A GRID WORLD AGENT WITH FAVORABLE INDUCTIVE BIASES

Anonymous authors

003

010 011

012

013

014

015

016

017

018

019

Paper under double-blind review

Abstract

We present a novel experiential learning agent with causally-informed intrinsic reward that is capable of learning sequential and causal dependencies in a robust and data-efficient way within grid world environments. After reflecting on stateof-the-art Deep Reinforcement Learning algorithms, we provide a relevant discussion of common techniques as well as our own systematic comparison within multiple grid world environments. Additionally, we investigate the conditions and mechanisms leading to data-efficient learning and analyze relevant inductive biases that our agent utilizes to effectively learn causal knowledge and to plan for rewarding future states of greatest expected return.

021 1 INTRODUCTION

Grid world environments come in many forms and have been studied extensively in the history of Artificial Intelligence with some notable examples such as Wumpus World (Bryce, 2011), Minigrid (Chevalier-Boisvert et al., 2024), and Tileworld (Pollack & Ringuette, 1990). However, the creation of intelligent grid world agents capable of learning effectively and in a data-efficient way has posed significant challenges. Current Reinforcement Learning (RL) agents struggle in some instances due to sequential dependencies, partial observability (Wang et al., 2023a), continual learning (primacy bias Kim et al. (2024), stability-plasticity dilemma Anand & Precup (2024)), relatively high-dimensional state spaces, compared to more traditional RL tasks – even when a single value per cell is provided rather than per pixel – and sometimes non-deterministic effects of actions.

Sequential dependencies usually raise the training data demand exponentially depending on the 032 combinatorics, and, ultimately, the number of the arising options, unless the agent is capable of 033 learning causal representations that transfer well because chaining them is favorable for reaching an 034 intended outcome. Partial observability, on the other hand, may require the model to have access to information from prior states, which typically corresponds to the previously observed values in a grid world outside the agent's current field of view. In terms of input dimensionality, there is 037 a trade-off between the observation window being too small for learning an effective policy and 038 the agent's observation window being more high-dimensional, thereby demanding more training data, even though this might be the least problematic. Additionally, if the agent is able to represent values outside of its observation window, a learned policy needs to consider not only the observation 040 window itself but also how it spatially relates to the remembered information beyond it. 041

With these considerations in mind, we introduce Non-Axiomatic Causal Explorer (NACE), a novel
 experiential learning agent, which leverages causal reasoning and intrinsic reward signals to enable
 more efficient learning as well as possesses learning mechanisms with the involved inductive biases.
 NACE is designed to induce causal rules from temporal and spatial local changes in the grid, which
 are often (but not always) caused by the agent. It utilizes these rules to plan for and reach future states
 of maximum uncertainty in order to effectively learn more (causal rules) about the environment,
 thereby improving predictability-based intrinsic reward formulations.

To illustrate the effectiveness of our approach, we provide a comprehensive discussion of state-of the-art Deep Reinforcement Learning (DRL) techniques as well as our own systematic comparison
 within multiple grid world environments demonstrating our agent's remarkable improvement in data
 efficiency, achieving similar performance with about 1000 samples where DRL algorithms typically
 require 1 million samples, representing a 1000-fold reduction in data requirements. We also thoroughly investigate the conditions and mechanisms under which learning in grid world environments

can become more data-efficient and attempt to answer the question about which inductive biases
 can lead to close-to-optimal learning speeds in the case where the agent is not pre-equipped with
 particular interaction rules between grid cell types but has the inductive biases to build and use them
 effectively. Lastly, we analyze useful inductive biases applicable among a wide range of grid worlds
 and their generality to raise the data efficiency of learning in such domains.

060

2 RELATED WORK

Current RL techniques such as value-based (e.g., Deep Q-Learning (Mnih et al., 2013)) and policy gradient-based (e.g., Proximal Policy Optimization (Schulman et al., 2017)) typically require millions of training iterations to solve grid-world environments (Zhang et al., 2020b), struggling to capture causal dependencies necessary for efficient planning and transferability.

Exploration improvements, such as intrinsic rewards based on information gain (Zhao et al., 2023), 067 prediction errors (Burda et al., 2018), or visitation counts (Zheng et al., 2021; Wang et al., 2023b), 068 enhance sample efficiency but often lack structured reasoning for generalization. In contrast, Tsi-069 vidis et al. (2021) propose Theory-Based RL (TBRL), exemplified by EMPA (Exploring, Modeling, and Planning Agent), which integrates Bayesian causal modeling, structured exploration, and heuris-071 tic planning to generalize efficiently across tasks with minimal training. Similarly, GALOIS (Cao 072 et al., 2022) addresses generalization by synthesizing interpretable, hierarchical programs with strict 073 cause-effect logic, though its reliance on predefined program sketches limits flexibility in loosely 074 structured environments. As a model-based RL alternative, DreamerV3 (Hafner et al., 2023) learns 075 latent state dynamics and improves behavior through imagination, enabling generalization across 076 diverse tasks with minimal domain-specific adjustments. However, its reliance on learning both la-077 tent representations and their dynamics reduces sample efficiency compared to methods that assume predefined representations. 078

Symbolic approaches like STRIPS and Behavior Trees (BTs) (Guo et al., 2023; Colledanchise & Ögren, 2018) handle human-defined causal knowledge but lack adaptive learning capabilities. Similarly, POMDPs (Spaan, 2012) focus on probability updates rather than causal discovery, while causal networks (Pearl, 1995) and structure-learning methods (Zheng et al., 2018) face challenges with ambiguity and scalability.

While TBRL, GALOIS, and DreamerV3 provide solutions for structured or model-based learning, their reliance on predefined logic, priors, or latent space learning introduces additional complexity.
Our approach proposes lightweight "empirical" causal relations learned from recurring cause-effect patterns, supporting efficient real-time learning in grid worlds without predefined program sketches, Bayesian modeling, or latent space learning. This ensures adaptability while retaining interpretability and transferability.

090 091

091 2.1 SELECTED RL TECHNIQUES IN GRID WORLDS

DRL often struggles with sample efficiency, requiring substantial interactions with environments.
 This paper examines foundational algorithms, scalable architectures, and exploration-focused methods that address these challenges.

096 Foundational methods include Deep Q-Network (DQN) Mnih et al. (2015), which combines deep 097 neural networks with Q-learning to handle large state spaces but struggles with sparse rewards; 098 Advantage Actor-Critic (A2C) Mnih et al. (2016), which reduces variance in updates through synchronized parallel actors but is limited by its on-policy nature; Trust Region Policy Optimization 099 (TRPO) Schulman et al. (2015), which ensures stable policy updates with trust region constraints 100 but is computationally intensive; and Proximal Policy Optimization (PPO) Schulman et al. (2017), 101 which refines TRPO with clipped objectives for improved data utilization and computational effi-102 ciency. 103

Scalable architectures such as Importance Weighted Actor-Learner Architectures (IMPALA) Espeholt et al. (2018) address multi-task learning by leveraging distributed architectures with off-policy corrections, offering scalability but facing synchronization challenges, whereby exploration-focused methods aim to address sparse rewards and complex state spaces. Count-Based Exploration (COUNT) Bellemare et al. (2016) uses pseudo-counts for better exploration but is computationally

demanding in large state spaces. Random Network Distillation (RND) Burda et al. (2018) incentivizes novelty through prediction errors, however, it depends on high-quality state representations. Curiosity-Driven Exploration (CURIOSITY) Pathak et al. (2017) rewards prediction errors of action outcomes, fostering intrinsic motivation, while Rewarding Impact-Driven Exploration (RIDE)
Raileanu & Rocktäschel (2020) focuses on impactful actions but may struggle with ambiguous state changes. Adversarially Motivated Intrinsic Goals (AMIGO) Campero et al. (2021) generates adversarial goals to guide exploration, requiring robust goal-generation mechanisms for effectiveness.

Furthermore, model-based RL techniques can learn a model of the environment and improve their behavior through imagined future scenarios Sekar et al. (2020). This framework is particularly ad-vantageous in high-dimensional and variable state spaces. DreamerV3 Hafner et al. (2023) builds on this principle by introducing a general algorithm designed to address a diverse range of tasks with minimal domain-specific adjustments. By learning latent state dynamics and planning through imagination, DreamerV3 offers broad applicability while reducing the need for extensive tuning or specialized configurations. However, the additional complexity of learning latent space representa-tions together with their dynamics leads to less sample efficiency compared to when representations are already present and only the dynamics need to be learned.

3 NON-AXIOMATIC CAUSAL EXPLORER

NACE is our proposed experiential learning technique with causality-informed intrinsic reward and
 strong inductive biases for grid world environments to boost sample efficiency. Here, we provide
 formal descriptions of NACE. For a comprehensive list of symbols used, refer to Appendix C.

131 3.1 STATES AND RULE REPRESENTATION

State in NACE is a tuple $s = (s_{\text{spatial}}, s_{\text{internal}})$ consisting of a two-dimensional value array $s_{\text{spatial}} \in \mathbb{N}^{m \times n}$ and a one-dimensional value array $s_{\text{internal}} \in \mathbb{N}^k$, as shown in Figure 1. The two-dimensional array reflects the spatial structure in the grid world, including remembered cells beyond the current view, while the one-dimensional array is used for internal values, such as inventory items (e.g., keys).

2D Spatial Array

	0	0		Key
	0	х		Ball
Mental map Inventory		iventory		

1D Internal Array

Figure 1: State components

Each rule is of the form (*preconditions*, *action*) \Rightarrow *consequence* where the precondition can hold a conjunction of cell value constraints spatially relative to the cell value of the consequence, and the consequence predicts one particular cell's value as well as the values of the one-dimensional array at the next timestep as depicted in Figure 2. Examples of created rules are provided in Appendix G.



Figure 2: Rule schema

Each rule tracks evidence using counters for w_+ and w_- similar to (Wang, 2013), which measure the accuracy of the rule's predictions. Positive evidence (w_+) is accumulated whenever a perfectly matching rule predicts correctly, while negative evidence (w_-) increases with incorrect predictions. Tracking of evidence helps the agent refine its causal knowledge by prioritizing more reliable rules.

167 3.2 INDUCTIVE BIASES

It is well-known that favorable inductive biases can enhance sample efficiency. Below are inductive biases that are incorporated in NACE and relevant for many grid world environments:

- 1. **Temporal Locality**: NACE constructs rules based solely on the current and previous state, modeling relevant dependencies locally in time.
- 2. **Causal Representation**: NACE's knowledge representation is centered around the aforementioned causal rules which can be chained and are independent of the objective.
- 3. **Spatial Equivariance**: Ability to model causal dependencies between grid cells independently of the specific location of the cells considered in the dependency. This means learned rules in NACE can be applied at any location.
- 4. **State Tracking**: Ability to effectively track state outside of the field of view of the agent based on the recorded or estimated locations. NACE explicitly keeps track of a bird's-eye view map by recording observations into it, updating the values that are within its observability window.
 - 5. Attentional Bias: Relevant dependencies tend to involve values that have either observably changed or a different value than predicted. Only rules that show a change from the previous to the current timestep, or differ from the predicted value, are considered for rule formation, evidence updating, and prediction.
- Additional discussions on inductive biases as well as ablation studies can be found in Appendix B.
- 189 190 191

192

200

201 202

203 204

205

206

207

208

210 211

166

171 172

173 174

175

176

177

178

179

180

181

182

183

186

187 188

3.3 CURIOSITY MODEL

This section outlines the mechanism which helps NACE systematically acquire missing causal knowledge about the environment. The key principle is realized by making the agent plan to reach a state which it is most unfamiliar with. The familiarity is judged by whether existing rules match well to the situation, whereby matching is a matter of degree dependent on how many rule conditions match the cells in the known state. This motivates the following formalism:

• *Match* value of a rule *r* is evaluated relative to consequence cell *c*:

$$M(r,c) = \frac{\text{Number of matched preconditions}}{\text{Total number of preconditions}}$$

• Cell match value of a cell c dependent on all m existing rule match values:

$$C(c) = \max(0, M(r_1, c), \dots, M(r_m, c))$$

• State match value S(s) of a state s is the average C(c) of its cells with C(c) > 0. This value, as we will see, is the secondary "explorative" objective in the planning process that guides the agent's decisions:

$$S(s) = \sum_{c \in X(s)} \frac{C(c)}{|X(s)|} \text{ where } X(s) = \{c \in s | C(c) > 0\},$$

212 3.4 NACE ARCHITECTURE

213

Figure 3 illustrates the high-level architecture of NACE, which consists of several interconnected components that work together to enable learning and decision-making. The related pseudocode is provided in Appendix D.



$$w_{-}(r) = \begin{cases} w_{-}(r) + 1 & \text{if } c_t \in M_{\text{mismatched}_t}^{\text{prediction}} \\ w_{-}(r) & \text{otherwise} \end{cases}$$

Finally, rules r for which $w_{-}(r) > w_{+}(r)$ become inactive, and for two rules r_{1}, r_{2} if their preconditions match (including the action) but the postconditions are different, only the rule with the higher truth expectation is selected, which is calculated according to:

$$w(r) = w_{-}(r) + w_{+}(r), \quad frequency(r) = \frac{w_{+}(r)}{w(r)}, \quad confidence(r) = \frac{w(r)}{w(r)+1}$$
$$f_{exp}(r) = (frequency(r) - \frac{1}{2}) * confidence(r) + \frac{1}{2}$$

This not only allows the system to find the relevant preconditions under which a consequence happens when the action is utilized but also gives the system tolerance to nondeterministic effects and enables accounting for uncertainty. A brief analysis of this can be found in Appendix A.

3. **Planner**: NACE makes use of depth- and width-bounded Breadth-First-Search algorithm with a combined search objective consisting of two components: it searches for states resulting from the different action sequences for futures that lead to the max. expected return or, if not existing, the lowest state match value. Hence, it applies a key RL principle to maximize the expected long-term return (Sutton et al., 1999), with the policy determined by the considered action sequence: $\pi(t) = a_t$ for t = 1, 2, ..., n whereby n is smaller-or-equal (dependent on where the optimum is found) to the maximum planning horizon:

$$\pi(t) = \begin{cases} \arg \max_{\pi} V^{\pi}(s_0) & \text{if } V^{\pi}(s_0) = \mathbb{E}\left[\sum_{t=0}^n \gamma^t R(s_t) \mid s_0 = s, \pi\right] > 0\\ \arg \min_{\pi} S(s_n) < 1 & \text{otherwise} \end{cases}$$

According to this definition, if no return greater than zero can be obtained for any considered action sequence, the system instead plans for a future state of lowest state match value, whereby $(\forall t : (0 \le t < n) \rightarrow S(s_t) = 1) \land S(s_n) < 1$, meaning the action sequence is constrained to be planned in such a way that state match value is 1 except for the last action where it is minimized for the resulting state.

Such constraint maximizes the agent's chance to reach the state of minimum state match while ensuring the low match value is not a consequence of predicting further from states where the knowledge was already not fully applicable.

Due to the amount of possible options, the planning algorithm dominates the asymptotics of NACE. It has the computational complexity of O(|V|+|E|) where V is the set of nodes, and E is the set of edges of the search graph. Constant-bounded search depth and width can be achieved by pruning of branches by expected return and state match value, however, bounded search depth can negatively affect performance, as analyzed in Appendix A.

4. **Predictor**: When the planner queries for the predicted state from a given state and an action, the role of the predictor is to construct the predicted state by applying all knowledge to the given state in the following way: initializing with the cell values from the given state, where for each cell we utilize only the rule r with M(r,c) = 1 and maximum $f_{exp}(r)$, meaning the rule preconditions match perfectly to the given state, the action that has been considered, and r has the highest truth expectation among the rule candidates.

In this case the postcondition cell value of the rule is applied to the corresponding cell at position (x, y) in the predicted state, while else the cell keeps the value from the previous state. Hence, for utilized rules $r^* = ((\overline{c}_t^1 \wedge ... \wedge \overline{c}_t^k \wedge \overline{v}_t \wedge \overline{a}_t) \Rightarrow (\overline{c}_{t+1} \wedge \overline{v}_{t+1} \wedge R(r^*)))$, where \overline{c}_{t+1} and \overline{v}_{t+1} constrains the cell value and value array of the consequence:

$$c_{t+1,x,y} = \begin{cases} c_{t+1} & \text{if } r^* = \underset{r|M(r,c_{t,x,y})=1}{\operatorname{argmax}} f_{\exp}(r) \\ c_{t,x,y} & \text{otherwise} \end{cases}$$

Now, while s_{t+1} is a composition of the cells at all locations at time t + 1, the reward associated with s_{t+1} is the average of the reward of each of the N utilized rules:

$$R(s_{t+1}) = \frac{1}{N} \sum_{i=1}^{N} R(r_i^*)$$

Hereby, the average was chosen since if the reward prediction of all the used rules aligns with the observed reward, their average will also align, while the sum would overestimate the outcome.

³²⁴ 4 EXPERIMENTS IN MINIGRID

To evaluate the effectiveness of NACE compared to other DRL techniques, we conducted a series of experiments in Minigrid Chevalier-Boisvert et al. (2024), a 2D grid world environment featuring diverse and procedurally generated scenarios (Hardware setup and further test environments in Appendices F and H). We focus on Minigrid levels that feature partial observability (using the default observation format, which provides values per grid cell rather than per pixel), challenging the agent to operate with limited information about its surroundings. The selected environments (see Table 1) are categorized based on the specific challenges they present:

- 1. Static: fixed start & goal locations.
- 2. **Dynamic:** randomized start, goal, and obstacle positions.
 - 3. Dynamic with sequential dependencies: tasks requiring specific action sequences (e.g., a door that needs a key or switch to be opened).

Environment	Туре
MiniGrid-Empty-16x16-v0	1
MiniGrid-DistShift2-v0	1
MiniGrid-LavaGapS7-v0	2
MiniGrid-SimpleCrossingS11N5-v	0 2
MiniGrid-Unlock-v0	3
MiniGrid-DoorKey-8x8-v0	3

Table 1: Environments with corresponding types

In each environment, we recorded the average reward, episode length, and standard deviation every 100 timesteps, whereby each timestep incorporates the observed state, action taken, and obtained reward. The following sections present and discuss some representative results for each category, using the selected RL techniques mentioned in Section 2.1. (Configuration and hyperparameter details are in Appendix E). Additionally, Behavior Trees (BTs) and hard-coded policies are employed as performance upper bounds in non-stationary and static environments.

4.1 STATIONARY ENVIRONMENTS

In this category, because the start and goal locations are fixed, the primary challenge for the agent is to consistently learn and optimize navigation strategies over repeated episodes.



MiniGrid-Empty-16x16-v0

Table 2: Final values for MiniGrid-Empty-16x16-v0

^{MiniGrid-Empty-16x16-v0: This environment features a large, static grid where the agent must navigate from a fixed start location to a fixed goal. Due to the limited observation window and the sparse reward—only granted upon reaching the goal—the task can present some difficulties. In this case, NACE, following its innate strategy, first learns to move effectively within the grid by exploring its immediate surroundings. Then it systematically expands its exploration, hitting observed walls out of curiosity, and finally exploring initially invisible parts of the map until the goal object is found and moved into. Due to its intrinsic inductive bias, it explores the area systematically and associates reward with the goal location within about 10³ timesteps. In contrast, other techniques}

like DreamerV3, TRPO, IMPALA, RIDE, and AMIGO, although capable of eventually learning the task, require over 10^5 timesteps to perform comparatively well (as seen in Figure 4, and Table 2).

MiniGrid-DistShift2-v0: In this case, the fixed start and goal locations are accompanied by stationary lava obstacles, which the agent must navigate around to reach the goal. DQN and DreamerV3 perform quite well, achieving a near-optimal policy with an average reward of 0.96, closely mirroring the performance of the BT. Although reaching a slightly lower average reward of 0.87, NACE was three orders of magnitudes more sample-efficient. The next-best policies were found by AMIGO and PPO with an average reward of 0.78 and 0.76, while the other techniques were below 0.5, all of them being much less sample-efficient than NACE (as in Figure 5 and Table 3).



MiniGrid-DistShift2-v0

MiniGrid-DistShift2-v0

4.2 DYNAMIC ENVIRONMENTS

Given that the start and goal locations, along with obstacle positions, are randomized in each episode, these environments require the agent to continuously adapt to new and unpredictable conditions.



MiniGrid-LavaGapS7-v0: In this environment, the agent must navigate around randomly placed 426 lava obstacles to reach a fixed goal, requiring adaptability due to the varying paths between episodes. 427 The 5x5 free space, which is mostly covered by the agent's observation window, is complicated 428 by dynamically spawning lava, unlike the stationary obstacles in *MiniGrid-DistShift2-v0*. From the 429 mean episode rewards (as seen in Figure 6, and Table 4), it is clear that PPO and NACE find a similar 430 effective strategy, whereby NACE takes about 10^3 timesteps while PPO takes 3×10^5 timesteps to 431 reach a mean reward value of around 0.8, while the optimal policies, as BT shows, are between 6.9 and 1.0, a range only DreamerV3 (0.953) managed to enter. Additionally, PPO exhibits more instability in learning and greater sensitivity to initialization, as evidenced by a higher standard deviation. Following these, AMIGO reached an average reward of only 0.69, while the remaining techniques performed poorly, despite the fact that this level is practically fully observable.

MiniGrid-SimpleCrossingS11N5-v0: Here the agent faces a large grid with multiple intersec-tions and potential dead ends. The randomized layout in each episode also forces the agent to develop a robust exploration strategy. As Figure 7, and Table 5 show, DreamerV3, IMPALA, COUNT, RIDE, CURIOSITY, RND achieved near-optimal policies since their intrinsic reward mechanisms seem to be particularly helpful in this large environment where the observable window covers only a small part. NACE and AMIGO found reasonable policies with average rewards of 0.88 and 0.78respectively, while the remaining techniques scored below 0.5. Again, NACE's strength lies in its data efficiency, driven by its inductive biases, even though it does not converge to the optimal policy.



Figure 7: Learning curves in MiniGrid-SimpleCrossingS11N5-v0

Table 5: Final values for *MiniGrid-SimpleCrossingS11N5-v0*

4.3 DYNAMIC ENVIRONMENTS WITH SEQUENTIAL DEPENDENCIES

In these environments, the need to perform actions in a specific sequence adds complexity and tests the agent's ability to plan and execute multi-step strategies.



480 MiniGrid-Unlock-v0: In this scenario, the agent must first locate and pick up a key before unlocking a door to reach the goal and obtain the reward. This sequential dependency adds a layer of complexity that challenges the agent's ability to plan ahead. Even though it is a single sequential dependency, the DRL techniques that learned the fastest initially, DreamerV3 and PPO, demands almost a million timesteps to converge to a similarly effective policy as NACE, which achieves this within just 10³ steps again (as seen in Figure 8, and Table 6). Additionally, while PPO shows more instability in learning, it is far less chaotic than AMIGO. IMPALA reached the optimal policy, and

486 did so after about 2 million steps, performing similarly well as COUNT and AMIGO in the end. It 487 is also visible, in our runs, that TRPO did not exceed a mean episode reward of 0.6, while A2C and 488 DQN completely failed to learn any effective policy. 489

490 **MiniGrid-DoorKey-8x8-v0:** This environment introduces an additional layer of sequential de-491 pendency by requiring the agent to navigate through an unlocked door to reach a goal in a separate room. While passing through the door adds complexity, the primary challenge lies in the sparse re-492 ward structure, as no reward is given for merely using the door, since only reaching the final goal is 493 rewarded. As presented in Figure 9, Table 7, DreamverV3 and COUNT nearly achieved the optimal 494 policy with an average reward of 0.975 and 0.96. AMIGO reached 0.87, but within 10^7 timesteps, 495 which is below the average reward which NACE reached within only 10^4 steps. Overall the re-496 sults suggest poor combinatorial scaling of the involved DRL techniques, while NACE, on average, 497 required a similar amount of timesteps as for *MiniGrid-Unlock-v0* to learn an effective policy. 498



515 The observed sample efficiency of NACE originates from explicitly exploiting the cell-based grid 516 world state observations for creating transition rules. This represents a strong inductive bias, which 517 makes NACE less generic than DreamerV3. However it can nevertheless be valuable in broader 518 applications where mapping high-dimensional input to a similar discrete world representation is feasible. This mapping, dependent on the problem domain, can be implemented with the appropri-519 ate choice of feature extraction techniques. Such approaches are commonly employed in robotics, 520 where methods like Simultaneous Localization and Mapping and object detection models are in-521 tegrated to construct semantic maps for operating mobile robots Zhang et al. (2020a). However 522 we acknowledge this demands a considerable engineering effort, while DreamerV3 can update its 523 perceptual representations dynamically via gradient-based updates. Additional discussions, such as 524 about NACE's sub-optimality due to representational limitations, can be found in Appendix A. 525

5 CONCLUSION

500

501

502

504

506

507

511

514

526

527

528 We introduced NACE, an experiential learning agent designed to enhance data efficiency in grid 529 world environments by leveraging causally-informed intrinsic rewards and inductive biases. We 530 compared NACE with state-of-the-art DRL techniques, demonstrating that while these techniques 531 are able to eventually achieve near-optimal policies, they often require significantly more data, es-532 pecially as task complexity increases due to factors such as sequential dependencies. NACE, by 533 contrast, extends the RL framework to empirically support causal relations, enabling effective learn-534 ing and decision-making even in low data settings without relying on pre-defined causal models. Our causality-informed curiosity model, combined with the outlined inductive biases, facilitates system-536 atic exploration and learning requiring significantly fewer timesteps. We hope that future work in 537 the field will strike new compromises regarding the inclusion of inductive biases, leading to highly sample-efficient DRL that retains the ability to converge to optimal policies. Moving forward, we 538 plan to generalize NACE to handle three-dimensional and continuous spaces, as well as explore neural implementations of NACE, further advancing the capabilities of learning agents.

540	Reproducibility Statement
541	• We utilized onen course implementations of the selected DBL elections from public
542	• We duffized open-source implementations of the selected DKL algorithms from public repositories (not including our technique):
545	repositories (not including our teeninque).
544	- AMIGO was from here: https://github.com/facebookresearch/
040 540	adversarially-motivated-intrinsic-goals
040 547	 BT is here: https://github.com/andreneco/minigrid_bt
540	- DQN, A2C, TRPO, and PPO were established on Stable Baselines3 (SB3)'s
540 540	baselines repository (Raffin et al., 2021): https://stable-baselines3.
550	readthedocs.10/
550	- DreamerV3 was from here: https://github.com/qxcv/dreamerv3
552	 All the other were from here: https://github.com/sparisi/cbet
552	• We used the the Minigrid package for the environments in our comparison, which is avail-
554	able here: https://github.com/Farama-Foundation/Minigrid
555	• For NACE we provide a stand-alone zip archive for reviewers to reproduce our results.
556	which is runnable on a regular computer with Python interpreter. It includes a README.txt
557	in the NACE folder, as well as scripts to generate the tables and the plots present in the
558	paper.
559	
560	REFERENCES
561	
562	Nishanth Anand and Doina Precup. Prediction and control in continual reinforcement learning.
563	Advances in Neural Information Processing Systems, 36, 2024.
564	Marc Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Remi Munos.
565	Unifying count-based exploration and intrinsic motivation. Advances in neural information pro-
566	cessing systems, 29, 2016.
567	Daniel Bruce Wumpus world in introductory artificial intelligence Journal of Computing Sciences
568	in Colleges 27(2):58–65, 2011
569	<i>in concepts</i> , 27(2).50-65, 2011.
570	Yuri Burda, Harrison Edwards, Amos Storkey, and Oleg Klimov. Exploration by random network
571	distillation. arXiv preprint arXiv:1810.12894, 2018.
572	A Campero, R Raileanu, H Küttler, JB Tenenbaum, T Rocktäschel, and E Grefenstette. Learning
573	with amigo: Adversarially motivated intrinsic goals. In ICLR 2021-9th International Conference
574	on Learning Representations. International Conference on Learning Representations, 2021.
575	Vichi Coo, Zhiming Li, Tionnoi Vang, Hoo, Zhang, Van Zhang, Vi Li, Jianya Hoo, and Yang Liu
576	Galois: boosting deep reinforcement learning via generalizable logic synthesis Advances in
577	Neural Information Processing Systems 35:19930–19943 2022
578	
579	Maxime Chevalier-Boisvert, Bolun Dai, Mark Towers, Rodrigo Perez-Vicente, Lucas Willems,
580	Salem Lahlou, Suman Pal, Pablo Samuel Castro, and Jordan Terry. Minigrid & miniworld: Mod-
581	ular & customizable reinforcement learning environments for goal-oriented tasks. Advances in
502	wearai myormation r tocessing systems, 50, 2024.
583	Michele Colledanchise and Petter Ögren. Behavior trees in robotics and AI: An introduction. CRC
504 505	Press, 2018.
585	Dorit Dor and Uri Zwick, Sokohan and other motion planning problems, Computational Converter
587	13(4)·215–228 1999
588	
589	Lasse Espeholt, Hubert Soyer, Remi Munos, Karen Simonyan, Vlad Mnih, Tom Ward, Yotam
590	Doron, Vlad Firoiu, Tim Harley, Iain Dunning, et al. Impala: Scalable distributed deep-rl with im-
591	portance weighted actor-learner architectures. In <i>International conference on machine learning</i> ,
592	pp. 1407–1410. PMLK, 2018.
593	Huihui Guo, Fan Wu, Yunchuan Qin, Ruihui Li, Keqin Li, and Kenli Li. Recent trends in task and motion planning for robotics: A survey. <i>ACM Computing Surveys</i> , 55(13s):1–36, 2023.

633

- Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains through world models. *arXiv preprint arXiv:2301.04104*, 2023.
 Woojun Kim, Yongjae Shin, Jongeui Park, and Youngchul Sung. Sample-efficient and safe deep reinforcement learning via reset deep ensemble agents. *Advances in Neural Information Processing*
- Systems, 36, 2024.
 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wiertransport of Martin Dischaillen, Dhaing stari with deer minformerent learning, arXiv present.
- stra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Belle mare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level
 control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pp. 1928–1937. PMLR, 2016.
- Deepak Pathak, Pulkit Agrawal, Alexei A Efros, and Trevor Darrell. Curiosity-driven exploration
 by self-supervised prediction. In *International conference on machine learning*, pp. 2778–2787.
 PMLR, 2017.
- Judea Pearl. From bayesian networks to causal networks. In *Mathematical models for handling partial knowledge in artificial intelligence*, pp. 157–182. Springer, 1995.
- Martha E Pollack and Marc Ringuette. Introducing the tileworld: Experimentally evaluating agent
 architectures. In *AAAI*, volume 90, pp. 183–189, 1990.
- Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah
 Dormann. Stable-baselines3: Reliable reinforcement learning implementations. Journal of
 Machine Learning Research, 22(268):1-8, 2021. URL http://jmlr.org/papers/v22/
 20-1364.html.
- Roberta Raileanu and Tim Rocktäschel. Ride: Rewarding impact-driven exploration for
 procedurally-generated environments. In *International Conference on Learning Representations*,
 2020.
- John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *International conference on machine learning*, pp. 1889–1897. PMLR, 2015.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Ramanan Sekar, Oleh Rybkin, Kostas Daniilidis, Pieter Abbeel, Danijar Hafner, and Deepak Pathak.
 Planning to explore via self-supervised world models. In *International conference on machine learning*, pp. 8583–8592. PMLR, 2020.
- Matthijs TJ Spaan. Partially observable markov decision processes. *Reinforcement learning: Stateof-the-art*, pp. 387–414, 2012.
- Richard S Sutton, Andrew G Barto, et al. Reinforcement learning. *Journal of Cognitive Neuroscience*, 11(1):126–134, 1999.
- Pedro A Tsividis, Joao Loula, Jake Burga, Nathan Foss, Andres Campero, Thomas Pouncy, Samuel J
 Gershman, and Joshua B Tenenbaum. Human-level reinforcement learning through theory-based
 modeling, exploration, and planning. *arXiv preprint arXiv:2107.12544*, 2021.
- Andrew Wang, Andrew C Li, Toryn Q Klassen, Rodrigo Toro Icarte, and Sheila A McIlraith. Learn ing belief representations for partially observable deep rl. In *International Conference on Machine Learning*, pp. 35970–35988. PMLR, 2023a.

653 654

655 656

657

658

659 660

661

662 663

664

665

666 667

668

669 670 671

672

676 677

678

679

680

682

683

684

685

686

687

688

689

690

- Kaixin Wang, Kuangqi Zhou, Bingyi Kang, Jiashi Feng, and YAN Shuicheng. Revisiting intrinsic reward for exploration in procedurally generated environments. In *The Eleventh International Conference on Learning Representations*, 2023b.
 - Pei Wang. Non-axiomatic logic: A model of intelligent reasoning. World Scientific, 2013.
 - Jiadong Zhang, Wei Wang, Xianyu Qi, and Ziwei Liao. Social and robust navigation for indoor robots based on object semantic grid and topological map. *Applied Sciences*, 10(24):8991, 2020a.
 - Tianjun Zhang, Huazhe Xu, Xiaolong Wang, Yi Wu, Kurt Keutzer, Joseph E Gonzalez, and Yuandong Tian. Bebold: Exploration beyond the boundary of explored regions. *arXiv preprint arXiv:2012.08621*, 2020b.
 - Qian Zhao, Jinhui Han, and Mao Xu. Boosting policy learning in reinforcement learning via adaptive intrinsic reward regulation. *IEEE Access*, 2023.
 - Lulu Zheng, Jiarui Chen, Jianhao Wang, Jiamin He, Yujing Hu, Yingfeng Chen, Changjie Fan, Yang Gao, and Chongjie Zhang. Episodic multi-agent reinforcement learning with curiosity-driven exploration. *Advances in Neural Information Processing Systems*, 34:3757–3769, 2021.
 - Xun Zheng, Bryon Aragam, Pradeep K Ravikumar, and Eric P Xing. Dags with no tears: Continuous optimization for structure learning. *Advances in neural information processing systems*, 31, 2018.

APPENDIX A ADDITIONAL STUDIES AND DISCUSSIONS

The performance results indicate that NACE often exhibits sub-optimal outcomes. To analyze this, we present contributing factors from a conceptual perspective, examine the impact of hyperparameter choices, and assess robustness to non-determinism arising from random action consequences.

• **Representational Limitations**: NACE's rule-based framework captures only spatially relative dependencies from one timestep to the next. It does not exploit the inherent structural statistics of environment generation, which are leveraged by various DRL techniques. While these structural dependencies are most apparent in static environments where locations remain constant, they are also present in dynamic environments. For example, the goal location consistently appears in the bottom-right corner not only in *MiniGrid-Empty* levels but also in *MiniGrid-SimpleCrossingS11N5-v0*. NACE's inability to utilize these broader environmental patterns limits its performance compared to methods that can.

Additionally, NACE's rules are tied to an action, meaning agent-external changes that are not caused by NACE need to be learned for each action separately, considerably lowering its sample efficiency by a factor of the amount of actions. The mechanism could be extended to incorporate learning of rules without an action as a precondition, leaving it to evidence collection whether the action is considered dependent on the truth expectation of the alternative rules.

 Study of Reduced Planning Horizon: NACE's estimation of expected returns relies heav-691 ily on its planning horizon. Short planning horizons can significantly reduce performance, 692 especially in tasks requiring long-term planning. To quantify this effect, we examine two 693 cases: MiniGrid-DoorKey-8x8-v0, which demands longer-horizon planning, and MiniGrid-DoorKey-6x6-v0, which is less demanding in this regard. As shown in Table 8, running NACE with a planning horizon of only 8 steps in MiniGrid-DoorKey-8x8-v0 results in con-696 vergence to an average return of 0.48, whereas extending the horizon to 100 steps improves 697 the average reward to 0.92. In contrast, in *MiniGrid-DoorKey-6x6-v0*, NACE maintains an average reward of 0.93 regardless of the planning horizon. A similar pattern is observed in *MiniGrid-Empty-16x16-v0*, where the average reward drops from 0.91 to 0.41 699 when the planning horizon is reduced from 100 to 8 steps. These results highlight NACE's 700 dependence on adequate planning horizons for effective rule chaining and the significant performance degradation that occurs when the planning horizon is too short.

702	Environment	Planning Horizon, Average Reward
703	MiniGrid-DoorKey-6x6-v0	8 steps, 0.93
704	MiniGrid-DoorKey-6x6-v0	100 steps, 0.93
705	MiniGrid-DoorKey-8x8-v0	8 steps, 0.48
706	MiniGrid-DoorKey-8x8-v0	100 steps, 0.92
707	MiniGrid-Empty-16x16-v0	8 steps, 0.41
708	MiniGrid-Empty-16x16-v0	100 steps, 0.91
709	Table 8. NACE Perform	ance with different planning horizons
710	Table 6. NACE FEITOIII	lance with different plaining nonzons
711		
712	Robustness to Non-Determinisr	n : NACE's rule representation incorporates uncertainty
713	handling through evidence counte	rs, enabling it to cope with non-deterministic state tran-
714	sitions. To assess this capability, v	we modify the environment to invoke unintended actions $\frac{1}{2}$
715	with certain probabilities. In <i>Mi</i>	<i>niGrid-Empty-10x10-v0</i> , when 10% of actions result in
716	hasia talaranga ta non datarminis	n However when the probability of unintended actions
717	increases to 20% NACE fails to c	omplete the task within the maximum allowed time in all
718	episodes Higher tolerance to not	n-determinism can be achieved by increasing the default
719	truth expectation threshold for ru	le usage above the default value of 0.5 However this
720	adjustment reduces sample efficie	ncv. as it requires the agent to confirm each rule multiple
721	times before utilizing it.	,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,
722	C	
723		EFFE
724	APPENDIX B ABLATION STUDY	: EFFECTS OF OMITTING KEY INDUCTIVE
725	BIASES IN NACE,	AND INDUCTIVE BIASES IN DRL
726		
727	B.1 CAUSAL RULE REPRESENTATION	
728	The causal rule representation is foundation	al to NACE's operation and cannot be omitted. However
729	we analyze the effects of reducing the n	anning horizon which limits the depth of chaining in
730	Appendix A.	aning nonzon, which mints the depth of channing, m
731	- ppendin - in	
732	B 2 TEMPORAL LOCALITY AND ATTER	NTIONAL BLAS
733	D.2 TEMPORAL EDGALITY AND ATTEM	CHORAE DIAS
734	Omitting these biases with larger environm	ent sizes is infeasible due to the combinatorial explosion
735	of potential rules as we will now analyze.	
736		
737	• For an environment of size $w \times h$	<i>i</i> , the number of possible rule preconditions for a single
738	timestep is $2^{w\cdot n}$, as each particula	r cell can either be considered or not be considered in the
739	precondition of a rule.	
740	• For a time window of du	ration d, this expands to $2^{w \cdot h \cdot d}$, leading to
741	18446744073709551616 possible	e rule preconditions for an 8×8 grid within a sin-
742	gle timestep.	
743	 NACE is tied to Markov Assump 	tion particularly within the observational window, as all
744	rule construction and updating c	onsiders only the previous and current state. However,
745	its bird view map representation	also contains values from observations of previous time
746	steps which are currently out of v	iew of the agent, which brings us to the next point, state
747	tracking.	
748		
749	B.3 STATE TRACKING	
750		
751	without state tracking, NACE lacks mem	ory of prior observations and memory of observations
752	which are outside of its field-of-view. This	results in oscillatory behavior caused by the exploration
753	sualegy of the agent, as it can only utilize	

In our experiments across 10 runs in MiniGrid-Empty-8x8-v0, this led the agent to turn indefinitely due to the curiosity model assigning low match values to previously visited areas (due to the lack of state tracking they are always considered to be of unknown value)

756	which are now outside of the field-of-view of the agent. The closest such cell is immediately				
757	behind the agent with the default partial observation model in Minigrid, which explains the				
758	behavior.				
759	• State Treaking plays a critical role in ansuring purposeful exploration and desision making				
760	for the agent to know which places have been visited and what it has been observing at the				
761	particular locations, as well as which locations have yet to be observed.				
762	 Sequential demonstration of the demond on state the plane. An ensurely of this is released does 				
763	• Sequential dependencies often depend on state tracking. An example of this is when a door				
764	concurrently, demanding some form of spatial memory. Another form of state tracking lies				
765	in the observable inventory array, which when absent would need the modeling of long-				
766	range temporal dependencies (e.g. did the agent already pick up the key?) which would				
767	demand a suitable model structure to be learnable by the agent.				
768					
769	B.4 SPATIAL EOUIVARIANCE				
770					
771	The absence of spatial equivariance significantly impacts sample efficiency.				
772					
773	• Each rule must be learned independently for every location, meaning in an 8x8 grid, the agent has to learn 64 times the same set of rules. However, since particular arrangements of				
774	cell values will not re-appear through the environment generation, it can take significantly				
775	longer to learn the relevant knowledge without this bias				
776	. Hence for the second constraint and environment of size why heating increases the mentioned				
777	• Hence for the general case with an environment of size $w \times h$, this increases the required sample count at least by a factor of w, h harming significantly the sample efficiency of the				
778	sample count at least by a factor of $w \cdot n$, naming significantly the sample efficiency of the technique				
779	Companyed				
780	• Conceptually, we also would like to point out that the rule learning mechanisms do not allow to learn spatial equivariance retrospectively either while some DPL techniques de				
781	pendent on the model structure, could potentially acquire it				
782	pendent on the model structure, courd potentially acquire it.				
783	These results highlight the necessity of each inductive bias in ensuring the scalability, efficiency, and				
784	functionality of NACE.				
700					
700	B.5 WHICH INDUCTIVE BIASES ARE PRESENT IN THE DRL TECHNIQUES				
788	In the main mean and in a the induction biases of NACE. However, would like to point out that				
789	in the main paper we outlined the inductive blases of NACE. However we would like to point out that some of them are also inherent in the DPL techniques, complementing our discussion on inductive				
790	biases in DRL and NACE.				
791					
792	• Temporal Locality : The DRL methods perform best when the Markov Assumption is met,				
793	despite LSTM allowing to cope with partial observability, the need to capture long-range				
794	temporal dependencies makes sample efficient learning more difficult.				
795	• Causal representation: While not explicitly stated as a set of cause-effect relations.				
796	DreamerV3's learned dynamics model can predict the consequence states of actions, which				
797	is not the case for the model-free DRL methods. Such modeling is to some extent inde-				
798	pendent from the objective (what is rewarding), and allows an agent to train itself from				
799	simulated experience by predicting novel states, and to reach novel goals.				
800	• Spatial Equivariance: Clearly the DRL techniques do not have an explicit rule represen-				
801	tation, however the Convolution layers in the DRL policies allow for learned features to be				
802	identified at different locations, improving generalization.				
803	• State Tracking: Is not explicitly handled by the DRL techniques as a separate point, in-				
804	stead it is handled in the same way as non-local temporal dependencies in the LSTM-				
805	including policies, while NACE builds a bird view map explicitly, which can be considered				
806	to be a form of spatial memory.				
807					
001	• Attentional Bias: While NACE has a strong prior to which cells to consider based on				
808	• Attentional Bias: While NACE has a strong prior to which cells to consider based on observably changed values and prediction mismatches, the DRL policies with Convolution				

810 APPENDIX C NOTATION AND SYMBOLS 811

Symbol	Description
s	State, represented as a combination of a 2D grid (s_{grid}) and a 1D array (s_{array})
a	Action taken by the agent
r	Causal rule in the form (preconditions, action) \Rightarrow consequence
$c_{t,x,y}$	Cell value at position (x, y) in the 2D grid at time t
\overline{c}	Equality constraint on a cell value (e.g., $c_r = c$)
v_t	Value array at time t
\overline{v}	Equality constraint on value array (e.g., $v_r = v$)
M(r,c)	Match value of a rule r for cell c, based on the fraction of preconditions satisfied
C(c)	Cell match value for cell c, derived from the maximum match value across all rules
S(s)	State match value for state s, calculated as the average $C(c)$ for cells with $C(c) > 0$
$w_+(r)$	Positive evidence counter for rule r, incremented when predictions align with observations
$w_{-}(r)$	Negative evidence counter for rule r , incremented when predictions differ from observations
w(r)	Total evidence count for rule r, defined as $w(r) = w_+(r) + w(r)$
frequency(r)	Fraction of positive evidence for rule r, defined as $f(r) = \frac{w_+(r)}{w(r)}$
confidence(r)	Confidence factor for rule r, defined as $c(r) = \frac{w(r)}{w(r)+1}$
$f_{\exp}(r)$	Expected truth value for rule r, calculated as $f_{\exp}(r) = (f(r) - \frac{1}{2}) \cdot c(r) + \frac{1}{2}$
M_t^{change}	Set of cells with changes in observed values between timesteps $t-1$ and t
$M_{\text{mismatched},t}^{i}$	Set of cells where observed values differ from predicted values at timestep t
$M_{\text{mismatched }t}^{\text{prediction}}$	Set of cells where predictions differ from observations at timestep t
R(r)	Reward associated with rule r
R(s)	Reward associated with a state s, defined as the average reward of rules applied to generate s
V(s)	Value of state s, used in planning for maximizing long-term returns
$\pi(t)$	Planned action sequence or policy at timestep t
γ	Discount factor for future rewards

APPENDIX D PSEUDOCODE

838 839

840 841

842

The system can be described by the pseudocode:

843	
844 -	
845	Algorithm 1: Pseudocode of NACE
846	• Actual World: perceived_array = perceive_partial(world)
847	• Observer
848	
849	$s_t = update_bird_view(s_{t-1}, perceived_array)$
850	$calculate(M_{change}, M_{mismatched}^{observation}, M_{mismatched}^{prediction})$
851	Hypothesizer:
852	- Create new rules for which $w_+(r) = 1$.
853	- Update rule evidences according to $w_+(r)$ and $w(r)$.
854	- Choose rules r_1 with $w_{\perp}(r_1) > w_{\perp}(r_1)$ for which there does not exist a rule r_2 with
855	same precondition and action, but different postcondition with $f_{exp}(r_2) > f_{exp}(r_1)$.
856	Planner utilizing Predictor:
857	= DEC with Direction(V(z) > 0)
858	$a_1, \dots, a_n = BFS_witn_Predictor(V(s) > 0)$ $a^* = a^* = BFS_with_Predictor(min(C(a)) < 1)$
859	$u_1, \dots, u_n = DFS_with_Predictor(min(S(S)) < 1)$
860	transition function
861	If found($a_1 = a_2$); return $a_1 = a_2$
862	If found $(a_1^*,, a_n^*)$; return $a_1^*,, a_n^*$
863	Else Perform a random action

864 Appendix E HYPERPARAMETER DETAILS 865

866 **E.1** FOUNDATIONAL ALGORITHMS 867

- 868 E.1.1 CORE MODELS AND THEIR MECHANISMS
 - Deep Q-Network (DQN): DQN integrates deep neural networks with classical Q-learning, making it effective for handling large state spaces. To stabilize training, DQN uses experience replay and a separate target network. The Q-value update in DQN follows:

$$Q(s,a) \leftarrow Q(s,a) + \alpha \left(R(s) + \gamma \max_{a'} Q(s',a') - Q(s,a) \right)$$

A(s,a) = Q(s,a) - V(s)

 $\theta \leftarrow \theta + \alpha \nabla_{\theta} \log \pi_{\theta}(a|s) A(s,a)$

where:

870

871

872 873

874

875

876

877

878

879

882

883

884

885

886

887 888

889

890

891

892

893 894

895

896

897

899

900 901 902

903 904

905

906

907

908

909

910 911

912

913

914

915

916 917

- s, a: current state and action, - s', a': next state and action, - R(s): reward received, - γ : discount factor for future rewards, - α : learning rate. • Advantage Actor-Critic (A2C): A2C builds on the actor-critic framework, synchronizing multiple parallel learners to reduce variance in policy updates. It calculates an advantage function to evaluate actions relative to the current policy's value estimate, stabilizing training but requiring frequent environmental interactions due to its on-policy nature. **Advantage Function:** Policy Update: The policy is updated using the gradient: where: - Q(s, a): action-value function, - V(s): state-value function, - $\pi_{\theta}(a|s)$: policy parameterized by θ , - α : learning rate. • Trust Region Policy Optimization (TRPO): TRPO addresses stability in policy updates by enforcing a trust region constraint, ensuring small policy changes during optimization. This constraint is implemented via a KL-divergence bound, preventing drastic shifts in behavior but requiring computationally expensive second-order optimization. **Objective Function:**

$$\max_{\theta} \mathbb{E}_{s \sim \pi_{\theta_{\text{old}}}} \left[\frac{\pi_{\theta}(a|s)}{\pi_{\theta_{\text{old}}}(a|s)} A(s,a) \right]$$

Constraint:

$$\mathbb{E}_{s \sim \pi_{\theta_{\text{old}}}} \left[D_{\text{KL}}(\pi_{\theta_{\text{old}}} || \pi_{\theta}) \right] \leq \delta$$

where:

- $\pi_{\theta}(a|s)$: new policy,
- $\pi_{\theta_{\text{old}}}(a|s)$: previous policy,
- A(s, a): advantage function,
- $D_{\rm KL}$: KL-divergence,
- δ: trust region size.
- Proximal Policy Optimization (PPO): PPO refines TRPO by introducing a clipped surrogate objective, which simplifies computation and allows for multiple updates per batch. This approach improves data utilization while maintaining policy stability.

Clipped Surrogate Objective:

$$\max_{\theta} \mathbb{E}_{s,a} \left[\min \left(\frac{\pi_{\theta}(a|s)}{\pi_{\theta_{\text{old}}}(a|s)} A(s,a), \operatorname{clip} \left(\frac{\pi_{\theta}(a|s)}{\pi_{\theta_{\text{old}}}(a|s)}, 1-\epsilon, 1+\epsilon \right) A(s,a) \right) \right]$$

where:

918	$-\pi_{\theta}(a s)$: new policy.
919	$-\pi_{\theta_{\alpha}}(a s)$: old policy.
920	- A(s, a): advantage function
921	$-\epsilon$ clipping threshold
922	- c. enpping uneshold.
923 024	E.1.2 Hyperparameter Configuration for Foundational Algorithms
925	We utilize the Stahle Deceline 2 from much (Deffer at al. 2021) to train and such as foundational
926	algorithms leveraging its pre-implemented models and customizable configurations. All algorithms
927	use the same convolutional neural network architecture to process observations, ensuring consis-
928	tency across experiments. The hyperparameters for each algorithm were selected based on achiev-
929	ing the best average performance across all tasks, rather than optimizing for a single task, to ensure
930	generalizability. The details of the network architecture and training setup for each algorithm are
931	outlined below.
932	• Network Architecture: Observations $(7 \times 7 \times 3)$ from the Minigrid environment are pro-
933	cessed through four convolutional layers. Each layer is configured as follows:
934	- Kernel size: 2×2
935	– Activation: ReLU
937	- Increasing number of filters: 16, 32, 64, and 128
938	The output of the final convolutional lower is flattened and passed to a fully connected lower
939	with:
940	Outent dimension 120
941	- Output dimension: 128
942	- Activation: ReLU
943	• Training Setup for DQN:
944	– Learning rate: 0.0001
945	– Buffer size: 1,000,000
946	– Learning starts: 100
947	- Batch size: 32
948	- Soft update coefficient: 1
949	– Discount factor: 0.99
950	– Train frequency: 4
952	- Gradient steps: 1
953	– Target update interval: 10,000
954	– Exploration fraction: 0.1
955	– Initial exploration epsilon: 1.0
956	– Final exploration epsilon: 0.05
957	– Max gradient norm: 10.0
958	• Training Setup for A2C:
959	- Learning rate: 0.0007
960	– Number of steps: 5
961	- Discount factor: 0.99
962	- Entropy coefficient: 0.0
963	- Value function coefficient: 0.5
904 065	- Max gradient norm: 0.5
205	• Training Satur for TDDO:
967	
968	- Learning rate: 0.001
969	- Number of steps: 2048
970	- Batch size: 128
971	– Discount factor: 0.99
	– Conjugate gradient max steps: 15

972	- Conjugate gradient damping: 0.1
973	- Line search shrinking factor: 0.8
974	- Line search max iterations: 10
975	- Enc scatch max netations. 10
976	- Number of critic updates. 10
977	- Target KL divergence: 0.01
978	 Training Setup for PPO:
979	- Learning rate: 0.0003
980	– Number of steps: 2048
981	- Batch size: 64
982	– Number of epochs: 10
983	- Discount factor: 0.99
984	- Clin range: 0.2
985	- Chp lange. 0.2
986	- Entropy coefficients 0.5
987	- value function coefficient: 0.5
988	- Max gradient norm: 0.5
989	
990	E.2 MODEL-BASED ALGORITHM: DREAMERV3
991	DreamerV3 is a model-based RL algorithm designed to enhance sample efficiency by learning a
992	latent world model of the environment. It optimizes both the world model and the policy within the
993	latent space, reducing the computational demands of interacting with the environment.
994	
995	World Model: The latent dynamics model predicts future latent states z based on prior latent state
996	z_{t-1} , action a_{t-1} , and reward R_{t-1} . This model facilitates long-term planning without requiring explicit relieves in the actual environment
997	explicit follouis in the actual environment.
998	Policy Optimization: The policy maximizes expected rewards in the learned latent space by lever-
999	aging the dynamics model to simulate trajectories. Policy updates use gradient-based methods in-
1000	formed by imagined rollouts.
1001	Loss Function:
1002	$\mathscr{L}_{\text{DreamerV3}} = \mathscr{L}_{\text{Reconstruction}} + \mathscr{L}_{\text{Dynamics}} + \mathscr{L}_{\text{Policy}}$
1003	where:
1004	
1005	• $\mathscr{L}_{\text{Reconstruction}}$: Measures the accuracy of reconstructing environment observations,
1006	• Louranies: Captures consistency in latent state transitions.
1007	\mathscr{A} - Maximizes imagined rewards
1008	• $\mathcal{Z}_{\text{Policy}}$. Maximizes imagined rewards.
1009	Hyperpersenter Configuration for DreamarV3. The hyperpersenter configuration has been
1010	chosen to match the settings provided in https://github.com/gycy/dreamery3 To avoid
1011	redundancy and maintain brevity, we do not include the full configuration here due to its extensive
1012	nature.
1013	
1014	E.3 MODEL-FREE EXPLORATION AND SCALABILITY EXTENSIONS
1015	
1017	All experiments for the other model-free methods are based on the Torchbeast implementation of
1018	IMPALA (Espeholt et al., 2018), which has been modified to support intrinsic reward algorithms as
1010	described in Raileanu & Rocktäschel (2020) and Campero et al. (2021). The hyperparameters were
1020	selected following the configurations used in these references. For clarity, we first list the values
1020	snared across all algorithms, followed by the specific details unique to each one.
1021	
1022	E.J.1 SHARED HYPERPARAMETERS
1023	• Network Architecture: Observations ($7 \times 7 \times 3$ for Minigrid) are processed through three
1024	convolutional layers:
· v … v	

- Number of filters: 32 per layer

1026	- Kernel size: 3×3		
1027	- Stride: 2		
1028	- Surde: 2		
1029	- Padding: 1		
1030	- Activation: Exponential Linear Unit (ELU)		
1031	The output of the convolutional layers is passed to:		
1032	- An LSTM layer to address partial observability by maintaining temporal dependencies		
1033	and encoding sequences of observations.		
1034	- A fully connected layer for computing:		
1035	* Policy logits: Unnormalized scores for each action, converted to probabilities us-		
1037	ing a softmax function.		
1038	* Value estimates: Predictions of expected future returns, used in actor-critic meth-		
1039	ods.		
1040	• Training Setup:		
1041	- Number of actors: 40		
1042	- Number of actors: 40		
1043	- Number of burlets: 80		
1044	– Unroll length: 100		
1045	– Number of learner threads: 4		
1046	– Batch size: 32		
1047	– Discount factor: 0.99		
1048	– Learning rate: 0.0001		
1049	– Policy entropy loss: 0.0005		
1051	- Gradient clipping: Norm of 40		
1052	- Save interval: Every 20 minutes		
1053	Snecial Parameters (Only When Annlicable)		
1054			
1055	- Count reset probability: 0.001 (COUNT, RIDE)		
1056	– Hash bits: 128 (COUNT)		
1057			
1058	E.3.2 INTRINSIC REWARDS AND COEFFICIENTS		
1059	Intrinsic rewards address sparse rewards and inefficient exploration. Each algorithm applies scal-		
1061	ing coefficients to normalize its intrinsic rewards. Additionally, all techniques incorporate policies		
1062	enhanced with LSTMs to address partial observability by maintaining memory of past observations		
1063	and actions.		
1064	• IMDALA. No intrincic reward $(n = 0.0)$		
1065	• IMPALA: No intrinsic reward $(r_i = 0.0)$.		
1066	• COUNT: $r_i = 0.005$.		
1067	• RIDE: $r_i = 0.1$.		
1068	• CURIOSITY: $r_i = 0.1$.		
1069	$\mathbf{P} = \mathbf{P} \mathbf{P} \mathbf{P} \mathbf{P} \mathbf{P} \mathbf{P} \mathbf{P} \mathbf{P}$		
1070	• RND: $r_i = 0.1$.		
1071	• AMIGO: $r_i = 0.1$ (applies to the teacher's intrinsic rewards).		
1072			
1073	The formal definitions of the intrinsic rewards are:		
1075	COUNT: The interime and is based on state a lifetime of the state of t		
1076	visited states:		
1077	1		
1078	$r_i = \overline{N(s_0)},$		
4.0000	- · (~0)		

where $N(s_0)$ is the (pseudo)count of visits to state s_0 . Counts are never reset during training.

1080 **RIDE** (Rewarding Impact-Driven Exploration): The intrinsic reward combines state novelty and state-change impact: 1082

$$r_i = \|\phi(s) - \phi(s_0)\|_2 \cdot \frac{1}{N(s_0)}$$

1084 where ϕ is trained to minimize both forward and inverse dynamics prediction errors. Counts $N(s_0)$ are reset at the beginning of each episode. 1086

CURIOSITY: The intrinsic reward comes from the prediction error of a forward dynamics model f, which predicts the next state embedding $\phi(s_0)$ from the current embedding $\phi(s)$ and action a: 1089

 $r_i = \|f(\phi(s), a) - \phi(s_0)\|_2.$

RND (Random Network Distillation): The intrinsic reward is computed as the prediction error of a trainable network ϕ attempting to match the output of a fixed random network ϕ :

$$r_i = \|\phi(s_0) - \hat{\phi}(s_0)\|_2.$$

AMIGO: The teacher policy generates goals g for the agent, with rewards given as:

	$\int +1$	if s_t satisfies g_t
$r_i = v(s_t, g) = \langle$	0	otherwise.

1100 The total reward is a weighted sum of intrinsic and extrinsic rewards: 1101

 $r_t = \beta r_i + \alpha r_e$, with $\beta = 0.3, \alpha = 0.7$.

1103 ALGORITHM-SPECIFIC HYPERPARAMETERS AND ARCHITECTURES E.3.3 1104

• IMPALA (Baseline):

- Intrinsic reward: None.
- Loss: Policy gradient, baseline, and entropy losses.
- 1108 • COUNT: 1109

1087

1088

1090 1091

1092

1093 1094 1095

1099

1102

1105

1106

1107

1110

1111 1112

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123 1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

- Intrinsic reward: State visitation counts.
- Uses count reset probability: p = 0.001.
- CURIOSITY:
 - Intrinsic reward: Forward prediction error.
 - Modules:
 - * State embedding model: Encodes observations into 256-dimensional embeddings.
 - * Forward dynamics model: Predicts next state embedding given current embedding and action.
 - * Inverse dynamics model: Predicts action given two successive state embeddings.
 - Loss weights:
 - * Forward dynamics loss: 10.0.
 - * Inverse dynamics loss: 0.1.

• RIDE:

- Intrinsic reward: Product of state visitation counts and the norm of state embedding changes.
 - Modules: Same as CURIOSITY.

• RND:

- Intrinsic reward: Prediction error between random target network and predictor network embeddings.
- Modules:
 - * Random target network: Produces fixed embeddings for observations.
 - * Predictor network: Trained to predict random target embeddings.
 - Loss weight: 0.1.

1134	• AMIGO:			
1136				
1137	- Intrinsic reward: Teacher-generated rewards.			
1138	– Teacher-specific parameters:			
1139	* Intrinsic reward coefficient (β) : 0.3			
1140	 Extrinsic reward coefficient (p): 0.3. 			
1141	* Examine reward coefficient (α). 0.1.			
1142	* Generator balon size: 150.			
1143	* Generator entropy cost: 0.05.			
1144	* Generator threshold: -0.5 .			
1146				
1147				
1148	APPENDIX F HARDWARE AND RUNTIME			
1149				
1150	In this section, we describe the hardware setup used to run the techniques and provide runtime			
1151	characteristics, including the duration of a single run for each technique. We specify the CPU and			
1152	GPU types and indicate whether the GPU was utilized for the corresponding models. While we			
1153	computational cost, but rather on sample efficiency.			
1154	NACE.			
1156	NACE:			
1157	• CPU: Apple M2 with 24GB RAM			
1158	CI C. Apple M2 with 240D KAM			
1159	• GPU: Not utilized for this technique			
1160				
1161	• Runtime: Approximately 15 minutes runtime till convergence per environment per run with			
1162	respective seed.			
1164				
1165	Intrinsic reward models (COUNT, RIDE, CURIOSITY, RND, AMIGO) and IMPALA:			
1166	• CPU: Intel Core i7-9750H with 32GB RAM			
1168				
1169	• GPU: Geforce GTX-1660 Ti with 6GB RAM			
1170	• Runtime: Approximately 8 hours per run on average			
1172	nammer representation of the ran on a course			
1173	Baseline models (TRPO, PPO, A2C, DQN):			
1174				
1175	• CPU: 1 compute node with 64 cores and 512GB RAM in total			
1176				
1177	 GPU: NVIDIA Tesla A100 HGX GPU with 80GB RAM 			
1170				
1180	• Runtime: Approximately I hour per run on average			
1181 1182	DreamerV3:			
1183	• CPU: 1 compute node with 64 cores and 512GB RAM in total			
1185				
1186	GPU: NVIDIA Tesla A100 HGX GPU with 80GB RAM			
1107	• Runtime: Approximately 40 hours per run on average			

1188 APPENDIX G EXAMPLE ENVIRONMENT WITH LEARNED RULES



1239 Prior to moving to Minigrid NACE was first tested with internal levels.

1240

1241

• Level 1: Food collection. In this level, as depicted in Figure 11, the agent needs to collect food.





