CONSTRUCTOR
UNIVERSITY

Thesis Presentation by Jan Steinmüller

# SAFE-RL
# DUCKIETOWN

# TABLE OF CONTENTS

## 01

**Proposition**

Why do we need safe reinforcement learning?

## 02

**Plan**

How can this be implemented?

## 03

**Implementation**

How does it work

# TABLE OF CONTENTS

## 04

### Videos

How well does it work?

## 05

### Conclusions

What do we learn from this?

## 06
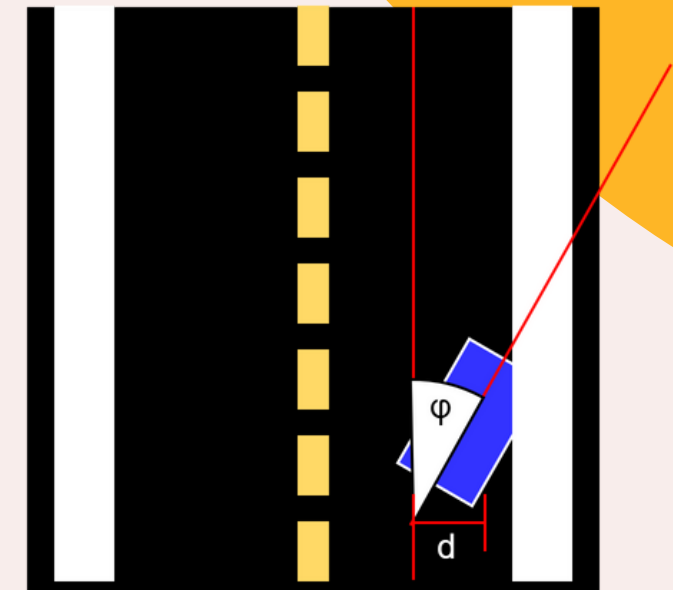
### Questions

Time for some questions

# PROPOSITION

WHY DO WE NEED SAFE REINFORCEMENT LEARNING?

- Reinforcement Learning does not need training data
- Model-free reinforcement learning algorithms
- Agent should not get stuck in local maxima because of the randomness factor
- Can adapt to different environments
- Starts out with completely random actions that could be unsafe
- Will only learn that an action is bad after it is executed

# PLAN

## HOW CAN THIS BE IMPLEMENTED

- Goal is lane following
- Should stay in its own lane at all times
- Reward can be based on distance and angle to the lane center
- Deep-Q-Networks are the deep reinforcement learning agent
  - Continuous state space
- Input is the lane pose given by the lane_filter node
  - lane_d is the distance to the center of the lane
  - lane_phi is the angle in radians
  - buffered and averaged in this implementation
- Safety layer between action selection and execution
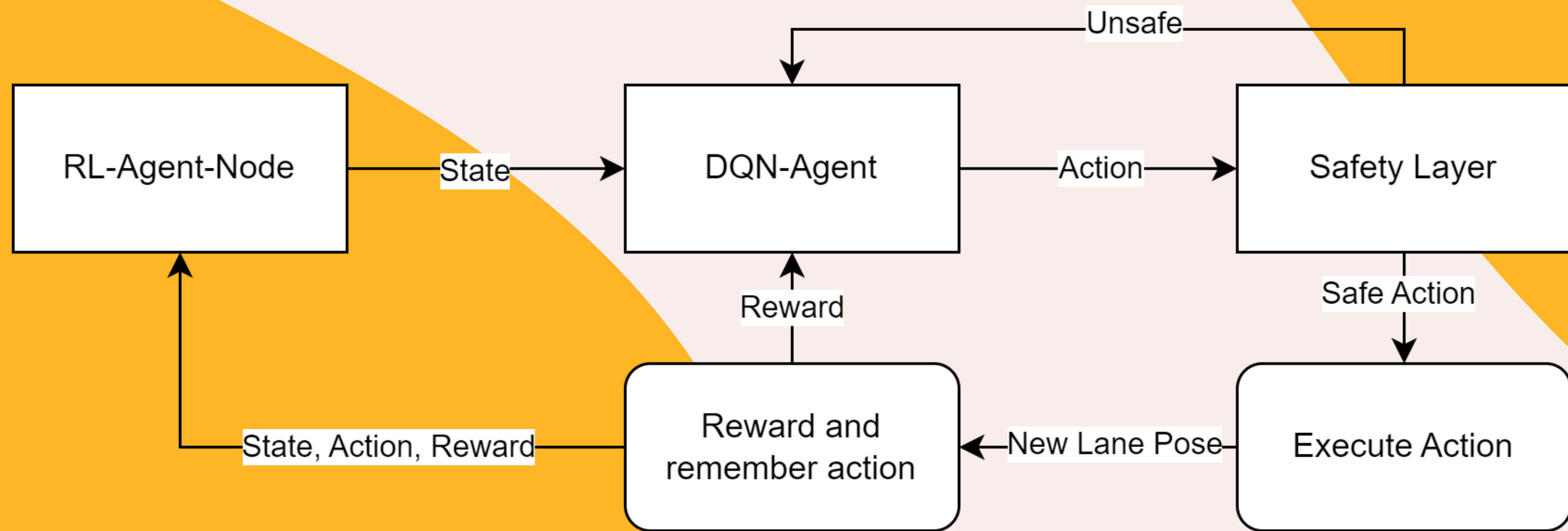
# DEEP-Q-NETWORKS

## HOW DOES DQN WORK?

- DQN is an extension of Q-learning
  - Q-learning has a discrete state table and action list
  - For a given state each action has a Q-Value
  - Q-Values are a combination of the expected reward and the long-term reward

  $$Q\_Value(state, action) = reward + discount\_factor * Q\_Value(state' + action')$$

- DQN introduces a neural network at the state input
  - Improves the scalability of the input layer
  - Continous input data
  - More resistance against noisy data
  - Replay data from executed actions to train the neural network
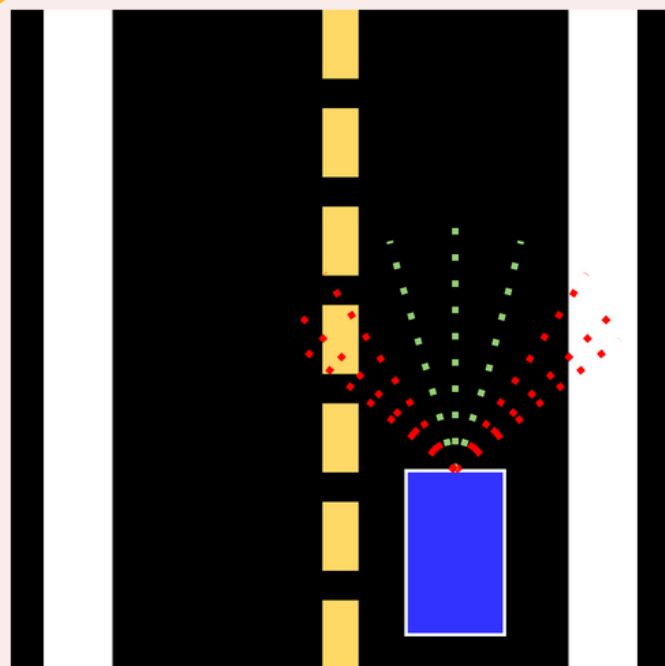  - $reward = 1 - (3 * lane\_d^2) - lane\_\varphi^2$
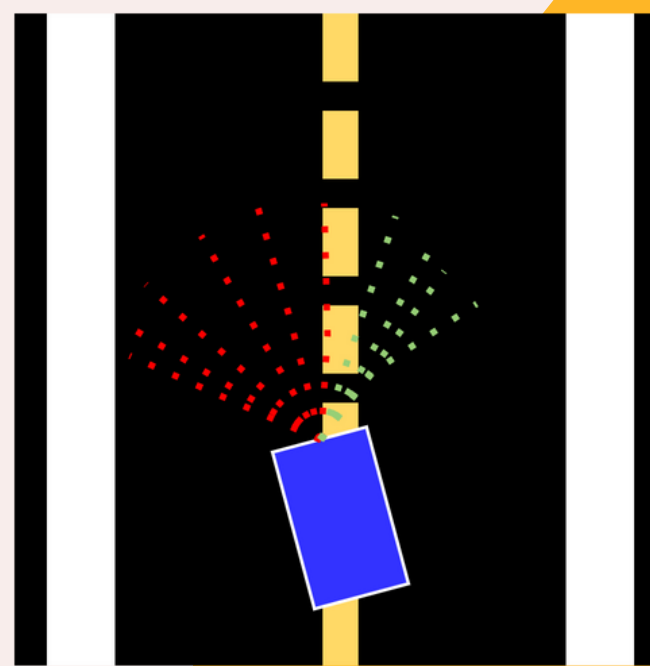
# IMPLEMENTATION

## HOW WAS IT IMPLEMENTED?

# Safety Layer

Predict the future state
with a data driven model
when selecting an action
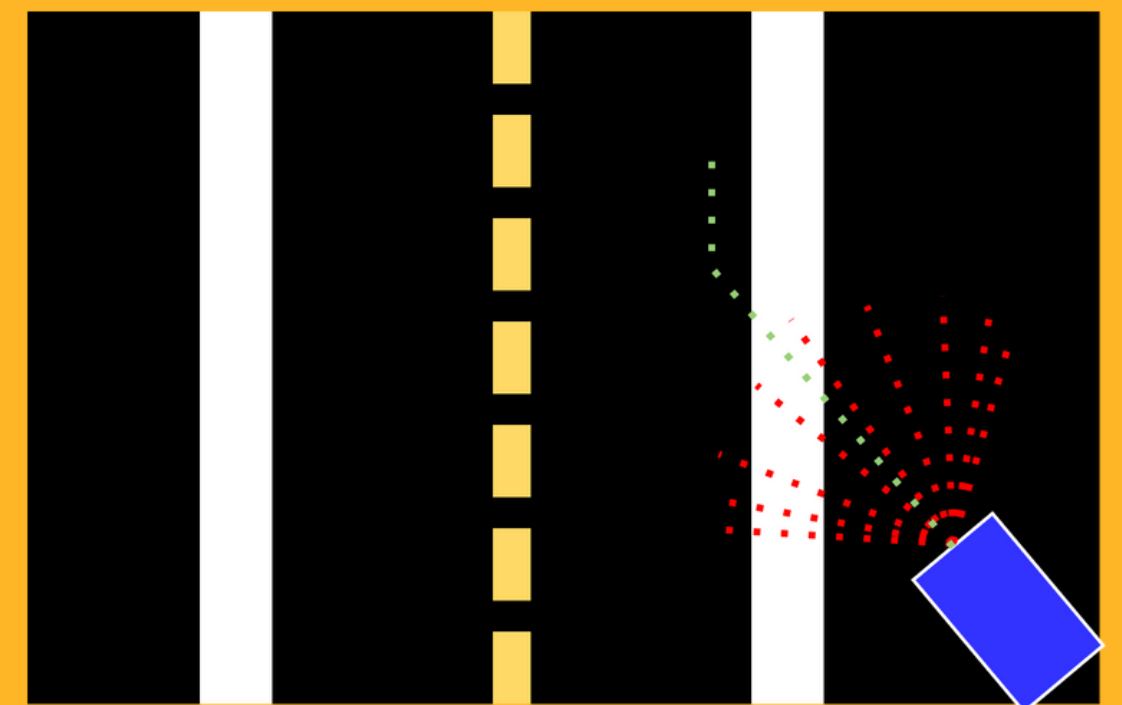
Check if the action is safe
and either optimize the
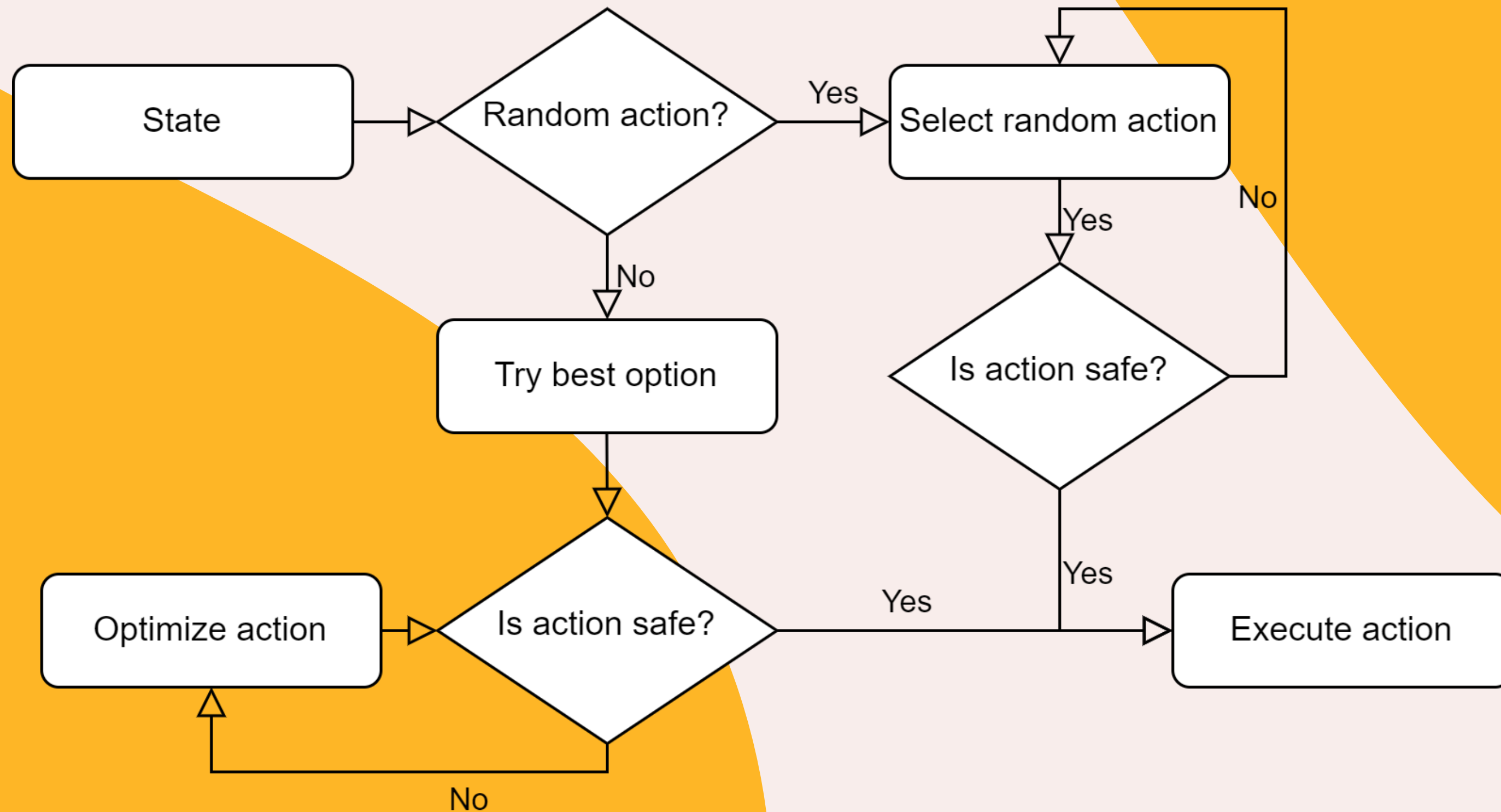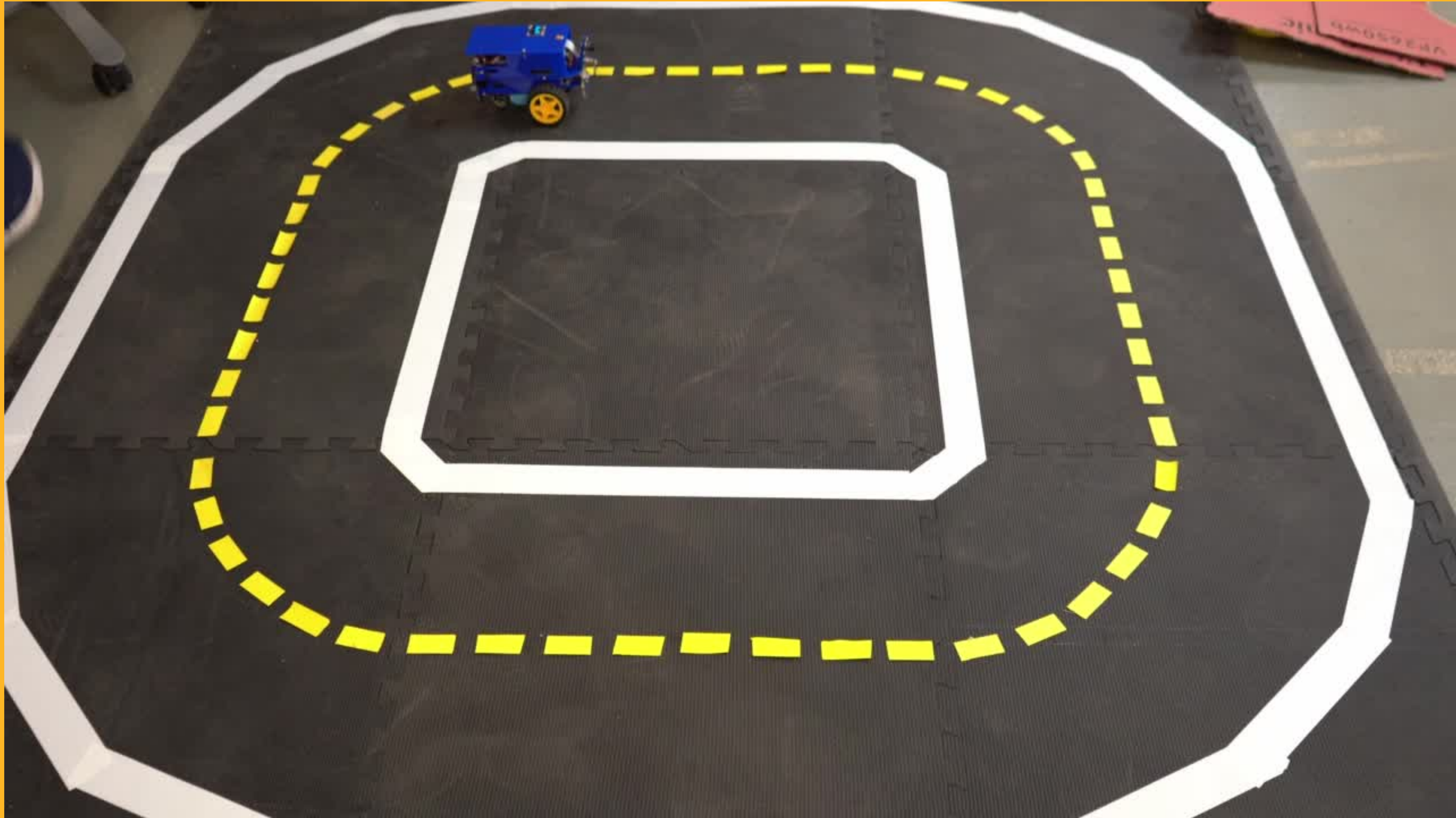action or consider another
action if it is unsafe

Recover from unsafe
states and learning which
unsafe action got us there

# ACT FUNCTION

**HOW ARE ACTIONS SELECTED**

Baseline with no safety layer

DQN with safety layer

# RESULTS

## WHAT DO WE LEARN FORM THE VIDEOS?

- There is a clear safety performance difference
- The safety layer is not perfect in the real world
- The robot mostly manages to recover by itself
  - When the lane gets out of the robot's view, it cannot recover
- The pose estimation by the linear regression model is not perfect
- The robot cannot execute too small velocities
- Too fast velocities make the robot overshoot easily
- The robot has trouble in corners
  - Most likely to the yellow lines being out of view
- Safety layer does not seem to impact execution time too much

# IMPROVEMENTS

**WHAT COULD BE ADJUSTED OR IMPROVED?**

- Adjust the safe lane_d distance

- Adjust angle and distance multiplier in reward

- Modify input data to the agent

- Try out different models for the pose estimation

# Any Questions?

# ACT FUNCTION

**Algorithm 1** act

1: **function** ACT(self, state)
2:     **if** np.random.rand() $\leq$ self.epsilon **then**
3:         $safe \leftarrow$ False
4:         $action\_list \leftarrow$ self.actions.copy()
5:         **while not** $safe$ **do**
6:             $action \leftarrow$ random.randrange(len(action_list))
7:             $predicted\_state \leftarrow$ self.predict_state_lr($state$, self.actions[action])
8:             $safe \leftarrow$ self.safety_layer.check_safety($predicted\_state$)
9:             $action\_list \leftarrow$ action_list.remove(action_list[action])
10:            **if** $safe$ **then**
11:                **Output** "Random action"
12:                **return** $action$
13:            **if** action_list = [] or action_list = None **then**
14:                **return** self.get_back_to_safety($state$)
15:        $q\_values \leftarrow$ self.model.predict($state$, verbose = 0)
16:        $action\_list \leftarrow$ self.actions.copy()
17:        // sorting actions by q_values from high to low
18:        $action\_list \leftarrow [x$ for $\_, x$ in sorted(zip($q\_values[0], action\_list$), reverse = True)]
19:        $action\_list \leftarrow$ action_list[:5]
20:        **for** action in action_list **do**
21:            $predicted\_state \leftarrow$ self.predict_state_lr($state$, action)
22:            $safe \leftarrow$ self.safety_layer.check_safety($predicted\_state$)
23:            **if** $safe$ **then**
24:                **Output** "Learned Action"
25:                **return** self.actions.index($action$)
26:            **else**
27:                self.iter $\leftarrow 0$
28:                $action \leftarrow$ self.optimize($state$, action)
29:                **if** action is not None **then**
30:                    **Output** "Optimized Action"
31:                    **return** $action$
32:        // If no safe action was found, we try to recover to a safe state
33:        **Output** "Using inverse model to get back to safety"
34:        **return** self.get_back_to_safety($state$)

# OPTIMIZE

**Algorithm 1** Optimize Action

---

    **function** OPTIMIZE(self, state, action)

2:        **if** self.iter $> 3$ **then**

            **return None**

4:        $new\_v \leftarrow \text{action}[0]/2$

        $new\_omega \leftarrow \text{action}[1]/2$

6:        $new\_action \leftarrow [new\_v, new\_omega]$

        $new\_state \leftarrow \text{self.predict\_state\_lr}(state, new\_action)$

8:        **Output** "Predicted state: ", new_state, " with action: ", new_action, ""

        // checking if new state is safe

10:      $safe \leftarrow \text{self.safety\_layer.check\_safety}(new\_state)$

        **if** $safe$ **then**

12:            **return** $new\_action$

        **else**

14:            **Output** "New action  is unsafe, trying again".format($new\_action$)

            self.iter $\leftarrow$ self.iter $+ 1$

16:            **return** self.optimize($state, new\_action$)

---

# RECOVER

**Algorithm 1** Get Back To Safety

**function** GET_BACK_TO_SAFETY(state)
    dt ← self.sleep_time
    // Getting back to a safe state
    // If we are angled away from the lane
    **if** (state[1] > 0 and state[0] > 0) or (state[1] < 0 and state[0] < 0) **then**
        // Turning towards lane
        theta ← -1.5 * state[1] / dt
        v ← 0
    **else**
        theta ← state[1] / dt
        v ← abs(state[0]) * np.sin(state[1]) * dt

    // Getting new action
    new_action ← [v, theta]
    OUTPUT "Using inverse model to get back to safety"
    **return** new_action