

---

# CSE 326 Analysis and Design of Algorithms

Dr.Walid Gomaa

---

Name	ID
Mohamed Abdelmonem Makram	120220055
Abdelrahman Ahmed Shaheen	120220228
Abdelrhman Mohamed Eldenary	120220253
Anas Ihab Badr	120220360



Computer Science Engineering Department  
Egypt-Japan University of Science and Technology

# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Our Thought Process</b>	<b>2</b>
2.1	Alpha Go . . . . .	2
2.2	Adapting Alpha Go Architecture . . . . .	2
<b>3</b>	<b>Implementation</b>	<b>2</b>
3.1	Policy-Value Network . . . . .	3
3.2	Monte Carlo Tree Search (MCTS) . . . . .	3
3.3	Training Loop . . . . .	4
3.4	Key Hyperparameters . . . . .	4
<b>4</b>	<b>Results</b>	<b>4</b>
<b>5</b>	<b>Conclusion and Future Directions</b>	<b>5</b>

# 1 Introduction

This week we aimed to understand how the popular Alpha Go Deep Learning Model works, and try to implement a simplified version that utilizes a **Policy Function**, **Value Function**, and a **Monte Carlo Tree Search** (MCTS) algorithm. We tried to learn Reinforcement learning in 1 day which was not an ideal choice. We managed to get a working demo. We illustrate the thought process, the implementation, and the results in this report.

## 2 Our Thought Process

### 2.1 Alpha Go

Alpha Go is a computer program developed by DeepMind Technologies, a subsidiary of Alphabet Inc. It plays the board game Go and was the first program to defeat a professional human player, the reigning world champion, and several other top players without handicaps on a full-sized board. The program uses deep neural networks and reinforcement learning to improve its performance.

### 2.2 Adapting Alpha Go Architecture

In our coloring problem, we noticed that the problem actually had a massive state recursive tree ([Link](#)), and brute forcing while providing a correct solution for  $S(4)$ , could not be scaled to  $S(5)$  or  $S(6)$  due to memory constraints (Maybe if we had more memory...). So we thought “What also has a massive state tree?” and it was the GO game which has approximately  $2.1 \times 10^{170}$  legal board positions.

So if it uses RL and MCTS to accurately traverse the states tree, so could we. One challenge was how we would assign the rewards or what to use when estimating the value of a state. Our solution was to generate random states and see their outcomes. If the outcome was a valid coloring, we would assign a reward of 0, and if it was invalid, we would assign a reward of  $-1$  to all the states in the path. Additionally, during the training phase, we would assign a reward of  $+1$  to the states that led to maximum depth we required it to reach, favoring the states that led to this outcome.

## 3 Implementation

We adopted a simplified version of the Alpha Go architecture to solve the Schur Number coloring problem. Our architecture consists of three main components: a

**Policy-Value Network**, a **Monte Carlo Tree Search (MCTS)** module, and a **training loop** that mimics self-play reinforcement learning.

### 3.1 Policy-Value Network

Our neural network takes as input the current coloring state represented as four binary bitmasks, one for each color class (Red, Blue, Green, Cyan), encoded into a  $4 \times 64$  tensor. The network has the following architecture:

- **2 Convolutional Layers:** Each with 64 filters of kernel size 3, followed by Batch Normalization and ReLU activation. These layers capture local spatial features within each bitmask.
- **Flattening:** The output is flattened into a vector of size  $64 \times 64$ .
- **Fully Connected Layer:** A dense layer with 256 units and ReLU activation to create a shared latent representation.
- **Output Heads:**
  - **Policy Head:** Outputs a probability distribution over 4 actions (corresponding to the 4 colors), using a Softmax layer.
  - **Value Head:** Predicts a scalar value estimating the expected outcome from the current state, using a Tanh activation to bound the output between -1 and 1.

### 3.2 Monte Carlo Tree Search (MCTS)

We implemented a basic version of MCTS to guide action selection during self-play:

- **Selection:** Traverse the tree using the PUCT formula to balance exploration and exploitation.
- **Expansion:** For a selected node, we expand valid children using the policy output from the neural network.
- **Simulation:** A random rollout is performed to the end of the episode using valid actions.
- **Backpropagation:** The simulation result is backpropagated to update visit counts and node values.

The search outputs a policy distribution over actions, based on visit counts at the root.

### 3.3 Training Loop

We simulate self-play episodes using MCTS-guided decisions. For each episode:

- We begin from a hardcoded initial state (covering some known valid subsets).
- At each time step, we perform MCTS with the current state and collect the state, policy, and outcome.
- When the episode ends (either by reaching a maximum element or by failure), we assign a reward of  $+1$  for success and  $-1$  for failure.

These collected triplets  $(s, \pi, v)$  are used to train the network using a combination of cross-entropy loss for the policy head and mean squared error for the value head.

### 3.4 Key Hyperparameters

- MCTS Simulations per Move: 100–500
- Max Element ( $z$ ): up to 20
- Training Epochs: 10
- Optimizer: Adam, with a learning rate of 0.001
- Batch Size: 32

## 4 Results

We tried to train the model for a maximum depth of 20 that’s to do random rollouts and see if the model can learn to color the graph till maximum depth of 20. The results were quite disappointing as the model could only reach a maximum depth of 27. We don’t understand enough theory to actually debug if this is a bad policy, or a bad value function. Or if it was because the data was not enough to generalize and actually learn the “coloring game”, as we would call it, since to train such a network to generalize we probably need more ground truth data. The results are shown in figure 1.

```

Episode 90/100
Episode 91/100
Episode 92/100
Episode 93/100
Episode 94/100
Episode 95/100
Episode 96/100
Episode 97/100
Episode 98/100
Episode 99/100
Episode 100/100
Epoch 1/10, Policy Loss: 1.3148, Value Loss: 0.2276
Epoch 2/10, Policy Loss: 1.2695, Value Loss: 0.0829
Epoch 3/10, Policy Loss: 1.2476, Value Loss: 0.0460
Epoch 4/10, Policy Loss: 1.2389, Value Loss: 0.0405
Epoch 5/10, Policy Loss: 1.2149, Value Loss: 0.0345
Epoch 6/10, Policy Loss: 1.2024, Value Loss: 0.0494
Epoch 7/10, Policy Loss: 1.1599, Value Loss: 0.0454
Epoch 8/10, Policy Loss: 1.1915, Value Loss: 0.0459
Epoch 9/10, Policy Loss: 1.1891, Value Loss: 0.0435
Epoch 10/10, Policy Loss: 1.1882, Value Loss: 0.0417
Model weights saved to schur_model_weights.pth
(venv) [anaconda3@netpan1c project]$ python3 inferprok.py
Searching for coloring...
z-6: Assigned to Green
z-7: Assigned to Green
z-8: Assigned to Green
z-9: Assigned to Green
z-10: Assigned to Cyan
z-11: Assigned to Cyan
z-12: Assigned to Cyan
z-13: Assigned to Cyan
z-14: Assigned to Cyan
z-15: Assigned to Cyan
z-16: Assigned to Cyan
z-17: Assigned to Cyan
z-18: Assigned to Cyan
z-19: Assigned to Cyan
z-20: Assigned to Blue
z-21: Assigned to Blue
z-22: Assigned to Green
z-23: Assigned to Green
z-24: Assigned to Green
z-25: Assigned to Blue
z-26: Assigned to Blue
z-27: Assigned to Red
called to find valid coloring at z-28

```

Figure 1: The model trying to assign numbers to valid colors.

## 5 Conclusion and Future Directions

In this report, we presented a simplified version of the Alpha Go architecture to tackle the Schur Number coloring problem. While our implementation successfully demonstrated the use of MCTS and a neural network for state evaluation, the results were not as promising as we had hoped. The model struggled to generalize and reach deeper states effectively.

**Future Directions:** We would also try to implement the code on our own hopefully because we can't trust the code we generated to be 100% correct. Also, we aim to generate our own ground truth data using the CUDA code we wrote previously. Now, the model would have a better chance to learn the coloring game. Also, this week's implementation was using a linear policy and value function, we need to implement a more complex one using a CNN, or an LSTM or a transformer architecture to capture long-term dependencies in the state space. It will not be easy with minimal experience in reinforcement learning, but we are willing to learn and improve.