DT-Guided DRL: A Transition from Utility-based Decision Theory to Deep **Reinforcement Learning**

Zelin Wan¹, Nithin Alluru¹, Jin-Hee Cho¹, Mu Zhu², Ahmed H. Anwar³, Charles

Kamhoua³, Munindar P. Singh⁴,

¹Virginia Tech

²Chinese Academy of Sciences

³DEVCOM Army Research Laboratory

⁴North Carolina State University

{zelin, nithin, jicho}@vt.edu, zhumu@cnic.cn, {ahmed.h.hemida.ctr, charles.a.kamhoua.civ}@army.mil, mpsingh@ncsu.edu

Abstract

This paper presents Decision Theory-guided Deep 1 Reinforcement Learning, called DT-guided DRL, to 2 3 address the cold start problem in DRL. By incorporating decision theory, DT-guided DRL improves 4 agents' initial performance, robustness, learning ef-5 ficiency, and reliability in complex environments. 6 We examine three problem contexts: cart pole, 7 maze navigation, and inverted pendulum. Ex-8 perimental results show that DT-guided DRL en-9 hances initial guidance and exploration. During 10 early training, it achieves at least 184% higher ac-11 cumulated reward in the cart pole problem and at 12 least 279% higher accumulated reward in the maze 13 problem. It also provides faster convergence in all 14 problem contexts. This approach leverages human-15 informed knowledge, offering a foundation for fur-16 ther research. 17

Introduction 1 18

Decision theory, particularly utility theory [Fishburn et al., 19 1979], has long been a cornerstone in fields like economics, 20 management science, and cybersecurity [Rădulescu and oth-21 ers, 2020], offering a structured framework for decision-22 making under uncertainty [North, 1968]. However, its lim-23 itations in handling complex environments, due to oversim-24 plified models and assumptions, restrict its broader appli-25 cability. In contrast, Deep Reinforcement Learning (DRL) 26 has demonstrated remarkable success in learning intricate be-27 haviors through neural networks that model complex state-28 action dynamics, particularly in some foundational robot con-29 trol tasks such as balancing (Inverted Pendulum), navigation 30 (Maze), and coordination (Cart Pole). Yet, the long train-31 ing times and initial random exploration of DRL pose signif-32 icant challenges for practical deployment, especially in criti-33 cal, real-world scenarios like robotics where early failures are 34 unacceptable [Dulac-Arnold et al., 2019]. 35

For researchers and practitioners in robotics, integrating 36 DRL offers transformative potential for enabling autonomous 37

systems to learn complex tasks. However, poor early performance and extensive training hinder its adoption in robotics, 39 where safety and efficiency are paramount. Addressing these 40 limitations is crucial for unlocking DRL's full potential in 41 real-world applications.

38

42

43

44

45

46

47

48

49

67

77

To address these challenges, we propose Decision Theoryguided Deep Reinforcement Learning, called DT-guided DRL, which incorporates utility functions from decision theory to guide early exploration. This approach mitigates the cold start problem and ensures consistent performance during training, particularly beneficial in robotics where trial-anderror learning can be costly or dangerous.

Our proposed DT-guided DRL approach aims to facilitate 50 the transition from decision theory to DRL and mitigate the 51 cold start problem. By integrating utility functions into the 52 DRL framework, our method allows researchers and practi-53 tioners to leverage existing decision theory knowledge while 54 benefiting from DRL's advanced capabilities, easing adoption 55 through a more gradual transition. DT-guided DRL uses util-56 ity functions to guide exploration and action selection in the 57 early learning stages, steering the agent toward promising re-58 gions of the state space and avoiding poor initial performance. 59

We evaluate the effectiveness of our approach using three 60 benchmark environments: the classic cart pole, maze naviga-61 tion tasks, and inverted pendulum. Those serve as proof-of-62 concept domains to compare DT-guided DRL against stan-63 dard DRL techniques and relevant baselines, including trans-64 fer learning [Zhu and others, 2023], sample efficiency [Lin et 65 al., 2021], and imitation learning [Gros and others, 2020]. 66

This work presents the following key contributions:

- 1. We introduce a novel integration of utility theory, a fun-68 damental component of decision theory, with the DRL 69 framework to enhance its applicability to robotics. This 70 is the first work to systematically merge these paradigms 71 in robotics, establishing a structured transition from 72 decision-theoretic models to DRL. This integration en-73 ables informed early-stage exploration, reducing ineffi-74 cient random search and accelerating learning in rein-75 forcement learning environments. 76
- 2. We experimentally demonstrate that DT-guided DRL ef-

fectively mitigates the cold start problem, a major chal-78 lenge in reinforcement learning. Through extensive evalu-79 ations of cart poles, maze navigation, and inverted pendu-80 lum, representing robotic control, navigation, and balanc-81 ing, we show that our method achieves significantly higher 82 initial performance, faster convergence, and greater sam-83 ple efficiency than standard DRL techniques and relevant 84 baselines, including transfer learning, imitation learning, 85 and sample-efficient RL. 86

3. Our approach leverages human-informed knowledge to
provide structured guidance during early exploration, improving learning stability and robustness in complex environments. This ensures that reinforcement learning can be
more reliably applied to real-world robotics, where trialand-error approaches are costly or unsafe.

4. By bridging decision theory and DRL, our work lays 93 the foundation for a new research direction at the in-94 tersection of both fields. The proposed framework not 95 only enhances reinforcement learning efficiency but also 96 opens new avenues for integrating expert knowledge and 97 decision-theoretic principles into autonomous robotic sys-98 tems, making AI-driven decision-making more reliable, 99 interpretable, and scalable. 100

101 2 Related Work

This section overviews existing approaches for enhancing
 DRL efficiency by reducing convergence time, focusing on
 prior-knowledge-enhanced initialization, transfer learning, imitation learning, and *sample efficiency.*

Prior-knowledge enhanced initialization improves DRL 106 efficiency by leveraging prior knowledge. Silva and Gom-107 bolay [2021] integrated domain-specific rules into RL, boost-108 ing performance without extensive data and highlighting its 109 real-world potential. Wang et al. [2024] accelerated Hybrid 110 Electric Vehicle (HEV) energy management by initializing 111 neural networks with expert knowledge, improving learning 112 speed and reducing energy consumption. Similarly, Xu and 113 others [2020] applied warm-start Q-learning to HEVs, reduc-114 ing iterations and enhancing fuel efficiency. Wexler and oth-115 ers [2022] introduced confidence-constrained learning to mit-116 igate performance degradation in Warm-Start RL, balancing 117 policy gradient and constrained learning. While warm starts 118 in DRL and RL accelerate training, they may introduce biased 119 exploration, limiting full-state-action space exploration. 120

Transfer Learning (TL) has been applied to address the 121 cold start problem in DRL. Approaches like the Hierar-122 chical Deep Reinforcement Learning Network (H-DRLN) 123 for Minecraft have shown superior performance by reusing 124 learned skills across tasks [Tessler and others, 2017]. Sim-125 ilarly, Reinforcement Learning Building Optimizer with 126 Transfer Learning (ReLBOT) applied algorithms from data-127 rich buildings to new ones, mitigating cold start in smart 128 buildings [Genkin and McArthur, 2022]. However, TL's suc-129 cess depends on selecting a pre-training task closely aligned 130 with the target task. Mismatched knowledge transfer can re-131 duce performance in the target domain. 132

Imitation Learning (IL), such as Reduction-based Active
 Imitation Learning (RAIL) algorithm, reduces expert queries

using active independently and identically distributed (i.i.d.) 135 learning, showing effectiveness in tasks like Cart-pole [Gros 136 and others, 2020]. Deep Deterministic Policy Gradient with 137 Imitation Learning (DDPG-IL) combines DDPG with imita-138 tion learning for autonomous driving, improving convergence 139 and performance [Zou et al., 2020]. However, it relies heavily 140 on high-quality demonstrations and human expertise, strug-141 gles with complex spaces, and faces ethical concerns and gen-142 eralization challenges, adding complexity to its application. 143

Sample Efficiency (SE) in DRL has advanced signifi-144 cantly. "JueWu-MC" used a hierarchical approach with hu-145 man demonstrations to address partial observability in open-146 world games, outperforming baselines in NeurIPS MineRL 147 competitions [Lin et al., 2021]. Option Machines (OMs) 148 leveraged high-level instructions for action selection, ex-149 celling in single-task, multi-task, and zero-shot learning [den 150 Hengst and others, 2022]. A DRL method for cloud-native 151 Service Function Chains (SFC) caching tackled cold-start 152 problems using Graph Neural Networks (GNNs), optimiz-153 ing latency and request acceptance under high-load condi-154 tions [Zhang and others, 2022]. Ma and others [2022] in-155 troduced a freshness discounted factor in prioritized experi-156 ence replay (PER), improving learning efficiency. Random-157 ized Ensemble Double Q-learning (REDQ) [Chen and oth-158 ers, 2021] achieved state-of-the-art sample efficiency on Mu-159 JoCo [Todorov et al., 2012] using a high update-to-data ra-160 tio and an ensemble of Q-functions. Dropout Q-learning 161 (DroQ) [Hiraoka and others, 2022], a variant of REDQ, im-162 proved computational and memory efficiency while main-163 taining similar sample efficiency. Reinforcement Learning 164 with Policy-Driven Regularization (RLPD) [Ball and others, 165 2023] used symmetric sampling and Layer Normalization to 166 efficiently leverage offline data in online RL, outperforming 167 prior methods without pre-training. Lastly, DRL trained a 168 quadruped robot to walk on varied terrains in 20 minutes by 169 optimizing actor-critic algorithms, achieving real-world sam-170 ple efficiency without simulation [Smith et al., 2023]. 171

Despite advancements, sample efficiency in RL remains 172 a challenge. Methods like REDQ [Chen and others, 2021], 173 DroQ [Hiraoka and others, 2022], and RLPD [Ball and others, 2023] rely on random exploration, leading to poor early 175 performance, especially in real-world applications. They 176 struggle with sparse rewards, complex dynamics, generalization issues, and balancing exploration with exploitation. 178

Our DT-guided DRL technique addresses these issues by leveraging decision theory, using utility functions to guide learning and ensure acceptable early performance. This structured exploration mitigates random exploration and sparse reward challenges, offering a promising direction for improving RL efficiency in data-limited real-world scenarios.

3 Decision Theory-Guided DRL

In DRL, balancing exploration and exploitation is crucial for finding optimal solutions. Traditional methods like dynamic ϵ -greedy exploration, which start with random exploration, often delay convergence. Our DT-guided DRL approach integrates decision theory to provide informed action distribu-190

tions from the start, address the cold start problem¹, and reduce the risk of local optima.

Unlike existing works [Chen and others, 2021; Hiraoka and
others, 2022; Ball and others, 2023], DT-guided DRL provides effective initial guidance, avoiding poor performance
from random exploration. It operates efficiently on small
datasets, enhancing generalizability across diverse RL scenarios.

199 3.1 Problem Formulation Using DRL

We demonstrate our technique using the Cart Pole environment from Gymnasium [Hsiao and others, 2022; Manrique Escobar and others, 2020] and maze problems common in decision theory [Dayan and Daw, 2008], employing the Proximal Policy Optimization (PPO) algorithm [Schulman and others, 2017] for its superior performance in our experiments.

A DRL agent operates within a Markov Decision Process (MDP) as follows:

- **State** (*S*): In a cart pole, the state includes cart position, velocity, pole angle, and angular velocity. For maze problems, it represents the agent's *x* and *y* coordinates.
- Action (*A*): In a cart pole, actions involve pushing the cart left or right; in maze problems, actions include moving up, down, left, or right.
- **Transition Probabilities** (T(s'|s, a)): The probability of transitioning from state s to s' given action a.
- **Reward** (R(s)): In the cart pole, +1 is given for each step the pole remains balanced; in maze problems, small penalties apply for non-exit steps, and +1 is awarded for reaching the exit.
- **Policy** (π): A mapping from states to actions.

222 **3.2** Problem Formulation Using Decision Theory

To guide the DRL agent effectively, the utility function must be tailored to the specific problem or environment.

Cart Pole Environment. In the cart pole problem, the utility function focuses on the pole's angle. Positive utility is assigned for pushing left when the angle is negative and pushing right when positive, aiming to keep the pole upright. Formally, a utility is defined as:

$$U(s_{\text{pole}_angle}, a) = \begin{cases} -\frac{s_{\text{pole}_angle}}{0.209}, & \text{if } a \text{ is push left} \\ \frac{s_{\text{pole}_angle}}{0.209}, & \text{otherwise.} \end{cases}$$
(1)

Here, $s_{\text{pole.angle}}$ represents the pole angle, ranging between (-0.418, 0.418) radians. Since the game ends if the angle exceeds $\pm 12^{\circ}$ (or ± 0.209 radians), the utility is normalized by dividing by 0.209.

Maze Environment. In the maze, the agent can move up, down, left, or right, with its position represented as $(P_{s,x}, P_{s,y})$. The utility function incorporates the distance from the agent's current position to the exit:

$$U(s,a) = \begin{cases} \frac{1}{\operatorname{dis}(P_{s'}, P_{\operatorname{exit}})}, & \text{if no obstacles in } a \text{'s direction} \\ 0, & \text{otherwise,} \end{cases}$$
(2)

where s' is the anticipated new state after action a, and 238 $dis(P_{s'}, P_{exit})$ is the distance to the exit, given by: 238

$$\operatorname{dis}(P_{s'}, P_{\operatorname{exit}}) = \sqrt{(P_{s',x} - P_{\operatorname{exit},x})^2 + (P_{s',y} - P_{\operatorname{exit},y})^2}.$$
(3)

3.3 Integrating DT with DRL for Discrete Actions 240

The effectiveness of DT-guided DRL lies in combining the 241 DT agent's action probabilities with the neural network (NN) 242 output. Then fine-tunes that network via PPO updates; this 243 two-stage loop repeats until convergence. As shown in Figure 1, the process is as follows: 245

- 1. Utility Function: Create a DT agent with a problemspecific utility function (e.g., Eq. (1) for cart pole or Eq. (2) for the maze). Convert these utility values into probabilities using a softmax layer, adjusting the temperature for determinism.
- 2. Deep Neural Network (DNN): Build a DNN that inputs the environment's current state and outputs action probabilities. 253
- 3. **Integration**: Combine the NN's outputs with the DT agent's probabilities (reverse softmax), summing them and adjusting a dynamic weight w, starting at w = 1 and gradually reducing to w = 0 with a linear decay schedule during training.
- 4. **Final Softmax**: Apply a softmax layer (temperature 1) to the combined outputs for a valid action probability distribution. 261
- Training: Train the DNN using RL algorithms like PPO or DQN, updating weights based on rewards from interactions with the environment.

Utility vs. Reward Function. Unlike reward shaping, 265 which modifies the environment's primary reward structure, 266 our utility functions act as an auxiliary heuristic to guide early 267 exploration. Rather than altering the original reward sig-268 nals that define task objectives, the utility biases the agent's 269 action selection during initial training, accelerating learning 270 while preserving the problem's fundamental nature. This 271 distinction ensures the canonical MDP formulation remains 272 intact while effectively mitigating the cold start challenge. 273 Although reward shaping can guide exploration by modi-274 fying reward signals, it alters the original MDP objectives, 275 whereas our approach provides early decision-theoretic guid-276 ance without changing the inherent reward structure. 277

3.4 Integrating DT with DRL for Continuous Actions

Integrating decision theory (DT) with DRL in continuous action spaces poses challenges due to the infinite action set. To address this, we adapt DT-guided DRL by incorporating decision-theoretic guidance and experience replay.

Integrating DT with DRL follows this process: 284

237

278

¹The cold start problem arises from limited initial data, leaving the learning agent with insufficient knowledge to make informed decisions and resulting in poor early performance.



Figure 1: Procedures for generating solutions by a DT-guided DRL agent: S_t is the state at round t, and $Prob_{DT}(a)$ or $Prob_{DT_PPO}(a)$ are the action probability.

- Utility Function: Develop a DT agent with a problemspecific utility function mapping continuous observations to actions. For example, in the inverted pendulum, the utility function computes optimal torque from state variables (e.g., pole angle, cart position) to balance.
- Initial Interaction: The DT agent interacts with the environment for multiple episodes, selecting actions to achieve task objectives (e.g., balancing the pendulum).
- 293 3. Experience Collection: Record state transitions, actions, rewards, and next states from the DT agent's interactions
 294 to capture decision-making patterns informed by the util 296 ity function.
- 4. Replay Buffer Integration: Insert the collected experiences into the DRL agent's replay buffer, enriching the initial training data with informed decisions and improving sample efficiency.
- 5. DRL Model: Build a DRL model that inputs continuous states and outputs continuous actions or action distribution
 parameters (e.g., mean and variance for Gaussian policies in Soft Actor-Critic or SAC).
- 805 6. Neural Network Training: Train DRL by sampling minibatches from the replay buffer (with DT and DRL experiences) and updating policy and value networks.
- 7. Policy Improvement: The DRL agent refines its policy by learning from DT-guided experiences and interactions, allowing for better generalization and exploring more relevant actions.
- 312 By initializing the replay buffer with decision-theory-

driven experiences, the DRL agent benefits from an informed 313 warm start, reducing random exploration and accelerating 314 convergence in continuous action environments. 315

4 Experimental Setup 316

317

320

321

334

4.1 Environment & Parameterization

Our study evaluates DT-guided DRL in both discrete and continuous action environments. 318

Discrete Action Environments

We consider two scenarios as follows.

- **Cart Pole:** The agent balances a pole on a cart by pushing it left or right. The state space includes cart position, velocity, pole angle, and angular velocity. The agent earns +1 per timestep, and the pole stays upright, with episodes ending if the pole tilts beyond $\pm 12^{\circ}$, the cart moves out of bounds or after 100 timesteps. 327
- **Maze:** The agent navigates from the top-left to the bottomright of the maze. The state space represents the agent's x/y position, and the action space includes moving up, down, left, or right. Reaching the exit rewards +1, while each step incurs a small penalty. Episodes end upon reaching the exit or exceeding the step limit. 333

Continuous Action Environment

The Inverted Pendulum problem involves balancing a pole 335 on a cart by applying continuous forces [Tarkhov and others, 2023]. The state space includes cart position, velocity, 337 Table 1: NEURAL NETWORK CONFIGURATION PARAMETERS

Parameter	Value (PPO)	Value (SAC)
Discount Factor (γ)	0.99	0.99
Learning Rate	0.0003	0.0003
Replay Buffer Size	2048	20000
Batch Size	64	256
PPO's Clipping Parameter	0.2	NA
SAC's Target Smoothing Coeffi-	NA	0.005
cient (τ)		
SAC's Temperature Parameter (α)	NA	0.036
Network Architecture	64 x 64	256 x 256
Activation Function	Tanh	ReLU

pole angle, and angular velocity. The action space is the continuous force applied to the cart. Table 1 summarizes key
configurations, including discount factor, learning rate, replay buffer size, batch size, and network architecture. These
settings are based on Stable-Baselines3 defaults [Raffin and
others, 2021], ensuring rigorous and reproducible results.

344 4.2 Comparing Schemes

345 Discrete Action Environments

- ³⁴⁶ Under this environment, an agent takes an action as follows:
- Decision Theory (DT) Agent [Fishburn *et al.*, 1979]: Uses
 DT-based utility functions for action selection.
- Standard Proximal Policy Optimization (PPO) [Schulman and others, 2017]: Uses the PPO algorithm without decision-theoretic guidance.
- **Transfer Learning (TL) PPO** [Genkin and McArthur, 2022]: Starts training on a 3x3 maze and transfers learning to larger mazes to evaluate TL effectiveness.
- Sample Efficiency (SE) PPO [Ball and others, 2023]: Enhances sample efficiency through techniques like an increased replay buffer and freshness-prioritized experience replay.
- **Imitation Learning (IL) PPO** [Gros and others, 2020]: An agent is pre-trained with expert demonstrations before continuing standard PPO training.
- **DT-guided PPO (Ours)**: Integrates DT into the action selection process during training. As outlined in Section 3, the method integrates the neural network's action probabilities with those from the utility function, employing a dy-
- namic weight that gradually shifts from DT guidance to anadaptive policy.

368 Continuous Action Environment

- ³⁶⁹ Under this environment, an agent behaves as follows:
- **Standard Soft Actor-Critic (SAC)** Haarnoja and others [2018]: Uses the SAC algorithm without decision-theoretic guidance.
- **DT-guided SAC** (Ours): Initializes the replay buffer with experiences collected from the DT agent.
- The source code is accessible at https://github.com/ Wan-ZL/DT-DRL.



Figure 2: Accumulated rewards under DT, PPO, SE PPO, IL PPO, and DT-guided PPO over 500 training episodes in the Cart Pole problem.



Figure 3: Accumulated rewards under DT, PPO, TL PPO, SE PPO, and DT-guided PPO over 500 training episodes in the Maze problem.

5 Simulation Results & Analysis

5.1 Cart Pole Problem

Comparative Performance Analysis

Figure 2 presents the accumulated rewards for five agents: 380 DT, PPO, SE PPO, IL PPO, and DT-guided PPO in the Cart Pole simulation. 382

Our key observations and findings are as follows. First, the 383 DT-guided PPO (red curve) outperforms other PPO agents 384 early in training, starting with a higher reward. This advan-385 tage comes from integrating DT utility values with the neu-386 ral network's output. Since the neural network begins with 387 zero-initialized outputs and uses Tanh as the activation func-388 tion, the DT utility dominates the initial decisions, making 389 the early behavior of the DT-guided PPO similar to the DT 390 agent (blue curve). Second, the DT agent's reward remains 391 static, showing its limitation in adapting beyond its prede-392 fined utility function. In contrast, the DT-guided PPO steadily 393 improves, reflecting the neural network's increasing ability 394 to learn from the environment and refine its actions, result-395 ing in progressively higher rewards. Finally, while all PPO 396 agents converge to similar reward levels, the DT-guided PPO 397 converges faster, demonstrating accelerated learning. This is 398 due to the DT-guided PPO's informed starting point, which 399 reduces the need for exploration and speeds up the learning 400

377 378

process. This early performance boost is crucial in real-world
applications where each failed attempt may be costly or dangerous. In particular, our data shows that DT-guided DRL
reaches near-optimal behavior significantly faster, demonstrating improved sample efficiency in the most critical initial
phase of training.

407 Overall, Figure 2 confirms that integrating decision the408 ory into the PPO framework mitigates random exploration
409 early in training and accelerates convergence to optimal per410 formance in the cart pole task.

411 5.2 Maze Problem

412 **Comparative Performance Analysis**

Figure 3 shows the accumulated rewards for five agents, including DT, PPO, TL PPO, SE PPO, and DT-guided PPO, in
the Maze simulation.

The results reveal several key insights: (1) DT-guided 416 PPO consistently outperforms all other agents, demonstrating 417 the effectiveness of combining decision theory's structured, 418 utility-based guidance with PPO's flexible neural learning. 419 This integration enables more informed decision-making, es-420 pecially in large mazes. (2) The maze environment poses 421 a sparse reward challenge, where rewards are only received 422 upon exit. DT-guided PPO mitigates this by using decision-423 theoretic guidance for early exploration, reducing inefficient 424 425 random search, and enabling more focused navigation. (3) DT-guided PPO converges faster than other agents. Its initial 426 utility-driven exploration accelerates learning by reducing the 427 time to discover optimal strategies, allowing it to reach near-428 optimal performance in fewer episodes. 429

430 Sensitivity Analysis – Effect of Maze Size (*m*)

Figure 4 illustrates the dynamics of accumulated rewards as a function of maze size (m), ranging from 3 to 8. Each subplot delineates the trajectory of accumulated reward over training episodes, with the x-axis representing the episode count and the y-axis quantifying the reward.

Upon examination, we discern several patterns as follows. 436 (1) The DT-guided PPO agent consistently surpasses all other 437 agents (DT, PPO, TL PPO, and SE PPO) regarding accumu-438 lated rewards across all maze sizes. This superior perfor-439 mance suggests that DT-guided PPO leverages the systematic 440 approach of decision theory and the flexible learning capabil-441 ities inherent in the PPO's neural network structure. Such an 442 443 integration gives the agent a robust navigational strategy in the maze, which becomes increasingly advantageous as the 444 maze's complexity escalates. This is crucial in larger mazes 445 where the agent must contend with more intricate challenges. 446 (2) As the maze size escalates, the differential in performance 447 between the DT-guided PPO and other PPO agents becomes 448 more pronounced, especially in later episodes. This trend 449 could be attributed to the DT component providing a more ef-450 fective heuristic in the early stages of exploration, guiding the 451 agent through larger state spaces more efficiently. The DT's 452 structured approach potentially reduces the agent's search 453 space by giving zero utility to actions toward obstacles. Thus, 454 it mitigates the challenges posed by a larger maze's complex-455 ity and helps to maintain higher performance levels than the 456 pure PPO agent. (3) A general trend observed is the decline 457

in the accumulated reward for all agents with increasing maze 458 size, with the DT agent demonstrating considerable fluctua-459 tion when $m \geq 7$. This fluctuation could stem from the am-460 plified complexity and inherent difficulties of larger mazes, 461 where a sole reliance on decision theory might not yield effi-462 cient pathfinding consistently. Without the capacity to learn 463 and adapt from interactions with the environment, a purely 464 decision-theoretic model may fail. Conversely, a pure PPO 465 approach may struggle with sparse rewards, as exits become 466 harder to reach and positive rewards become less frequent as 467 the maze expands. 468

Despite these challenges, DT-guided PPO resists the sparse reward problem, likely due to initial guidance from the decision theory component directing the agent toward more rewarding trajectories in larger mazes. The hybrid model's integration of structured decision-making with adaptive learning effectively manages maze complexities across scales. 469

The DT-guided PPO's consistent outperformance across maze sizes illustrates the value of combining structured decision-making with empirical adaptive neural networks, especially when dealing with problems of increasing size and complexity. The data suggests this hybrid approach could be a promising direction for developing robust solutions in complex, dynamic environments.

5.3 Inverted Pendulum Problem

Comparative Performance Analysis

Figure 5 shows the accumulated rewards for the SAC and DTguided SAC agents in the Inverted Pendulum simulation.

482

483

493

494

509

The key observations are: (1) Both agents display low rewards (around 4 to 6) during the initial 150 episodes, reflecting the random exploration phase. (2) The DT-guided SAC (red curve) demonstrates an early advantage over the SAC agent (blue curve), as the decision-theoretic function provides informed action guidance, enabling more effective exploration compared to SAC's purely random approach.

Sensitivity Analysis – Impact of Decision-Theoretic Initialization on SAC Performance

Figure 6 compares the performance of SAC* and DT-guided 495 SAC* variants under different initialization strategies. In 496 standard SAC, a fixed number of random steps is used initially to promote exploration. SAC* removes this random 498 initialization, while DT-guided SAC* replaces it with 1 or 2 episodes of interactions from a decision-theoretic (DT) agent, 500 providing structured guidance early in training. 501

Key observations are: (1) DT-guided SAC* consistently 502 outperforms SAC*, achieving higher accumulated rewards. 503 The initial DT-driven interactions guide the agent away from 504 inefficient exploration and toward promising actions. (2) 505 Both the 1-episode and 2-episode variants perform similarly, 506 suggesting that the DT agent provides consistent guidance 507 and that additional episodes yield diminishing returns. 508

6 Conclusion & Future Work

We proposed Decision Theory-guided Deep Reinforcement 510 Learning (DT-guided DRL) to address the cold start problem 511 by integrating decision-theoretic principles to improve initial performance and accelerate convergence. Experiments on 513



Figure 4: Effect of maze size (m): Accumulated rewards under DT, PPO, TL PPO, SE PPO, and DT-guided PPO over 500 training episodes.



Figure 5: Accumulated reward of SAC and DT-guided SAC over 500 training episodes in the Inverted Pendulum problem.

cart pole, maze, and inverted pendulum tasks show that DT-514 guided DRL consistently outperforms conventional agents, 515 achieving higher early-stage rewards and faster learning. 516

In the cart pole and maze tasks, DT-guided DRL demon-517 strated strong initial performance and robustness, particu-518 larly under sparse rewards and large-scale navigation. In the 519 continuous-action inverted pendulum task, it enhanced early 520 learning by leveraging decision-theoretic interactions. 521

While effective in providing structured exploration and im-522 proved early learning, the approach depends on problem-523 specific utility functions, which may not generalize easily. 524 Although tested on simplified benchmarks, the results clearly 525 demonstrate DT-guided DRL's value in mitigating the cold 526 start problem. Future work will explore scaling to complex 527 domains, such as multi-joint robots, legged locomotion, and 528



Figure 6: Accumulated reward of SAC* and DT-guided SAC* with one or two episodes of decision theory agent interactions during the training in the Inverted Pendulum problem.

vision-based tasks.

We also plan to investigate automated ways of deriving 530 utility functions through inverse RL or meta-learning, reduc-531 ing the reliance on manual heuristics. Combining DT-guided 532 exploration with intrinsic-motivation approaches [Aubret et al., 2023] may further enhance sample efficiency. Evaluating DT-guided DRL in large-scale, high-dimensional tasks and 535 in safety-critical real-world robotics platforms forms the next 536 crucial step in validating its broader efficacy. 537

In conclusion, DT-guided DRL represents a significant step 538 towards more efficient and reliable reinforcement learning. 539 We hope this work inspires further research at the intersection 540 of decision theory and deep reinforcement learning. 541

529

542 **References**

- Arthur Aubret, Laetitia Matignon, and Salima Hassas. An
 information-theoretic perspective on intrinsic motivation
- in reinforcement learning: A survey. *Entropy*, 25(2):327,2023.
- ⁵⁴⁷ Philip J Ball et al. Efficient online reinforcement learning
- with offline data. In *International Conference on Machine Learning*, pages 1577–1594. PMLR, 2023.
- Xinyue Chen et al. Randomized ensembled double q learning: Learning fast without a model. *International Conference on Learning Representations*, 2021.
- Peter Dayan and Nathaniel D Daw. Decision theory, rein forcement learning, and the brain. *Cognitive, Affective, & Behavioral Neuroscience*, 8(4):429–453, 2008.
- 556 Floris den Hengst et al. Reinforcement learning with op-
- tion machines. In Proceedings of the Thirty-First Interna-
- tional Joint Conference on Artificial Intelligence, IJCAI-
- 22, pages 2909–2915. International Joint Conferences on
 Artificial Intelligence Organization, 2022.
- Gabriel Dulac-Arnold, Daniel Mankowitz, and Todd Hester.
 Challenges of real-world reinforcement learning. *ICML Workshop on Real-Life Reinforcement Learning*, 2019.
- Peter C Fishburn, Peter C Fishburn, et al. Utility theory for
 decision making. Krieger NY, 1979.
- 566 Mikhail Genkin and JJ McArthur. Using reinforcement 567 learning with transfer learning to overcome smart building
- cold start. In 2022 International Conference on Computational Science and Computational Intelligence (CSCI),
 pages 713–718. IEEE, 2022.
- Timo P Gros et al. Tracking the race between deep reinforcement learning and imitation learning. In *Quantitative Eval*
- 573 uation of Systems: 17th International Conference, QEST
- 574 2020, Vienna, Austria, August 31–September 3, 2020, Pro-
- *ceedings 17*, pages 11–17. Springer, 2020.
- Tuomas Haarnoja et al. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic
- actor. In *International Conference on Machine Learning*,
 pages 1861–1870. Pmlr, 2018.
- Takuya Hiraoka et al. Dropout q-functions for doubly effi cient reinforcement learning. *International Conference on Learning Representations*, 2022.
- Jen-Yueh Hsiao et al. Unentangled quantum reinforce ment learning agents in the OpenAI gym. *arXiv preprint arXiv:2203.14348*, 2022.
- ⁵⁸⁶ Zichuan Lin, Junyou Li, Jianing Shi, Deheng Ye, Qiang ⁵⁸⁷ Fu, and Wei Yang. Juewu-mc: Playing minecraft with
- Fu, and Wei Yang. Juewu-mc: Playing minecraft with
 sample-efficient hierarchical reinforcement learning. *arXiv preprint arXiv:2112.04907*, 2021.
- Jue Ma et al. Fresher experience plays a more important
 role in prioritized experience replay. *Applied Sciences*,
 12(23):12489, 2022.
- Camilo Andres Manrique Escobar et al. A parametric study
 of a deep reinforcement learning control system applied to

the swing-up problem of the cart-pole. *Applied Sciences*, 595 10(24):9013, 2020. 596

- D Warner North. A tutorial introduction to decision theory. 597 IEEE Transactions on Systems Science and Cybernetics, 598 4(3):200–210, 1968. 598
- Roxana Rădulescu et al. Multi-objective multi-agent decision making: a utility-based analysis and survey. *Autonomous Agents and Multi-Agent Systems*, 34(1):10, 2020.
- Antonin Raffin et al. Stable-baselines3: Reliable reinforcement learning implementations. *Journal of Machine Learning Research*, 22(268):1–8, 2021. 605
- John Schulman et al. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017. 607
- Andrew Silva and Matthew Gombolay. Encoding human domain knowledge to warm start reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 5042–5050, 2021.
- Laura Smith, Ilya Kostrikov, and Sergey Levine. Demonstrating a walk in the park: Learning to walk in 20 minutes with model-free reinforcement learning. *Robotics: Science and Systems (RSS) Demo*, 2(3):4, 2023.
- DA Tarkhov et al. Optimal control selection for stabilizing the inverted pendulum problem using neural network method. *Optical Memory and Neural Networks*, 32(Suppl 2):S214– S225, 2023.
- Chen Tessler et al. A deep hierarchical approach to lifelong learning in minecraft. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 31, 2017. 622
- Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In 2012 IEEE/RSJ International Conference on Intelligent Robots and Systems, pages 5026–5033. IEEE, 2012. 626
- Hanchen Wang, Ziba Arjmandzadeh, Yiming Ye, Jiangfeng Zhang, and Bin Xu. Flexnet: A warm start method for deep reinforcement learning in hybrid electric vehicle energy management applications. *Energy*, 288:129773, 2024. 630
- Benjamin Wexler et al. Analyzing and overcoming degradation in warm-start reinforcement learning. In 2022 *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 4048–4055. IEEE, 2022. 634
- Bin Xu et al. Learning time reduction using warm-start methods for a reinforcement learning-based supervisory control in hybrid electric vehicle applications. *IEEE Transactions on Transportation Electrification*, 7(2):626–635, 2020.
- Jiayin Zhang et al. Cold-start aware cloud-native service 639 function chain caching in resource-constrained edge: A reinforcement learning approach. *Computer Communications*, 195:334–345, 2022. 642
- Zhuangdi Zhu et al. Transfer learning in deep reinforcement
 learning: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- Qijie Zou, Kang Xiong, and Yingli Hou. An end-to-end learning of driving strategies based on ddpg and imitation learning. In 2020 Chinese Control and Decision Conference (CCDC), pages 3190–3195. IEEE, 2020. 649