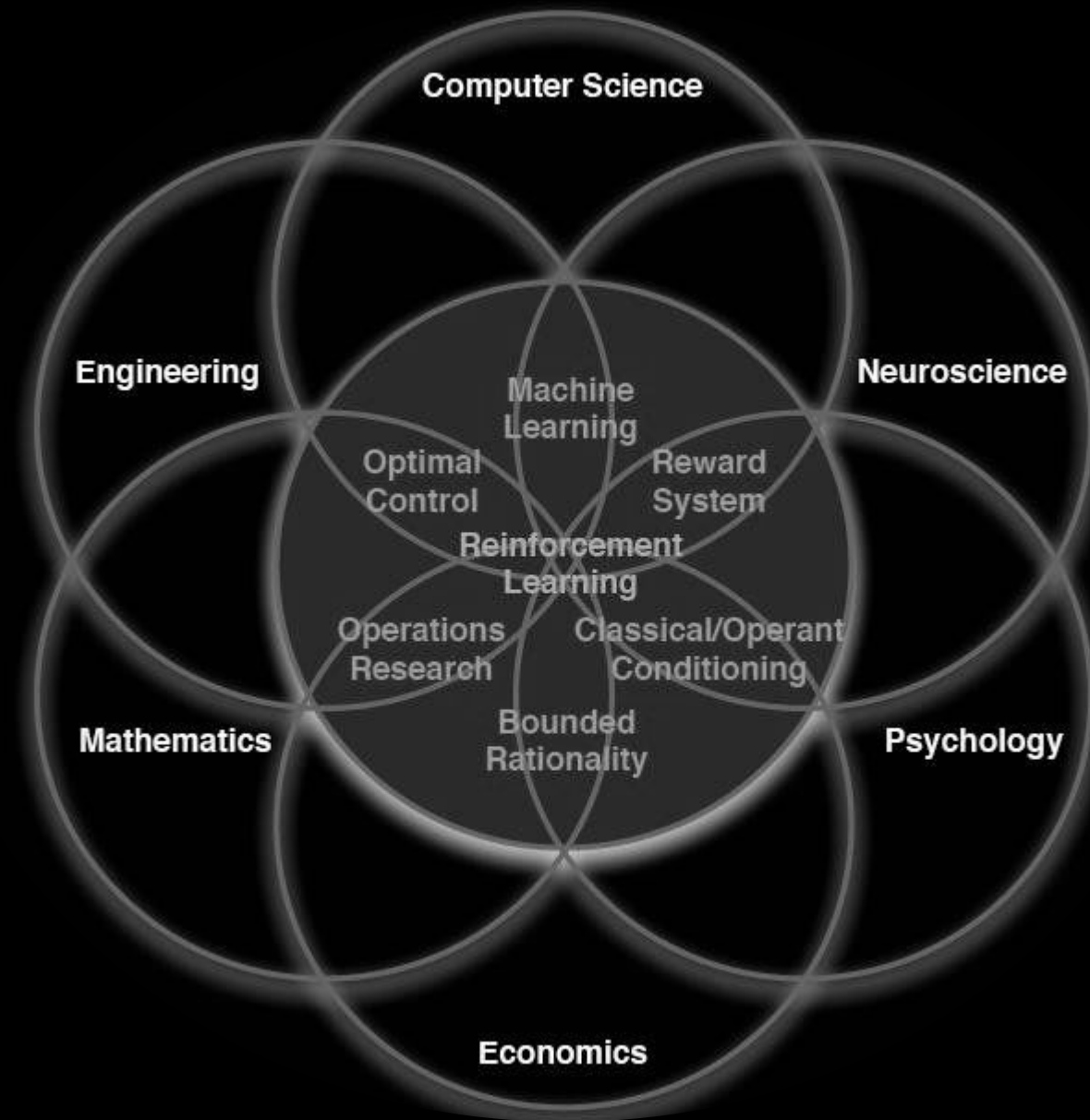
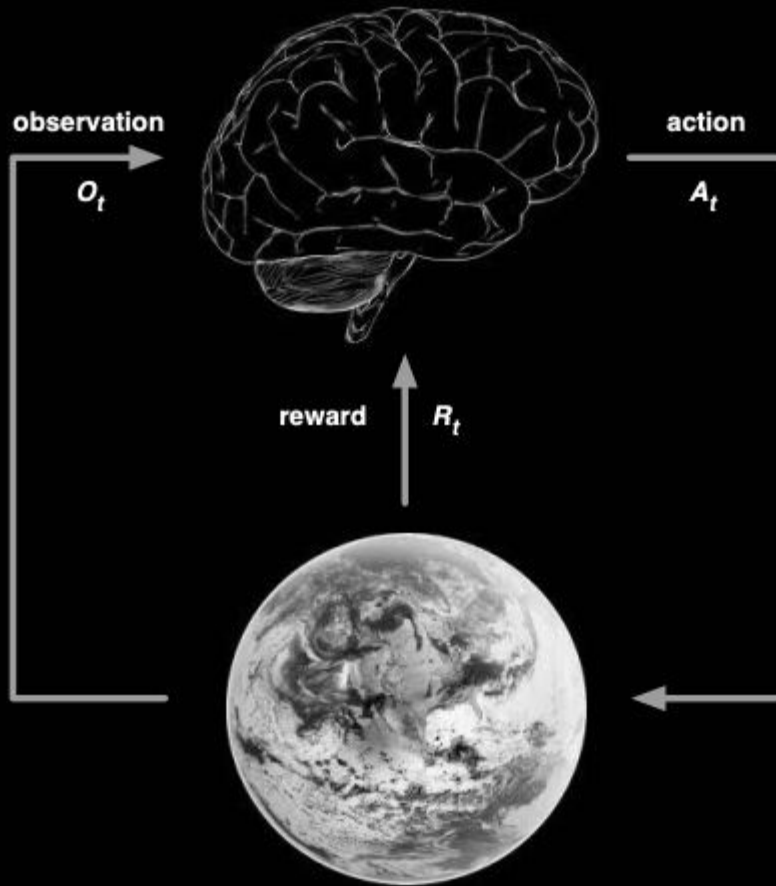


Reinforcement Learning

Overview of RL and Deep RL





R_t feedback at time t

SEQUENTIAL PARTIALLY OBSERVABLE ENVIRONMENT

Agent Goal: select actions to maximise total future reward

At each step t :

the agent:

Executes action A_t

Receives observation O_t

Receives scalar reward R_t

The environment:

Receives action A_t

Emits observation O_{t+1}

Emits scalar reward R_{t+1}

t increments at env. step

$$H_t = O_1, R_1, A_1, \dots, A_{t-1}, O_t, R_t$$

$$S_t = f(H_t)$$

environment state S_t^e - real state

agent state S_t^a - the agent's representation of state

TOOLS USED TO CHOOSE ACTIONS

Policy π function mapping state to action

Deterministic: $a = \pi(s)$

Stochastic: $\pi(a|s) = P[A_t = a | S_t = s]$

Value $v_\pi(s)$ a prediction of future reward

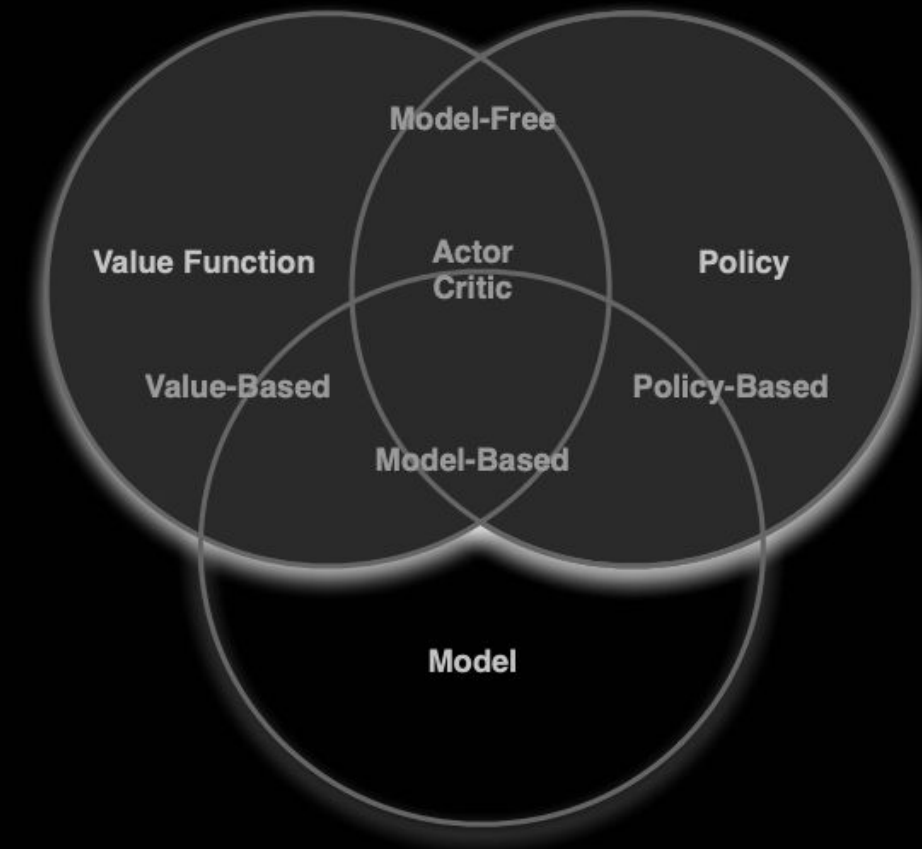
$$v_\pi(s) = E_\pi [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s]$$

Model predicts what the next state P and the next reward R

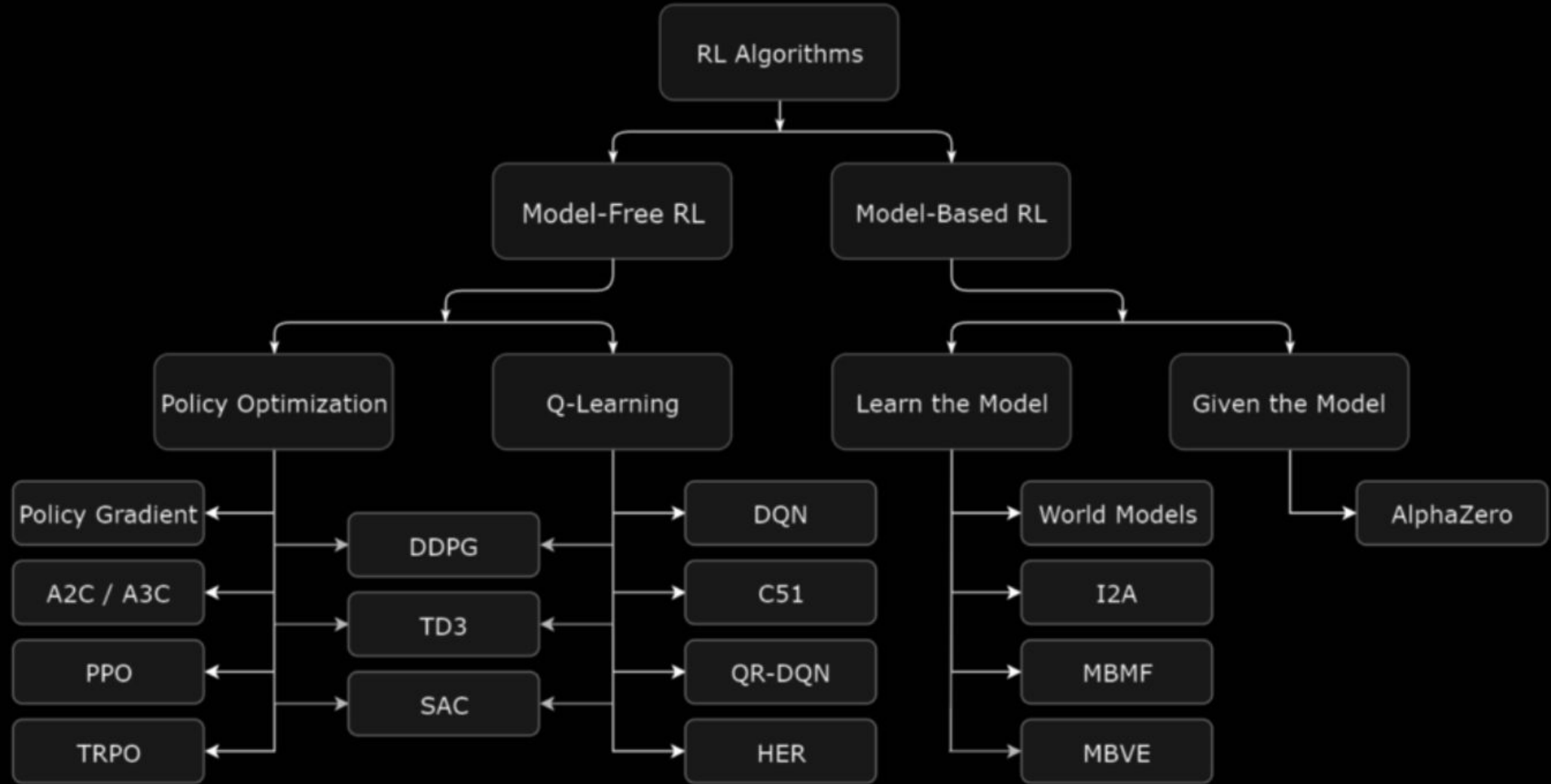
$$P_{ss'}^a = P[S_{t+1} = s' | S_t = s, A_t = a]$$

$$R_s^a = E[R_{t+1} | S_t = s, A_t = a]$$

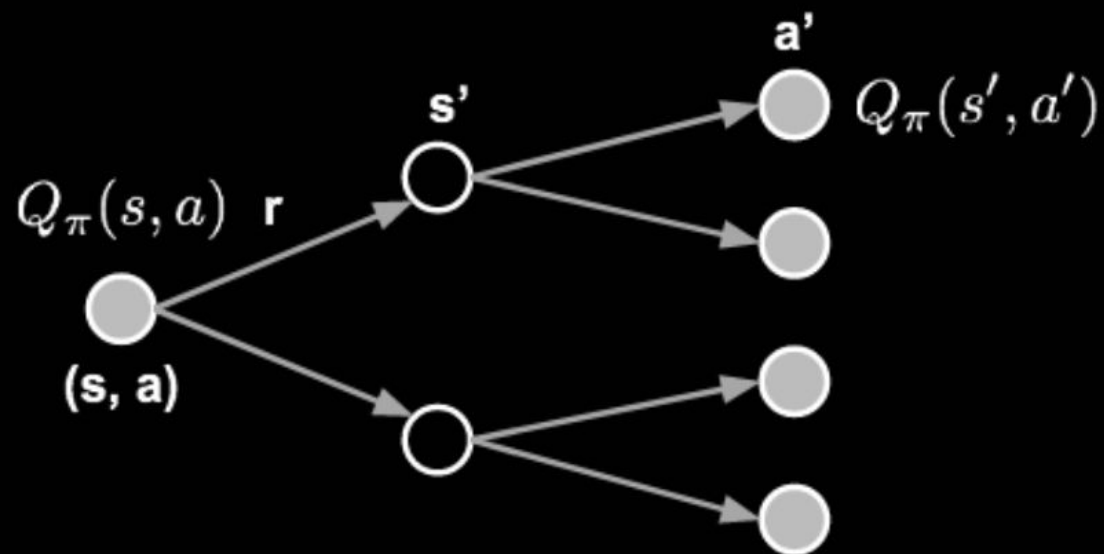
DIFFERENT TYPES OF RL



EXPLORATION VS EXPLOITATION



Q learning to AlphaZero



$$V_\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) Q_\pi(s, a)$$

$$Q_\pi(s, a) = R(s, a) + \gamma \sum_{s' \in \mathcal{S}} P_{ss'}^a V_\pi(s')$$

$$V_\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \left(R(s, a) + \gamma \sum_{s' \in \mathcal{S}} P_{ss'}^a V_\pi(s') \right)$$

$$Q_\pi(s, a) = R(s, a) + \gamma \sum_{s' \in \mathcal{S}} P_{ss'}^a \sum_{a' \in \mathcal{A}} \pi(a'|s') Q_\pi(s', a')$$

Q-learning: Learn function $Q : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$

Require:

States $\mathcal{X} = \{1, \dots, n_x\}$

Actions $\mathcal{A} = \{1, \dots, n_a\}$, $A : \mathcal{X} \Rightarrow \mathcal{A}$

Reward function $R : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$

Black-box (probabilistic) transition function $T : \mathcal{X} \times \mathcal{A} \rightarrow \mathcal{X}$

Learning rate $\alpha \in [0, 1]$, typically $\alpha = 0.1$

Discounting factor $\gamma \in [0, 1]$

procedure QLEARNING($\mathcal{X}, A, R, T, \alpha, \gamma$)

 Initialize $Q : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$ arbitrarily

while Q is not converged **do**

 Start in state $s \in \mathcal{X}$

while s is not terminal **do**

 Calculate π according to Q and exploration strategy (e.g. $\pi(x) \leftarrow \arg \max_a Q(x, a)$)

$a \leftarrow \pi(s)$

$r \leftarrow R(s, a)$ ▷ Receive the reward

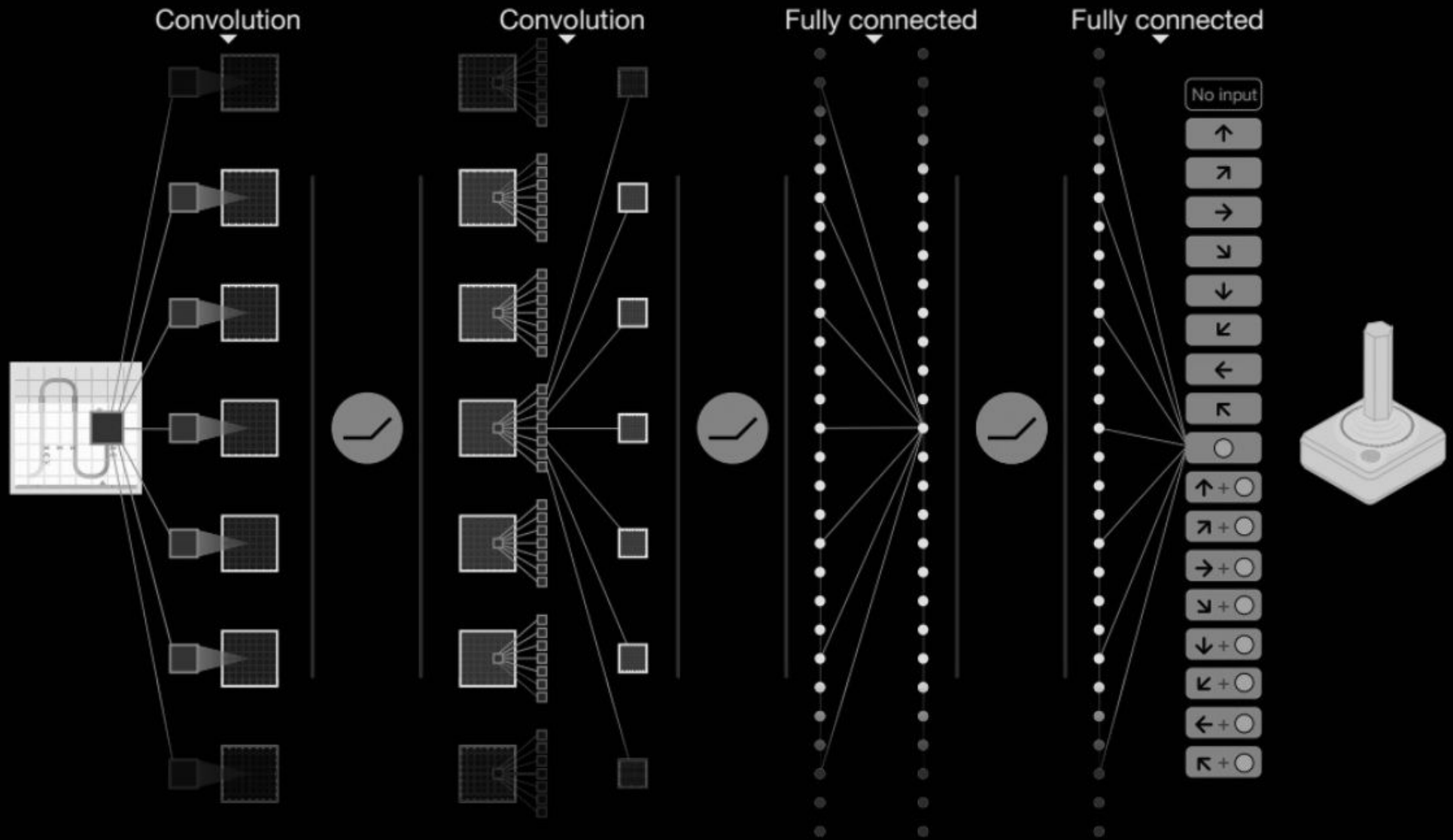
$s' \leftarrow T(s, a)$ ▷ Receive the new state

$Q(s', a) \leftarrow (1 - \alpha) \cdot Q(s, a) + \alpha \cdot (r + \gamma \cdot \max_{a'} Q(s', a'))$

$s \leftarrow s'$

return Q

PROBLEM: DOES NOT HAVE GENERALITY



Algorithm 1: deep Q-learning with experience replay.

Initialize replay memory D to capacity N

Initialize action-value function Q with random weights θ

Initialize target action-value function \hat{Q} with weights $\theta^- = \theta$

For episode = 1, M **do**

 Initialize sequence $s_1 = \{x_1\}$ and preprocessed sequence $\phi_1 = \phi(s_1)$

For $t = 1, T$ **do**

 With probability ε select a random action a_t

 otherwise select $a_t = \operatorname{argmax}_a Q(\phi(s_t), a; \theta)$

 Execute action a_t in emulator and observe reward r_t and image x_{t+1}

 Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$

 Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in D

 Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from D

 Set $y_j = \begin{cases} r_j & \text{if episode terminates at step } j+1 \\ r_j + \gamma \max_{a'} \hat{Q}(\phi_{j+1}, a'; \theta^-) & \text{otherwise} \end{cases}$

 Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ with respect to the network parameters θ

 Every C steps reset $\hat{Q} = Q$

End For

End For

EXPERIENCE REPLAY

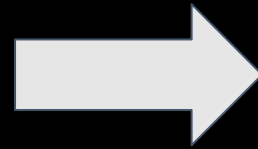
SAMPLE TARGET NETWORK

ALPHA GO

User Player knowledge
and learning to play

Knows the model of GO

Not generalizable

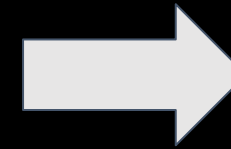


ALPHA GO Zero

learns to play from
scratch

Knows the model of GO

Not generalizable

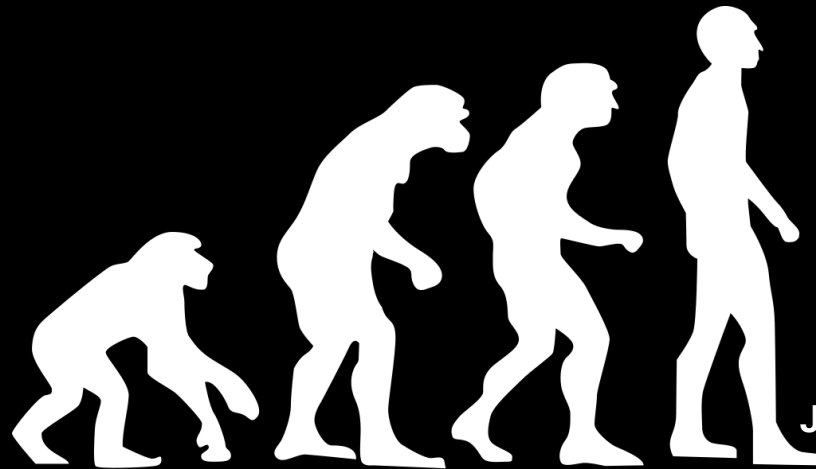


ALPHA Zero

learns to play from
scratch

Learns the model for any
game

Generalizable to minimax
games



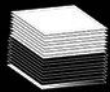
The training pipeline for AlphaGo Zero consists of three stages, executed in parallel

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



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

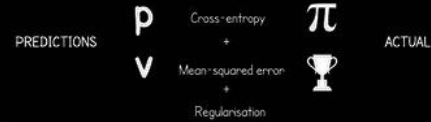
A TRAINING LOOP

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



After every 1,000 training loops, evaluate the network

EVALUATE NETWORK

Test to see if the new network is stronger

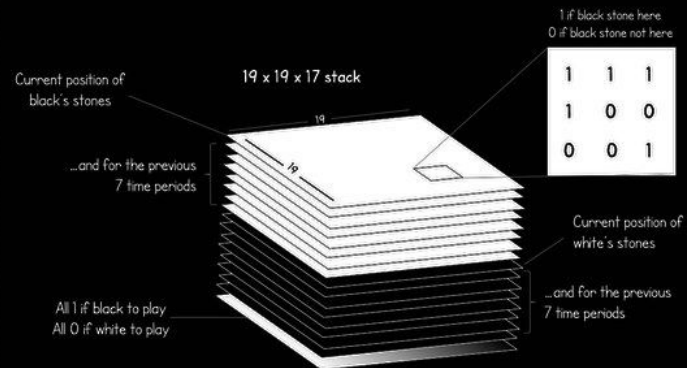
Play 400 games between the latest neural network and the current best neural network

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'



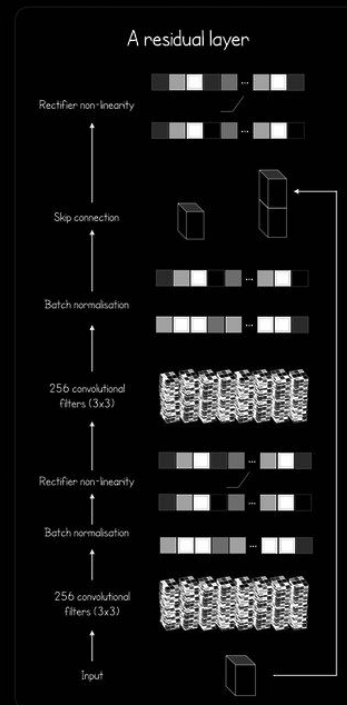
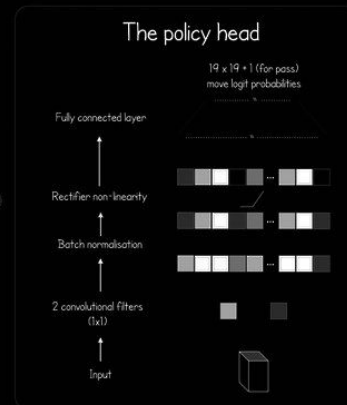
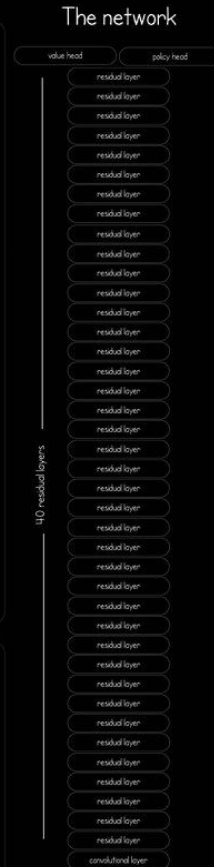
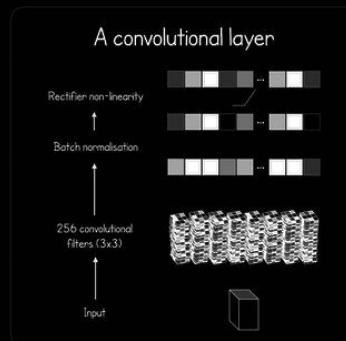
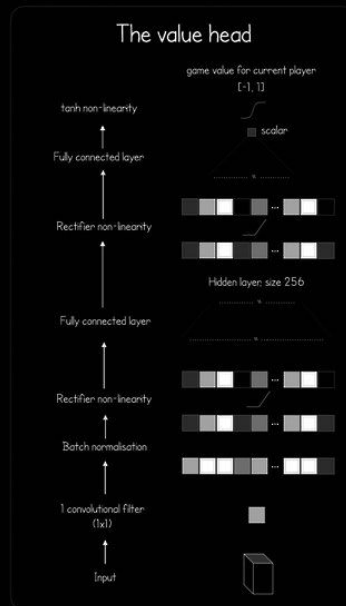
This stack is the input to the deep neural network

THE DEEP NEURAL NETWORK ARCHITECTURE

How AlphaGo Zero assesses new positions

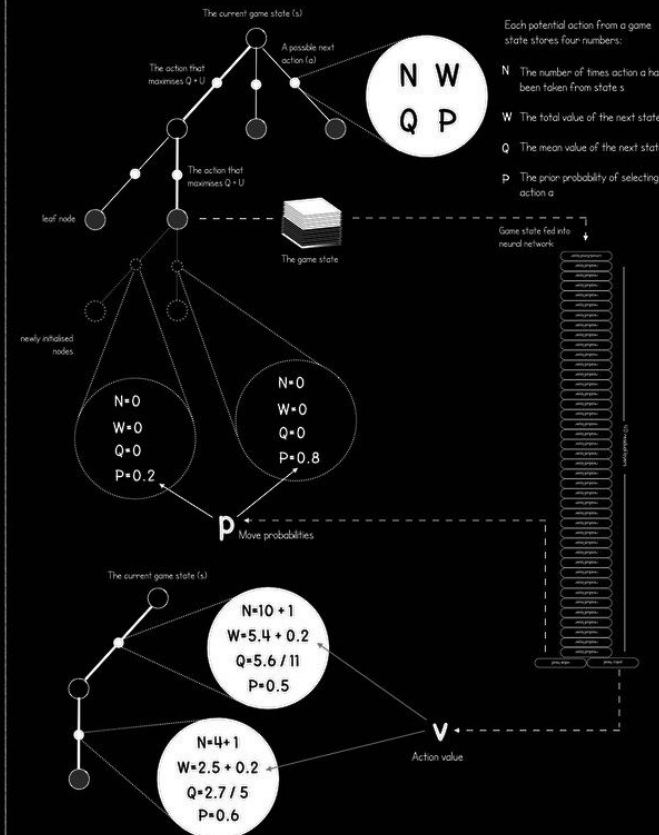
The network learns 'tabula rasa' (from a blank slate)

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



MONTE CARLO TREE SEARCH (MCTS)

How AlphaGo Zero chooses its next move



... then select a move

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

Deterministically (for competitive play)

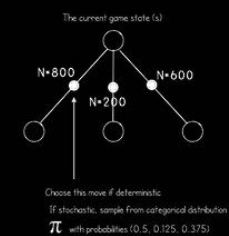
Choose the action from the current state with greatest N

Stochastically (for exploratory play)

Choose the action from the current state from the distribution

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

where τ is a temperature parameter, controlling exploration



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:

$$N \rightarrow N + 1$$

$$W \rightarrow W + v$$

$$Q = W / N$$

Other points

- The sub-tree from the chosen move is retained for calculating subsequent moves
- The rest of the tree is discarded



APPLICATIONS



Artificial General Intelligence



How Much Information is the Machine Given during Learning?

► “Pure” Reinforcement Learning (**cherry**)

- The machine predicts a scalar reward given once in a while.

► A few bits for some samples

► Supervised Learning (**icing**)

- The machine predicts a category or a few numbers for each input
- Predicting human-supplied data
- 10→10,000 bits per sample

► Self-Supervised Learning (**cake génoise**)

- The machine predicts any part of its input for any observed part.
- Predicts future frames in videos
- Millions of bits per sample



