

# Rocky

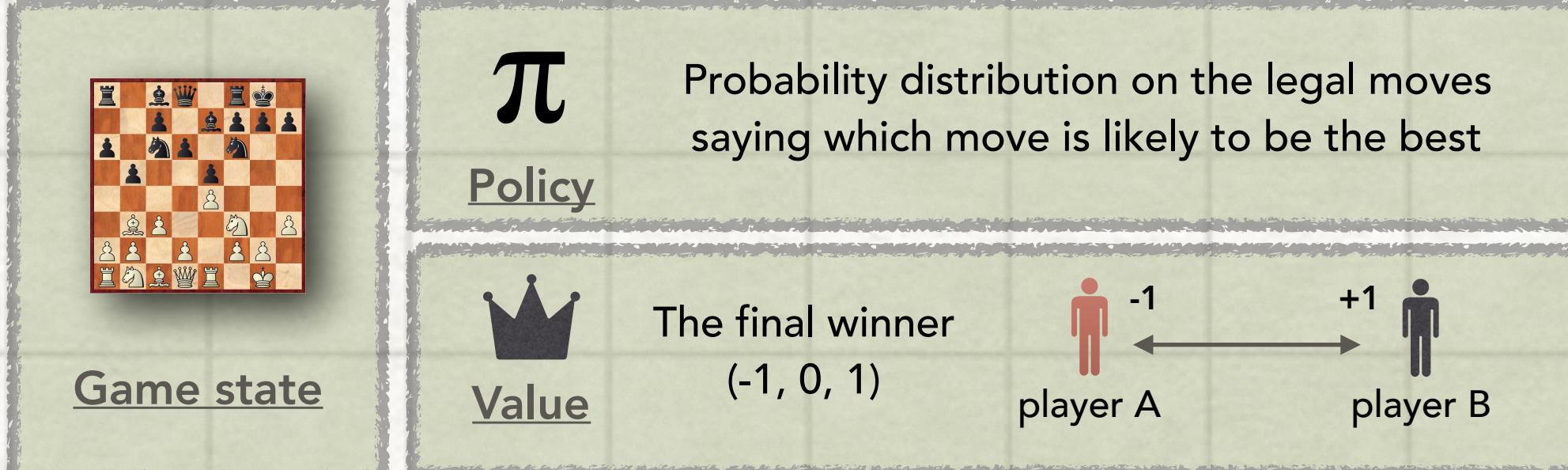
## AlphaZero at home

Gergely Halász & Gergely Papp

### I. Self-play

#### (data generation)

Play  $m$  self-plays that you believe are enough to have a good picture about the game in general. While playing, save the following information in each game:



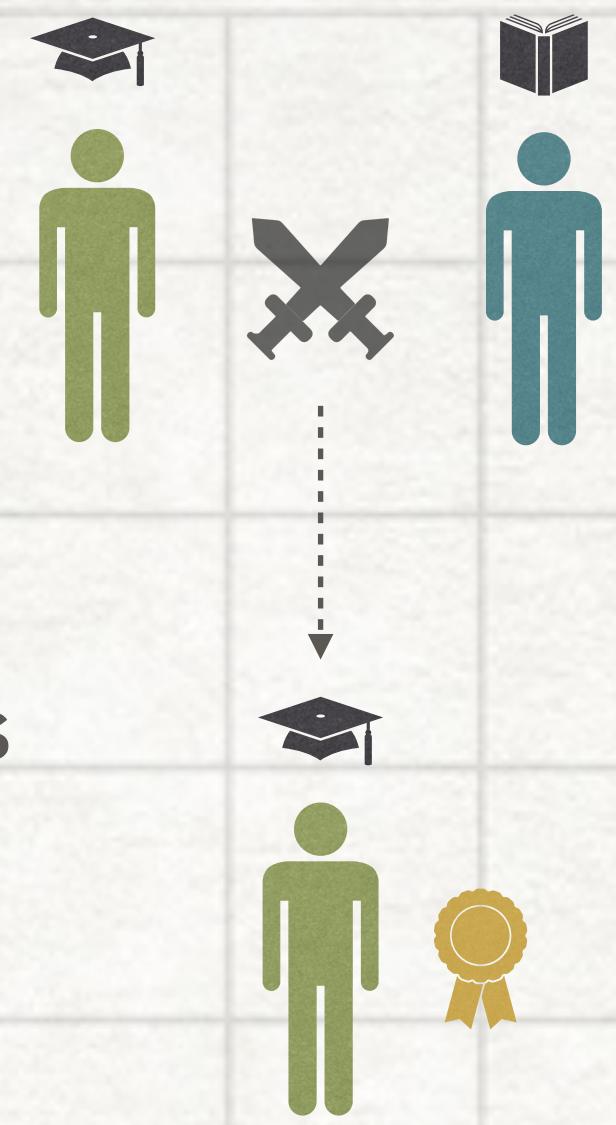
### III. Evolution



Make the fresh neural network play  $n$  games against the reigning one. Choose  $n$  to be a number which is enough to decide if a new prototype is better or not, and if it wins more than 55% of the games, replace the reigning network with the new one.

### Grand picture

0. Let teacher be a random agent
1. Teacher plays some games
2. Student learns from teacher's plays
3. Student challenges their teacher
4. Winner takes teacher's place



### II. Learning

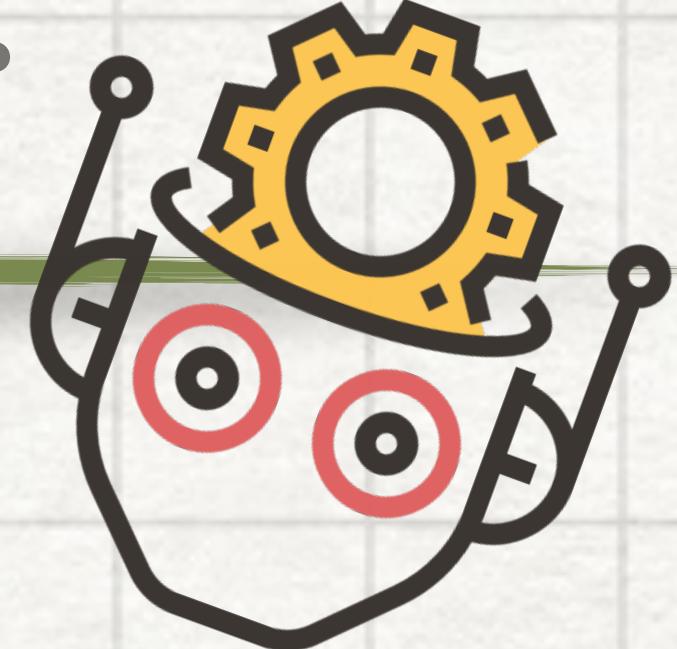
#### (optimising neural network)

Make a copy of your best performing neural network, and learn from the games generated in the last  $k$  self-played games ( $k \geq m$ ).

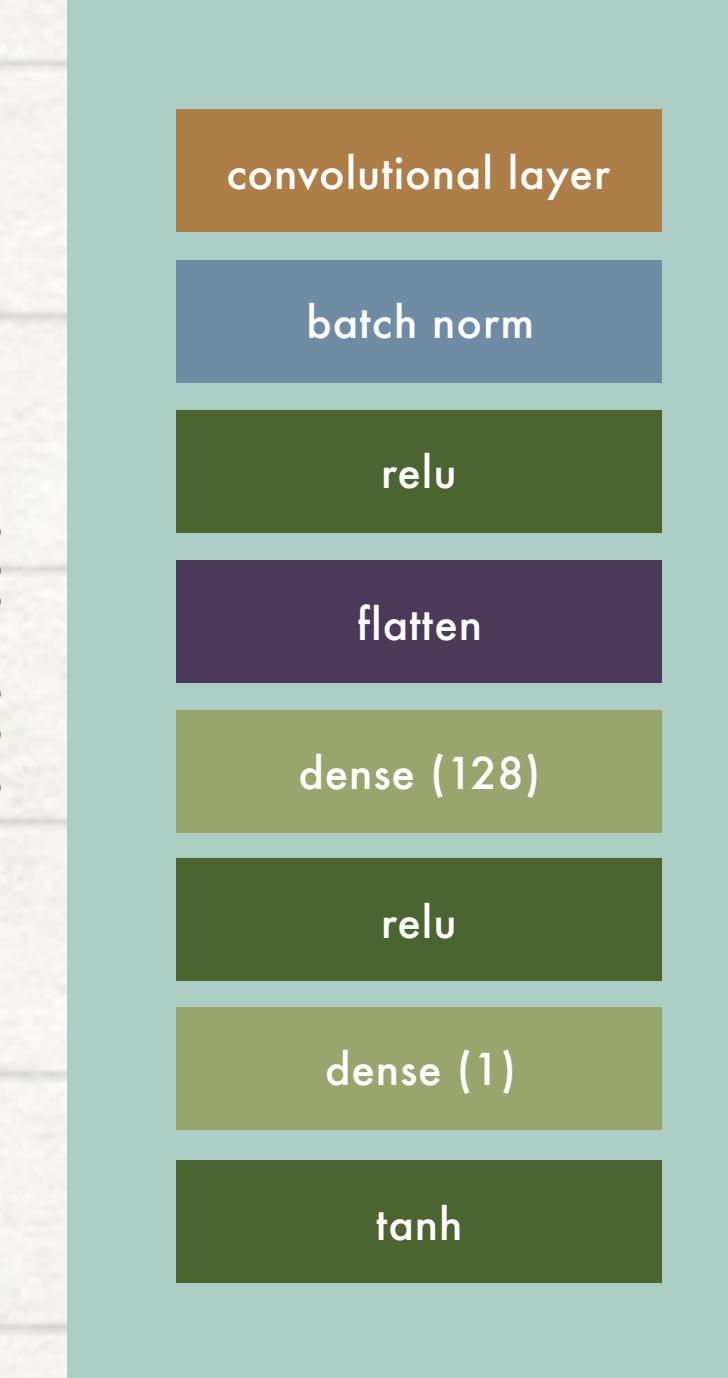
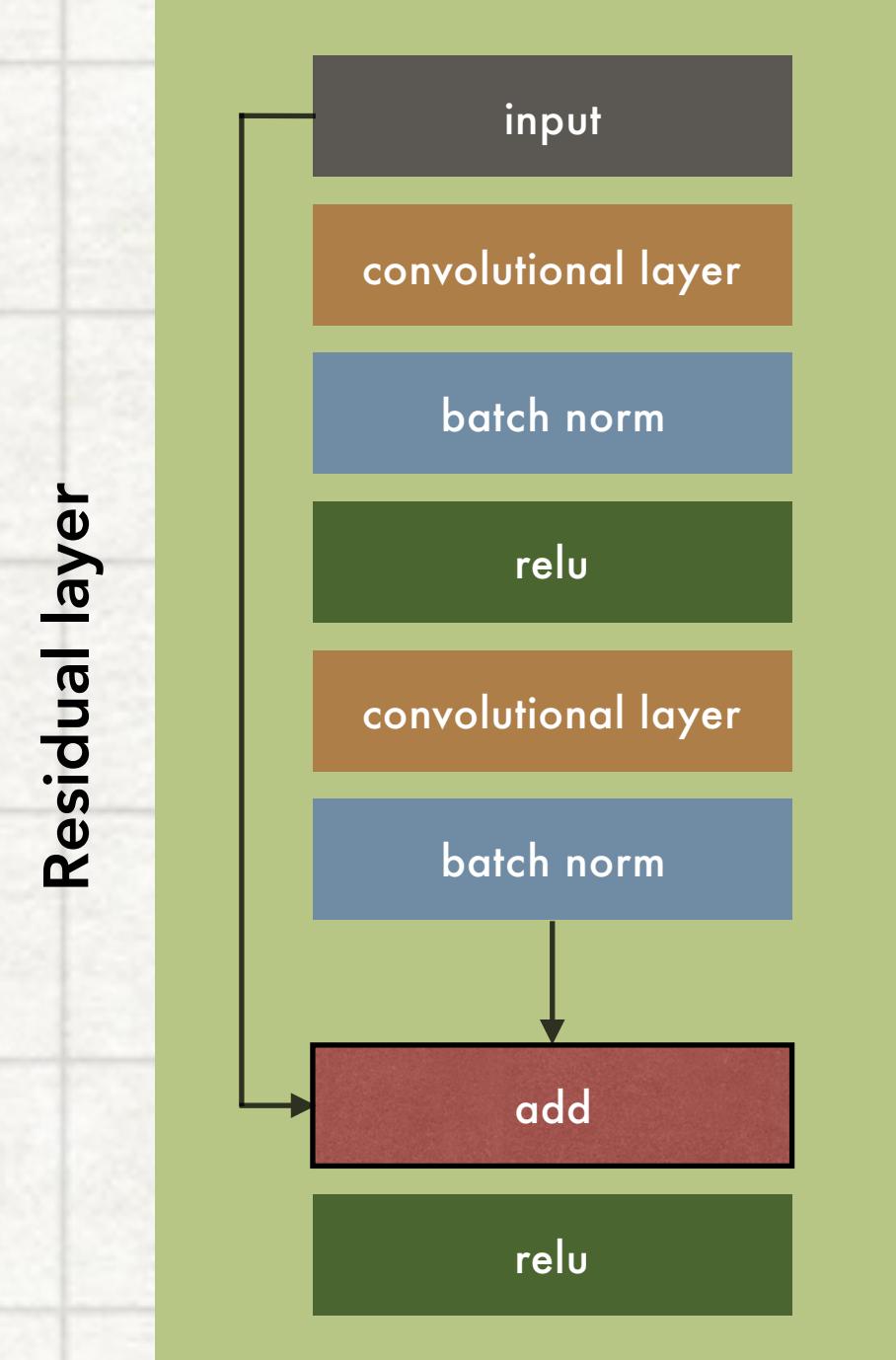
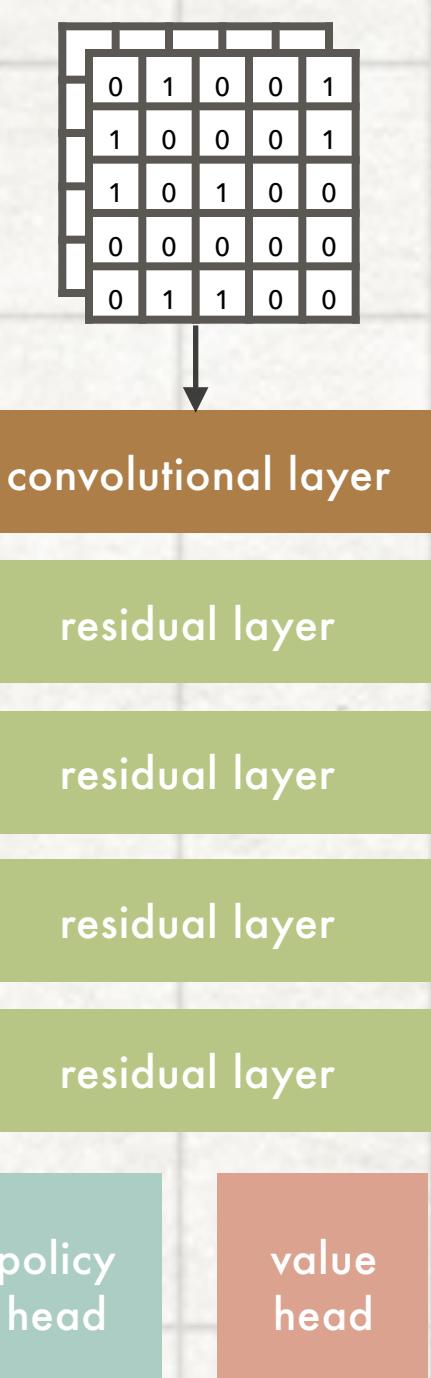
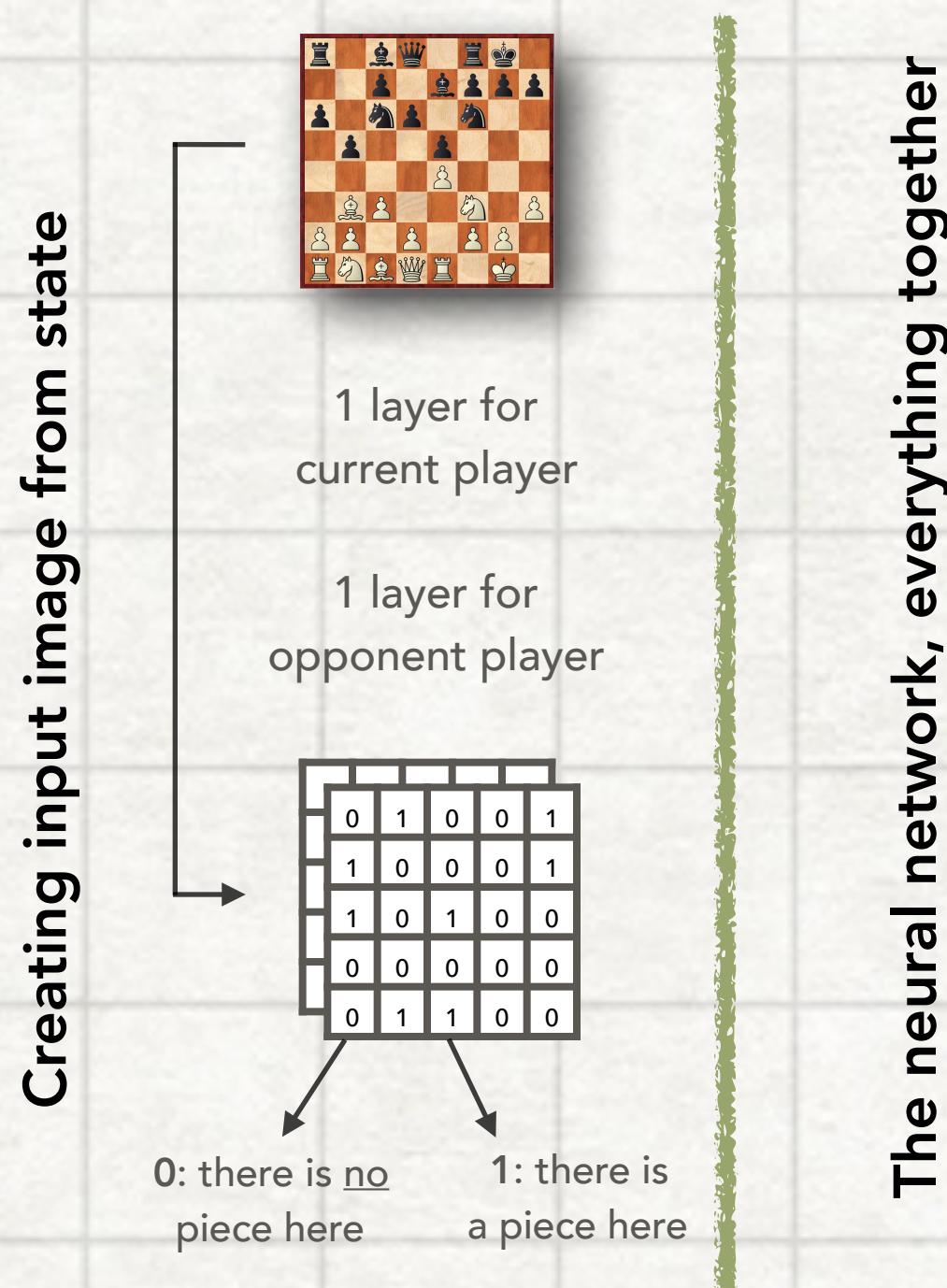
$$\text{loss(state)} = \text{cross\_entropy} + \text{mean_squared_error} + \text{regularisation}$$

$\pi$  Retrieved from MCTS by analysing how many times a move was played out of all  
 $\text{crown}$  A value based on which player won the game. It's either -1, 0 or 1.

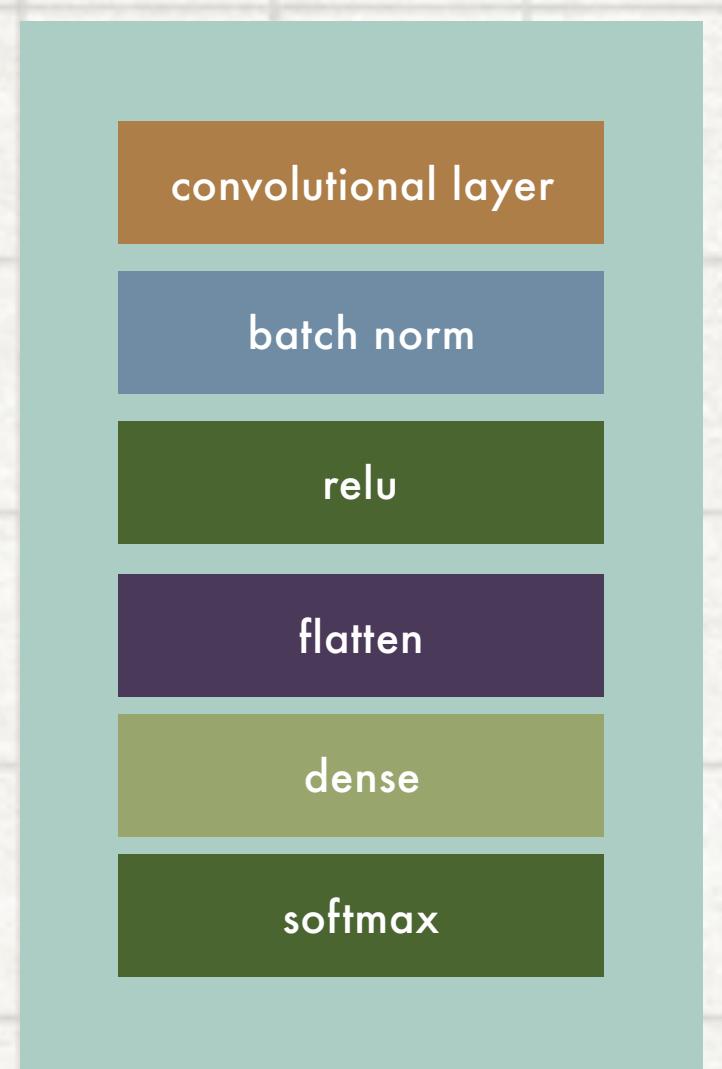
### IV. Go back to step I.



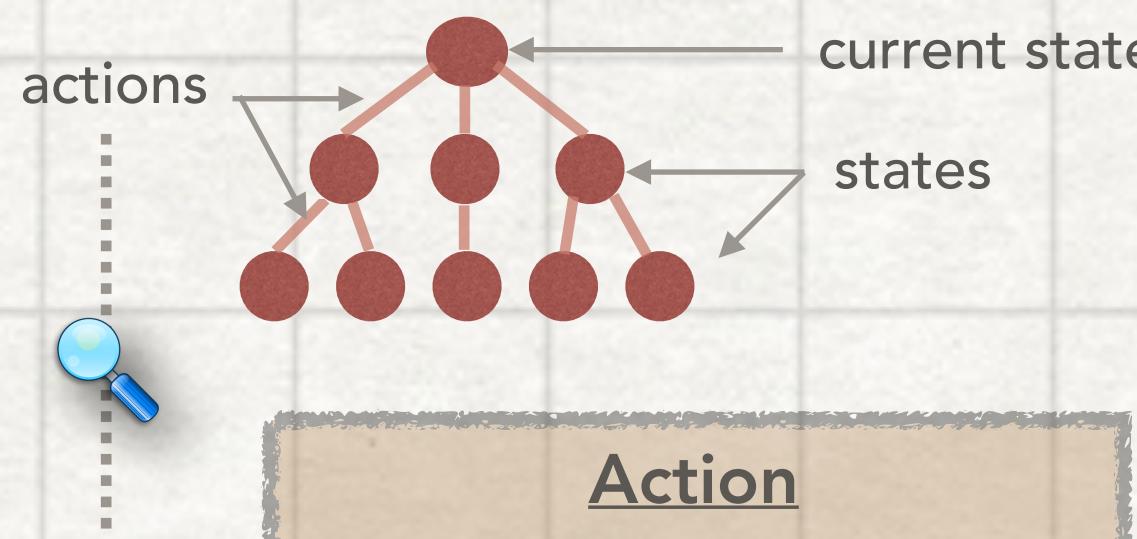
### The neural network



value head



### Monte Carlo Tree Search (MCTS)



#### I. Selection

Choose action, that maximises

$$Q: \text{The mean value of the state}$$

$$U: \text{Value controlling exploration}$$

Repeat until a leaf node reached

#### II. Evaluation

Run the neural network on the leaf state and get predictions on

$$p: \text{Probabilities on next moves}$$

$$v: \text{The outcome of the game}$$

$$p, v$$

#### III. Backpropagation

Update the traversed actions as:

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

$$U = p + C * \sqrt{\frac{N}{\sum N_i}}$$

Run MCTS on current state for a couple of times ( $j$ ). Then, choose your action as

$$\frac{N^{1/\tau}}{\sum N_i^{1/\tau}}$$

stochastic way  
early game

$\tau$  temperature

$$\text{deterministic way}$$
  
late & competitive game  
 $\text{argmax}([N_1, N_2, \dots])$

+ Some tiny dirichlet noise in order to achieve different plays

### Achievements

#### Environment

Game: Connect4

CPU: Intel Core™ 2 Quad, 2.40 GHz

GPU: -

Simulations: 100 (training), 1200 (testing)

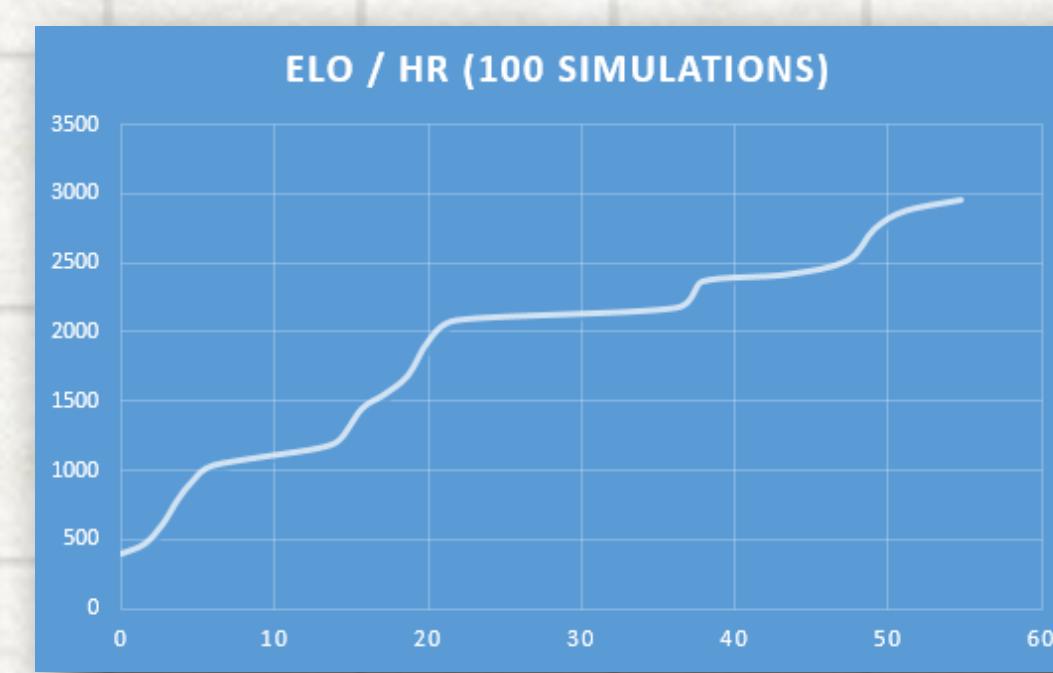
Training time: 24 hours before testing, ~60 hours in total

Iterations: 16 before testing, ~35 in total



#### Results

The engine was tested (unofficially) against a small group of data scientists and physicists. Interestingly, after 24 hours of training Rocky managed to beat the human team 5 to 1.



### Further work

We also tested what if we insert temporal difference methods into the algorithm, but it has not achieved anything outstanding. The idea is to reduce the difference between the predicted value and real outcome if a risky move was played in the next couple of states.