
Construing the Reward Function: Maximum Likelihood Inverse Reinforcement Learning with Limited Cognitive Resources

Sam Liang

Advisors: Mark K. Ho and Thomas Griffiths

Princeton University

saml@princeton.edu

Abstract

Humans lack the cognitive resources to plan optimally when faced with tasks. We cannot consider every factor of a task when determining a plan to solve the task and achieve our goal. Yet, many work on inferring the plan of an expert assume that they are planning optimally. This is unrealistic when humans are the experts. When applying these algorithms to the real world, they may perform poorly, because they fail to take into account that humans do not plan optimally. We would like inference algorithms to learn the optimal plan for a task even when they only have access to sub-optimal data. I present a revised Maximum Likelihood Inverse Reinforcement Learning algorithm that assumes the expert may not be acting optimally, using construals to model an expert acting sub-optimally in which they consider a subset of features when planning. It aims to learn the optimal plan from sub-optimal experts. Performance between the original and construed Maximum Likelihood algorithm on an expert that considers a subset of features is presented.

Contents

1	Introduction	4
1.1	Markov Decision Process	4
1.2	Learning From Demonstration Terminology	4
2	Methodology	6
2.1	Research Goals	6
2.2	Gridworld MDPs, Trajectories, and Policies	6
3	Related Work	7
3.1	Imitation Learning	7
3.1.1	Implementation Details	8
3.2	Intention Learning	9
3.2.1	Implementation Details	10
4	Construal	11
4.1	Value of Representation	11
5	Construed Maximum Likelihood IRL Algorithm	12
5.1	Construing the Reward Function	12
5.2	Loss Function	12
5.3	Implementation Details	13
6	Experiments and Evaluations	14
6.1	Imitation Learning Classifier and Maximum Likelihood IRL	14
6.2	Construed Maximum Likelihood IRL	19
7	Conclusion	21
7.1	Imitation and Inverse Reinforcement Learning	21
7.2	Construed Maximum Likelihood IRL Algorithm	22
7.3	Future Work	22
8	Code	22
9	Honor Code	22
References		23
10	Appendix: Algorithms, Policy Graphs, and State Value Maps	24
10.1	Maximum Likelihood IRL Algorithm	24
10.2	Classifier and MLIRL Results on Gridworld 1	25

10.3 Classifier and MLIRL Results on Gridworld 2	38
10.4 Classifier and MLIRL Results on Gridworld 3	49
10.5 Classifier and MLIRL Results on Gridworld 4	63
10.6 Classifier Training with 1-hot Vector Representation	69
10.7 CMIRL Results on Gridworld 1	71
10.8 CMIRL Results on Gridworld 5	75
10.9 CMIRL Results on Gridworld 6	79
10.10 CMIRL Results on Gridworld 7	83

1 Introduction

Learning from demonstration is learning how to do something by examining an expert perform the task. It can be divided into two categories: imitation learning and intention learning [6]. Most work in imitation and intention learning assumes the expert is planning optimally [8]. However, this assumption may lead to learning a sub-optimal plan, because the expert is not planning optimally, when we want to learn the optimal plan.

In this paper, I present a construal version of the Maximum Likelihood Inverse Reinforcement Learning algorithm that seeks to address this issue. It eliminates the assumption that the expert must be acting optimally by introducing construals [4] into the algorithm. A construal is a subset of all the details or features of a task, and the idea is that humans form construals when planning for a task [4]. Equally, this paper seeks to understand imitation and intention learning. Thus, I also implement a classifier trained with supervised learning and the Maximum Likelihood algorithm as baselines to compare performance between each other and the construed Maximum Likelihood algorithm.

1.1 Markov Decision Process

A Markov Decision Process (MDP) is a model of the expert's environment and task. It defines a set of states S ; a set of initial states S_0 ; a set of actions A that the expert can take in each state; a transition function $T : S \times A \times S \rightarrow [0, 1]$ that defines the probability of going to state i when taking action a in state j ; and a reward function $R : S \times A \times S \rightarrow \mathbb{R}$ that defines the reward of taking an action in state i and arriving in state j . The expert traverses through the MDP by taking actions in each state and arriving in another state until it reaches the goal state, mimicking the expert performing an action in their environment and potentially changing environment conditions before accomplishing their goal in the task.

1.2 Learning From Demonstration Terminology

Formally, in learning from demonstration, an agent takes in as input examples of the expert traversing through the MDP and tries to learn a policy π from those examples that allows it to act and perform like the expert in the MDP [6]. A policy π is a mapping from states to a probability distribution over actions. It specifies with what probability an agent takes an action in a certain state. With a policy, an agent will be able to decide on which actions to take as it traverses

the MDP. Furthermore, an optimal policy is one that maximizes the expected discounted reward which is how much reward on average the expert will acquire if they take this action in a state. As mentioned, most work in imitation and intention learning assume that the expert is using an optimal policy which may not be the case. Furthermore, the examples of the expert traversing through the MDP are called a set of trajectories D . A trajectory $d_i \in D$ is a sequence of state-action pairs $d_i = \langle (s_0, a_0), (s_1, a_1), \dots, (s_n, a_n) \rangle$ [6] starting from an initial state s_0 , ending at the goal state, and specifying in order which states the expert visited and what actions they took in those states. In both imitation and intention learning, once trajectories are acquired, a modification of them through a process called featurization is required as the data must only contain numbers. A state-feature vector function $\phi : S \rightarrow \mathbb{R}^n$ converts states to state-feature vectors [6] and an action-feature vector function $\varphi : A \rightarrow \mathbb{R}$ converts actions to numbers (labels). A state-feature vector is a vector of features, usually numbers, that identifies the state. For example, in Figure 1, some gridworld features are {lightgreen, green, red, lightblue, and blue}. A red state could be featurized as $\langle x\text{-coordinate}, y\text{-coordinate}, a \text{ number corresponding to its color} \rangle$ or as a 1-hot vector where each index corresponds to a color.

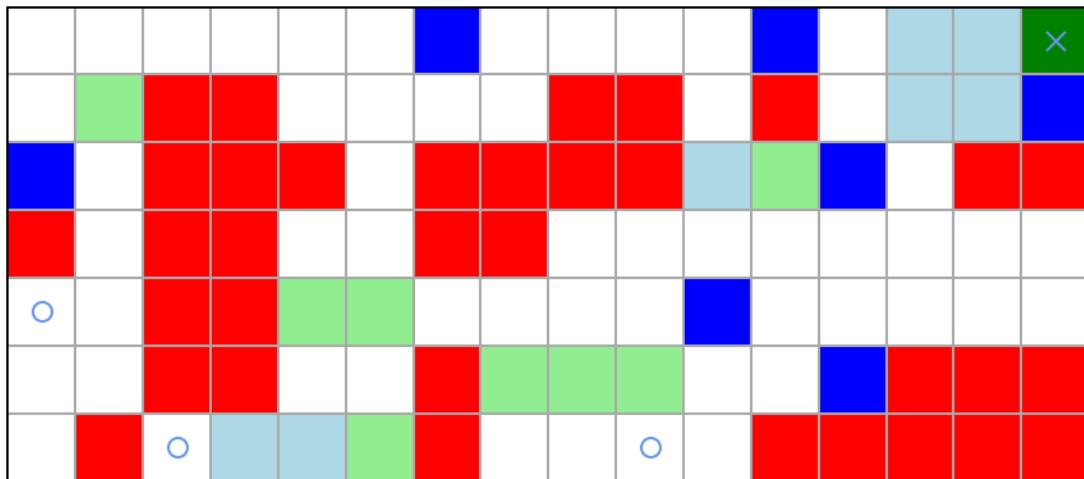


Figure 1: A gridworld example consisting of states with different features indicated by its color. Created with MSDM [5].

2 Methodology

2.1 Research Goals

The main goals for this project were two-fold: first to understand imitation learning and IRL, and secondly to extract the optimal policy from watching a sub-optimal expert.

2.2 Gridworld MDPs, Trajectories, and Policies

In this paper, all algorithms, models, and experiments are performed on gridworlds which are created through the MSDM codebase [5]. See Figure 1 for an example gridworld. Gridworlds are MDPs and have all the characteristics of MDPs: states are square tiles; actions are up, left, right, down, and stay; start states are indicated by a blue circle on a white tile; and the goal state is indicated by a blue "x" on the green tile. The goal of the expert is to reach the goal state from the starting states. All gridworlds in this paper have a discount rate of 0.99 and step cost of -1. The discount rate is how much future rewards are discounted and the step cost is the cost of taking any action.

In gridworlds, each action and state has a reward when performing the action and landing on the state. In my MSDM [5] gridworlds, the only action reward is the step cost, and state rewards are associated with its color. In Figure 1, states that are red have a reward of -500. Furthermore, gridworld features are characteristics of states. In a gridworld, some features are the (x, y) coordinates and color of a state. These are the only features I consider in the gridworlds and use in the featurization process. A note on the 1-hot vector featurization process in this paper is that white is not a feature and green (goal state) can be a feature. This is because a vector of zeros will correctly identify a white state. Green is considered the same as white when it has a reward of 0, because the goal state is a known global attribute. If it has a nonzero reward, it is considered a feature.

Meanwhile, the trajectories dataset contains trajectories generated by an expert executing their plan on a gridworld. Trajectories begin at a start state and end at the goal state. All trajectories were generated from scratch by repeatedly running the expert's policy on the gridworld using MSDM [5]. Policies are computed using MSDM's [5] Maximum Entropy Reinforcement Learning (RL) algorithm.

3 Related Work

While much work in imitation learning and inverse reinforcement learning model the expert as acting optimally, there is research on experts acting sub-optimally like humans. The Receding Horizon IRL [6] algorithm introduces a horizon parameter as a form of sub-optimal behavior. The horizon specifies how far into the future an expert plans which is done by limiting how many steps into the future from the current state that the expert takes into account when calculating the action which maximizes the expected discounted future reward from the current state [6]. This certainly addresses humans' cognitive resources as they do not plan infinitely ahead. While I did not use this algorithm as the baseline IRL algorithm, construals can be introduced as another form of sub-optimal behavior alongside the horizon. Meanwhile, online Bayesian goal inference [8] models the expert as a boundedly-rational planning agent that is limited by how far ahead they can plan (horizon idea), and they also continue to update their previously computed plan on-the-fly as they execute it and traverse the MDP. There is also the Bayesian Theory of Mind [2] model that uses a partially observable MDP (POMDP) to outline a human expert's environment. A POMDP is an MDP that has unobservable states and in their case, unobservable states are those that are not in the human's line-of-sight. They model an expert as someone who forms beliefs based on prior knowledge and what they can currently perceive in the POMDP and plans based on their beliefs [2]. The delineated related work all touch upon sub-optimal expert behavior, but they model the problem differently than I do. I use the idea of construals, maintain an infinite horizon planning scheme with a fully observable MDP, and do not recompute the policy.

3.1 Imitation Learning

In imitation learning, the agent learns a function that maps states to actions, $F : S \rightarrow A$ [6]. This is akin to the agent memorizing the correct action to take in all states, because to map states to actions, each state must be uniquely featurized. It is commonly setup as a supervised learning problem in which we model the agent as a classifier and train it to learn this function [6]. The classifier accepts a state as input and returns action scores. The classifier is trained on a dataset of featurized trajectories $X = \{\phi(s) \mid (s, a) \in \bigcup_{d_i \in D} d_i\}$ [6]. After training the classifier, the learned policy is built by deriving a probability distribution over actions from the action scores over all the states in the MDP.

3.1.1 Implementation Details

I implemented a classifier with the Python deep learning library PyTorch [7]. The classifier has a convolutional layer and a fully connected layer. The convolutional layer consists of three convolutional blocks performing 1D convolution. The first block consists of sixteen 3x1 filters, a BatchNorm layer, and a ReLU layer. The second block consists of thirty-two 3x1 filters, a BatchNorm layer, and a ReLU layer. The last block consists of sixty-four 3x1 filters, a BatchNorm layer, and a ReLU layer. The output is flattened before being passed into the final two fully connected layers, the first with a ReLU layer and the second with no extra layers. The final output is scores over the actions. See Figure 2 for an architecture visualization.

I use a {x-coordinate, y-coordinate, state color} feature representation. White is excluded as a color. To train the classifier, I used Cross Entropy Loss and Stochastic Gradient Descent (SGD) with a learning rate of 0.1, weight decay of 0.0001, and momentum of 0.9. I also decreased the learning rate by 0.1 every 10 epochs. For results, see Table 1. For the specific implementation, see Section 8.

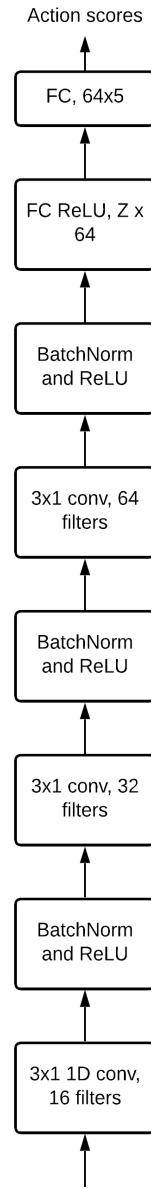


Figure 2: The Classifier Architecture

3.2 Intention Learning

Intention learning, also called inverse reinforcement learning (IRL) when working with MDPs, tries not to memorize which actions to take in which states, but

tries to learn the expert’s reward function [6]. The reward function is like an explanation of why the expert takes some action in some state, because we derive policies from the reward function by using any RL policy algorithm. Thus, if we can learn the same reward function as the expert, we will be able to compute the same policy as the expert, and the agent will have successfully learned from the expert. Furthermore, if the expert is planning optimally, then learning a reward function close to the expert’s means the agent will be able to plan optimally in any MDP that has the same features and rewards, making intention learning more generalizable. There are many IRL algorithms. In this paper, I implement the Maximum Likelihood IRL (MLIRL) algorithm [1]. In MLIRL, the reward function $r_\theta(s, a)$ is parameterized by a vector of reward weights θ in the following way:

$$r_\theta(s, a) = \theta^T \phi(s, a, s') \quad (1)$$

where ϕ returns a feature vector [1]. Each weight in θ corresponds to one feature in the MDP. The data preparation is similar to imitation learning. We have a dataset of trajectories D that the algorithm learns from. However, we are no longer training a classifier’s parameters, but rather, learning θ . MLIRL tries to find the θ that maximizes the log likelihood of seen trajectories (trajectories in D). Mathematically, the log likelihood is

$$L(D|\theta) = \prod_{i=1}^N \prod_{(s,a) \in d_i} \pi_\theta(s, a)^{w_i} = \sum_{i=1}^N \sum_{(s,a) \in d_i} w_i \log \pi_\theta(s, a) \quad (2)$$

where $L(D|\theta)$ is the log likelihood of the trajectories given the current reward weights θ and w_i is a weight encoding of the frequency of trajectory i in the dataset [1]. The solution to the MLIRL problem is the θ that maximizes $L(D|\theta)$: $\theta^* = \operatorname{argmax}_\theta L(D|\theta)$ [1]. After learning the reward weights, the learned policy is built by passing in the learned reward function into the RL algorithm used during the algorithm. The learned reward function is derived from the reward weights according to Equation (1). See Figure 7 for the full algorithm.

3.2.1 Implementation Details

In my implementation of MLIRL, I use a 1-hot vector encoding the color of the state as the feature representation. The 1-hot vector excludes white as a color. Because of this featurization, the reward weight is equivalent to the actual reward of a feature by Equation 1. π_θ is computed using MSDM’s [5] Maximum Entropy

RL algorithm. w_i was selected to be $\frac{1}{N}$ for all i where N is the number of trajectories. In the original paper, the trajectories are all unique so any duplicates are factored into w_i . I allow for duplicate trajectories to be processed, so I do not weight a trajectory during the loss calculation and reward weights update. In my implementation, w_i acts as a scaling factor, because without it, the gradients are too large. Reward weights θ were learned using PyTorch’s [7] SGD with a learning rate of 1, weight decay of 0, and momentum of 0.9. No weight decay is used because learning the reward weights correctly means the agent will be able to generalize well. The loss function is the negative log likelihood $-L(D|\theta)$ instead of the log likelihood, because PyTorch minimizes the loss. Minimizing $-L(D|\theta)$ is equivalent to maximizing $L(D|\theta)$. For results, see Table 1. For the specific implementation, see Section 8.

4 Construal

When humans plan for a task, we do not consider all details relevant to the task to come up with a plan. We lack the computational resources required to do so. Although we face this limitation, we still plan to achieve our goals by planning as optimally as possible with our given resources. This can be interpreted as humans selecting the most important details of the task to use during planning. To model this sub-optimal behavior, I use the idea of construals [4]. A construal c is a subset of all the details, or cause-effect relationships, of a task [4] that the expert considers during planning. More importantly, because we operate on gridworlds, a construal can be considered as a subset of all the features of the gridworld. Following from the logic established, experts form construals when planning for a task, because they have limited cognitive resources and cannot consider all task features [4]. Thus, they choose the most important features to still solve the task well [4]. For example, in Figure 1, there are five features: red, blue, lightgreen, lightblue, and green. A construal c from the set of all possible construals C is any subset of these five features in this gridworld. It can be the set {red, blue}, {lightblue}, {red, green, lightgreen}, or {red, blue, lightgreen, lightblue, and green} among others.

4.1 Value of Representation

How does the expert choose the best construals? Each construal has a value of representation (VOR) [4] which is a numerical value that indicates the importance of this construal. The VOR is defined in the following way:

$$VOR = U(\pi_c) - |c| \quad (3)$$

where π_c is the policy derived from construal c , $U(\pi_c)$ is the utility of π_c , and $|c|$ is the size of set c [4]. The utility of π_c can be thought of as the expected future discounted reward over the initial starting states when following the policy. A higher utility means the policy is better. Once the VOR for all construals are calculated, the best construals are the ones with the highest VOR value.

5 Construed Maximum Likelihood IRL Algorithm

5.1 Construing the Reward Function

What does it mean for an expert to disregard features not in their construal? Because I adopt the reward function representation $r_\theta(s, a) = \theta^T \phi(s, a, s')$ from MLIRL, I chose to define this as the expert disregarding and ignoring the rewards associated with those features. Thus, when an expert construes their reward function, they zero out the weights in θ of the features not in their construal. This leads to a reward function where the reward of taking an action and landing in another state changes because the reward weights of those features of the state are now 0, so they do not contribute to the reward anymore. They are ignored. Hence, the Construed Maximum Likelihood IRL algorithm (CMLIRL) adds the assumption that the expert construes their reward function in this manner to introduce construals into MLIRL [3].

In order for the reward weights to be learned, an extra assumption that the expert chooses their construals according to a differentiable probability distribution based on the VOR must be included [3]. This is so that the gradients can be computed. For example, if they only chose the best construal which is equal to a max operation over the construals, then the gradient cannot be calculated because max is not differentiable [3].

5.2 Loss Function

Starting from MLIRL, we still find the $\theta^* = \operatorname{argmax}_\theta L(D|\theta)$ to maximize the likelihood of seen trajectories [1]. However, the likelihood depends on all the construals now. This is modelled through this probability [3]:

$$P(D|\theta) = \sum_{\pi_c} P(D, \pi_c|\theta) = \sum_{\pi_c} P(D|\pi_c)P(\pi_c|\theta) = \sum_{\pi_c} \sum_{i=1}^N \sum_{(s,a) \in d_i} \log \pi_{c,\theta}(s, a) P(\pi_c|\theta) \quad (4)$$

This is the loss function used to update the reward weights. We sum over all the construed policies because the likelihood of a trajectory depends on which construed policy the expert uses which we do not know. $P(\pi_c|\theta)$ is the expert's probability distribution over construals from the algorithm assumption.

5.3 Implementation Details

Putting everything together, the algorithm is outlined below. `Construed_policies` is a mapping from construals to policies derived from those construals.

Algorithm 2 Construed Maximum Likelihood IRL Algorithm

Input: MDP, features ϕ , trajectories $\{d_1, d_2, \dots, d_N\}$, number of epochs M , learning rate α , SoftMax function σ , utility function U .

Initialize: Choose random initial reward weights θ .

```

for  $i \in [1, M]$  do
    Initialize VOR and construed_policies  $\delta$ 
    Compute set of all construals  $C$ 
    for  $c \in C$  do
        Compute  $\pi_{c,\theta}$ 
         $u = U(\pi_{c,\theta})$ 
         $\delta[c] = \pi_{c,\theta}$ 
         $VOR[c] = u - |c|$ 
    endfor
     $L = 0$ 
    for  $\pi_{c,\theta} \in \delta$  do
         $L \leftarrow L + \sum_{i=1}^N \sum_{(s,a) \in d_i} \sigma(VOR)_c \log \pi_{c,\theta}(s, a)$ 
    endfor
     $\theta \leftarrow \theta + \alpha \nabla L$ 
endfor
```

In my implementation of Algorithm 2, I imposed a further assumption that the expert uses construals of size 2. This mimics the behavior that the expert can only consider 2 details of the task. This is not a necessary assumption for the algorithm but makes it so that computing all possible construals C is computationally easier during algorithm execution. I use a 1-hot vector encoding the color of the state except white as the feature representation. I use the initial value as the utility function. The initial value is the expected future discounted reward from start states. $\pi_{c,\theta}$ is computed using MSDM’s [5] Maximum Entropy RL algorithm. Moreover, I assume the expert chooses construals according to a SoftMax distribution $\sigma(VOR)$ where $\sigma(VOR)_c$ is the probability associated with construal c . Reward weights θ were learned using PyTorch’s [7] SGD with a variable learning rate, weight decay of 0, and momentum of 0.9.

6 Experiments and Evaluations

6.1 Imitation Learning Classifier and Maximum Likelihood IRL

Evaluations are done using a policy’s initial value as the quantitative measure. The initial value is defined as the following [3]:

$$\sum_{s_0 \in S_0} P(s_0)V(s_0) \quad (5)$$

where S_0 is the set of initial states, $P(s_0)$ is the probability of starting at s_0 , and $V(s_0)$ is the expected future discounted reward at s_0 . A more positive initial value is better. We want the classifier’s and MLIRL’s learned policies’ initial values to match the expert’s initial value.

Moreover, policy graphs which are the learned policies and expert’s policies plotted on their respective gridworlds and their state value maps are generated for qualitative comparison. State value maps display the expected future discounted reward on every state using the specified policy: $V_\pi(s)$ for all $s \in S$. See Sections 10.2 to 10.5 for these results.

For MLIRL, the learned reward weights are also compared to the actual rewards. The goal is to see how well the classifier and MLIRL can learn an expert’s policy and how well their learned policies apply to novel gridworlds.

Table 1 presents the initial values of the policies that the classifier and MLIRL

learned from four different gridworlds (GW) against the initial values of the expert policy. This tests how well the classifier and MLIRL can learn from the expert. Initial values are calculated with MSDM [5]. GW1 has 5 features {lightgreen, green, red, lightblue, and blue}. GW2 and GW3 have the same 5 features {lightgreen, green, red, black, and blue}. GW4 has 8 features {lightgreen, grey, purple, magenta, yellow, black, orange, and red}. For each gridworld example, both the classifier and MLIRL algorithm were trained from the same trajectories generated from the expert who acts optimally, using batches of size 128 for 50 epochs.

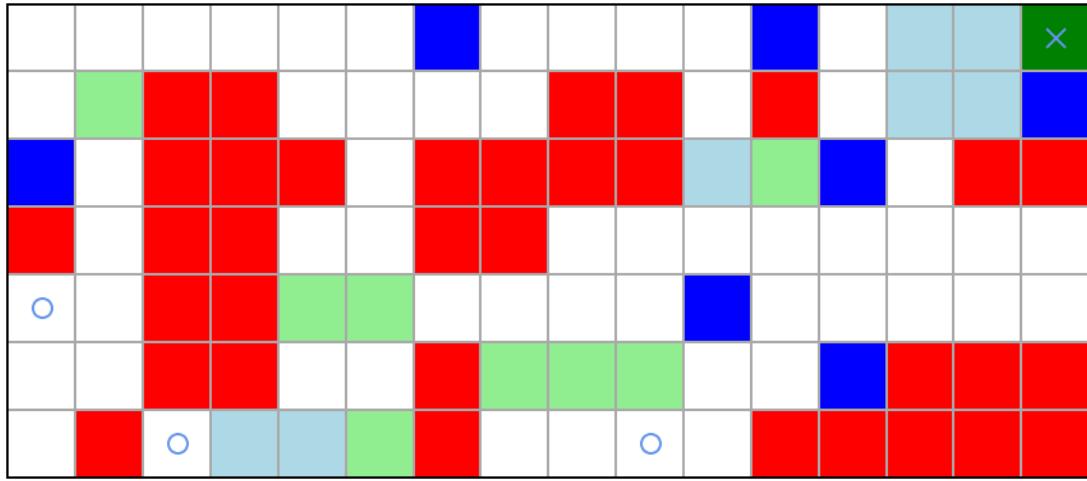


Figure 3: Gridworld 1. Rewards are: {lightgreen: 0, green: 10, red: -500, lightblue: 0, and blue: -5}. Created with MSDM [5].

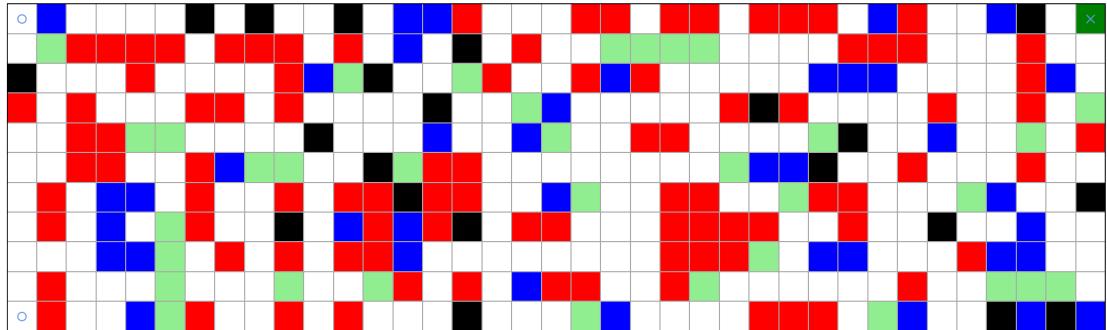


Figure 4: Gridworld 2. Rewards are: {lightgreen: 0, green: 10, red: -12, black: -5, and blue: -1}. Created with MSDM [5].

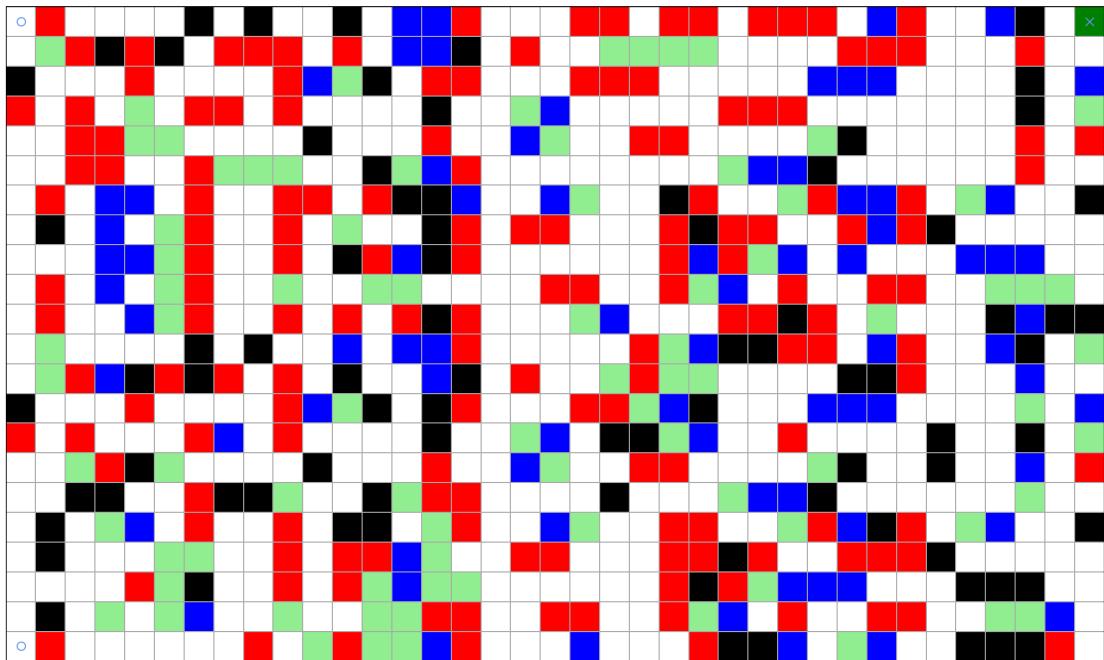


Figure 5: Gridworld 3. Rewards are: {lightgreen: 0, green: 10, red: -12, black: -5, and blue: -1}. Created with MSDM [5].

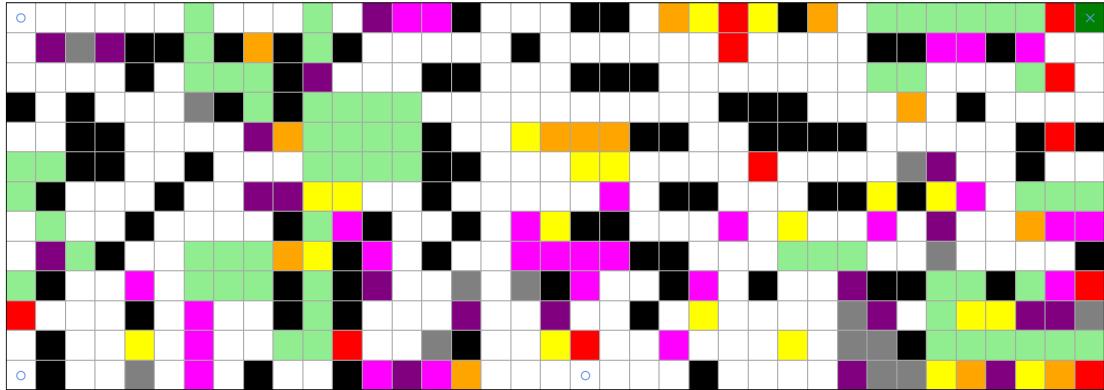


Figure 6: Gridworld 4. Rewards are: {lightgreen: 0.5, grey: -1, purple: -2, magenta: -5, yellow: -7, black: -10, orange: -13, and red: -20}. Created with MSDM [5].

Policy	10000 Trajectories				100000 Trajectories			
Classifier	-42.59	-68.92	-81.34	-63.07	-27.84	-64.24	-79.76	-
Maximum Likelihood	-38.29	-63.63	-71.33	-60.43	-38.41	-63.64	-71.22	-
Optimal Expert	-24.21	-62.72	-70.41	-58.63	-24.21	-62.72	-70.41	-

Table 1: Initial values of my implementations compared to the expert policy’s initial values. For a given cell, the initial values are for the learned policy of gridworld 1, 2, 3, and 4, respectively.

Reward Weights	GW1
Learned	-0.11, -0.06, -4.93, 0.89, -20.19
Actual	0, 0, -5, 10, -500
Reward Weights	GW2
Learned	-0.02, -1.01, -4.93, 5.45, -12.03
Actual	0, -1, -5, 10, -12
Reward Weights	GW3
Learned	0.00, -0.94, -4.95, 7.14, -11.95
Actual	0, -1, -5, 10, -12
Reward Weights	GW4
Learned	0.54, -1.05, -2.09, -4.95, -7.33, -9.83, -8.41, -17.51
Actual	0.5, -1, -2, -5, -7, -10, -13, -20

Table 2: Learned rewards compared to the actual rewards. In GW1, from left to right, the rewards correspond to lightgreen, lightblue, blue, green, red. In GW2 and GW3, the rewards correspond to lightgreen, blue, black, green, red. In GW4, the rewards correspond to lightgreen, grey, purple, magenta, yellow, black, orange, and red.

From Table 1, we see that the classifier benefits from seeing more trajectories per epoch. In each gridworld example, the initial values of the classifier’s policy became more positive by training on 100000 trajectories. This occurs because the classifier is a deep neural network, so it has 10000+ parameters that need to be learned. Seeing more new trajectories helps it learn. On the other hand, MLIRL does not benefit from seeing more trajectories. It’s initial values do not change from 10000 to 100000 trajectories. In these experiments, the MLIRL deals with 5 reward weights, because a 1-hot vector encoding on only the colors is used. It

only needs to learn 5 reward weights instead of 10000+, so it needs less data to train on for good results.

Between the classifier and MLIRL, performance is similar with MLIRL doing better than the classifier. There is an anomaly in gridworld 1 on 100000 trajectories: the classifier does better than MLIRL. This offers insight into how MLIRL works. This anomaly occurs because of the -500 reward associated with the red feature. Because for a given set of trajectories, there are multiple reward weights that can maximize its likelihood, the MLIRL will never learn the correct reward weight of -500 unless the random weights are initialized close to -500. I used a Gaussian distribution with mean 0 and standard deviation of 1 for initializing reward weights, so the MLIRL algorithm never learned a reward weight close to -500 for that feature. From Table 2, it learned a reward weight of -20.19. Thus, it thinks that its not as bad to step on a red state, when in reality it should never step on red. This leads to a more negative initial value. The classifier does not suffer from this issue, because it does not care about the weights and tries to copy the expert in every state.

From Table 2, most of the learned rewards match closely to the actual rewards. This means the MLIRL does well in learning the expert’s reward function and generates an optimal policy similar to the expert’s. This is reflected in the MLIRL initial values being close to the expert’s initial values.

Additionally, the differences between gridworlds 1, 2, and 3, are the size of the gridworld and the state arrangement. Gridworld 1 is the smallest and gridworld 3 is the largest. They have the same number of features. From the results, the classifier and MLIRL perform well even as the size increases. Meanwhile, the difference between the first three gridworlds and gridworld 4 is the number of features. Gridword 4 has 8 features and the others have 5. Again, it shows that they both still perform well as the number of features increases. We can conclude that as the number of features increases, the number of reward weights also increases, so the MLIRL algorithm should benefit from training on more trajectories.

Moving on to novel gridworlds, Table 3 shows the initial values of the learned policies applied to novel gridworlds. This tests if the learned policies can generalize to other gridworlds. Note, these novel gridworlds can only vary in size and state configuration. The number of features must remain the same as I do not deal with how to handle unseen features. Each learned policy derived from the same gridworld is tested on the same novel gridworld. Learned policies from different gridwords are tested on different novel gridworlds for a total of four novel gridworlds. The novel gridworlds are in Sections 10.2 to 10.5 as well as the policy

graphs and state value maps.

Policy	GW1	GW2	GW3	GW4
Classifier	-2241.32 -2179.14	-200.11 -208.63	-262.47 -183.63	-436.25 -
Maximum Likelihood	-48.06 -48.42	-38.18 -38.20	-61.26 -61.05	-65.24 -
Optimal Expert	-24.02	-37.75	-60.25	-64.43

Table 3: Initial values of the learned policies on new, unseen gridworlds. For a given cell, the initial values are for the learned policy of that gridworld trained with 10000 or 100000 trajectories applied to unseen gridworlds, respectively. Each policy trained from the same gridworld is applied to one novel gridworld, so the expert policy row only has one initial value.

Results demonstrate that the MLIRL is generalizable, while the classifier is not. On every novel gridworld, the classifier performs poorly. This is because the classifier does not know what to do when it sees a state-feature representation that it has never seen before. For example, using our current featurization, the classifier will not know what to output when it sees a state-feature representation with (x, y) coordinates that it has never seen before during training. In an effort to make the classifier more generalizable, I used the MLIRL 1-hot vector featurization for the classifier and trained it on Figure 3. See subsection 10.6 for the results. With this featurization, the classifier does not learn at all. It's initial value is -4084.27 compared to the expert's -24.21. This reaffirms that the classifier is memorizing which action to take and must featurize every state uniquely. With this featurization, a blue state at coordinates $(1, 1)$ and a blue state at coordinates $(4, 3)$ are treated the same. But, if the expert takes different actions at these blue states, the classifier will be confused because the gradients will push in different directions. Thus, it never learns with this featurization and cannot be made more general this way. On the other hand, the MLIRL does well on novel gridworlds, because it learned the expert's reward function well. Since a policy is derived from the reward function, it will do well on any novel gridworld that has the same features and rewards associated with those features.

6.2 Construed Maximum Likelihood IRL

Evaluations are done using a policy's initial value as the quantitative measure and policy graphs and state value maps as the qualitative measure. Learned and actual rewards are also compared. The goal is to see if CMLIRL can learn the actual rewards when learning from trajectories of an expert that construes their reward function.

Table 4 compares the initial values of CMLIRL and MLIRL and Table 5 compares the learned rewards of CMLIRL and MLIRL with the actual rewards. Trajectories are from a construed expert. Each trajectory is generated by the expert after they sample a construal according to a SoftMax distribution over their VOR. The same 10000 trajectories are used for training the CMLIRL and MLIRL for 10 epochs with batches of size 128. See Sections 10.7 to 10.10 for the policy graphs, state value maps, GW5, GW6, and GW7.

Policy	GW1	GW5	GW6	GW7
CMLIRL	-25.48	-193.63	-580.05	-407.83
MLIRL	-90.09	-130.21	-135.35	-191.22
Optimal	-24.21	-70.90	-79.26	-104.13

Table 4: Initial values of the policies that CMLIRL and MLIRL learned.

Reward Weights	GW1
CMLIRL	0.59, -2.89, -4.99, 0.16, -40.13
MLIRL	-0.09, -0.08, -4.78, -0.18, -14.09
Actual	0, 0, -5, 10, -500

Reward Weights	GW5
CMLIRL	-1.37, 0.49, -0.19, 0.31, 0, -9.65, 0.55, -11.14
MLIRL	-0.04, 0.05, -0.01, 0.04, 0.01, -9.66 -0.01 -11.17
Actual	0.5, -1, -2, -5, -7, -10, -13, -50

Reward Weights	GW6
CMLIRL	-5683.33, 58.57, 335.17, 50.06, -2159.2, 98.96
MLIRL	-14.42, 0.44, -1.19, 0.38, -1.15, -0.47
Actual	-3.8, -4, -5, -4.75, -3, -5

Reward Weights	GW7
CMLIRL	0.38, -1.82, 0.05, -2.01, 0.26, 1.13
MLIRL	-0.09, -1.96, 0.02, -1.83, 0.05, -0.27
Actual	-3.8, -4, -5, -4.75, -3, -5

Table 5: CMLIRL, MLIRL, and actual rewards compared. From left to right in each GW, the rewards correspond to the same features for all three.

The initial experiment showed promising results. On GW1, CMLIRL outperformed MLIRL and closely matched the optimal policy’s initial value. Looking

at the rewards though, it only learns a more negative reward for the last feature that is closer to -500. Subsequent experiments on GW5, GW6, and GW7 show that CMLIRL fails to learn the actual rewards and performs worse than MLIRL. Surprisingly, MLIRL performs decently.

Looking at the learned reward weights for GW5, we see that MLIRL pushes all reward weights close to 0 except for two reward weights. There are two reasons for this: MLIRL treats all ignored features as if they were white squares and the SoftMax distribution for the construals was highly skewed. In this gridworld, there was one construal that had a probability of nearly 1. This means that although the expert is construing, they are almost always picking 1 construal each time they traverse the gridworld. Because of this, they consistently visit the states with the ignored features which the MLIRL algorithm thinks then is a white square because the expert does not avoid it. Thus, this pushes the reward weights to zero. Why does CMLIRL fail to learn the true rewards? Initially, I thought that the issue was the highly skewed SoftMax distribution. If the SoftMax policy heavily leans towards a select few construals where many construals have probabilities close to 0 of being chosen, then certain features may not be updated, because they are always masked out.

I tried to make the SoftMax distribution less skewed. However, looking at GW6 where 5 construals had a probability greater than 0.1 of being chosen, CMLIRL still performs poorly. The large reward weights are due to the learning rate used for CMIRL. Finally, I removed most of the white states from GW6 to put more importance on the states with rewards associated to them and the construals, but the CMLIRL still did not perform well.

7 Conclusion

7.1 Imitation and Inverse Reinforcement Learning

In conclusion, I have showed the differences between imitation learning and IRL through implementation of a classifier and MLIRL algorithm. From the results, both perform well in learning from the expert, but MLIRL performance depends on initialization as learning the reward function is a non-convex problem, whereas the classifier can copy the expert exactly. Furthermore, MLIRL is generalizable to any MDP with the same features, and will perform almost optimally on MDPs with the same rewards as well, if it learned well from an expert acting optimally.

7.2 Construed Maximum Likelihood IRL Algorithm

To combat the problem of trying to learn the real rewards from a sub-optimal expert, I presented the CMLIRL algorithm that uses the idea of construals to model a sub-optimal expert. Results show that CMLIRL fails to solve this problem. Further investigation is required to improve and iterate upon this version of CMLIRL.

7.3 Future Work

Many researchers have already explored other ways of modelling sub-optimal behavior like those in Section 3. As mentioned, there must be more work done on this version of CMLIRL. It is also imperative to speed up policy generation from construals as calculating a policy for all possible construals is computationally expensive. Similarly, one can integrate the idea of construals into other algorithms that have other forms of sub-optimal human behavior included already to see how they will perform. Finally, researching the best ways to featurize states is also important. For example, how can we make intention learning generalize better? Overall, IRL with sub-optimal behavior is a new and growing field with potential and many avenues of to explore and build upon.

8 Code

The code for the implementations of the classifier, MLIRL, and CMLIRL is available at https://github.com/iamsamliang/IRL_IW in algorithms.py and dataset.py.

9 Honor Code

This project represents my own work, in accordance with the University regulations.

/s/ Sam Liang

References

- [1] M. Babes, V. N. Marivate, K. Subramanian, and M. L. Littman. Apprenticeship learning about multiple intentions. In *ICML*, pages 897–904, 2011. URL https://icml.cc/2011/papers/478_icmlpaper.pdf.
- [2] C. L. Baker, J. Jara-Ettingera, R. Saxe, and J. B. Tenenbaum. Rational quantitative attribution of beliefs, desires and percepts in human mentalizing. *Nat Hum Behav*, 1, 0064, 2017. URL <https://doi.org/10.1038/s41562-017-0064>.
- [3] M. Ho. email, Dec. 2021.
- [4] M. K. Ho, D. Abel, C. G. Correa, M. L. Littman, J. D. Cohen, and T. L. Griffiths. Control of mental representations in human planning. *CoRR*, abs/2105.06948, 2021. URL <https://arxiv.org/abs/2105.06948>.
- [5] M. K. Ho, C. G. Correa, and D. Ritter. Models of Sequential Decision Making (msdm), 5 2021. URL <https://github.com/markkho/msdm>.
- [6] J. MacGlashan and M. L. Littman. Between imitation and intention learning. In *IJCAI*, pages 3692–3698, 2015. URL <http://ijcai.org/Abstract/15/519>.
- [7] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019. URL <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>.
- [8] T. Zhi-Xuan, J. L. Mann, T. Silver, J. B. Tenenbaum, and V. K. Mansinghka. Online bayesian goal inference for boundedly-rational planning agents. *CoRR*, abs/2006.07532, 2020. URL <https://arxiv.org/abs/2006.07532>.

10 Appendix: Algorithms, Policy Graphs, and State Value Maps

All graphs created with MSDM [5].

10.1 Maximum Likelihood IRL Algorithm

Algorithm 1 Maximum Likelihood IRL

Input: MDP \(\mathcal{r}\), features \(\phi\), trajectories \(\{\xi_1, \dots, \xi_N\}\), trajectory weights \(\{w_1, \dots, w_N\}\), number of iterations \(M\), step size for each iteration (\(t\)) \(\alpha_t\), \(1 \leq t < M\).

Initialize: Choose random set of reward weights \(\theta_1\).

for $t = 1$ **to** M **do**

 Compute Q_{θ_t} , π_{θ_t} .

$$L = \sum_i w_i \sum_{(s,a) \in \xi} \log(\pi_{\theta_t}(s, a)).$$

$$\theta_{t+1} \leftarrow \theta_t + \alpha_t \nabla L.$$

end for

Output: Return $\theta_A = \theta_M$.

Figure 7: Maximum Likelihood IRL Algorithm. Taken from [1].

10.2 Classifier and MLIRL Results on Gridworld 1

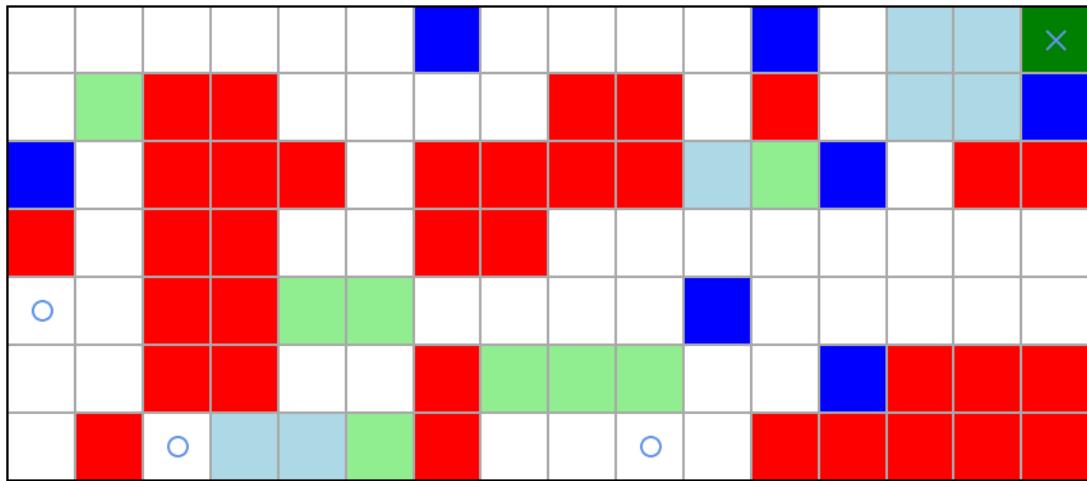


Figure 8: Gridworld 1

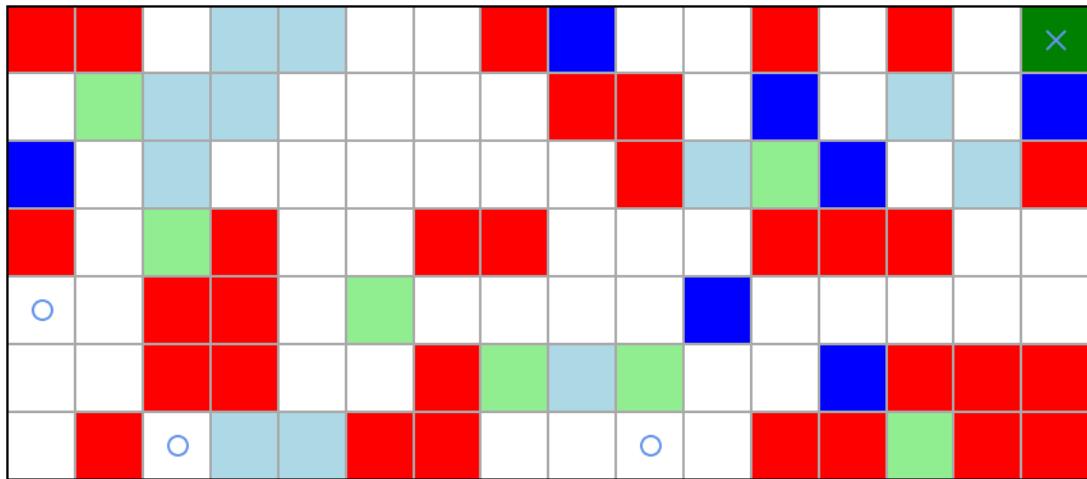


Figure 9: Novel Gridworld for Gridworld 1

Trained on 10000 trajectories with batch size of 128 for 50 epochs.

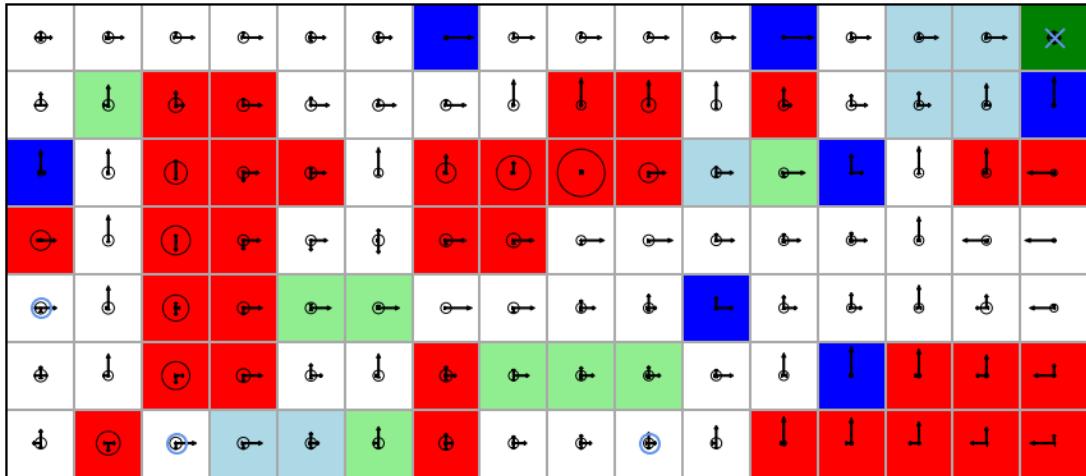


Figure 10: Classifier Policy. The arrows within a state represent the action to take: left, right, up, and down. The stay action is a circle. The length of the arrow corresponds to the probability of taking this action in this state.

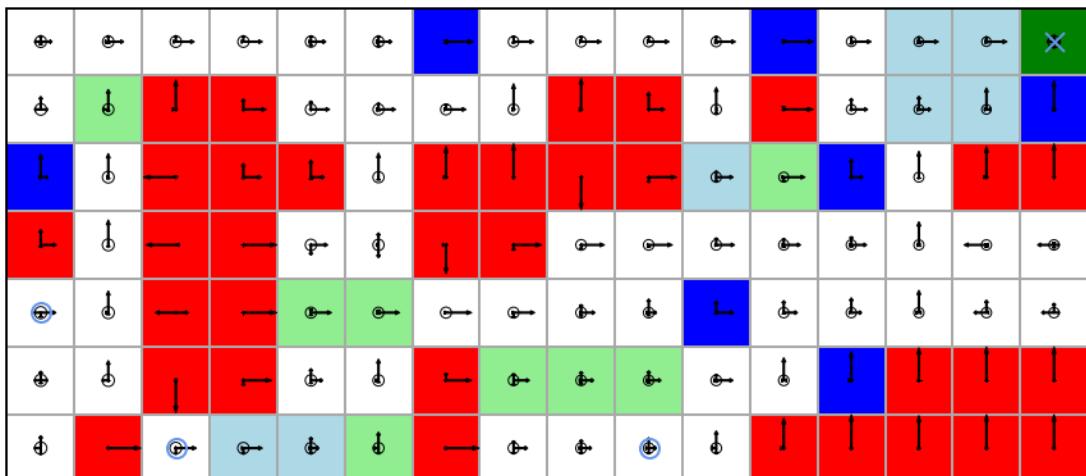


Figure 11: MLIRL Policy

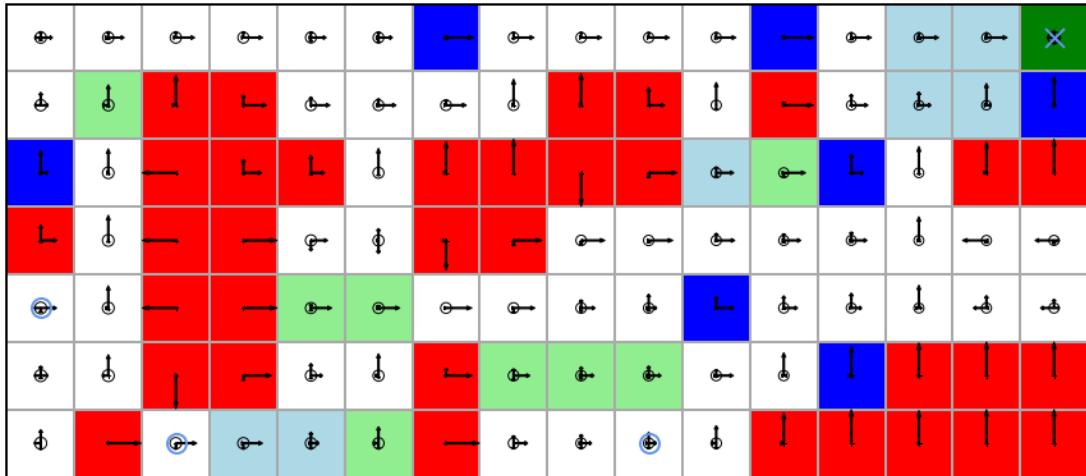


Figure 12: Expert Policy

-37.74	-35.62	-33.57	-31.28	-29.52	-26.66	-18.72	-16.46	-14.67	-10.51	-5.67	1.58	3.48	5.64	7.90	X0
-38.93	-37.60	-667.70	-210.51	-33.16	-28.34	-26.60	-21.46	-107.42	-182.28	-9.60	-108.44	1.27	2.99	3.28	6.55
-42.30	-39.90	-2238.62	-887.11	-159.06	-34.87	-344.21	-889.40	-2384.64	-324.14	-10.65	-8.20	-0.27	1.17	-91.20	-627.67
-451.12	-45.38	-2985.69	-473.76	-48.60	-37.91	-687.62	-222.13	-22.64	-12.46	-9.63	-7.06	-4.34	-1.27	-12.18	-14.26
-60.00	-52.06	-2149.19	-366.55	-38.03	-33.89	-31.76	-29.40	-20.69	-15.83	-10.57	-8.11	-6.03	-4.87	-12.46	-16.03
-83.67	-98.63	-920.58	-211.87	-38.65	-35.98	-426.81	-58.40	-24.46	-18.44	-15.52	-12.58	-9.93	-48.22	-230.56	-497.54
-117.39	-955.79	-40.18	-41.81	-40.00	-38.93	-913.55	-117.93	-32.82	-20.08	-22.16	-64.65	-156.14	-609.49	-936.36	-1324.86

Figure 13: Classifier State Value Map

-31.72	-29.28	-26.79	-24.89	-22.80	-20.94	-13.98	-11.88	-10.05	-8.25	-6.32	1.75	3.81	5.85	8.02	0.00
-33.22	-32.31	-33.01	-24.40	-22.20	-19.94	-17.39	-15.12	-14.94	-16.79	-22.52	0.02	2.50	4.09	5.57	8.12
-34.83	-33.97	-53.06	-526.14	-24.53	-23.37	-20.98	-19.07	-32.88	-17.61	-15.50	-11.94	0.32	1.29	3.18	0.35
-40.24	-35.71	-89.57	-31.75	-29.60	-26.62	-36.48	-18.43	-16.11	-14.24	-11.89	-8.49	4.86	-2.16	-28.60	-81.84
-48.06	-40.37	-231.43	-29.62	-27.70	-25.02	-22.50	-20.60	-18.17	-16.76	-11.83	-9.03	-6.73	-5.60	-21.03	-48.35
-61.59	-61.02	-57.89	-37.83	-34.09	-35.75	-23.50	-21.38	-19.46	-17.38	-14.62	-11.22	-9.23	-7.50	-23.62	-50.86
-90.70	-48.46	-46.11	-44.86	-42.62	-51.41	-25.48	-22.55	-21.07	-10.49	-18.17	-14.07	-17.21	-509.97	-526.83	-553.98

Figure 14: MLIRL State Value Map

-27.58	-25.21	-22.93	-21.25	-19.71	-18.15	-11.61	-9.66	-7.97	-6.34	-4.70	2.22	3.96	5.94	8.06	0.00
-28.69	-26.99	-24.13	-19.99	-18.27	-16.20	-13.51	-10.98	-9.07	-7.22	-6.23	1.07	2.81	4.30	5.74	8.19
-30.03	-28.49	-29.21	-520.25	-18.91	-17.50	-14.68	-11.88	-11.20	-8.08	-6.89	-5.33	1.22	2.72	4.42	2.10
-32.96	-29.55	-30.25	-20.76	-19.96	-18.78	-17.10	-11.90	-10.30	-8.57	-6.89	-4.83	-2.44	0.83	-0.86	-2.82
-32.32	-30.97	-31.66	-20.09	-19.28	-17.66	-15.94	-14.73	-12.95	-11.76	-7.72	-5.55	-3.43	-1.07	-2.40	-4.02
-34.31	-32.55	-25.99	-21.95	-20.67	-19.31	-17.02	-16.18	-14.68	-12.91	-10.64	-7.52	-5.05	-2.31	-3.38	-4.98
-35.91	-26.05	-26.33	-23.91	-22.44	-21.27	-18.40	-17.54	-16.30	-10.68	-12.78	-8.96	-11.00	-503.53	-504.34	-505.93

Figure 15: Expert State Value Map

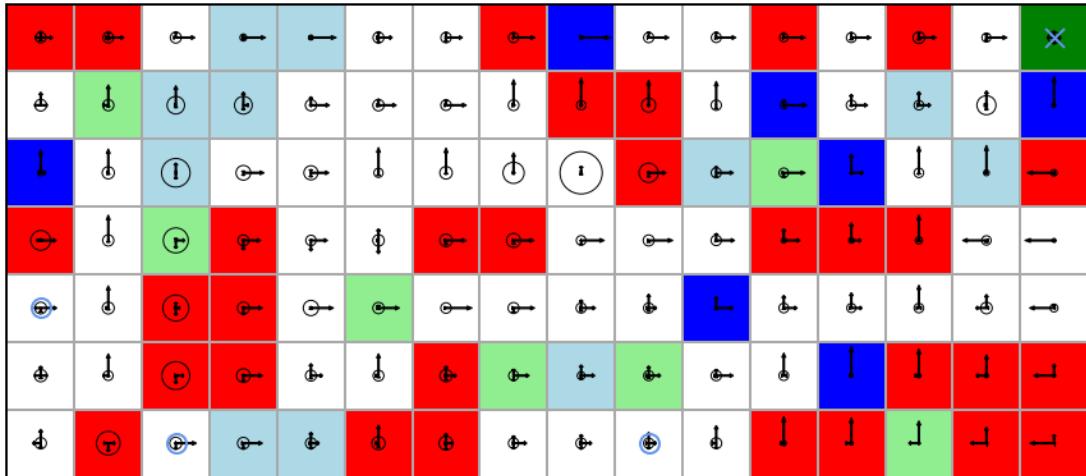


Figure 16: Classifier Policy on Novel Gridworld

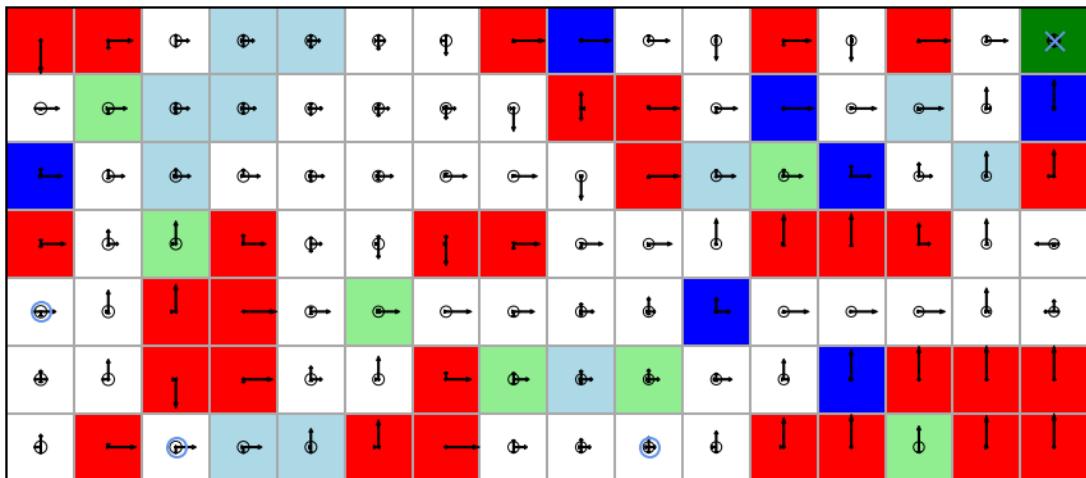


Figure 17: MLIRL Policy on Novel Gridworld

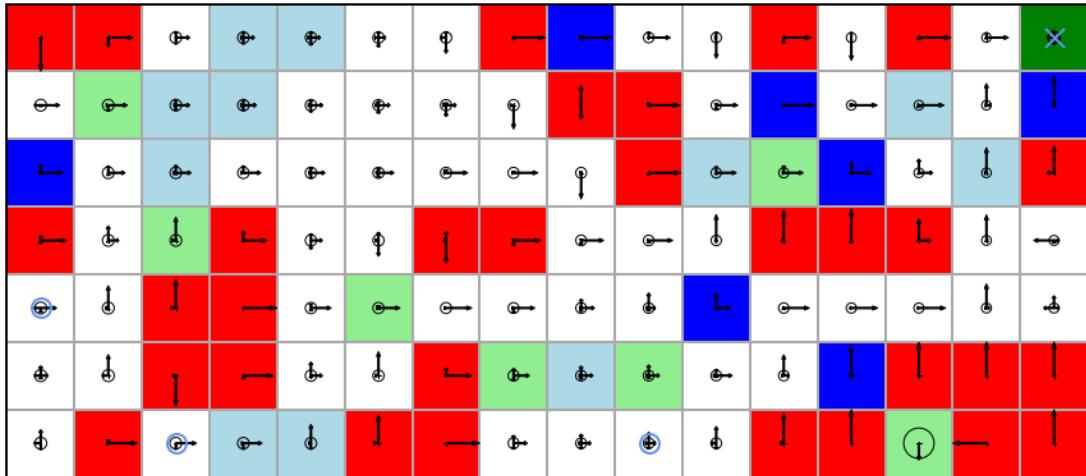


Figure 18: Expert Policy on Novel Gridworld

-4222.91	-3052.06	-2296.12	-2181.02	-2208.82	-2239.35	-2298.06	-1846.01	-1585.71	-1591.71	-1613.77	-1141.47	-911.73	-437.29	-34.08	0.00
-4281.56	-3657.48	-2383.39	-2137.39	-2199.89	-2234.86	-2282.82	-2318.55	-1687.65	-1736.82	-1532.41	-873.71	-693.67	-559.89	-116.30	0.15
-4117.64	-3631.14	-2588.44	-2083.05	-2054.40	-2161.98	-2267.18	-2274.15	-2000.74	-1488.85	-1145.09	-902.45	-672.79	-584.70	-169.22	-312.50
-3962.08	-3583.38	-3005.07	-1993.76	-1656.40	-1780.84	-2180.28	-1729.06	-1562.65	-1582.73	-1607.13	-1311.74	-992.10	-759.69	-1230.28	-1216.92
-3426.37	-3514.92	-3266.27	-1835.69	-1496.97	-1506.68	-1486.66	-1502.11	-1525.49	-1548.67	-1570.52	-1579.98	-1396.96	-1261.04	-1229.04	-1216.33
-3357.39	-3451.21	-2617.78	-1777.55	-1559.88	-1569.32	-1866.96	-1499.81	-1492.22	-1507.29	-1524.88	-1540.08	-1394.79	-1308.19	-1448.34	-1696.16
-3292.07	-3268.22	-1820.93	-1847.61	-1880.80	-1881.97	-2392.39	-1535.01	-1466.50	-1476.67	-1498.02	-1568.83	-1548.96	-1870.87	-1911.80	-2348.61

Figure 19: Classifier State Value Map on Novel Gridworld

-75.78	-68.71	-67.44	-67.21	-67.83	-70.39	-78.75	-37.44	-28.56	-27.09	-24.04	-36.95	-40.71	5.71	8.14	0.00
-67.78	-66.62	-65.56	-64.45	-63.53	-64.07	-65.16	-66.71	-99.12	-14.00	-11.03	-3.52	-1.72	2.07	5.62	8.13
-67.31	-65.42	-63.62	-60.92	-57.49	-58.13	-59.69	-57.73	-53.90	-13.02	-10.46	-8.33	-1.25	1.23	3.41	1.02
-70.27	-67.51	-68.43	-52.45	-45.77	-44.03	-46.91	-25.48	-21.76	-17.97	-13.96	-12.10	-16.88	-0.57	0.88	-10.79
-78.01	-70.00	-80.99	-37.54	-35.42	-31.14	-27.38	-25.38	-22.86	-20.82	-16.32	-15.22	-11.78	-8.88	-1.84	-8.26
-75.52	-74.12	-64.82	-41.55	-38.75	-40.04	-28.34	-26.28	-24.44	-22.56	-20.60	-17.73	-15.45	-12.71	-3.68	-10.99
-81.13	-52.42	-46.03	-44.38	-41.41	-56.63	-30.17	-27.52	-26.19	-20.85	-23.97	-20.71	-33.82	-513.91	-506.24	-514.37

Figure 20: MLIRL State Value Map on Novel Gridworld

-30.80	-28.63	-27.71	-26.26	-24.35	-22.25	-20.10	-16.46	-10.53	-8.98	-7.47	-0.86	1.13	7.01	8.30	0.00
-30.10	-28.47	-27.09	-25.14	-22.83	-20.28	-17.65	-15.49	-15.22	-6.99	-5.84	0.95	2.65	4.31	5.90	8.21
-30.13	-28.23	-26.33	-23.85	-21.48	-18.81	-15.88	-13.95	-12.51	-7.60	-6.50	-5.05	1.44	2.97	4.31	2.32
-30.92	-29.04	-27.63	-22.88	-21.16	-19.68	-17.33	-12.87	-11.30	-9.52	-7.90	-6.06	4.44	1.81	2.76	0.90
-32.07	-30.49	-28.79	-20.66	-19.86	-18.10	-16.37	-15.16	-13.42	-12.05	-7.38	-3.32	-1.66	-0.24	1.08	-0.51
-33.88	-32.10	-26.59	-21.97	-20.87	-19.44	-17.29	-16.42	-14.84	-12.78	-9.75	-5.83	-3.18	-1.32	0.07	-1.50
-35.50	-26.99	-26.08	-23.95	-22.44	-20.71	-18.49	-17.66	-16.34	-10.91	-12.19	-7.49	-9.15	-100.00	-100.00	-502.49

Figure 21: Expert State Value Map on Novel Gridworld

Trained on 100000 trajectories with batch size of 128 for 50 epochs.

\oplus	$\oplus\rightarrow$	$\oplus\rightarrow$	$\oplus\rightarrow$	$\oplus\rightarrow$	$\oplus\rightarrow$	$\oplus\rightarrow$	\rightarrow	$\oplus\rightarrow$	$\oplus\rightarrow$	$\oplus\rightarrow$	$\oplus\rightarrow$	\rightarrow	$\oplus\rightarrow$	$\oplus\rightarrow$	$\oplus\rightarrow$	\times
\ominus	\downarrow	\leftarrow	\leftarrow	\leftarrow	\leftarrow	\leftarrow	\rightarrow	\uparrow	\uparrow	\uparrow	\uparrow	\rightarrow	\leftarrow	\leftarrow	\leftarrow	\downarrow
L	\uparrow	\leftarrow	\leftarrow	\leftarrow	\leftarrow	\leftarrow	\uparrow	\uparrow	\uparrow	\circ	\circ	\leftarrow	\leftarrow	L	\uparrow	\leftarrow
\odot	\uparrow	\leftarrow	\odot	\leftarrow	φ	ϕ	\odot	\odot	\leftarrow	\leftarrow	\leftarrow	\leftarrow	\leftarrow	\leftarrow	\uparrow	\leftarrow
$\ominus\odot$	\uparrow	\leftarrow	\odot	\leftarrow	\odot	\odot	\leftarrow	\leftarrow	\leftarrow	\leftarrow	\leftarrow	L	\leftarrow	\leftarrow	\uparrow	\leftarrow
\oplus	\uparrow	\leftarrow	\odot	\leftarrow	\leftarrow	\leftarrow	\leftarrow	\leftarrow	\leftarrow	\leftarrow	\leftarrow	\leftarrow	\leftarrow	\leftarrow	\leftarrow	\leftarrow
\odot	\odot	\odot	\odot	\odot	\odot	\odot	\odot	\odot	\odot	\odot	\odot	\odot	\odot	\odot	\odot	\odot

Figure 22: Classifier Policy

\oplus	$\oplus\rightarrow$	$\oplus\rightarrow$	$\oplus\rightarrow$	$\oplus\rightarrow$	$\oplus\rightarrow$	$\oplus\rightarrow$	\rightarrow	$\oplus\rightarrow$	$\oplus\rightarrow$	$\oplus\rightarrow$	$\oplus\rightarrow$	\rightarrow	$\oplus\rightarrow$	$\oplus\rightarrow$	$\oplus\rightarrow$	\times
\ominus	\downarrow	\leftarrow	L	L	\leftarrow	\leftarrow	\uparrow	\uparrow	\uparrow	L	\uparrow	\rightarrow	\leftarrow	\leftarrow	\leftarrow	\downarrow
L	\uparrow	\leftarrow	L	L	\leftarrow	\leftarrow	\uparrow	\uparrow	\uparrow	\uparrow	\uparrow	\leftarrow	\leftarrow	L	\uparrow	\leftarrow
L	\uparrow	\leftarrow	\leftarrow	\leftarrow	φ	ϕ	\downarrow	\leftarrow	\leftarrow	\leftarrow	\leftarrow	\leftarrow	\leftarrow	\leftarrow	\uparrow	\leftarrow
$\ominus\odot$	\uparrow	\leftarrow	\leftarrow	\leftarrow	\odot	\odot	\odot	\odot	\odot	\odot	\odot	L	\leftarrow	\leftarrow	\uparrow	\leftarrow
\oplus	\uparrow	\leftarrow	L	\leftarrow	\leftarrow	\leftarrow	\leftarrow	\leftarrow	\leftarrow	\leftarrow	\leftarrow	\leftarrow	\leftarrow	\leftarrow	\leftarrow	\leftarrow
\odot	\leftarrow	\odot	\odot	\odot	\odot	\odot	\odot	\odot	\odot	\odot	\odot	\odot	\odot	\odot	\odot	\odot

Figure 23: MLIRL Policy

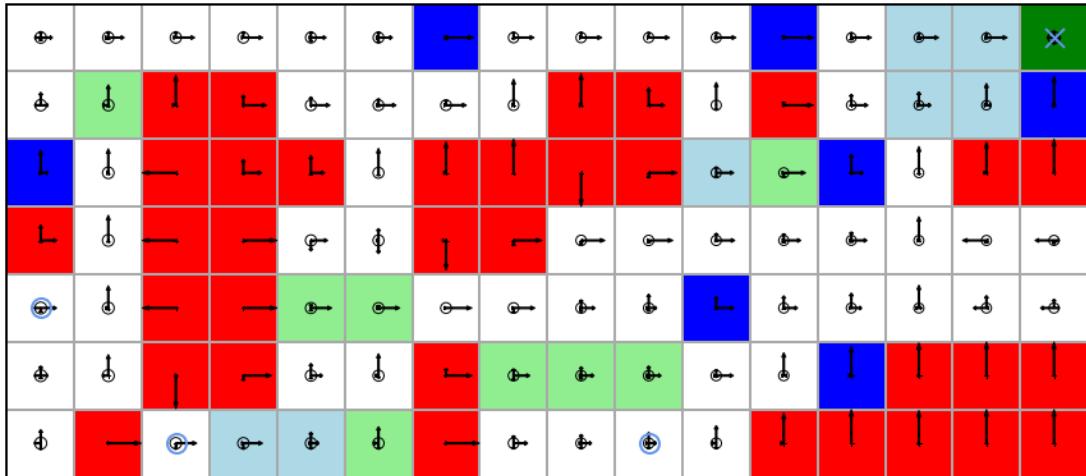


Figure 24: Expert Policy

-29.39	-27.04	-24.77	-22.80	-21.16	-19.17	-12.32	-10.29	-8.58	-6.80	-4.90	2.09	3.92	5.91	8.04	☒0
-30.56	-28.97	-502.85	-138.59	-21.26	-17.94	-15.21	-12.02	-163.71	-254.24	-7.22	-153.66	2.56	4.10	5.39	8.15
-32.39	-30.74	-1281.54	-807.22	-155.07	-21.74	-141.42	-174.12	-1384.59	-496.77	-7.88	-5.91	0.89	2.30	-156.64	-668.17
-307.24	-32.27	-1696.69	-357.22	-27.75	-23.07	-832.34	-210.86	-12.37	-9.60	-7.63	-5.50	-3.13	-0.08	-10.16	-11.61
-37.07	-34.55	-1881.43	-508.73	-23.30	-20.79	-18.89	-17.23	-14.47	-12.79	-8.56	-6.40	-4.39	-2.47	-7.61	-9.09
-40.32	-39.45	-1128.44	-274.14	-24.76	-22.74	-359.88	-21.74	-16.47	-14.16	-11.87	-8.85	-6.50	-50.53	-202.54	-465.23
-45.67	-903.91	-30.17	-28.07	-26.41	-25.15	-906.26	-28.08	-18.89	-16.28	-14.56	-75.58	-212.42	-662.23	-982.34	-1381.69

Figure 25: Classifier State Value Map

-31.78	-29.34	-26.84	-24.94	-22.83	-20.97	-13.99	-11.90	-10.07	-8.26	-6.33	1.78	3.83	5.88	8.04	0.00
-33.28	-32.39	-33.14	-24.46	-22.25	-19.98	-17.43	-15.17	-15.02	-16.92	-22.78	0.03	2.52	4.12	5.58	8.13
-34.89	-34.05	-53.50	-526.21	-24.60	-23.44	-21.06	-19.15	-33.11	-17.78	-15.65	-12.05	0.33	1.28	3.17	0.33
-40.34	-35.80	-89.83	-31.90	-29.74	-26.71	-36.78	-18.51	-16.18	-14.30	-11.95	-8.54	4.90	-2.21	-29.23	-83.83
-49.30	-40.43	-230.27	-29.74	-27.79	-25.10	-22.57	-20.67	-18.23	-16.82	-11.88	-9.08	-6.79	-5.70	-21.49	-49.53
-61.30	-60.64	-58.52	-37.99	-34.24	-35.95	-23.58	-21.44	-19.51	-17.43	-14.66	-11.27	-9.31	-7.62	-24.10	-52.05
-89.65	-48.76	-46.39	-45.13	-42.87	-51.89	-25.57	-22.60	-21.13	-10.95	-18.24	-14.14	-17.33	-510.10	-527.35	-555.20

Figure 26: MLIRL State Value Map

-27.58	-25.21	-22.93	-21.25	-19.71	-18.15	-11.61	-9.66	-7.97	-6.34	-4.70	2.22	3.96	5.94	8.06	0.00
-28.69	-26.99	-24.13	-19.99	-18.27	-16.20	-13.51	-10.98	-9.07	-7.22	-6.23	1.07	2.81	4.30	5.74	8.19
-30.03	-28.49	-29.21	-520.25	-18.91	-17.50	-14.68	-11.88	-11.20	-8.08	-6.89	-5.33	1.22	2.72	4.42	2.10
-32.96	-29.55	-30.25	-20.76	-19.96	-18.78	-17.10	-11.90	-10.30	-8.57	-6.89	-4.83	-2.44	0.83	-0.86	-2.82
-32.32	-30.97	-31.66	-20.09	-19.28	-17.66	-15.94	-14.73	-12.95	-11.76	-7.72	-5.55	-3.43	-1.07	-2.40	-4.02
-34.31	-32.55	-25.99	-21.95	-20.67	-19.31	-17.02	-16.18	-14.68	-12.91	-10.64	-7.52	-5.05	-2.31	-3.38	-4.98
-35.91	-26.05	-26.33	-23.91	-22.44	-21.27	-18.40	-17.54	-16.30	-10.68	-12.78	-8.96	-11.00	-503.53	-504.34	-505.93

Figure 27: Expert State Value Map

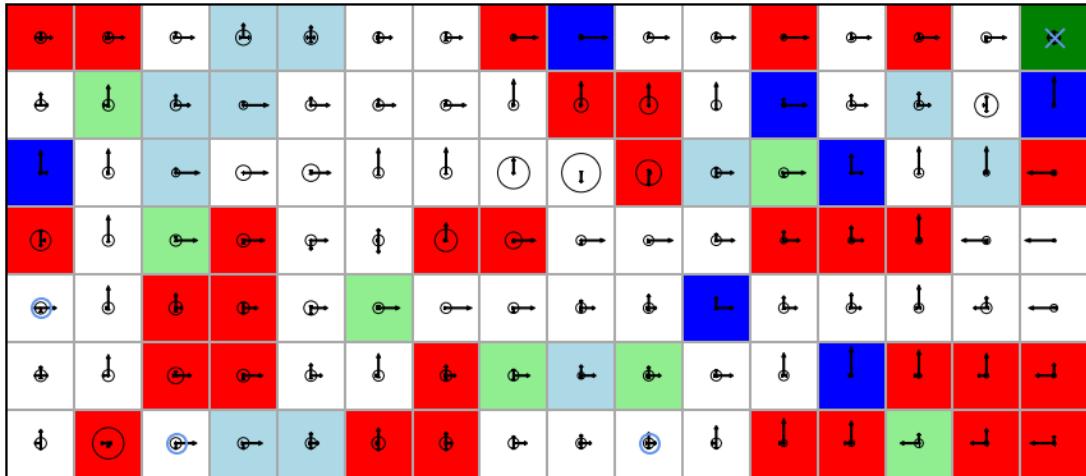


Figure 28: Classifier Policy on Novel Gridworld

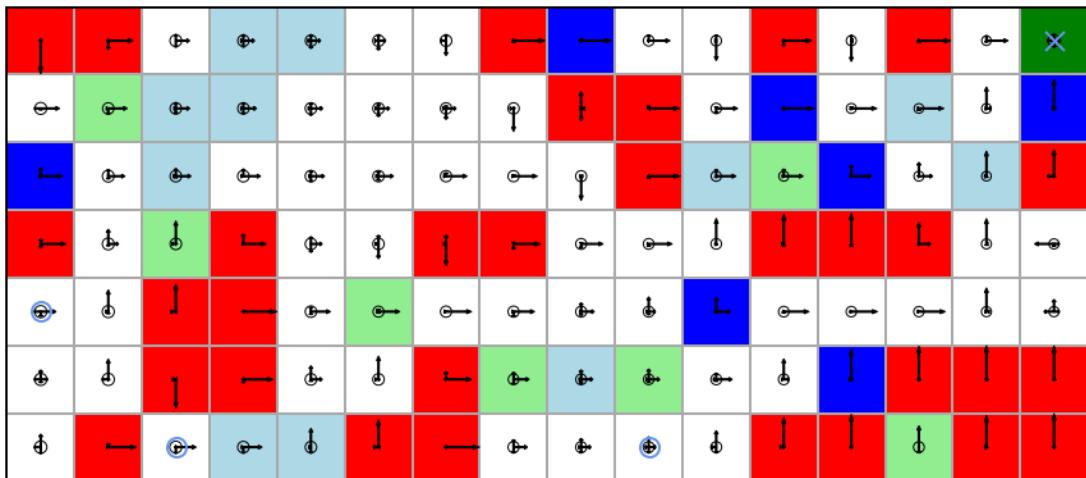


Figure 29: MLIRL Policy on Novel Gridworld

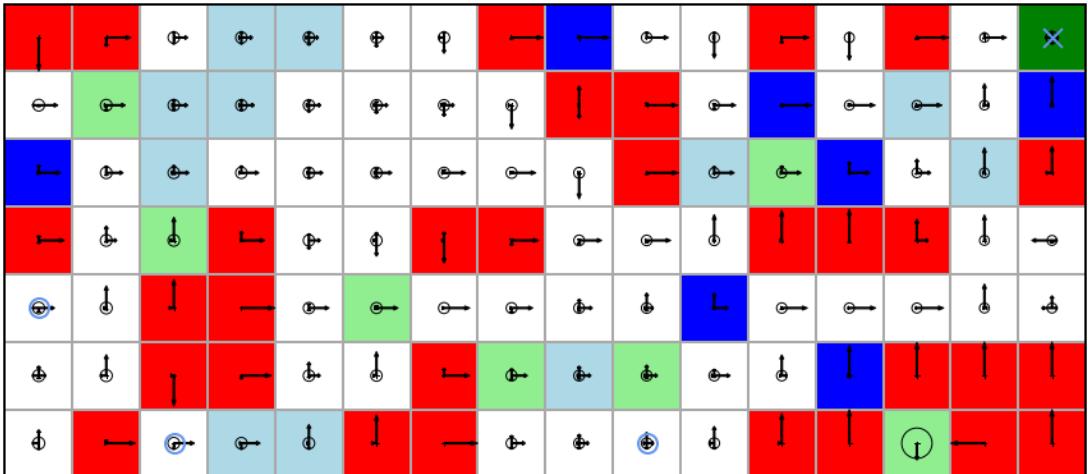


Figure 30: Expert Policy on Novel Gridworld

-3911.06	-2714.02	-2004.50	-1885.27	-1902.38	-2007.62	-2070.16	-1612.09	-1448.13	-1436.18	-1456.37	-981.61	-843.92	-363.40	-92.82	X0
-3940.77	-3322.27	-2023.23	-1923.32	-1952.87	-2008.06	-2055.67	-2089.13	-1630.21	-1663.87	-1377.97	-851.89	-677.70	-578.27	-254.78	-5.65
-3748.73	-3303.55	-1889.33	-1883.06	-1905.69	-1952.54	-2062.01	-2051.96	-1655.48	-1975.06	-1098.15	-887.73	-672.08	-605.26	-312.06	-445.96
-3681.25	-3256.88	-2315.82	-1880.13	-1622.19	-1699.39	-2260.58	-1769.85	-1600.94	-1627.06	-1652.52	-1393.14	-1058.11	-759.60	-1235.75	-1222.55
-3110.55	-3193.50	-2654.59	-2023.88	-1536.26	-1526.00	-1519.99	-1537.65	-1566.97	-1594.15	-1620.25	-1633.48	-1442.01	-1266.70	-1235.22	-1220.49
-3045.03	-3124.51	-2752.56	-1877.43	-1599.78	-1599.96	-1847.51	-1513.40	-1535.17	-1553.46	-1573.72	-1590.98	-1441.55	-1328.29	-1430.84	-1673.63
-2977.99	-3483.98	-1907.29	-1935.24	-1970.25	-2058.61	-2425.71	-1490.67	-1501.19	-1509.56	-1541.24	-1630.51	-1661.09	-2027.83	-2006.71	-2438.62

Figure 31: Classifier State Value Map on Novel Gridworld

-76.44	-69.26	-67.98	-67.78	-68.45	-71.08	-79.56	-37.89	-28.97	-27.51	-24.44	-37.77	41.62	5.69	8.14	0.00
-68.28	-67.13	-66.06	-64.98	-64.10	-64.69	-65.82	-67.43	-100.48	-14.14	-11.12	-3.61	-1.80	2.04	5.62	8.13
-67.79	-65.90	-64.09	-61.42	-58.03	-58.71	-60.28	-58.30	-54.43	-13.12	-10.52	-8.38	-1.29	1.22	3.41	1.01
-70.80	-68.04	-69.04	-52.92	-46.22	-44.49	-47.38	-25.70	-21.95	-18.13	-14.07	-12.22	-17.12	-0.60	0.86	-10.98
-74.09	-70.55	-81.98	-37.84	-35.70	-31.35	-27.56	-25.54	-23.01	-20.97	-16.46	-15.37	-11.93	-9.01	-1.88	-8.39
-76.10	-74.71	-65.69	-41.87	-39.05	-40.36	-28.52	-26.42	-24.57	-22.69	-20.73	-17.88	-15.64	-12.89	-3.74	-11.15
-81.73	-52.88	-46.07	-44.72	-41.71	-57.26	-30.34	-27.67	-26.33	-26.99	-24.13	-20.88	-34.17	-514.15	-506.33	-514.57

Figure 32: MLIRL State Value Map on Novel Gridworld

-30.80	-28.63	-27.71	-26.26	-24.35	-22.25	-20.10	-16.46	-10.53	-8.98	-7.47	-0.86	1.13	7.01	8.30	0.00
-30.10	-28.47	-27.09	-25.14	-22.83	-20.28	-17.65	-15.49	-15.22	-6.99	-5.84	0.95	2.65	4.31	5.90	8.21
-30.13	-28.23	-26.33	-23.85	-21.48	-18.81	-15.88	-13.95	-12.51	-7.60	-6.50	-5.05	1.44	2.97	4.31	2.32
-30.92	-29.04	-27.63	-22.88	-21.16	-19.68	-17.33	-12.87	-11.30	-9.52	-7.90	-6.06	4.44	1.81	2.76	0.90
-32.07	-30.49	-28.79	-20.66	-19.86	-18.10	-16.37	-15.16	-13.42	-12.05	-7.38	-3.32	-1.66	-0.24	1.08	-0.51
-33.88	-32.10	-26.59	-21.97	-20.87	-19.44	-17.29	-16.42	-14.84	-12.78	-9.75	-5.83	-3.18	-1.32	0.07	-1.50
-35.50	-26.99	-26.08	-23.95	-22.44	-20.71	-18.49	-17.66	-16.34	-16.91	-12.19	-7.49	-9.15	-100.00	-100.00	-502.49

Figure 33: Expert State Value Map on Novel Gridworld

10.3 Classifier and MLIRL Results on Gridworld 2

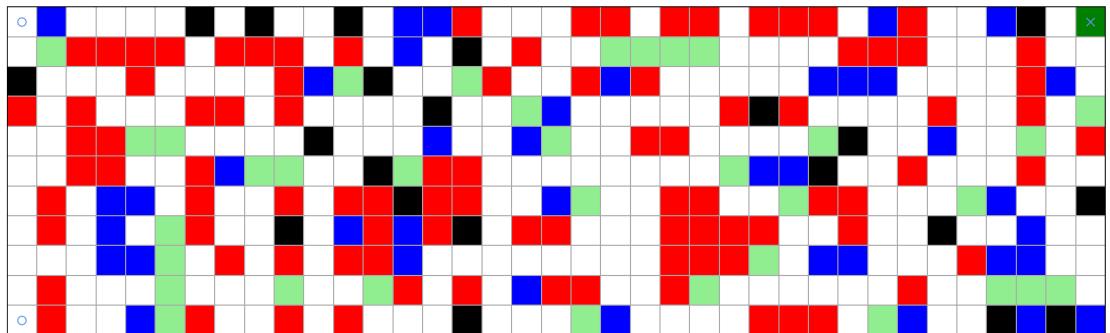


Figure 34: Gridworld 2

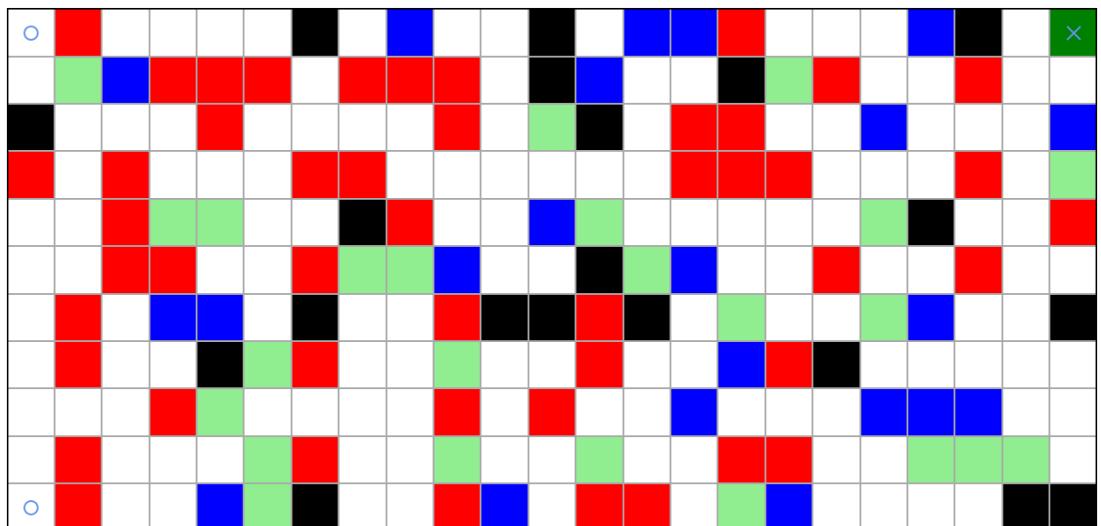


Figure 35: Novel Gridworld for Gridworld 2

Trained on 10000 trajectories with batch size of 128 for 50 epochs.

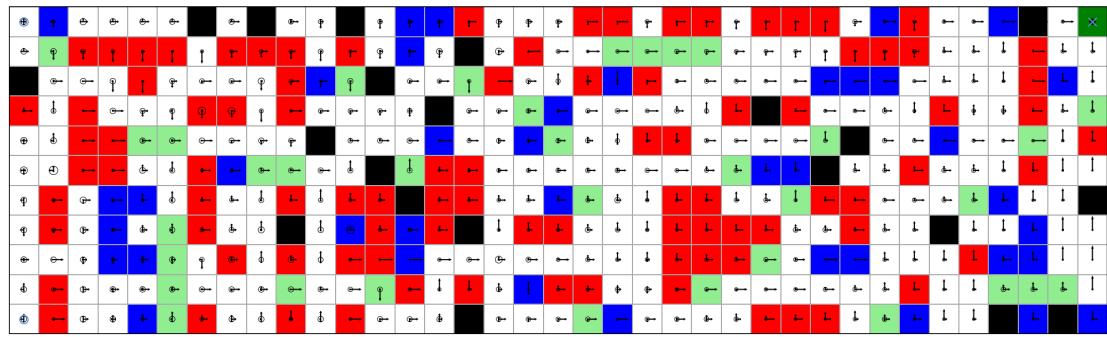


Figure 36: Classifier Policy

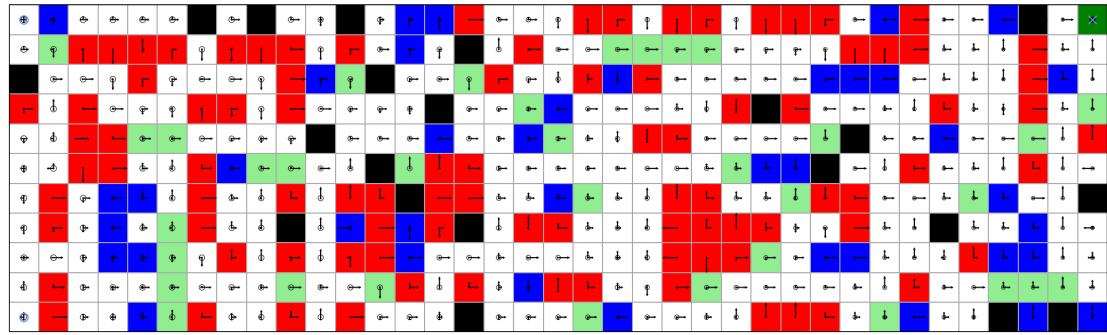


Figure 37: MLIRL Policy

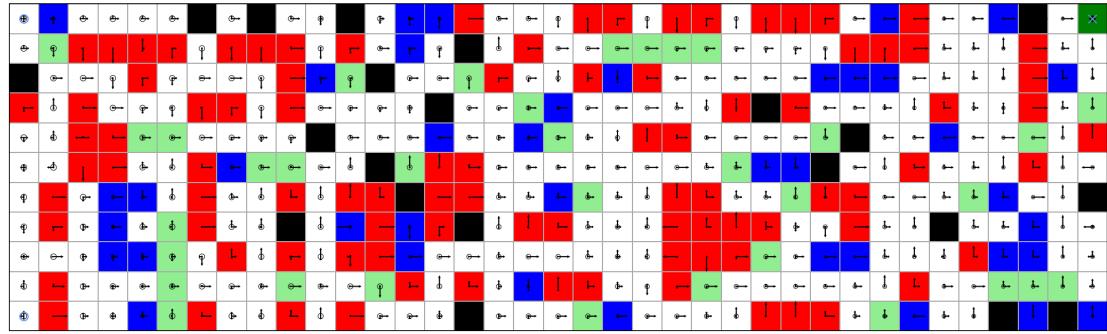


Figure 38: Expert Policy



Figure 39: Classifier State Value Map

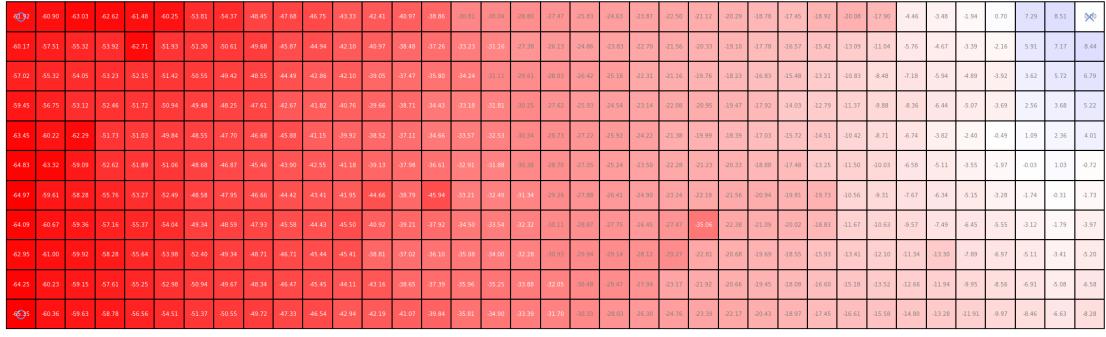


Figure 40: MLIRL State Value Map



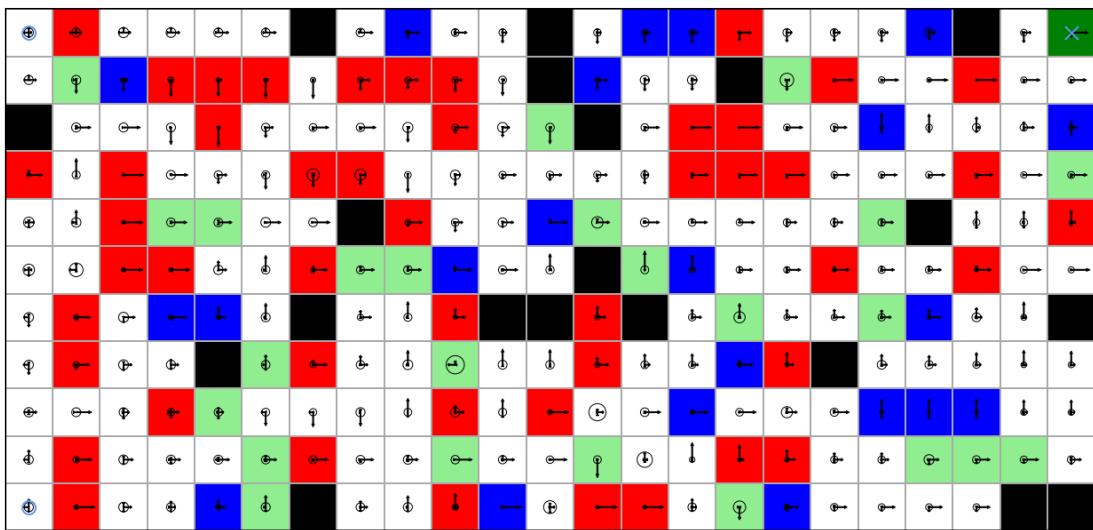


Figure 42: Classifier Policy on Novel Gridworld

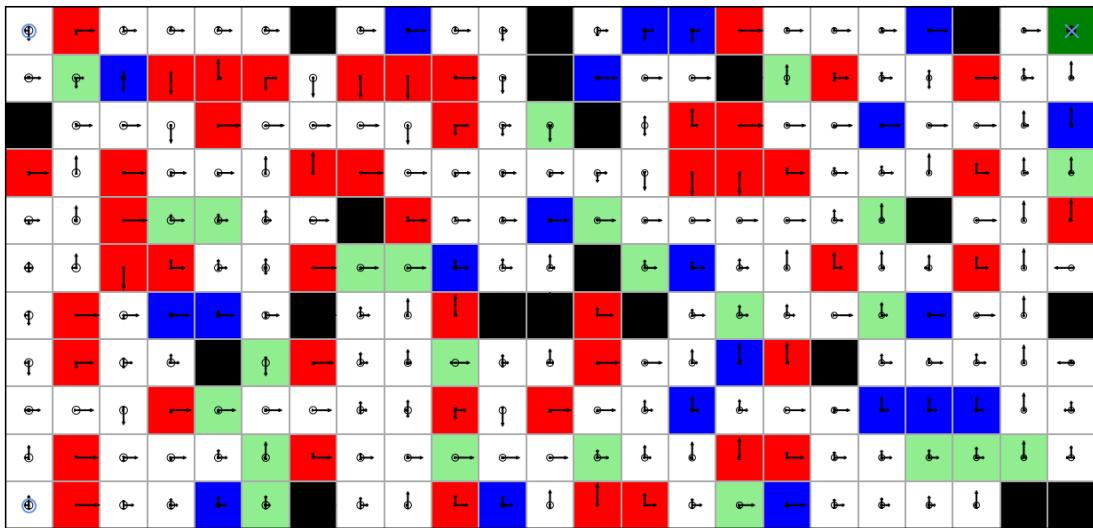


Figure 43: MLIRL Policy on Novel Gridworld

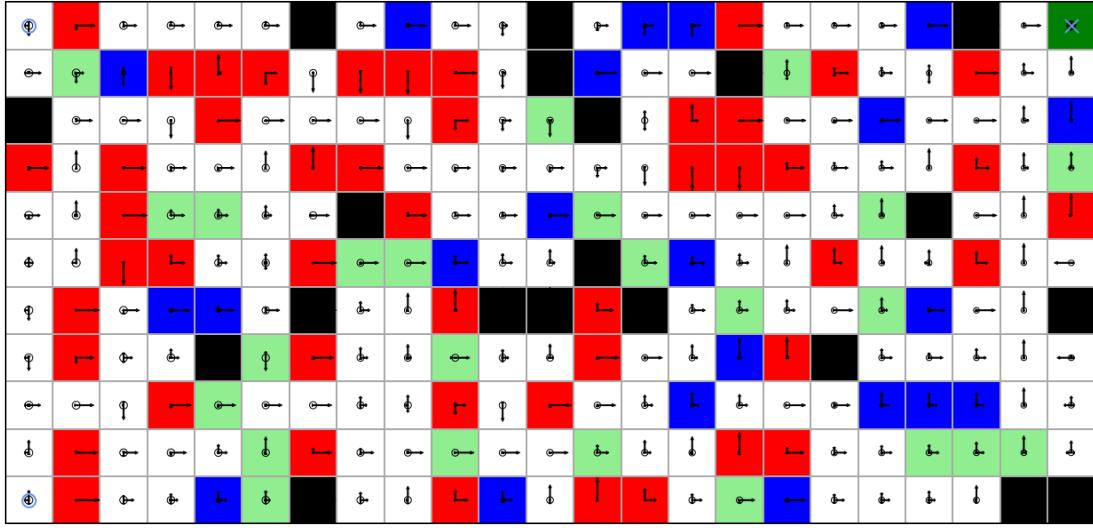


Figure 44: Expert Policy on Novel Gridworld

-200.42	208.95	-205.39	397.93	197.11	-198.83	196.26	-195.82	394.89	-186.86	-183.72	185.16	-184.75	380.98	-176.90	-160.22	152.24	-141.86	-329.08	-115.83	-85.35	-69.92	⊗
-193.65	192.01	-194.64	-195.11	-207.52	-197.86	199.24	-206.28	-202.19	-192.70	-182.57	-184.60	-185.99	-183.38	-179.91	-163.03	157.18	-138.46	-336.27	-135.51	-123.76	122.58	-121.67
-189.50	189.42	-190.55	-191.92	-194.64	-196.76	-198.95	-197.89	-196.38	-184.32	-182.05	-183.07	-185.60	-186.68	-175.58	-160.16	158.38	-154.07	-447.11	-146.17	-147.15	141.15	-145.21
-191.67	188.55	-193.01	-193.24	-194.63	-196.18	204.46	-203.22	-196.91	-183.22	-182.27	-183.66	-185.45	-188.69	-194.66	-185.89	175.04	-174.30	-374.31	-175.37	-164.26	162.24	-163.26
-192.31	188.91	-196.38	-195.26	-195.84	-197.92	199.80	-196.29	-185.71	-183.07	-182.24	-182.72	-183.59	-184.58	-185.59	-182.75	179.99	-177.37	-375.67	-173.46	-173.65	159.85	-164.83
-194.29	193.94	-207.10	-195.21	-194.76	-196.08	-190.74	-186.80	-184.62	-182.03	-182.65	-183.01	-182.89	-183.58	-183.96	-181.94	-183.26	-174.13	-372.16	-171.60	-159.10	152.82	-152.28
-198.61	-204.06	-200.59	-197.44	-195.32	-195.66	-190.80	-184.12	-183.56	-189.36	-190.25	-192.15	-189.60	-183.53	-182.23	-181.11	178.65	-175.90	-369.43	-166.89	-164.61	154.80	-161.31
-202.44	208.55	-206.98	-207.49	-200.79	-202.30	-190.44	-184.82	-183.20	-188.79	-194.74	-199.39	-194.83	-188.22	-185.34	-188.67	180.50	-172.76	-372.10	-169.65	-165.74	158.40	-166.33
-205.92	211.11	-213.80	-214.07	-213.17	-216.39	221.00	-204.89	-189.70	-200.14	-202.23	-209.74	-209.66	-198.27	-198.29	-201.06	202.46	-204.04	-206.70	-219.79	-212.35	184.17	-188.39
-202.63	213.63	-211.83	-214.65	-216.41	-220.09	212.95	-209.47	-209.96	-221.16	-224.05	-233.27	-236.05	-207.07	-200.41	-213.52	235.68	-257.96	-280.58	-320.58	-320.57	312.84	-302.24
-198.81	211.75	-210.12	-212.22	-214.96	-217.44	212.74	-207.52	-208.10	-223.09	-231.04	-232.44	-226.17	-213.53	-212.41	-237.15	-306.81	-333.98	-351.03	-363.04	-371.94	-376.17	-375.91

Figure 45: Classifier State Value Map on Novel Gridworld

-36.31	-35.87	-35.77	-34.95	-33.72	-32.46	-26.09	-26.41	-24.25	-22.93	-21.66	-19.68	-18.59	-16.84	-16.06	-5.73	-4.74	-3.34	-1.90	0.69	7.29	8.51	0.00
-34.22	-32.57	-32.18	-30.44	-34.60	-22.56	-21.69	-21.08	-19.82	-19.58	-18.36	-17.22	-15.42	-14.01	-12.69	-6.34	-5.28	-3.11	-1.79	-0.19	5.84	7.17	8.46
-31.66	-30.25	-28.84	-27.52	-22.66	-21.67	-20.59	-19.43	-18.34	-16.97	-16.07	-15.08	-13.89	-13.24	-14.39	-3.87	-2.85	-1.43	0.99	2.69	4.19	5.48	6.80
-33.02	-31.58	-27.19	-26.28	-24.80	-23.18	-22.07	-18.14	-17.21	-16.12	-14.78	-13.16	-11.57	-10.41	-8.47	-7.51	-3.31	-2.06	-0.55	1.30	2.93	3.93	4.20
-35.74	-33.65	-28.46	-27.71	-27.06	-26.26	-27.76	-21.73	-17.17	-16.00	-14.53	-11.83	-10.27	-8.90	-7.40	-6.01	-4.66	-3.36	-2.29	0.16	1.37	2.61	3.00
-37.87	-35.66	-35.58	-29.58	-29.22	-28.87	-21.31	-20.43	-19.24	-17.04	-15.91	-15.01	-11.47	-10.54	-8.74	-7.39	-6.11	-4.63	-3.95	-5.90	0.25	1.29	-0.34
-39.43	-35.35	-34.59	-32.34	-30.07	-28.94	-22.61	-21.70	-20.61	-19.33	-17.06	-16.03	-17.37	-11.45	-10.29	-8.77	-7.64	-6.81	-5.48	-3.65	-1.49	-0.06	-1.42
-38.59	-36.31	-35.56	-35.04	-29.69	-28.71	-24.41	-23.25	-22.61	-24.20	-23.42	-23.31	-14.42	-13.35	-11.81	-10.19	-9.20	-7.55	-6.28	-4.71	-2.89	-1.58	-3.73
-37.24	-35.53	-34.47	-29.58	-28.31	-27.09	-25.73	-24.68	-23.86	-22.67	-21.74	-17.39	-16.28	-15.06	-13.40	-12.73	-11.96	-10.74	-7.98	-6.32	-4.35	-3.23	-4.99
-38.83	-33.87	-32.96	-31.34	-29.77	-28.64	-25.51	-24.40	-22.76	-21.15	-20.07	-18.69	-17.47	-16.48	-15.96	-14.20	-12.71	-11.56	-10.04	-8.34	-6.56	-4.91	-6.30
-40.35	-34.99	-34.28	-33.37	-31.74	-30.93	-26.52	-25.52	-24.32	-22.68	-21.32	-20.10	-18.46	-17.85	-17.55	-16.60	-14.19	-12.81	-11.42	-9.88	-8.35	-6.36	-7.56

Figure 46: MLIRL State Value Map on Novel Gridworld

-36.34	35.63	-35.56	-34.76	33.55	-32.32	25.97	-26.31	-24.15	-22.85	-21.59	-19.58	-18.51	-16.79	-16.03	-5.63	-4.65	-3.25	-1.83	0.73	7.31	8.53	0.00
-33.87	-32.27	-31.88	-30.15	-34.42	-22.28	-21.42	-20.83	-19.59	-19.40	-18.19	-17.09	-15.30	-13.90	-12.58	-6.24	-5.18	-3.05	-1.75	-0.16	5.88	7.20	8.49
-31.33	-29.94	-28.55	-27.25	-22.38	-21.40	-20.32	-19.18	-18.09	-16.74	-15.85	-14.87	-13.77	-13.13	-14.31	-3.76	-2.74	-1.32	1.08	2.75	4.24	5.52	6.84
-32.68	-31.28	-26.91	-26.00	-24.52	-22.89	-21.83	-17.89	-16.96	-15.88	-14.55	-12.96	-11.38	-10.24	-8.31	-7.40	-3.21	-1.96	-0.46	1.38	2.98	3.97	4.25
-35.37	-33.37	-28.23	-27.49	-26.84	-26.06	-27.51	-21.45	-16.90	-15.73	-14.29	-11.63	-10.09	-8.74	-7.24	-5.88	-4.54	-3.25	-2.19	0.23	1.42	2.67	3.06
-37.44	-35.32	-35.25	-29.36	-29.02	-28.70	-21.05	-20.17	-18.99	-16.79	-15.67	-14.79	-11.31	-10.38	-8.58	-7.25	-5.97	-4.51	-3.82	-5.81	0.32	1.35	-0.28
-38.98	-34.99	-34.25	-32.03	-29.79	-28.65	-22.31	-21.41	-20.34	-19.06	-16.81	-15.80	-17.20	-11.28	-10.13	-8.62	-7.50	-6.66	-5.34	-3.52	-1.40	0.02	-1.33
-38.09	-35.92	-35.21	-34.70	-29.37	-28.41	-24.08	-22.92	-22.29	-23.96	-23.19	-23.10	-14.22	-13.16	-11.63	-10.02	-9.04	-7.41	-6.14	-4.58	-2.77	-1.47	-3.59
-36.75	-35.10	-34.05	-29.20	-27.96	-26.75	-25.39	-24.35	-23.54	-22.42	-21.49	-17.18	-16.07	-14.86	-13.21	-12.54	-11.78	-10.57	-7.84	-6.19	4.22	-3.09	-4.83
-38.32	-33.43	-32.52	-30.94	-29.41	-28.29	-25.17	-24.07	-22.47	-20.91	-19.82	-18.46	-17.26	-16.27	-15.74	-13.99	-12.51	-11.36	-9.85	-8.17	-6.41	-4.75	-6.13
-40.06	-34.57	-33.87	-32.96	-31.36	-30.57	-26.18	-25.19	-23.99	-22.43	-21.07	-19.86	-18.25	-17.64	-17.33	-16.37	-13.97	-12.59	-11.20	-9.68	-8.17	-6.18	-7.38

Figure 47: Expert State Value Map on Novel Gridworld

Trained on 100000 trajectories with batch size of 128 for 50 epochs.

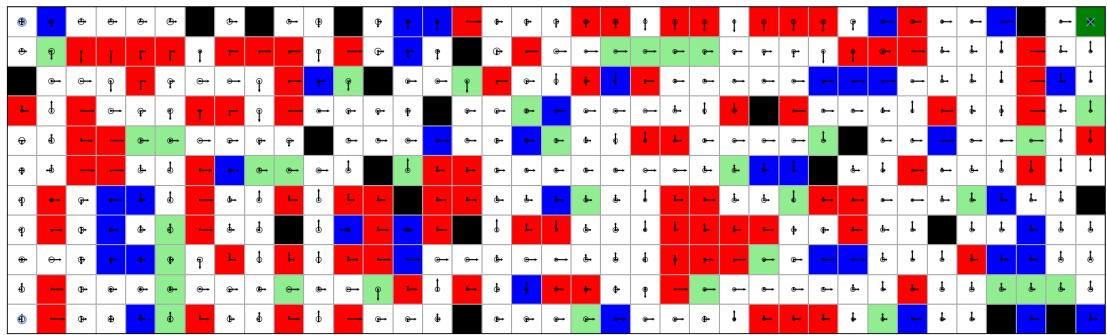


Figure 48: Classifier Policy

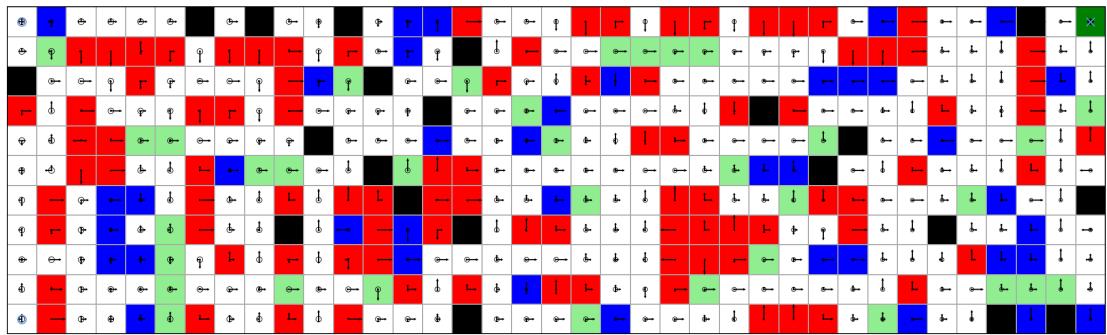


Figure 49: MLIRL Policy

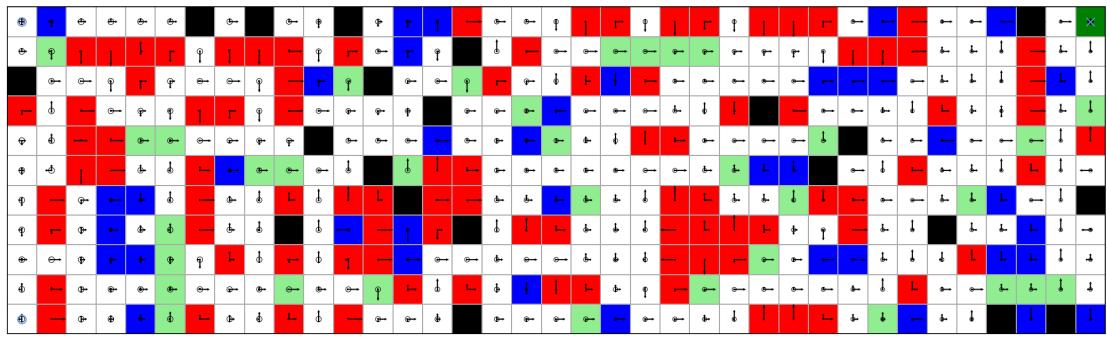


Figure 50: Expert Policy



Figure 51: Classifier State Value Map

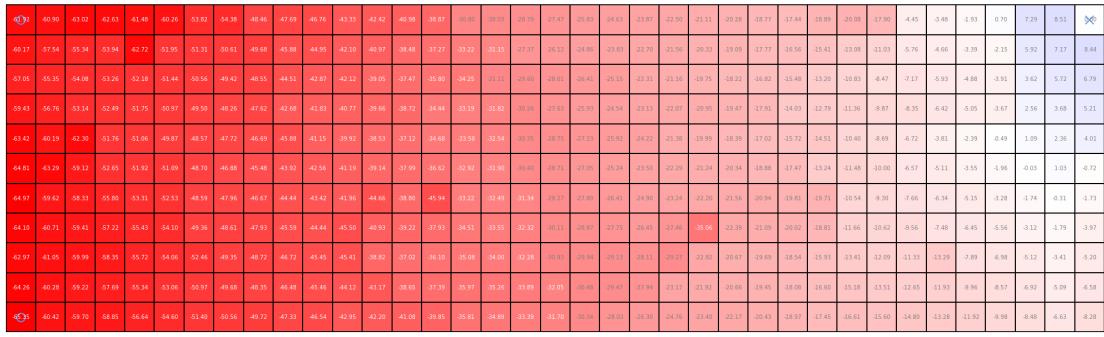


Figure 52: MLIRL State Value Map

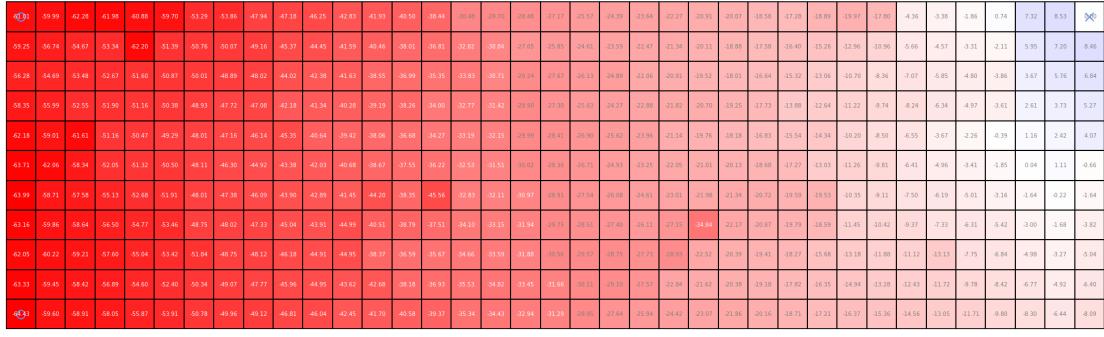


Figure 53: Expert State Value Map

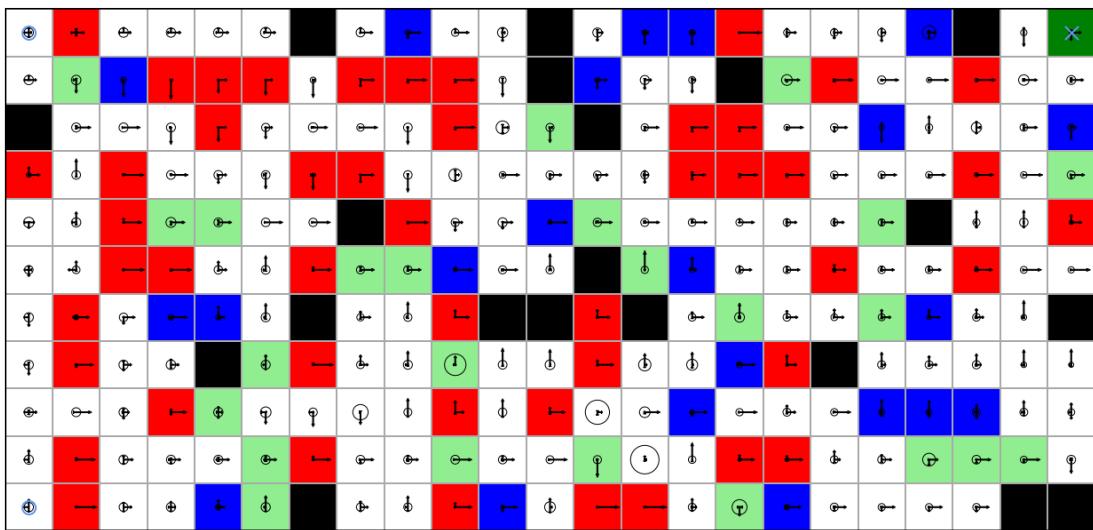


Figure 54: Classifier Policy on Novel Gridworld

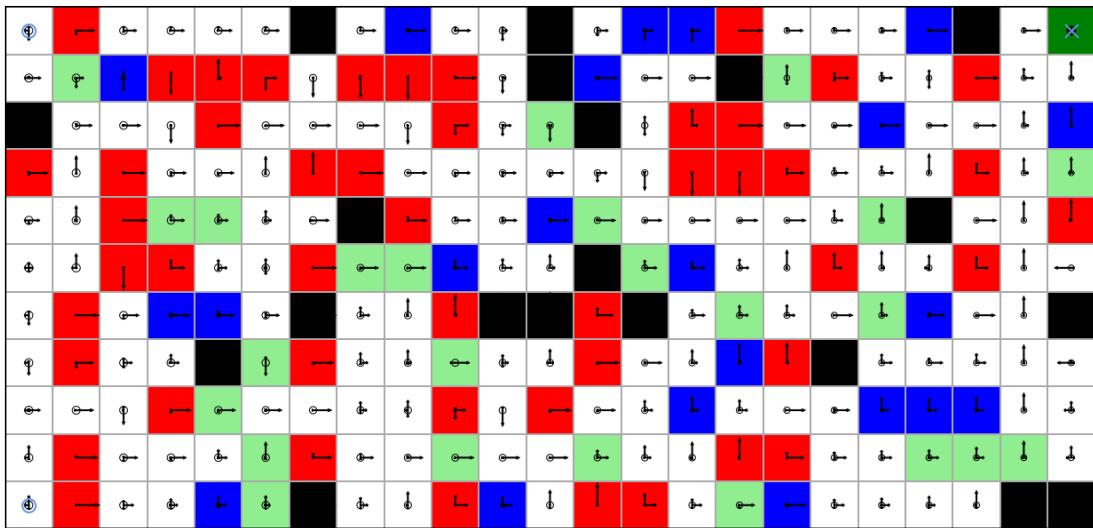


Figure 55: MLIRL Policy on Novel Gridworld

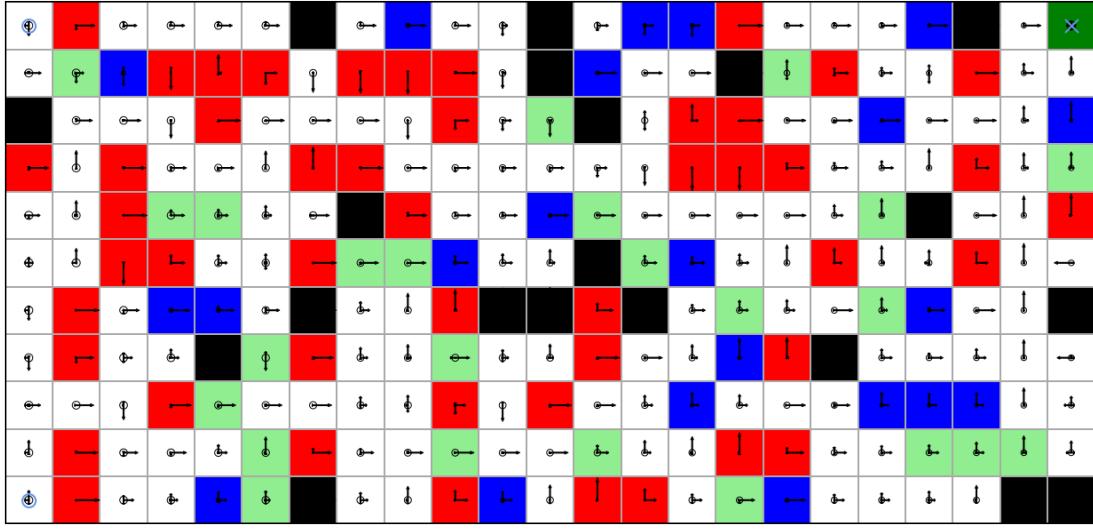


Figure 56: Expert Policy on Novel Gridworld

-2070	214.40	-214.66	-210.48	211.53	-214.02	211.77	-212.03	212.30	206.11	-207.17	216.91	-217.84	-215.00	204.51	-147.37	146.41	-139.89	-131.24	-128.86	-122.76	108.89	⊗
-206.94	207.55	-208.02	-211.80	-223.58	-212.77	213.59	-222.13	-217.14	-206.59	-205.74	-218.09	-220.95	-219.69	-215.99	-172.98	159.37	-138.89	-131.94	-128.88	-116.51	110.40	-101.12
-206.01	-206.96	-208.75	-210.45	-212.70	-213.27	-212.42	-213.08	-214.24	-206.03	-206.52	-209.81	-220.15	-225.89	-218.41	-202.18	-195.63	-193.85	-191.26	-160.49	-176.66	-187.26	-197.39
-207.94	205.38	-211.10	-212.08	-213.75	-215.47	219.67	-217.35	-215.65	-208.76	-207.37	-209.06	-212.04	-217.35	-223.04	-216.27	-206.63	-208.36	-211.93	-216.61	-208.07	-212.21	-213.49
-201.64	203.13	-215.93	-215.28	-215.39	-218.05	-219.94	-216.70	-205.17	-205.51	-205.32	-206.36	-207.30	-208.37	-209.04	-205.81	-201.48	-199.96	-200.15	-201.50	-199.29	-186.62	-209.78
-202.31	203.86	-224.47	-213.48	-214.03	-215.64	-208.86	-207.79	-206.37	-204.63	-205.65	-206.23	-206.18	-207.10	-204.74	-195.62	-194.34	-183.86	-175.58	-168.69	-146.24	132.20	-130.76
-205.58	-214.08	-214.81	-214.15	-213.61	-214.14	-206.72	-204.61	-204.83	-214.00	-215.58	-215.50	-209.91	-203.82	-196.09	-193.53	-184.42	-179.04	-165.80	-156.66	-148.01	-134.48	-157.54
-209.50	-217.17	-217.81	-220.04	-215.83	-216.41	-204.12	-203.87	-203.88	-219.46	-219.12	-219.71	-207.22	-203.51	-195.12	-192.30	-180.70	-169.02	-167.21	-158.48	-151.89	-141.74	-165.83
-212.99	218.33	-221.19	-220.72	-222.02	-223.40	-227.68	-215.55	-206.75	-219.91	-221.02	-215.75	-212.45	-198.20	-197.63	-198.34	-198.79	-196.01	-195.71	-205.15	-219.44	191.15	-225.33
-209.57	217.48	-218.11	-220.79	-222.97	-226.40	219.07	-217.87	-219.15	-227.81	-230.71	-236.29	-238.74	-204.37	-200.04	-236.56	-236.18	-242.49	-280.17	-317.14	-321.25	-321.71	-321.77
-206.57	-216.18	-218.16	-218.22	-221.11	-223.60	-222.52	-215.32	-216.62	-231.79	-233.57	-234.54	-229.14	-217.80	-218.95	-257.55	-298.39	-310.80	-323.75	-331.01	-336.24	-336.67	-335.03

Figure 57: Classifier State Value Map on Novel Gridworld

-36.35	-35.89	-35.80	-34.97	-33.75	-32.49	-26.13	-26.44	-24.27	-22.95	-21.68	-19.67	-18.58	-16.83	-16.03	-5.72	-4.74	-3.34	-1.90	0.69	7.29	8.51	0.00
-34.26	-32.62	-32.22	-30.48	-34.62	-22.57	-21.70	-21.10	-19.84	-19.60	-18.38	-17.23	-15.43	-14.01	-12.69	-6.34	-5.28	-3.12	-1.80	-0.19	5.84	7.17	8.46
-31.71	-30.29	-28.88	-27.56	-22.67	-21.68	-20.60	-19.44	-18.35	-16.99	-16.09	-15.10	-13.91	-13.25	-14.38	-3.88	-2.85	-1.43	0.99	2.69	4.19	5.48	6.80
-33.06	-31.62	-27.23	-26.33	-24.83	-23.19	-22.08	-18.15	-17.22	-16.13	-14.79	-13.18	-11.58	-10.42	-8.47	-7.50	-3.32	-2.06	-0.55	1.30	2.93	3.92	4.20
-35.79	-33.70	-28.52	-27.77	-27.12	-26.32	-27.77	-21.73	-17.18	-16.00	-14.54	-11.84	-10.28	-8.91	-7.40	-6.02	-4.66	-3.37	-2.30	0.15	1.36	2.61	3.00
-37.92	-35.72	-35.60	-29.63	-29.28	-28.92	-21.33	-20.45	-19.26	-17.05	-15.92	-15.02	-11.49	-10.56	-8.75	-7.40	-6.11	-4.63	-3.95	-5.90	0.25	1.28	-0.35
-39.46	-35.37	-34.61	-32.36	-30.10	-28.96	-22.62	-21.71	-20.63	-19.34	-17.07	-16.04	-17.38	-11.46	-10.30	-8.78	-7.65	-6.82	-5.48	-3.65	-1.50	-0.06	-1.42
-38.60	-36.32	-35.58	-35.06	-29.71	-28.73	-24.42	-23.25	-22.62	-24.21	-23.43	-23.30	-14.43	-13.36	-11.82	-10.20	-9.21	-7.56	-6.29	-4.72	-2.89	-1.59	-3.73
-37.24	-35.55	-34.49	-29.60	-28.33	-27.10	-25.73	-24.69	-23.87	-22.68	-21.74	-17.41	-16.29	-15.07	-13.41	-12.74	-11.97	-10.76	-8.00	-6.34	-4.36	-3.23	-4.99
-38.83	-33.89	-32.98	-31.36	-29.80	-28.66	-25.51	-24.40	-22.77	-21.17	-20.08	-18.70	-17.49	-16.50	-15.97	-14.21	-12.72	-11.57	-10.05	-8.35	-6.57	-4.91	-6.30
-40.35	-35.01	-34.30	-33.39	-31.76	-30.95	-26.53	-25.53	-24.32	-22.70	-21.33	-20.11	-18.48	-17.87	-17.57	-16.61	-14.20	-12.81	-11.43	-9.89	-8.37	-6.38	-7.57

Figure 58: MLIRL State Value Map on Novel Gridworld

-36.34	35.63	-35.56	-34.76	33.55	-32.32	25.97	-26.31	-24.15	-22.85	-21.59	19.58	-18.51	16.79	-16.03	-5.63	-4.65	-3.25	-1.83	0.73	7.31	8.53	0.00
-33.87	-32.27	-31.88	-30.15	-34.42	-22.28	-21.42	-20.83	-19.59	-19.40	-18.19	-17.09	-15.30	-13.90	-12.58	-6.24	-5.18	-3.05	-1.75	-0.16	5.88	7.20	8.49
-31.33	-29.94	-28.55	-27.25	-22.38	-21.40	-20.32	-19.18	-18.09	-16.74	-15.85	-14.87	-13.77	-13.13	-14.31	-3.76	-2.74	-1.32	1.08	2.75	4.24	5.52	6.84
-32.68	-31.28	-26.91	-26.00	-24.52	-22.89	-21.83	-17.89	-16.96	-15.88	-14.55	-12.96	-11.38	-10.24	-8.31	-7.40	-3.21	-1.96	-0.46	1.38	2.98	3.97	4.25
-35.37	-33.37	-28.23	-27.49	-26.84	-26.06	-27.51	-21.45	-16.90	-15.73	-14.29	-11.63	-10.09	-8.74	-7.24	-5.88	-4.54	-3.25	-2.19	0.23	1.42	2.67	3.06
-37.44	-35.32	-35.25	-29.36	-29.02	-28.70	-21.05	-20.17	-18.99	-16.79	-15.67	-14.79	-11.31	-10.38	-8.58	-7.25	-5.97	-4.51	-3.82	-5.81	0.32	1.35	-0.28
-38.98	-34.99	-34.25	-32.03	-29.79	-28.65	-22.31	-21.41	-20.34	-19.06	-16.81	-15.80	-17.20	-11.28	-10.13	-8.62	-7.50	-6.66	-5.34	-3.52	-1.40	0.02	-1.33
-38.09	-35.92	-35.21	-34.70	-29.37	-28.41	-24.08	-22.92	-22.29	-23.96	-23.19	-23.10	-14.22	-13.16	-11.63	-10.02	-9.04	-7.41	-6.14	-4.58	-2.77	-1.47	-3.59
-36.75	-35.10	-34.05	-29.20	-27.96	-26.75	-25.39	-24.35	-23.54	-22.42	-21.49	-17.18	-16.07	-14.86	-13.21	-12.54	-11.78	-10.57	-7.84	-6.19	4.22	-3.09	4.83
-38.32	-33.43	-32.52	-30.94	-29.41	-28.29	-25.17	-24.07	-22.47	-20.91	-19.82	-18.46	-17.26	-16.27	-15.74	-13.99	-12.51	-11.36	-9.85	-8.17	-6.41	-4.75	-6.13
-40.06	-34.57	-33.87	-32.96	-31.36	-30.57	-26.18	-25.19	-23.99	-22.43	-21.07	-19.86	-18.25	-17.64	-17.33	-16.37	-13.97	-12.59	-11.20	-9.68	-8.17	-6.18	-7.38

Figure 59: Expert State Value Map on Novel Gridworld

10.4 Classifier and MLIRL Results on Gridworld 3

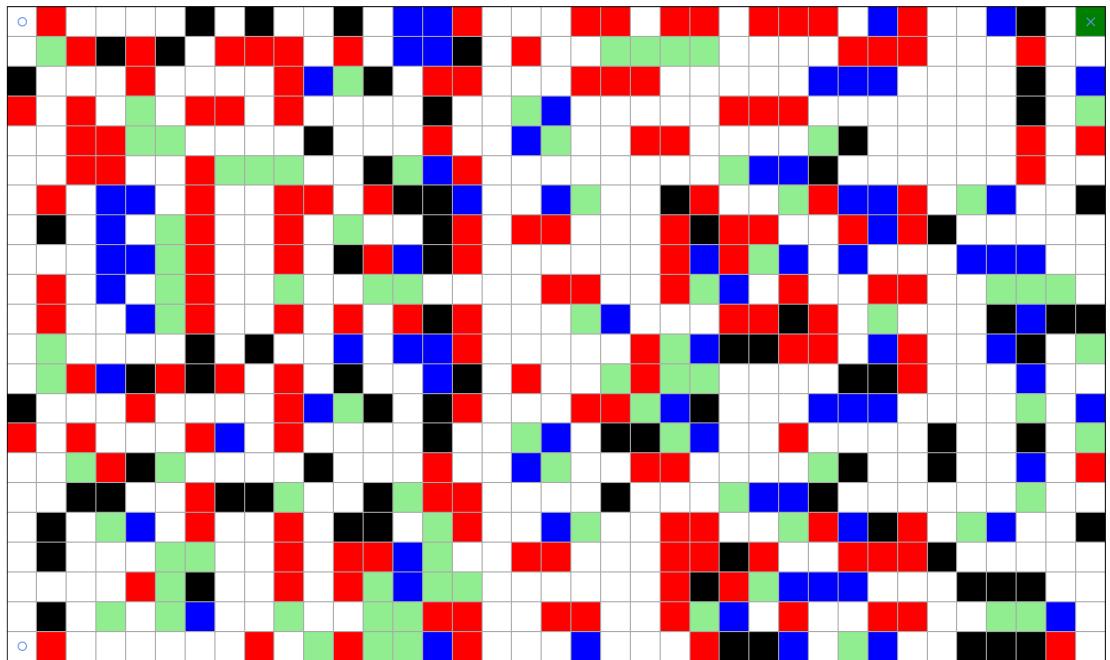


Figure 60: Gridworld 3

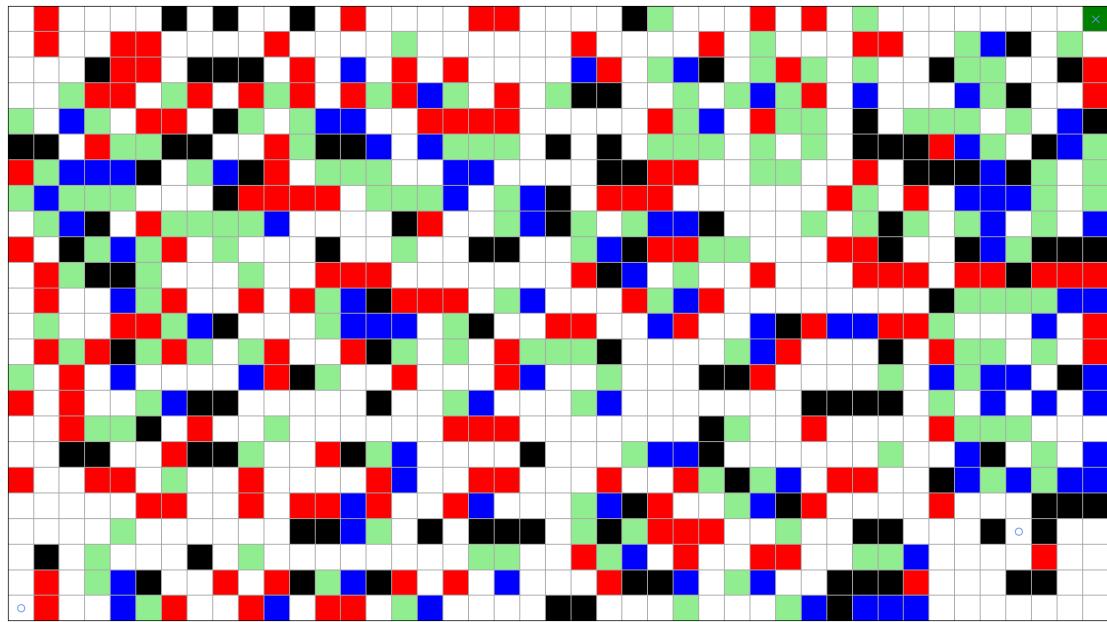


Figure 61: Novel Gridworld for Gridworld 3

Trained on 10000 trajectories with batch size of 128 for 50 epochs.

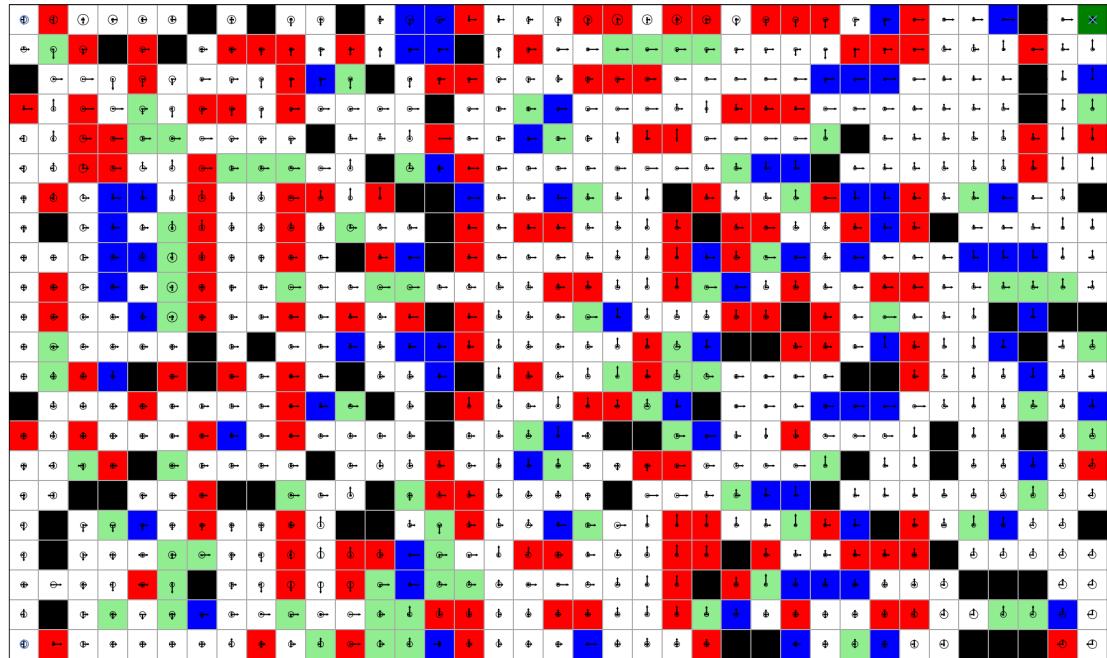


Figure 62: Classifier Policy

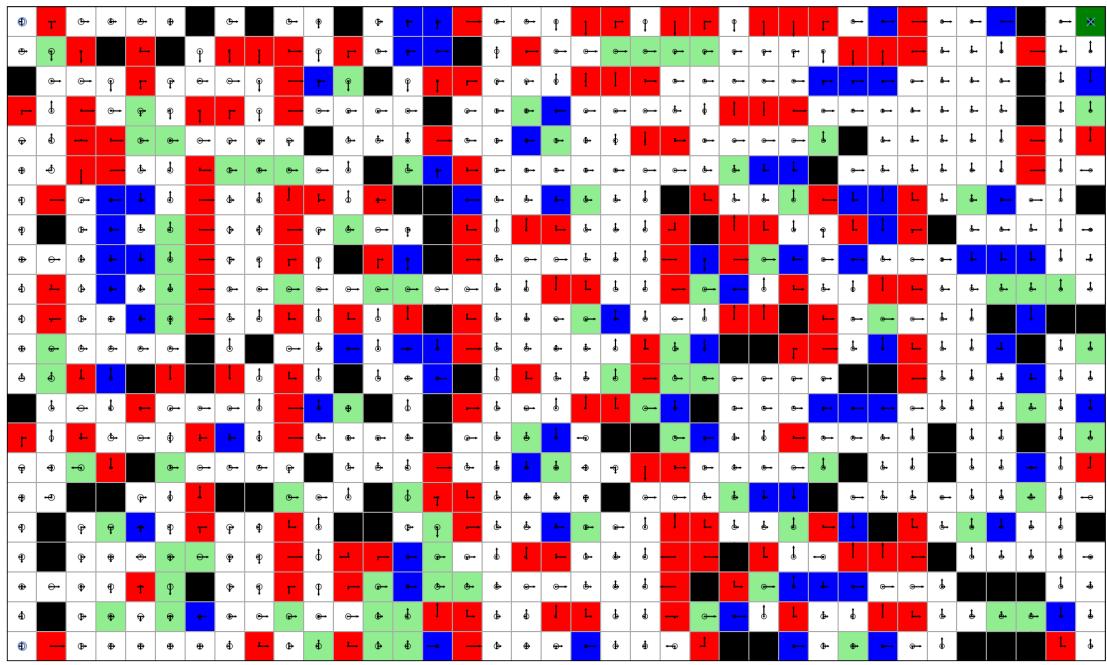


Figure 63: MLIRL Policy

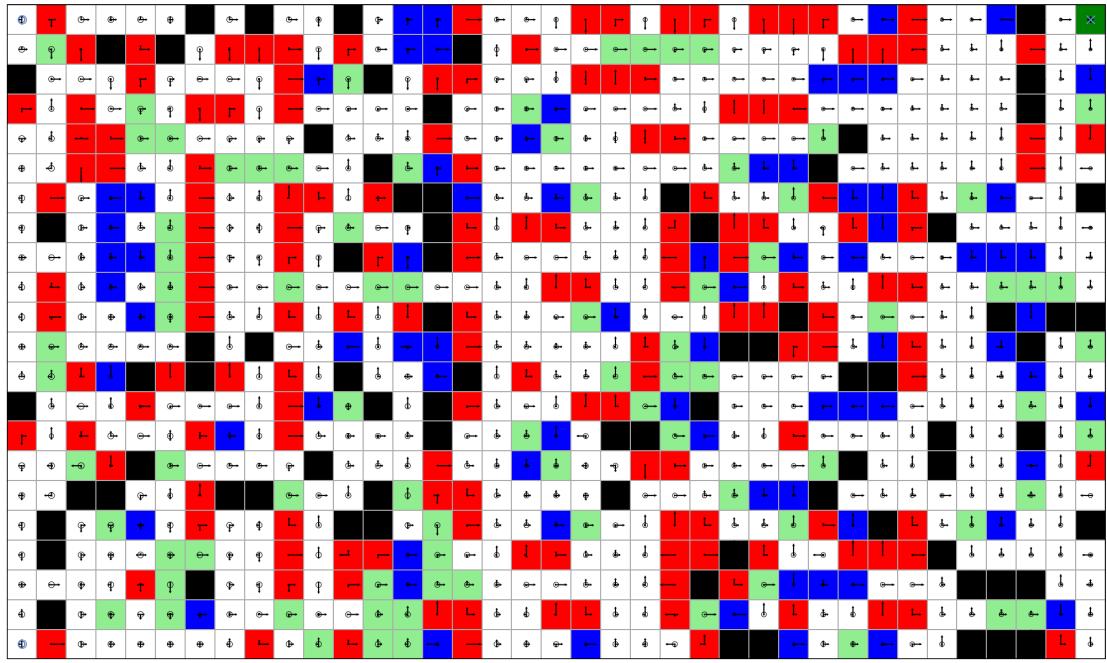


Figure 64: Expert Policy

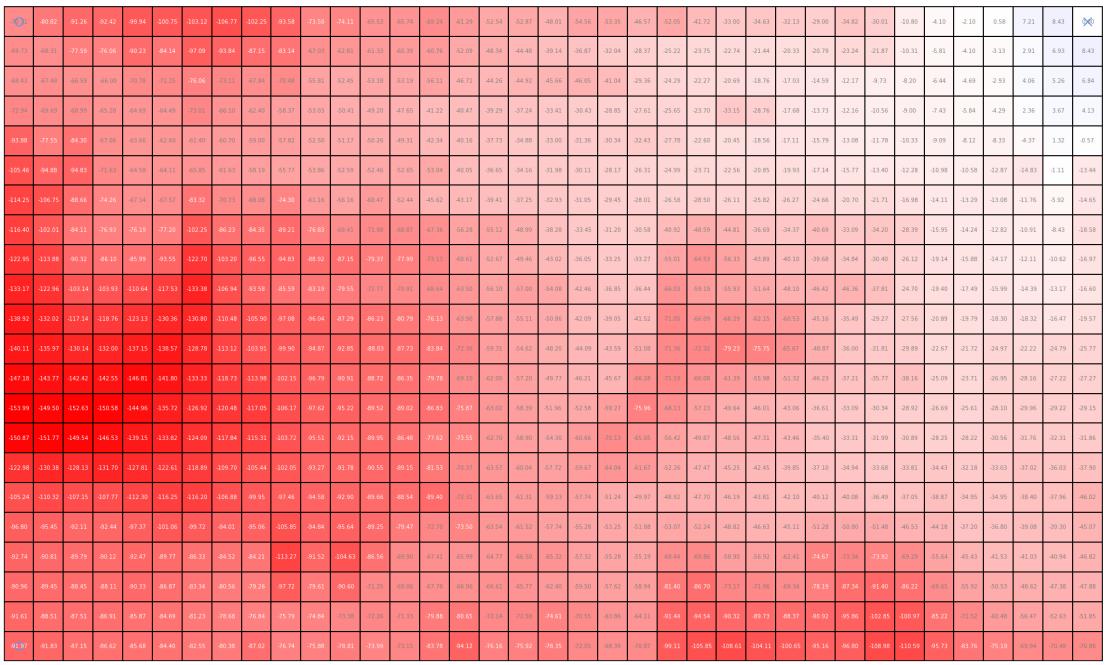


Figure 65: Classifier State Value Map

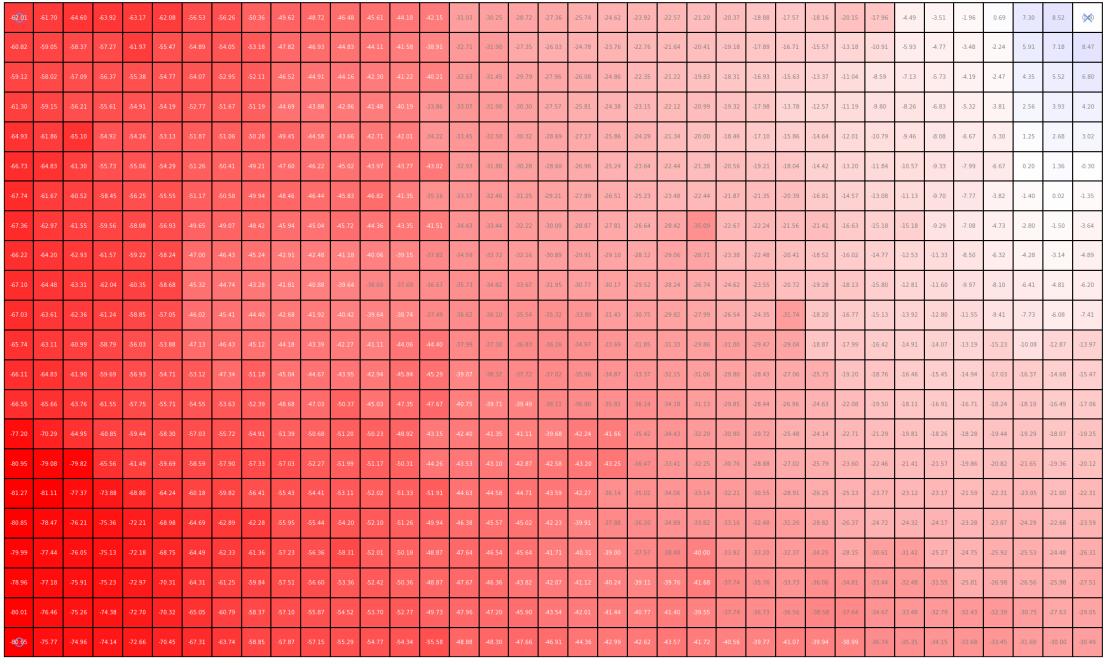


Figure 66: MLIR State Value Map

43.39	41.17	44.12	43.43	42.68	41.08	56.09	55.88	49.98	49.15	48.25	49.05	49.19	49.77	41.85	30.73	29.51	29.45	27.11	25.52	24.41	23.78	22.20	21.00	20.17	-18.68	17.40	18.10	20.04	17.85	4.41	3.41	-1.89	0.74	7.32	8.53	X
60.29	58.55	57.88	56.78	43.49	55.08	54.43	53.58	52.73	47.57	46.48	44.41	43.69	41.18	38.85	32.30	31.57	27.08	29.39	24.56	23.55	22.54	21.42	20.21	18.99	17.70	-16.53	15.39	13.03	20.79	9.84	-4.76	-3.42	-2.20	5.94	7.20	8.49
58.61	57.53	56.62	55.90	44.92	54.32	53.62	52.51	51.48	46.58	44.47	43.73	42.89	40.92	39.83	33.20	31.59	29.46	27.66	25.84	24.63	22.32	21.00	19.61	18.11	-16.78	-15.45	-13.20	12.89	8.46	-7.02	5.64	-4.11	-2.40	4.40	5.57	6.85
67.49	58.58	55.74	55.14	54.44	43.71	42.31	31.23	30.75	44.26	41.46	42.45	41.09	38.30	33.50	32.77	31.53	29.76	27.26	25.54	24.14	22.92	21.00	20.75	19.37	17.79	-13.58	-12.39	11.03	9.66	-8.14	-6.72	5.23	2.63	3.99	4.26	
43.88	40.94	44.00	54.44	53.78	52.66	51.42	50.61	49.83	49.02	44.16	43.26	42.31	41.62	37.85	33.03	32.51	29.95	28.37	26.87	25.58	24.04	21.06	19.73	18.23	16.89	15.66	14.95	11.86	10.66	9.34	7.91	-6.57	-5.21	1.32	3.74	3.08
45.75	43.79	40.79	55.23	54.57	53.81	50.79	49.94	48.74	47.15	45.79	44.60	43.57	43.37	42.64	32.54	31.39	29.88	28.24	26.83	24.93	23.35	22.26	21.11	20.23	18.96	17.82	-14.26	13.04	-11.70	10.44	-9.20	-7.68	-6.57	8.27	3.43	0.22
46.87	40.98	59.93	57.89	55.73	55.05	50.69	50.13	49.46	48.03	46.02	49.41	46.40	40.94	34.76	32.97	32.08	30.07	28.84	27.55	26.19	24.02	23.26	22.37	21.41	21.11	20.16	-18.65	14.40	-12.92	-11.08	-9.58	-7.61	-3.67	1.31	0.10	1.26
46.54	42.24	40.91	58.94	57.53	56.40	49.10	48.51	47.86	45.45	44.57	49.20	43.91	42.90	41.11	34.03	33.74	31.84	29.73	28.53	27.48	26.31	28.17	34.91	22.42	22.01	21.32	21.21	-16.45	-15.03	-15.01	0.11	4.51	4.58	2.69	4.40	3.52
45.39	43.46	42.25	40.91	48.62	47.69	46.47	45.88	44.73	42.45	42.02	40.72	39.62	38.71	37.47	34.25	33.71	31.79	30.55	29.57	28.75	27.77	28.77	28.47	23.13	22.23	20.16	-18.27	15.78	14.85	12.11	0.35	4.14	4.14	3.02	4.76	
46.22	43.67	42.55	41.25	39.65	48.08	44.83	44.25	42.83	41.38	40.45	40.23	39.28	37.21	36.25	35.31	34.43	33.20	31.58	30.42	29.81	29.14	27.06	26.41	24.35	23.28	20.46	19.03	17.80	15.51	12.60	11.33	0.77	-7.92	-4.25	-4.66	6.06
46.03	42.58	42.46	40.32	38.06	46.42	45.53	44.93	43.93	42.24	42.48	39.99	39.22	38.32	37.07	34.22	30.63	30.12	34.91	31.41	31.05	30.38	29.52	27.70	26.27	24.10	30.45	17.92	16.51	14.96	13.77	12.58	11.37	0.22	-7.55	5.91	7.25
44.59	41.93	39.95	57.87	55.29	53.22	46.03	45.94	44.64	43.52	42.92	41.70	40.64	43.55	43.91	31.95	34.05	36.91	25.86	26.57	33.30	31.49	31.00	29.53	30.63	29.11	29.78	18.61	17.73	16.17	14.65	13.65	12.95	15.01	8.90	12.67	13.79
44.99	43.62	40.81	58.72	56.15	44.14	42.62	40.84	35.73	44.17	41.43	42.42	45.35	44.84	38.83	37.93	37.31	36.82	35.38	34.99	33.63	29.63	28.08	26.72	25.44	18.92	18.50	16.21	15.21	14.67	16.78	18.12	34.44	15.26			
45.36	44.36	42.58	40.55	56.97	55.08	53.99	53.08	51.88	48.10	46.45	49.82	44.45	46.83	47.21	40.39	39.21	39.06	31.65	36.49	35.41	35.14	33.79	33.15	29.51	28.11	26.63	24.29	21.76	20.20	17.81	16.64	18.41	17.97	17.90	16.29	16.81
75.58	68.44	43.35	59.94	58.61	47.51	58.41	55.19	54.29	50.93	50.13	50.68	49.70	48.40	47.69	41.91	40.87	40.64	49.21	41.76	41.23	35.01	34.02	31.88	30.12	29.36	25.10	23.79	22.38	-20.91	-19.51	-17.68	-17.05	18.99	-17.76	18.98	
73.54	77.48	78.26	43.84	49.36	28.89	37.98	57.31	56.76	56.48	51.72	51.46	58.63	49.78	43.71	43.61	42.17	42.34	42.07	42.70	42.75	36.00	32.98	31.85	30.13	28.46	26.64	25.43	23.26	22.34	21.08	21.21	-19.55	-20.52	21.29	19.04	19.84
79.62	79.69	75.87	72.14	44.90	42.71	49.38	59.16	55.83	54.68	53.86	52.57	53.48	40.77	41.40	44.16	44.63	44.18	43.05	41.74	35.87	34.57	33.67	32.69	30.10	28.50	28.89	24.79	19.44	22.78	22.81	21.05	21.98	22.68	20.64	21.94	
79.50	76.69	74.74	73.83	29.49	47.31	43.60	42.11	41.58	55.39	54.87	43.63	63.52	40.87	49.35	45.83	45.82	44.46	41.64	39.36	37.39	35.34	34.47	33.38	32.73	32.08	30.83	28.46	26.02	24.89	23.97	23.81	22.95	23.51	23.89	22.27	23.19
78.61	75.59	74.53	73.61	70.54	47.22	43.77	41.46	40.58	58.59	55.76	57.72	55.39	49.56	48.28	47.07	45.95	45.09	41.16	39.87	38.51	37.08	38.01	39.57	38.65	32.77	32.93	33.83	27.79	29.71	31.01	24.63	24.37	25.55	25.11	29.04	25.81
77.76	75.74	74.47	73.77	71.42	48.78	43.04	40.75	39.96	54.84	50.97	52.71	52.78	49.72	48.28	47.08	45.71	43.25	41.56	40.61	40.73	38.62	39.23	41.26	37.75	38.17	37.29	31.62	31.17	29.41	26.60	29.13	25.51	17.04			
78.85	75.61	73.81	72.89	71.16	48.77	43.59	59.79	57.63	54.63	55.21	53.86	43.04	62.11	49.11	47.35	46.60	45.28	42.97	41.49	40.32	40.24	40.89	39.98	37.01	36.26	38.12	37.37	34.24	33.00	32.21	32.05	31.37	30.27	32.12	28.51	
74.31	73.50	72.65	71.12	48.69	45.60	42.49	48.18	47.22	66.49	54.64	54.11	53.67	54.91	48.29	47.76	47.05	-6.32	43.78	42.45	42.07	41.60	40.93	40.58	39.47	38.32	35.21	18.91	31.71	31.01	31.21	29.58	31.99				

Figure 67: Expert State Value Map

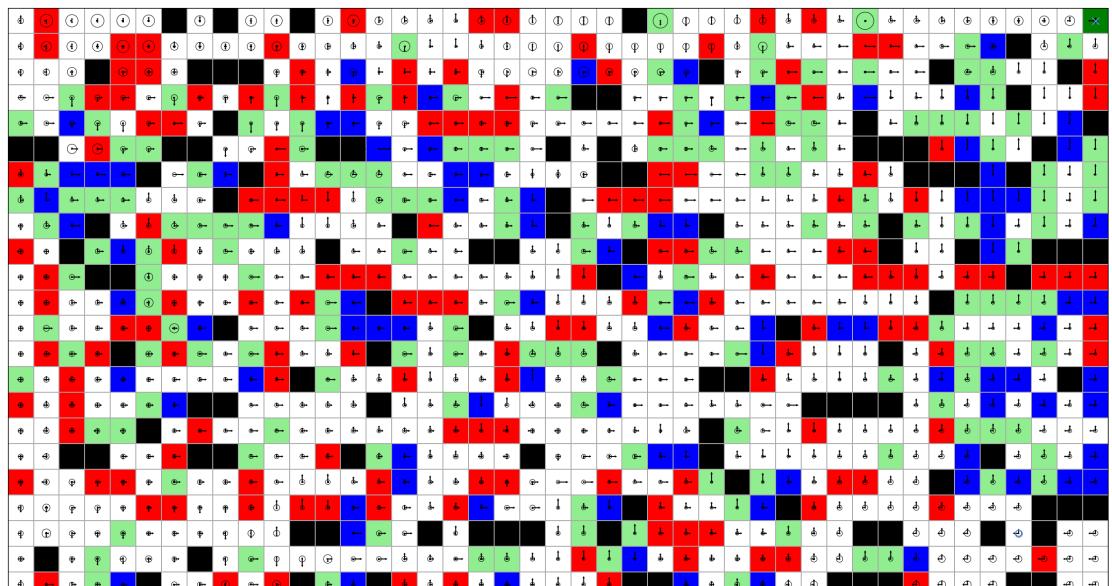


Figure 68: Classifier Policy on Novel Gridworld

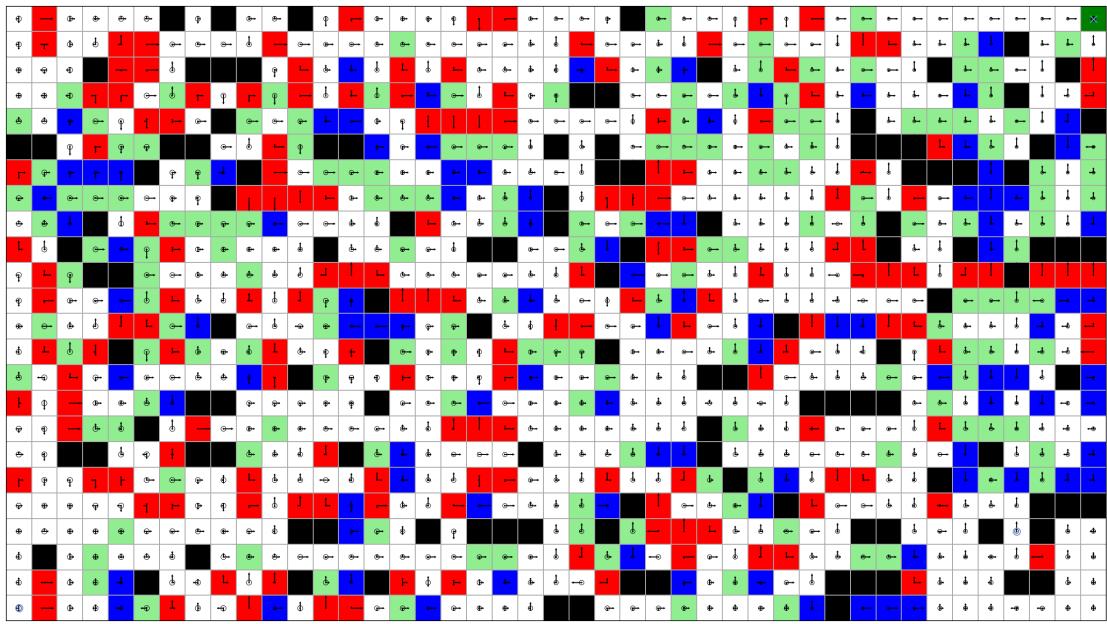


Figure 69: MLIRL Policy on Novel Gridworld

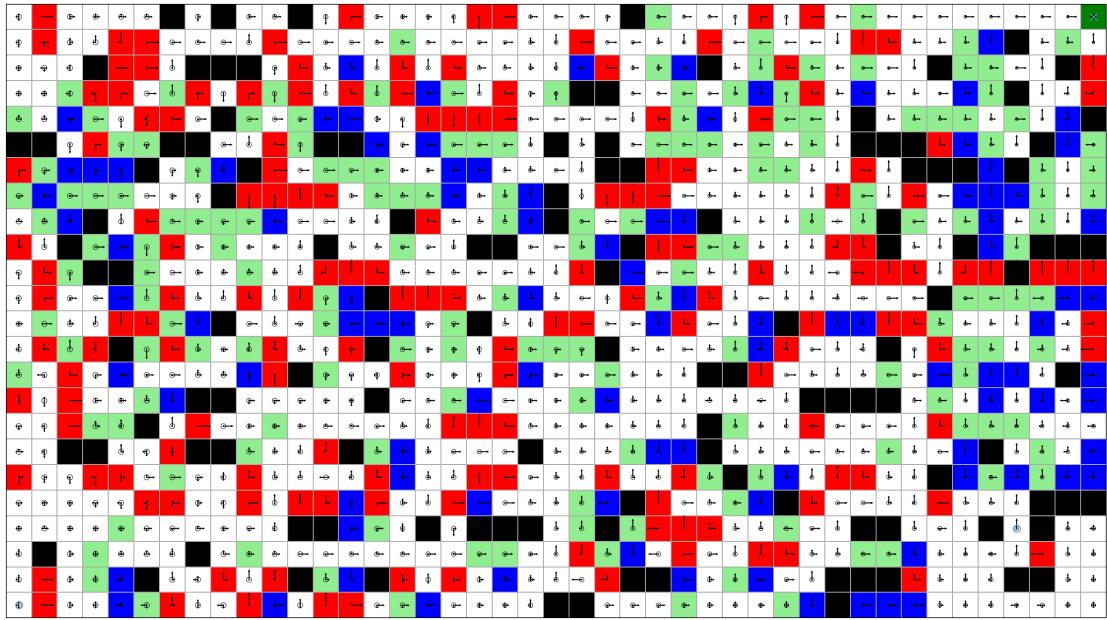


Figure 70: Expert Policy on Novel Gridworld

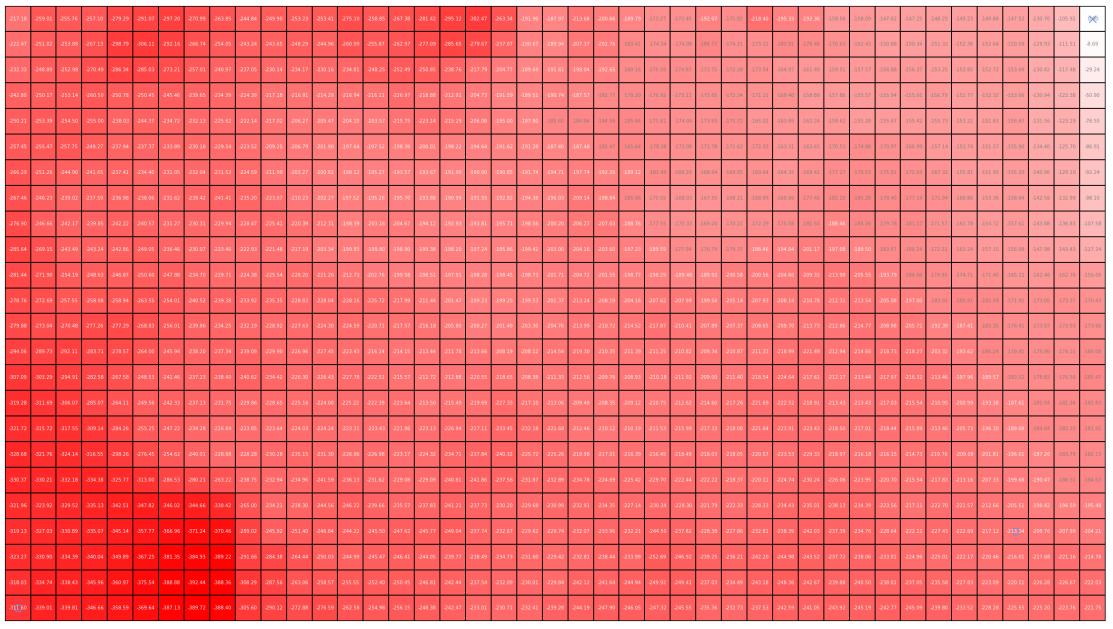


Figure 71: Classifier State Value Map on Novel Gridworld

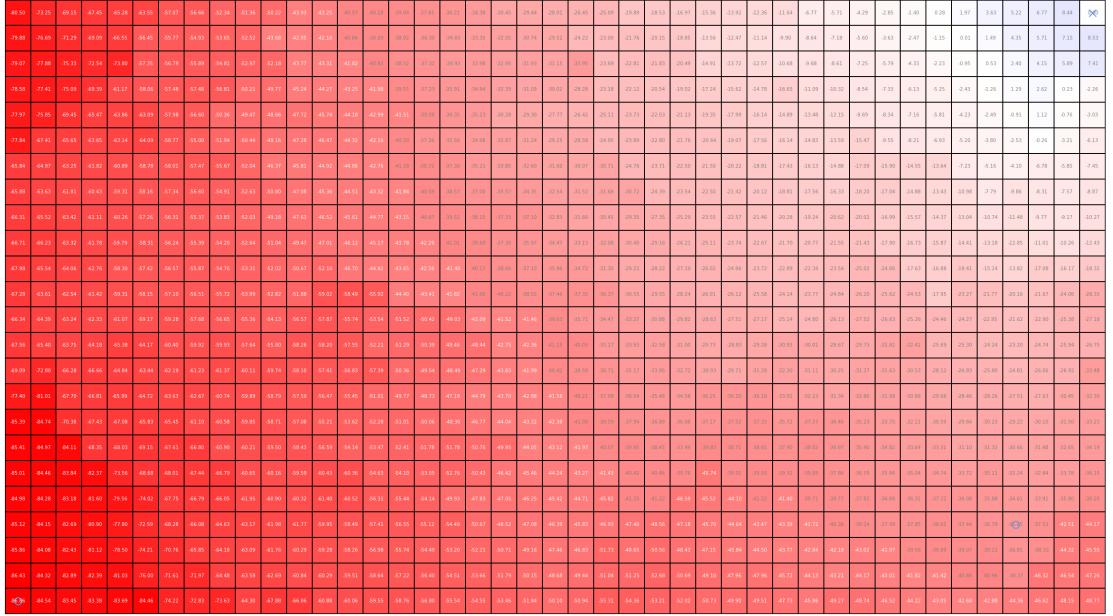


Figure 72: MLIR State Value Map on Novel Gridworld

Figure 73: Expert State Value Map on Novel Gridworld

Trained on 100000 trajectories with batch size of 128 for 50 epochs.

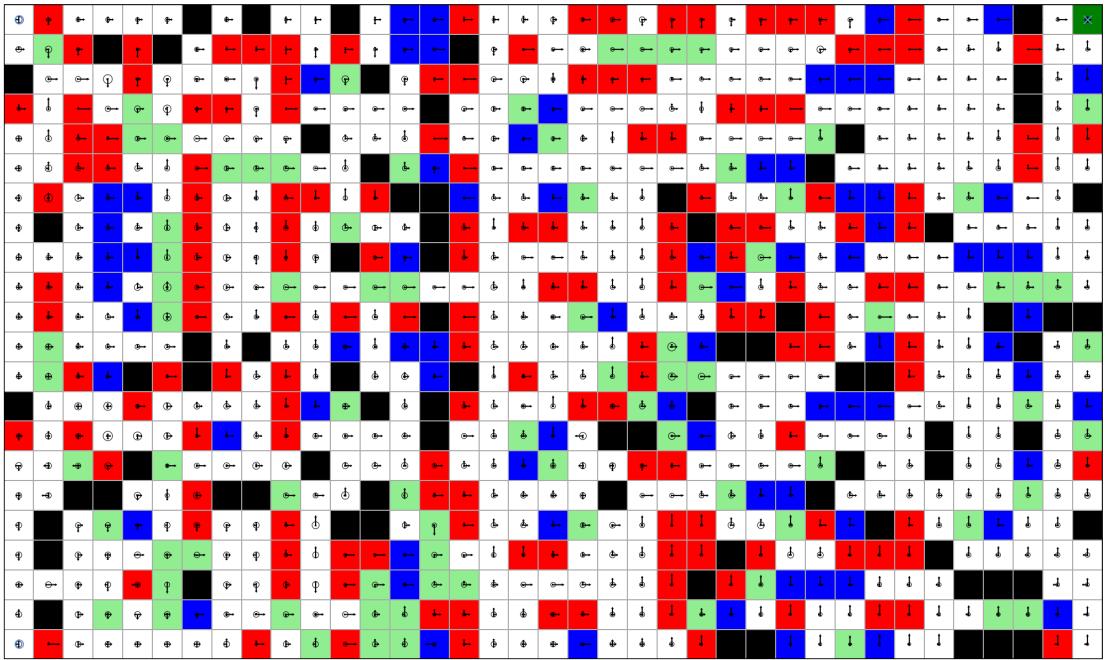


Figure 74: Classifier Policy

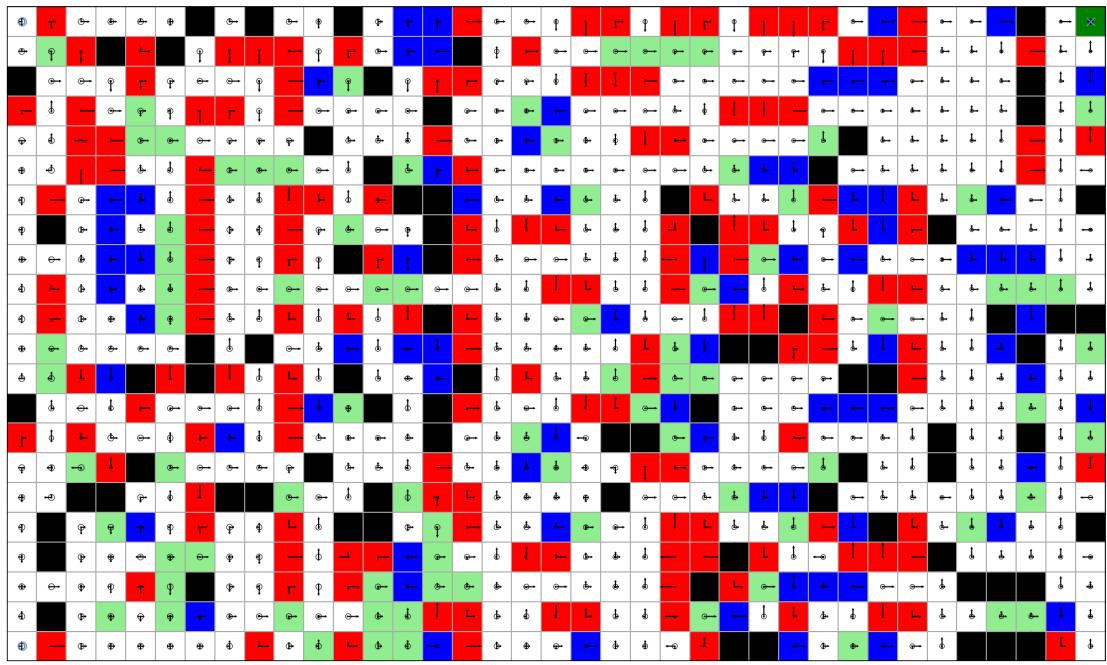


Figure 75: MLIRL Policy

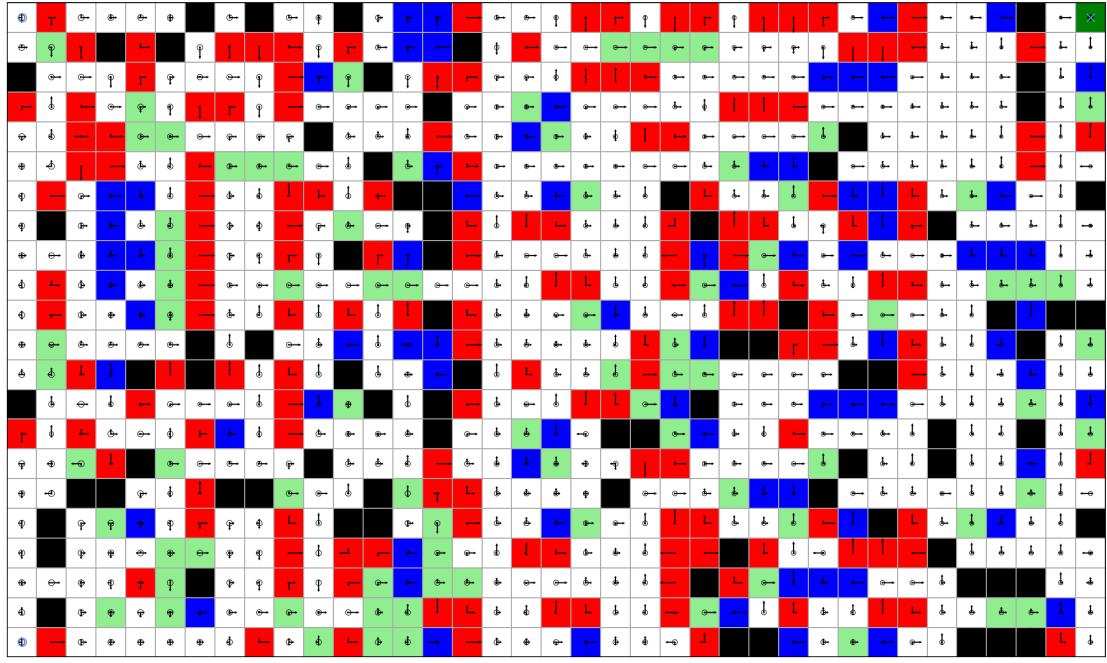


Figure 76: Expert Policy

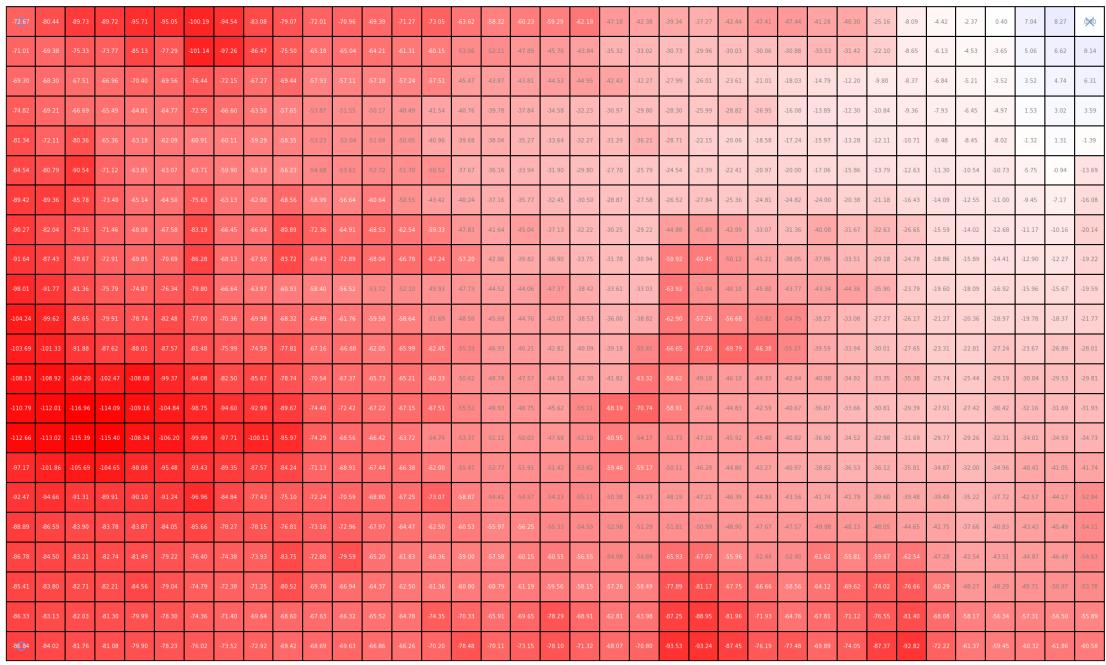


Figure 77: Classifier State Value Map

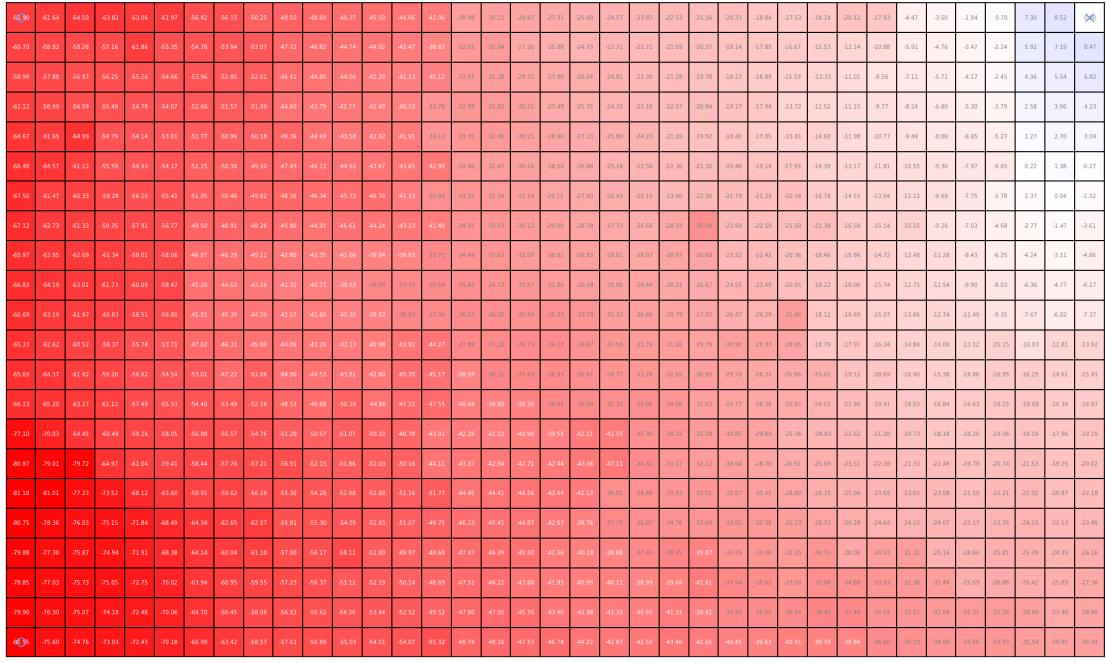


Figure 78: MLIR State Value Map

43.39	41.17	44.12	43.43	42.68	41.08	56.09	55.88	49.98	49.15	48.25	49.05	49.19	43.77	41.85	30.73	29.51	29.45	27.11	25.52	24.41	23.78	22.39	21.03	20.17	-18.68	17.40	18.10	20.04	17.85	4.41	3.41	-1.89	0.74	7.32	8.53	X	
60.29	58.55	57.88	56.78	43.49	55.08	54.43	53.58	52.73	47.57	46.48	44.41	43.69	41.18	38.85	32.38	31.57	27.98	29.39	24.56	23.55	22.54	21.42	20.21	18.99	17.70	-16.53	15.39	13.03	20.79	9.84	-4.76	-3.42	-2.20	5.94	7.20	0.49	
58.61	57.53	56.62	55.90	44.92	54.32	53.62	52.51	51.48	46.58	44.47	43.73	42.89	40.82	39.83	33.20	31.59	29.46	27.66	29.84	24.63	22.32	21.09	19.61	18.11	-16.78	-15.45	-13.29	12.89	-8.46	-7.02	-5.64	-4.11	-2.40	4.40	5.57	0.85	
67.49	58.58	55.74	55.14	54.44	43.71	42.31	31.23	30.75	44.26	41.46	42.45	41.09	38.30	33.50	32.77	31.53	29.76	27.26	25.54	24.14	22.92	21.80	20.75	19.37	17.79	-13.58	-12.39	11.03	9.66	-8.14	-6.72	5.23	2.63	3.99	4.26		
43.88	40.94	44.00	54.44	53.78	52.66	51.42	50.61	49.83	49.02	44.16	43.26	42.31	41.62	37.85	33.03	32.51	29.95	28.37	26.87	25.58	24.04	21.96	19.73	18.23	16.89	15.66	14.95	11.86	10.66	9.34	7.91	6.57	5.21	1.32	3.74	3.08	
45.75	43.79	40.79	55.23	54.57	53.81	50.79	49.94	48.74	47.15	45.79	44.60	43.57	43.37	42.64	32.54	31.38	29.88	28.24	26.83	24.93	23.35	22.26	21.11	20.23	18.96	17.82	-14.26	13.04	-11.70	10.44	-9.20	-7.68	-6.57	8.27	3.43	0.22	
46.87	40.98	59.93	57.89	55.73	55.05	50.68	50.13	49.46	48.03	46.02	49.41	46.40	40.94	34.76	32.97	32.08	30.07	28.84	27.55	26.19	24.02	23.26	22.37	21.41	21.11	20.16	-18.65	14.40	-12.92	-11.08	-9.58	-7.61	-3.67	1.31	0.10	1.26	
46.54	42.24	40.91	58.94	57.53	56.40	49.10	48.51	47.86	45.45	44.57	49.28	43.91	42.90	41.11	34.03	33.74	31.84	29.73	28.53	27.48	26.31	28.27	34.91	22.42	22.03	21.32	21.21	16.45	15.03	15.01	9.11	4.51	4.58	2.69	4.40	3.52	
45.39	43.46	42.25	40.91	48.62	47.69	46.47	45.88	44.73	42.45	42.02	40.72	39.62	38.71	37.47	34.25	33.71	31.79	30.55	29.57	28.75	27.77	28.77	28.47	23.13	22.23	20.16	18.27	15.78	14.85	12.11	11.14	8.35	6.14	4.14	3.02	4.76	
44.22	43.67	42.55	41.25	39.65	48.08	44.83	44.25	42.83	41.38	40.45	40.23	39.28	37.21	36.28	35.31	34.43	33.28	31.58	30.42	29.81	29.14	27.06	26.41	24.35	23.28	20.46	19.03	17.80	15.51	12.60	11.33	8.77	7.92	4.25	4.66	6.06	
46.03	42.58	42.46	40.32	38.06	46.42	45.53	44.93	43.93	42.24	42.48	39.99	39.22	38.32	37.07	34.22	30.63	30.12	34.91	31.41	31.05	30.38	29.52	27.70	26.27	24.10	30.45	17.92	16.51	14.96	13.77	12.58	11.37	8.22	7.55	5.91	7.25	
44.59	41.93	39.95	57.87	55.29	53.22	46.03	45.94	44.64	43.52	42.92	41.70	40.64	43.55	43.91	31.95	34.05	36.91	25.86	26.57	33.30	31.49	31.00	29.53	30.63	29.11	29.78	18.61	17.73	16.17	14.65	13.65	12.95	15.01	8.90	12.67	13.79	
44.99	43.62	40.81	58.72	56.15	44.14	42.62	40.84	35.73	44.05	44.17	41.43	42.42	45.35	44.94	38.83	37.83	37.31	36.82	35.38	33.99	33.63	29.63	28.38	26.72	25.44	18.92	18.50	16.21	15.21	14.67	16.78	18.12	14.44	15.26			
45.36	44.36	42.58	40.55	56.97	55.08	53.99	53.08	51.88	48.10	46.45	49.82	44.45	46.83	47.21	40.39	39.21	39.06	31.65	36.49	35.41	35.14	33.79	33.15	29.51	28.11	26.63	24.29	21.76	18.41	17.97	17.90	16.29	16.81				
75.58	68.44	43.35	59.94	58.61	37.51	58.41	55.19	54.29	50.83	50.13	50.68	49.70	49.40	47.69	41.91	40.87	40.64	39.21	41.76	41.23	35.01	34.02	31.88	30.12	29.36	25.10	23.79	22.38	-20.91	-19.51	-17.68	-17.05	18.99	-17.76	18.98		
73.54	77.48	78.26	43.84	49.36	28.89	37.98	57.31	56.76	56.48	51.72	51.46	58.63	49.78	43.71	43.61	42.17	42.34	42.07	42.70	42.75	36.00	32.98	31.85	30.13	28.46	26.64	25.43	23.26	22.34	21.08	21.21	19.55	-20.52	21.29	19.04	19.84	
79.62	79.69	75.87	72.14	44.90	42.71	49.38	59.16	55.83	54.68	53.86	52.57	53.48	40.77	51.40	44.16	44.63	44.18	43.05	41.74	35.87	34.57	33.67	32.69	29.50	28.50	28.89	24.79	19.44	22.78	22.81	21.05	21.98	22.68	20.64	21.94		
79.50	76.69	74.74	73.83	29.49	47.31	43.60	42.11	41.58	55.39	54.87	43.63	63.52	40.87	49.35	45.83	45.82	44.46	41.64	39.36	37.39	35.34	34.47	33.38	32.73	32.08	30.83	28.46	26.02	24.39	23.97	23.81	22.95	23.51	23.89	22.27	23.19	
78.61	75.59	74.53	73.61	70.54	47.22	43.77	41.46	40.58	58.59	55.76	57.72	55.39	49.56	48.28	47.07	45.95	45.09	41.16	39.87	38.51	37.08	38.01	29.57	28.35	27.77	32.77	32.93	33.83	27.79	29.71	31.01	24.65	24.37	25.55	25.11	29.04	25.81
77.76	75.74	74.47	73.77	71.42	48.78	43.04	40.75	39.96	54.84	50.97	52.71	52.78	49.72	48.28	47.08	45.71	43.25	41.56	40.61	39.73	38.62	39.23	41.26	37.75	35.17	37.29	31.62	31.17	31.17	29.41	26.60	29.13	25.51	37.04			
78.85	75.61	73.81	72.89	71.16	48.77	43.59	59.79	57.63	54.43	55.21	53.86	43.04	62.11	49.11	47.35	46.60	45.28	42.97	41.49	40.32	40.24	40.89	39.98	37.01	36.26	40.18	37.17	34.24	33.00	32.21	32.05	31.37	30.27	32.12	30.54		
74.31	73.50	72.65	71.12	48.69	45.60	42.49	58.18	47.22	66.49	54.64	54.11	53.67	54.91	48.29	47.76	47.05	-6.32	43.78	42.45	42.07	41.60	49.71	40.58	39.47	38.32	35.21	18.91	31.71	31.01	31.21	29.58	31.99					

Figure 79: Expert State Value Map

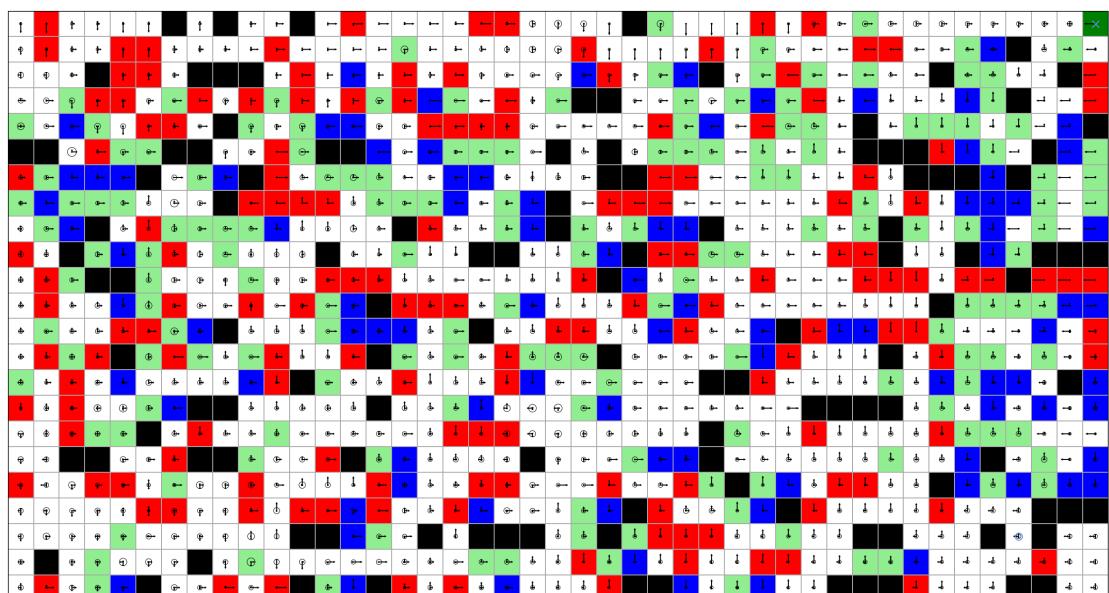


Figure 80: Classifier Policy on Novel Gridworld

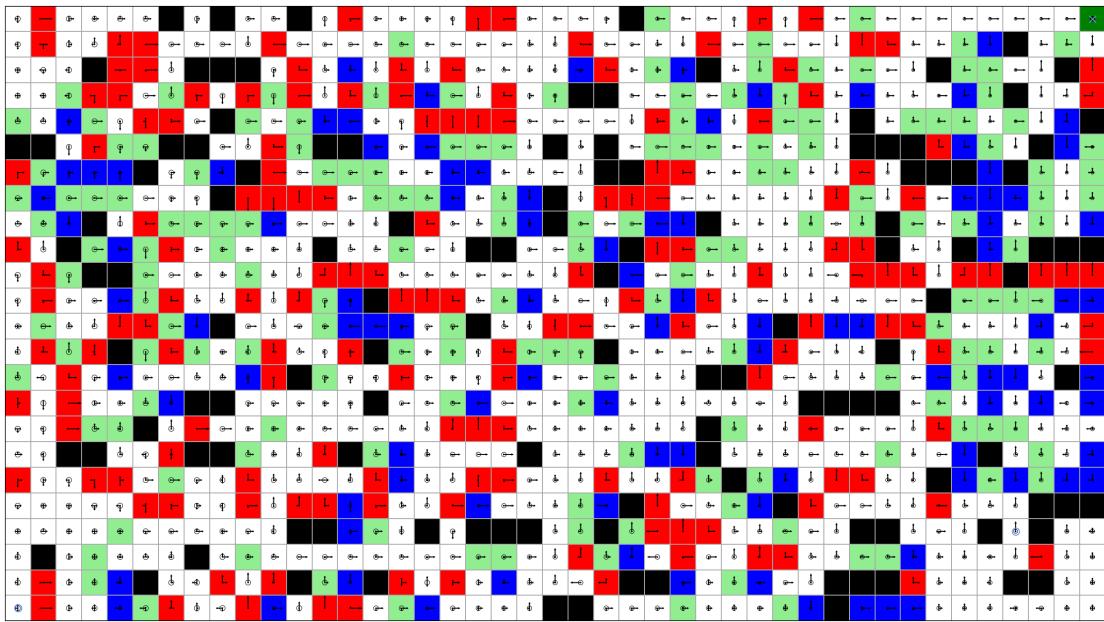


Figure 81: MLIRL Policy on Novel Gridworld

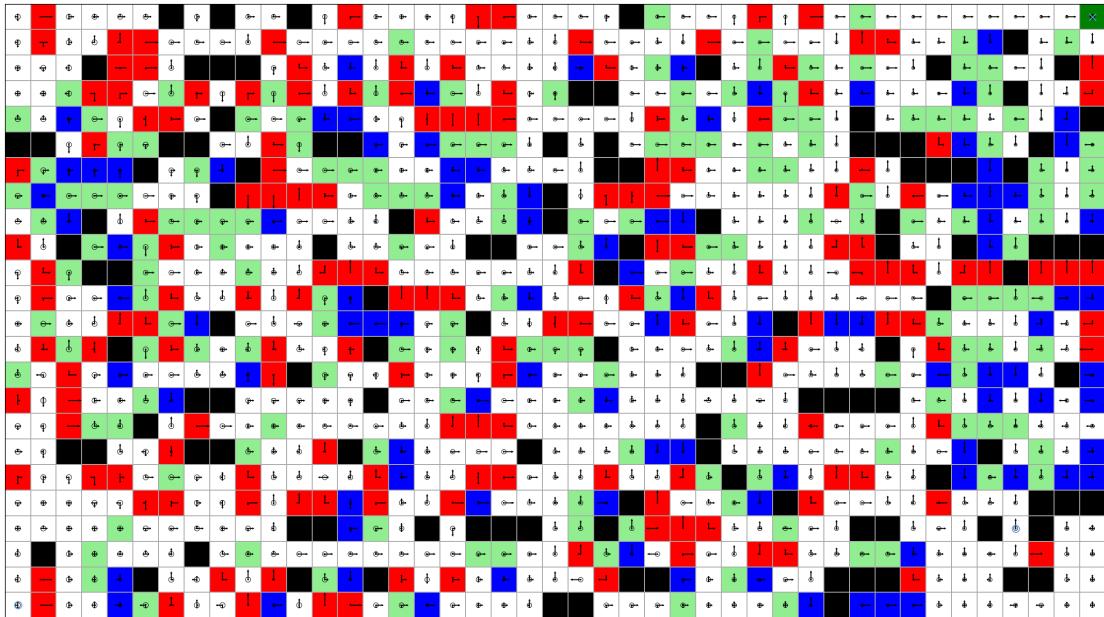


Figure 82: Expert Policy on Novel Gridworld

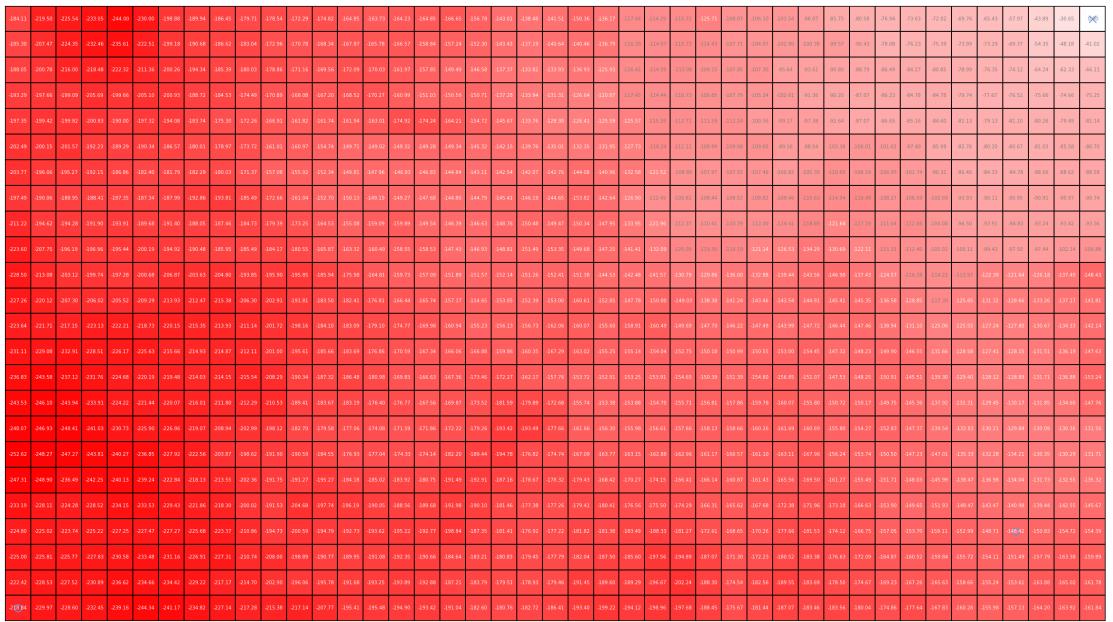


Figure 83: Classifier State Value Map on Novel Gridworld

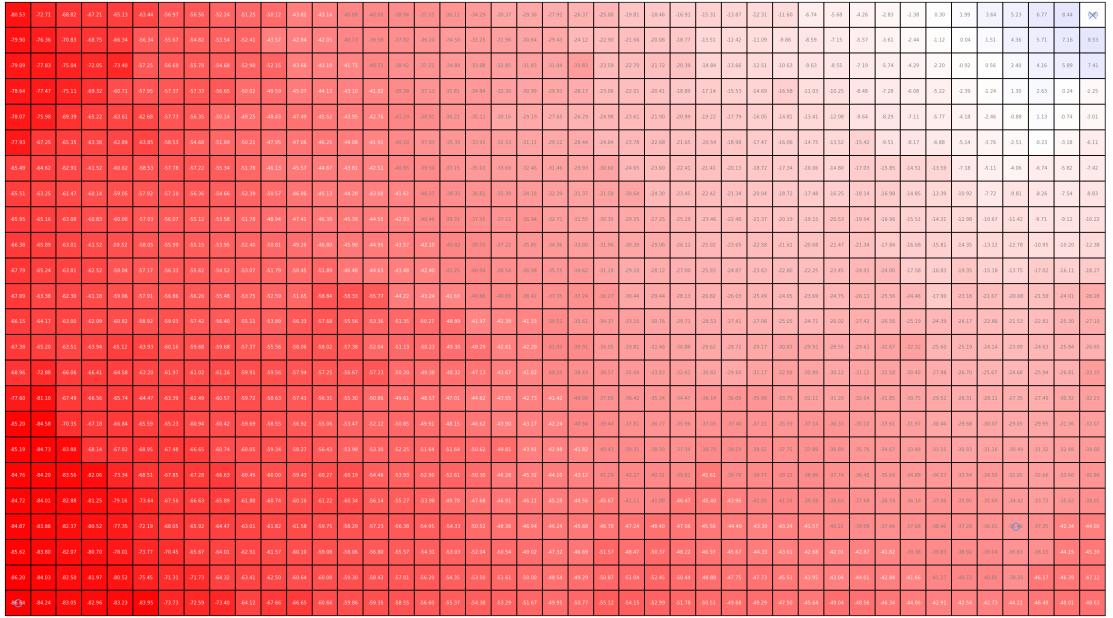


Figure 84: MLIR State Value Map on Novel Gridworld

79.15	43.47	48.03	44.51	44.14	42.87	56.53	48.11	31.79	50.83	43.69	43.61	42.72	39.69	29.59	33.03	39.00	29.09	21.48	28.13	24.78	39.49	18.27	18.72	21.14	33.71	32.17	11.46	6.62	3.58	4.15	2.71	3.28	3.36	2.05	3.70	3.27	4.00	8.65	2.00			
79.49	44.36	69.94	47.94	45.98	55.81	59.22	44.36	53.11	51.99	43.14	42.67	43.87	44.37	44.37	39.13	39.24	39.48	24.51	31.35	31.09	30.36	29.14	22.68	21.45	13.05	18.57	21.35	32.27	21.99	8.72	6.44	7.02	3.45	3.51	2.35	4.14	0.11	1.57	4.39	5.14	7.11	8.54
79.62	45.24	71.44	51.92	44.44	58.86	56.21	45.31	48.25	52.87	43.14	43.69	42.46	43.29	43.11	38.05	36.75	38.14	21.55	31.54	30.73	31.54	31.54	21.35	21.45	13.51	18.25	18.47	13.51	21.16	0.61	4.01	2.11	0.04	3.62	7.46	4.19	5.91	7.42	8.54			
77.14	25.95	73.55	42.05	40.39	37.51	56.52	42.87	42.20	49.87	42.27	49.53	41.72	41.09	42.42	40.75	39.95	35.54	30.55	32.05	29.03	22.95	21.20	20.19	20.68	20.99	22.75	24.55	16.45	10.89	12.12	4.96	7.75	3.99	5.32	2.36	1.15	1.39	3.97	0.27	2.27		
79.52	24.38	69.08	44.50	43.38	42.27	37.26	35.92	49.72	48.84	48.30	47.08	45.10	43.51	43.34	46.08	39.13	35.81	28.73	29.85	29.89	27.37	28.13	28.10	23.37	24.89	20.78	19.01	21.53	25.89	24.40	13.27	13.94	6.53	4.10	2.07	3.66	4.09	2.37	4.83	1.19	0.39	2.96
79.39	46.25	64.45	42.43	42.48	43.30	56.03	44.19	31.25	49.83	43.53	46.63	45.81	43.66	44.47	33.96	30.73	34.93	33.58	32.20	39.80	20.85	28.70	24.63	23.37	22.40	21.85	20.36	18.78	17.39	15.80	14.60	13.37	15.28	2.38	8.05	4.76	3.05	2.67	4.14	3.13	4.00	
54.93	44.07	42.31	45.99	46.18	58.82	37.24	46.72	54.84	51.36	45.15	45.17	44.27	43.43	42.11	40.65	38.62	31.67	34.89	33.37	30.15	31.18	29.02	30.38	28.44	33.09	29.20	21.21	19.94	18.54	21.78	15.89	24.05	36.89	35.72	14.98	31.40	7.08	0.01	31.04	4.67	3.78	2.35
34.97	42.35	69.94	43.63	55.18	37.42	56.69	45.86	54.18	51.93	55.12	46.41	44.73	43.88	42.69	41.23	40.01	37.79	36.81	35.05	30.11	31.18	30.42	34.08	33.24	21.21	21.13	21.88	18.34	17.31	14.09	33.04	18.05	14.78	32.26	13.79	9.63	16.16	0.17	7.61	4.75		
45.41	44.62	42.32	49.32	49.32	49.16	34.52	50.57	44.67	53.09	51.73	49.46	46.90	45.89	44.58	44.19	40.74	36.93	37.42	36.79	34.63	32.81	31.41	29.11	29.61	27.01	16.94	13.75	21.16	18.96	20.27	24.79	18.07	15.30	22.45	22.66	23.13	4.62	9.61	20.12			
53.83	45.26	42.49	42.01	49.65	37.68	55.11	44.65	51.46	52.32	52.36	49.81	46.76	45.55	45.35	43.23	42.72	41.41	30.45	30.46	30.96	32.55	30.93	34.75	31.51	30.51	29.89	23.84	22.38	21.41	20.46	22.28	21.27	17.93	18.82	23.45	24.20	20.98	42.01	10.93	20.31	22.10	
47.39	44.71	43.29	42.02	53.55	56.71	55.85	35.18	54.01	52.53	52.31	49.98	51.61	49.98	46.25	41.11	42.28	40.82	20.83	38.21	35.87	31.40	34.75	30.34	28.87	27.89	26.78	21.71	24.48	21.44	22.03	23.22	22.63	22.32	22.08	21.42	16.64	33.20	15.95	13.69	16.93	24.02	38.11
46.49	43.94	41.79	48.48	58.57	37.41	56.38	45.77	44.99	53.27	52.11	51.36	54.46	37.92	50.19	43.38	42.06	41.38	30.37	39.75	39.30	37.93	38.78	36.70	36.70	30.21	27.93	20.40	25.58	23.26	23.03	23.03	21.54	18.95	21.43	23.06	26.17						
35.54	43.82	42.49	43.58	46.32	58.81	59.33	48.92	35.98	54.59	53.39	55.81	37.21	53.12	52.94	49.98	49.48	41.81	41.09	49.98	39.15	35.31	34.88	32.89	39.53	29.48	28.26	27.12	24.86	24.36	25.79	26.08	23.63	24.21	21.01	21.26	22.68	23.33	26.36				
34.71	44.57	43.61	43.36	44.12	43.37	39.62	38.14	48.68	66.77	45.08	37.54	36.65	31.63	37.92	49.08	46.47	41.83	42.24	41.62	40.81	39.71	34.75	33.35	33.43	30.43	28.63	26.61	31.04	31.01	25.44	23.03	25.68	22.63	26.43	25.68	24.49						
48.10	43.29	40.49	43.81	43.95	42.10	51.39	43.48	48.85	59.32	49.85	47.49	45.95	47.93	47.70	40.72	46.27	48.42	45.95	47.93	47.70	47.75	47.75	47.75	47.75	31.13	32.22	31.22	31.22	31.22	31.22	31.22	31.22	31.22	31.22	31.22	31.22						
75.40	23.02	66.71	45.26	45.14	43.87	52.79	49.02	59.12	58.05	26.36	55.81	49.95	52.37	49.12	48.18	49.55	44.22	42.35	42.34	41.03	37.37	36.95	36.95	31.41	31.49	30.54	29.12	28.72	27.92	27.17	31.21	31.21	31.21	31.21	31.21	31.21	31.21					
83.30	42.98	69.95	44.59	66.38	44.89	56.63	42.39	59.08	58.14	58.00	56.38	54.58	22.97	33.14	40.39	49.85	47.72	46.23	43.51	47.78	41.43	41.72	38.19	37.46	36.43	35.83	37.11	34.16	35.25	36.87	38.08	34.83	33.73	30.24	29.60	29.87	29.73	32.42				
35.72	43.23	42.31	47.46	47.21	48.37	66.36	48.09	49.13	50.51	58.79	57.73	55.98	43.50	52.83	51.78	51.18	52.12	50.22	49.42	43.53	42.59	42.48	40.98	39.39	37.05	26.93	30.73	30.73	37.48	37.05	32.57	30.93	34.93	33.29	33.73	31.71	29.54	26.26	31.93	33.79		
34.29	42.69	42.02	49.46	71.14	47.39	57.35	48.72	49.08	59.95	59.68	59.75	59.69	53.99	53.97	52.48	52.19	49.99	45.88	48.92	42.70	42.73	40.81	39.36	39.23	39.23	39.23	39.23	39.23	39.23	39.23	39.23	39.23	39.23	39.23	39.23	39.23	39.23	39.23				
33.25	42.45	43.31	79.43	37.14	77.24	56.63	44.05	45.32	43.24	42.06	46.71	59.85	56.44	34.30	43.24	49.35	47.28	48.51	47.71	44.88	44.06	45.26	40.74	38.72	46.15	41.08	41.08	34.39	33.36	30.45	30.45	31.64	34.44	24.37	33.75	37.75						
83.39	42.24	40.70	79.90	25.77	70.91	57.21	43.29	43.09	42.49	42.21	42.19	52.14	37.75	56.74	52.09	54.47	37.95	46.23	43.51	47.78	41.43	41.72	37.46	36.43	35.83	37.11	34.16	30.24	29.60	29.87	29.73	32.42	42.49	42.49								
34.19	42.24	40.49	79.11	29.45	72.44	69.45	45.01	43.44	42.21	43.30	59.55	27.52	56.28	45.70	37.08	52.58	51.56	25.09	48.58	40.70	48.37	31.37	48.38	44.03	42.38	42.37	41.71	40.37	39.11	38.00	31.43	43.79	43.79									
34.01	42.69	40.98	48.43	78.95	74.06	20.33	39.99	43.78	42.81	43.96	40.13	59.55	48.76	37.90	44.48	55.68	53.88	53.80	42.13	49.54	40.11	48.87	50.86	50.05	49.55	47.38	47.38	43.63	42.73	43.68	42.73	43.68	43.68	43.68								
40.55	42.74	43.55	41.46	41.11	42.40	72.58	33.79	72.65	43.55	47.12	46.11	40.11	39.37	58.66	48.06	56.19	34.87	53.89	42.89	31.19	49.48	40.48	47.11	45.29	48.67	46.37	45.95	43.68	42.52	42.36	42.34	41.81	46.13	47.16	40.20							

Figure 85: Expert State Value Map on Novel Gridworld

10.5 Classifier and MLIRL Results on Gridworld 4

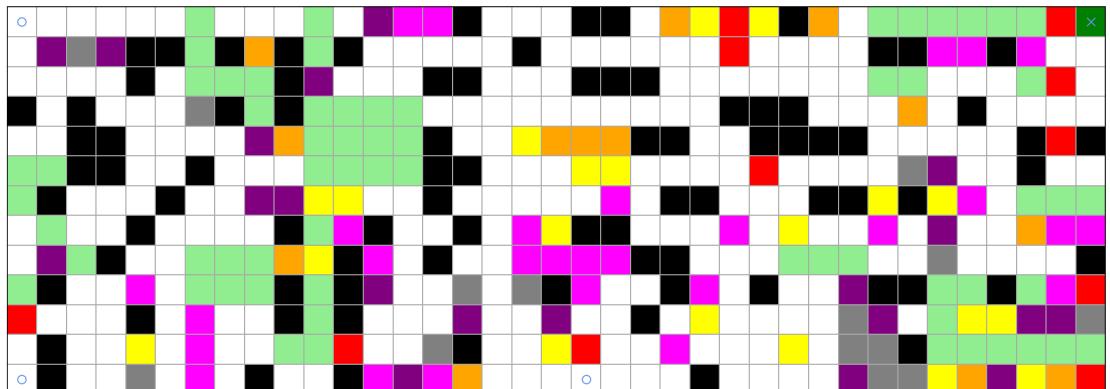


Figure 86: Gridworld 4

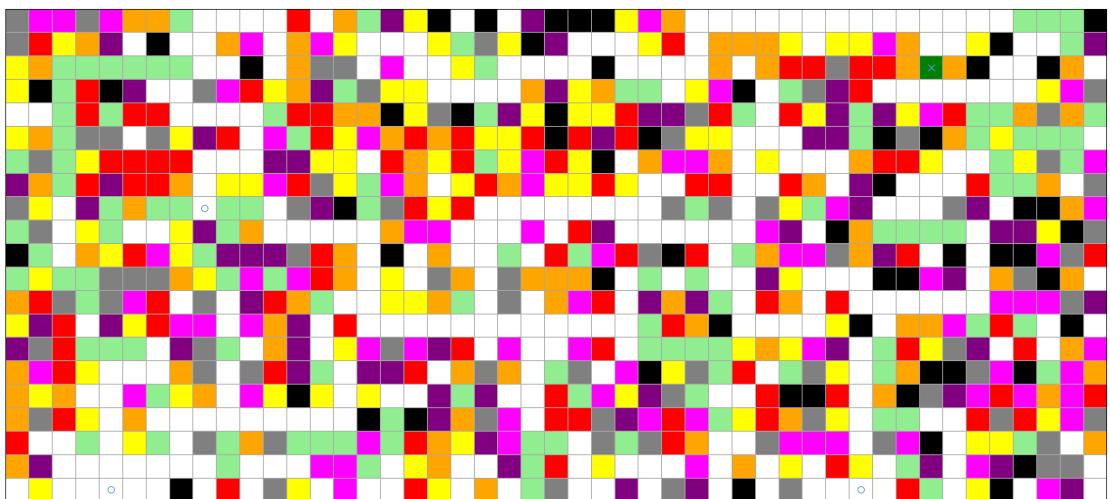


Figure 87: Novel Gridworld for Gridworld 4

Trained on 10000 trajectories with batch size of 128 for 50 epochs.

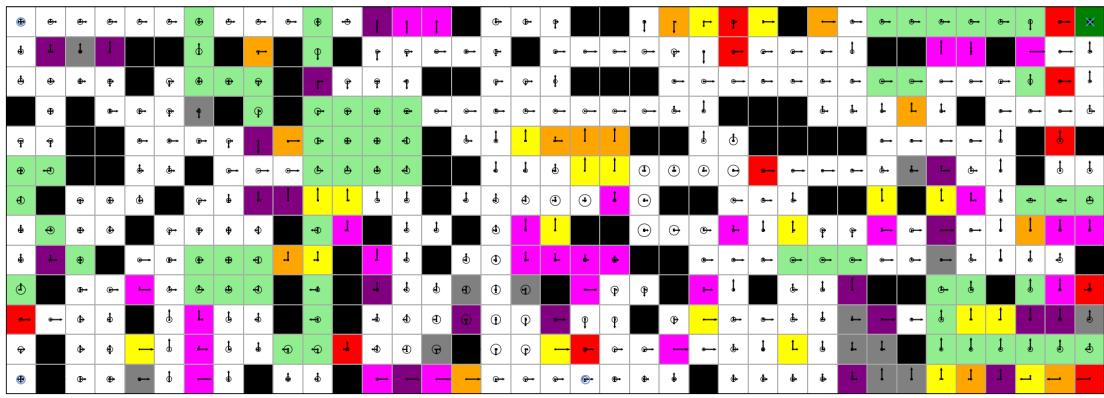


Figure 88: Classifier Policy

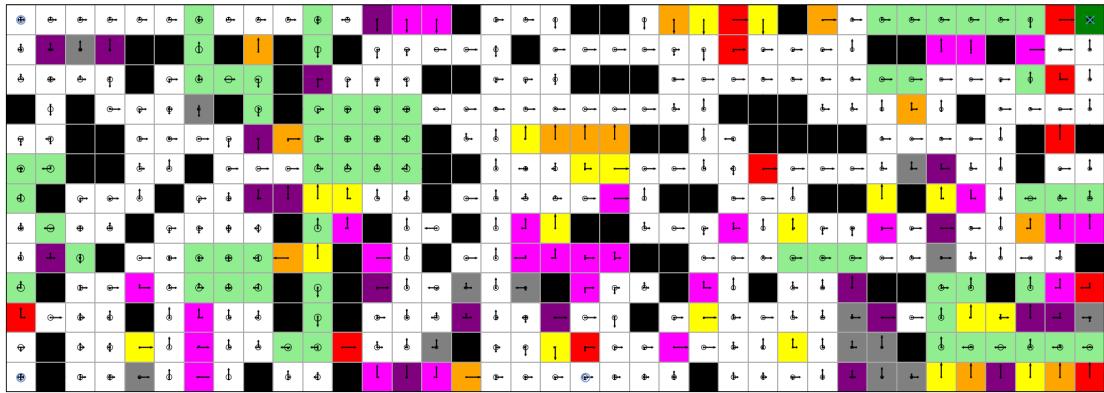


Figure 89: MLIRL Policy

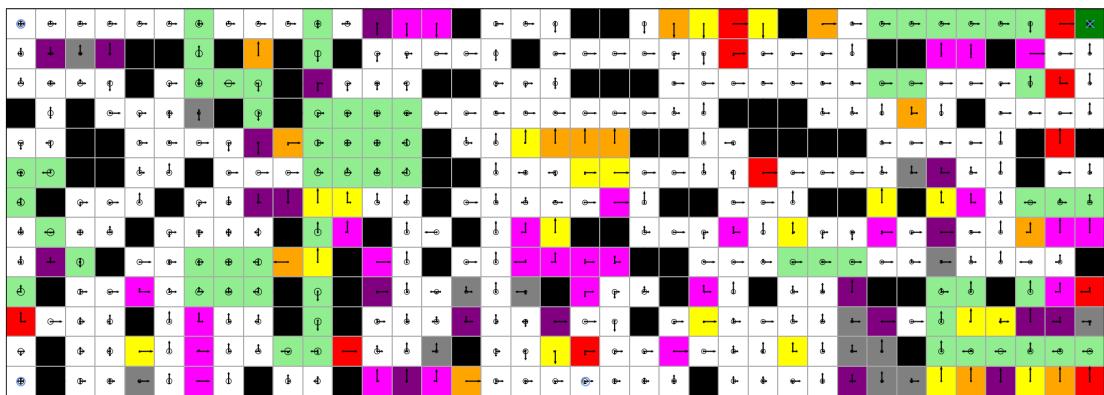


Figure 90: Expert Policy

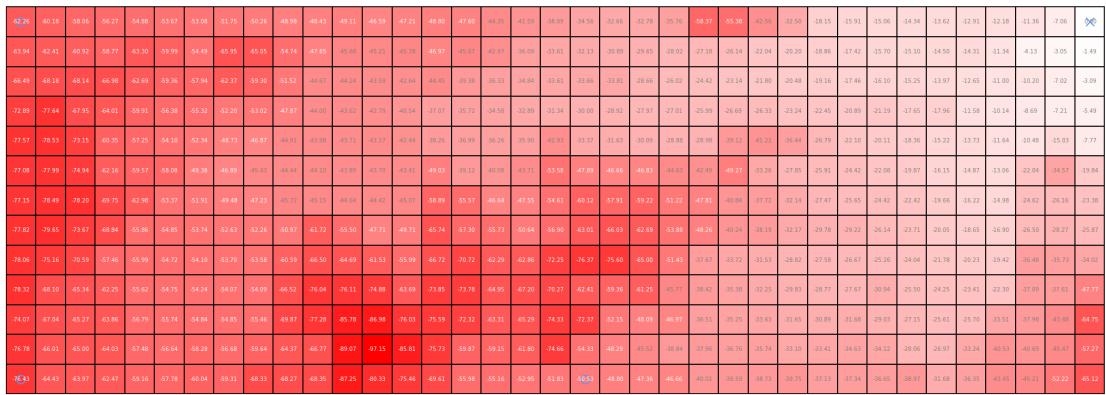


Figure 91: Classifier State Value Map

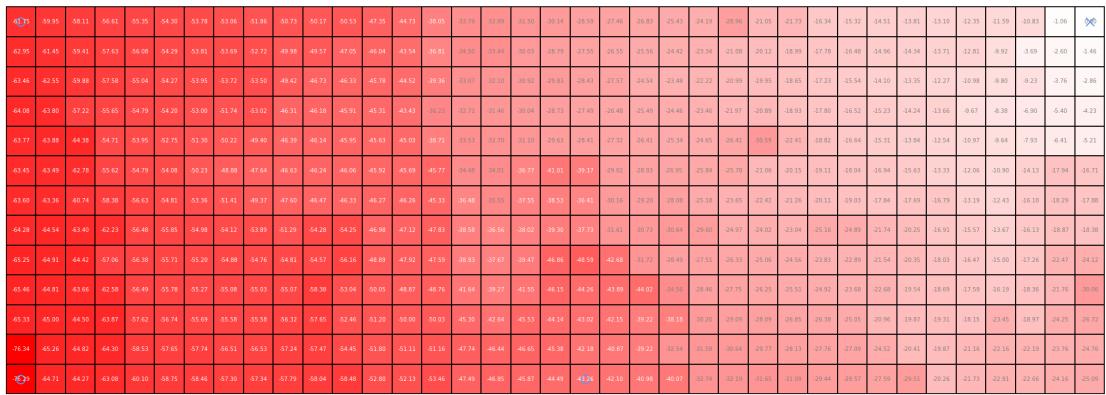


Figure 92: MLIRL State Value Map

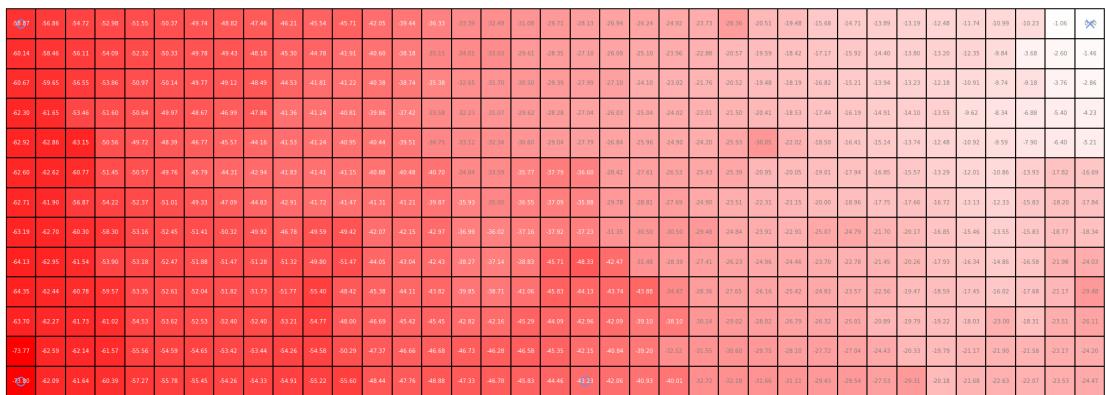


Figure 93: Expert State Value Map

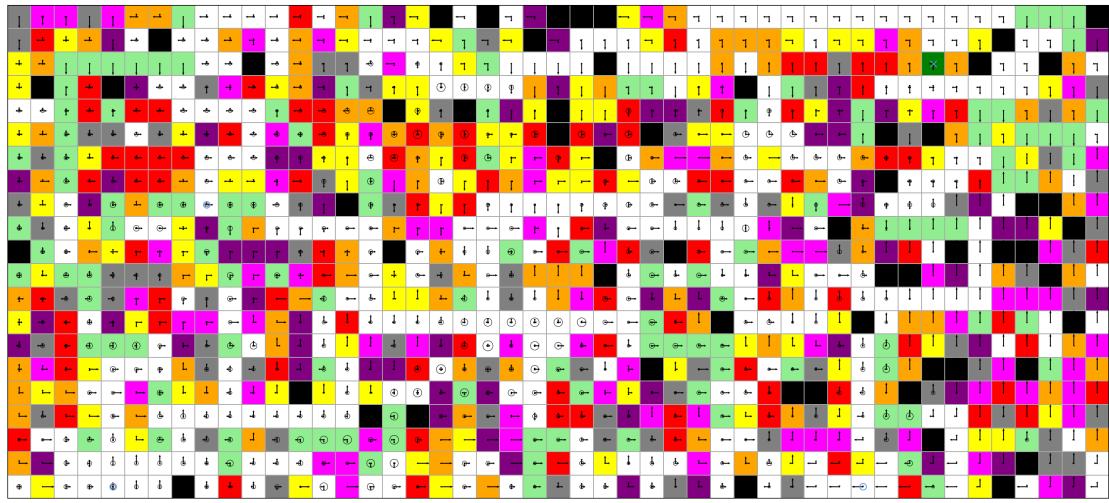


Figure 94: Classifier Policy on Novel Gridworld

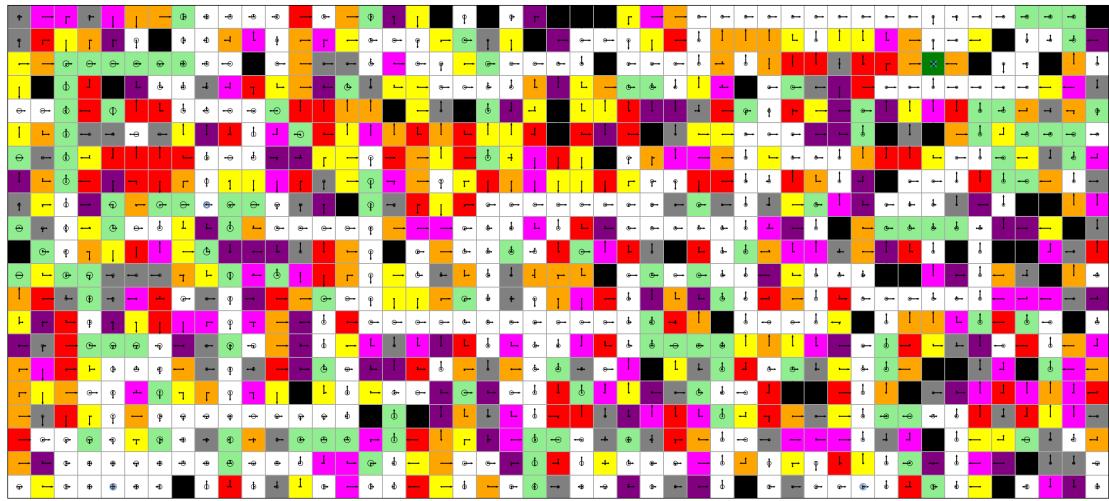


Figure 95: MLIRL Policy on Novel Gridworld

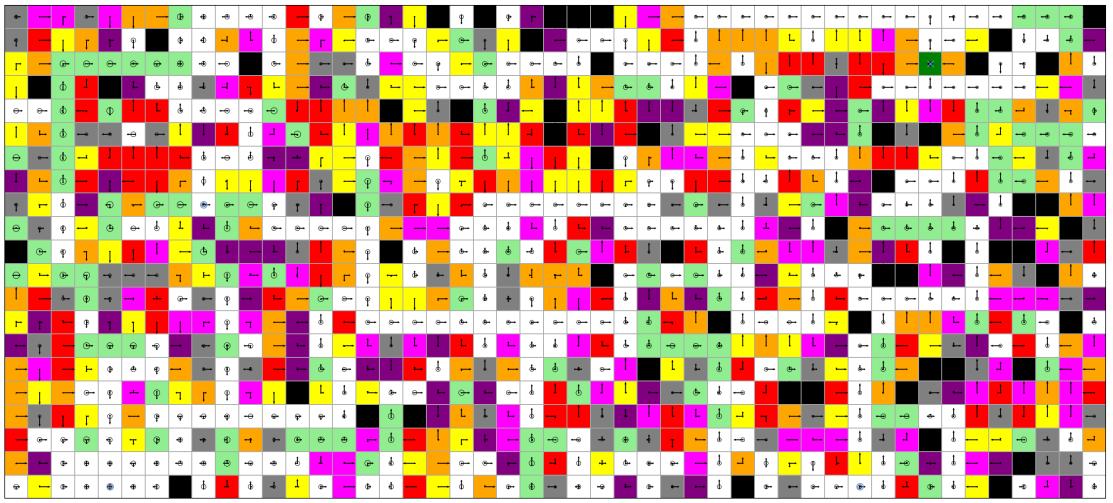


Figure 96: Expert Policy on Novel Gridworld

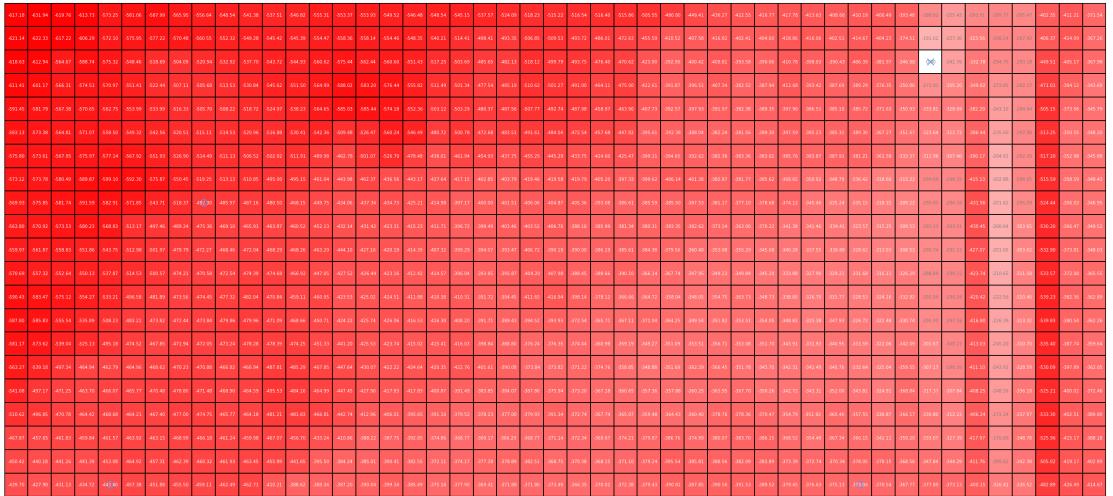


Figure 97: Classifier State Value Map on Novel Gridworld

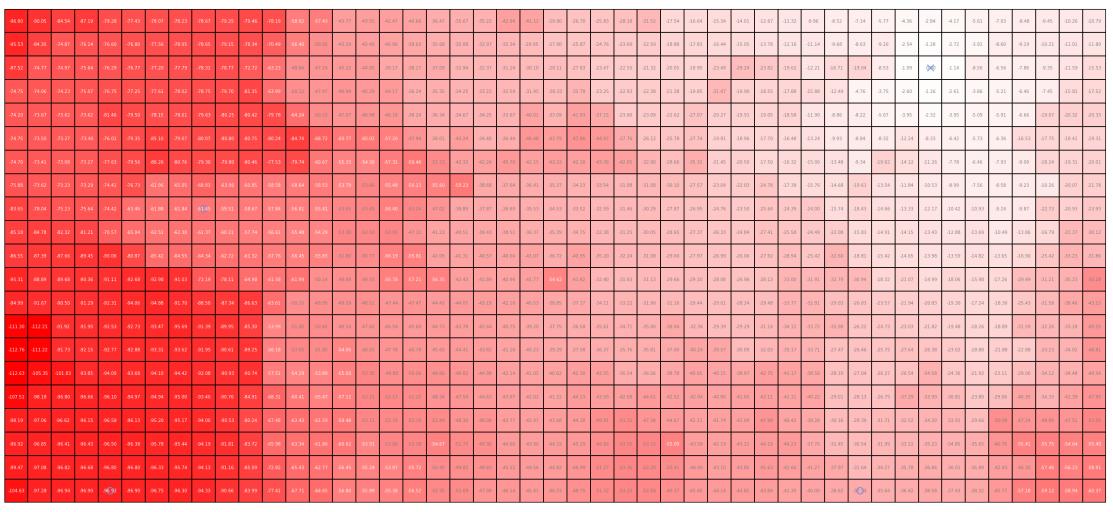


Figure 98: MLIRL State Value Map on Novel Gridworld

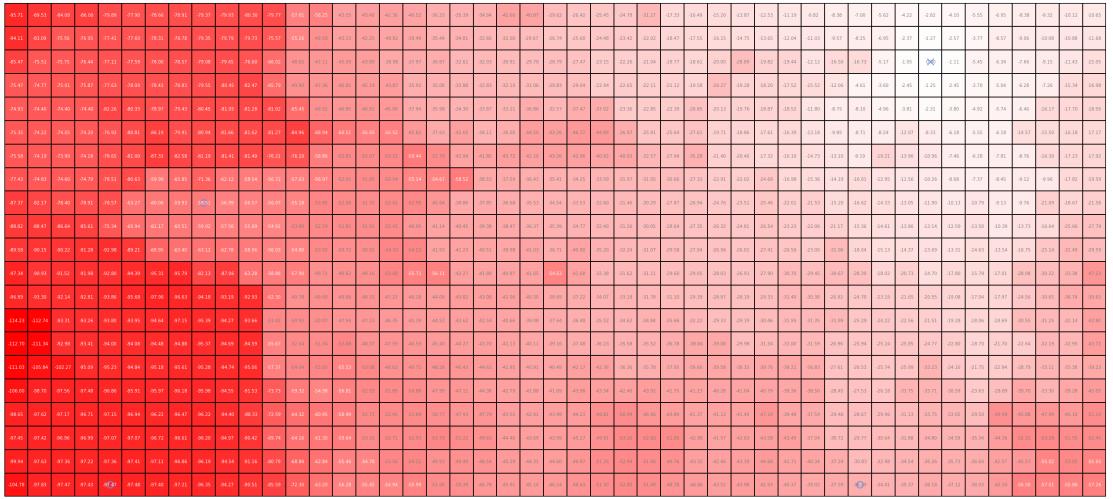


Figure 99: Expert State Value Map on Novel Gridworld

10.6 Classifier Training with 1-hot Vector Representation

Trained on 10000 trajectories with batch size of 128 for 50 epochs on Figure 3. All other hyperparameters remain the same.

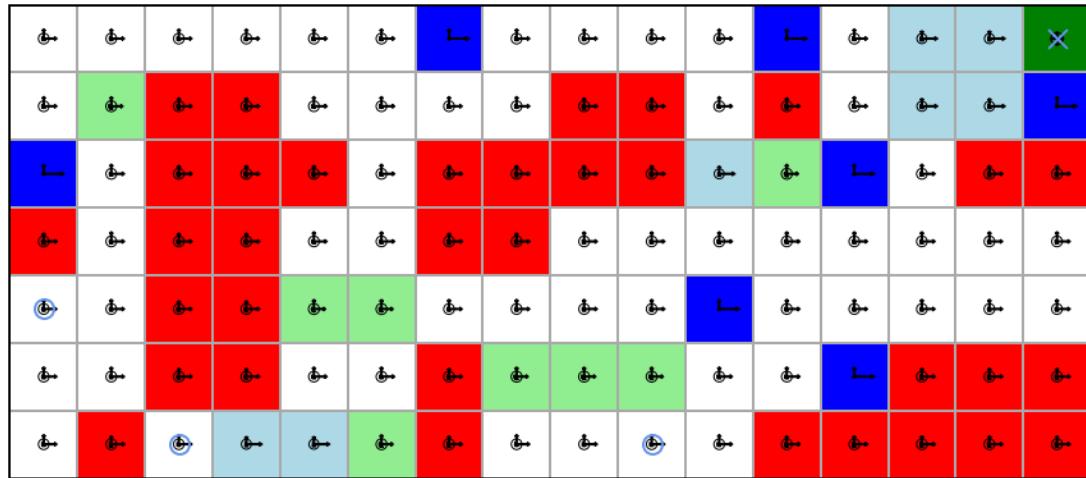


Figure 100: Graph of the classifier's learned policy.

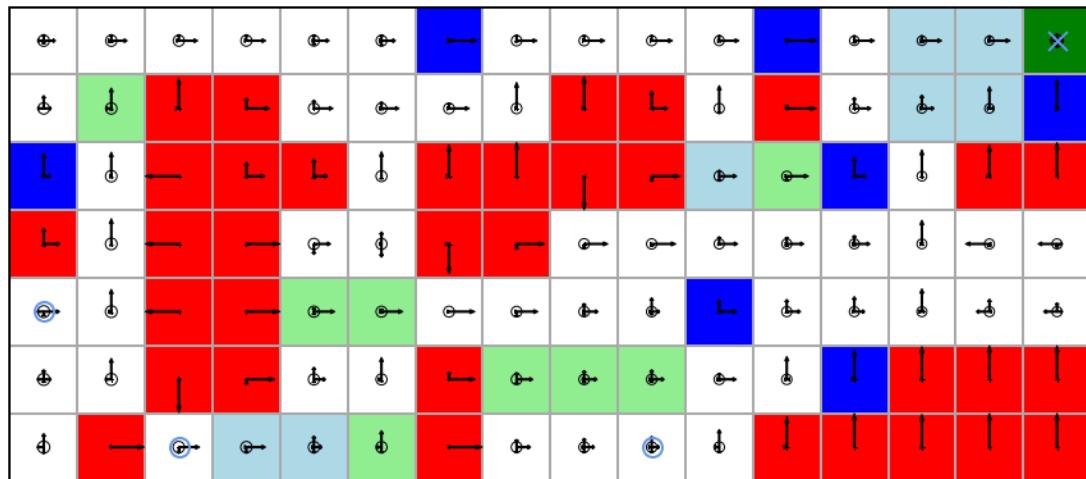


Figure 101: Graph of the expert's policy.

-1574.56	-1528.98	-1423.67	-1203.16	-1010.36	-935.71	-871.16	-853.79	-734.19	-507.47	-302.31	-182.61	-140.16	-90.96	-42.60	0.00
-2055.07	-2309.11	-2299.03	-1782.16	-1404.72	-1336.84	-1456.69	-1638.78	-1531.97	-1050.92	-814.34	-543.91	-389.28	-360.43	-333.13	-187.02
-2901.22	-3278.10	-3472.76	-2989.26	-2370.03	-2464.53	-2664.97	-2542.34	-2196.71	-1504.08	-1066.97	-1084.73	-1061.94	-1370.28	-1604.07	-1622.22
-3873.76	-4126.69	-4215.13	-3544.24	-3077.98	-3134.37	-3216.25	-2658.76	-2132.69	-1636.63	-1304.12	-1377.78	-1536.68	-1824.60	-2132.61	-2163.72
-4489.35	-4646.83	-4616.42	-3800.84	-3178.77	-3064.18	-3013.65	-2488.75	-1972.26	-1767.10	-1757.10	-1952.19	-2211.46	-2474.88	-2624.78	-2648.37
-4686.36	-4905.07	-4689.53	-3858.08	-3281.65	-3255.82	-2933.47	-2501.95	-2294.54	-2410.42	-2753.52	-3326.76	-4082.96	-4444.09	-4576.53	-4595.70
-4923.82	-4713.91	-4339.92	-3686.47	-3366.50	-3447.19	-3211.56	-2774.11	-2895.27	-3423.56	-4329.40	-5222.59	-5828.90	-6227.97	-6316.71	-6331.12

Figure 102: Classifier's state value map.

-27.58	-25.21	-22.93	-21.25	-19.71	-18.15	-11.61	-9.66	-7.97	-6.34	-4.70	2.22	3.96	5.94	8.06	0.00
-28.69	-26.99	-24.13	-19.99	-18.27	-16.20	-13.51	-10.98	-9.07	-7.22	-6.23	1.07	2.81	4.30	5.74	8.19
-30.03	-28.49	-29.21	-520.25	-18.91	-17.50	-14.68	-11.88	-11.20	-8.08	-6.89	-5.33	1.22	2.72	4.42	2.10
-32.96	-29.55	-30.25	-20.76	-19.96	-18.78	-17.10	-11.90	-10.30	-8.57	-6.89	-4.83	-2.44	0.83	-0.86	-2.82
-32.32	-30.97	-31.66	-20.09	-19.28	-17.66	-15.94	-14.73	-12.95	-11.76	-7.72	-5.55	-3.43	-1.07	-2.40	-4.02
-34.31	-32.55	-25.99	-21.95	-20.67	-19.31	-17.02	-16.18	-14.68	-12.91	-10.64	-7.52	-5.05	-2.31	-3.38	-4.98
-35.91	-26.05	-26.33	-23.91	-22.44	-21.27	-18.40	-17.54	-16.30	-16.08	-12.78	-8.96	-11.00	-503.53	-504.34	-505.93

Figure 103: Expert's state value map.

10.7 CMIRL Results on Gridworld 1

Trained on 10000 trajectories with batch size of 128 for 10 epochs.

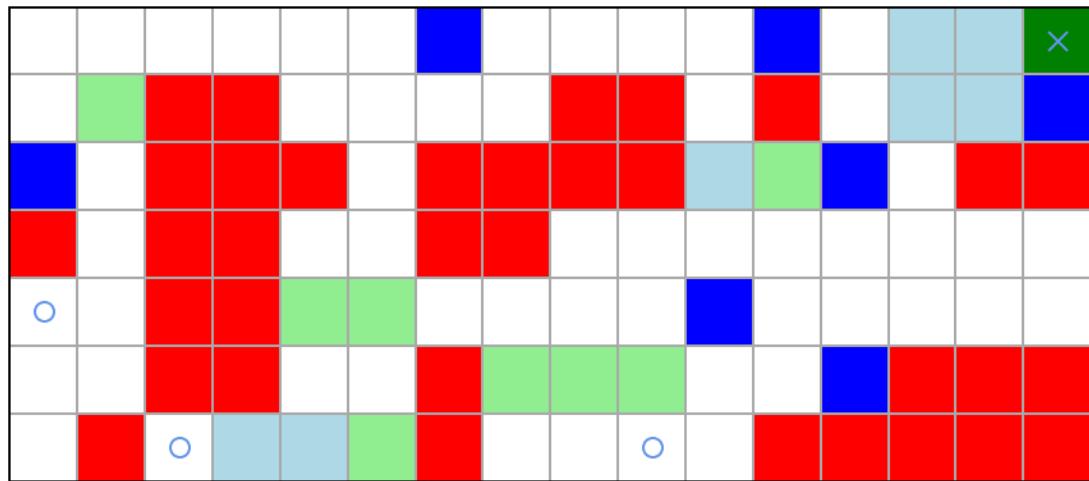


Figure 104: Gridworld 1

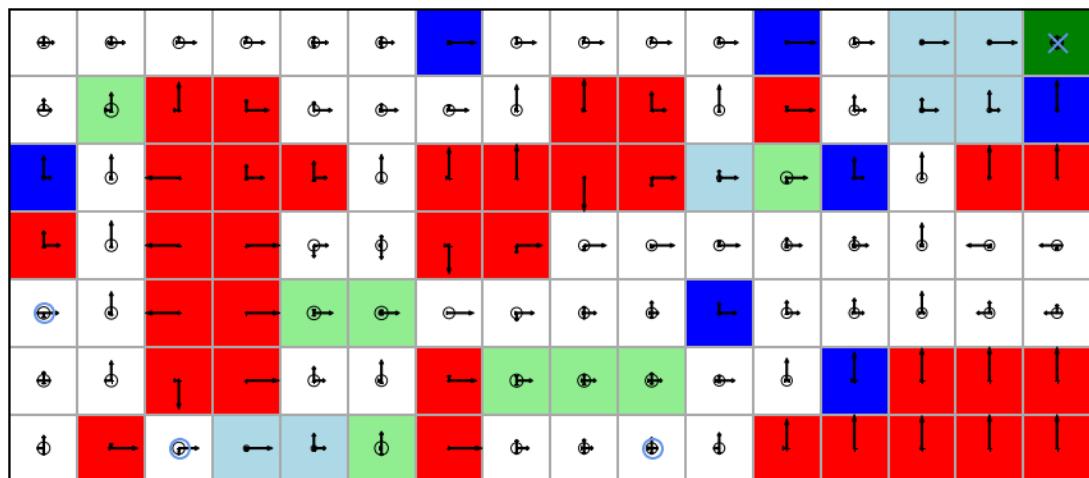


Figure 105: CMIRL Policy

Figure 106: MLIRL Policy

Figure 107: Expert Policy

-28.50	-25.84	-22.91	-20.97	-19.30	-17.61	-10.87	-8.80	-7.01	-5.27	-3.53	3.48	5.29	7.30	8.73	0.00
-29.64	-27.95	-24.51	-19.70	-17.89	-15.72	-12.81	-10.19	-8.17	-6.09	-4.96	1.86	3.80	5.04	5.66	8.36
-31.00	-29.47	-30.25	-520.27	-18.91	-17.37	-14.10	-11.10	-11.68	-7.63	-6.03	4.70	2.00	3.32	4.49	2.27
-33.87	-30.55	-31.54	-22.55	-21.75	-20.12	-19.12	-12.96	-10.71	-8.49	-6.49	-4.67	-2.23	1.16	-1.35	-8.05
-34.02	-32.10	-34.52	-22.32	-21.53	-19.91	-17.91	-16.58	-14.23	-12.37	-7.71	-5.53	-3.37	-0.98	-2.85	-6.42
-35.70	-33.84	-28.43	-23.49	-22.66	-21.64	-18.87	-17.99	-16.34	-14.15	-11.20	-7.61	-5.09	-2.27	-3.83	-7.37
-37.36	-28.79	-26.21	-24.96	-23.93	-23.35	-20.13	-19.25	-17.95	-16.02	-13.88	-9.33	-11.06	-503.54	-504.81	-508.31

Figure 108: CMIRL State Value Map

-45.06	-42.56	-39.68	-37.08	-33.01	-30.18	-21.80	-19.29	-17.04	-14.68	-11.65	0.47	3.63	5.78	8.00	0.00
-48.11	-50.56	-64.10	-39.67	-35.66	-32.83	-31.15	-30.11	-37.48	-49.04	-73.77	-3.53	1.69	3.61	5.10	7.98
-50.84	-52.90	-144.19	-547.27	-44.74	-44.54	-44.24	-46.95	-107.94	-50.47	-44.42	-34.54	-2.65	-3.83	-1.42	-6.13
-66.98	-58.52	-251.40	-72.72	-64.97	-55.94	-102.38	-41.51	-36.53	-34.09	-29.21	-21.28	-13.34	-12.89	-112.19	-245.81
-116.45	-81.36	-448.43	-64.58	-58.15	-51.57	-45.86	-41.17	-36.25	-34.05	-26.02	-21.14	-18.32	-21.81	-81.74	-166.89
-167.75	-175.15	-165.38	-93.97	-81.77	-93.15	-46.42	-39.31	-35.85	-32.77	-28.37	-24.10	-24.10	-26.44	-90.02	-174.62
-262.58	-125.96	-117.81	-116.19	-111.21	-146.65	-50.66	-39.72	-37.36	-36.00	-36.99	-32.25	-40.30	-533.81	-601.12	-685.13

Figure 109: MLIRL State Value Map

-27.58	-25.21	-22.93	-21.25	-19.71	-18.15	-11.61	-9.66	-7.97	-6.34	-4.70	2.22	3.96	5.94	8.06	X
-28.69	-26.99	-24.13	-19.99	-18.27	-16.20	-13.51	-10.98	-9.07	-7.22	-6.23	1.07	2.81	4.30	5.74	8.19
-30.03	-28.49	-29.21	-520.25	-18.91	-17.50	-14.68	-11.88	-11.20	-8.08	-6.89	-5.33	1.22	2.72	4.42	2.10
-32.96	-29.55	-30.25	-20.76	-19.96	-18.78	-17.10	-11.90	-10.30	-8.57	-6.89	-4.83	-2.44	0.83	-0.86	-2.82
-32.32	-30.97	-31.66	-20.09	-19.28	-17.66	-15.94	-14.73	-12.95	-11.76	-7.72	-5.55	-3.43	-1.07	-2.40	-4.02
-34.31	-32.55	-25.99	-21.95	-20.67	-19.31	-17.02	-16.18	-14.68	-12.91	-10.64	-7.52	-5.05	-2.31	-3.38	-4.98
-35.91	-26.05	-26.13	-23.91	-22.44	-21.27	-18.40	-17.54	-16.30	-16.08	-12.78	-8.96	-11.00	-503.53	-504.34	-505.93

Figure 110: Expert State Value Map

10.8 CMIRL Results on Gridworld 5

Trained on 10000 trajectories with batch size of 128 for 10 epochs.

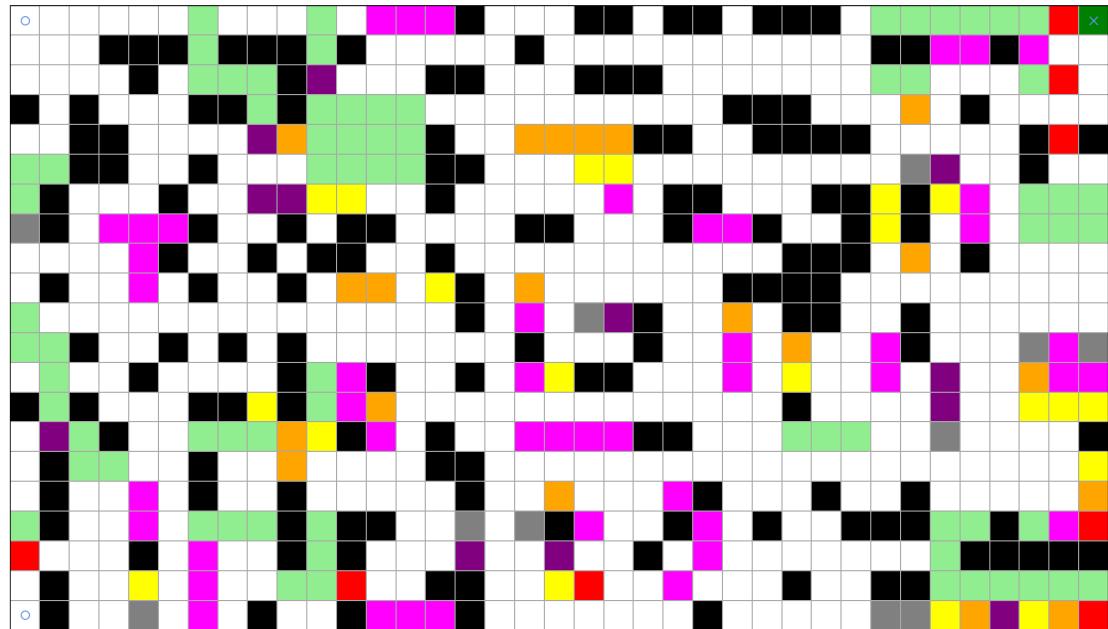


Figure 111: Gridworld 5

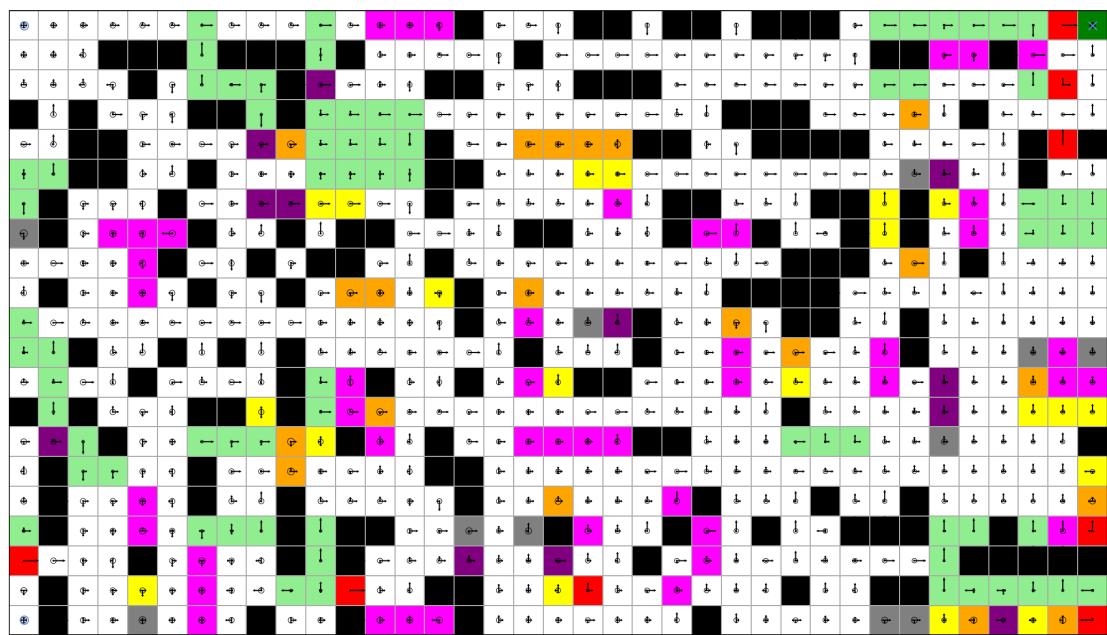


Figure 112: CMIRL Policy

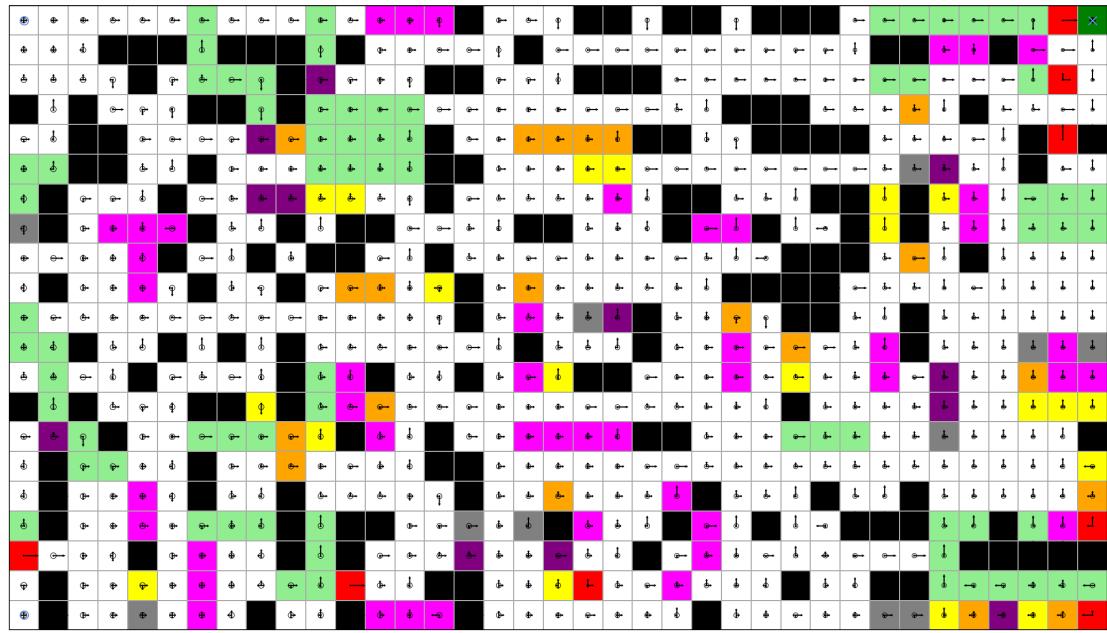


Figure 113: MLIRL Policy

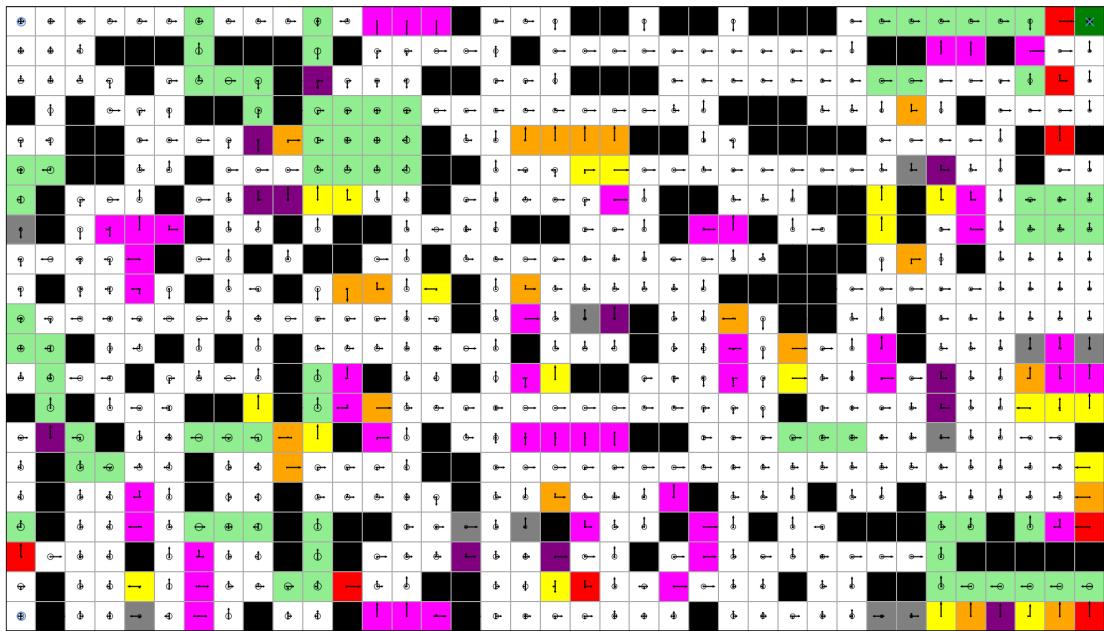


Figure 114: Expert Policy

0	122.95	128.95	125.95	129.27	129.15	120.85	121.21	122.75	122.94	124.14	125.72	121.07	122.45	71.25	42.65	48.23	45.49	49.20	38.66	38.09	36.98	35.75	31.85	32.65	32.39	27.06	24.22	21.75	21.31	20.89	17.07	15.63	15.13	11.13	0	
120.65	121.01	122.81	125.76	120.08	128.93	129.65	120.45	121.21	120.95	120.81	120.15	121.77	121.21	100.11	95.85	95.21	56.48	44.30	39.01	37.93	37.08	36.11	35.87	31.63	32.79	31.61	30.25	28.00	24.12	27.71	23.80	19.01	11.29	8.27	-2.97	3.48
120.48	121.01	122.47	127.15	121.50	123.51	129.45	124.51	121.45	120.76	117.45	116.78	121.88	121.47	116.15	119.54	124.48	94.96	74.57	83.42	41.54	56.92	35.91	34.84	33.62	32.49	31.43	30.40	27.31	22.40	18.93	16.53	13.94	12.13	41.44	-4.01	2.88
120.47	120.81	127.87	120.63	125.42	129.22	129.74	127.81	129.13	121.97	117.10	116.54	121.35	118.36	118.78	118.93	116.68	110.37	90.83	88.05	45.96	38.02	36.73	35.81	34.48	31.71	35.97	35.94	35.94	36.12	24.43	18.97	12.24	-10.49	8.46	-3.75	4.26
122.78	120.93	127.63	128.08	128.05	125.97	126.41	126.83	128.14	120.42	117.72	115.78	116.81	120.85	127.44	129.61	122.08	107.89	88.16	45.87	42.09	36.24	35.18	33.84	32.36	31.80	35.18	31.18	29.03	25.18	16.58	13.12	11.72	8.60	8.79	5.24	
137.68	124.48	100.35	156.33	144.19	137.33	138.65	132.49	132.79	129.98	119.94	113.05	109.44	109.87	120.80	119.79	118.98	114.01	89.82	77.60	82.28	37.44	34.89	31.83	32.69	31.45	30.38	29.39	28.46	27.53	23.95	18.59	13.12	13.07	17.25	28.74	17.14
157.40	147.89	169.48	171.60	168.55	126.97	131.53	131.17	130.51	127.87	119.08	108.80	109.35	108.87	106.13	106.93	99.30	85.22	71.67	66.03	42.13	36.17	35.24	33.81	32.42	31.49	30.38	33.66	30.38	29.48	25.22	19.75	15.55	17.68	26.02	18.24	
158.16	145.43	149.59	171.42	172.93	179.57	122.3	131.81	132.20	127.89	120.00	116.52	104.59	123.54	123.23	123.29	123.43	96.70	72.24	64.48	57.54	47.75	50.31	44.84	37.31	32.73	32.60	40.47	40.93	34.07	32.44	25.05	19.92	19.66	24.01	18.81	
158.15	140.30	141.59	116.33	144.05	105.36	139.42	133.05	124.81	124.79	126.78	115.90	125.58	125.53	96.95	93.72	87.46	74.18	62.31	63.54	58.00	50.90	49.61	48.11	43.63	44.38	33.40	45.66	52.88	42.30	40.81	34.45	29.93	20.27	21.22	23.02	20.50
149.35	141.94	154.23	154.42	154.71	131.71	134.94	137.28	138.14	144.81	145.37	140.87	127.80	121.84	127.14	91.45	89.53	71.97	45.37	69.21	54.73	54.14	54.85	53.65	61.63	42.17	49.68	61.21	58.79	49.32	46.63	31.45	26.87	21.89	22.94	24.17	23.81
148.71	147.44	148.11	145.91	142.16	136.21	135.24	136.29	116.95	117.46	114.64	114.64	122.50	120.30	108.30	88.04	87.60	41.06	47.26	41.93	58.49	66.19	73.84	83.00	94.24	71.09	40.57	32.96	51.88	31.16	41.00	34.50	28.02	24.21	25.84	27.74	27.01
146.18	146.77	149.36	144.06	142.04	136.03	135.99	135.44	136.22	120.36	127.89	122.23	120.44	97.61	92.21	89.64	89.24	40.04	47.83	43.84	42.49	37.91	47.73	52.01	42.09	32.60	47.60	48.88	48.13	34.45	38.47	34.21	30.46	29.29	34.71	37.52	37.71
144.13	145.58	145.82	146.35	145.10	139.49	136.81	139.28	140.16	131.91	119.49	112.67	100.59	98.30	87.48	87.61	85.14	80.52	75.64	79.17	81.44	91.31	92.61	88.13	89.09	70.88	40.88	39.31	36.66	44.04	38.95	36.28	41.54	48.18	34.09	30.98	
159.16	163.28	160.86	163.18	181.38	131.56	143.93	159.05	159.28	159.00	140.05	128.10	121.59	120.87	106.29	107.59	106.98	105.70	109.16	96.72	92.82	91.42	91.14	90.74	89.74	82.04	68.05	60.15	58.60	56.36	46.89	42.81	42.99	50.92	64.11	65.87	63.95
234.49	234.88	232.36	215.17	218.12	196.03	272.63	170.48	174.44	185.66	146.68	128.51	127.54	124.89	111.87	113.69	115.98	114.98	110.64	104.57	97.75	92.20	89.49	88.11	85.79	78.18	41.78	59.21	36.79	33.76	49.41	48.06	48.68	56.93	49.37	71.50	71.98
238.70	233.31	216.28	255.65	250.76	240.23	270.11	168.37	169.17	156.70	137.57	127.39	126.73	121.81	116.19	117.38	117.70	114.35	120.43	99.80	87.70	85.34	84.31	79.35	69.89	59.52	64.97	55.12	53.34	51.03	51.03	52.71	59.18	49.70	73.27	63.23	
226.54	257.84	259.39	262.65	263.61	262.28	196.49	169.35	168.81	149.65	320.40	157.27	121.98	122.21	121.63	121.57	122.97	118.70	109.91	99.35	91.33	87.70	81.30	75.57	68.33	42.40	68.75	55.28	64.13	52.64	53.04	55.52	41.39	49.46	37.08	33.78	
265.30	219.10	240.73	216.29	267.02	169.09	441.12	168.31	177.52	131.97	410.11	120.98	121.93	121.40	121.88	121.27	122.97	123.40	121.98	105.11	98.03	92.80	89.91	89.57	76.65	48.07	42.13	43.49	15.82	34.43	34.11	33.81	43.81	27.47	27.72	22.63	
239.13	200.57	242.56	205.95	275.10	276.02	200.39	226.23	180.39	329.07	120.31	125.50	124.93	124.65	122.61	122.01	121.81	111.30	120.57	98.82	93.28	90.71	84.26	77.57	48.96	43.26	43.32	40.15	56.13	54.12	53.82	55.70	73.49	22.29	31.94	22.09	
267.49	243.27	263.36	270.75	274.49	278.19	300.79	375.16	260.94	232.08	148.28	116.06	122.89	123.37	121.81	121.62	118.14	105.27	108.39	94.83	91.31	86.81	72.77	40.17	44.95	44.46	41.63	56.51	54.98	46.19	42.93	48.90	36.44	37.34			
243.95	245.56	269.33	273.36	276.79	281.19	280.18	269.81	185.64	161.08	150.33	145.72	144.22	151.99	120.62	119.84	118.06	113.20	101.37	98.04	98.97	98.58	85.04	77.63	74.42	72.36	70.37	74.81	41.51	41.16	72.21	41.11	106.44	114.04	117.59	120.00	

Figure 115: CMIRL State Value Map

65	85.81	84.64	43.36	37.70	32.35	32.32	30.43	30.13	30.78	41.50	43.57	38.29	31.99	33.31	43.37	41.01	37.84	36.43	34.57	33.45	32.83	31.61	30.23	29.46	27.98	26.90	21.65	20.58	19.77	19.27	18.73	17.79	-16.31	15.76	4.13	50	
66.99	86.55	85.81	44.88	43.02	43.23	42.64	43.38	33.63	37.32	76.51	75.10	69.21	63.18	55.01	51.94	51.60	37.02	35.20	33.58	32.65	31.75	30.74	29.51	28.37	27.05	25.74	24.31	22.54	20.90	25.00	22.41	18.68	-11.58	6.37	2.94	3.47	
67.24	87.53	89.02	42.02	41.72	39.89	37.14	40.43	30.76	37.79	-71.39	-49.85	-49.18	-49.28	-47.64	-59.81	-45.65	-49.93	-42.47	-57.22	-54.40	-31.73	-30.74	-29.55	-28.29	-27.11	-20.00	-24.93	-23.02	-21.22	-19.47	-17.28	-14.90	-22.32	-11.63	4.04	2.88	
67.67	87.66	91.59	33.22	43.20	42.87	42.59	33.07	33.78	37.03	-70.49	-49.47	-49.42	-39.08	-39.50	-49.15	-45.21	-39.37	-49.97	-41.31	-35.29	-32.69	-31.67	-30.85	-29.25	-28.14	-29.65	-29.84	-30.42	-32.01	-23.45	-19.46	-13.15	-11.38	-8.32	-0.82	4.26	
69.68	85.25	42.36	41.41	53.33	33.15	43.06	32.81	42.04	37.51	-72.35	-40.32	-40.58	-33.57	-39.86	-42.21	-45.86	-31.35	-72.10	-40.11	-46.30	-35.49	-33.11	-32.98	-34.07	-32.93	-32.54	-30.71	-30.07	-29.87	-25.80	-17.61	-14.10	-12.60	-10.44	-6.86	0.24	
72.25	89.41	104.63	44.94	53.84	33.55	39.58	39.60	37.68	43.07	-75.54	-71.38	-70.95	-71.24	-42.70	-42.76	-43.34	-40.48	-76.51	-69.68	-49.27	-37.32	-34.98	-34.01	-31.38	-31.98	-30.92	-29.94	-29.01	-28.06	-24.45	-19.23	-15.87	-13.92	-12.47	-28.36	-17.22	
102.11	104.45	106.16	112.69	97.24	22.17	91.66	91.52	90.31	88.42	43.52	71.12	75.48	77.03	79.33	79.39	79.11	76.01	70.46	64.81	52.59	40.98	38.18	35.23	34.15	32.93	32.01	30.92	34.19	30.88	30.04	25.35	19.46	-16.20	18.94	-23.65	18.72	
111.17	114.32	114.37	112.33	108.04	110.30	82.05	92.03	82.42	80.83	90.74	49.72	49.25	79.27	73.31	70.79	78.44	75.18	64.00	59.83	54.87	48.40	35.56	49.64	44.21	37.91	34.09	33.11	37.39	40.05	39.95	34.49	32.51	25.45	19.14	20.33	22.20	19.68
113.70	115.51	116.08	118.08	117.18	15.55	42.03	92.62	93.08	105.41	88.77	82.93	45.31	30.67	77.00	35.12	72.55	64.08	38.06	56.41	62.53	39.11	48.30	47.21	43.13	44.13	33.88	48.13	50.89	39.54	38.03	34.07	29.18	20.52	21.60	22.59	20.99	
111.77	-114.16	115.42	118.40	114.43	107.71	96.31	96.28	331.87	-113.53	117.13	109.60	43.73	46.76	48.11	37.47	39.03	47.48	49.52	55.85	42.86	33.04	50.01	49.83	46.57	53.11	46.88	49.74	49.20	43.40	-43.31	34.10	26.12	-22.30	-23.24	-23.81	-22.82	
110.32	-111.61	112.71	111.47	109.38	103.11	101.79	101.46	104.13	104.81	104.93	102.60	45.29	46.85	46.43	47.81	37.85	27.51	21.49	42.44	61.90	45.24	16.55	43.07	66.85	41.95	27.09	62.23	49.28	48.55	39.89	32.08	-22.27	-24.42	-25.28	-25.84	-24.97	
109.61	110.11	111.37	110.54	108.76	102.71	102.13	102.27	104.95	111.19	99.81	91.96	48.76	43.74	30.94	29.21	29.31	42.31	49.75	58.58	75.29	74.60	80.00	41.13	25.21	41.25	54.88	31.13	31.75	32.73	29.24	-27.87	-30.47	-30.89	19.99			
108.48	109.89	111.56	112.42	112.57	110.20	100.71	111.14	112.02	104.55	104.07	102.99	91.05	45.80	34.69	44.24	44.47	30.31	70.35	47.14	30.97	79.98	81.86	78.13	73.06	43.78	36.14	34.68	30.81	35.01	33.38	35.94	40.63	31.68	39.67			
113.21	117.27	120.38	121.93	129.10	129.59	123.79	124.92	133.23	119.80	112.87	112.04	98.66	90.83	89.99	90.79	89.54	89.04	49.82	43.20	30.05	93.76	93.73	70.37	76.71	71.51	40.82	54.61	52.75	49.82	42.60	39.55	38.57	43.99	36.43	35.69	32.53	
131.18	129.90	144.39	119.48	136.63	139.70	139.77	141.18	143.36	133.33	117.52	112.39	104.32	94.14	83.99	95.42	97.09	46.32	42.64	48.06	32.79	78.37	73.02	71.06	66.62	59.61	53.57	52.89	31.09	-48.21	-44.84	43.58	43.22	-48.85	41.42	-43.01	0.96	
131.46	144.24	146.14	147.17	143.95	140.68	139.77	140.36	141.08	128.39	126.22	105.41	104.17	27.64	101.03	86.78	87.23	86.84	43.89	75.98	88.63	66.77	45.64	61.42	56.76	53.43	32.08	20.40	48.39	46.43	46.31	46.68	52.01	40.72	43.94	11.21		
133.11	152.38	153.16	153.98	151.04	148.62	149.87	159.43	149.63	121.41	108.03	104.96	102.63	101.21	102.24	99.98	99.52	100.90	93.47	30.80	75.37	72.15	68.93	43.73	59.77	56.77	54.51	52.15	39.66	39.32	47.75	48.09	49.54	54.20	43.34	45.56	28.71	
138.11	157.80	159.47	159.04	156.65	150.72	145.81	140.21	139.14	-111.14	107.74	104.53	102.57	102.21	102.43	101.81	100.44	100.66	95.74	45.66	76.65	73.32	72.01	44.23	46.98	57.17	55.30	64.74	61.25	49.80	49.17	48.58	49.93	55.69	43.45	47.22	45.16	
160.89	161.80	162.54	164.02	164.04	159.32	153.24	144.84	139.77	113.51	107.29	107.68	103.41	102.86	102.93	101.23	99.48	47.63	49.05	44.44	75.53	75.12	24.42	47.65	42.31	58.60	56.16	45.97	42.68	50.20	49.40	49.11	49.09	66.18	43.37	22.51	11.60	
174.55	165.61	167.37	170.54	170.96	166.82	162.26	151.32	126.96	116.50	108.16	107.01	107.29	105.63	103.81	100.58	88.42	90.07	94.10	83.59	30.79	79.80	74.21	70.11	62.55	40.07	37.54	47.61	16.55	14.10	31.45	50.01	56.67	44.21	71.11	44.03	43.83	
147.79	167.27	167.89	169.69	171.51	169.97	168.77	160.34	137.72	113.13	112.23	121.47	118.81	115.54	112.38	88.53	86.16	84.49	80.24	82.61	80.00	30.34	79.26	48.87	41.67	40.37	40.41	40.76	38.08	38.86	41.62	42.03	47.77	48.09				

Figure 116: MLIRL State Value Map

65	59.44	16.23	33.24	51.59	39.39	49.77	48.93	47.63	46.50	45.75	46.53	41.83	39.52	36.11	32.87	31.07	30.77	29.23	27.70	26.55	25.80	24.51	23.05	22.23	20.51	19.77	15.47	14.48	13.60	12.95	12.24	11.49	10.73	9.97	1.06	50
42.59	45.64	58.40	54.74	42.19	43.35	49.76	49.48	45.31	45.46	45.00	42.08	46.75	48.54	33.10	34.00	33.61	29.21	27.07	24.68	35.68	24.49	23.83	22.20	20.61	19.69	18.31	17.03	15.73	14.20	13.60	13.61	12.29	9.82	3.68	2.60	1.46
42.59	43.57	59.08	64.31	50.51	49.78	49.59	49.06	48.68	44.78	47.07	41.51	40.74	39.36	38.35	38.61	36.30	35.28	31.93	27.64	24.66	23.68	22.20	20.61	19.46	18.11	16.75	15.15	13.92	13.21	12.16	10.80	9.73	-8.16	-3.76	2.86	
44.09	43.41	35.64	50.95	49.85	49.17	47.08	46.86	48.08	41.81	41.49	41.08	49.23	38.31	34.52	32.81	31.81	30.17	28.85	27.58	26.57	25.58	24.56	23.54	21.83	20.46	18.47	17.38	16.15	14.89	14.08	13.54	9.61	8.34	4.87	3.39	4.23
44.58	44.53	44.84	49.70	48.96	47.76	45.42	41.77	41.49	41.21	40.73	39.31	35.65	33.61	32.71	31.15	29.59	28.33	27.39	26.49	25.41	24.75	24.42	22.74	22.07	18.45	16.40	15.14	13.73	12.48	10.92	9.59	7.90	6.40	5.21		
64.33	64.31	43.62	50.66	49.79	49.05	44.25	43.11	42.68	41.66	41.40	40.77	40.98	34.76	34.74	30.78	30.80	30.20	27.																		

10.9 CMIRL Results on Gridworld 6

Trained on 10000 trajectories with batch size of 128 for 10 epochs.

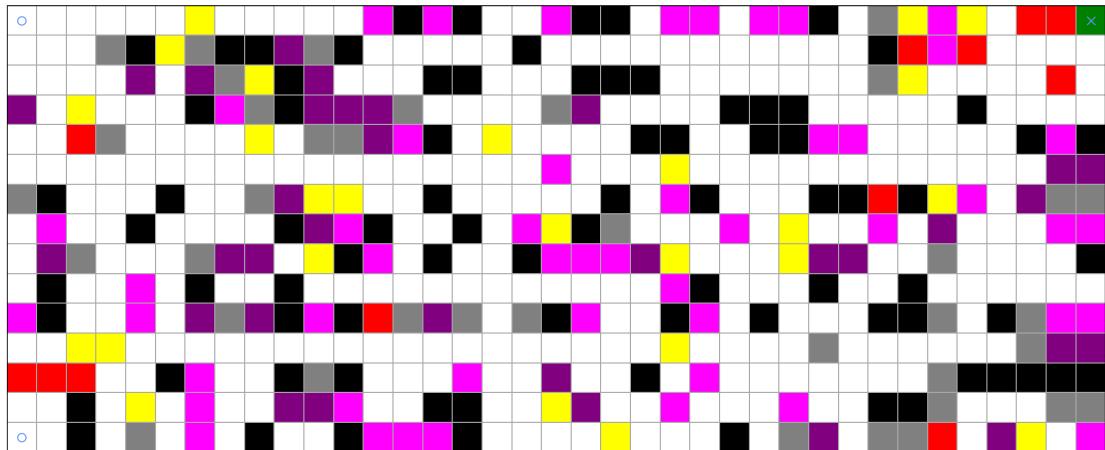


Figure 118: Gridworld 6

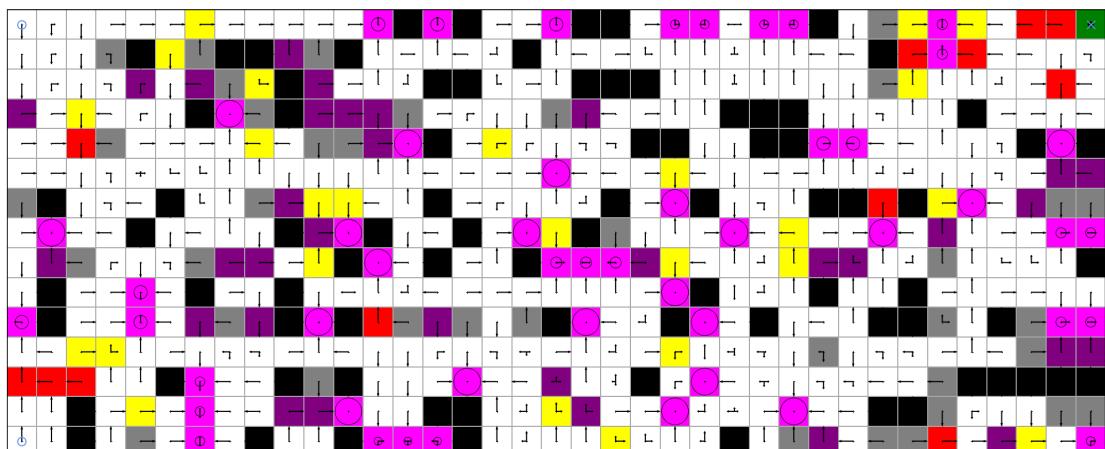


Figure 119: CMIRL Policy

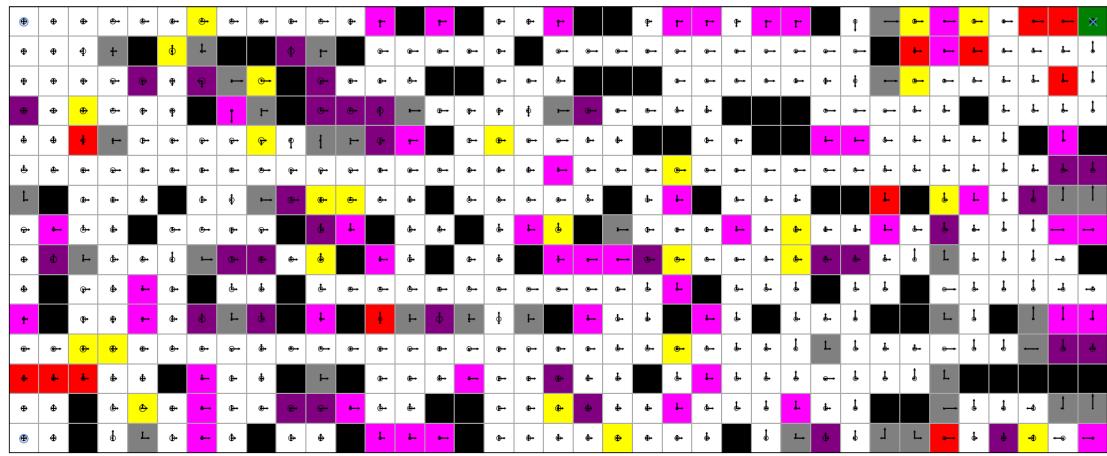


Figure 120: MLIRL Policy

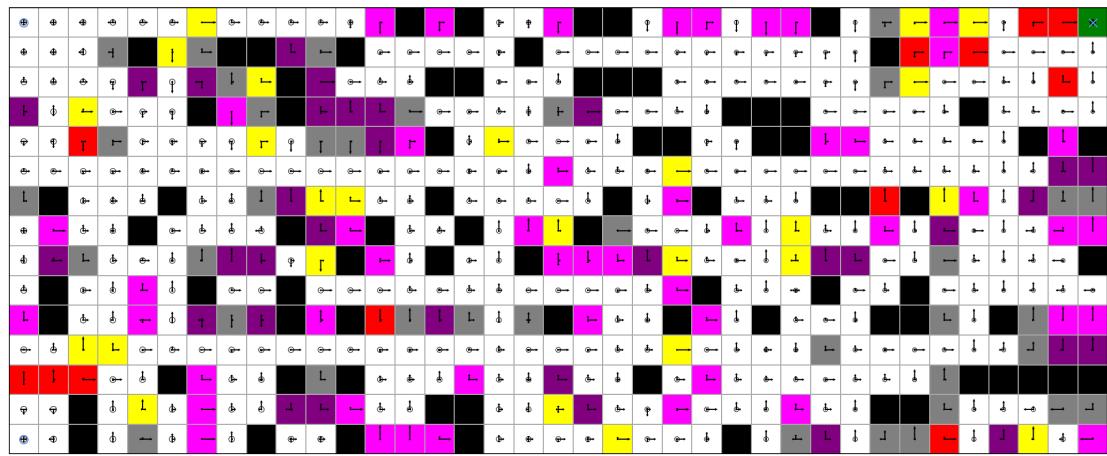


Figure 121: Expert Policy



Figure 122: CMIRL State Value Map

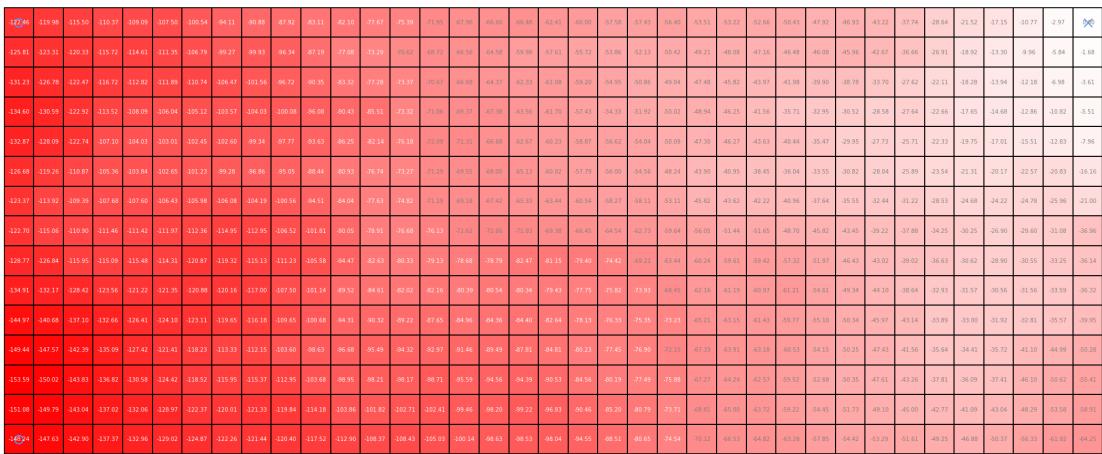


Figure 123: MLIRL State Value Map

79.34	77.82	76.25	47.17	41.47	33.85	53.67	52.33	50.88	49.77	49.39	49.75	42.60	41.76	39.73	38.78	37.58	35.31	31.80	30.57	29.68	28.23	26.55	25.75	24.38	23.31	22.39	21.71	20.69	21.67	15.38	9.37	8.12	4.79	3.20			
79.62	78.64	76.79	70.61	44.27	39.65	58.32	54.51	53.28	52.03	47.99	43.60	42.38	41.18	39.91	38.80	37.43	32.90	31.58	30.47	29.77	28.18	26.81	25.65	24.03	23.05	22.89	20.81	20.20	22.21	37.53	31.94	6.87	-5.67	-4.28	-2.96	3.48	
78.71	76.39	71.42	43.90	57.53	56.51	57.81	58.35	55.18	50.71	49.42	44.85	43.94	43.76	41.49	39.24	39.03	39.39	33.63	31.79	31.16	27.94	26.61	25.25	23.67	22.39	20.81	19.37	19.90	16.51	11.07	-9.87	-8.60	-7.24	4.22	4.06	3.90	
73.94	72.63	42.85	57.97	56.42	53.82	54.11	52.73	54.68	50.43	50.98	46.33	45.31	41.72	40.64	39.56	38.07	37.20	35.69	31.67	30.40	29.22	27.80	26.97	25.57	23.11	29.09	27.73	18.46	25.18	13.77	12.24	8.70	-8.48	-7.34	-6.04	4.33	
44.62	44.37	40.22	56.79	55.57	44.32	42.51	51.73	49.68	49.00	40.07	41.16	48.29	45.13	42.11	41.59	37.51	36.37	31.03	31.89	32.48	33.88	27.73	26.51	25.51	24.02	22.58	23.95	5.71	13.96	14.61	13.28	11.69	10.06	8.76	7.47	5.55	
42.43	40.45	58.33	56.96	55.32	53.63	51.98	50.21	48.75	47.75	46.75	40.86	44.93	43.71	42.46	41.21	39.83	38.94	35.62	38.55	33.41	32.29	25.63	24.53	23.22	21.77	20.44	19.19	18.03	16.71	15.58	14.28	12.95	11.66	12.96	13.77	8.70	
43.42	40.09	58.95	57.96	56.87	14.27	52.90	51.88	50.29	48.88	48.34	46.72	46.13	45.55	42.50	41.29	40.05	38.68	37.43	38.11	34.58	33.74	29.52	25.24	24.14	23.11	22.16	20.83	19.41	18.06	16.84	15.49	13.04	12.91	13.99	18.09	14.04	
47.65	40.43	59.75	59.28	57.24	15.46	43.91	53.28	54.41	54.75	55.26	51.21	47.28	46.88	46.41	42.23	41.28	39.95	38.73	37.34	33.55	32.42	31.11	29.84	25.51	24.44	23.81	24.88	23.94	22.84	21.29	18.10	16.23	14.31	16.69	19.36	19.68	
48.66	45.81	48.51	59.86	49.69	46.87	25.09	54.40	56.09	56.18	52.58	42.32	49.23	48.35	46.97	44.37	42.64	44.59	44.51	43.81	48.55	33.80	30.52	29.23	27.51	25.69	28.73	26.21	24.98	24.00	22.55	28.63	17.16	15.83	17.49	19.04	21.00	
48.58	43.53	42.31	40.99	59.29	27.99	38.40	57.98	55.92	51.66	50.75	49.97	49.05	48.05	46.66	45.22	44.24	44.22	42.96	43.91	49.85	39.46	34.03	29.65	28.49	27.28	30.52	29.39	16.89	23.34	21.45	19.76	18.36	17.34	19.57	19.86	21.61	
68.63	43.82	42.75	43.90	40.89	59.26	40.11	58.54	57.42	56.26	52.49	51.52	50.53	49.92	48.15	46.36	45.79	46.01	44.91	42.77	41.87	40.90	39.11	31.53	30.03	29.39	31.22	29.96	29.46	26.89	25.55	20.71	19.70	18.63	19.73	21.03	22.78	
49.44	48.08	43.66	42.05	40.78	59.77	48.88	37.94	56.74	55.56	54.50	53.66	52.57	51.26	49.79	48.42	47.22	46.35	44.99	43.85	42.57	41.18	34.40	33.35	32.23	31.54	32.78	29.78	18.17	26.24	24.35	22.78	21.07	22.48	23.99	27.01	19.18	
73.43	70.22	44.47	43.35	42.11	41.21	49.50	58.53	58.52	58.27	37.05	56.14	54.07	53.25	42.31	31.36	49.38	48.68	48.24	45.68	44.81	44.11	41.92	40.50	34.61	33.53	31.10	32.11	30.34	18.77	27.22	25.91	24.11	22.19	21.65	20.30	31.25	24.59
78.17	76.36	46.51	44.01	43.68	45.76	40.74	59.87	59.41	45.19	40.81	55.00	54.01	53.53	53.07	51.59	50.41	50.15	50.38	45.34	44.30	43.03	38.76	31.61	35.31	35.17	32.51	31.19	30.29	28.58	27.11	28.31	26.07	27.71	29.49	71.02	37.01	
75.55	49.24	45.95	47.49	46.88	41.94	41.02	41.28	45.64	45.15	40.91	34.84	54.49	54.55	52.89	51.60	51.02	50.53	48.90	42.44	31.21	39.89	38.51	30.62	31.46	37.67	32.22	31.45	33.60	31.46	29.42	27.06	28.32	30.79	38.15	37.52		

Figure 124: Expert State Value Map

10.10 CMIRL Results on Gridworld 7

Trained on 10000 trajectories with batch size of 128 for 10 epochs.

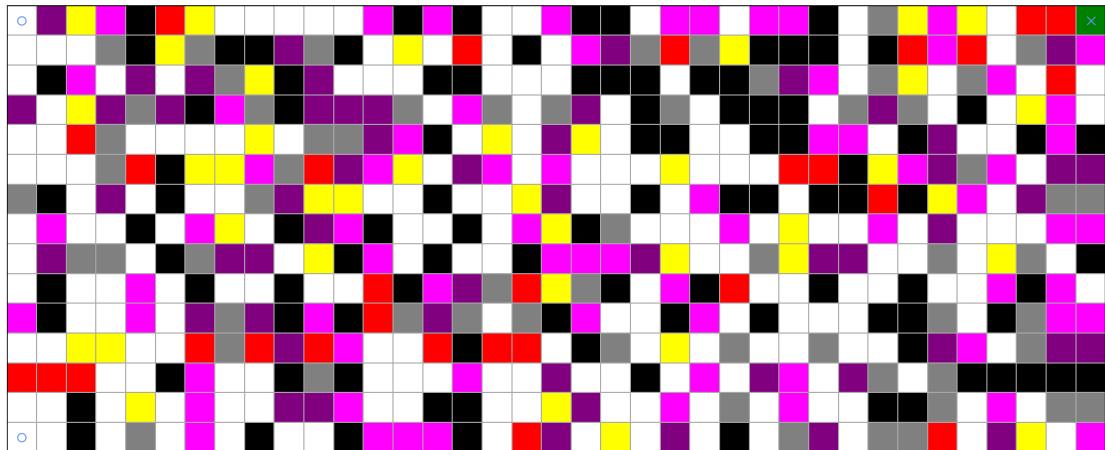


Figure 125: Gridworld 7

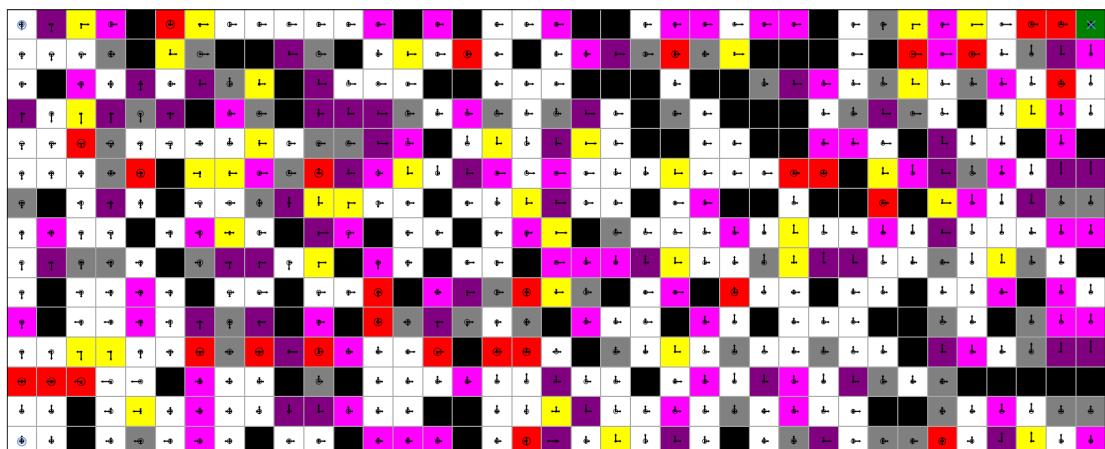


Figure 126: CMIRL Policy

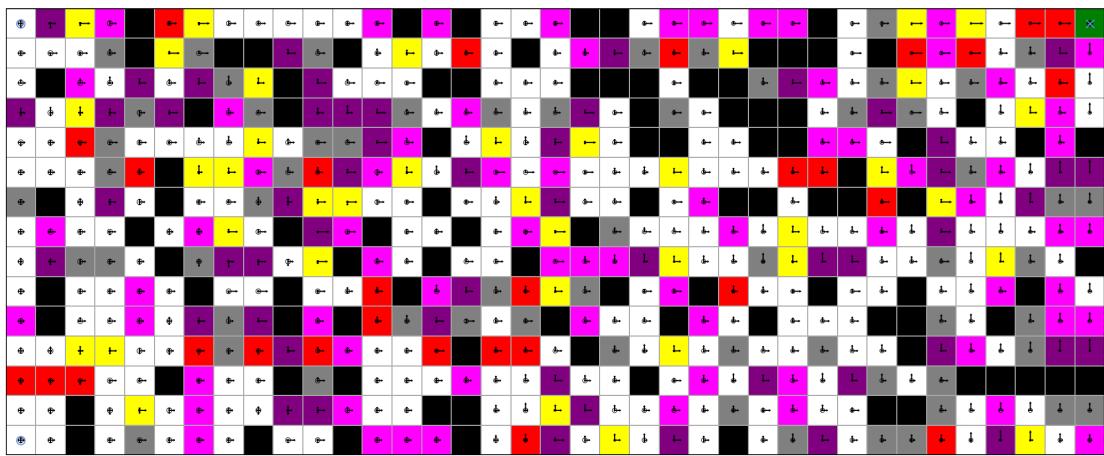


Figure 127: MLIRL Policy

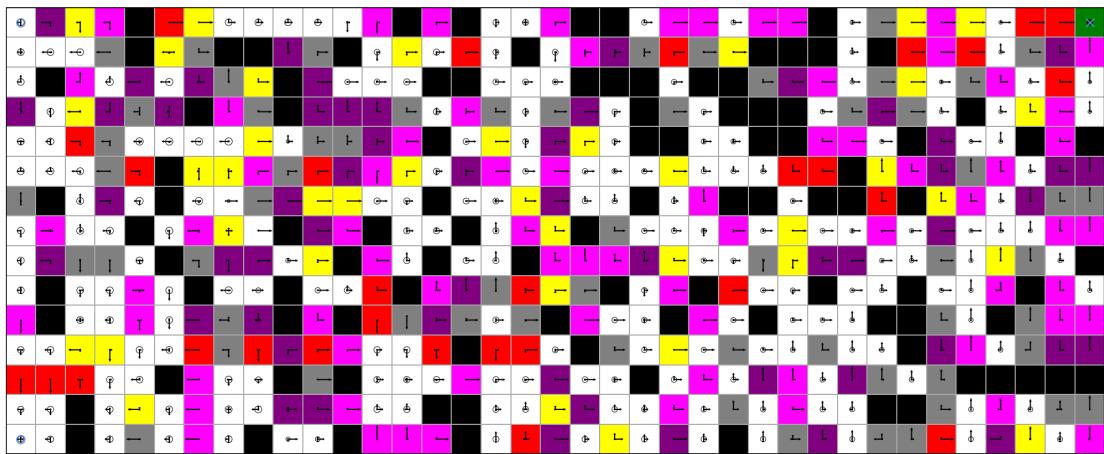


Figure 128: Expert Policy

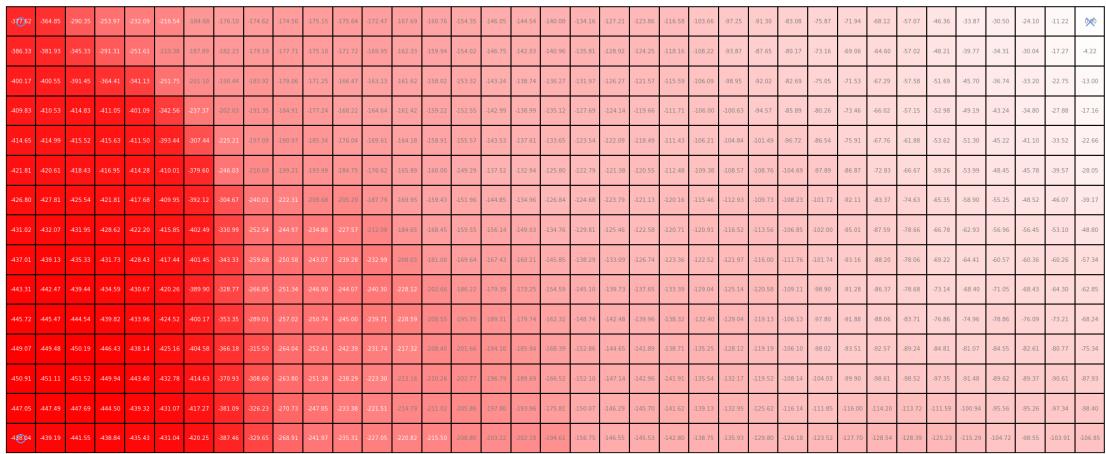


Figure 129: Classifier State Value Map

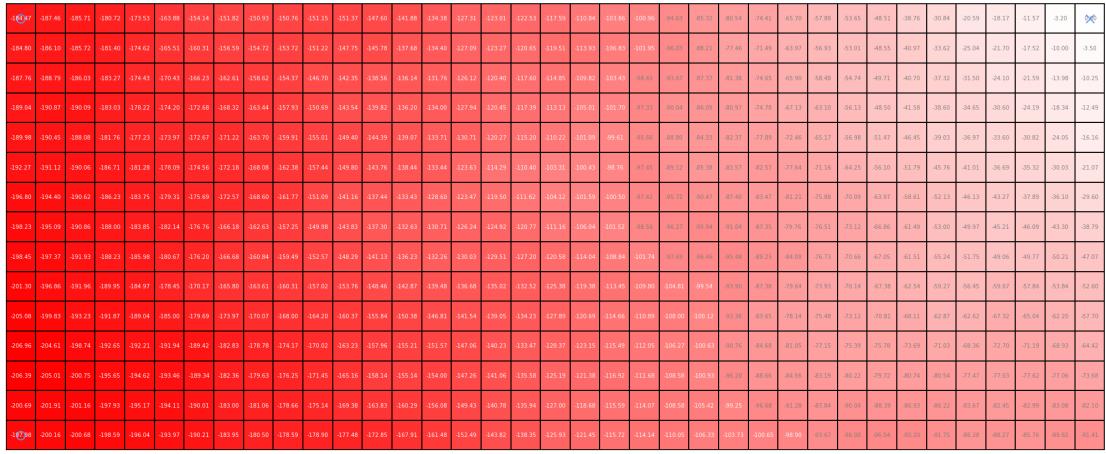


Figure 130: MLIR State Value Map

109.96	-106.31	109.62	111.50	114.08	99.16	93.94	93.39	93.94	92.55	91.75	89.69	88.24	89.65	83.92	80.28	79.57	80.25	75.27	73.19	67.65	66.35	40.78	55.03	54.35	40.58	42.39	55.68	17.45	32.23	26.45	20.48	14.71	-13.68	7.27	3.13	X			
109.97	-106.30	109.41	109.66	110.78	109.87	97.68	94.94	93.64	93.15	90.16	86.76	86.84	85.03	84.54	79.02	78.28	76.93	76.31	75.08	71.16	67.33	62.14	57.31	51.57	47.63	45.80	39.88	38.87	35.25	-29.01	-22.83	-16.70	-15.57	12.59	7.23	3.14			
109.98	-106.58	-109.52	111.41	111.76	115.20	103.39	97.41	94.50	90.57	90.24	85.94	85.57	85.09	81.37	77.35	76.82	75.99	75.11	71.09	47.22	43.47	42.73	59.03	59.97	61.19	-40.49	-49.40	79.99	34.58	28.71	27.77	22.70	-16.69	-15.90	9.18	3.55			
109.99	-107.24	-129.69	110.12	112.71	118.95	117.40	116.19	102.21	97.85	94.68	90.85	86.99	47.02	45.96	25.64	41.08	76.85	76.19	72.93	47.75	47.01	43.29	48.32	57.31	31.17	49.17	40.15	44.28	39.35	34.14	29.23	28.21	23.38	-22.59	-26.58	10.20	8.85		
110.00	109.47	-109.45	109.83	114.61	118.60	-118.59	118.15	116.41	105.03	99.55	86.91	83.46	81.14	86.33	82.56	81.46	75.83	75.12	72.67	47.14	48.50	42.62	38.54	37.78	36.88	32.38	40.03	40.45	57.11	-35.21	32.11	26.78	25.36	-21.96	21.47	-14.59	5.02		
110.01	109.47	-109.74	110.01	114.79	-117.98	116.69	109.78	103.76	99.91	94.88	90.55	89.66	84.50	84.19	79.87	74.07	73.43	47.54	68.96	68.17	65.10	59.24	58.38	57.88	57.38	32.04	-37.11	33.00	-37.52	-36.98	31.43	-26.58	-25.13	-25.18	-19.79	14.12			
110.02	110.77	-110.18	110.59	111.52	116.07	-113.96	115.11	116.85	99.65	95.35	90.16	89.04	84.71	44.15	40.34	79.32	78.65	73.06	48.31	47.71	47.36	45.18	44.44	59.48	58.88	66.26	33.08	50.99	27.51	43.97	41.88	36.69	-31.04	-31.21	26.38	24.44	-19.21		
110.03	112.41	-112.84	112.09	113.23	113.81	112.81	113.36	110.74	103.49	98.35	93.80	88.31	44.82	44.35	-83.79	80.17	-79.71	-77.72	-72.19	48.49	48.20	43.68	42.61	40.95	56.73	55.37	49.43	-49.49	27.48	41.07	-40.04	34.62	33.61	-32.27	-31.72	-29.37	-23.08		
110.04	112.45	-113.36	112.25	112.12	114.50	108.88	111.97	111.42	100.92	100.02	95.06	41.57	45.73	45.05	82.77	81.58	80.53	42.08	-77.50	-72.30	49.16	44.35	59.11	57.78	54.15	41.53	40.93	40.14	33.01	-21.92	-40.37	36.03	34.92	-33.29	-22.84	-14.99	10.20		
110.05	112.87	-112.00	111.87	113.65	110.73	109.18	108.78	111.91	105.41	95.27	94.95	84.66	89.53	85.58	86.52	82.28	81.54	80.72	74.55	49.91	66.05	65.35	60.91	54.85	50.19	49.63	49.89	44.94	33.89	42.74	28.79	37.13	-35.88	-30.13	-17.74	-15.58	-9.31		
110.06	109.05	-109.30	111.68	111.54	-107.87	107.92	108.33	112.36	107.82	99.45	94.77	91.23	84.00	85.84	87.22	81.74	80.47	75.19	-71.31	45.25	44.66	53.82	58.71	53.12	52.31	-48.57	-47.73	-45.58	-25.31	-44.07	-42.57	-37.71	-36.78	-26.70	-41.69	-41.28	-29.14		
110.07	107.65	-107.77	108.56	106.75	107.11	-107.63	107.83	111.77	105.16	95.48	90.04	84.46	43.97	43.60	42.94	43.06	76.29	73.38	72.68	48.97	44.19	43.97	38.98	34.12	31.47	40.31	-49.16	-47.43	-46.66	-40.00	-46.77	42.53	-37.98	-42.70	-40.26	-48.75	-31.04		
110.08	102.68	-102.80	106.57	-106.33	106.75	-107.49	110.86	110.34	104.66	92.26	47.50	44.06	49.33	42.83	41.69	75.66	74.99	74.21	49.18	48.40	47.61	43.43	42.41	38.12	36.91	31.39	30.82	32.20	-47.61	-39.14	30.86	-47.61	-43.37	-47.72	-40.09	-30.99	-46.00		
110.09	102.46	-103.10	106.98	-106.53	-110.55	110.61	111.98	108.96	96.18	-80.09	44.38	43.82	41.52	39.61	76.57	75.87	75.56	-72.62	-68.93	48.31	37.86	42.77	41.95	57.08	56.93	34.11	37.10	32.95	33.79	51.75	49.41	-47.37	47.69	72.98	-74.37	-39.36			
110.10	102.41	-103.04	106.04	-106.50	-110.26	110.51	112.98	105.88	94.27	49.50	50.03	84.43	84.36	81.50	77.68	76.75	77.07	-74.07	-71.97	49.04	48.51	43.87	42.08	39.98	58.14	59.13	54.08	-14.03	34.11	33.51	-49.71	-48.46	-50.13	-33.87	39.06	-14.03			

Figure 131: Expert State Value Map