# Solving Rubik's Cube Without Tricky Sampling

Yicheng Lin[1,*] and Siyu Liang[2]

[1]*School Of Basic Medical Sciences, Fudan University*
[2]*School Of Mechanical Engineering, Tongji University*

December 2, 2024

## Abstract

The Rubik's Cube, with its vast state space and sparse reward structure, presents a significant challenge for reinforcement learning (RL) due to the difficulty of reaching rewarded states. Previous research addressed this by propagating cost-to-go estimates from the solved state and incorporating search techniques. These approaches differ from human strategies that start from fully scrambled cubes, which can be tricky for solving a general sparse-reward problem. In this paper, we introduce a novel RL algorithm using policy gradient methods to solve the Rubik's Cube without relying on near solved-state sampling. Our approach employs a neural network to predict cost patterns between states, allowing the agent to learn directly from scrambled states. Our method was tested on the 2x2x2 Rubik's Cube, where the cube was scrambled 50,000 times, and the model successfully solved it in over 99.4% of cases. Notably, this result was achieved using only the policy network without relying on tree search as in previous methods, demonstrating its effectiveness and potential for broader applications in sparse-reward problems.

**Keywords:** Rubik's cube, Reinforcement learning, Sparse reward

## 1. Introduction

The Rubik's Cube is a classic combinatorial puzzle that presents unique and significant challenges for artificial intelligence and machine learning [1–3]. With an enormous number of possible states, only a single state represents the solved configuration, making it particularly difficult for reinforcement learning algorithms to tackle. The challenge lies in the fact that, during training, the moves taken by the agent's policy are exceedingly unlikely to reach the solved state. This lack of reinforcement for long stretches of exploration often leaves the agent struggling to discover any rewarded state, making it difficult to learn an effective strategy.

This highlights a key issue in sparse reward problems within reinforcement learning, where the absence of frequent feedback severely hinders the agent's ability to learn [4–6]. Without enough guidance, it becomes much harder for the agent to develop useful policies. The Rubik's Cube serves as an ideal testbed for addressing these challenges, and developing reinforcement learning algorithms capable of solving this puzzle could offer broader insights into how to handle sparse-reward environments in other domains.

---

*Author e-mail addresses: `linyc20@fudan.edu.cn`, `cedric_liang@163.com`

While classical methods for solving the Rubik's Cube have existed for decades [7], it wasn't until the work of Agostinelli et al. [3], who developed the DeepCube algorithm, that an AI system using reinforcement learning [8–13] was able to solve the Rubik's Cube from any starting configuration. In their follow-up work, DeepCubeA [2], they implement Deep Approximate Value Iteration (DAVI), which iteratively estimates the cost-to-go for states near the solved state and propagates this information outward to states further away. This process involves scrambling a solved cube multiple times and collecting trajectories that start from the solved configuration. While this approach allows DeepCube to effectively solve the Rubik's Cube, it contrasts with the way humans typically solve the puzzle. Humans start with a fully scrambled cube, without prior knowledge of states close to the solved state. We observe states far from the solution and gradually develop an understanding of the underlying structure and relationships between states, even when they are distant from the solved state.

In this paper, we address this issue by developing a new reinforcement learning approach that uses policy gradient methods to solve the Rubik's Cube without sampling states from a solved configuration. Instead, we continuously sample states from a fully scrambled cube and build up rewards based on the underlying distance patterns between states. Unlike methods mentioned above, which rely on search methods like Monte Carlo Tree Search (MCTS), our approach requires no such search techniques. Using this method, we successfully solved the 2x2x2 Rubik's Cube with a success rate of 99.4% across 50,000 test cases.

## 2. Methods

### 2.1. State and action spaces

A Rubik's Cube consists of a number of stickers, each uniquely associated with a specific position on the cube. In general, for any given Rubik's Cube, all stickers can be encoded as elements of a set of numbers. The state of the Rubik's Cube is then defined as a vector, where each element corresponds to the encoded value of a sticker, recorded in a specific order (Figure 2.1a). This vector is denoted as $s$, with its dimension $N$ representing the total number of stickers. An action performed on a Rubik's Cube involves scrambling the cube hence rearranging specific stickers according to a predefined rule. Using the state vector representation, an action $a$ applied to a given state $s$ can be defined as a function $a(s) = \mathbf{\Gamma}s$, where $\mathbf{\Gamma}$ is a permutation matrix of size $N \times N$, representing the specified rearrangement rule (Figure 2.1b). The action space $\mathcal{H}$ of a Rubik's Cube is defined as the set of all possible actions, each determined by a specific rearrangement rule. Applying an action $a_{[*]}$ in $\mathcal{H}$ to a given state $s_s$ results in a new state $s_t = a_{[*]}(s_s)$. Starting from an initial state $s_0$, repeated applications of actions will cause the Rubik's Cube to transition through a sequence of states. The set of all states that can be reached from $s_0$ is defined as the state space derived from $s_0$, denoted as $S$. Consequently, the topological structure of the Rubik's Cube can be represented as a state graph, where each node corresponds to a state, and edges represent actions connecting these states. Solving the Rubik's Cube can thus be equivalently formulated as a pathfinding problem in the state graph: given a pair of states $(s_s, s_t)$, the goal is to identify a feasible sequence of actions $\langle a_i \rangle_M$ that transforms $s_s$ into $s_t$, with $s_t = a_M(\ldots a_i(\ldots a_2(a_1(s_s))))$, where $a_i \in A, i = 1, 2, \ldots, M$, with $M$ the length of the action sequence (Figure 2.1c).
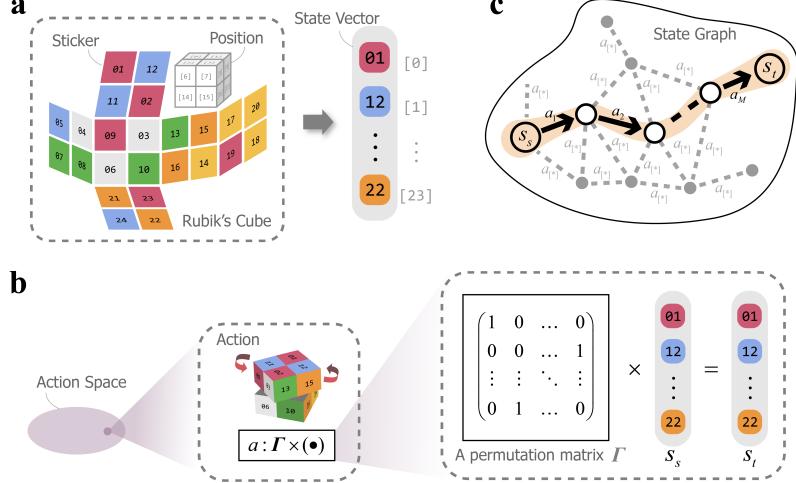
Figure 2.1. General Representation of Rubik's cube problem. a. The state of the cube can be represented by a vector, with each sticker encoded as a number. b. Scrambles are modeled using a permutation matrix applied to the state vector. c. The cube's topological structure is visualized as a state graph, where nodes represent states and edges represent transitions.
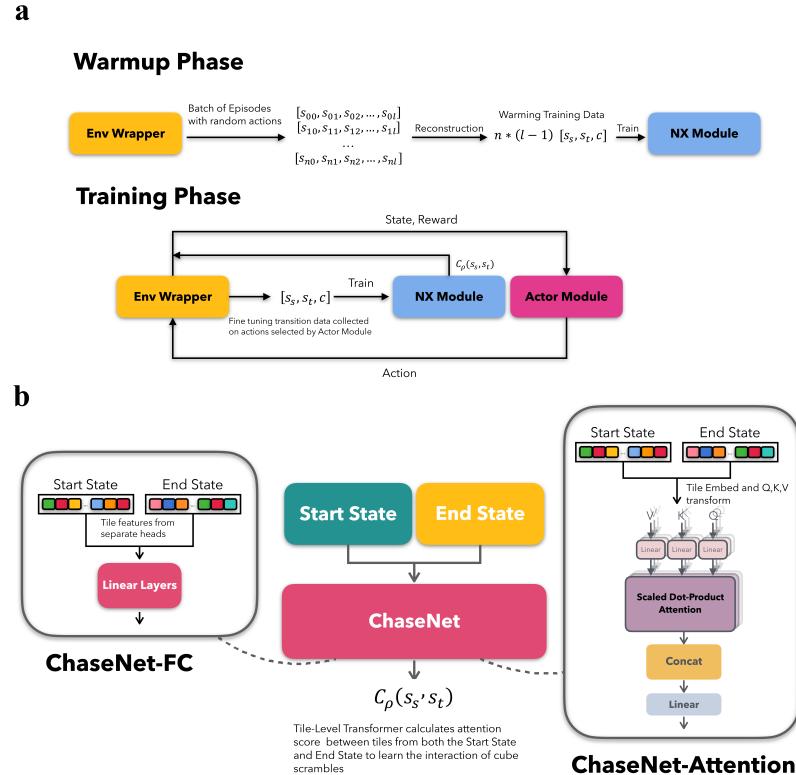


Figure 2.2. The NX Module a. The NX Module training process is divided into a warmup phase and a training phase. b. Two architectural variants of ChaseNet: ChaseNet-FC (fully connected) and ChaseNet-Attention (attention-based).

## 2.2. Cost of states pairs

We define the cost between state pair $(s_s, s_t)$ as the minimal number of scrambles needed from the start state $s_s$ to the target state $s_t$. It can also be represented as the length of the optimal path which linked the state pair $(s_s, s_t)$ in the state graph.

## 2.3. NX Module

The core objective of the NX Module is to collect states far from the solved state and train a model to accurately estimate the cost between any pair of states based on their sticker representation. In this module, a neural network $C_\rho(s_s, s_t)$ termed ChaseNet, parameterized by $\rho$, was trained to estimate the cost between pairs of states $(s_s, s_t)$ using samples generated from scrambled cubes. During the warmup phase, state pairs $(s_s, s_t)$ are created by applying random actions sampled from the action space, with the restriction that no action is repeated more than three times in a row, as this would return the cube to a prior state, producing incorrect labels. We emphasize that the $s_s$ collected in the trajectories are randomly scrambled, without any restriction on their distance from the solved $s_g$, unlike in previous works. In the training phase, state pairs are generated by applying actions sampled from the policy network $p_\theta(a|s)$, allowing the model to refine its predictions based on learned policies.

---

**Algorithm 1** Warmup for NX-Module. **Input**: $B$: Dataset Size, $K$: Maximum number of twists, $J$: training iterations. **Output**: $\rho$: the trained ChaseNet parameters.

---

$\rho \leftarrow$ INITIALIZENETWORKPARAMETERS
**for** $j = 1, 2, \ldots, J$ **do**
    INITIALIZEDATASET $D \leftarrow \emptyset$
    **while** $|D| < B$ **do**
        Generate a random initial scrambled state $s_s$
        $s_{current} \leftarrow s_s$                                       ▷ Set current satate
        **for** $i = 1, 2, \ldots, K$ **do**
            Apply a random scramble to get the next state $s_i$
            Compute the cost $y_i = i$ (the number of twists applied)
            $D \leftarrow D \cup (s_s, s_i, y_i)$                          ▷ Update dataset
            $s_{current} \leftarrow s_i$                              ▷ Update current state
        **end for**
    **end while**
    $\rho \leftarrow \mathrm{Train}(C_\rho, X, y)$ with each $X_i = (s_s, s_i)$ in $D$
**end for**

---

## 2.4. Env Module

The environment is wrapped into the Env Module. During interaction with the agent, it takes in the action $a_i$ and returns the consequent state $s_{i+1}$ and the reward $r_i$, with $i$ the index within the episode. The reward is formulated by the following:

$$r_i = -\log_b C_\rho(s_{i+1}, s_g) \tag{2.1}$$

where $s_g$ is the resolved state and $b = 1.2$ is a base specified to rescale the direct predicted cost of state pair $(s_s, s_g)$. Specifically, when the goal state is reached, a fixed reward of 100 is assigned.

## 2.5. Actor Module

The Actor Module refines the policy network $p_\theta(a|s)$ with Proximal Policy Optimization (PPO) [14] algorithm. Specifically, given the reward returned from the Env Module which is formulated using the cost of the current state to the solved state, the objective function is formulated as following:

$$L\left(s, a, \theta_k, \theta\right) = min\left(\frac{p_\theta(a \mid s)}{p_{\theta_k}(a \mid s)} A^{p_{\theta_k}}(s, a), \ g\left(\epsilon, \ A^{p_{\theta_k}}(s, a)\right)\right) \tag{2.2}$$

where

$$g\left(\epsilon, A\right) = \begin{cases} (1+\epsilon)A & A \geq 0 \\ (1-\epsilon)A & A < 0 \end{cases} \tag{2.3}$$

The advantage $A^{p_\theta}(a|s)$ is the difference between the Q-value $Q^{p_\theta}(s,a)$, the expected return of selection action a in state s, and the value $V^{\varphi_k}(s)$, the predicted return of state s by the critic network parameterized by $\varphi$ in the $k$-th iteration. $\epsilon$ is a hyperparameter that we set at 0.2, as used in the previous paper [14].

---

**Algorithm 2** Training for Actor Module and NX module finetune. **Input**: $C_\rho$: Trained ChaseNet model parameterized by $\rho$, $p_\theta$: Policy network parameterized by $\theta$, $J$:Training iterations, $B$: Batch size, $\epsilon$: PPO clipping parameter, $K$: Number of twists. **Output**: $\rho$: The trained ChaseNet parameters, $\theta$: The trained policy network parameters.

---

**for** $j = 1, 2, \ldots, J$ **do**
    Generate a random initial scrambled state $s_s$
    $s_{current} \leftarrow s_s$                           ▷ Set current state
    $B_f \leftarrow \emptyset$        ▷ Initialize storage buffer for states, actions, rewards and log probabilities
    $D \leftarrow \emptyset$              ▷ Initialize Dataset $D$ for ChaseNet fine tuning
    $i \leftarrow 0$                    ▷ Index of states in an episode
    **while** Episode not end **do**
        Select action $a_i$ $p_\theta(a|s_i)$
        Apply action $a_i$ to get next state $s_{i+1}$
        $r_i \leftarrow -\log_b C_\rho(s_{i+1}, s_g)$       ▷ Predict cost using ChaseNet with $s_g$ the resolved state
        Store $(s_i, a_i, r_i, \log p_\theta(a_i|s_i))$ in $B_f$
        $s_i \leftarrow s_{i+1}$                   ▷ Update state
        $D \leftarrow D \cup (s_s, s_i, i)$           ▷ Collect Data for ChaseNet fine tuning
        $\rho \leftarrow FINETUNE(C_\rho, X, y)$ with each $X_i = (s_s, s_i)$ in $D$
        **if** $|B_f| > B$ **then**
            Compute objective $L_{policy}(s, a, \theta_k, \theta)$ using formula 2.2 with $s$, $a$ sampled from $B_f$
            Update Policy network parameters $\theta$ by maximizing $L_{policy}$
        **end if**
        $i \leftarrow i + 1$
    **end while**
**end for**

---

## 2.6. Network Architectures

For our comparative study, we developed ChaseNet to predict the cost between state pairs, implementing two distinct neural network architectures: ChaseNet-FC, composed entirely of fully

connected layers, and ChaseNet-Attention, which utilizes a sticker-level transformer [15] architecture, as described below. ChaseNet-FC offers a straightforward architecture, with two linear heads that independently extract features from the flattened embedding vectors of the start state $s_s$ and end state $s_t$. These features are then concatenated and passed through additional linear layers to produce the final output cost.

ChaseNet-Attention used a sticker-level transformer that computes attention scores between stickers from both the start state $s_s$ and end state $s_t$. This enables the model to learn the interactions between cube scrambles more effectively. Specifically, batches of vector representations of $s_s$ and $s_t$ are combined to form the input tensor $X$, which is of shape $B \times M$, where $B$ denotes the batch size and $M = 2 \times N$ with $N$ the total number of stickers on the cube. The attention scores $A$ are computed between each position in the combined input tensor $X$:

$$A_{ij} = Attention\left(X_i, X_j\right), \; \forall i, j \in [1, 2, \ldots, M] \tag{2.4}$$

where $X_i$ and $X_j$ correspond to the $i$-th and $j$-th positions in the combined input tensor. This approach allows the model to capture the relationships and dependencies between every sticker's position in the start state and end state. By leveraging these attention scores and the transformer's ability to capture intricate dependencies between sticker positions, ChaseNet-Attention can possibly learn to predict the cost between states more effectively than the approach of ChaseNet-FC.

In the Actor Module, the policy network $p_\theta\left(a|s\right)$ consists of linear layers that extract features from a single state $s$ and output a probability distribution over the discrete action space for the current state.

## 3. Results

### 3.1. Comparison of performance of ChaseNet-FC and ChaseNet-Attention in the warmup phase

In this work, we focus on the 2x2x2 Rubik's Cube as a challenging yet computationally feasible problem for evaluating our approach. We monitored the loss curves of both ChaseNet-FC and ChaseNet-Attention across multiple epochs. With the training iteration set to $J = 1000$, the resulting loss curves are displayed in Figure 3.1a, providing a comparative view of the models' convergence behavior. ChaseNet-FC demonstrates a faster convergence rate than ChaseNet-Attention, likely due to its simpler architecture composed solely of linear layers, which enables more straightforward feature extraction and quicker optimization. In contrast, ChaseNet-Attention, which incorporates a sticker-level transformer for capturing more complex spatial dependencies, initially converges more slowly. However, after 1,000 iterations of warm-up training, ChaseNet-Attention achieves a lower final loss value, suggesting that its architecture, though more complex, ultimately captures richer state representations that improve prediction accuracy. We evaluated both ChaseNet-FC and ChaseNet-Attention on an independent test set to measure Spearman's correlation coefficient between their outputs and the true cost values. ChaseNet-FC achieved a coefficient of 0.834, while ChaseNet-Attention reached a higher coefficient of 0.901, indicating that the sticker-level transformer in ChaseNet-Attention provides a more accurate prediction of state-pair costs.

### 3.2. Solving 2x2x2 Rubik's Cube without searching

After the warm-up phase of training ChaseNet, we used its predictions to train the policy network via the Proximal Policy Optimization (PPO) algorithm. For comparative analysis, we evaluated the
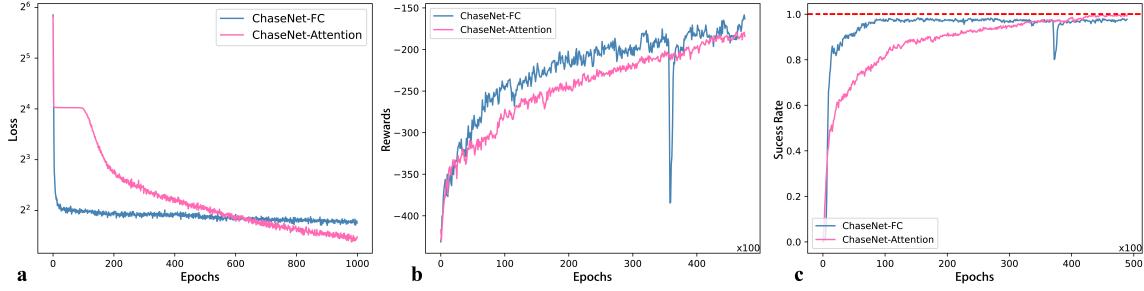
Figure 3.1. a. Warmup loss for ChaseNet-FC and ChaseNet-Attention. b. Average rewards during RL training using rewards from ChaseNet-FC and ChaseNet-Attention. c. Success rate during RL training using rewards from ChaseNet-FC and ChaseNet-Attention.

performance of the policy network trained using rewards from both ChaseNet-FC and ChaseNet-Attention. During testing, we scrambled the cube and measured the number of times the policy network successfully solved it. Notably, this testing phase involved no tree search; instead, we relied solely on the policy network to guide each move. Despite the absence of search, the policy network achieved a remarkable success rate. Here, the success rate is defined as the ratio of successful solves to the total number of 50 test attempts. Figure 3.1b-c shows the performance of both ChaseNet-FC and ChaseNet-Attention. During the final validation, we achieved a success rate of over 99.4% across 50,000 test cases.

## 4. Discussion

In this work, we presented a reinforcement learning approach for solving the Rubik's Cube without relying on sampling starting from solved states. Unlike methods like DeepCubeA, which sample states near the solution, our approach learns a solution policy directly from learning the cost patterns of fully scrambled states. By leveraging ChaseNet to accurately estimate state transition costs, the NX Module guides policy optimization effectively from random starting points, aligning with the way humans solve the puzzle by building heuristics from disordered states. Achieving a success rate over 99.4% on the 2x2x2 Rubik's Cube, this method demonstrates the potential to address sparse-reward challenges without complex sampling or search methods.

Despite these promising results, our approach currently has some limitations. We focused on the 2x2x2 Rubik's Cube to test feasibility within a manageable state space. Scaling to the 3x3x3 cube would require a much larger neural network to estimate costs across a more complex space, posing a key challenge for future work. Furthermore, while the 2x2x2 cube provided an initial testbed, broader testing across different sparse-reward environments will be essential to assess the generalizability and practicality of this method. Applying our approach to various domains could demonstrate its robustness for other sparse-reward tasks where rewards are infrequent or difficult to access.

In summary, our approach offers a new method for solving sparse-reward problems by optimizing policies from fully scrambled states, without relying on search or solved-state sampling. These results lay a foundation for developing more flexible and scalable methods, though further testing on complex puzzles and varied environments will be important to confirm its wider applicability.

# References

[1] Wolfgang Konen. Towards learning rubik's cube with n-tuple-based reinforcement learning, 2023.

[2] Alexander Shmakov Pierre Baldi Forest Agostinelli, Stephen McAleer. Solving the rubik's cube with deep reinforcement learning and search. *nature machine intelligence*, 1:356–363, 2019.

[3] Stephen McAleer, Forest Agostinelli, Alexander Shmakov, and Pierre Baldi. Solving the rubik's cube with approximate policy iteration. In *International Conference on Learning Representations*, 2019.

[4] Gautham Vasan, Yan Wang, Fahim Shahriar, James Bergstra, Martin Jagersand, and A. Rupam Mahmood. Revisiting Sparse Rewards for Goal-Reaching Reinforcement Learning. *arXiv e-prints*, page arXiv:2407.00324, June 2024.

[5] Martin Riedmiller, Roland Hafner, Thomas Lampe, Michael Neunert, Jonas Degrave, Tom Van de Wiele, Volodymyr Mnih, Nicolas Heess, and Jost Tobias Springenberg. Learning by Playing - Solving Sparse Reward Tasks from Scratch. *arXiv e-prints*, page arXiv:1802.10567, February 2018.

[6] Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob McGrew, Josh Tobin, Pieter Abbeel, and Wojciech Zaremba. Hindsight Experience Replay. *arXiv e-prints*, page arXiv:1707.01495, July 2017.

[7] Tomas Rokicki, Herbert Kociemba, Morley Davidson, and John Dethridge. The diameter of the rubik's cube group is twenty. *SIAM Review*, 56(4):645–670, 2014.

[8] D.P. Bertsekas and J.N. Tsitsiklis. Neuro-dynamic programming: an overview. In *Proceedings of 1995 34th IEEE Conference on Decision and Control*, volume 1, pages 560–564 vol.1, 1995.

[9] J. Bagnell, Sham M Kakade, Jeff Schneider, and Andrew Ng. Policy search by dynamic programming. In S. Thrun, L. Saul, and B. Schölkopf, editors, *Advances in Neural Information Processing Systems*, volume 16. MIT Press, 2003.

[10] Sham M. Kakade and John Langford. Approximately optimal approximate reinforcement learning. In *International Conference on Machine Learning*, 2002.

[11] Bruno Scherrer. Approximate Policy Iteration Schemes: A Comparison. *arXiv e-prints*, page arXiv:1405.2878, May 2014.

[12] Alessandro Lazaric, Mohammad Ghavamzadeh, and Rémi Munos. Analysis of classification-based policy iteration algorithms. *Journal of Machine Learning Research*, 17(19):1–30, 2016.

[13] Martin L. Puterman and Moon Chirl Shin. Modified policy iteration algorithms for discounted markov decision problems. *Management Science*, 24(11):1127–1137, 1978.

[14] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal Policy Optimization Algorithms. *arXiv e-prints*, page arXiv:1707.06347, July 2017.

[15] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention Is All You Need. *arXiv e-prints*, page arXiv:1706.03762, June 2017.