# From Game-Playing to Self-Driving: Comparing AlphaGo vs AlphaZero Approaches for Driving Controls

Ellen Xu, Robert J. Moss, Mykel J. Kochenderfer

{ellenjxu, mossr, mykel}@stanford.edu

#### **Stanford University**

## Abstract

We investigate approaches comparing AlphaGo-style methods, initialized from human policies, with AlphaZero-style learning from scratch for real-world control tasks. While AlphaZero achieved superior performance in Go without human data, we hypothesize that for human-centered environments, human policies can encode safety constraints and behavioral priors difficult to capture in reward functions alone. We evaluate these approaches on a realistic driving simulator using a PID controller as our human-level baseline. Our results show that human-guided Monte Carlo Tree Search (MCTS) significantly outperforms the PID controller, achieving 23% higher rewards. Importantly, standard AlphaZero Continuous fails to converge due to exploration instability. We explore two key components for stable convergence: guided exploration and guided rollouts. These findings suggest that human priors may provide crucial constraints for safe and efficient learning in real-world reinforcement learning applications.

## 1 Introduction

Deep reinforcement learning has achieved superhuman performance in games and robotics (Silver et al., 2018; Kober et al., 2013). In the game of Go, AlphaZero ultimately outperformed AlphaGo by learning completely from self-play, while AlphaGo first trained on human expert games before reinforcement learning. However, for real-world domains like autonomous driving or industrial robotics which are inherently human-centered, the optimal learning strategy remains largely unexplored.

Does leveraging priors from human policies outperform learning from scratch in human-centered control environments?

We hypothesize that real-world control tasks benefit from human initialization because human policies inherently encode safety constraints and preferred behaviors that are difficult to specify through reward functions alone. Unlike games which have well-defined rules and rewards, real-world control systems often operate in noisy, safety-critical environments. This leads to constraints on exploration which can hinder learning for reinforcement learning systems.

In this work, we investigate this hypothesis by comparing AlphaGo-style methods, which utilize human behavioral priors, with AlphaZero-style learning from scratch on a realistic driving control task. Specifically, we evaluate steering control using the benchmark Comma AI Controls Challenge (Comma AI, 2024), which simulates real-world vehicle dynamics with a GPT-based model trained on customer data.

# 2 Background

In this section, we provide background on key algorithm foundations used in our experiments.

### 2.1 Monte Carlo Tree Search

Monte Carlo Tree Search (MCTS) is an online search and planning algorithm (Coulom, 2007). Unlike exhaustive search, MCTS selectively explores promising parts of the search tree, and caches local statistics for each node. The four key stages of MCTS are: *selection* using Upper Confidence Bound for Trees (UCT), *expansion* of leaf nodes, *simulation* via random rollouts to terminal states, and *backpropagation* of rewards to update node statistics. After iteratively building up a search tree, a final action is selected from the current state (i.e., root node) and executed in the environment.

# 2.2 AlphaZero and AlphaGo

AlphaZero has achieved state-of-the-art, super-human performance in Chess, Shogi, and the game of Go (Moerland et al., 2018; Schrittwieser et al., 2020; Silver et al., 2018). The key innovation of AlphaZero is the augmentation of Monte Carlo tree search (MCTS) with a neural network generated policy. The augmented policy can learn to generalize from self-play, and progressively bootstraps the network from random initialization up to superhuman play, without requiring extra human input. AlphaGo, in contrast, first learned from human expert games through supervised learning before improving via self-play reinforcement learning.

# 2.3 Controls Challenge

Classical control theory provides mathematical frameworks for regulating system behavior, with Proportional-Integral-Derivative (PID) controllers remain widely used in practice due to their simplicity and effectiveness (Åström & Hägglund, 2006). However, autonomous vehicle control presents challenges including nonlinear dynamics, environment noise, and real-time performance requirements, which often make conventional control methods insufficient (Rajamani, 2011).

The Comma AI Controls Challenge (Comma AI, 2024) is a realistic driving simulator built using real-world driving data (illustrated in Figure 1). The goal is to output torque control to steer a car along a desired trajectory, while minimizing control effort. To simulate the car's steering responses, the benchmark uses a GPT-based simulator trained on real-world noisy dynamics from customer vehicles.



Figure 1: Control challenge diagram showing interaction between learned policy and simulator.

The controls challenge captures three critical aspects of real-world RL: (1) high-dimensional continuous state-action spaces, (2) sensitive dynamics where errors accumulate rapidly, and (3) realistic noise from a simulator trained on real vehicle data.

# 3 Related Work

Reinforcement learning (RL) approaches have shown promise in continuous control tasks (Lillicrap et al., 2016; Haarnoja et al., 2018; Schulman et al., 2017). These methods can learn control policies directly from interaction with the environment, which are generic solutions that can outperform classical hand-crafted controllers. Model-free RL methods learn a direct mapping from states to actions, and model-based methods utilize an environment simulator to derive a controller(Wang et al., 2019).

However, learning-based controllers face challenges in safety-critical applications like autonomous driving, where balancing safe exploration and exploitation becomes difficult. This has led to interest in hybrid methods, which combine model-based planning with learned components. These hybrid methods offer several advantages: they can separate the exploration and learning phases, provide safer exploration through model planning, and achieve greater data-efficiency than pure model-free approaches (Wang et al., 2019).

AlphaZero is a hybrid approach which combines neural networks with tree-based search. Recent work has extended AlphaZero beyond discrete games. AlphaZero Continuous (A0C) adapted the algorithm to continuous action spaces by replacing self-play with simulation and using progressive widening for action expansion (Moerland et al., 2018). Moss et al. (2024) further generalized these methods to partially observable environments.

In autonomous driving, Hoel et al. (2020) applied AlphaGo Zero for discrete decision-making in highway lane changes. More recently, Cusumano-Towner et al. (2025) demonstrated large-scale self-play reinforcement learning for driving policies, by simulating copies of its own agent for robust policy learning. These works illustrate that the core ideas of AlphaZero—self-play, MCTS search, and neural network learning—can be successfully applied to self-driving. However, the approaches focus on high-level planning rather than low-level continuous control, which presents unique challenges including noise sensitivity and safety requirements.

In previous work, Xu (2024) used model-free RL methods to learn a low-level autonomous vehicle controller on the controls challenge. While Proximal Policy Optimization (PPO) faced difficulties with efficient exploration in the large continuous action space, evolutionary methods were able to converge reliably, yet do not scale well to larger parameter counts (Schulman et al., 2017; Hansen, 2023). The challenge of safe and efficient exploration motivate our hybrid learning approach, which combines human behavior priors with planning over an internal model of vehicle dynamics.

# 4 Methods

We compare three approaches: (1) a tuned PID controller as our human baseline, (2) standard AlphaZero learning from scratch, and (3) AlphaGo Control with human guidance.

# 4.1 Baseline: PID Controller

We use a Proportional-Integral-Derivative (PID) controller as our human-level policy baseline. The controller parameters ( $K_p = 0.3, K_i = 0.05, K_d = -0.1$ ) were manually tuned for the environment task and encode human preferences for desired control behavior.

# 4.2 AlphaZero Implementation

We developed a parallelized AlphaZero Continuous (A0C) implementation combining MCTS with neural network learning. Our implementation differs from Moerland et al. (2018) since we use Q-value estimates rather than normalized visit counts to construct the target policy, and we employ a softmax temperature to control the sharpness of the distribution. The policy network learns from

MCTS search results using cross-entropy loss:

$$\mathcal{L}_{\text{policy}} = \mathbb{E}_{s \sim \mathcal{D}} \left[ -\sum_{a \in \mathcal{A}} \pi_{\text{MCTS}}(a|s) \log \pi_{\theta}(a|s) \right]$$

where  $\pi_{\text{MCTS}}(a|s) = \text{softmax}(Q(s, a)/\tau)$  is the target policy derived from Q-values with temperature  $\tau$ , and  $\pi_{\theta}(a|s)$  is the neural network policy. The temperature parameter prevents overconfident target policies when Q-value estimates have high variance.

The value network is trained to predict expected returns minimizing mean squared error:

$$\mathcal{L}_{\text{value}} = \mathbb{E}_{s \sim \mathcal{D}} \left[ (V_{\phi}(s) - z)^2 \right]$$

where z is the discounted return from MCTS rollouts and  $V_{\phi}(s)$  is the value network prediction.

The policy and value network are distilled from the MCTS exploration. The networks enable generalization and improves MCTS efficiency by (1) reducing search breadth (policy network biases exploration toward promising actions) and (2) reducing search depth (truncating rollouts and using value network leaf node estimates).

#### 4.3 AlphaGo Control

Standard MCTS planning performed poorly in the control environment due to action sensitivity in closed-loop feedback systems. Search depths beyond 2-3 steps caused rapid error accumulation, leading to extreme Q-value estimates and system divergence (Figure 4c). We experimented with two key modifications to address the instability:

**Guided exploration.** Instead of initializing with sampling from a random policy, we sample actions using Gaussian exploration around the PID policy output:

$$a_t \sim \mathcal{N}(\pi_{\text{PID}}(s_t), \sigma^2)$$

where  $\pi_{\text{PID}}$  is the PID controller policy and  $\sigma^2$  is the action variance, controlling the amount of exploration around the selected action.

*Rationale.* While discrete AlphaZero uses policy priors  $\pi_{\theta}(a|s)$  as probability weights for action selection, continuous action spaces have unbounded probability densities that cannot be directly used as weights. Our method achieves similar bias toward promising actions by sampling from a Gaussian centered on the expert policy, ensuring exploration remains within a reasonable action region while preserving the ability to discover improvements.

**Guided rollouts.** We replace random simulation rollouts used in standard AlphaZero with PIDguided rollouts. At each timestep, we explore alternative actions from the root node using the action sampling above, but for all descendant nodes, we follow the PID policy:

$$a_t \sim \begin{cases} \mathcal{N}(\pi_{\text{PID}}(s_t), \sigma^2) & \text{if root node} \\ \pi_{\text{PID}}(s_t) & \text{otherwise} \end{cases}$$

*Rationale.* This formulation is well-suited for policy improvement because: (1) it provides low-variance value estimates by leveraging the known PID policy for rollouts, in comparison with Monte Carlo rollouts which are high-variance, and (2) it enables direct policy improvement by comparing the value of alternative actions against the current policy, following the standard Q-learning paradigm where  $\pi'(s) = \arg \max_a Q^{\pi}(s, a)$  yields policy improvement when  $Q^{\pi}(s, a)$  estimates are accurate.

These modifications enable more stable policy improvement through  $Q^{\pi}(s, a)$  estimates. The approach preserves MCTS's ability to discover improved policies, while avoiding the exploration instability in standard tree search.

# 5 Results

The Comma AI Controls Challenge (Comma AI, 2024) is used as the test benchmark to evaluate learned policies.

# 5.1 Planning Performance

We first evaluate the planning capabilities of different MCTS variants. Human-guided MCTS achieved 23% lower cost than the PID baseline, while standard MCTS failed to converge due to exploration instability (Table 1).

Table 1: Planning method comparison. Results averaged over 10 seeded evaluation rollouts (200 steps each) with search depth d = 3 and simulations n = 10. MCTS planning requires 8 times longer execution time per training step (0.627s vs 0.075s for PID baseline), highlighting the tradeoff between planning accuracy and computational efficiency.

Algorithm	Total Cost ( $\pm$ std)	Runtime/Step
PID (baseline)	80.17 (± 95.04)	0.075s
Standard MCTS	6443.13 (± 4657.55)	0.605s
Human-guided MCTS	<b>61.70</b> (± 42.64)	0.627s

# 5.2 Policy Learning Results

After distilling MCTS exploration into neural networks, we evaluated the learned policies over 5,000 rollouts. AlphaGo Control significantly outperformed the standard AlphaZero (Table 2).

Table 2: The policy cost is calculated over 5,000 trajectory rollouts. The resulting mean and standard deviation are are reported across 10 random seeds to ensure statistical significance. AlphaZero training used 100 workers and 30 iterations.

Algorithm	Total Cost	Lat Accel Cost	Jerk Cost
Standard AlphaZero	$\begin{array}{l} 4954.00 \ (\pm \ 846.50) \\ \textbf{197.68} \ (\pm \ 8.50) \end{array}$	98.55 (± 17.48)	<b>26.86</b> (± 2.16)
AlphaGo Control		<b>3.36</b> (± 0.15)	29.64 (± 3.17)

Figure 2 shows that AlphaGo Control achieves stable learning with consistent reward improvement, while standard AlphaZero fails to converge. This may be due to the exploration-exploitation dilemma, where deviating from stable policies leads to error accumulation, preventing learning from random initialization.

# 5.3 Ablation Studies

We validated our design choices through ablation studies removing (A) guided exploration and (B) guided rollouts. Figure 3 shows that both components are essential to stable learning.

Qualitatively, ablation A (i.e., no guided exploration) learns oscillatory control behavior, while ablation B (no guided rollouts) learns conservative policies that produce minimal control actions (Figure 5). The oscillatory behavior may be due to poor initial actions from a randomly initialized policy. The conservative behavior in ablation B suggests that high-variance value estimates may lead to similar Q-values across actions, causing the policy to default to low-magnitude control outputs to minimize immediate penalties. Furthermore, the difference in policy behaviors suggest multiple local optima exist, and human priors can help select for behaviorally optimal solutions.



Figure 2: Learning curves comparing AlphaGo Control (human-guided) versus standard AlphaZero (from scratch). Left: episode rewards during training showing stable improvement for AlphaGo Control versus diverging returns for AlphaZero. Right: value network loss convergence.



Figure 3: Learning curves for ablation studies. The AlphaGo Control learning algorithm is shown in red. Without guided exploration or action sampling (green), random exploration fails to discover improvements in the large continuous action space. Without guided rollouts (blue), high-variance returns prevent effective value network learning, resulting in flat performance curves.

# 6 Conclusions

Our results suggest that human demonstrations can provide crucial constraints for safe exploration in sensitive control environments. Importantly, standard AlphaZero Continuous fails to converge due to exploration instability. We demonstrate in ablations two key components for stable convergence: (1) guided exploration and (2) guided rollouts, both achieved by using human priors. These components were used in the AlphaGo training framework to train a control policy, which showed greater stability over standard AlphaZero learning from scratch.

This finding suggests that human priors can be effective for safe exploration in real-world control tasks. Unlike games with well-defined rules, control systems require careful balance between exploration and safety in order to achieve stable learning, which human priors may naturally provide. While standard MCTS diverges rapidly in the control setting, human-guided search achieves stable improvement with 23% better performance than baseline PID control. While we used a PID controller as our human policy, this could be replaced with other heuristics or initializing a policy learned through imitation learning from human data. These can extend beyond driving controls to other human-centered domain where safety and exploration in low-level continuous state spaces is important.

Our work provides preliminary evidence that in real-world applications, behavior priors—such as those used in AlphaGo—lead to safer and more stable learning compared to from-scratch approach like AlphaZero. However, there are several limitations to our work. First, we focus on a single control task with limited scope and complexity. Second, computational constraints restricted our MCTS search depth, potentially limiting the full potential of from-scratch learning methods. Third, our comparison uses a relatively simple PID baseline rather than more sophisticated hand-crafted or imitation learning policies. Future work aims to scale up search and generalize to other tasks to better evaluate from scratch and human-guided methods.

# Appendix

Additional implementation details, training hyperparameters, and figures are included in the following sections.

#### 6.1 Problem Formulation

The controls challenge is formulated as a finite-horizon MDP  $\mathcal{M} = \{S, \mathcal{A}, T, R, \gamma, T_h\}$ , where  $S \subseteq \mathbb{R}^{n_s}$  is the continuous state space,  $\mathcal{A} \subseteq \mathbb{R}^{n_a}$  is the continuous action space, T(s' | s, a) represents the GPT-based transition model, R(s, a) is the deterministic reward function,  $\gamma \in [0, 1)$  is the discount factors, and  $T_h$  is the time horizon. The reward function is defined as:

$$R(s, a) = -|a_{\text{lateral}}(s, a) - a_{\text{desired}}| - \lambda_u ||a||$$

The objective is to find the policy that maximizes the expected cumulative discounted reward:

$$\pi^* = \arg \max_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=0}^{T_h - 1} \gamma^t R(s_t, a_t) \right]$$

The controller receives ego vehicle state  $s_t = \{v_{ego}, a_{ego}, roll\}$ , current and desired acceleration  $(a_{current} \text{ and } a_{desired})$ , and 20 seconds of future states and past context  $(s_{t-20:t} \text{ and } s_{t:t+20})$ . The controller output u is a continuous steering torque corresponding to a steering wheel input.

#### 6.2 Gym Environment Implementation

A custom Gym wrapper for the Control Challenge is used for training the RL algorithms. The key modifications from the Control Challenge (2.3) include: formulating the cost for each trajectory rollout into a dense reward for each timestep, support for batched environment observations and forward pass in the simulator, and utilizing a custom control state as a context window of past states. More specifically, the custom ControlState class is used as the internal simulator for MCTS with the following state representation: context window of state history, action history, and current lateral acceleration history  $o_{t-C:t}$ ,  $a_{t-C:t}$ ,  $s_{t-C:t}$  where C is the context window of past 20 timesteps.

#### 6.3 Network Architecture

Both actor and critic networks use multi-layer perceptrons with [256, 256] hidden layers and ReLU activation. The input to the network is 304-dimensional tensor, which includes:

- Current state: target/current lateral acceleration, roll, ego velocity/acceleration (5 dimensions)
- Context history: past 20 timesteps of states, actions, and lateral accelerations (100 dimensions)
- Future plan: next 5 seconds of target accelerations and predicted states (196 dimensions)
- PID terms: current error, integral error, derivative error (3 dimensions)

The actor outputs a Gaussian policy with learnable log standard deviation initialized to -2.

## 6.4 Training Hyperparameters

Key training hyperparameters include: learning rate  $1 \times 10^{-3}$ , batch size 10,000, 100 parallel workers, 30 training iterations. MCTS used a depth of 3, 10 simulations per action, action sampling variance 0.1, and discount factor of  $\gamma = 0.99$ . Full details are provided in Table 3.

Parameter	Default Value	Description
Training Parameters		
Max iterations	30	Maximum training iterations
Episode steps	200	Steps per episode during rollout
Workers	100	Number of parallel tree workers
Learning rate	$1 \times 10^{-3}$	Adam optimizer learning rate
Training epochs	10	Epochs per iteration
Batch size	10,000	Training batch size
Replay buffer size	100,000	Maximum replay buffer capacity
Network Architecture		
Hidden sizes	[256, 256]	Hidden layer dimensions for actor/critic
Log std	-2	Initial log standard deviation for actor
L2 regularization	$1 \times 10^{-4}$	L2 penalty coefficient
Value loss weight	0.5	Weight for value function loss
MCTS Parameters		
Exploration weight	0.1	UCB exploration constant
Discount factor ( $\gamma$ )	0.99	Reward discount factor
MCTS depth	3	Search tree depth
MCTS simulations	10	Number of simulations per action
Action variance	0.1	Variance for action sampling
UCB parameter (k)	2	UCB formula parameter
Dirichlet alpha ( $\alpha$ )	0.5	Dirichlet noise parameter
Policy Training		
Temperature $(\tau)$	1	Softmax temperature for policy targets
Entropy coefficient	0.001	Entropy regularization weight
Evaluation		
Evaluation episodes	5	Episodes for evaluation during training
Evaluation steps	200	Steps per evaluation episode

Table 3: Hyperparameters used in AlphaGo Control training

#### 6.5 Trajectory Analysis

Examples of control trajectories and qualitative analysis are provided. Figure 4 compare trajectories from the planning algorithms versus the baseline PID policy trajectories. Figures 5 and 6 compare the trajectories from ablation study.

# Acknowledgments

This work was supported by the Stanford Intelligent Systems Laboratory for additional compute and resources.



(a) Human-guided MCTS rollout

(b) PID rollout



(c) Standard MCTS divergence

Figure 4: Comparison of trajectories for MCTS planning and baseline PID controller. Desired tracking trajectory is in blue and current lateral acceleration trajectory in red. (a) Guided MCTS exploration around PID actions can lead to discovering better actions. (b) This improvement is particularly evident during rapid acceleration changes, where PID tends to overshoot. (c) Using standard MCTS as an online planner (without additional modifications for stability), the lateral acceleration diverges after a few timesteps. A search depth of 10 and 100 simulations per step were used for standard MCTS, and search depth of 3 and 10 simulations per step for human-guided MCTS which is more efficient.



Figure 5: Sampled trajectories for ablations. Top row: ablation A (no action sampling) produces oscillatory behavior. Bottom row: ablation B (no guided rollouts) produces near-zero control output.



Figure 6: Sampled trajectories for AlphaGo Control learned policy. The controller is able to properly track lateral acceleration, as shown in the red and blue trajectories.

## References

- Karl Johan Åström and Tore Hägglund. Advanced PID Control. ISA-The Instrumentation, Systems, and Automation Society, 2006.
- Comma AI. Controls Challenge v2. https://github.com/commaai/controls\_ challenge, 2024. Accessed: 2025-06-24.
- Rémi Coulom. Efficient Selectivity and Backup Operators in Monte-Carlo Tree Search. In Computers and Games. Springer, 2007.
- Marco Cusumano-Towner, David Hafner, Alex Hertzberg, Brody Huval, Aleksei Petrenko, Eugene Vinitsky, Erik Wijmans, Taylor Killian, Stuart Bowers, Ozan Sener, Philipp Krähenbühl, and Vladlen Koltun. Robust Autonomy Emerges from Self-Play. arXiv preprint arXiv:2502.03349, 2025.
- Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor. In *International Conference on Machine Learning (ICML)*, pp. 1861–1870. PMLR, 2018.
- Nikolaus Hansen. The CMA Evolution Strategy: A Tutorial, 2023. URL https://arxiv.org/ abs/1604.00772.
- Carl-Johan Hoel, Katherine Driggs-Campbell, Krister Wolff, Leo Laine, and Mykel J. Kochenderfer. Combining Planning and Deep Reinforcement Learning in Tactical Decision Making for Autonomous Driving. *IEEE Transactions on Intelligent Vehicles*, 5(2):294–305, 2020.
- Jens Kober, J Andrew Bagnell, and Jan Peters. Reinforcement learning in robotics: A survey. *The International Journal of Robotics Research*, 32(11):1238–1274, 2013.
- Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. In *International Conference on Learning Representations (ICLR)*, 2016.
- Thomas M. Moerland, Joost Broekens, Aske Plaat, and Catholijn M. Jonker. AOC: Alpha Zero in Continuous Action Space. *arXiv preprint arXiv:1805.09613*, 2018.

- Robert J. Moss, Anthony Corso, Jef Caers, and Mykel J. Kochenderfer. BetaZero: Belief-State Planning for Long-Horizon POMDPs using Learned Approximations. *Reinforcement Learning Journal (RLJ)*, 1:158–181, 2024.
- Rajesh Rajamani. Vehicle Dynamics and Control. Mechanical Engineering Series. Springer Science & Business Media, 2nd edition, 2011.
- Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Simon Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, Timothy Lillicrap, and David Silver. Mastering Atari, Go, chess and shogi by planning with a learned model. *Nature*, 588(7839):604–609, 2020.
- 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, Timothy Lillicrap, Karen Simonyan, and Demis Hassabis. A general reinforcement learning algorithm that masters chess, shogi, and Go through self-play. *Science*, 362(6419):1140–1144, 2018.
- Tingwu Wang, Xuchan Bao, Ignasi Clavera, Jerrick Hoang, Yifan Wen, Eric Langlois, Shunshi Zhang, Guodong Zhang, Pieter Abbeel, and Jimmy Ba. Benchmarking Model-Based Reinforcement Learning. arXiv preprint arXiv:1907.02057, 2019.
- Ellen Xu. RL for Car Controls. https://blog.comma.ai/rlcontrols/, 2024. Accessed: 2025-06-24.