

Q-Learning – How to?

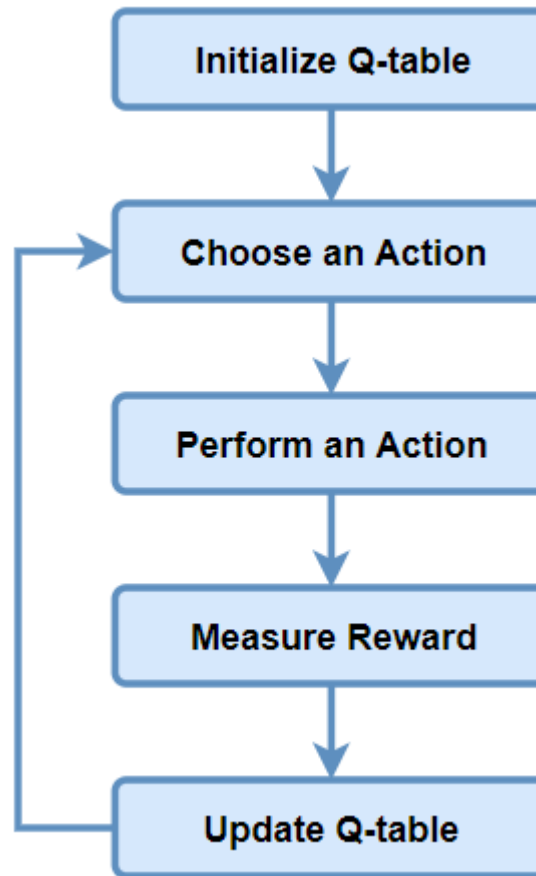
Institute of Automation and Information Systems

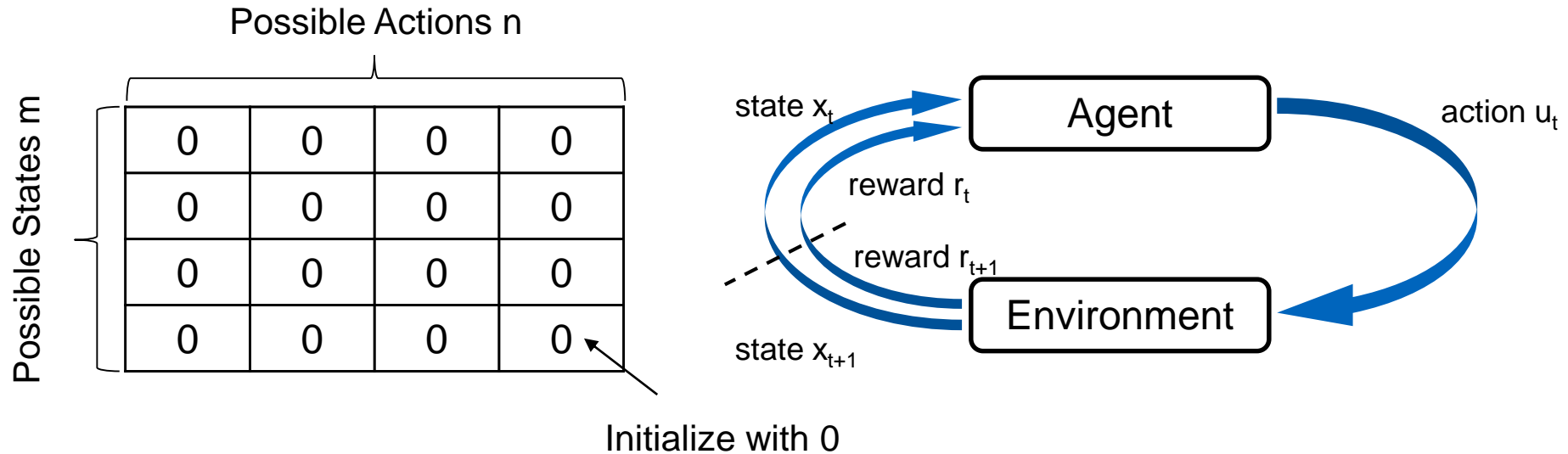
Technical University of Munich



- 1** **Q-Learning Process Steps**
- 2** **Exercise: Grid World**
- 3** **Exercise: Crane-Simulation**

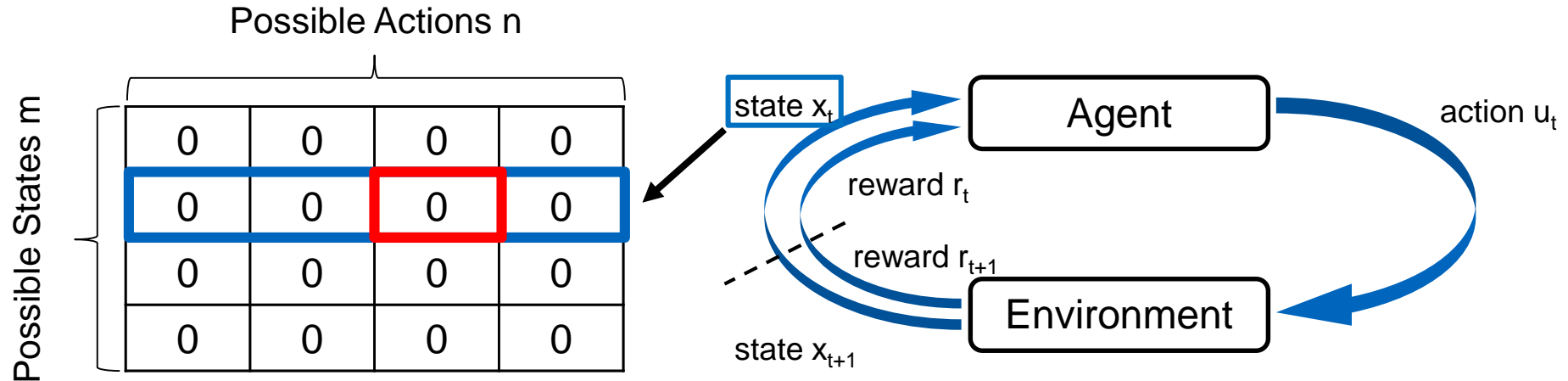
Q-Learning Process Steps





1. Initialize Q-Table

- Use discretization for environments with continuous state variables
- If the state is described by multiple state variables X_1, X_2, \dots, X_k (such as position, velocity and acceleration) with corresponding number of discretization intervals m_1, m_2, \dots, m_k , the Q-Table is of dimension $m_1 \times m_2 \times \dots \times m_k \times n$

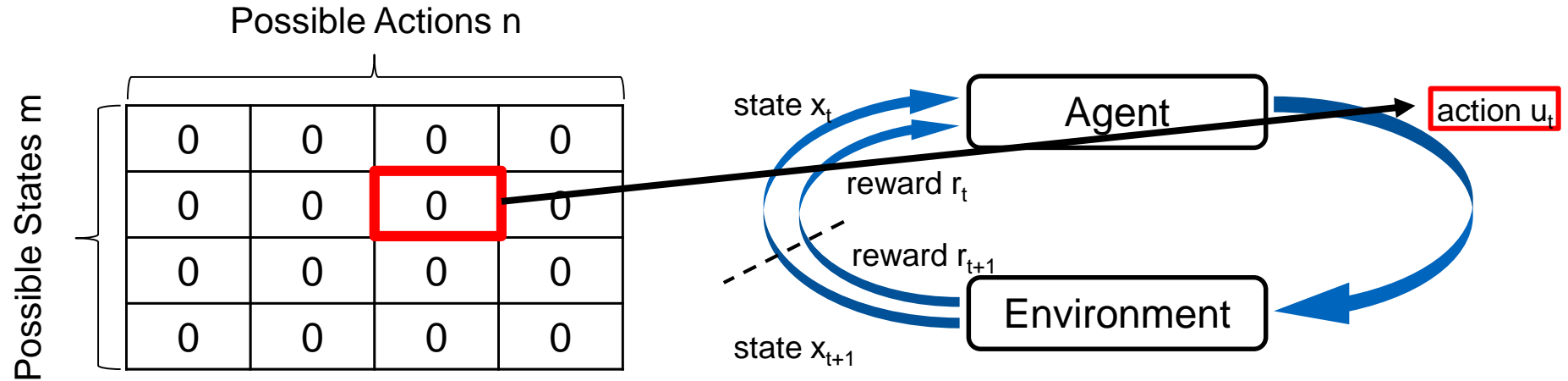


2. Choose an Action:

- Get the current state x_t
- Perform ϵ -greedy strategy:

$$\text{Action} = \begin{cases} \max_u Q(x, u), & R > \epsilon \\ \text{Random } u, & R \leq \epsilon \end{cases}, \text{ R is uniform random value in } [0,1] \text{ and } \epsilon \in [0,1]$$

→ when highest Q-Value occurs more than once, choose randomly from actions with this value

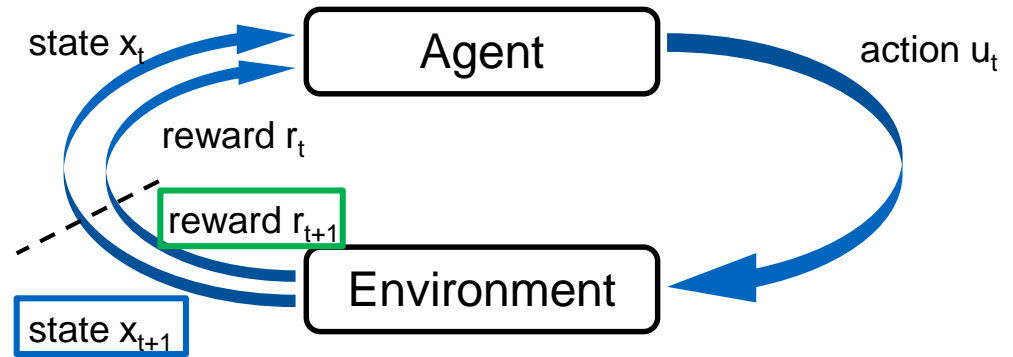


3. Perform an Action

Possible Actions n

Possible States m

0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0



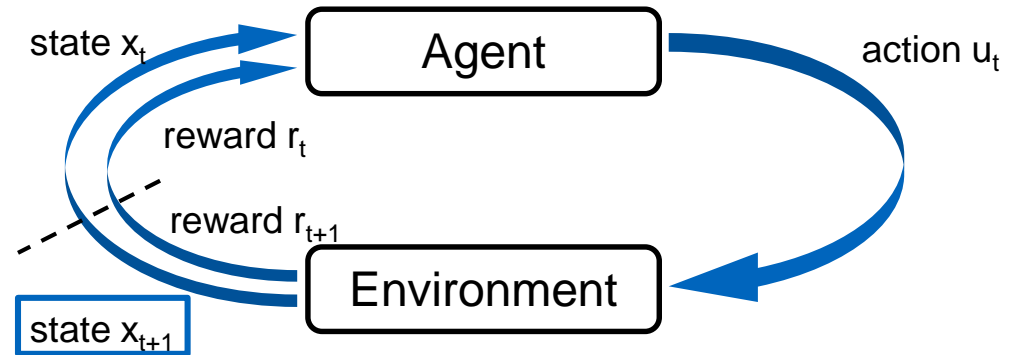
4. Measure reward

And remember new state x_{t+1}

Possible Actions n

Possible States m

0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0



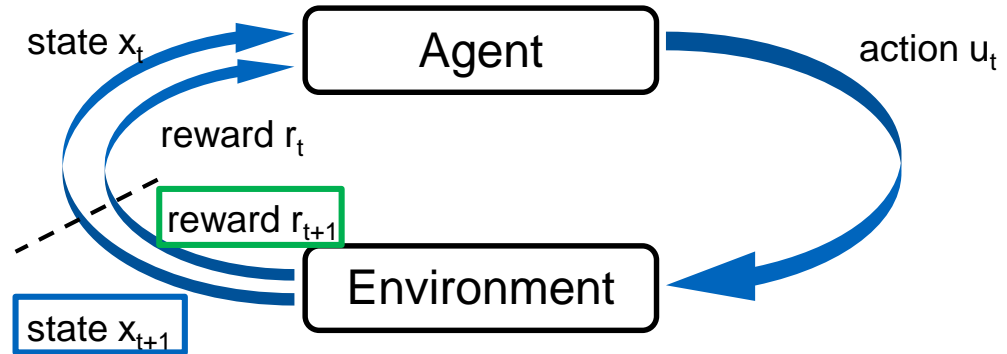
5. Update Q-Table

- Go to new state x_{t+1}
- Find the highest possible Q-value of state x_{t+1} (Q-value = expected future reward)

Possible Actions n

Possible States m

0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0



5. Update Q-Table

- Go to new state x_{t+1}
- Find the highest possible Q-value of state x_{t+1} (Q-value = expected future reward)
- Update Q-value for the original state-action pair:

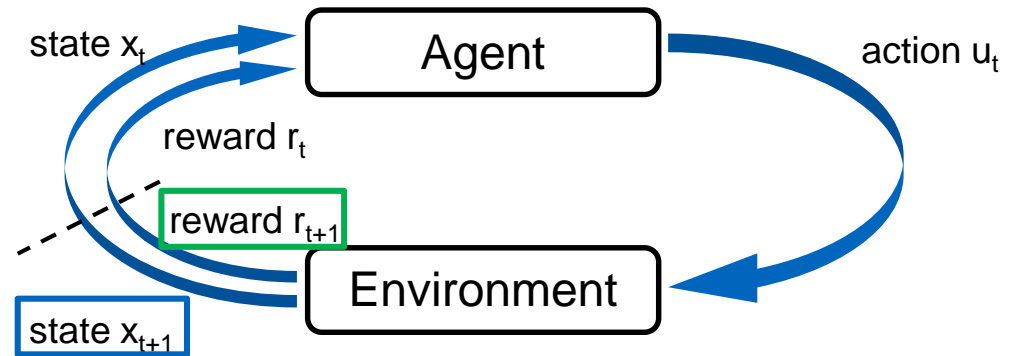
$$Q_{k+1}^{\pi}(x_t, u_t) = (1 - \alpha) \cdot Q_k^{\pi}(x_t, u_t) + \alpha \cdot (r_{t+1} + \gamma \cdot \operatorname{argmax}_u Q(x_{t+1}, u))$$

$r_{t+1} = 20, \quad \alpha = 10^{-3}, \quad \gamma = 0.95$

Possible Actions n

Possible States m

0	0	0	0
0	0	0.02	0
0	0	0	0
0	0	0	0

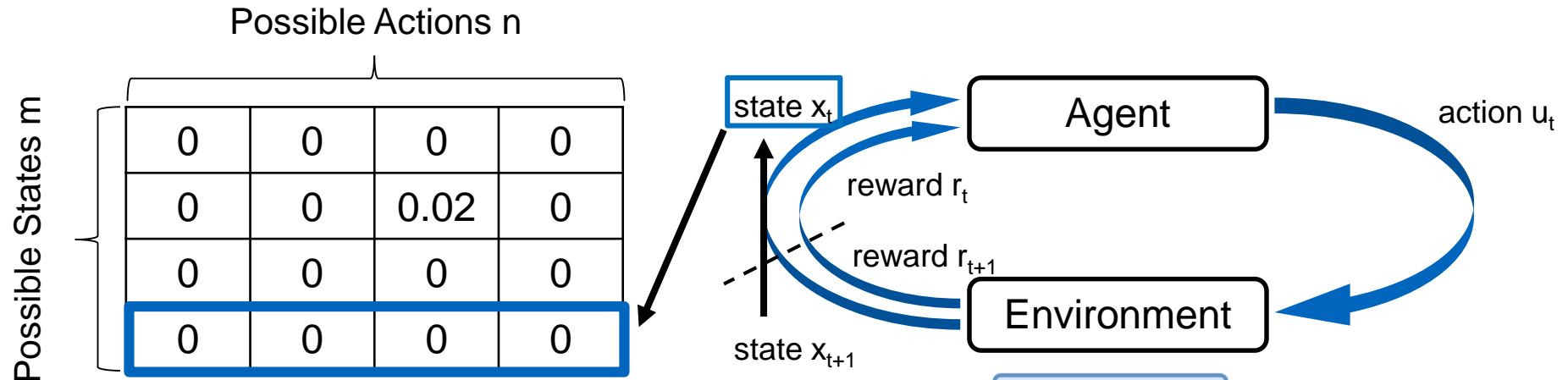


5. Update Q-Table

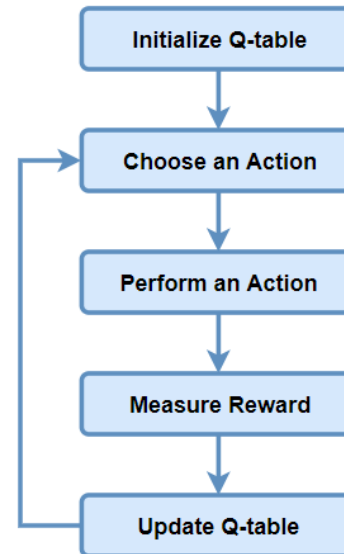
- Go to new state x_{t+1}
- Find the highest possible Q-value of state x_{t+1} (Q-value = expected future reward)
- Update Q-value for the original state-action pair:

$$Q_{k+1}^{\pi}(x_t, u_t) = (1 - \alpha) \cdot Q_k^{\pi}(x_t, u_t) + \alpha \cdot (r_{t+1} + \gamma \cdot \underset{u}{\operatorname{argmax}} Q(x_{t+1}, u)) = 0.02$$

$r_{t+1} = 20, \quad \alpha = 10^{-3}, \quad \gamma = 0.95$



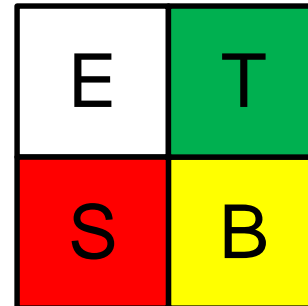
Go back to step 2 and continue



- 1 Q-Learning Process Steps
- 2 Exercise: Grid World
- 3 Exercise: Crane-Simulation

- Small Grid World: 2x2
- Go from a starting point (S) to a target location (T)
- Possible actions: up, down, left or right (crashing into the boundary → position remains the same, but counts as new step for reward calculation)
- Rewards:

x_{t+1}	$r_{t+1}(x_{t+1})$
E	-1
T	+100
B	-0.5
S	-1

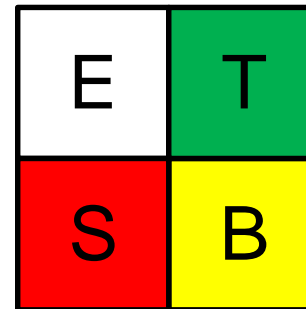


- Task 1:
 - Perform the first iteration of Q-Learning (Steps 1-5) by hand
 - Use learning rate $\alpha = 0.1$ and discount factor $\gamma = 1$
 - Neglect epsilon-greedy strategy

– Task 2:

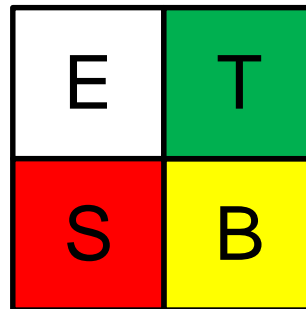
- After 5 completed episodes with epsilon-greedy strategy ($\epsilon = 0.05$) the Q-Table looks like shown below
- Compute the first iteration of Q-Learning for episode 6, starting again at (S) by hand
- As before: Neglect epsilon-greedy strategy

	Up	Down	Left	Right
E	0	-0.1	-0.1	0
T	0	0	0	0
S	-0.1	-0.1	-0.1	7.91
B	40.95	0	-0.1	0



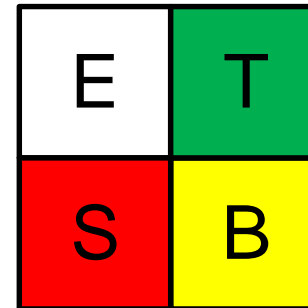
- Imagine an epsilon-greedy strategy where ϵ is computed as a function of the current episode:
→ $\epsilon(\text{episode}) = 0.9999^{\text{episode}}$
- Epsilon gets very small over time
- Task 3: What does the Q-Table look like after an infinite amount of episodes?

$$Q_{k+1}^{\pi}(x_t, u_t) = (1 - \alpha) \cdot Q_k^{\pi}(x_t, u_t) + \alpha \cdot \left(r_{t+1} + \gamma \cdot \underset{u}{\operatorname{argmax}} Q(x_{t+1}, u) \right)$$



- Task 4: What happens, if we apply the following reward structure to the Grid World problem?

x_{t+1}	$r(x_{t+1})$
E	-1
T	+100
B	+1
S	-1



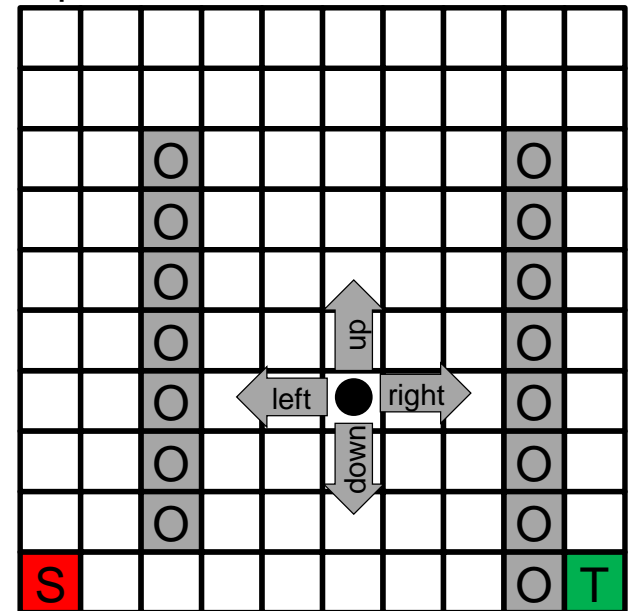
- 1 Q-Learning Process Steps**
- 2 Exercise: Grid World**
- 3 Exercise: Crane-Simulation**

- Quadratic Warehouse → discretized into a 10 x 10 grid
- A product has to be transported with a crane from a starting point (S) to a target location (T)
- When the crane crashes into an high rack storage (O), the game ends
- Possible actions: up, down, left or right (crashing into the boundary → position remains the same)
- Rewards:

x_{t+1}	$r_{t+1}(x_{t+1})$
White field / (S)	-1
(T)	+100
(O)	-100



Top view:



Tasks:

- Open `crane_simulation_Template.py`
- In class `CraneSim`: define the grid and all required variables
- In class `Q_Agent`: implement the epsilon-greedy policy and the update function for the Q-Table
- In function `play`: implement all necessary steps to run one episode of the simulation

