

# Generalised reinforcement learning with self play

INF552 Project

Kaladhar  
Abhinav  
Pallavi



# INTRODUCTION

- Reinforcement learning- software agents ought to take actions in an environment to maximize some notion of reward.
- Alpha zero algorithm
- Neural networks
- Self play

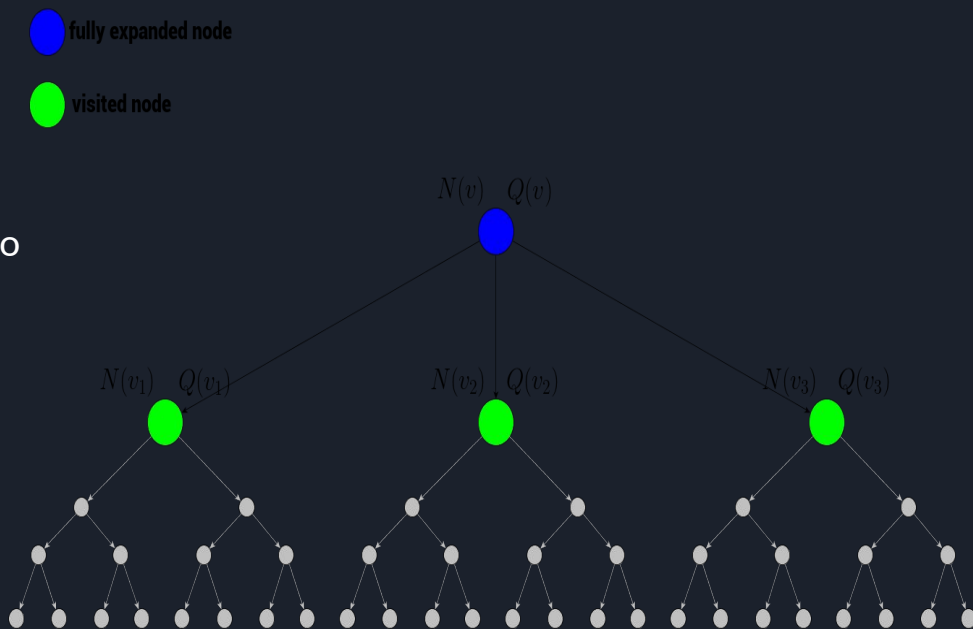


# ALPHAGO COMPONENTS

- MCTS
- Policy network
- Value Network
- Reinforcement learning used for training the network via self-plays

# MCTS

- Simulates the game many times
- Computes the most promising move
- Moves are chosen with respect to policy function
- Result is back propagated to the current game tree node.
- $Q(v)$  total simulation reward
- $N(v)$  total number of visits

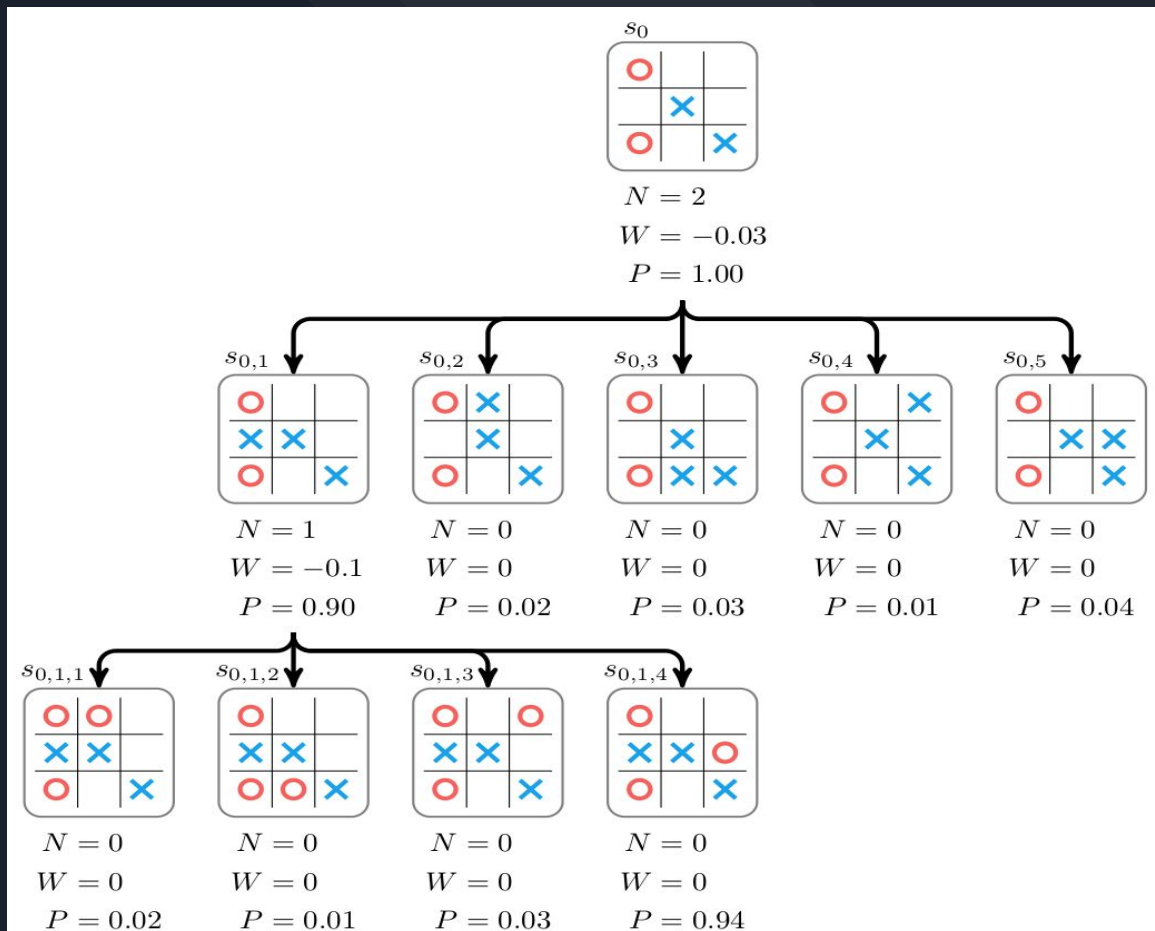




# Upper confidence bound

- Function that chooses next node among visited nodes to traverse through-
- Exploration component
- Exploitation component

# MCTS



# Policy and Value Network

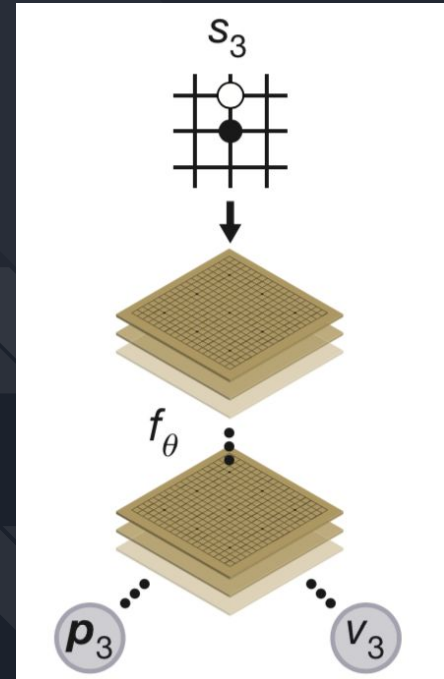
Both Policy and Value are determined by the same network.

The heads (output layer) are different for the Network. The value function is a continuous-valued function between  $[-1,1]$ . Also, Policy network provides the Probability for the next best move.

This makes the network robust as the loss function is applied to the 'common network' and the respective heads.

The Loss function for the network is given as follows :

SGD is used to reduce the loss for each batch.



$$l = \sum_t (v_\theta(s_t) - z_t)^2 - \vec{\pi}_t \cdot \log(\vec{p}_\theta(s_t))$$

$$U_i = \frac{W_i}{N_i} + cP_i \sqrt{\frac{\ln N_p}{1 + N_i}}$$

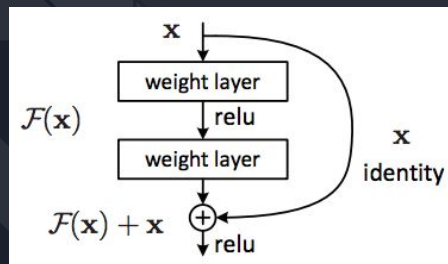
# Neural Network Architecture

Residual Networks are used to achieve more depth in terms of layers.

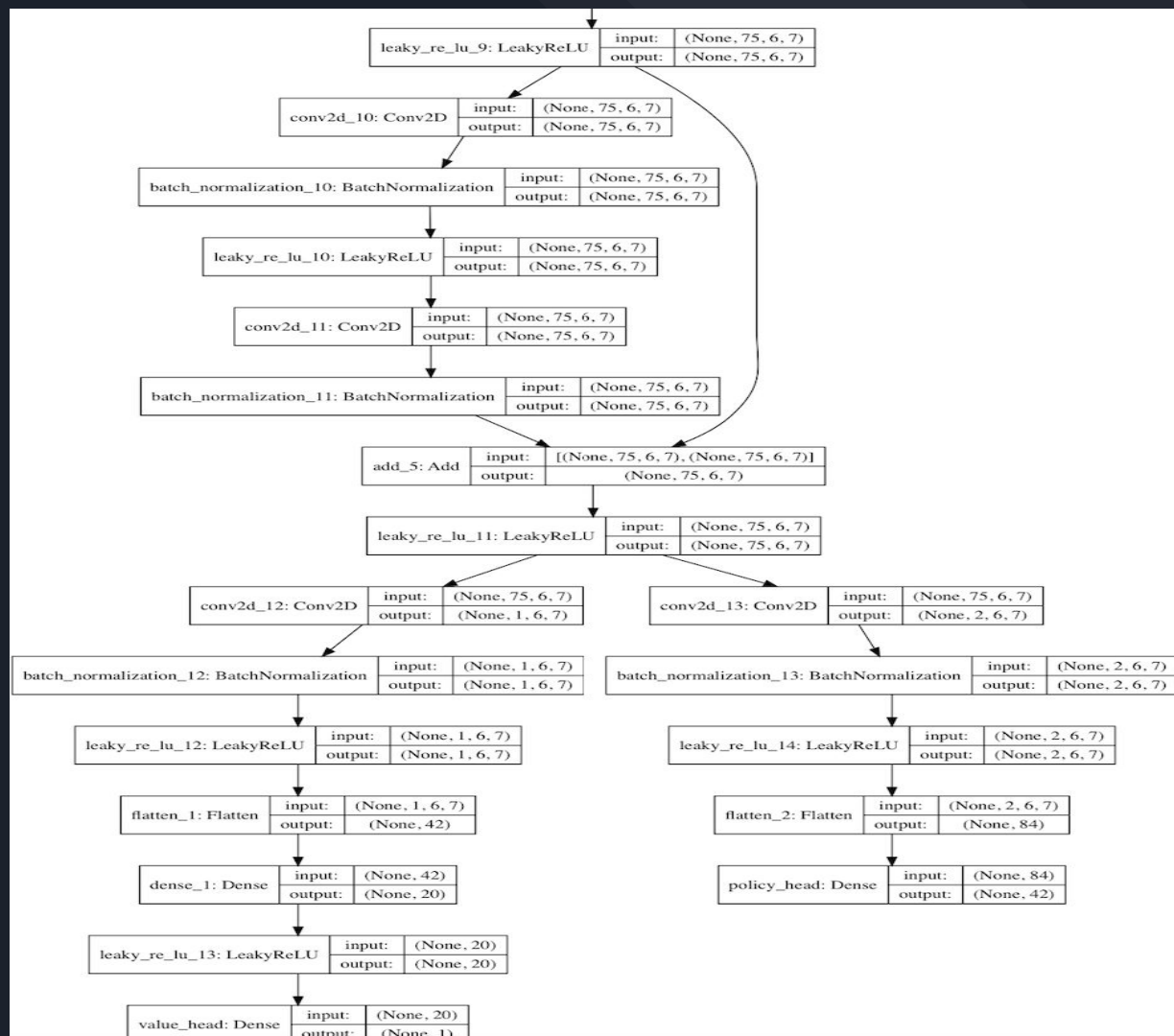
The utilization of RN's allowed for 20 RN's.

Vanishing Gradients problem is reduced and training error does not increase as the input takes a shortcut with the help of an identity mapping.

The actual residual block consists of two CNN's instead of a Dense Layer.







# Policy Iteration through Self-play

In each simulation of our algorithm, the `best_model` play a number of episodes (10 in the current experiments) of self-play using MCTS. The resulting training examples is of the form  $(s(t), \sim\pi(t), z(t))$ .

Then, the current best is updated neural network using the new training examples, to get a new neural network. We then play the old and new networks against each other for a number of games (40 in our experiments). If the new network wins more than a set threshold number of times, the network is updated and continue with the next iteration, resetting the MCTS tree. Else, we continue with the old network and the old MCTS tree, and conduct another iteration to augment our training examples further.

---

**Algorithm 2** Policy Iteration through Self-Play

---

```
1: procedure POLICYITERATIONSP
2:    $\theta \leftarrow \text{initNN}()$ 
3:    $\text{trainExamples} \leftarrow []$ 
4:   for  $i$  in  $[1, \dots, \text{numIters}]$  do
5:     for  $e$  in  $[1, \dots, \text{numEpisodes}]$  do
6:        $ex \leftarrow \text{executeEpisode}(\text{nn})$ 
7:        $\text{trainExamples.append}(ex)$ 
8:        $\theta_{\text{new}} \leftarrow \text{trainNN}(\text{trainExamples})$ 
9:       if  $\theta_{\text{new}}$  beats  $\theta \geq \text{thresh}$  then
10:         $\theta \leftarrow \theta_{\text{new}}$ 
11:   return  $\theta$ 
```

---

---

**Algorithm 3** Execute Episode

---

```
1: procedure EXECUTEEPISODE( $\theta$ )
2:    $\text{examples} \leftarrow []$ 
3:    $s \leftarrow \text{gameStartState}()$ 
4:   while True do
5:     for  $i$  in  $[1, \dots, \text{numSims}]$  do
6:        $\text{MCTS}(s, \theta)$ 
7:        $\text{examples.add}((s, \pi_s, \_))$ 
8:        $a^* \sim \pi_s$ 
9:        $s \leftarrow \text{gameNextState}(s, a^*)$ 
10:    if  $\text{gameEnded}(s)$  then
11:      //fill  $\_$  in examples with reward
12:       $\text{examples} \leftarrow \text{assignRewards}(\text{examples})$ 
13:    return  $\text{examples}$ 
```

---

Challenges faced..



# DEMO

0 <sub>0</sub>	0 <sub>1</sub>	0 <sub>2</sub>	1 <sub>3</sub>	0 <sub>4</sub>	1 <sub>5</sub>	0 <sub>6</sub>	1 <sub>7</sub>
1 <sub>8</sub>	0 <sub>9</sub>	0 <sub>10</sub>	0 <sub>11</sub>	1 <sub>12</sub>	0 <sub>13</sub>	1 <sub>14</sub>	0 <sub>15</sub>
0 <sub>16</sub>	0 <sub>17</sub>	0 <sub>18</sub>	0 <sub>19</sub>	0 <sub>20</sub>	1 <sub>21</sub>	0 <sub>22</sub>	1 <sub>23</sub>
1 <sub>24</sub>	0 <sub>25</sub>	0 <sub>26</sub>	0 <sub>27</sub>	1 <sub>28</sub>	0 <sub>29</sub>	0 <sub>30</sub>	0 <sub>31</sub>
0 <sub>32</sub>	-1 <sub>33</sub>	0 <sub>34</sub>	1 <sub>35</sub>	0 <sub>36</sub>	0 <sub>37</sub>	0 <sub>38</sub>	0 <sub>39</sub>
0 <sub>40</sub>	0 <sub>41</sub>	0 <sub>42</sub>	0 <sub>43</sub>	0 <sub>44</sub>	0 <sub>45</sub>	-1 <sub>46</sub>	0 <sub>47</sub>
0 <sub>48</sub>	-1 <sub>49</sub>	0 <sub>50</sub>	-1 <sub>51</sub>	0 <sub>52</sub>	-1 <sub>53</sub>	0 <sub>54</sub>	-1 <sub>55</sub>
-1 <sub>56</sub>	0 <sub>57</sub>	-1 <sub>58</sub>	0 <sub>59</sub>	-1 <sub>60</sub>	0 <sub>61</sub>	-1 <sub>62</sub>	0 <sub>63</sub>

