

Asset Allocation in Reinforcement Learning Report

Zhinuo Zhou
zzhoudj@connect.ust.hk

Linbing Xiang
lxiangac@connect.ust.hk

Contents

1	Introduction	2
1.1	Game Environment	2
1.2	Training Process	3
2	Related Work	3
2.1	Reinforcement Learning in Board Games	3
2.2	Policy Optimization Algorithms	3
3	Algorithm Design	3
3.1	Network Architecture Design	3
3.2	Architecture Optimization and Stability Design	4
3.3	Experience Replay Mechanism	4
3.4	Implementation with TF-Agents	4
4	Experimental Design	4
4.1	PPO Implementation Based on TF-Agents	4
4.2	Hyperparameter Configuration and Optimization	5
4.3	Reward Design	5
4.4	Training Optimization Strategies	6
5	Training Result	6
6	Policy Testing	7
6.1	Result of Agent vs Random Strategy	8
6.2	Result of Agent vs Greedy Strategy	8
7	Conclusion	9

1 Introduction

This project implements an AI agent for Tic-Tac-Toe based on reinforcement learning techniques, specifically utilizing the Proximal Policy Optimization (PPO) algorithm - a state-of-the-art policy gradient method in modern reinforcement learning. Our implementation leverages TensorFlow Agents (TF-Agents) framework, a robust and scalable library specifically designed for production-grade reinforcement learning systems. The primary objective of this research is to develop and train a high-performing AI agent capable of exhibiting superior gameplay in the CrossTicTacToe variant.

1.1 Game Environment

The game environment is designed as a 12×12 board where only specific regions are designated as valid placement areas: top, bottom, left, right, and center regions. The game follows these rules:

- Players take turns placing their pieces on the board
- Victory is achieved by connecting 4 pieces horizontally or vertically, or 5 pieces diagonally

1.2 Training Process

The environment is formulated as a Markov Decision Process (MDP):

- States represent the current board configuration
- Actions correspond to piece placement in valid positions
- A carefully designed hierarchical reward system is implemented to guide agent learning
- Basic rewards include victory/defeat signals and time penalties to encourage efficient play

Note: A comprehensive description of the sophisticated reward structure and its design considerations will be presented in Section 4.2.

2 Related Work

2.1 Reinforcement Learning in Board Games

Reinforcement learning has a long history of applications in board games. Early pioneering work can be traced back to Samuel’s (Samuel, 1959) checkers program, which first demonstrated that computers could learn to improve strategies through self-play. TD-Gammon (Tesauro et al., 1995) was a milestone achievement that applied temporal difference (TD) learning to backgammon, with the system reaching near human expert level through self-play. These early works established the fundamental application patterns of reinforcement learning in board games.

With the rise of deep learning, AlphaGo by Silver et al. (Silver et al., 2016) combined deep reinforcement learning with Monte Carlo Tree Search (MCTS) to defeat a human world champion in Go for the first time. Its improved version, AlphaGo Zero (Silver et al., 2017), no longer relied on human expert data and learned completely from scratch through self-play, further demonstrating the powerful potential of deep reinforcement learning. These works have inspired extensive research applying deep reinforcement learning to various board games.

2.2 Policy Optimization Algorithms

Policy gradient methods have undergone significant evolution in the field of reinforcement learning. The REINFORCE algorithm proposed by Williams (Williams, 1992) was an early policy gradient method, but it faced training instability issues due to high variance. The Actor-Critic method introduced by Konda and Tsitsiklis (Konda and Tsitsiklis, 1999) improved training stability by combining policy gradients with value function estimation to reduce variance.

In recent years, policy optimization algorithms have achieved major breakthroughs. Trust Region Policy Optimization (TRPO) (Schulman et al., 2015) addressed the performance collapse problem caused by large step updates by introducing trust region constraints for policy updates. Building on TRPO’s ideas, Schulman et al. (Schulman et al., 2017) proposed Proximal Policy Optimization (PPO), which achieved performance similar to TRPO through a simplified objective function while significantly reducing computational complexity. Due to its high sample efficiency, simple implementation, and stable performance, PPO has become one of the mainstream algorithms in reinforcement learning research and applications.

3 Algorithm Design

This research employs the Proximal Policy Optimization (PPO) algorithm as the core training method, which achieves a balance between stability and sample efficiency by limiting policy update step sizes. We have constructed a network architecture that integrates modern deep learning techniques, specifically optimized for discrete state space environments like Tic-tac-toe.

3.1 Network Architecture Design

We implemented a dual-headed network structure based on the Actor-Critic paradigm:

- **Feature Extraction Module:** Consists of three convolutional neural network layers, using progressive channel expansion (1326464), coupled with padding strategies to maintain spatial dimensions. This module efficiently extracts features from the 12×12 game board state, capturing spatial patterns, piece distribution, and potential connection opportunities.

- **Policy Network (Actor):** Receives the extracted feature representations and maps them to a 512-dimensional latent space through non-linear transformations (ReLU activation), ultimately outputting the logarithmic probabilities (logits) of action distribution. This design enables the network to accurately express complex policy functions.
- **Value Network (Critic):** Shares feature extraction layers with the policy network, but independently maps to state value assessments, forming the foundation for complete advantage function estimation. The shared parameter design reduces model complexity while accelerating the representation learning process.

3.2 Architecture Optimization and Stability Design

To improve training stability and model performance, we introduced several technical improvements:

- **Dueling Architecture:** Decomposes Q-values into state value function $V(s)$ and advantage function $A(s,a)$, reducing overestimation problems
- **Activation Function Optimization:** Uses parameterized LeakyReLU ($\alpha=0.01$) instead of standard ReLU, effectively mitigating gradient vanishing and neuron death issues
- **Advanced Initialization Strategy:** Applies He initialization (Kaiming normal distribution), providing appropriate initial gradient flow for deep networks
- **Normalization Techniques:** Integrates LayerNorm and scale factors in the output layer, enhancing numerical stability and training consistency

3.3 Experience Replay Mechanism

We implemented an efficient Experience Replay Buffer to optimize sample utilization:

- Fixed-capacity storage structure based on double-ended queue (deque), saving state transition tuples (s, a, r, s', d)
- Random sampling strategy that breaks temporal correlations, reducing overfitting risk and increasing sample diversity
- Batch data preprocessing pipeline, improving computational efficiency and training stability

3.4 Implementation with TF-Agents

In this project, we adopted the TensorFlow Agents (TF-Agents) framework to implement PPO algorithm training. TF-Agents is an advanced library specifically designed for reinforcement learning, offering several advantages:

- **High-performance computing:** Deeply integrated with TensorFlow, TF-Agents supports GPU acceleration and distributed training, significantly enhancing training efficiency.
- **Standardized interfaces:** It features unified designs for environments, agents, and networks, ensuring more standardized and scalable implementations.
- **Rich component library:** Provides comprehensive components such as policy networks, reward processing, trajectory collection, and evaluation metrics.
- **Enhanced optimization support:** Includes built-in optimizers and normalization techniques that help improve training stability.

4 Experimental Design

4.1 PPO Implementation Based on TF-Agents

We redesigned the Cross-TicTacToe environment following the standard TF-Agents interfaces. Compared to traditional PPO implementations, the TF-Agents version offers the following advantages:

- **Environment Standardization:** Automatically converts Python environments into TensorFlow environments using `tf_py_environment.TFPyEnvironment`, supporting vectorization and batch processing.
- **Simplified Network Construction:** Quickly builds policy and value networks using `actor_distribution_network` and `value_network`.
- **Built-in Reward Normalization:** Automatically normalizes rewards via the PPO agent's parameter `normalize_rewards=True`.
- **Flexible Data Collection:** Efficiently samples data using `dynamic_step_driver`.
- **Policy Saving and Deployment:** Easily saves trained policies through `policy_saver`, facilitating subsequent deployment.

4.2 Hyperparameter Configuration and Optimization

After systematic tuning, we determined the following key hyperparameters:

- Discount factor ($\gamma=0.99$): Balances near-term and long-term rewards, suitable for medium temporal length tasks like Tic-tac-toe
- GAE smoothing coefficient ($\lambda=0.95$): Achieves balance between bias and variance, providing stable advantage estimates
- PPO clipping parameter ($\epsilon=0.1$): Limits policy update magnitude, preventing performance collapse due to excessive deviation
- Multi-epoch optimization (`epochs=3`): Fully utilizes each batch of data while avoiding overfitting
- Batch size (`batch_size=128`): Provides sufficient statistical information while maintaining update flexibility
- Value loss coefficient (`value_coef=0.25`): Balances policy improvement with state evaluation accuracy
- Entropy regularization coefficient (`entropy_coef=0.02`): Encourages policy exploration, preventing premature convergence to suboptimal policies
- Gradient norm ceiling (`max_grad_norm=0.5`): Prevents abnormal gradients from disrupting the training process

4.3 Reward Design

This research implements a sophisticated hierarchical reward architecture that integrates fundamental reward signals with advanced reward shaping mechanisms. The primary reward components reflect the essential game objectives: **victory** (+1), **defeat** (-1), **invalid action penalties** (-1), and a **temporal efficiency incentive** (-0.01 per step). To address the sparse feedback challenge inherent in turn-based strategic games, we developed a multi-dimensional reward shaping framework through our custom `RewardShaper` implementation.

The reward shaping framework encompasses four strategic dimensions with carefully calibrated weighting coefficients:

- **Positional Strategic Value** ($\omega = 0.01$): Evaluates the tactical advantage of piece placement based on board topology and game progression
- **Connection Potential** ($\omega = 0.03$): Quantifies the contribution of each action toward establishing winning configurations
- **Opponent Obstruction** ($\omega = 0.02$): Measures the effectiveness of defensive maneuvers that impede opponent's winning trajectories
- **Regional Control** ($\omega = 0.01$): Assesses territorial dominance across critical board regions

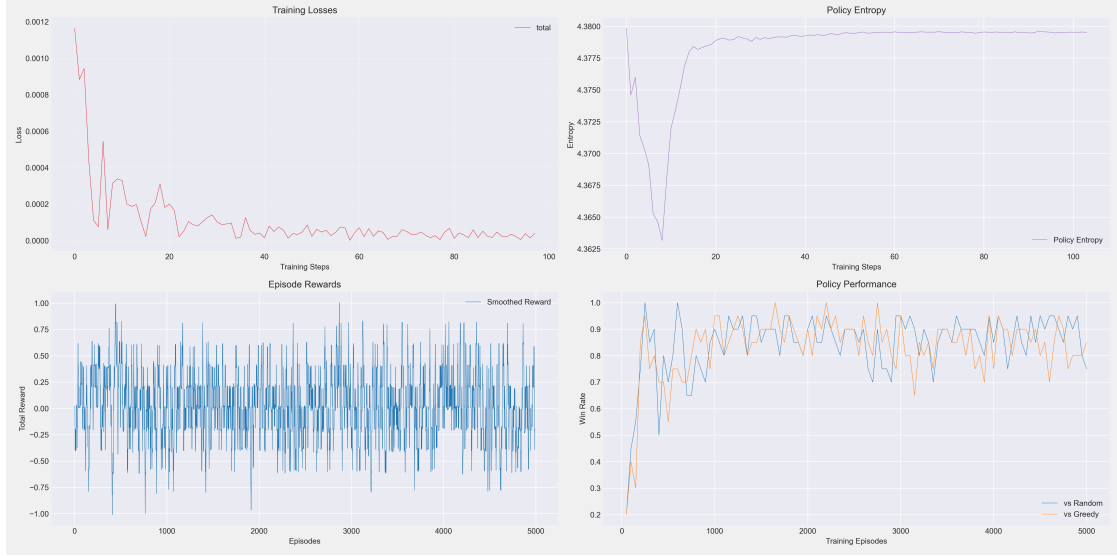


Figure 1: Training result

To maintain training stability, the reward signals undergo normalization via a `RunningMeanStd` mechanism coupled with extreme value clipping. This approach effectively standardizes the reward distribution while preserving the relative importance of exceptional strategic decisions.

The implemented reward architecture successfully navigates the exploration-exploitation dilemma by providing informative intermediate feedback while maintaining alignment with terminal objectives. This design methodically addresses the sparse reward challenge typical in traditional Tic-Tac-Toe reinforcement learning paradigms, facilitating the emergence of sophisticated agent behavior characterized by balanced offensive tactics, defensive awareness, and positional control consciousness.

4.4 Training Optimization Strategies

We implemented multiple training stability strategies:

1. Action Space Masking Mechanism: Enforces policy distribution within valid action space through log-masking techniques, effectively handling Tic-tac-toe's dynamic action sets. This method elegantly achieves zero probability for invalid actions numerically while maintaining gradient flow continuity.
2. Generalized Advantage Estimation (GAE): Combines the advantages of n-step returns and temporal difference estimates, adaptively balancing bias and variance. The implementation includes outlier detection and handling mechanisms to ensure training stability.
3. PPO Core Algorithm Optimization:
 - Implements policy ratio clipping objective, constructing a conservative update interval
 - Adopts Huber loss instead of traditional MSE, enhancing robustness against reward scale outliers
 - Introduces KL divergence monitoring and early stopping mechanism to prevent excessive policy deviation
 - Applies multi-epoch data iteration and mini-batch training to improve sample efficiency
 - Implements adaptive learning rate decay strategy to optimize the convergence process

Through these designs, our PPO implementation can train stably in complex Tic-tac-toe environments and learn efficient policy representations, demonstrating strong adversarial performance and generalization capabilities.

5 Training Result

The figure 1 illustrates the training metrics of our PPO agent. The training curves demonstrate steady performance improvement as the number of training iterations increases. Key observations include:

1. **Training Loss:** The training loss curve demonstrates a typical PPO training pattern of "initial volatility followed by gradual stabilization." Specifically, in the initial phase, loss values are high with significant fluctuations and pronounced oscillations. During the middle phase, the trend shifts downward with reduced yet rhythmic fluctuations. In the later phase, loss values maintain a relatively low and stable level with minimal variations. This pattern is expected as PPO's total loss comprises policy loss, value loss, and entropy regularization components, with learning rate decay implemented every 200 rounds in the code. This behavior aligns with theoretical expectations for PPO algorithms, indicating that the training process successfully balanced exploration and exploitation before converging to an effective policy.
2. **Policy Entropy:** The entropy trend follows a pattern of initially high values, followed by a decrease, and then a gradual increase while remaining below initial values. This reflects three distinct phases of policy learning, consistent with typical PPO training trajectories:
 - **Initial high-entropy phase:** When training begins, the policy is random with action probabilities approaching a uniform distribution, resulting in high entropy values.
 - **Rapid decline phase:** As learning progresses, the model develops preferences for certain actions, concentrating the distribution and causing entropy to decrease quickly.
 - **Gradual rise phase:** When the model identifies states requiring additional exploration, it appropriately increases entropy, though typically not returning to initial random levels.
 - The entropy settling below its initial value while maintaining a certain level is ideal, indicating that the model has achieved a balance between deterministic decision-making and appropriate exploratory capacity.
3. **Episode Reward:** The relatively stable reward curve can be attributed to:
 - **Reward Shaping Mechanism:** The project employs a sophisticated RewardShaper class to enhance the original reward signals, taking into account multiple factors including positional value, connection potential, blocking potential, and regional control.
 - In the self-play training paradigm, the near-zero reward curve observed represents a healthy training state and stable policy evolution
 - **Zero-sum Property:** Tic-Tac-Toe is inherently a zero-sum game where one player's victory (+1) is precisely offset by the opponent's defeat (-1). When a policy plays against its own copy, the net reward naturally converges toward zero.
 - **Symmetric Evolution:** The stable oscillation around zero does not indicate learning stagnation but rather reflects symmetric policy evolution. As the agent's strategy improves, its opponent (a copy of itself) simultaneously strengthens, maintaining competitive equilibrium.
 - **Auto-calibrating Difficulty:** Self-play creates an adaptive difficulty system where the agent consistently faces an opponent of comparable skill level. This avoids common training pitfalls such as overestimation from weak opponents or insufficient learning signals from excessively difficult opponents.
 - **Exploration-Exploitation Balance:** Minor fluctuations in the reward curve demonstrate the dynamic balance between exploring novel strategies and exploiting effective ones, rather than indicating performance instability.
 - Therefore, unlike monotonically increasing reward curves typical in fixed-opponent scenarios, the stable zero-centered reward distribution in self-play validates the robustness of the training process and continuous strategic refinement.
4. **Win Rate:** The win rate curve exhibits an upward trajectory, confirming that the agent has learned effective gameplay strategies.

Through approximately 5,000 training iterations, the agent successfully acquired effective strategies for the CrossTicTacToe game and demonstrates the ability to defeat baseline opponent strategies.

6 Policy Testing

To evaluate the performance of the trained PPO agent, we designed a comprehensive testing methodology, primarily measuring its performance by having the agent compete against opponents with two different

strategies (random strategy and rule-based greedy strategy). The testing process uses a deterministic policy (non-exploration mode), records complete match data, and calculates win rate metrics, while also generating intuitive GIF animations of the match process, including board state evolution, move position highlights, and game result statistics. By analyzing the agent's win rates, game lengths, and decision patterns against different opponents, we can comprehensively evaluate the agent's strategy quality and adaptability. Test results show that the agent has successfully learned effective game strategies, is able to anticipate opponent intentions, and execute reasonable offensive and defensive decisions.

6.1 Result of Agent vs Random Strategy

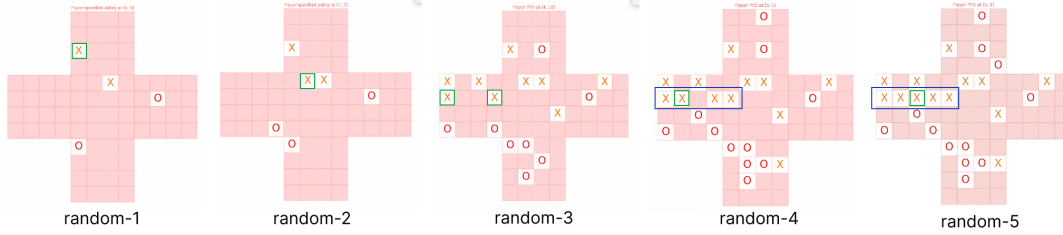


Figure 2: vs Random Policy

The test results demonstrate our agent's (X) successful performance against a random strategy opponent (O) in one game for five stages. Our algorithm effectively prioritized strategic positions, beginning with the center and corner placements in the early phases. In the fourth stage (random-4), the agent made a particularly strategic move by placing X in a position that simultaneously blocked the opponent's potential diagonal winning path while strengthening its own horizontal winning opportunity. This dual-purpose move exemplifies the agent's advanced tactical awareness.

This validation confirms our reinforcement learning algorithm has successfully captured the strategic essence of Tic-Tac-Toe. The agent exhibited consistent decision-making, proactive planning, and the ability to create and capitalize on winning opportunities rather than merely reacting to the opponent's moves. The clear victory against random play demonstrates the effectiveness of our training approach and the agent's acquired strategic capabilities.

6.2 Result of Agent vs Greedy Strategy

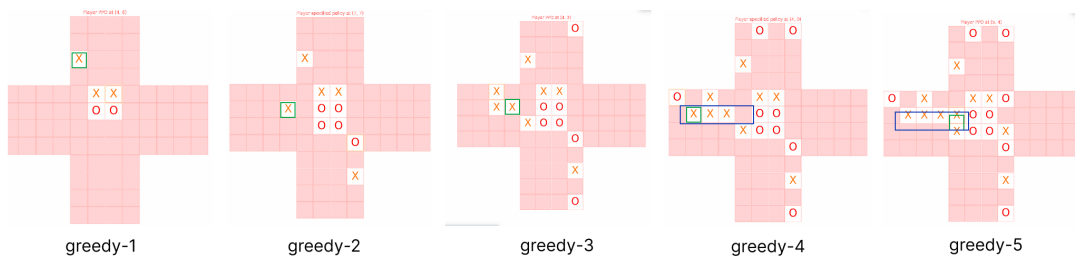


Figure 3: vs Greedy Policy

The image shows our agent (X) successfully defeating a greedy opponent (O) across five game stages. This matchup represents a significantly higher challenge than the random opponent, as the greedy algorithm consistently selects optimal immediate moves.

In the early stages (greedy-1 and greedy-2), our agent established control by securing the center position and strategically placing pieces to create multiple potential winning paths. The greedy opponent responded by blocking immediate threats but failed to anticipate our agent's long-term strategy.

By the middle game (greedy-3), our agent had created a fork situation - establishing two potential winning lines simultaneously. This forced the greedy opponent into a defensive position where it could only block one threat at a time.

The critical turning point came in greedy-4, where our agent executed an excellent tactical move that simultaneously defended against the opponent’s potential winning line while advancing its own winning strategy. This dual-purpose move demonstrates advanced planning capability beyond simple reactive play.

In the final stage (greedy-5), our agent completed its horizontal winning line in the top row, showcasing its ability to maintain strategic focus throughout the game despite the greedy opponent’s best efforts to block immediate threats.

This victory against a competent greedy opponent validates our reinforcement learning approach, demonstrating that our agent has developed sophisticated multi-move planning abilities and can successfully outmaneuver opponents who only optimize for immediate advantage.

7 Conclusion

In conclusion, this project not only demonstrates the successful application of PPO to the Cross Tic-Tac-Toe game, implemented efficiently using the TensorFlow Agents framework, but also provides insights into the practical considerations for implementing reinforcement learning in discrete action space environments. The methodologies developed herein contribute to the broader understanding of how strategic decision-making can emerge from reinforcement learning processes, with potential implications for more complex domains beyond board games.

Code is available in https://github.com/ZhinuoNunu/RL_ttt.

References

- Vijay Konda and John Tsitsiklis. 1999. Actor-critic algorithms. *Advances in neural information processing systems*, 12.
- A. L. Samuel. 1959. [Some studies in machine learning using the game of checkers](#). *IBM Journal of Research and Development*, 3(3):210–229.
- John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. 2015. Trust region policy optimization. In *International conference on machine learning*, pages 1889–1897. PMLR.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.
- David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. 2016. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489.
- David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. 2017. Mastering the game of go without human knowledge. *nature*, 550(7676):354–359.
- Gerald Tesauro et al. 1995. Temporal difference learning and td-gammon. *Communications of the ACM*, 38(3):58–68.
- Ronald J Williams. 1992. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8:229–256.