

# Lecture 17: Planning IV (Decision-Making II)

Professor Katie Driggs-Campbell

April 11, 2024

ECE484: Principles of Safe Autonomy



# Administrivia

- Upcoming due dates:
  - Final Presentations in class on 4/23 and 4/25
  - Final Video due 5/3
- Exam on 4/18 at 7pm
  - Email me ASAP about conflict exams
  - Make reservation in testing center for DRES accommodations
  - No cheatsheet will be needed
  - CA review session on Friday 4/12 (HW party time)
  - In-class review session on Tuesday 4/16
- Prof. DC OH by appointment next week (4/16)
  - Otherwise in 260 CSL

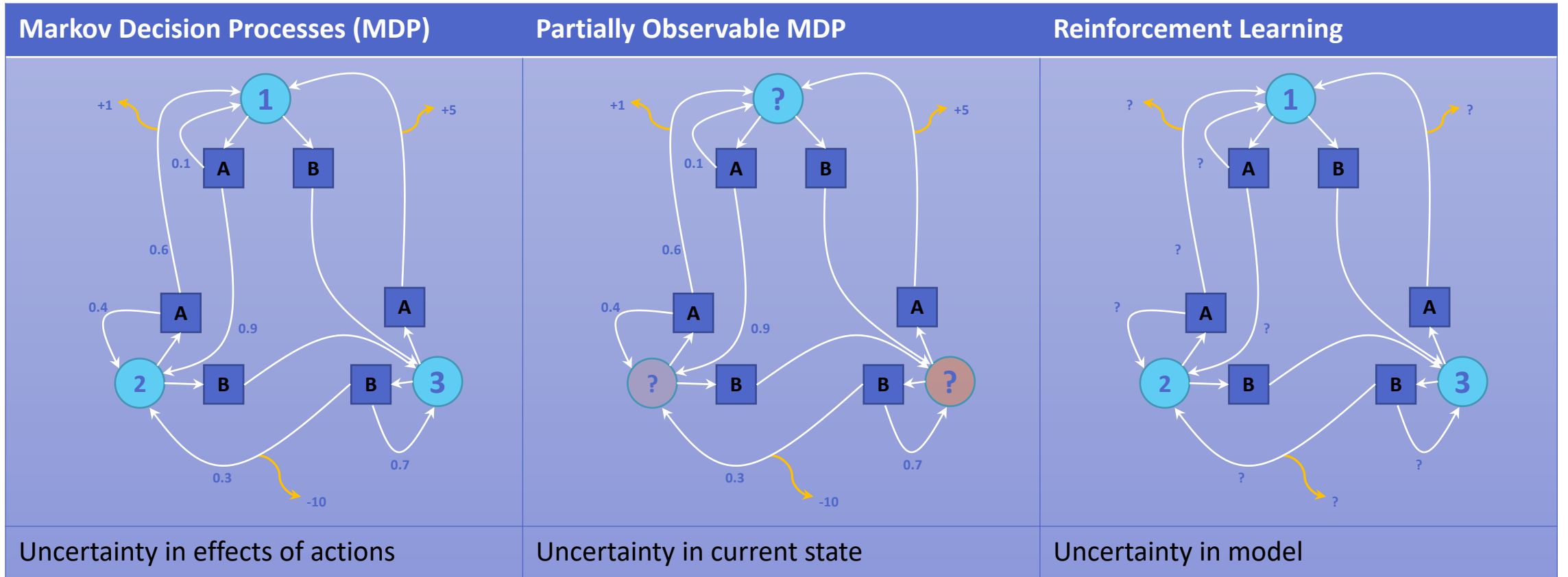


# Today's Plan

- Possible solutions for decision-making
- Markov Decision Processes
- MDP Policies and Value Iteration



# Markov Models



# Decision-Making Policies

- We want to devise a scheme that generates actions to optimize the future payoff *in expectation*



# Decision-Making Policies

- We want to devise a scheme that generates actions to optimize the future payoff *in expectation*
- Policy:  $\pi : x_t \rightarrow u_t$ 
  - Maps states to actions
  - Can be low-level reactive algorithm or a long-term, high-level planner
  - May or may not be deterministic



# Decision-Making Policies

- We want to devise a scheme that generates actions to optimize the future payoff *in expectation*
- Policy:  $\pi : \mathcal{X}_t \rightarrow \mathcal{U}_t$ 
  - Maps states to actions
  - Can be low-level reactive algorithm or a long-term, high-level planner
  - May or may not be deterministic
- Typically, we want a policy that optimizes future payoff, considering optimal actions over a planning (time) horizon



# MDP Policies

- Policies map states to actions

$$\pi: \mathcal{X} \rightarrow \mathcal{U}$$

- We want to find a policy that maximizes future pay off
  - Suppose  $T = 1$ :  $\pi_1(x) = \operatorname{argmax}_u r(x, u)$





# MDP Policies

- Policies map states to actions

$$\pi: x \rightarrow u$$

- We want to find a policy that maximizes future pay off
  - Suppose  $T = 1$ :  $\pi_1(x) = \operatorname{argmax}_u r(x, u)$

- We write the Value Function for given  $\pi$ :

$$V_1(x) = \gamma \max_u r(x, u)$$

- Generally, we want to find the sequence of actions that optimize the *expected cumulative discounted future payoff*



# Expected Cumulative Payoff

$$R_T = \mathbb{E} \left[ \sum_{\tau=0}^T \gamma^\tau r_{t+\tau} \right]$$



# Expected Cumulative Payoff

$$R_T = \mathbb{E} \left[ \sum_{\tau=0}^T \gamma^\tau r_{t+\tau} \right]$$

1. Greedy case:  $T = 1$   
→ Optimize next payoff



# Expected Cumulative Payoff

$$R_T = \mathbb{E} \left[ \sum_{\tau=0}^T \gamma^\tau r_{t+\tau} \right]$$

1. Greedy case:  $T = 1$   
→ Optimize next payoff
2. Finite Horizon:  $1 \leq T < \infty, (\gamma \leq 1)$   
→ Optimize  $R_T$  for set time window



# Expected Cumulative Payoff

$$R_T = \mathbb{E} \left[ \sum_{\tau=0}^T \gamma^\tau r_{t+\tau} \right]$$

1. Greedy case:  $T = 1$   
→ Optimize next payoff
2. Finite Horizon:  $1 \leq T < \infty$ , ( $\gamma \leq 1$ )  
→ Optimize  $R_T$  for set time window
3. Infinite Horizon:  $T = \infty$ , ( $\gamma < 1$ )  
→ Optimize  $R_\infty$  for all time

If  $|r| \leq r_{max}$ , discounting guarantees  $R_\infty$  is finite

$$R_\infty \leq r_{max} + \gamma r_{max} + \gamma^2 r_{max} + \dots = \frac{r_{max}}{1 - \gamma}$$



# Value Functions

For longer time horizons (T), we define  $V(x)$  recursively:

$$\text{Recall: } V_1(x) = \gamma \max_u r(x, u)$$



# Value Functions

- In the infinite time horizon, we tend to reach equilibrium:

$$V_{\infty}(x) = \gamma \max_u \left[ r(x, u) + \int V_{\infty}(x') p(x' | x, u) dx' \right]$$

- This is the *Bellman Equation*
  - Satisfying this is necessary and sufficient for an optimal policy



# Computing the (Approximate) Value Function

- Initial guess for  $\hat{V}$ 
  - $\hat{V}(x) \leftarrow r_{min}, \forall x$





# Computing the (Approximate) Value Function

- Initial guess for  $\hat{V}$ 
  - $\hat{V}(x) \leftarrow r_{min}, \forall x$
- Successively update for increasing horizons
  - $\hat{V}(x) \leftarrow \gamma \max_u [r(x, u) + \int \hat{V}(x') p(x' | x, u) dx']$
- Value iteration converges if  $\gamma < 1$



# Computing the (Approximate) Value Function

- Initial guess for  $\hat{V}$ 
  - $\hat{V}(x) \leftarrow r_{min}, \forall x$
- Successively update for increasing horizons
  - $\hat{V}(x) \leftarrow \gamma \max_u [r(x, u) + \int \hat{V}(x') p(x' | x, u) dx']$
- Value iteration converges if  $\gamma < 1$
- Given estimate  $\hat{V}(x)$ , policy is found:
  - $\pi(x) = \operatorname{argmax}_u [r(x, u) + \int \hat{V}(x') p(x' | x, u) dx']$



# Computing the (Approximate) Value Function

- Initial guess for  $\hat{V}$ 
  - $\hat{V}(x) \leftarrow r_{min}, \forall x$
- Successively update for increasing horizons
  - $\hat{V}(x) \leftarrow \gamma \max_u [r(x, u) + \int \hat{V}(x') p(x'|x, u) dx']$
- Value iteration converges if  $\gamma < 1$
- Given estimate  $\hat{V}(x)$ , policy is found:
  - $\pi(x) = \operatorname{argmax}_u [r(x, u) + \int \hat{V}(x') p(x'|x, u) dx']$
- Often, we use the discrete version:
  - $\pi(x) = \operatorname{argmax}_u [r(x, u) + \sum_{x'} \hat{V}(x') p(x'|x, u)]$



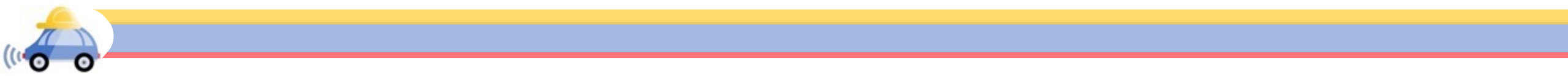
# A Simple MDP



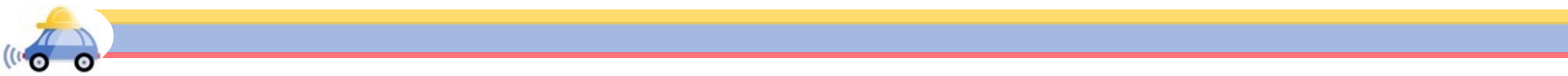
# Example: Value Iteration Setup



# Example: Value Iteration

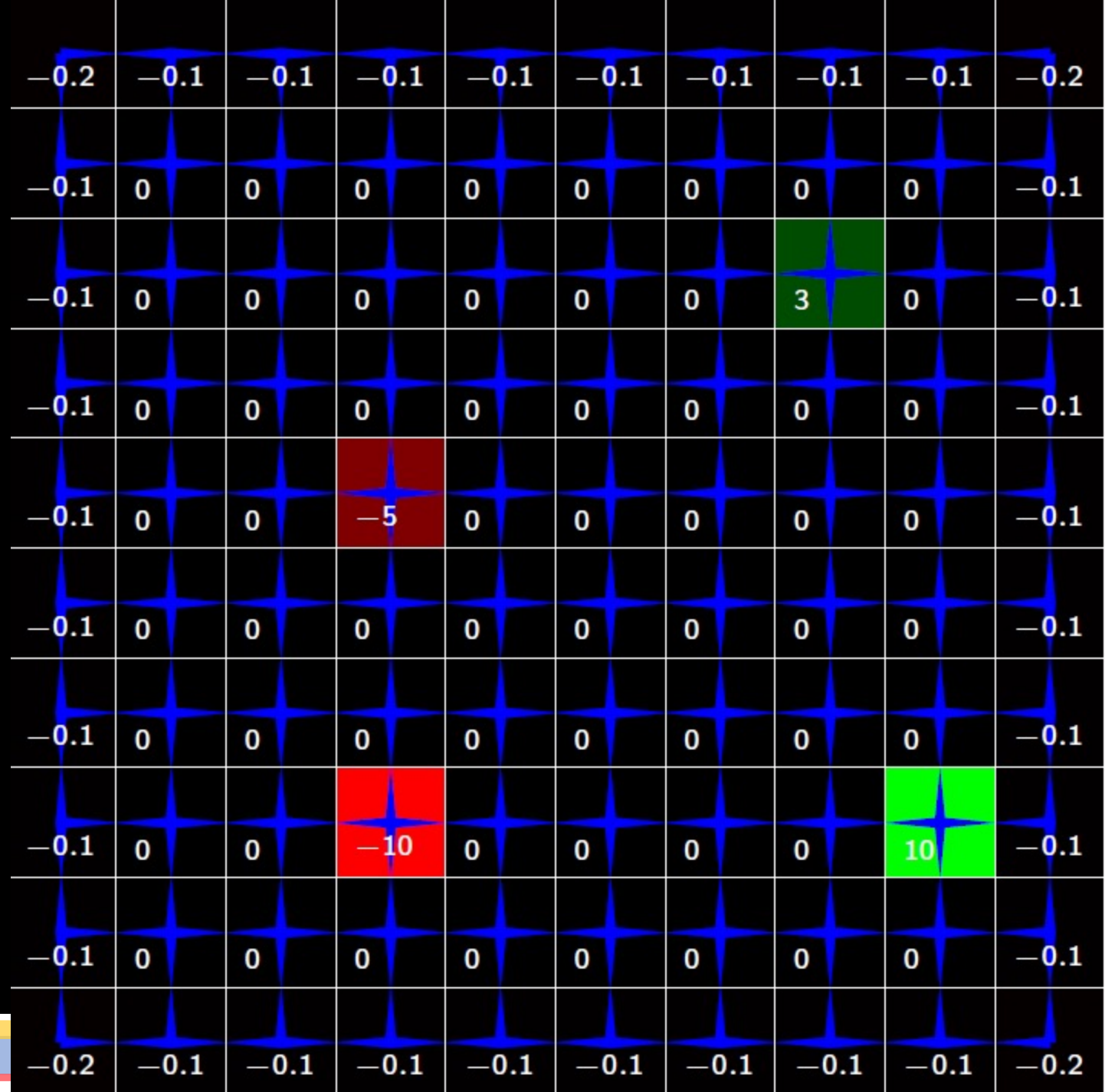


# Example: Rewards, Values, Policy



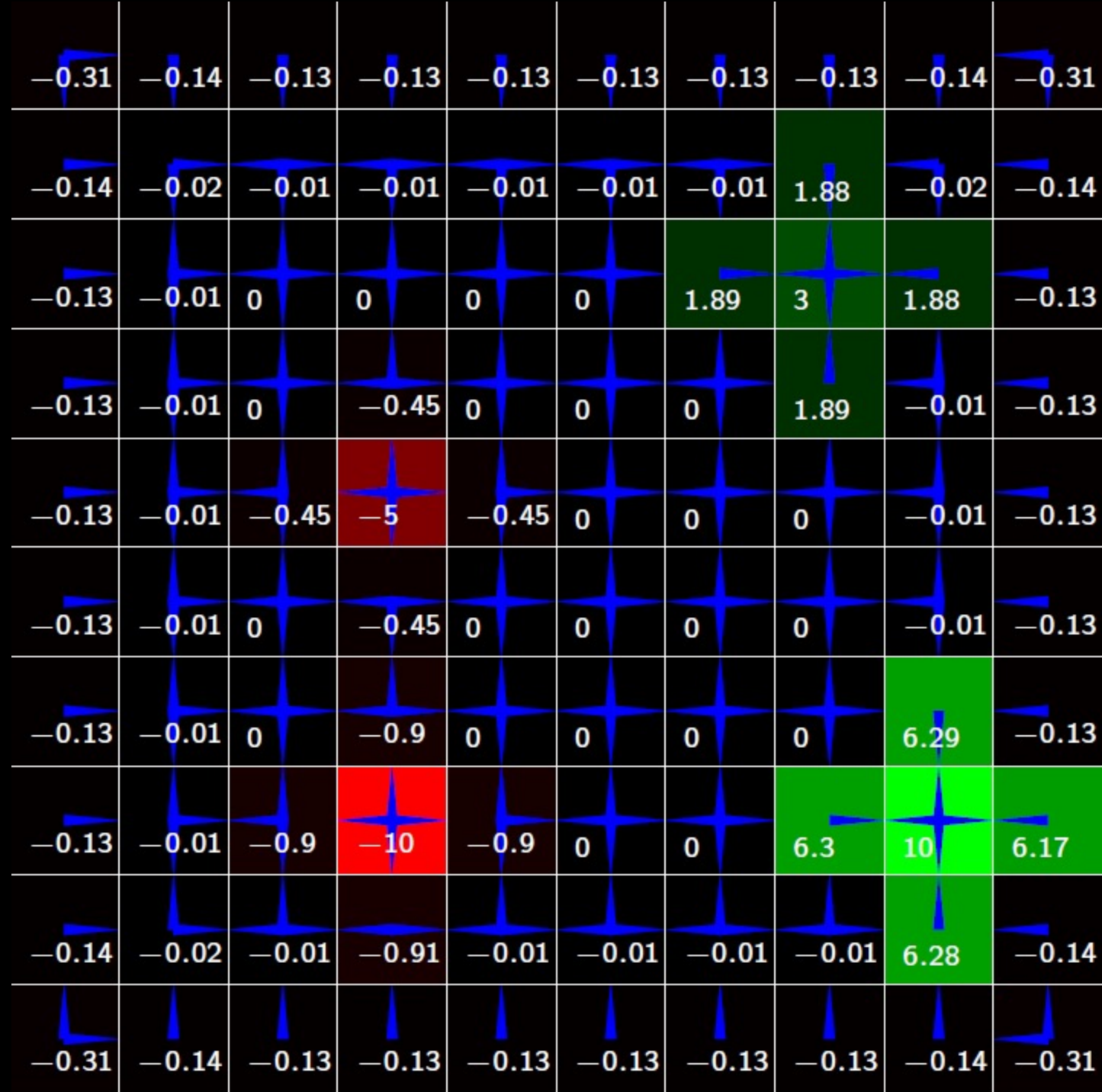
# Grid world example

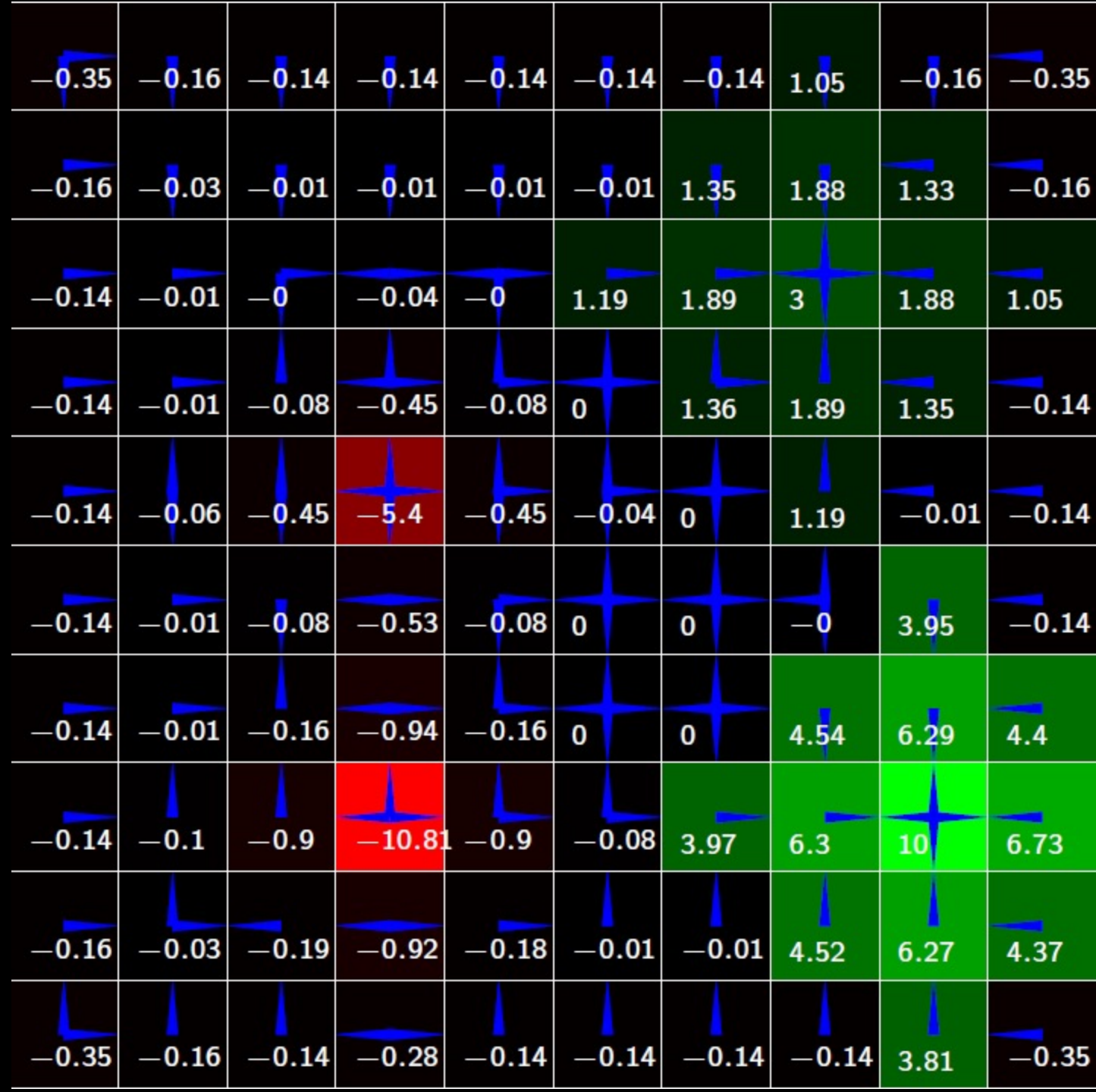
- States: cells in  $10 \times 10$  grid
- Actions: up, down, left, right
- Transition model: 0.7 chance of moving in intended direction, uniform in other directions
- Reward:
  - two states with cost
  - two terminal states with rewards
  - -1 for wall crash
- Discount is 0.9



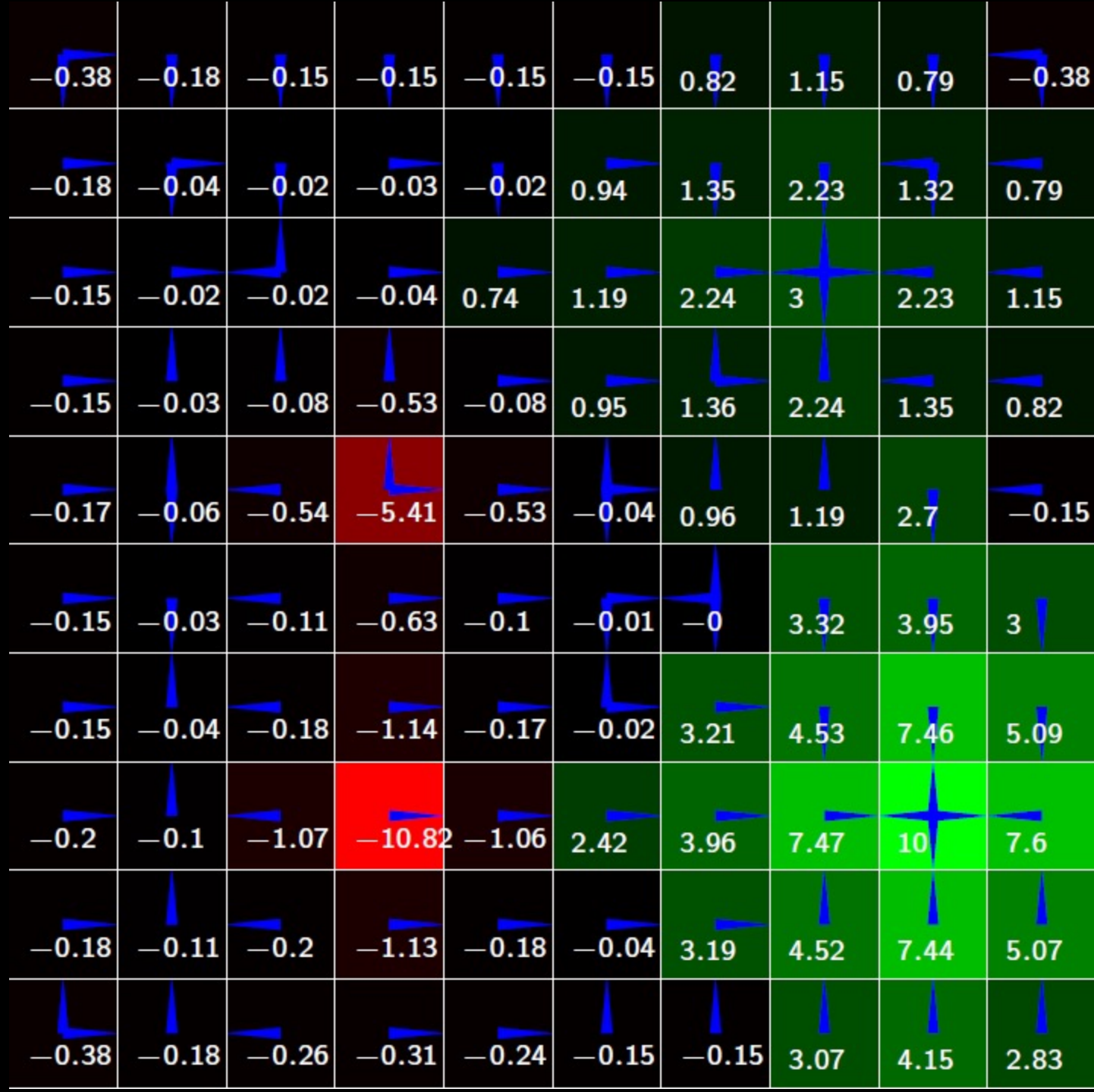


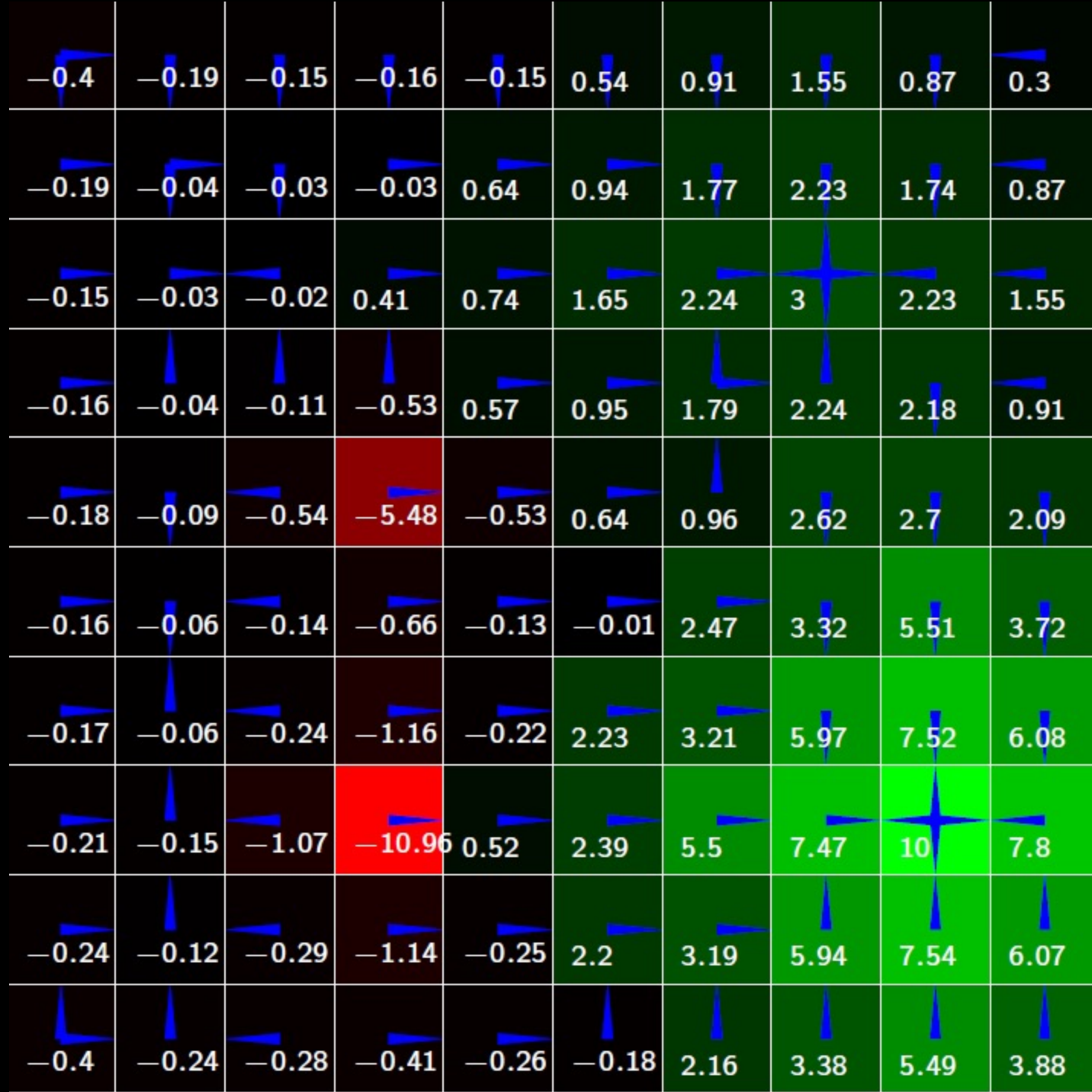


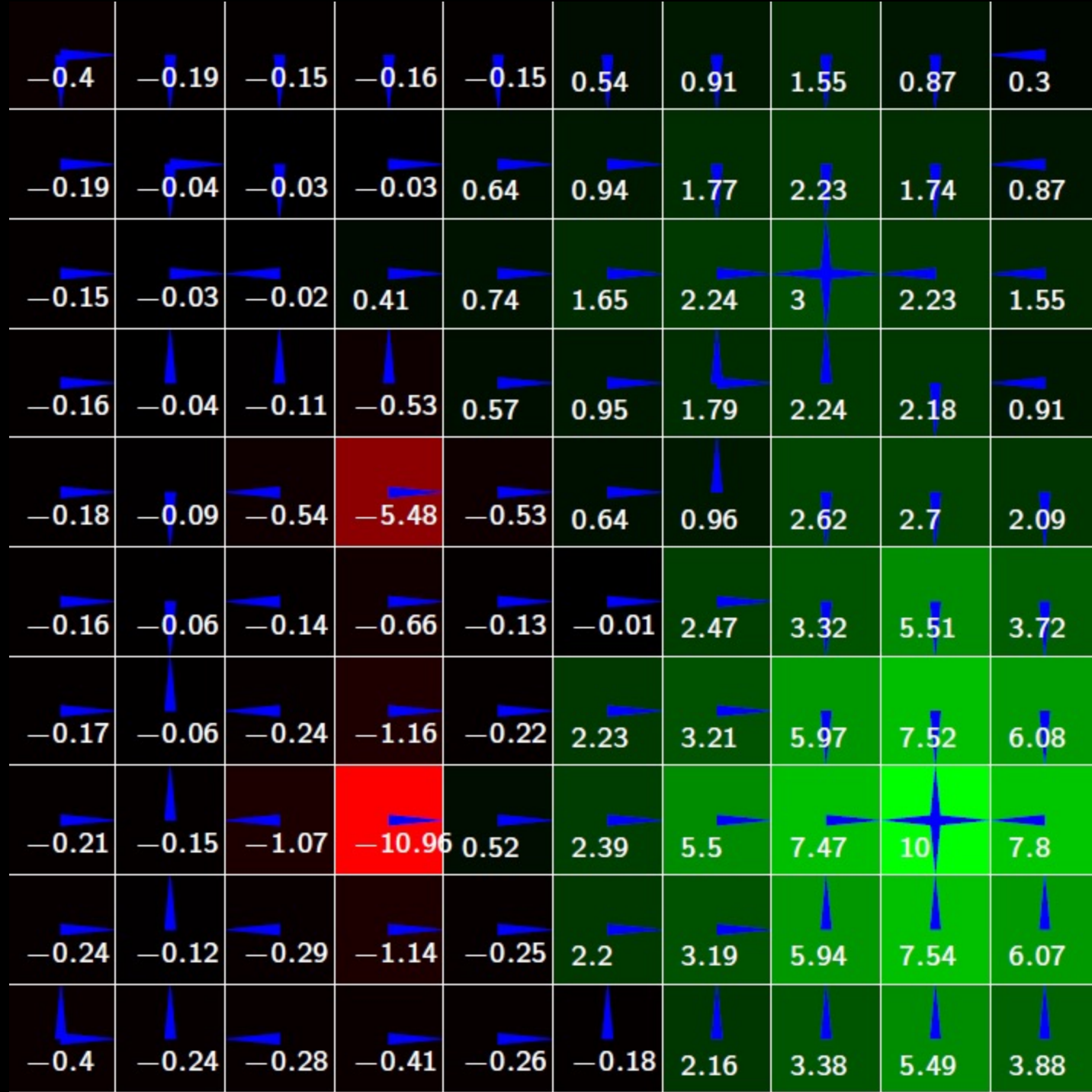




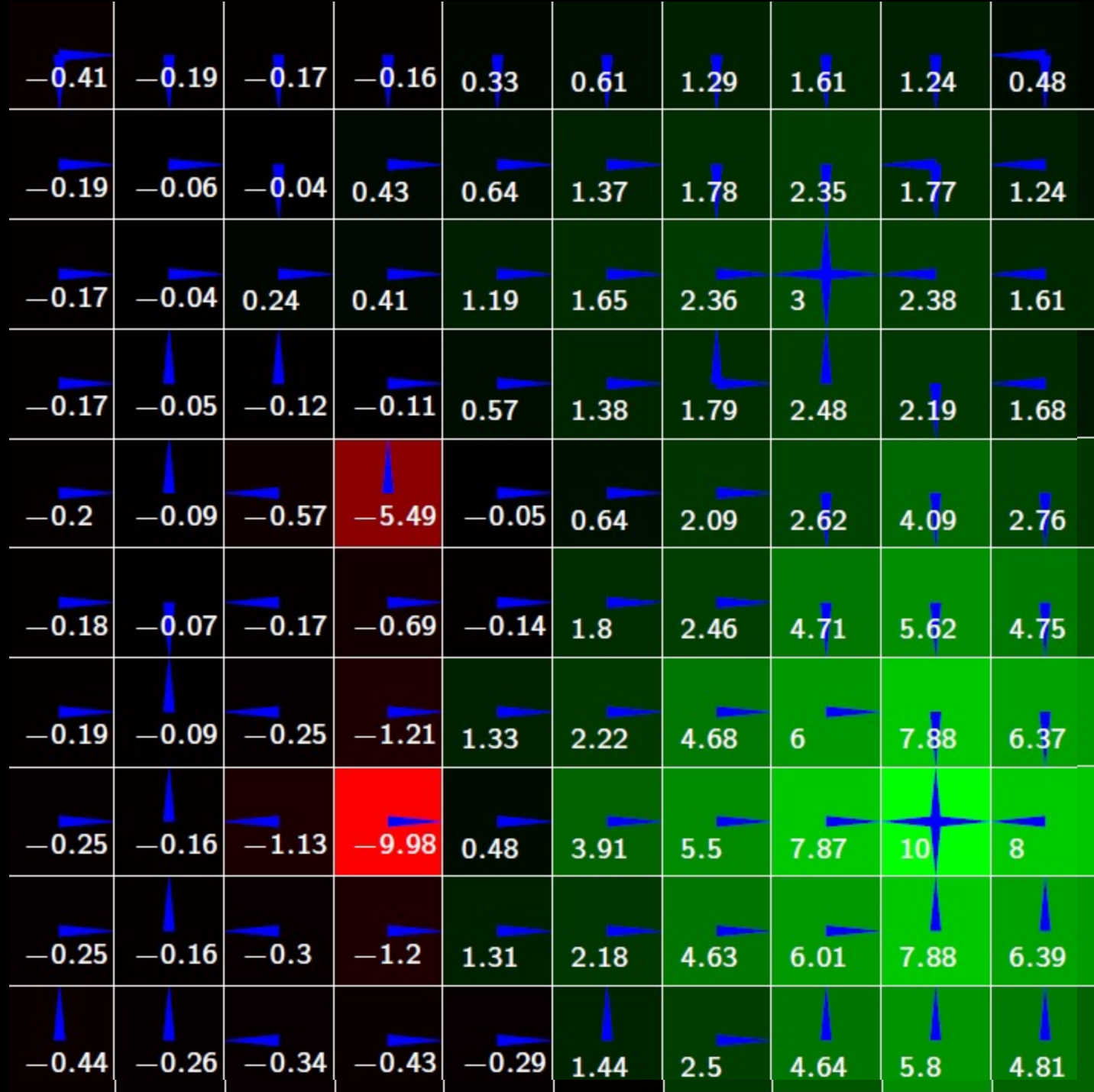


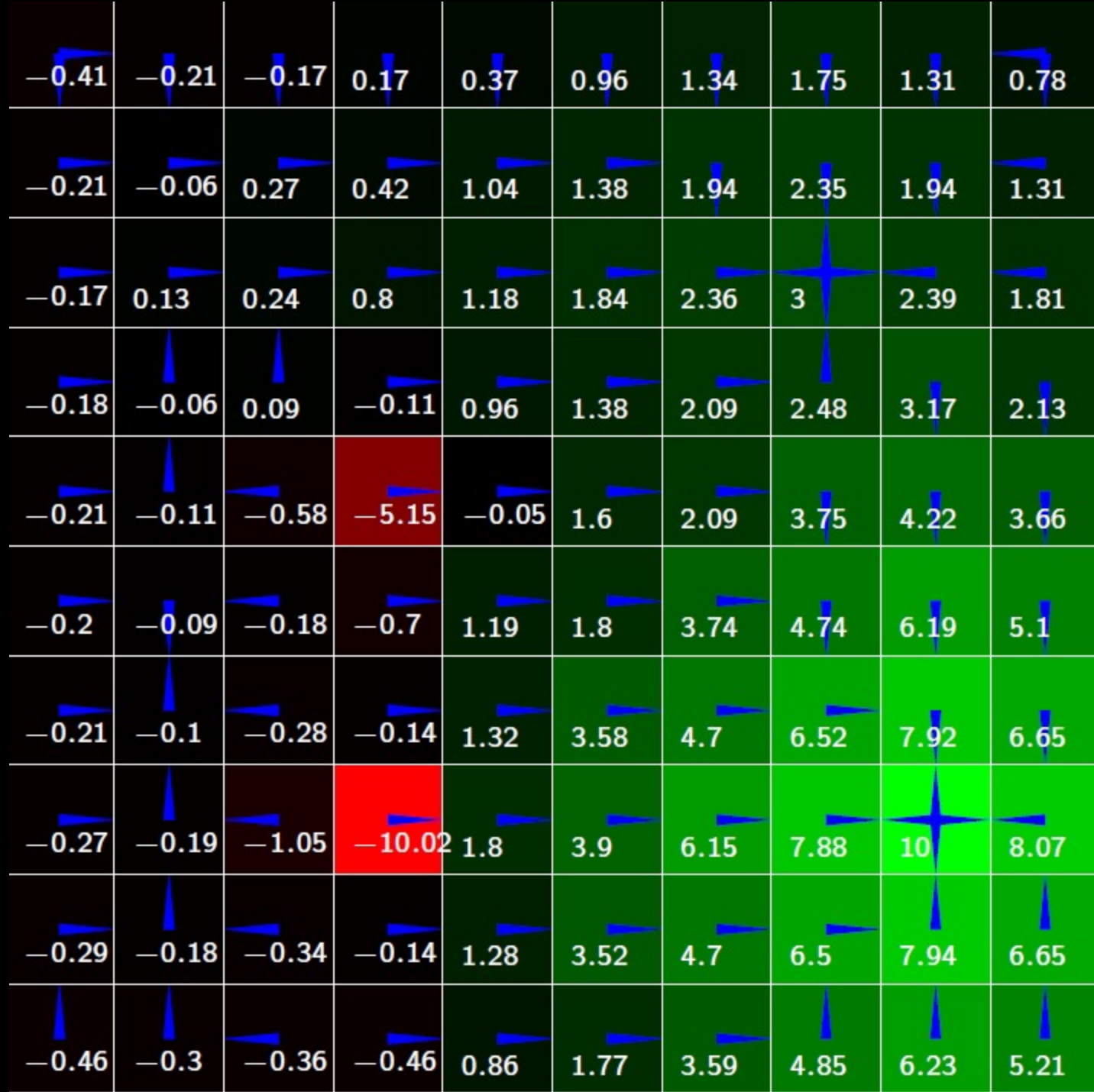




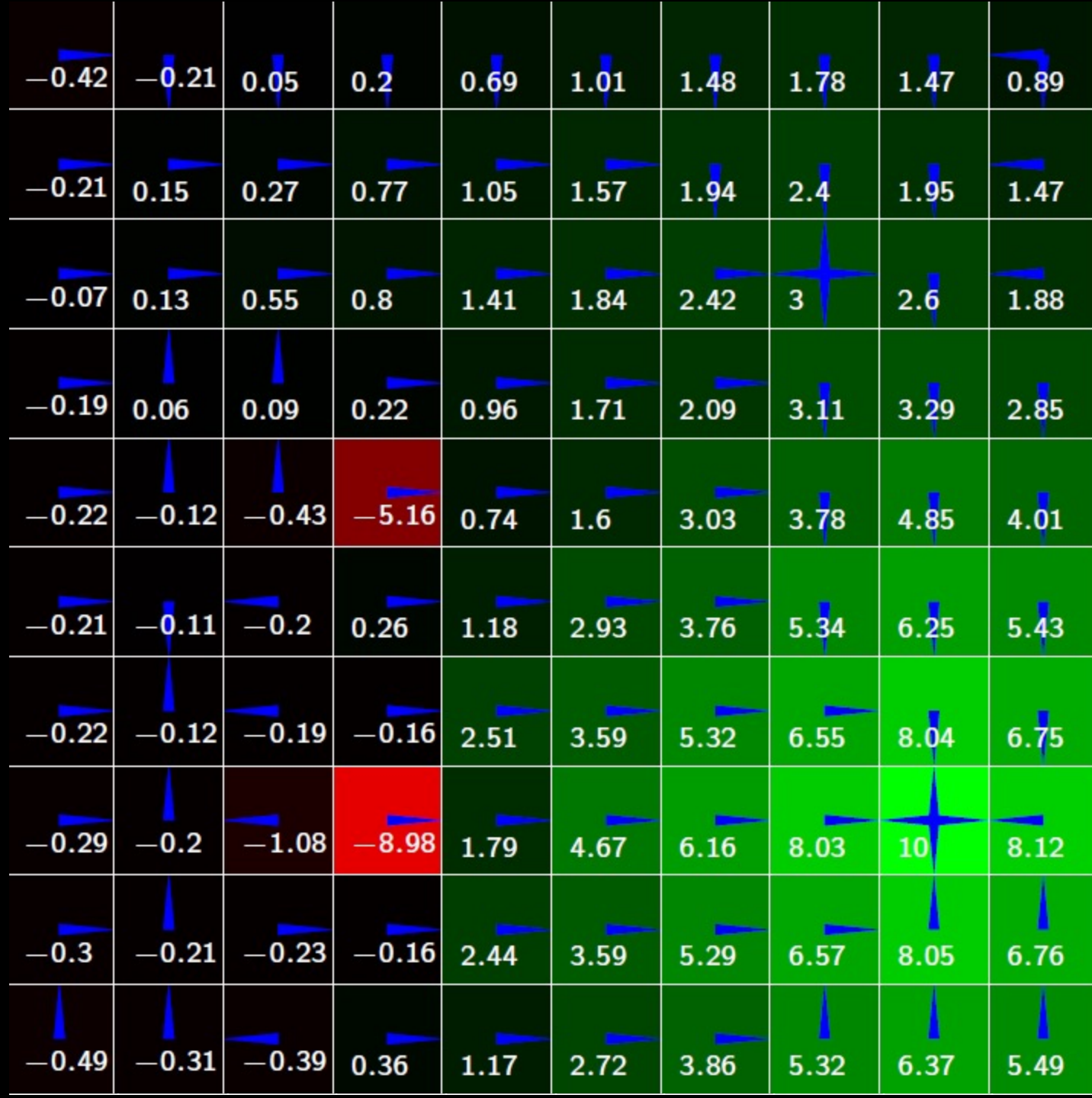


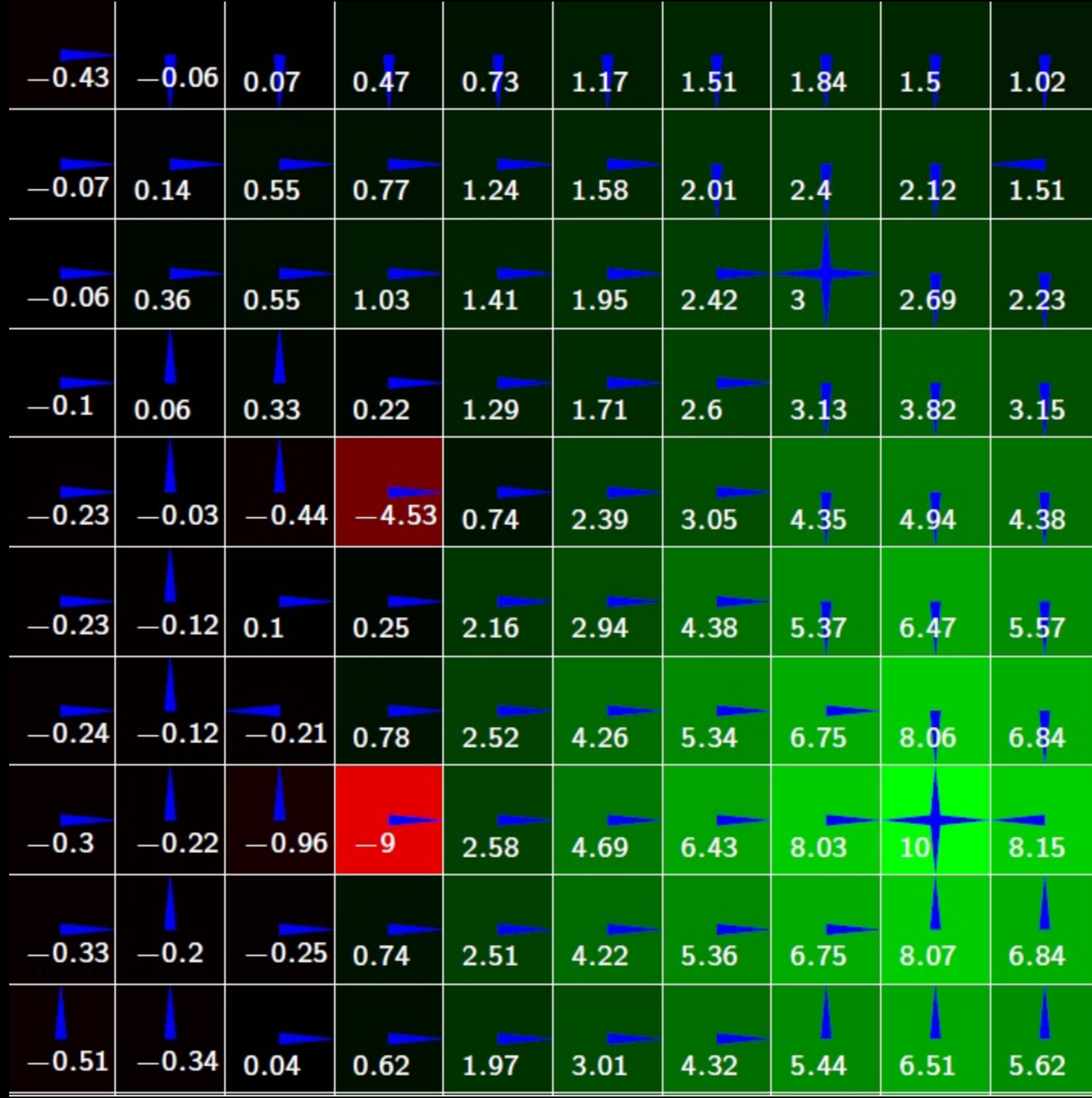








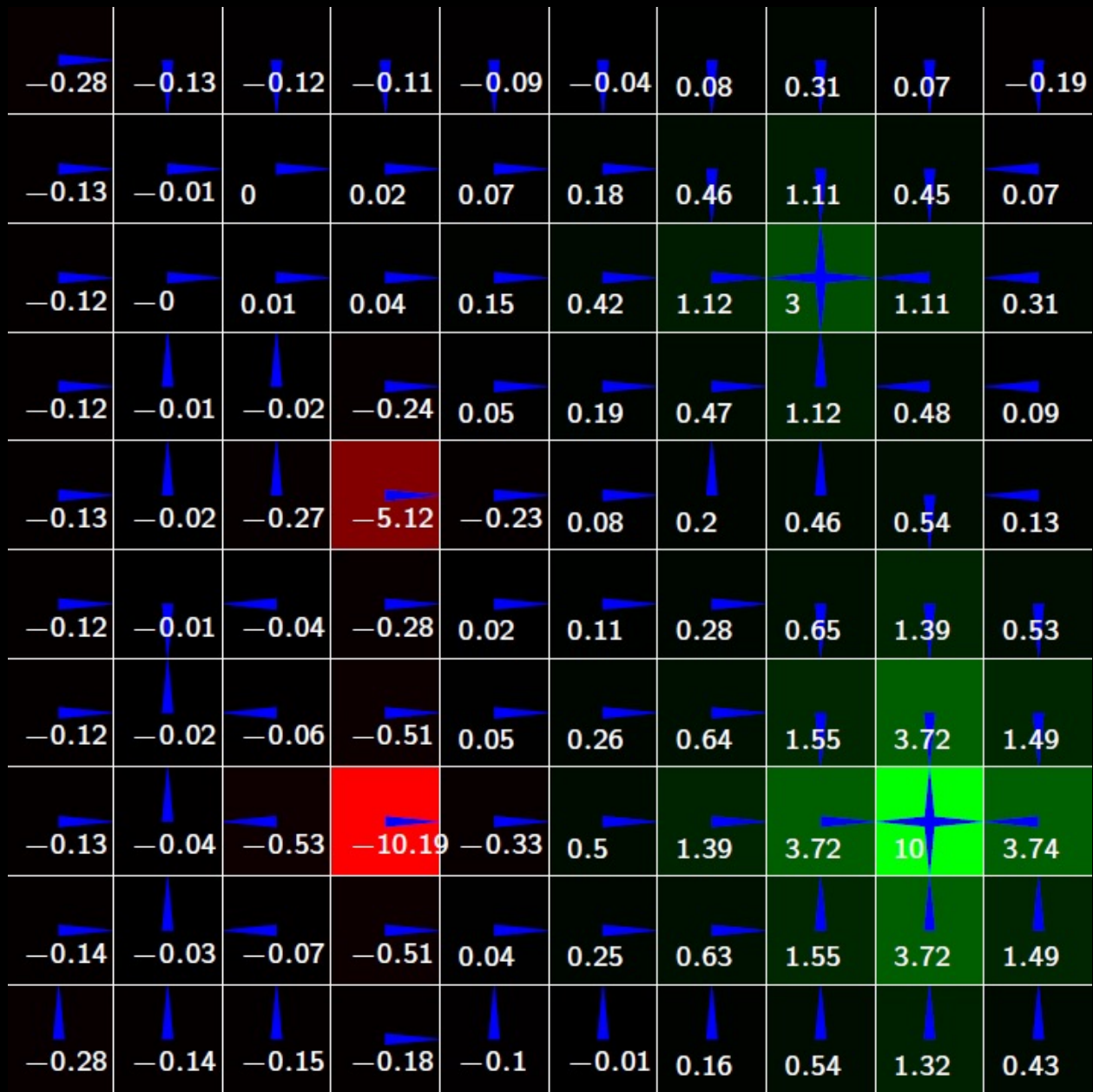
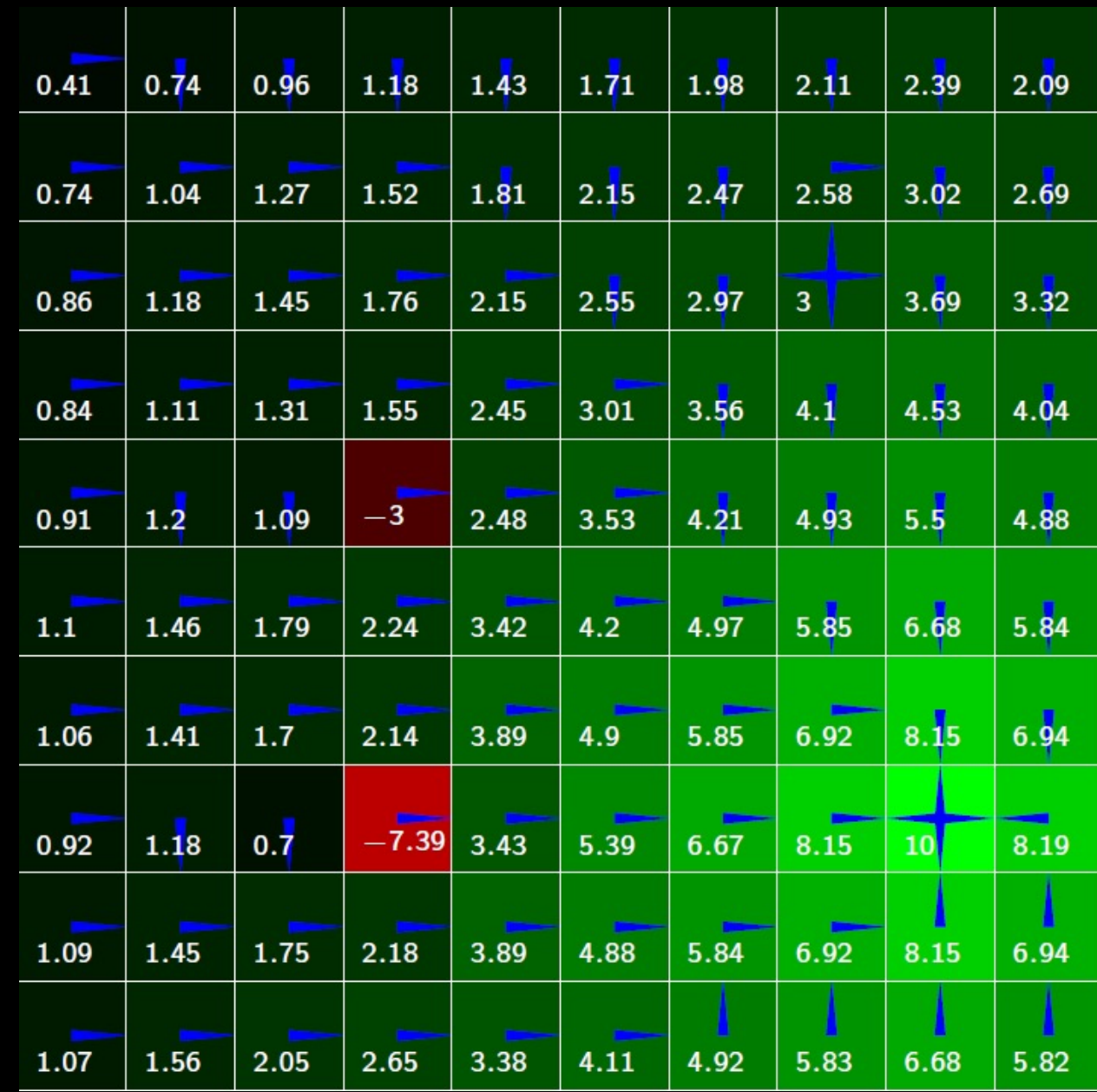




0.41	0.74	0.96	1.18	1.43	1.71	1.98	2.11	2.39	2.09
0.74	1.04	1.27	1.52	1.81	2.15	2.47	2.58	3.02	2.69
0.86	1.18	1.45	1.76	2.15	2.55	2.97	3	3.69	3.32
0.84	1.11	1.31	1.55	2.45	3.01	3.56	4.1	4.53	4.04
0.91	1.2	1.09	-3	2.48	3.53	4.21	4.93	5.5	4.88
1.1	1.46	1.79	2.24	3.42	4.2	4.97	5.85	6.68	5.84
1.06	1.41	1.7	2.14	3.89	4.9	5.85	6.92	8.15	6.94
0.92	1.18	0.7	-7.39	3.43	5.39	6.67	8.15	10	8.19
1.09	1.45	1.75	2.18	3.89	4.88	5.84	6.92	8.15	6.94
1.07	1.56	2.05	2.65	3.38	4.11	4.92	5.83	6.68	5.82

Converged!



$\gamma = 0.9$  $\gamma = 0.5$ 

# Decision-Making Summary

- Given an MDP, we defined **Expected Cumulative Payoff**, which plays a key role in **optimizing actions over planning horizons**
- Used **value iteration** to determine the “value” of a particular state, which helps us determine the best action to take considering future payoff
- We generally assumed the transition and reward function are known exactly – but what if we don’t have access to this information?



# Return to Safety

How would you incorporate desirable behaviors and safety requirements into your AV using the various decision-making frameworks?



# Responsibility Sensitive Safety

developed by Intel / MobilEye

Instead of looking at absolute safety, introduce a safety notion that depends on *responsibility*

→ AV should never be responsible for an accident

## RSS Rules:

1. Keep a safe longitudinal distance from the car ahead.
2. Keep a safe lateral distance from the cars on your sides.
3. Respect right-of-way rules (multiple geometries, traffic lights, pedestrians, unstructured roads).
4. Be cautious of occluded areas.



# RSS Example: Safe Following Distance

The longitudinal distance between a car ( $c_r$ ) that drives behind another car ( $c_f$ ) is safe w.r.t. a response time  $\rho$  if:

- $c_f$  applies at most  $a_{max}^{brake}$
  - $c_r$  will apply at most  $a_{max}^{accel}$  during response time
  - After  $\rho$ ,  $c_r$  will brake by at least  $a_{min}^{brake}$  until full stop
- $c_r$  will not collide with  $c_f$





# RSS Example: Safe Following Distance

The longitudinal distance between a car ( $c_r$ ) that drives behind another car ( $c_f$ ) is safe w.r.t. a response time  $\rho$  if:

- $c_f$  applies at most  $a_{max}^{brake}$
  - $c_r$  will apply at most  $a_{max}^{accel}$  during response time
  - After  $\rho$ ,  $c_r$  will brake by at least  $a_{min}^{brake}$  until full stop
- $c_r$  will not collide with  $c_f$

Remarks:

1. This is basic reachability!
2. The safe distance depends on a set of parameters that can be determined by regulation.
3. The parameters can be different for a robotic car and a human driver.
4. The parameters can be different for different road conditions.



# RSS Example: Safe Following Distance

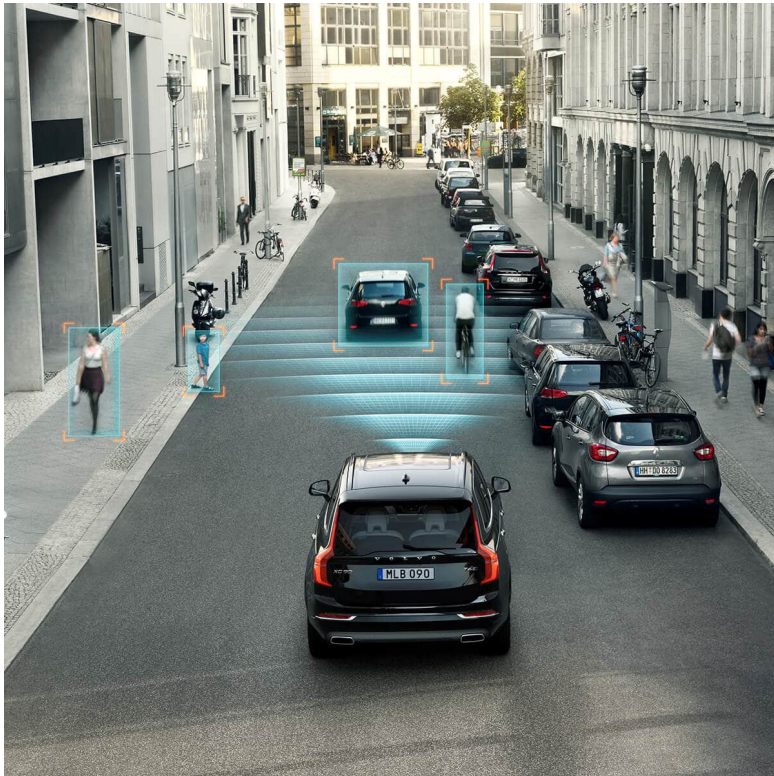
- Let  $v_r, v_f$  be the longitudinal velocities of the cars
- The minimal safe distance is:

$$d_{min} = \left[ v_r \cdot \rho + \frac{1}{2} a_{max}^{accel} \cdot \rho^2 + \frac{(v_r + \rho \cdot a_{max}^{accel})^2}{2a_{min}^{brake}} - \frac{v_f^2}{2a_{max}^{brake}} \right]_+$$



# Bonus MP4: Designing Decision-Logic

## Task 1: Emergency Braking



## Task 2: Stopping at a Red Light



# Bonus MP4: Designing Decision-Logic

## Validation

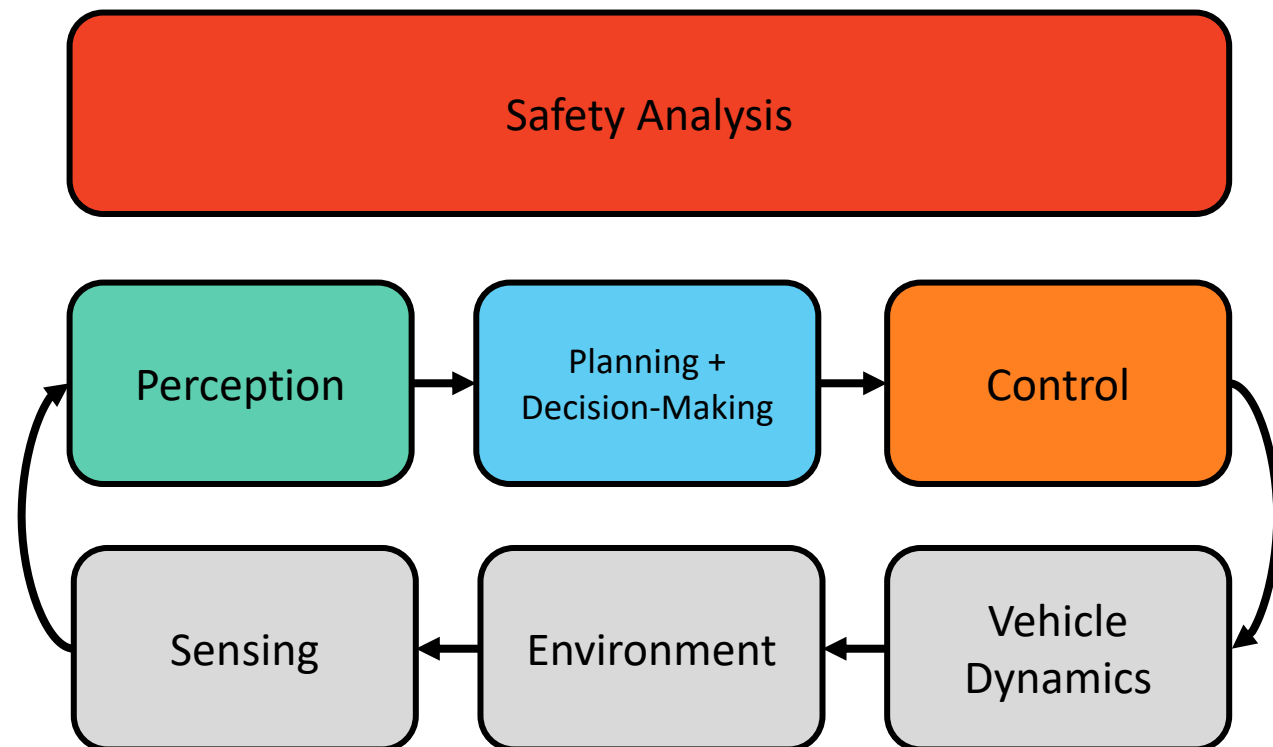
- Define test cases
- Gather samples through simulation
- Collect statistical evidence of safety

## Verification

- Define uncertainty
- Compute reachable sets to assess safety and provide guarantees



# Course Recap



# Why are we here?

Know



Components of an autonomous system and safety standards.

→ How to use software modules for perception, planning, control, ROS, OpenCV, ...

Do



Code and analyze algorithms for perception, localization, planning, control, & verification

→ Plan, propose, organize and execute a team project

Understand



Models, algorithms, data, biases, assumptions for building trustworthy autonomous systems

→ Theoretical properties of algorithms and their limitations

Get Inspired



Become the Isaac Newton of Autonomy

→ “To do things right, first you need love, then technique.” – Antoni Gaudí

