

Lecture 15: Decision-Making

Professor Katie Driggs-Campbell

March 30, 2021

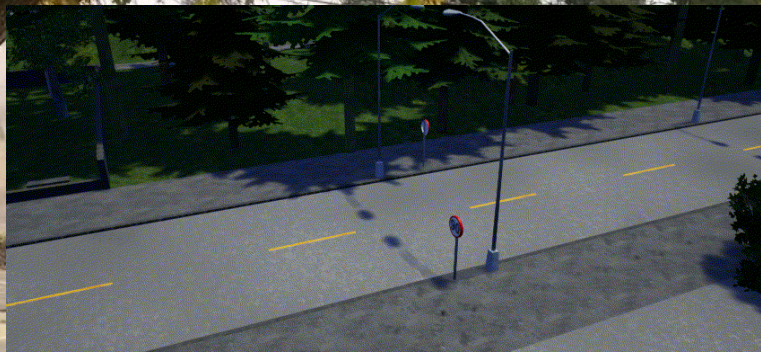
ECE484: Principles of Safe Autonomy



Administrivia

- Get started on MP3 and your project!
- Will go over oral exam protocol on Thursday
- Upcoming guest lecture attendance is “worth” double
- Participation grades (10% of total grade):
 - $P \approx \frac{1}{3}\{\text{attendance}\} + \frac{1}{3}\{\text{guest lecture participation}\} + \frac{1}{3}\{\text{team assessment}\}$
 - Stop by OH or make appointment to check attendance grade
 - For guest lecture participation, you can either send in questions beforehand (via Google forms to be posted on discord) or ask during class
 - For team assessment, we will post a Google form for collecting feedback on your teammates



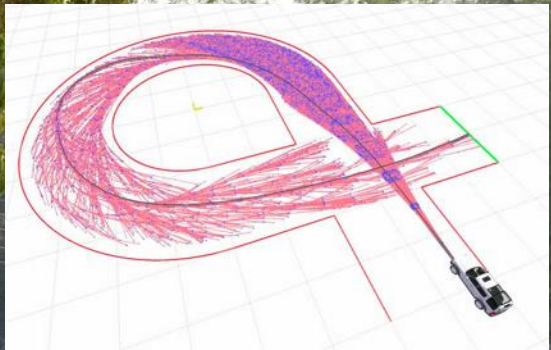


Simulation & Validation ✓

Sensors ✓



Perception ✓



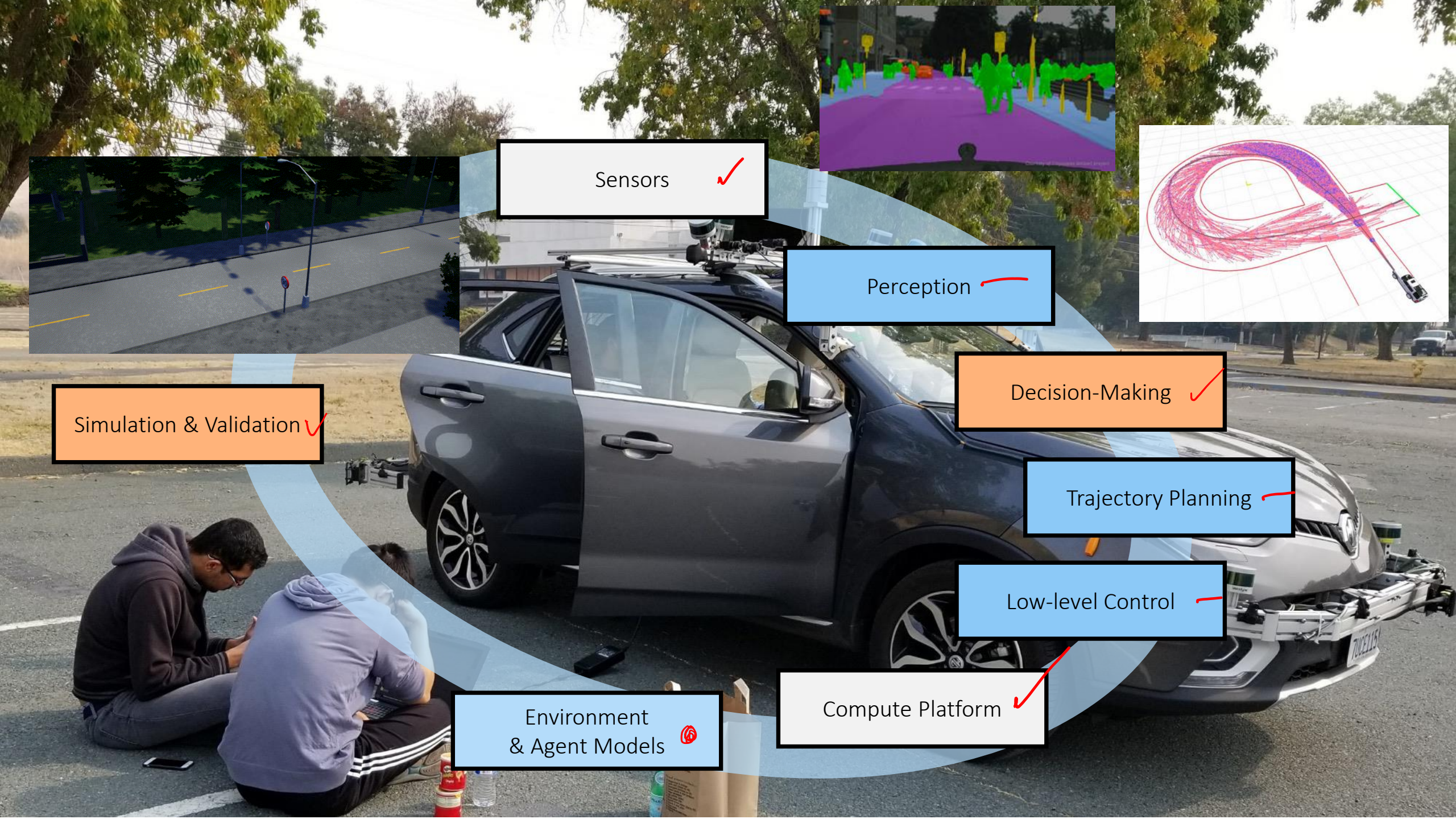
Decision-Making ✓

Trajectory Planning ✓

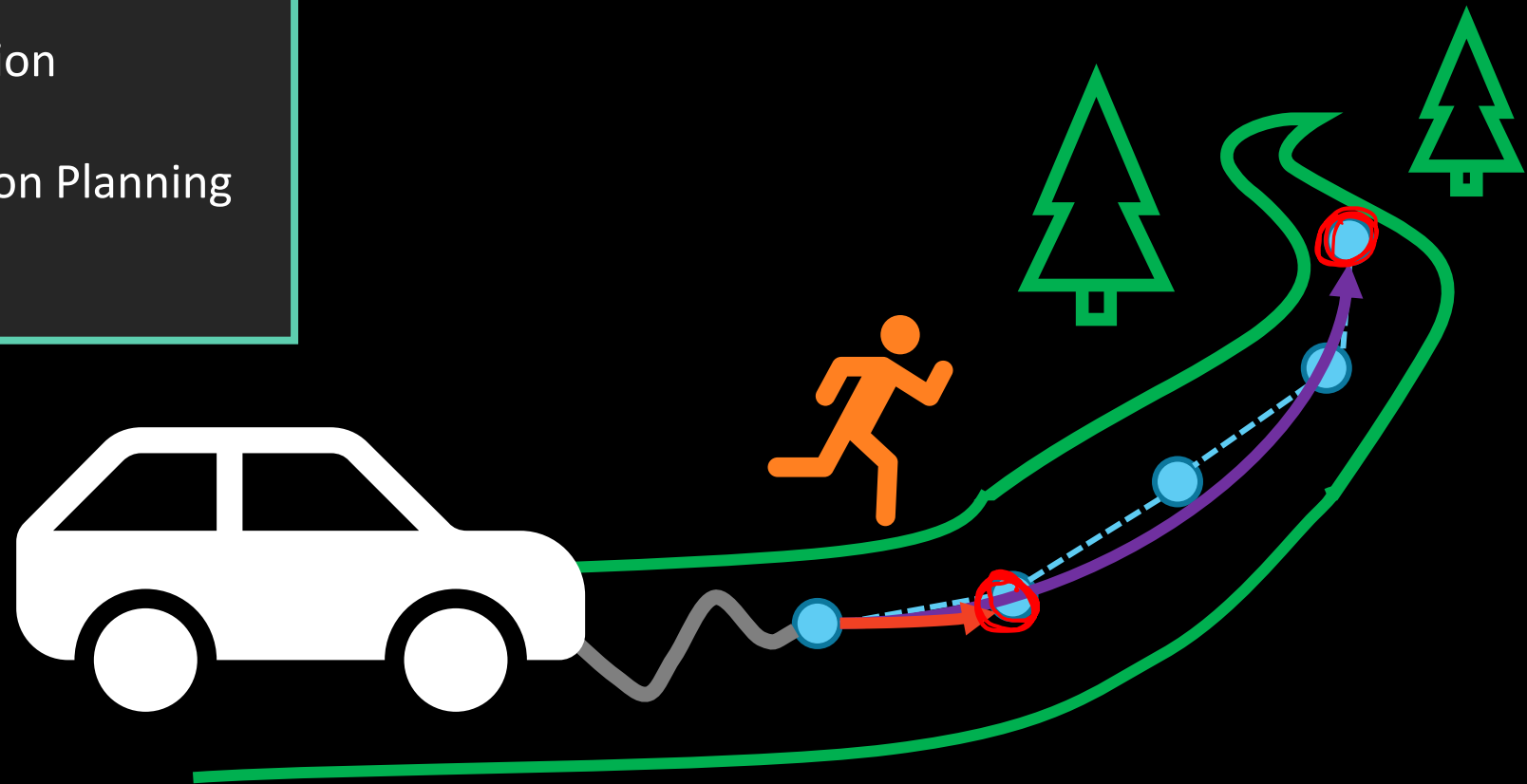
Low-level Control ✓

Environment & Agent Models Ⓜ

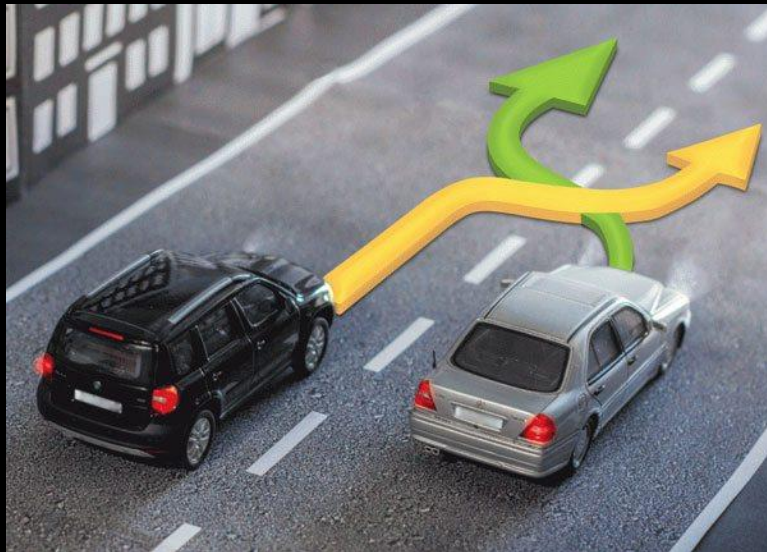
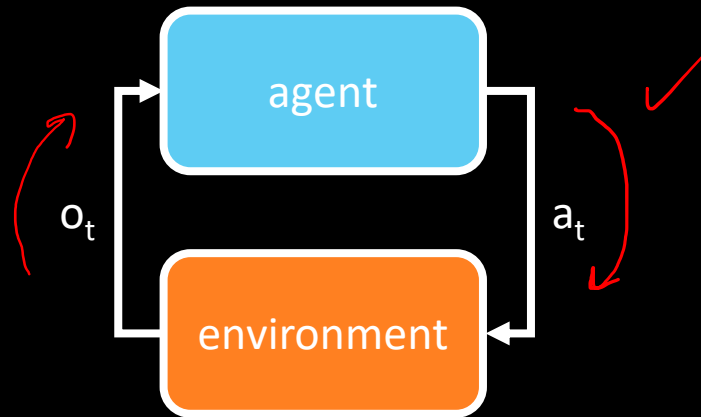
Compute Platform ✓



- Vehicle Modeling
- Localization
- Detection & Recognition
- Control
- Trajectory/Path/Motion Planning
- **Decision Making**
- Final topic: Safety!



High-Level Decision-Making



Coordinated drones
Posted by Tech Insider
8,033,416 Views

They used coding and algorithms so the drones didn't crash into each other

16k 485 Share

BEST

u/Skizm · 2mo

```
if(goingToCrashIntoEachOther)
{ dont(); }
```

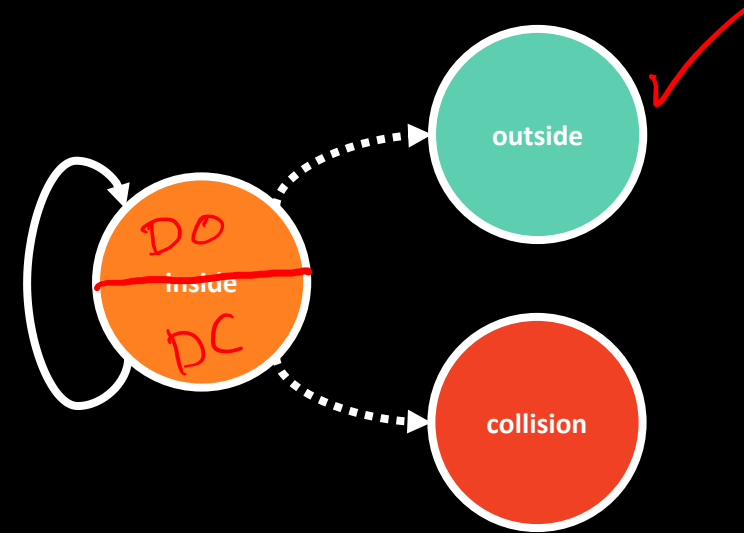
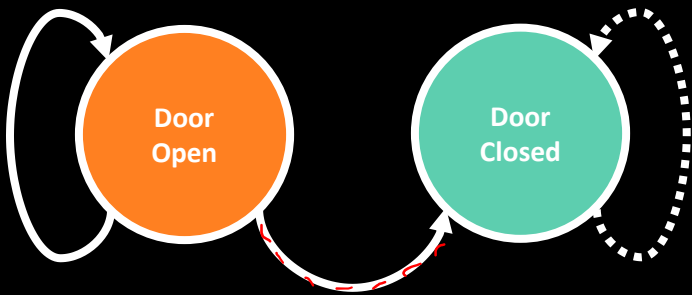
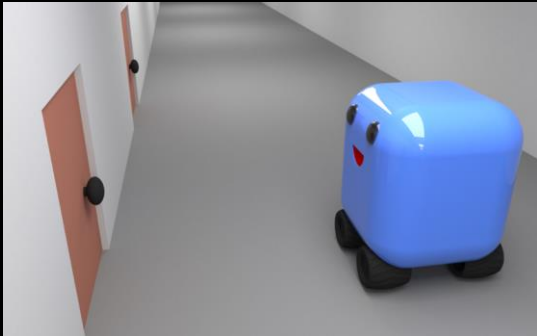
skiz never

as a robotics major i can confirm this is 100% how coding works



From Filtering to Decision-Making

Recall: Filtering allows us to recursively update our belief about some state



Decision-making helps us reason about what actions we should take



Today's Plan

- Introduction to decision-making
- Markov Decision Processes
- MDP Policies and Value Iteration
- Simple Example
 - Will post worked out complicated example

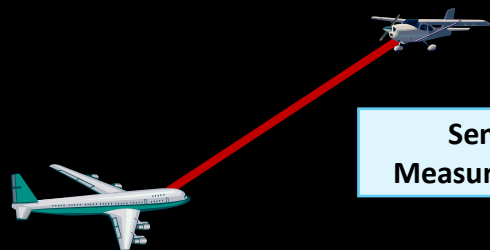


Today's Plan

- Introduction to decision-making
- Markov Decision Processes
- MDP Policies and Value Iteration
- Simple Example



Traffic Alert and Collision Avoidance System (TCAS)



Sensor
Measurements

```
IF (ITF.A LT G.ZTHR)
THEN IF (ABS(ITF.VMD) LT G.ZTHR)
THEN SET ZHIT;
ELSE CLEAR ZHIT;
ELSE IF (ITF.ADOT GE P.ZDTHR)
THEN CLEAR ZHIT
ELSE
ITF.TAUV = -ITF.A/ITF.ADOT;
IF (ITF.TAUV LT TVTHR AND
((ABS(ITF.VMD) LT G.ZTHR) OR
(ITF.TAUV LT ITF.TRTRU))
THEN SET ZHIT
ELSE CLEAR ZHIT
IF (ZHIT EQ $TRUE AND
ABS(ITF.ZDINT) GT P.MAXZDINT
THEN CLEAR ZHIT
```

Resolution
Advisory



Surveillance

Advisory Logic

Display



PROCESS Reversal_modeling;

```

Default modeled separation for current RA is 0 if current RA is negative;
Set own altitude and own rate to own tracked altitude and own tracked rate;

IF (own does not follow his RAs)
  THEN Model separation achieved assuming RA not followed;
  IF (current RA is a climb RA)
    THEN CLEAR flag indicating the sense of the RA after a reversal;
    ELSE SET flag indicating the sense of the RA after a reversal;
  IF (modeled separation achieved by continuing current RA greater than 1.2 *
      P.CROSSTHR)
    THEN CLEAR reversal flag in ITF
  ELSE
    <Begin own is assumed to follow its RA>
    IF (current RA is positive)
      THEN model response to current RA;
      <model maximum displayable rate for climb if current rate exceeds
      maximum displayable rate or minimum displayable rate for descent if
      current rate is less than minimum displayable rate>

      IF (tracked response lags modeled response in RA direction AND
          time since RA less than a parameter time AND
          own's rate has not changed by more than P.MODEL_ZD since the
          RA was first issued)
        THEN set own altitude and own rate to modeled altitude and rate
            for use in reversal modeling;
        Model separation achieved by continuing current RA;
        Set delay time to greater of pilot delay time remaining for last advisory against a
        new threat, and the pilot quick reaction time;

      IF (considering a reversal from a descend RA to a climb RA)
        THEN set own goal rate to greater of own tracked rate (or maximum
        displayable rate, whichever is less) and nominal climb rate;
        ELSE IF (own too close to ground to descend)
          THEN set own goal rate to zero;
          ELSE set own goal rate to lesser of own tracked rate (or minimum
          displayable rate, whichever is greater) and nominal descent
          rate;

      IF (vertical chase, low VMD geometry was not the reason for considering
          reversal)
        THEN IF (intruder causing crossing OR (intruder level AND own crossing
            from above) OR intruder rate and own modeled rate are opposite in sign)
          THEN use outer rate bound to model intruder;
          ELSE use inner rate bound to model intruder;
        ELSE use intruder's tracked vertical rate to model intruder;
      CALL MODEL_SEP
      IN (delay, goal rate, own altitude, own rate, acceleration response, sense after
          reversal, intruder altitude, modeled intruder rate, ITF entry)
      OUT (predicted separation for sense reversal);

      IF (Predicted separation for sense reversal is not positive OR
          modeled separation achieved by continuing current RA GE G.ALIM)
        THEN CLEAR reversal flag in ITF;
    <End own is assumed to follow its RA>

```

END Reversal_modeling;

RESOLUTION HIGH-LEVEL LOGIC

6-P22

PROCESS Reversal_modeling;

```

NOMINAL_SEP = 0;
Z = G.ZDOWN;
ZD = G.ZDOWN;
DELAY = 0;

IF (G.OWN_FOLLOW EQ FALSE)
  THEN CALL MODEL_SEP
  IN (DELAY, ZD, Z, ZD, P.VACCEL, OWNTENT(7), ITF.ZINT, ITF.ZDINT, ITF.entry)
  OUT (NOMINAL_SEP);
  IF (OWNTENT(7) EQ $TRUE)
    THEN NEW_SENSE = $FALSE;
    ELSE NEW_SENSE = $TRUE;
  IF (NOMINAL_SEP GT 1.2 * P.CROSSTHR)
    THEN CLEAR ITF.REVERSE;
  ELSE
    <Begin own is assumed to follow its RA>
    IF (OWNTENT(5,6) EQ '00')
      THEN DELAY = MAX(P.TV1 - (G.TCUR - G.TPOSRA), 0);
      IF (OWNTENT(7) EQ $FALSE)
        THEN ZDGOAL = MAX(MIN(G.ZDOWN, P.MAXDRATE), P.CLMRT);
        ELSE ZDGOAL = MIN(MAX(G.ZDOWN, P.MINDRATE), P.DESRT);
      CALL PROJECT_VERTICAL_GIVEN_ZDGOAL
      IN ((G.TCUR - G.TPOSRA), G.ZTV, G.ZDTV, ZDGOAL, P.TV1, P.VACCEL)
      OUT (ZPROJ, ZDPROJ);
      IF (((OWNTENT(7) EQ $FALSE AND ZPROJ GT G.ZDOWN AND
          (G.ZDOWN GE G.ZDTV - P.MODEL_ZD)) OR
          (OWNTENT(7) EQ $TRUE AND ZPROJ LT G.ZDOWN AND
          (G.ZDOWN LE G.ZDTV + P.MODEL_ZD))) AND
          G.TCUR - G.TPOSRA LT P.MODEL_T)
        THEN Z = ZPROJ;
        ZD = ZDPROJ;
      CALL MODEL_SEP
      IN (DELAY, ZDGOAL, Z, ZD, P.VACCEL, OWNTENT(7),
          ITF.ZINT, ITF.ZDINT, ITF.entry)
      OUT (NOMINAL_SEP);
      IF (OWNTENT(7) EQ $TRUE)
        THEN NEW_SENSE = $FALSE;
        ELSE NEW_SENSE = $TRUE;
      DELAY = MAX(P.TV1 - (G.TCUR - G.TLASTNEWRA), P.QUICKREAC);

      IF (NEW_SENSE EQ $FALSE)
        THEN ZDGOAL = MAX(P.CLMRT, MIN(G.ZDOWN, P.MAXDRATE));
        ELSE IF (G.NODESCENT EQ $TRUE)
          THEN ZDGOAL = 0;
          ELSE ZDGOAL = MIN(P.DESRT, MAX(G.ZDOWN,
              P.MINDRATE));
      IF (G.REV_CONSDRD EQ FALSE)
        THEN IF ((ITF.INT_CROSS EQ $TRUE) OR (ITF.ZDINT EQ 0 AND
            ITF.RZ GT 0) OR (ITF.ZDINT * G.ZDMODEL LT 0))
          THEN MZDINT = ITF.ZDOUTR;
          ELSE MZDINT = ITF.ZDINR;
        ELSE MZDINT = ITF.ZDINT;

      CALL MODEL_SEP
      IN (DELAY, ZDGOAL, Z, ZD, P.RACCEL, NEW_SENSE, ITF.ZINT, MZDINT, ITF.entry)
      OUT (ZMP);

      IF (ZMP LE 0 OR NOMINAL_SEP GE G.ALIM)
        THEN CLEAR ITF.REVERSE;
    <End own is assumed to follow its RA>
END Reversal_modeling;

```

RESOLUTION LOW-LEVEL LOGIC

6-P23



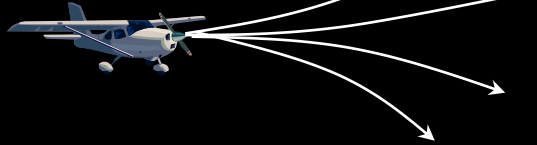
Why is it hard?

State Uncertainty



Imperfect sensor information leads to uncertainty in position and velocity of aircraft

Dynamic Uncertainty



Variability in pilot behavior makes it difficult to predict future trajectories of aircraft

Multiple Objectives



System must carefully balance both safety and operational considerations



Decision-Making Methods

1. Explicit programming
 - Ex: if/then statements
 - Heavy burden on designer



Heuristic Method for Lane Changing: MOBIL

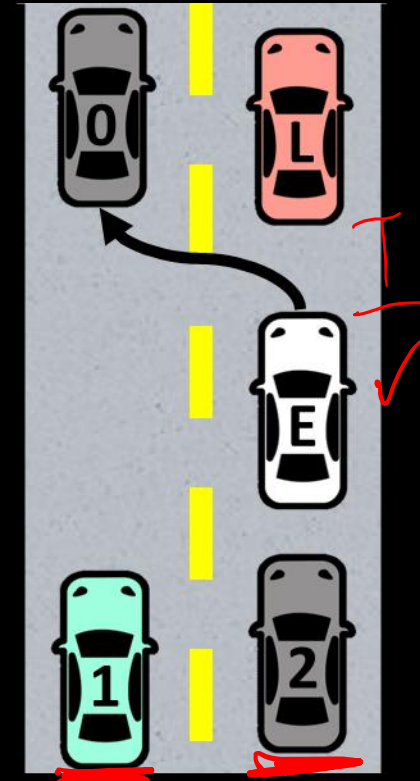
- Safety criterion:

$$\tilde{a}_E \geq -b_{safe}$$

- Decision rule:

$$\tilde{a}_E - a_E + p(\tilde{a}_1 - a_1 + \tilde{a}_2 - a_2) > \Delta a_{th}$$

- Politeness factor, p : 0.35
- Safe braking limit, b_{safe} : 2 m/s^2
- Acceleration threshold: 0.1 m/s^2
- Look-ahead horizon: 30m



Decision-Making Methods

1. Explicit programming
 - Ex: if/then statements
 - Heavy burden on designer
2. Supervised learning
 - Ex: imitation learning
 - Generalizing is often a challenge
3. Optimization / optimal control
 - Ex: MPC
 - Requires a high-fidelity model and lots of computation
4. Planning
 - Given a **stochastic model**, how to algorithmically determine best policy?
5. Reinforcement Learning
 - If model is unknown (or very complex), learn policy through experience



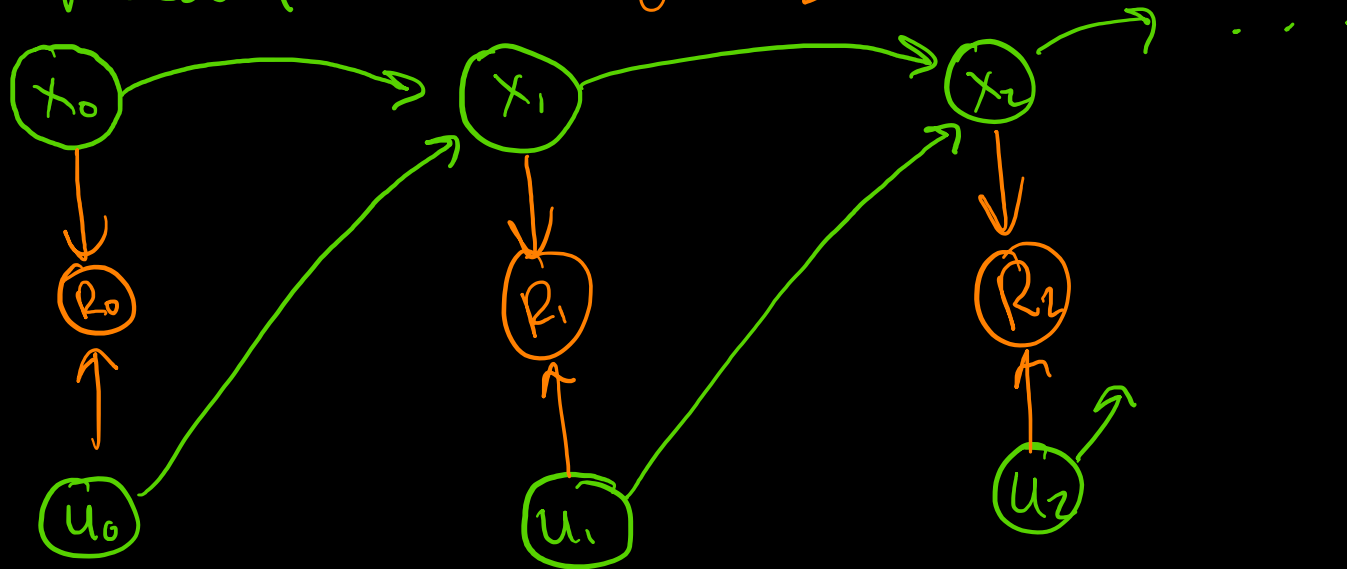
Today's Plan

- Introduction to decision-making
- **Markov Decision Processes**
- MDP Policies and Value Iteration
- Simple Example



Markov Decision Processes (MDPs)

(S, A, T, γ, R)
model design

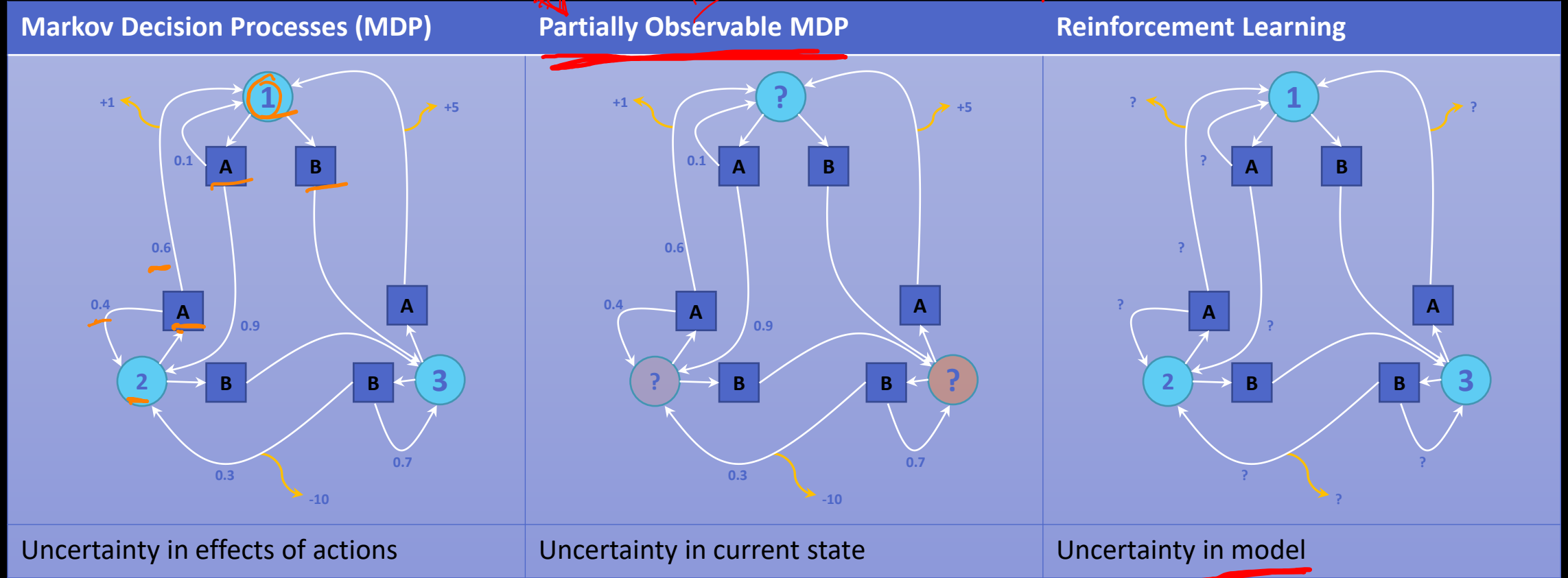


Uncertainty in Motion

- **Markov Decision Processes (MDPs)** model the AV and environment assuming full observability
 - $P(z|x)$: *deterministic* and bijective ✓
 - $P(x'|x, u)$: may be nondeterministic
 - Must incorporate uncertainty into the planner and generate actions for each state
- A policy for action selection is defined for all states



Markov Models



Markov Assumptions and Common Violations

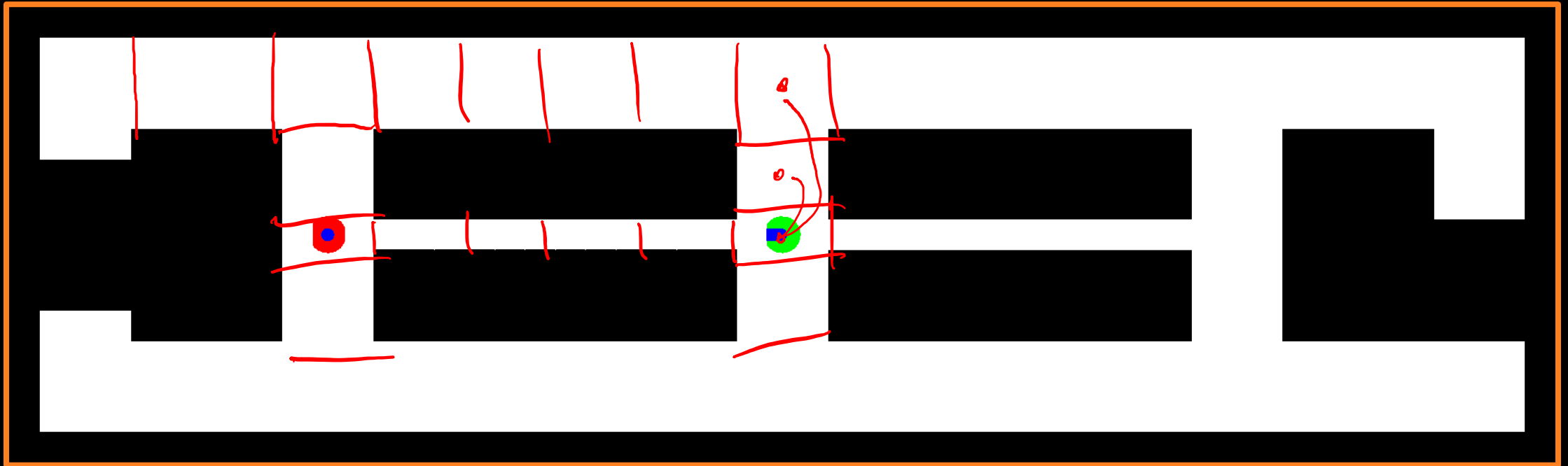
Markov Assumption postulates that past and future data are independent if you know the current state.

What are some common violations?

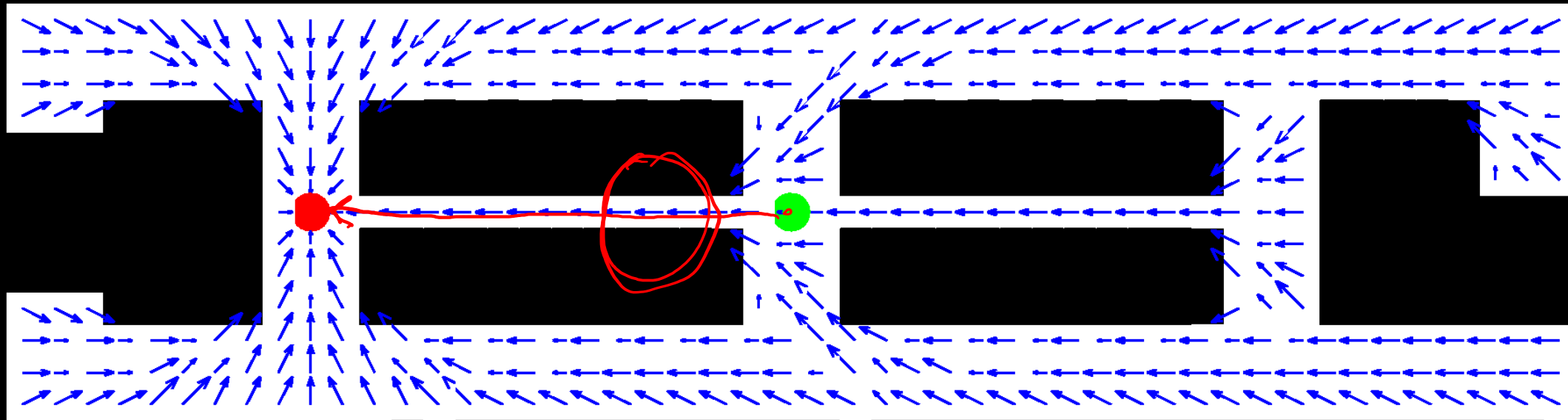
- Unmodeled dynamics in the environment not included in state
 - E.g., moving people and their effects on sensor measurements in localization
- Inaccuracies in the probabilistic model
 - E.g., error in the map of a localizing agent or incorrect model dynamics
- Approximation errors when using approximate representations
 - E.g., discretization errors from grids, Gaussian assumptions
- Variables in control scheme that influence multiple controls
 - E.g., the goal or target location will influence an entire sequence of control commands



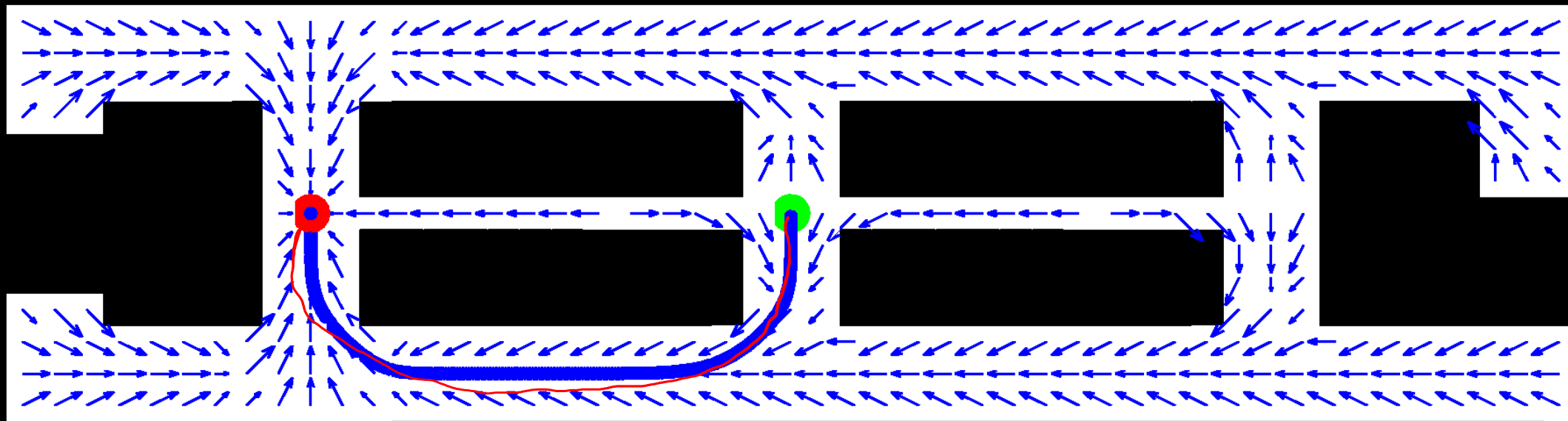
Grid World Example



Deterministic actions



Nondeterministic actions



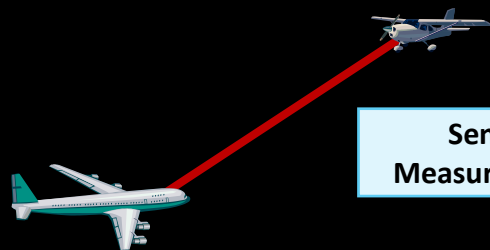
Defining Values

- Actions are driven by goals
 - E.g., reach destination, stay in lane
- Often, we want to reach goal while optimizing some cost
 - E.g., minimize time / energy consumption, obstacle avoidance
- We express both costs and goals in a single function, called the payoff function

ex) $r(x,u) = \begin{cases} 100 & \text{if reach goal} \\ -10 & \text{if collision} \\ -1 & \text{ow} \end{cases}$



Traffic Alert and Collision Avoidance System (TCAS)



Sensor
Measurements

```
IF (ITF.A LT G.ZTHR)
THEN IF (ABS(ITF.VMD) LT G.ZTHR)
THEN SET ZHIT;
ELSE CLEAR ZHIT;
ELSE IF (ITF.ADOT GE P.ZDTHR)
THEN CLEAR ZHIT
ELSE
ITF.TAUV = -ITF.A/ITF.ADOT;
IF (ITF.TAUV LT TVTHR AND
((ABS(ITF.VMD) LT G.ZTHR) OR
(ITF.TAUV LT ITF.TRTRU)))
THEN SET ZHIT
ELSE CLEAR ZHIT
IF (ZHIT EQ $TRUE AND
ABS(ITF.ZDINT) GT P.MAXZDINT)
THEN CLEAR ZHIT
```

Resolution
Advisory



Surveillance

Advisory Logic

Display



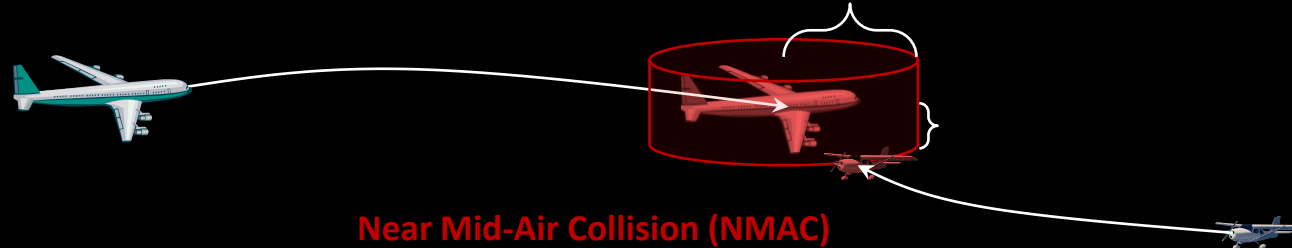
ACAS X: Simplified MDP



State space	Action space
<ul style="list-style-type: none">• Relative altitude• Own vertical rate• Intruder vertical rate• Time to lateral NMAC• State of advisory	<ul style="list-style-type: none">• Clear of conflict• Climb > 1500 ft/min• Climb > 2500 ft/min• Descend > 1500 ft/min• Descend > 2500 ft/min
Dynamic model	Reward model
<ul style="list-style-type: none">• Head-on, constant closure• Random vertical acceleration• Pilot response delay (5 s)• Pilot response strength (1/4 g)• State of advisory	<ul style="list-style-type: none">• NMAC (-1)• Alert (-0.01)• Reversal (-0.01)• Strengthen (-0.009)• Clear of conflict (0.0001)



ACAS X: Simplified MDP



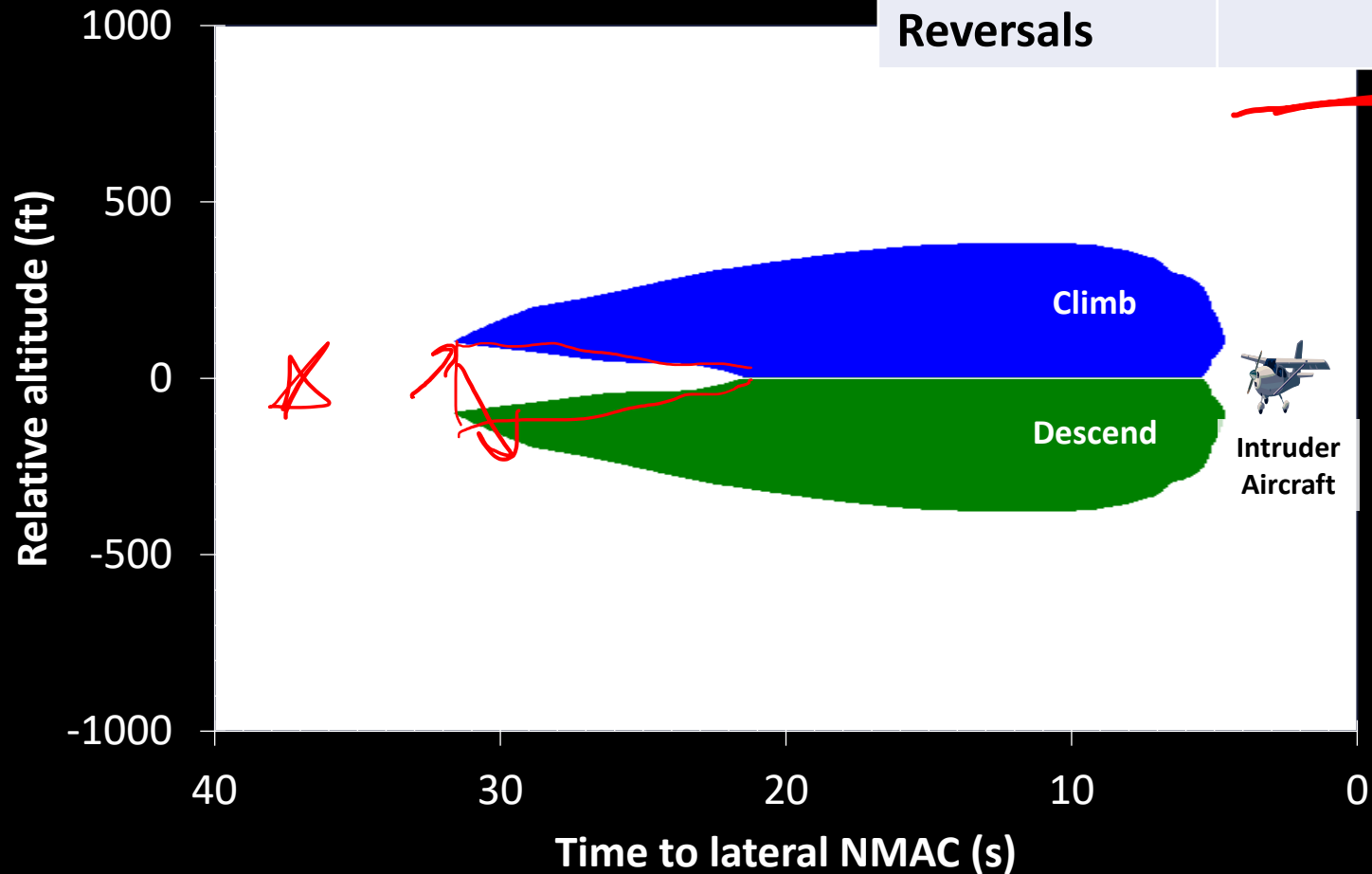
State space	Action space
<ul style="list-style-type: none"> Relative altitude Own vertical rate Intruder vertical rate Time to lateral NMAC State of advisory 	<ul style="list-style-type: none"> Clear of conflict Climb > 1500 ft/min Climb > 2500 ft/min Descend > 1500 ft/min Descend > 2500 ft/min
Dynamic model	Reward model
<ul style="list-style-type: none"> Head-on, constant closure Random vertical acceleration Pilot response delay (5 s) Pilot response strength (1/4 g) State of advisory 	<ul style="list-style-type: none"> NMAC (-1) Alert (-0.01) Reversal (-0.01) Strengthen (-0.009) Clear of conflict (0.0001)



Optimized Logic

Both Own and Intruder Level

Metric	TCAS	ACAS X
NMACs	169	3
Alerts	994,317	690,406
Strengthens	40,470	92,946
Reversals	197,315	9,569



Today's Plan

- Introduction to decision-making
- Markov Decision Processes
- **MDP Policies and Value Iteration**
- Simple Example



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
- Typically, we want a policy that optimizes future payoff, considering optimal actions over a planning (time) horizon



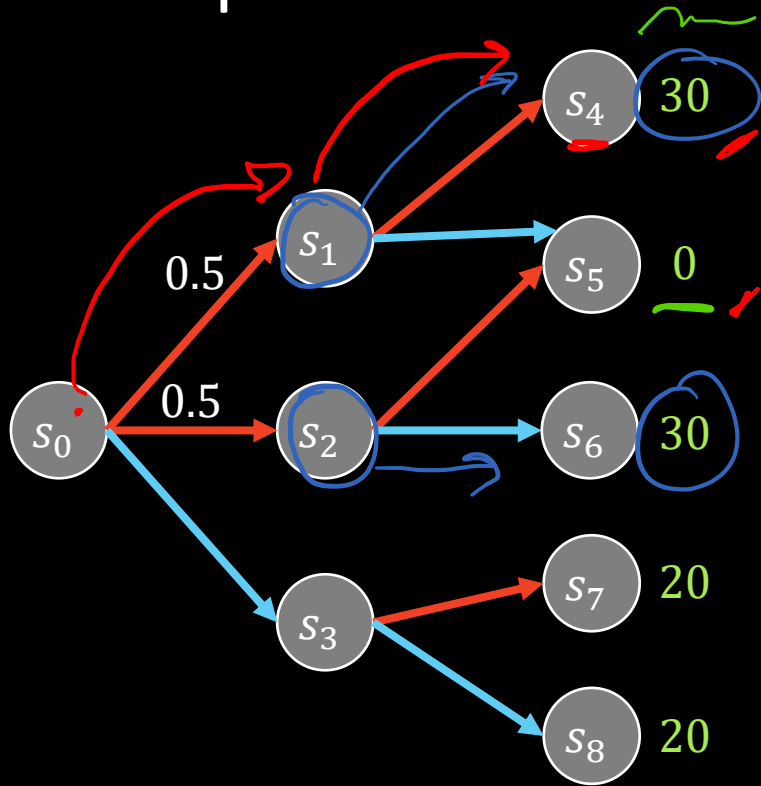
Open vs. Closed Loop Planning

- Closed-Loop Planning: accounts for future information in planning. This creates a reactive plan (policy) that can react to different outcomes over time
- Open-Loop Planning: path planning algorithms develop a static sequence of actions



Open Loop vs. Closed Loop Planning

\uparrow
 a_1, a_2
 \downarrow



$t = 0$

$t = 1$

$t = 2$

open-loop plans: expected R

$$u(a_1, a_1) = 0.5 \cdot 30 + 0.5 \cdot 0 = 15$$

$$u(a_1, a_2) = 15$$

$$u(a_2, a_1) = 20$$

$$u(a_2, a_2) = 20$$

~



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 π :

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

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



Expected Cumulative Payoff

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

Handwritten annotations: "payoff" with an arrow pointing to $r_{t+\tau}$, and "discount factor" with an arrow pointing to γ^τ . The summation index τ and the upper limit T are also circled in blue.

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

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

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



Value Functions

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

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

$T=2$: pick an action that max sum $V_1 + 1$ -step [~]payoff

$$\pi_2(x) = \operatorname{argmax}_u \left[r(x, u) + \int V_1(x') p(x'|x, u) dx' \right]$$

$$V_2(x) = \gamma \max_u \left[\right]$$

For finite T :

$$\pi_T(x) = \operatorname{argmax}_u \left[r(x, u) + \int V_{T-1}(x') p(x'|x, u) dx' \right]$$

$$V_T(x) = \gamma \max_u \left[\right]$$



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}
 - $\hat{V}(x) \leftarrow r_{min}, \forall x$
- 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:
 - $\pi(x) = \operatorname{argmax}_u [r(x, u) + \sum'_x \hat{V}(x') p(x'|x, u)]$



Today's Plan

- Introduction to decision-making
- Markov Decision Processes
- MDP Policies and Value Iteration
- Simple Example



Example: Create an MDP

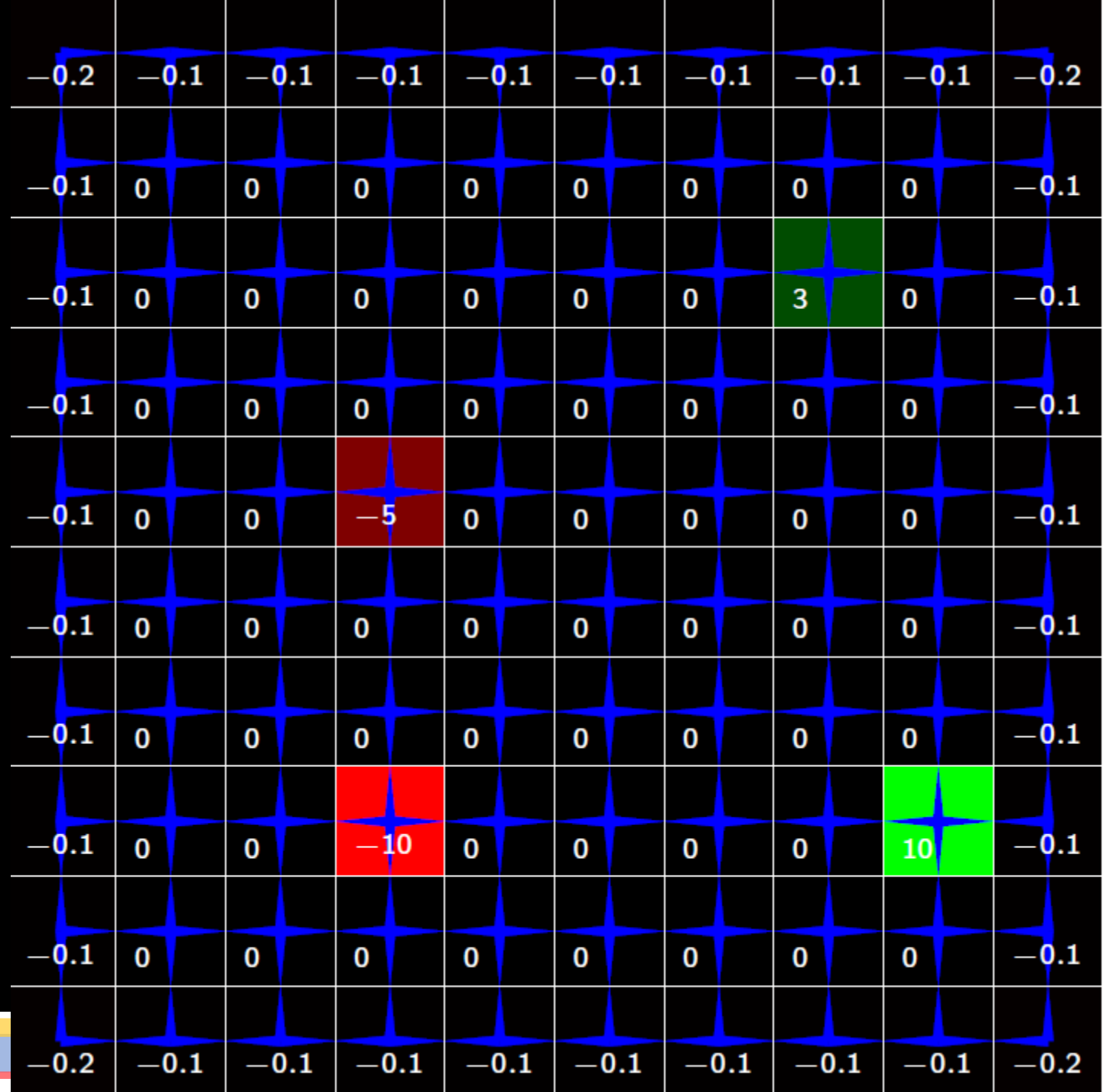


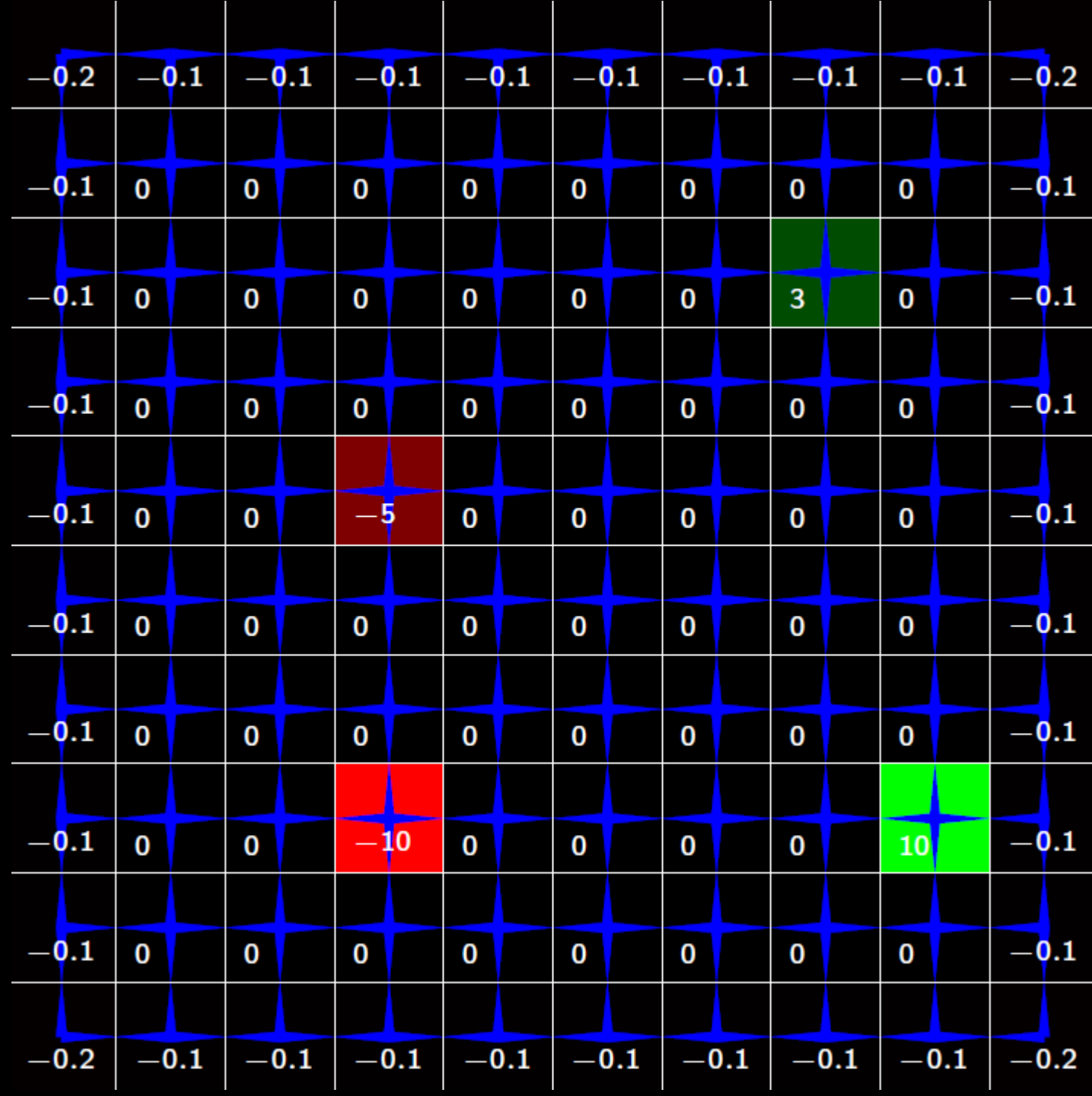
Example: Value Iteration

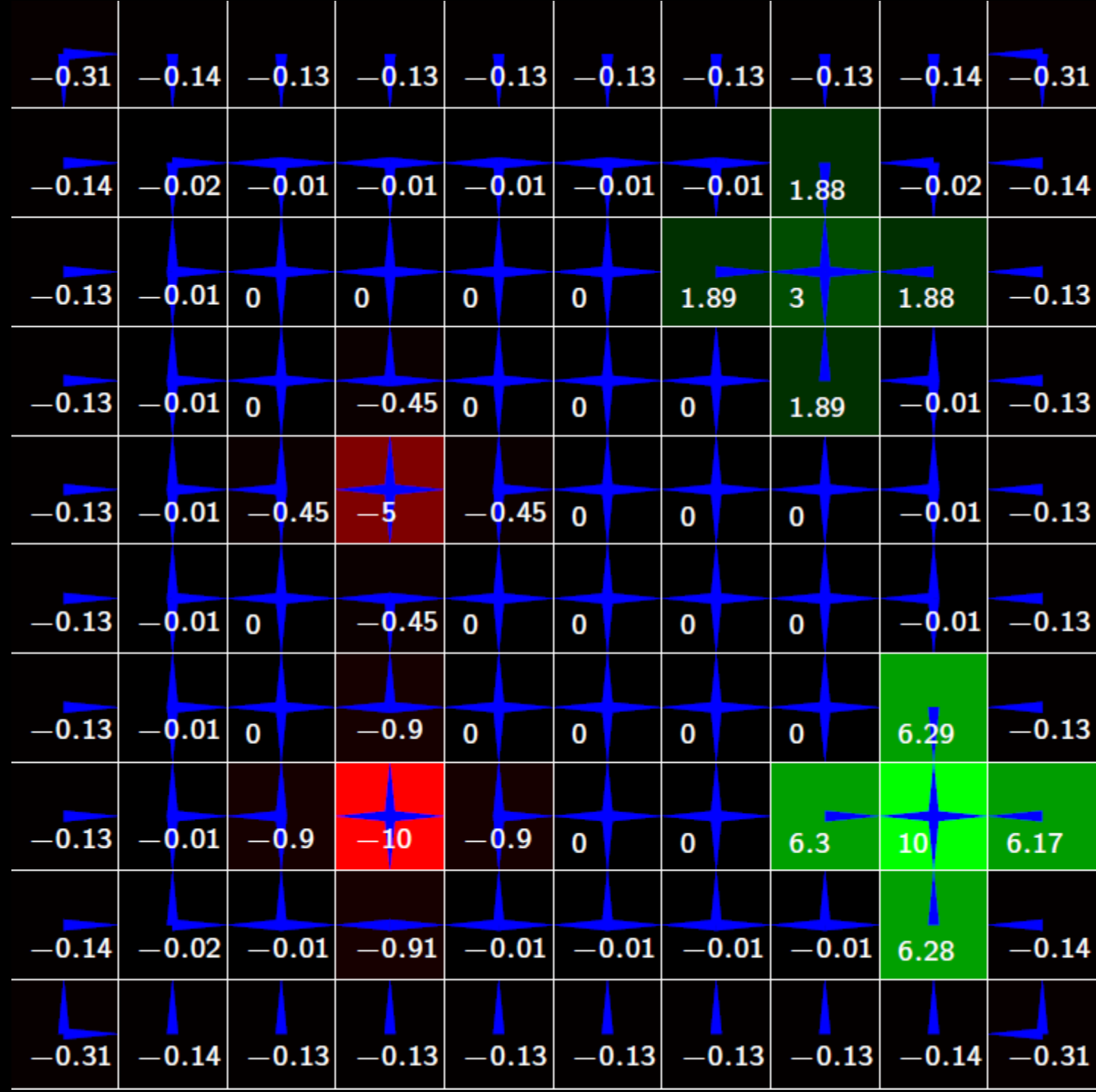


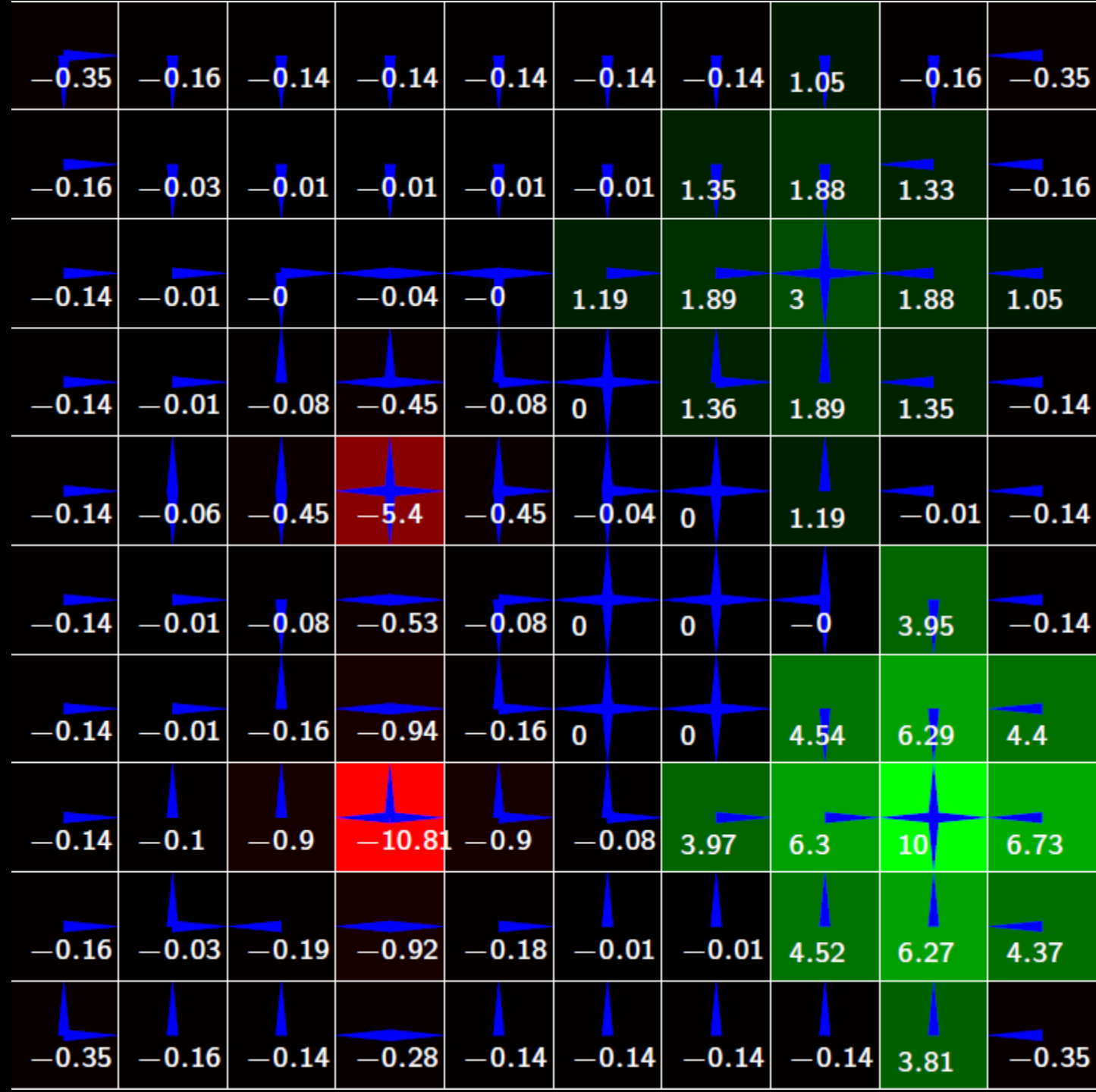
Grid world example

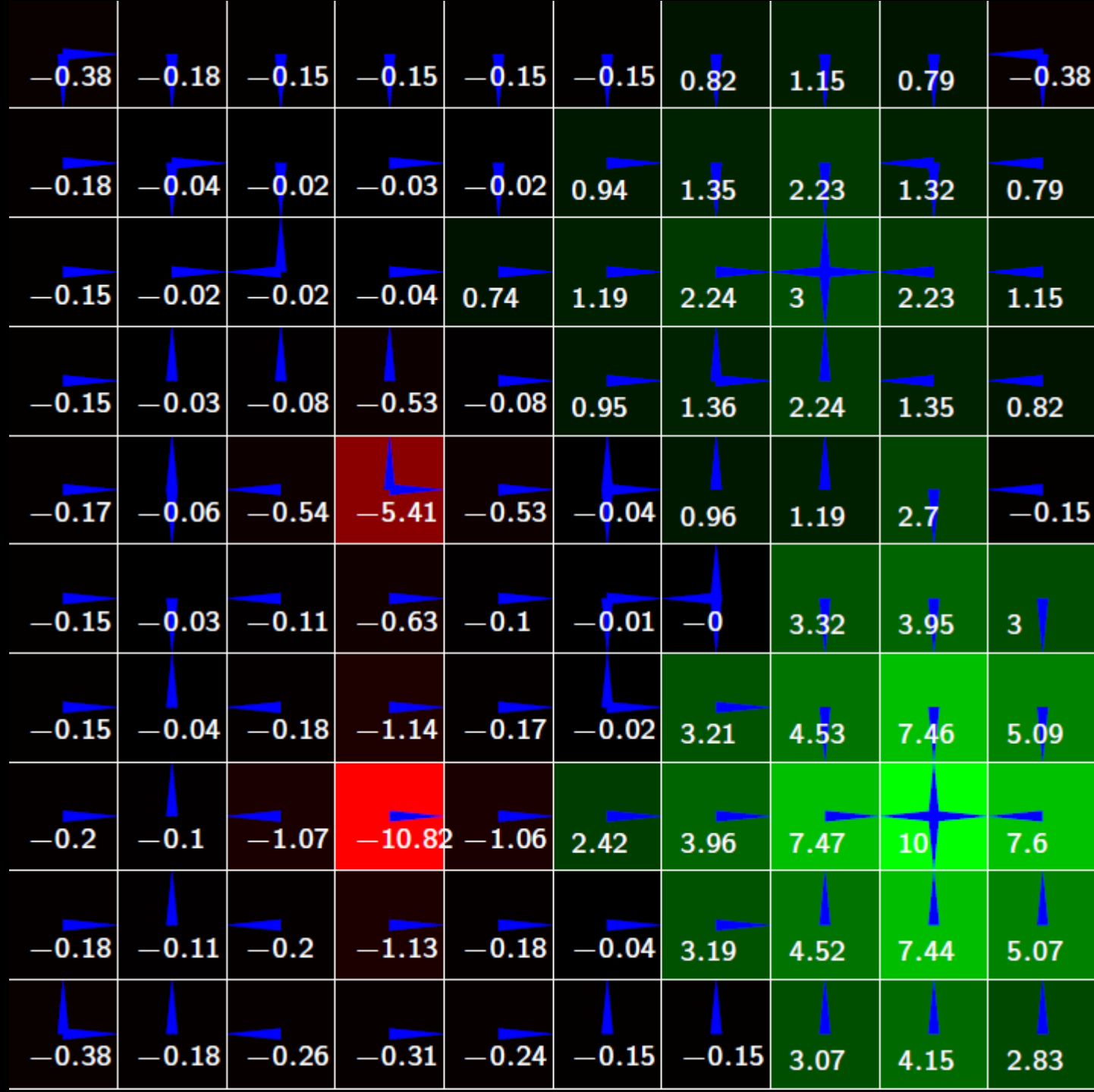
- States: cells in 10×10 grid
- Actions: up, down, left, right
- Transition model: 0.7 chance of moving in intended direction, uniform in other directions
- Reward:
 - two states with cost
 - two terminal states with rewards
 - -1 for wall crash
- Discount is 0.9

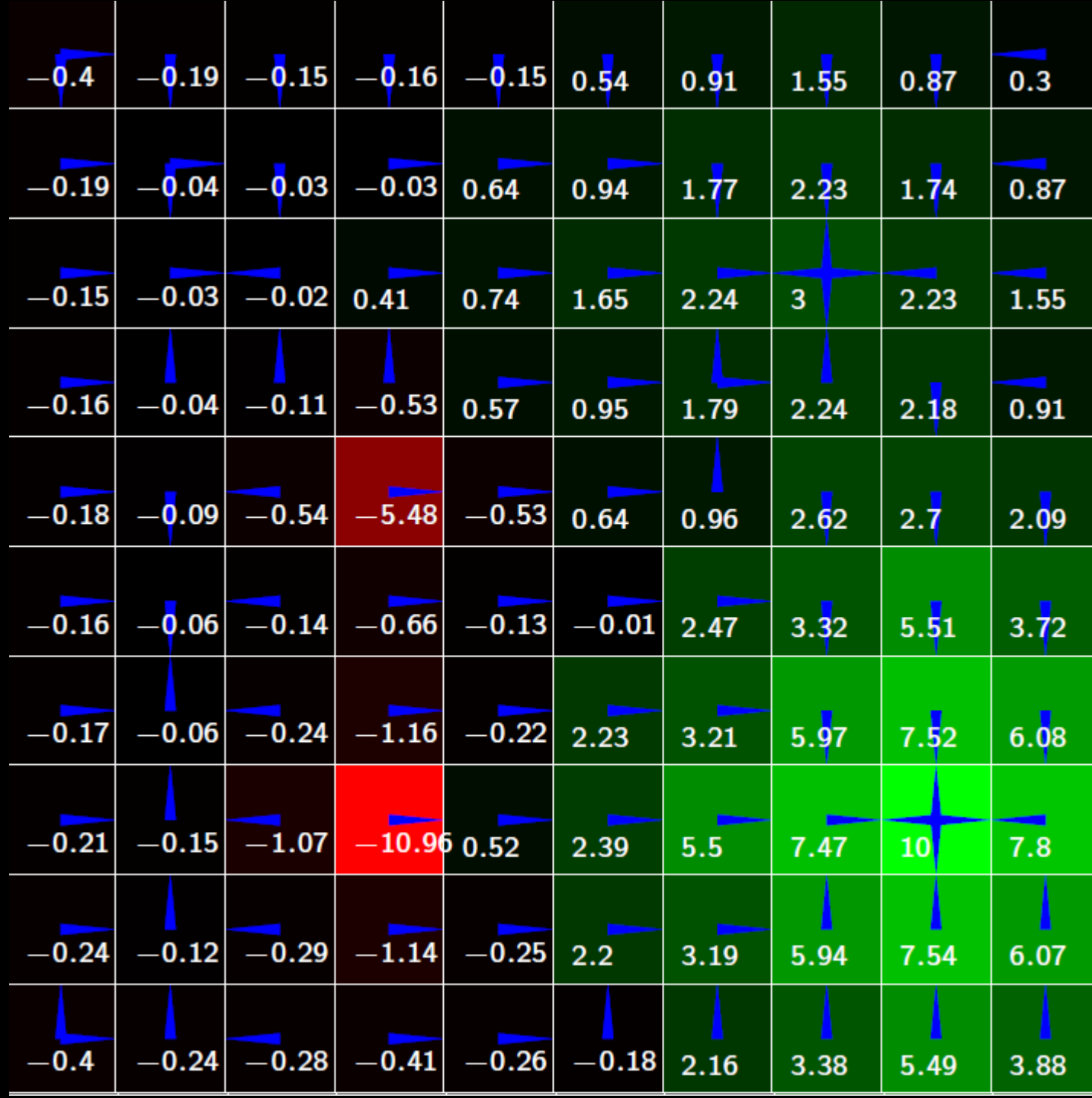


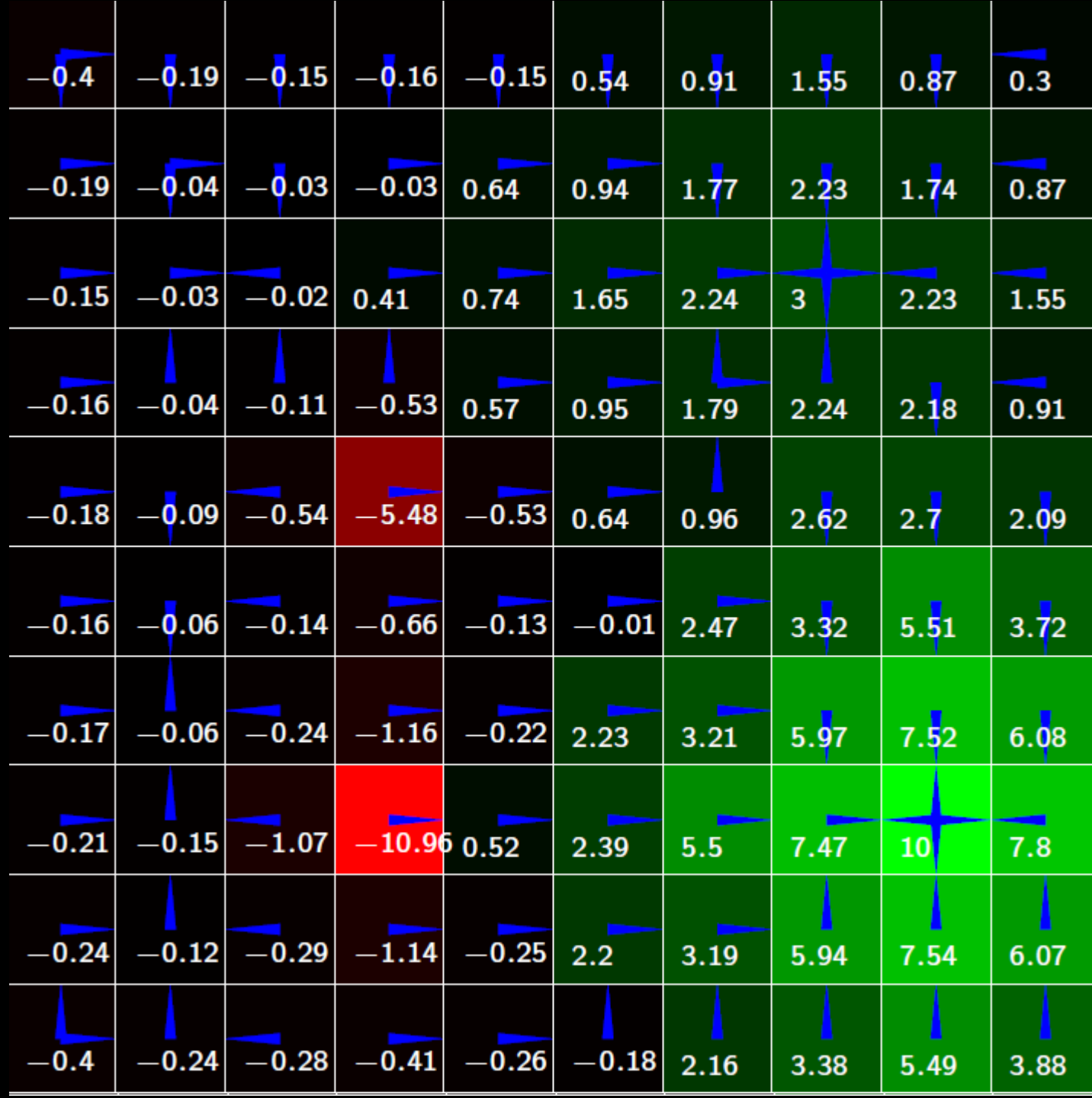


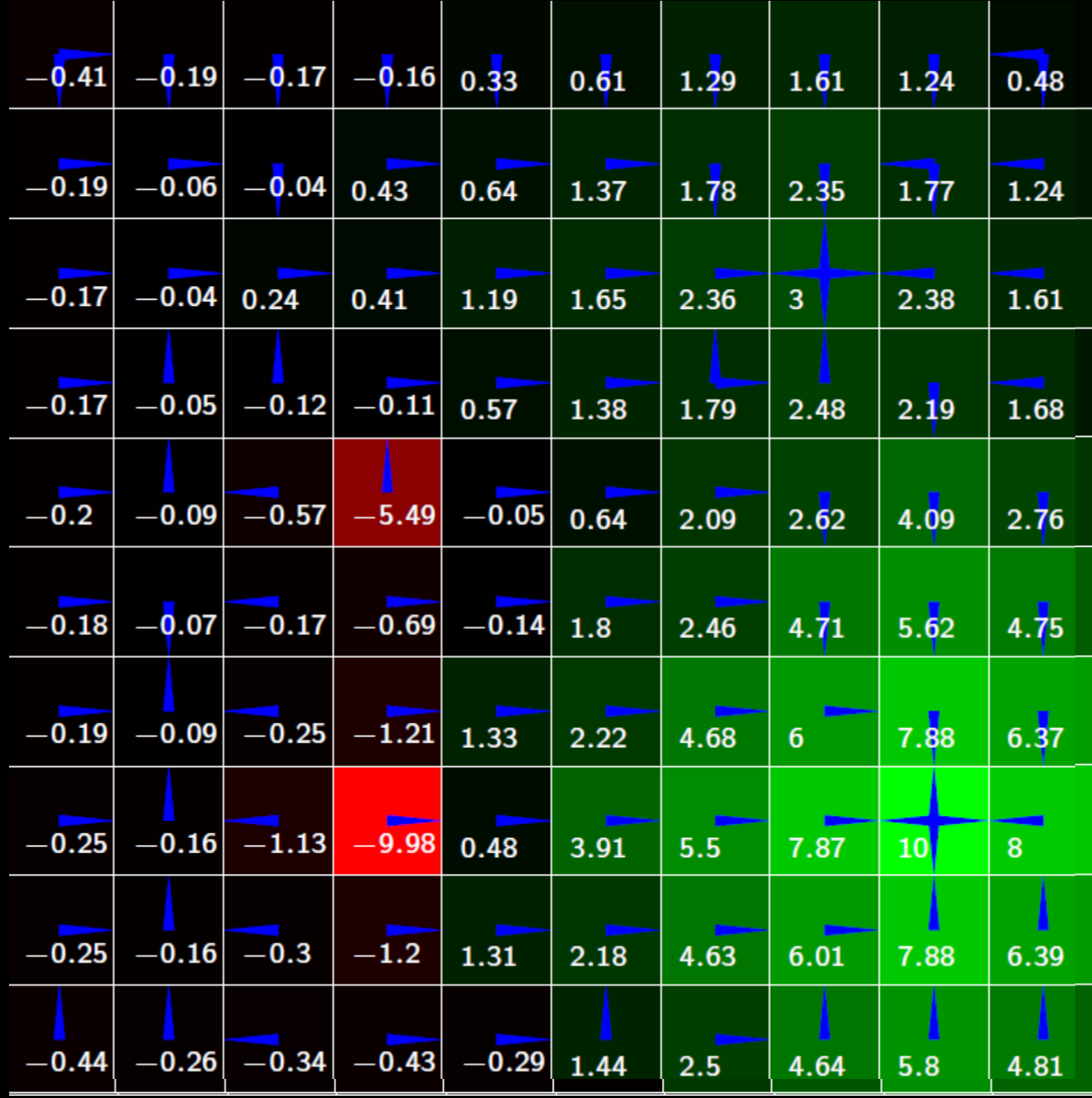


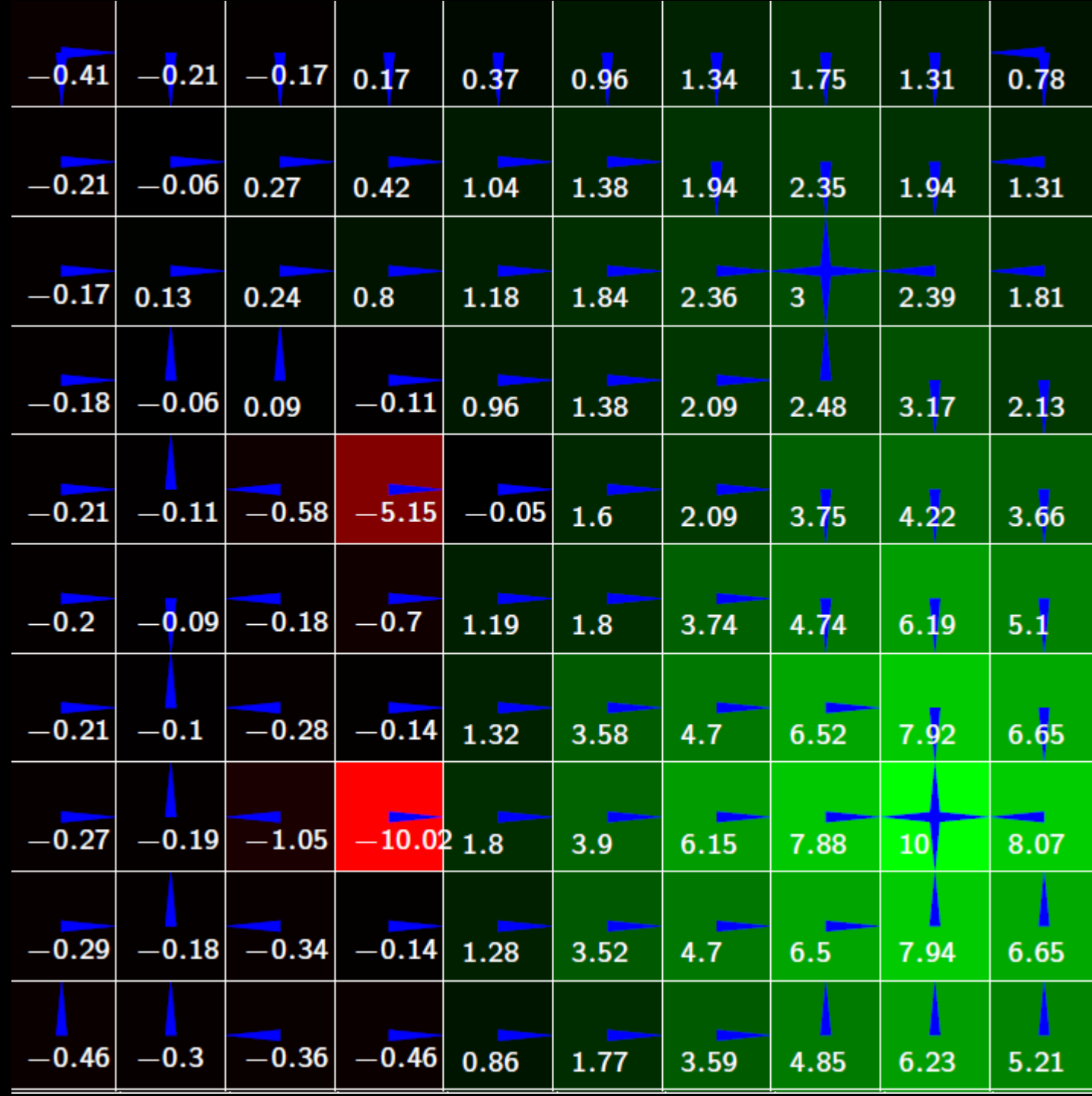


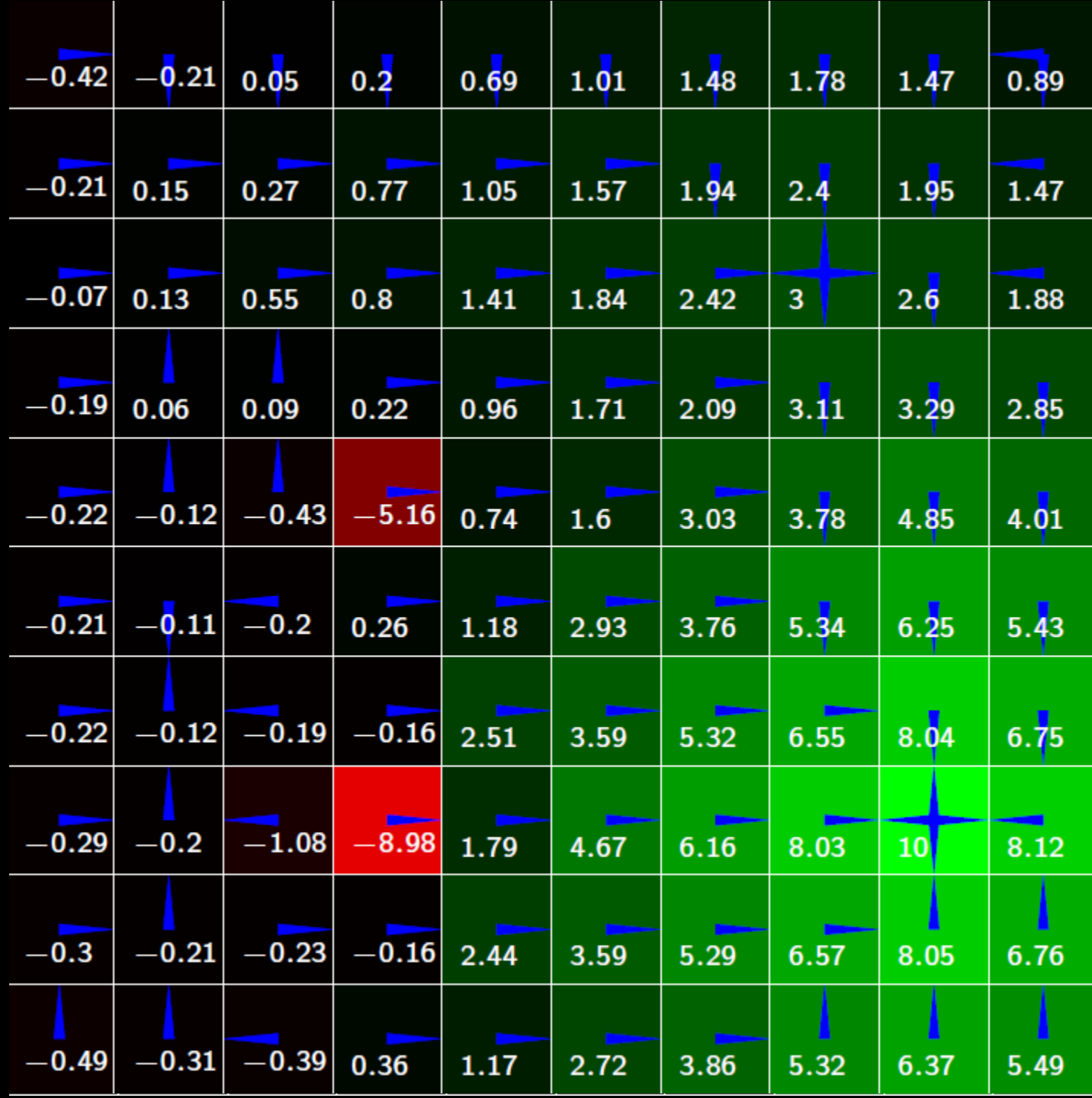


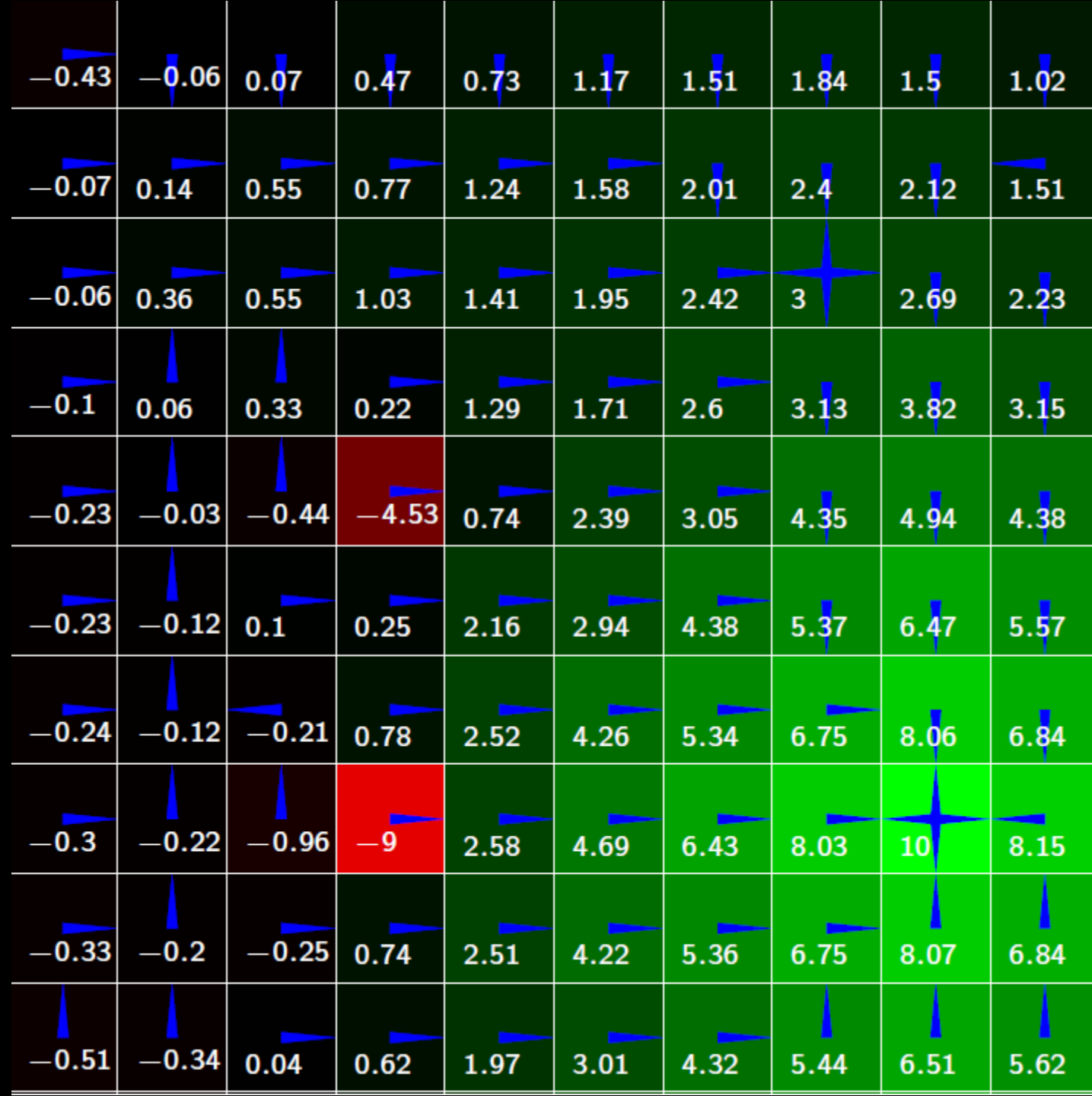


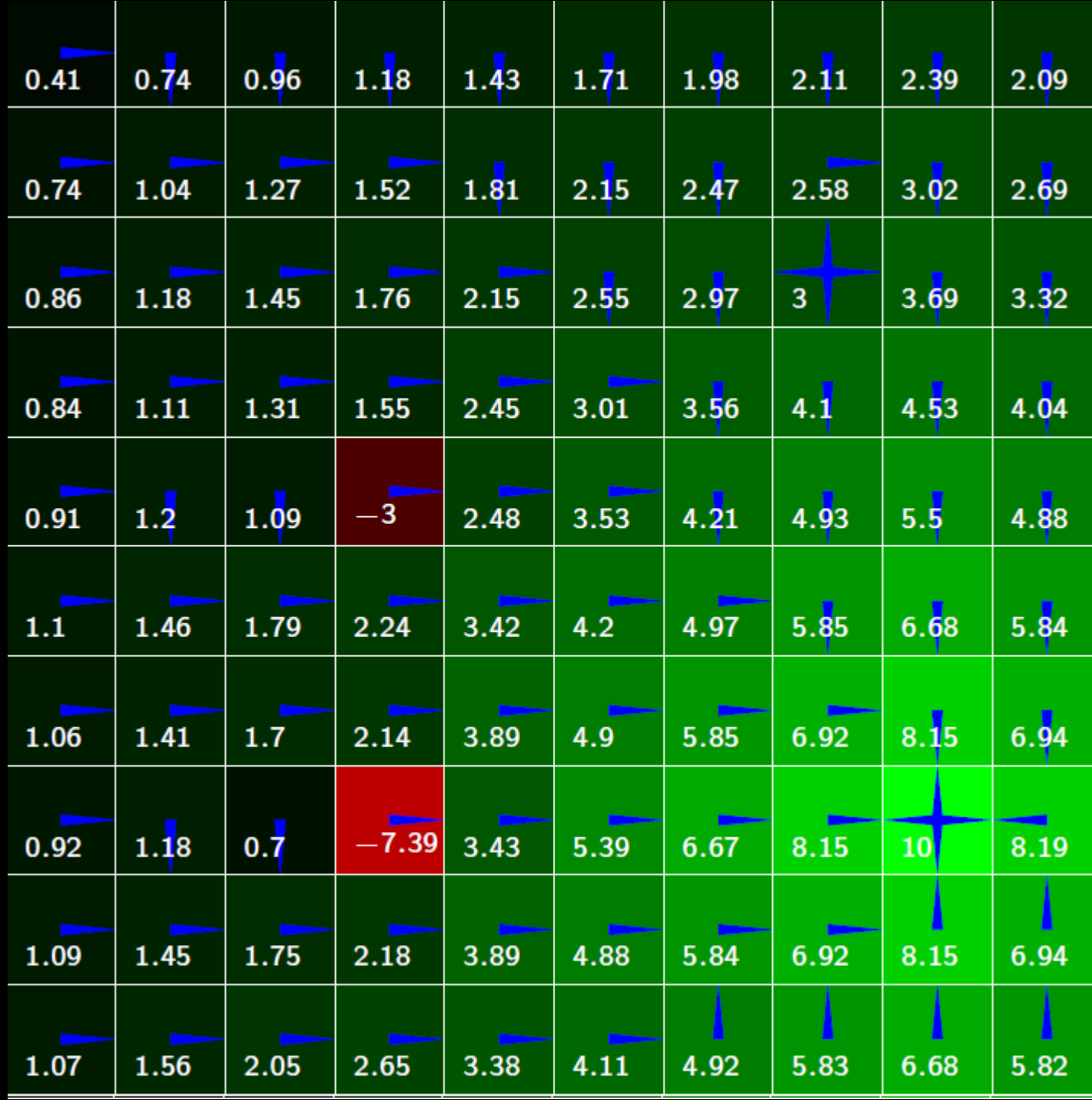




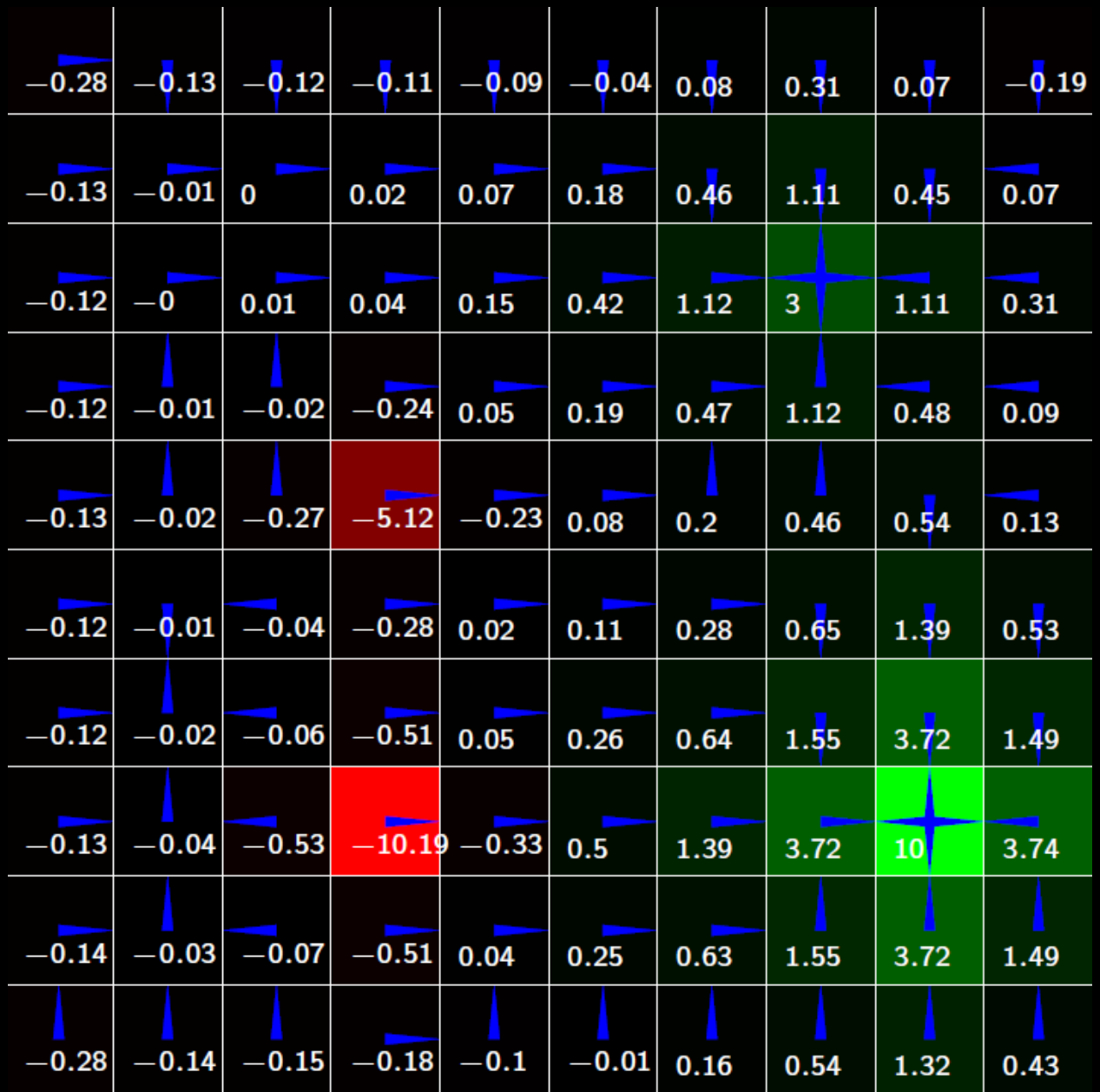
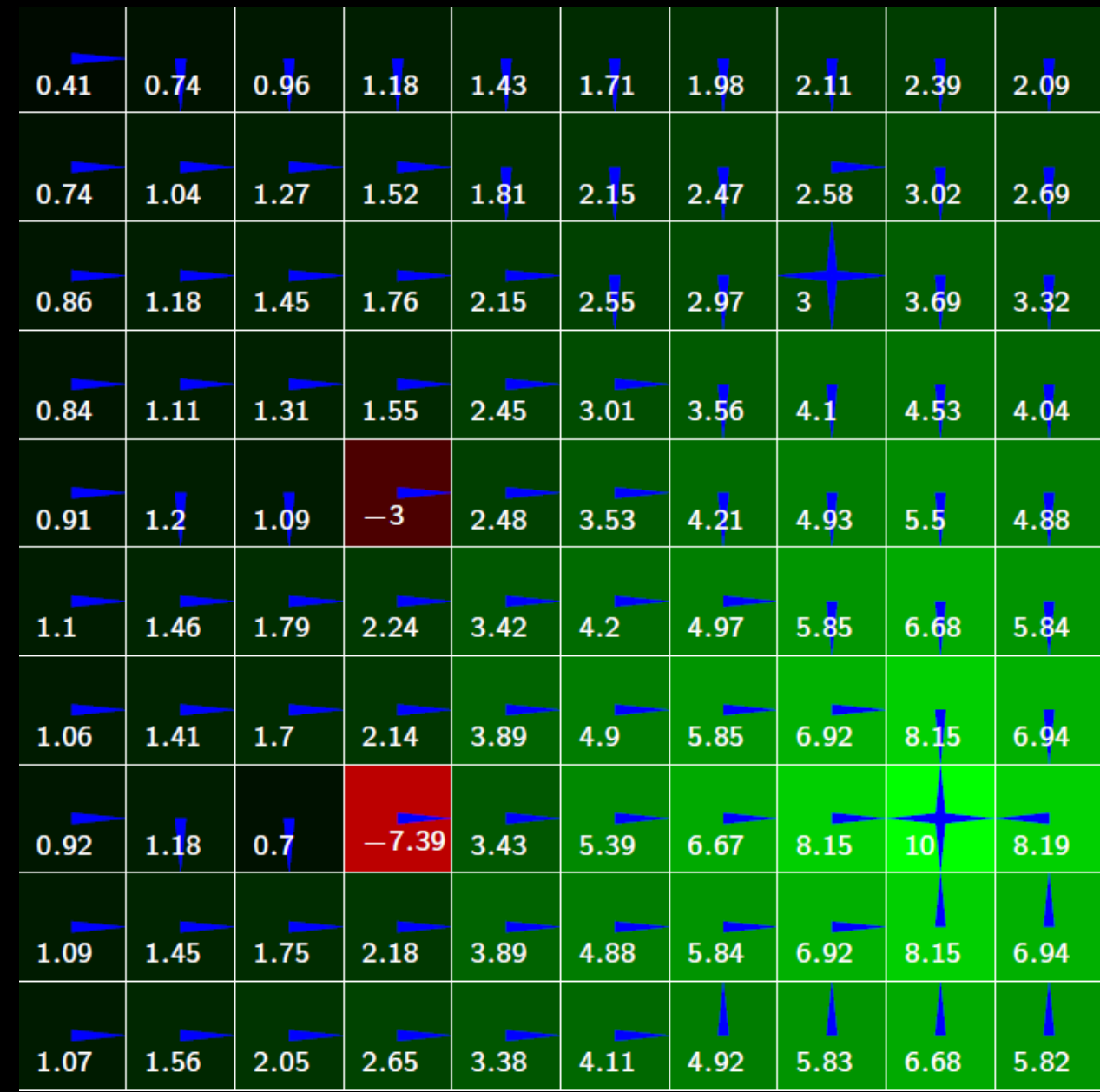








Converged!

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

Summary

- Discussed **decision-making (planning)** schemes and how they fit into the AV stack
- Defined the **MDP** model for decision-making, including **goals, costs, payoff, and policies**
- 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?
 - Will post notes on basic Q-learning for RL!



Extra Slides



Models of Optimal Behavior

- In the finite-horizon model, the agent should optimize expected reward for the next H steps: $E[\sum_{t=0}^H r_t]$
 - Continuously executing H -step optimal actions is known as receding horizon control
- In the infinite-horizon discounted model, agent should optimize:
 $E[\sum_{t=0}^H \gamma^t r_t]$
 - Discount factor is between 0 and 1, can be thought of as interest rate (reward now is worth more than reward later)
 - Keeps the utility of an infinite sequence finite



Challenges

- Value iteration and Policy iteration are both standard, and no agreement on which is better in theory
- In practice, value iteration is preferred over policy iteration as the latter requires solving linear equations, which scales \sim cubically with the size of the state space
- Real-world applications face challenges:
 1. Curse of modeling: Where does the (probabilistic) environment model come from?
 2. Curse of dimensionality: Even if you have a model, computing and storing expectations over large state-spaces is impractical



Mods to Dynamic Programming

Structured Dynamic Programming

- $R(s, a)$ and $U(s)$ can also be represented using a decision tree
- Structured value iteration and structured policy iteration performs updates on leaves of the decision trees instead of all the states
- Structured dynamic programming algorithms improve efficiency by aggregating states, and additive decomposition of reward and value functions

Approximate Dynamic Programming

- For large or continuous spaces, ADP is concerned with finding approx. optimal policies
- This is an active area of research that is conceptually similar to reinforcement learning
- Some approximation methods are:
 - Local approximation relies on the idea that close states have similar values (builds on kNN)
 - Global approximation uses a fixed set of parameters to approximate the value function over the entire state space, generally based on linear regression



Online Methods

- Online methods compute optimal action from current state
 - Expand tree up to some horizon
 - States reachable from the current state is typically small compared to full state space
- Heuristics and branch-and-bound techniques allow search space to be pruned
- Monte Carlo methods provide approximate solutions



Forward Search

Provides optimal action from current state s up to depth d

```
1: function SELECTACTION( $s, d$ )
2:   if  $d = 0$ 
3:     return (NIL, 0)
4:   ( $a^*, v^*$ )  $\leftarrow$  (NIL,  $-\infty$ )
5:   for  $a \in A(s)$ 
6:      $q \leftarrow R(s, a)$ 
7:     for  $s' \in S(s, a)$ 
8:       ( $a', v'$ )  $\leftarrow$  SELECTACTION( $s', d - 1$ )
9:        $q \leftarrow q + \gamma T(s' | s, a)v'$ 
10:    if  $q > v^*$ 
11:      ( $a^*, v^*$ )  $\leftarrow$  ( $a, q$ )
12:   return ( $a^*, v^*$ )
```

Time complexity is $O((|S| \times |A|)^d)$

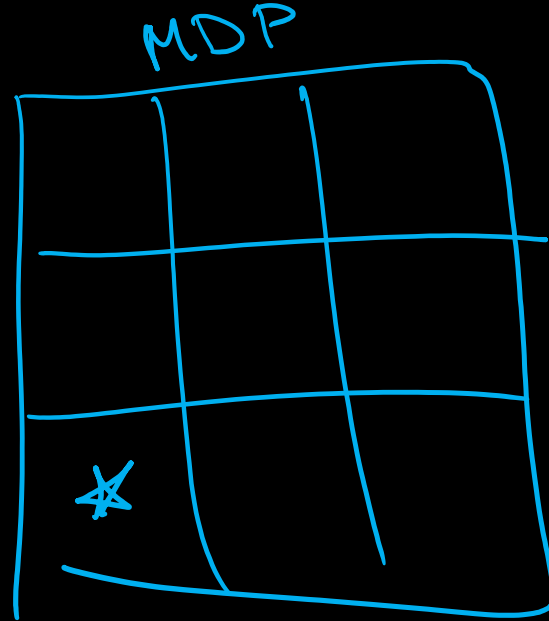
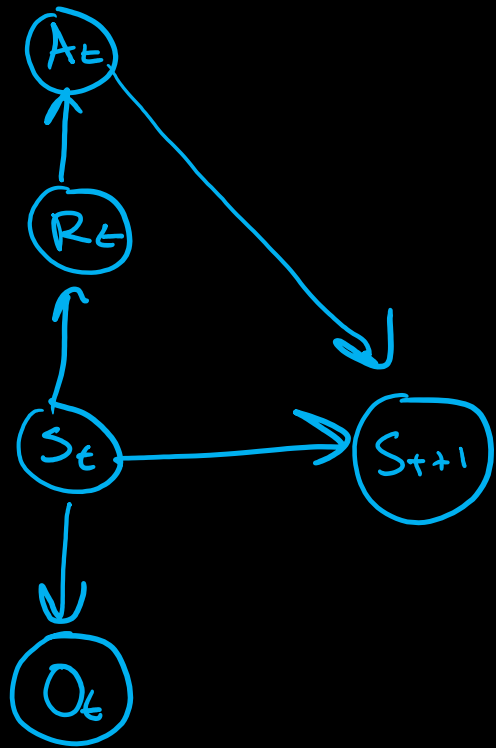


MDP Policy Summary

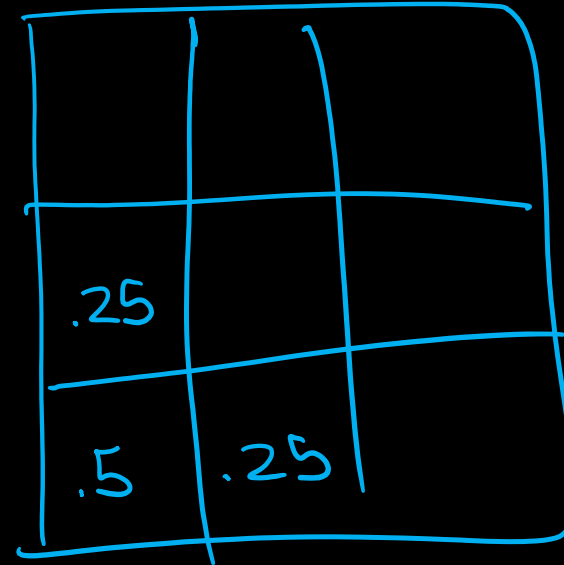
- MDPs represent sequential decision making problems using a transition and reward function
- Optimal policies can be found using dynamic programming
- Problems with large or continuous state spaces can be solved approximately using function approximation
- We generally assumed the transition and reward function are known exactly. On Wednesday, we'll relax this assumption.



Partially Observable MDPs



belief state
MDP



$$\pi : \underbrace{O_{0:t-1}, a_{0:t-1}} \rightarrow a_t$$



POMDP Executions

```
function POMDPPolicyExecution( $\pi$ )  
     $b \leftarrow$  initial belief state  
    loop  
        Execute action  $a = \pi(b)$   
        Observe  $o$  and reward  $r$   
         $b \leftarrow$  UpdateBelief( $b, a, o$ )
```



Alpha vectors

one step horizons

$$U^*(s) = \max_a R(s, a) \quad \rightarrow \quad U^*(b) = \max_a \sum_s b(s) R(s, a)$$

an alpha vector α_a represents $R(\cdot, a)$
 \bar{b} is the vector of beliefs

$$U^*(b) = \max_a (\alpha_a^T \bar{b})$$



alpha vector example

Imagine we have an exam tomorrow, but there is a non-negligible chance I forget about the exam. You can choose to either study or take the evening off.

- If you study and there is an exam, you ace it (R=100)
- If you study and there is no exam, you get nothing (R=0)
- If you relax and there is an exam, you fail and are stressed (R=-100)
- If you relax and there is no exam, you are very happy (R=100)

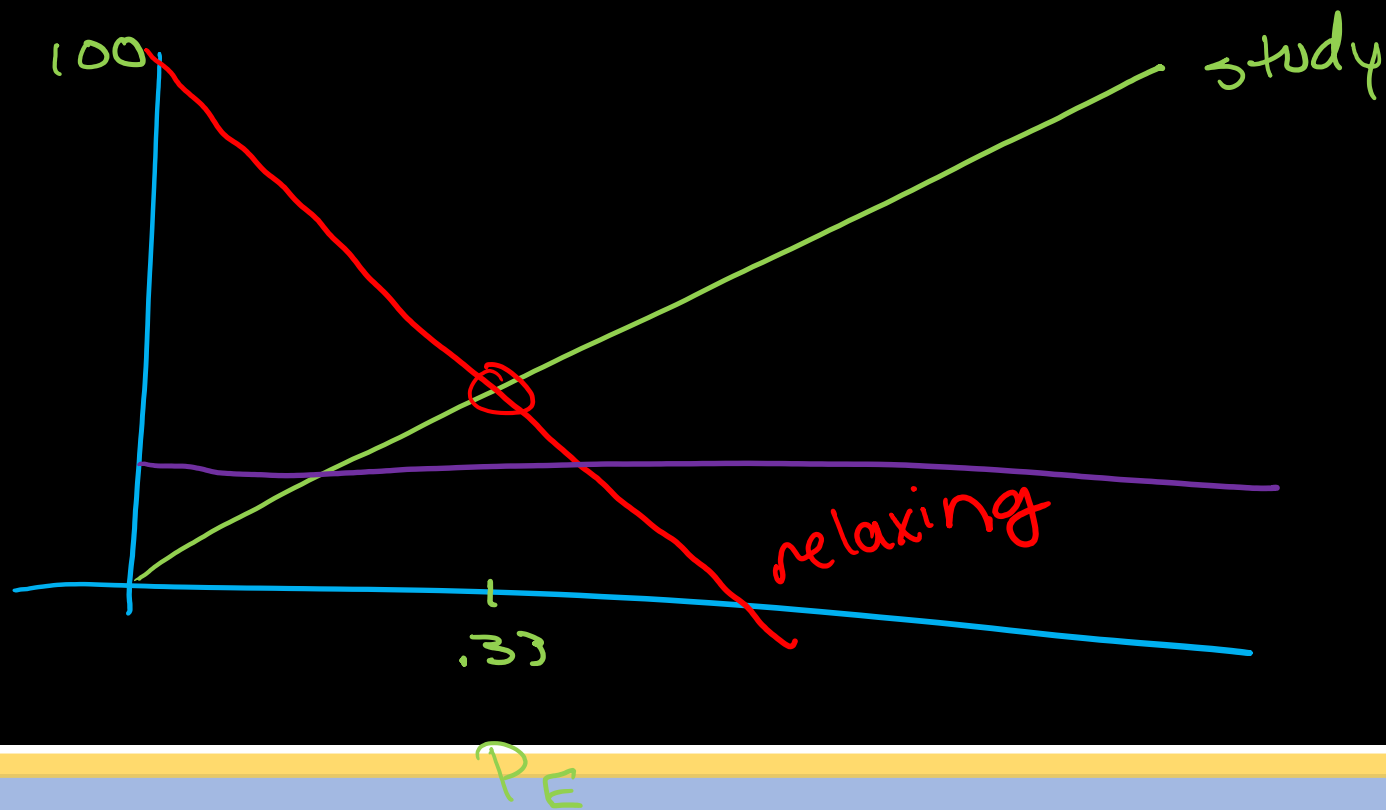
$$\bar{b} = \begin{pmatrix} P_E \\ 1 - P_E \end{pmatrix} \quad \alpha_{\text{study}} = \begin{pmatrix} 100 \\ 0 \end{pmatrix} \quad \alpha_{\text{relax}} = \begin{pmatrix} -100 \\ 100 \end{pmatrix}$$



alpha vector example

$$\alpha_{\text{study}} b = \alpha_{\text{relax}} b$$

$$100 \cdot P_E = -100 P_E + 100(1 - P_E) \rightarrow P_E = .33$$



dropout
 $R = 30$



Why are POMDPs hard to solve?

- Combinatorial explosions!
 - H-step conditional plans: $(|O|^{h-1})/(|O|-1)$, so the number of policies is $A^{(|O|^{h-1})/(|O|-1)}$
 - For a two action two observation problem, there are 2^{63} six step conditional plans
- Instead of solving exactly, we can approximate value iteration and/or solve offline
- There are many great solvers available

