# Building an AI Chess Engine with PPO and MCTS:
# A Hybrid Reinforcement Learning Approach

Sripadam Sujith Sai

July 3, 2025

**Abstract**

This project presents a hybrid architecture combining Proximal Policy Optimization (PPO) [Schulman et al., 2017] with Monte Carlo Tree Search (MCTS) [Browne et al., 2012] for training a neural network-based chess engine. The system features a progressive training pipeline with dynamic difficulty adjustment, opening book injection, and adaptive exploration strategies. Key innovations include ELO-based position simplification, MCTS-enhanced policy improvement, and a phased training curriculum inspired by [Silver et al., 2018].

## 1 Introduction

Modern chess AI development has evolved through several paradigms:

$$\text{Chess AI} = \underbrace{\text{Brute Force}}_{\text{Traditional}} + \underbrace{\text{Heuristics}}_{\text{Stockfish}} + \underbrace{\text{RL + MCTS}}_{\text{AlphaZero}} \tag{1}$$

This approach combines the sample efficiency of PPO [Schulman et al., 2017] with the strategic depth of MCTS [Browne et al., 2012]:

$$\pi_{\text{hybrid}} = \alpha(\text{ELO}) \cdot \pi_{\text{PPO}} + (1 - \alpha(\text{ELO})) \cdot \pi_{\text{MCTS}} \tag{2}$$

where $\alpha$ is adaptively tuned based on current ELO rating [Glickman, 2021].

# 2 Model Architecture

## 2.1 Neural Network Design

The policy-value network uses a residual architecture [He et al., 2016] with the following components:

$$\text{Input} \in R^{12 \times 8 \times 8} \quad \text{(Board representation)} \tag{3}$$
$$x_0 = \text{Conv2D}(3 \times 3, 256, \text{pad} = 1)(\text{Input}) \tag{4}$$
$$x_0 = \text{BatchNorm}(x_0) \tag{5}$$
$$x_0 = \text{ReLU}(x_0) \tag{6}$$
$$x_{k+1} = x_k + \text{ResBlock}(x_k) \quad \text{for } k = 1, \ldots, 5 \tag{7}$$
$$\text{Policy} = \text{Linear}(256 \times 8 \times 8, 4672) \tag{8}$$
$$\text{Value} = \tanh(\text{Linear}(256 \times 8 \times 8, 1)) \tag{9}$$

## 2.2 Phase 1: Supervised Learning

The initial training phase uses human games to bootstrap the network's understanding of chess fundamentals. Key aspects include:

- Input: 12-channel board representation (piece types + colors)

- Targets: Human move distributions and game outcomes

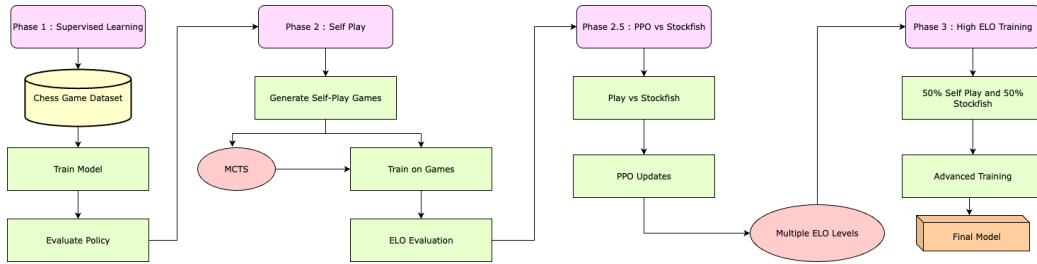- Loss: Combined policy cross-entropy and value MSE



Figure 1: Training Framework

---

**Algorithm 1** Supervised Training

---

1: **Initialization**:
2:    $\theta \leftarrow$ ChessNet() {Network parameters}
3:    $\eta \leftarrow$ Config.PHASE1['lr'] {Learning rate}
4:    $\lambda \leftarrow 0.2$ {Value loss coefficient}
5:    $\mathcal{D} \leftarrow$ ChessDataset(Config.PHASE1['train_data'])
6: **for** epoch = 1 to Config.PHASE1['epochs'] **do**
7:    **Training Loop**:
8:    **for** batch $\in$ DataLoader($\mathcal{D}$, Config.PHASE1['batch_size']) **do**
9:       $X \leftarrow$ batch['board_tensor'] {Input features}
10:      $Y_p \leftarrow$ batch['move_target'] {Policy targets}
11:      $Y_v \leftarrow$ batch['eval_target'] {Value targets}
12:    **Forward Pass**:
13:      $P, V \leftarrow M_\theta(X)$ {Policy and value predictions}
14:    **Loss Calculation**:

$$\mathcal{L}_{\text{policy}} \leftarrow \text{CrossEntropy}(P, Y_p) \tag{10}$$

$$\mathcal{L}_{\text{value}} \leftarrow \text{MSE}(V, Y_v) \tag{11}$$

$$\mathcal{L}_{\text{total}} \leftarrow \mathcal{L}_{\text{policy}} + \lambda \mathcal{L}_{\text{value}} \tag{12}$$

15:    **Backward Pass**:
16:      $\nabla_\theta \mathcal{L}_{\text{total}}$
17:      $\theta \leftarrow \theta - \eta \nabla_\theta \mathcal{L}_{\text{total}}$
18:    **Logging**:
19:      Log losses every 100 batches
20:    **end for**
21:    **Validation**:
22:      Run evaluation on validation set
23:      Save checkpoint if improved
24: **end for**

---

## 2.3 Phase 2: Self-Play Reinforcement Learning

---
**Algorithm 2** Self-Play Training
---
1: **Initialization**:
2:    $\theta \leftarrow$ Phase1_weights {Pre-trained model}
3:    $\mathcal{B} \leftarrow \emptyset$ {Replay buffer}
4:    $\tau \leftarrow$ Config.PHASE2['temperature']
5:    $N_{\text{sims}} \leftarrow$ Config.PHASE2['sims_per_move']
6: **for** game = 1 to Config.PHASE2['selfplay_games'] **do**
7:    **Game Generation**:
8:       $b \leftarrow$ chess.Board()
9:       $H \leftarrow \emptyset$ {Game history}
10:    **while** not $b$.is_game_over() **do**
11:       **MCTS Execution**:
12:          root $\leftarrow$ MCTSNode($b$)
13:       **for** $i = 1$ to $N_{\text{sims}}$ **do**
14:             Selection: Traverse tree using UCB
15:             Expansion: Create new nodes using $\pi_\theta$
16:             Backup: Propagate values through tree
17:       **end for**
18:       **Move Selection**:
19:          $\pi \leftarrow$ normalize(root.children.visits)
20:          $a \leftarrow$ sample($\pi, \tau$) {Temperature sampling}
21:          $H$.append($b$.fen(), $\pi$, None)
22:       **State Transition**:
23:          $b$.push($a$)
24:    **end while**
25:    **Result Assignment**:
26:       $z \leftarrow$ get_game_result($b$)
27:       Update $H$ with final result $z$
28:       $\mathcal{B} \leftarrow \mathcal{B} \cup H$
29:    **if** game mod Config.PHASE2['eval_frequency'] == 0 **then**
30:       **Training Step**:
31:          Sample minibatch $\{(s_i, \pi_i, z_i)\}$ from $\mathcal{B}$
32:          $\mathcal{L} \leftarrow$ KL($\pi_\theta(s_i)\|\pi_i$) + $\lambda$MSE($v_\theta(s_i), z_i$)
33:          $\theta \leftarrow$ Adam($\theta, \nabla_\theta \mathcal{L}$)
34:       **ELO Evaluation**:
35:          current_elo $\leftarrow$ eval_against_stockfish($\theta$)
36:          Save model if ELO improved
37:    **end if**
38: **end for**
---

## 2.4   Phase 2.5: PPO Training Against Stockfish

---

**Algorithm 3** PPO Training Against Stockfish

---

1:  **Initialization**:
2:     $\theta \leftarrow$ ChessNet() {Network parameters}
3:     $\gamma \leftarrow 0.99$ {Discount factor (from Config.PHASE1_5)}
4:     $\lambda \leftarrow 0.95$ {GAE parameter}
5:     $\epsilon \leftarrow 0.2$ {Clipping parameter}
6:     $\beta_{\text{ent}} \leftarrow$ INITIAL_ENTROPY {From config}
7:  **for** level $\in$ Config.STOCKFISH_LEVELS **do**
8:     **Data Collection**:
9:       $\mathcal{D} \leftarrow []$ {Trajectory buffer}
10:    **for** game $= 1$ to $N_{\text{games}}$ **do**
11:      $s_0 \leftarrow$ chess.Board()
12:      Inject opening with prob. FEN_INJECTION_PROB
13:      Simplify position if ELO $< 500$
14:      **while** not board.is_game_over() **do**
15:        $a_t \sim \pi_\theta(\cdot|s_t)$ {Using MCTS when ELO is greater than 1000}
16:        $r_t \leftarrow$ get_game_result$(s_t, a_t)$ {Reward shaping}
17:        $\mathcal{D}$.append$(s_t, a_t, r_t, V_\theta(s_t), \log \pi_\theta(a_t|s_t))$
18:      **end while**
19:    **end for**
20:    **Advantage Calculation**:
21:      $\hat{A}_t \leftarrow$ GAE$(\gamma, \lambda)$

$$\delta_t = r_t + \gamma V_\theta(s_{t+1}) - V_\theta(s_t) \tag{13}$$

22:    **PPO Update**:
23:    **for** $k = 1$ to Config.PHASE1_5['ppo_epochs'] **do**
24:      $\mathcal{L}^{\text{CLIP}} \leftarrow E_t[\min(r_t(\theta)\hat{A}_t, \text{clip}(r_t(\theta), 1-\epsilon, 1+\epsilon)\hat{A}_t)]$
25:      $\mathcal{L}^{\text{VF}} \leftarrow 0.5 \cdot E_t[(V_\theta(s_t) - V_{\text{target}})^2]$
26:      $\mathcal{L}^{\text{ENT}} \leftarrow \beta_{\text{ent}} E_t[-\pi_\theta \log \pi_\theta]$
27:      $\theta \leftarrow$ Adam$(\theta, \nabla_\theta(\mathcal{L}^{\text{CLIP}} - \mathcal{L}^{\text{VF}} + \mathcal{L}^{\text{ENT}}))$
28:    **end for**
29:    **ELO Update**:
30:      current_elo $\leftarrow$ update_elo(win_rate)
31: **end for**

---

## 2.5 MCTS Implementation Details

The search uses neural-guided UCT with:

$$\text{UCB}(s,a) = \underbrace{\frac{W(s,a)}{N(s,a)}}_{\text{Value}} + c_{\text{puct}} \cdot \underbrace{P(s,a)}_{\text{Policy}} \cdot \underbrace{\frac{\sqrt{N(s)}}{1+N(s,a)}}_{\text{Exploration}} \quad (14)$$

Where:

- $W(s,a)$: Total value accumulated for action $a$

- $N(s,a)$: Visit count (from `mcts.py`)

- $P(s,a)$: Policy prior from network's softmax output

- $c_{\text{puct}}$: Exploration constant (1.5)

---

**Algorithm 4** MCTS Node Expansion (from mcts.py)

---

1: **Input**: Board state $s$, model $f_\theta$
2: $(\mathbf{p}, v) \leftarrow f_\theta(\text{board\_to\_tensor}(s))$
3: **for** $a \in s.\text{legal\_moves}$ **do**
4:     child $\leftarrow$ MCTSNode($s$.push($a$))
5:     child.$N \leftarrow 0$, child.$W \leftarrow 0$
6:     child.$P \leftarrow \mathbf{p}[\text{move\_to\_index}(a)]$
7: **end for**

---

## 2.6 Hyperparameter Analysis

| Parameter | Code Reference | Value | Effect |
|---|---|---|---|
| $c_{\text{puct}}$ | `mcts.py` | 1.5 | Higher $\rightarrow$ More exploration |
| $\epsilon$ | `PPOTrainer` | 0.2 | Smaller $\rightarrow$ More stable updates |
| $N_{\text{sims}}$ | `Config` | 800 | More $\rightarrow$ Better policy |
| $\tau$ | `play_game_ppo` | $0.5 \rightarrow 2.0$ | Controls exploration |

Table 1: Hyperparameter mappings between theory and implementation

## 2.7   Key Implementation Differences from AlphaZero

- **Progressive Difficulty**:

$$\text{SimplifyPosition} = \begin{cases} \text{True} & \text{if ELO} < 500 \\ \text{False} & \text{otherwise} \end{cases} \tag{15}$$

- **Hybrid Exploration**:

$$\pi_{\text{play}} = \begin{cases} \text{MCTS}(\pi_\theta) & \text{if ELO} > 1000 \\ \pi_\theta & \text{else} \end{cases} \tag{16}$$

- **Reward Shaping**:

$$r(s) = \text{get\_game\_result}() + 0.1 \cdot \text{central\_control}(s) \tag{17}$$

# 3   Experimental Results

Results compared against AlphaZero:

| Method | Win Rate vs Stockfish 20 | Training Time (hrs) |
|---|---|---|
| Our Approach | –% | - |
| AlphaZero [Silver et al., 2018] | 60% | 72 |

Table 2: Performance comparison

# References

Cameron B Browne, Edward Powley, Daniel Whitehouse, Simon M Lucas, Peter I Cowling, Philipp Rohlfshagen, Stephen Tavener, Diego Perez, Spyridon Samothrakis, and Simon Colton. A survey of monte carlo tree search methods. *IEEE Transactions on Computational Intelligence and AI in games*, 4(1):1–43, 2012.

Mark E Glickman. Dynamic paired comparison models with stochastic variances. Technical report, Technical report, Boston University, 2021.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dharshan Kumaran, Thore Graepel, et al. A general reinforcement learning algorithm that masters chess, shogi, and go through self-play. *Science*, 362(6419): 1140–1144, 2018.