

# Reinforcement Learning I - How Agents Learn by Trial and Error

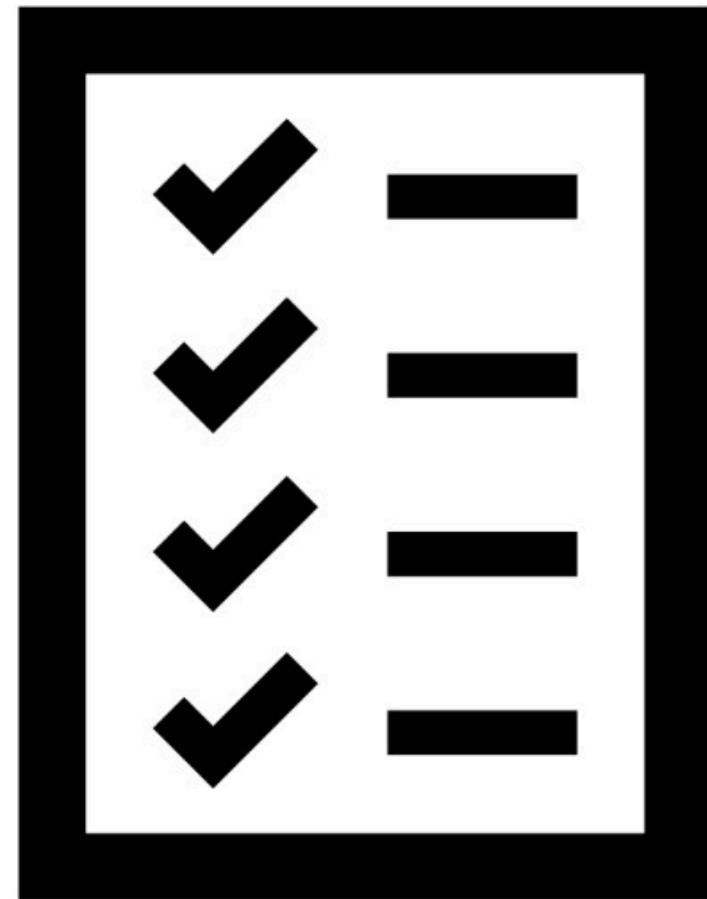
**Design IT.  
Create Knowledge.**

[www.hpi.de](http://www.hpi.de)



# Attendance List

[kisz@hpi.de](mailto:kisz@hpi.de)  
[hpi.de/kisz](http://hpi.de/kisz)





<https://hpi.de/ki-servicezentrum/>

## Overview

Gefördert durch:



Bundesministerium  
für Forschung, Technologie  
und Raumfahrt

# Four AI Service Centres in Germany



KI-Servicezentren

**Goal:**  
**Reduce barriers to the implementation  
of AI applications in society and the  
economy**

# AI Service Centre Berlin-Brandenburg

Structure

Research

Education

AISC

Infrastructure

Consulting

Gefördert durch:



Bundesministerium  
für Forschung, Technologie  
und Raumfahrt

## • AI methods research

## • AI operations research

Research

Gefördert durch:



Bundesministerium  
für Forschung, Technologie  
und Raumfahrt

### PubMedCLIP: How Much Does CLIP Benefit Visual Question Answering in the Medical Domain?

Sedigheh Eslami, Christoph Meinel, Gerard de Melo

Hasso Plattner Institute / University of Potsdam

{sedigheh.eslami, christoph.meinel, gerard.demelo}@hpi.de

#### Abstract

Contrastive Language-Image Pre-training (CLIP) has shown remarkable success in learning with cross modal supervision from extensive amounts of image-text pairs collected online. Thus far, the effectiveness of CLIP has been investigated primarily in general-domain multimodal problems. In this work, we evaluate the effectiveness of CLIP for the task of Medical Visual Question Answering (MedVQA). We present PubMedCLIP, a fine-tuned version of CLIP for the medical domain based on PubMed articles. Our experiments conducted on two MedVQA benchmark datasets illustrate that PubMedCLIP achieves superior results improving the overall accuracy up to 3% in comparison to the state-of-the-art. Model Agnostic Meta-Learning (MAML) networks pre-trained only on visual data. The PubMedCLIP model with different back-ends, the source code for pre-training them and reproducing our MedVQA pipeline is publicly available at <https://github.com/sarathESL/PubMedCLIP>.

#### 1 Introduction

Medical visual question answering (MedVQA) seeks answers to natural language questions about a given medical image. The development of MedVQA has considerable potential to benefit health care systems, as it may aid clinicians in interpreting medical images and obtaining more accurate diagnoses by consulting a second opinion. Thus, it has become a very active area of research, with competitive benchmarks and yearly competitions (Abacha et al., 2021). Yet, visual question answering in the medical domain in particular remains non-trivial as we suffer from a general lack of large balanced training data, in part due to privacy concerns. To solve the multimodal task of MedVQA, a system must understand both medical images and textual questions and infer the associations between them sufficiently well to produce a correct answer (An-

Findings of the Association for Computational  
May 2-6, 2023 ©2023 Associa

### Exploring Paracrawl for Document-level Neural Machine Translation

Yusser Al Ghassis<sup>1,2</sup>, Jingyi Zhang<sup>1</sup> and Josef van Genabith<sup>1,2</sup>

<sup>1</sup>German Research Center for Artificial Intelligence (DFKI),  
Saarland Informatics Campus, Saarbrücken, Germany

<sup>2</sup>Department of Language Science and Technology, Saarland University, Germany

<sup>1</sup>Hasso-Plattner-Institut (HPI), Potsdam, Germany  
yusser.al.ghassis@dfki.de, Jingyi.Zhang@hpi.de

#### Abstract

Document-level neural machine translation (NMT) has outperformed sentence-level NMT on a number of datasets. However, document-level NMT is still not widely adopted in real-world translation systems mainly due to the lack of large-scale general-domain training data for document-level NMT. We examine the effectiveness of using Paracrawl for learning document-level translation. Paracrawl is a large-scale parallel corpus crawled from the Internet and contains data from various domains. The official Paracrawl corpus was released as parallel sentences (extracted from parallel web-pages) and therefore previous works only used Paracrawl for learning sentence-level translation. In this work, we extract parallel paragraphs from Paracrawl parallel webpages using automatic sentence alignments and we use the extracted parallel paragraphs as parallel documents for training document-level translation models. We show that document-level NMT models trained with only parallel paragraphs from Paracrawl can be used to translate real documents from TED, News and Europarl, outperforming sentence-level NMT models. We also perform a targeted pronoun evaluation and show that document-level models trained with Paracrawl data can help context-aware pronoun translation. We release our data and code here<sup>1</sup>.

#### 1 Introduction

The Transformer translation model (Vaswani et al., 2017), which performs sentence-level translation based on attention networks, has achieved great success and significantly improved the state-of-the-art in machine translation. Compared to sentence-level translation, document-level translation (Xu et al., 2021; Bao et al., 2021; Jauegi Uusane et al., 2020; Mu et al., 2020; Mansf et al., 2019; Tu et al., 2018; Mansf and Haffari, 2018) performs translation at document-level and can potentially fur-

<sup>1</sup><https://github.com/YusserW/Exploring-Paracrawl-for-Document-level-Neural-Machine-Translation>

Proceedings of the 17th Conference of the European Chapter of  
May 2-6, 2023 ©2023 Associa

### Efficient Parallelization Layouts for Large-Scale Distributed Model Training

Johannes Hagemann

Aleph Alpha / Hasso Plattner Institute

johannes.hagemann@student.hpi.de

Samuel Weinhardt

Aleph Alpha

samuel.weinhardt@aleph-alpha.com

Konstantin Döbler

Hasso Plattner Institute

konstantin.doebler@hpi.de

Maximilian Schall

Hasso Plattner Institute

maximilian.schall@hpi.de

Gerard de Melo

Hasso Plattner Institute

gerard.demelo@hpi.de

#### Abstract

Efficiently training large language models requires parallelizing across hundreds of hardware accelerators and invoking various compute and memory optimizations. When combined, many of these strategies have complex interactions regarding the final training efficiency. Prior work tackling this problem did not have access to the latest set of optimizations, such as FLASHATTENTION or sequence parallelism. In this work, we conduct a comprehensive ablation study of possible training configurations for large language models. We distill this large study into several key recommendations for the most efficient training. For instance, we find that using a micro-batch size of 1 usually enables the most efficient training layouts. Larger micro-batch sizes necessitate activation checkpointing or higher degrees of model parallelism and also lead to larger pipeline bubbles. Our most efficient configurations enable us to achieve state-of-the-art training efficiency results over a range of model sizes, most notably a Model FLOPs utilization of 70.5% when training a LLaMA 13B model.

#### 1 Introduction

The number of parameters and computational resources spent on training deep neural networks is growing rapidly [1, 3, 14]. The largest models consisting of hundreds of billions of parameters do not even fit onto a single hardware accelerator. Thus, training these models requires various ways of reducing the memory requirements, such as ZeRO [16], activation checkpointing [2], and 3D-parallel data, tensor, and pipeline parallel training [13]. 3D parallelism, in particular, has been demonstrated to be effective for the training of Transformer based large language models (LLMs) with hundreds of billions of parameters [13].

However, training these models efficiently with 3D parallelism requires significant domain expertise and extensive manual effort to determine the ideal configurations. These configurations not only need to combine data, model, and pipeline parallelism most efficiently, but also consider complex interactions with other memory and compute optimizations. FLASHATTENTION [5] in particular has had a notable impact since its release, enabling us to train models at previously impossible degrees of training efficiency. In light of these developments, we conduct a systematic study via a large-scale training efficiency sweep of these interactions. We consider up to 256 GPUs and LLaMA models with up to 65 billion parameters.

Workshop on Advancing Neural Network Training at 37th Conference on Neural Information Processing Systems (NeurIPS 2023).



Zugangsangefrage  
[aisc.hpi.de](http://aisc.hpi.de)

# Infrastructure

Gefördert durch:



Bundesministerium  
für Forschung, Technologie  
und Raumfahrt



- **Free access**
- No production operation
  - Data should be **anonymised** or **synthesised**
  - No **hosting** of products
- **Reporting & publication** by users
- **Old rights** remain with users
- **New rights** remain with users
  - Granting of rights of use for research and teaching

## Training

- 64 NVIDIA H100 GPU

## Inference

- 40 NVIDIA A30 GPU

## ARM Server

- Ampere Altra Max  
M128-30 CPU
- 2 x NVIDIA L40 GPUs

## GPU Server

- AMD Epyc CPU
- 8 x NVIDIA L40S GPU

## Edge

- ARMv8 CPU
- NVIDIA Jetson AGX Module

## Neuromorph

- 288 SpiNNaker2 Chips

## Storage

- 1.5 PB NVMe

## Network

- 400 Gb/s Infiniband
- 200 Gb/s Ethernet



Newsletter

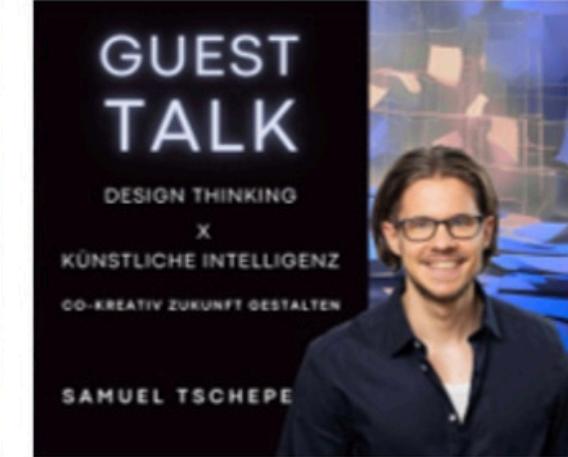
## Education

### Talks



[tele-task.de/series/1463](http://tele-task.de/series/1463)

- Guest talks about research and innovation

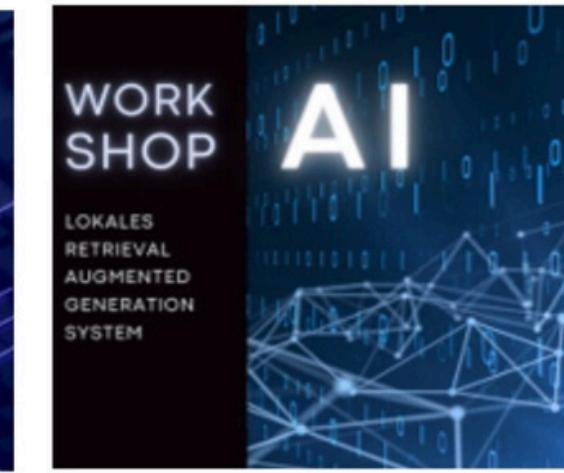


### Work shops



[aimaker.community](http://aimaker.community)

- Practical topics
- Example topics: Speech2summary, Docker for ML, semantic search

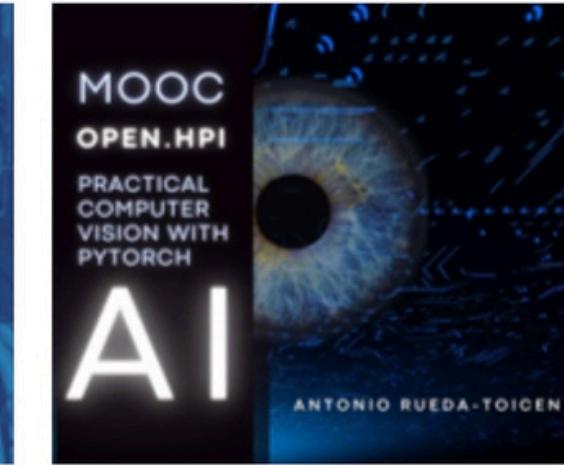
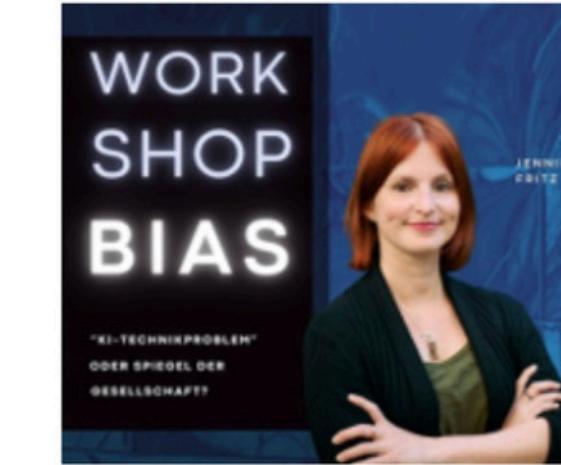


### MOOCs



[open.hpi.de/channels/ai-service-center](http://open.hpi.de/channels/ai-service-center)

- ChatGPT: What does generative AI mean for our society?
- Profitable AI
- Understanding and avoiding AI biases



Gefördert durch:



Bundesministerium  
für Forschung, Technologie  
und Raumfahrt



Sprechstunde buchen

# Consulting

Gefördert durch:



Bundesministerium  
für Forschung, Technologie  
und Raumfahrt

## AI consultation hours

- Answering questions:
  - about AI infrastructure
  - about AI models & frameworks
  - about AI use cases

## AI pilot projects

- Co-development of a prototype
- Application every three months
- Selection criteria, e.g. AI maturity, public benefit
- Publication of results

## Cooperation

- Jointly organized network meetings

Kontast: Normal Kontast: Invertiert Kontast: Graustufen A A A

Ki-Prototyp Leichte-Sprache-Assistent (LeiSA)

Vereinfachte Texte mit LLMs! [Zum Projekt auf GitHub](#)

Originaltext

Das KI-Servicezentrum Berlin-Brandenburg ist ein Projekt am Hasso-Piattner-Institut. Unser Ziel ist, den Einsatz von künstlicher Intelligenz im Wirtschaft und Gesellschaft zu fördern.

Unser Schwerpunkt liegt auf der Unterstützung von Start-ups, kleinen und mittelständischen Unternehmen, öffentlichen Einrichtungen, Bildungseinrichtungen und der Forschung. Wir bieten Interne Beratungsleistungen, Bildungsangebote und Infrastruktur. Unser Anwendungsbereich reicht von der Produktion bis zur Vermarktung.

Leichte Sprache

„Unser Ziel“  
Wir möchten, dass viele Menschen und Unternehmen KI-Anwendungen nutzen können.

„Wer kann sich am ehesten wenden?“  
Startende  
Kleine und mittelständische Unternehmen  
Start-ups  
Öffentliche Einrichtungen

„Was können wir tun?“  
Kostenlose Beratung  
Bildungsangebote  
Infrastruktur, um KI-Anwendungen zu entwickeln

„Beispiele für unsere Arbeit“  
Ein Unternehmen möchte eine App entwickeln, die Menschen bei der Einkaufsliste hilft.  
Eine Forschungsgruppe möchte ein System bauen, das Menschen bei der Krankenversorgung unterstützt.

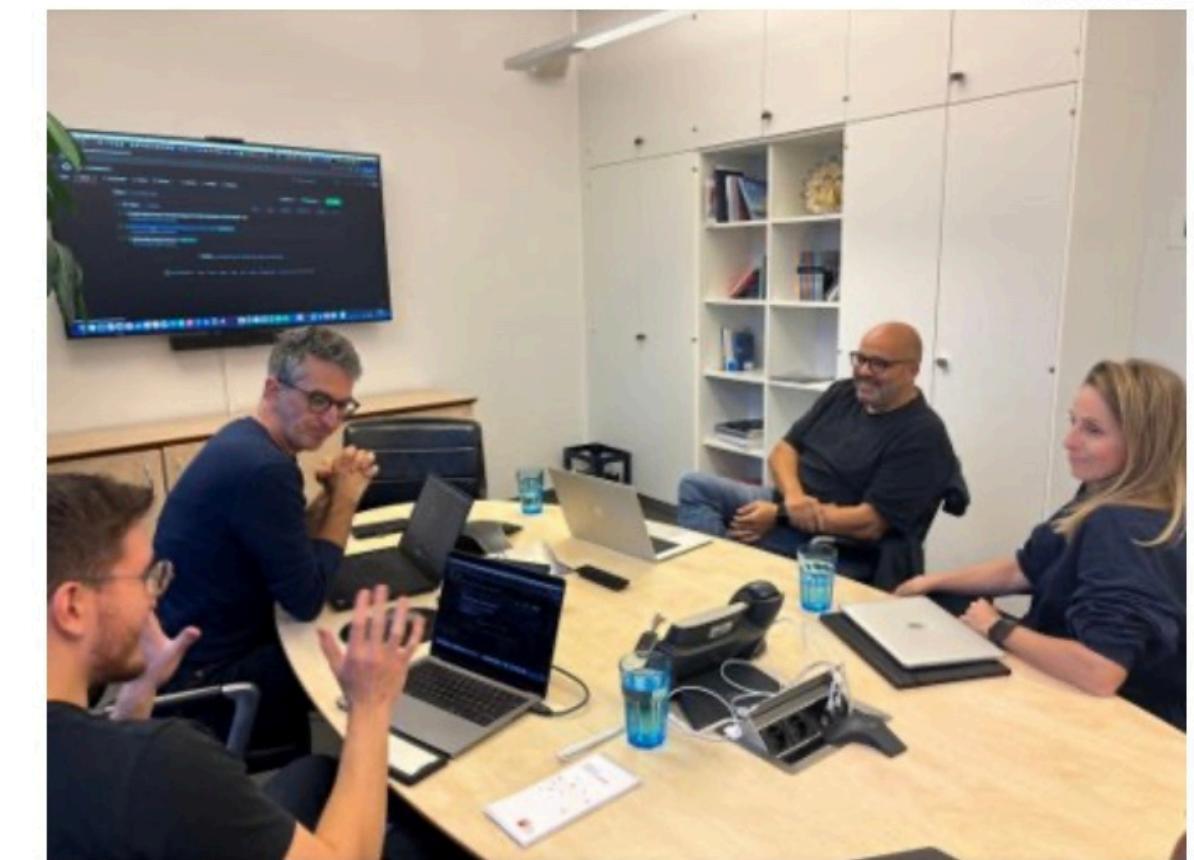
Wir helfen Ihnen gerne!

Hinweis zur Übersetzung: Dieser Text wurde mit einer KI erstellt.  
Bitte verarbeiten Sie keine personenbezogenen Daten mit dem Tool. Falls der Prototyp einmal nicht verfügbar sein sollte, kontaktieren Sie uns gern, wir kümmern uns darum.

[Vereinfachen!](#) [Löschen](#)

Der KI-Prototyp ist ein gemeinsames Projekt der Digitalagentur Brandenburg GmbH und des KI-Servicezentrum Berlin-Brandenburg.

[github.com/aihpi/leichte-sprache](https://github.com/aihpi/leichte-sprache)



## Previous AI pilot projects

- Generation of math problems
- Generation of easy language
- Generation of upcycling suggestions
- Reduction of food waste
- Dating by handwriting



[Jetzt bewerben!](#)



<https://tts.aisc.hpi.de/>

user: event

pw: aisc@hpi2022!

## Prototypes

Gefördert durch:



Bundesministerium  
für Forschung, Technologie  
und Raumfahrt

# Speech processing

[Processing →](#)    [← Scenarios](#)

## 🎙 Voice Dubbing Demonstrator

### 🎙 Select Reference Audio

Choose input method:

Record Audio

Use Predefined Audio

Click the button below to start recording your audio.

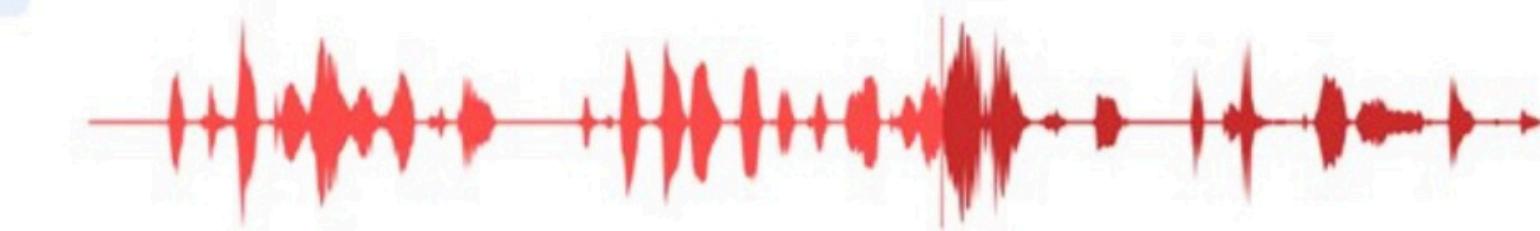


## 🎙 Voice Dubbing Demonstrator

### 🔍 At an Event

### 🤖 Cloned Voice (3/3)

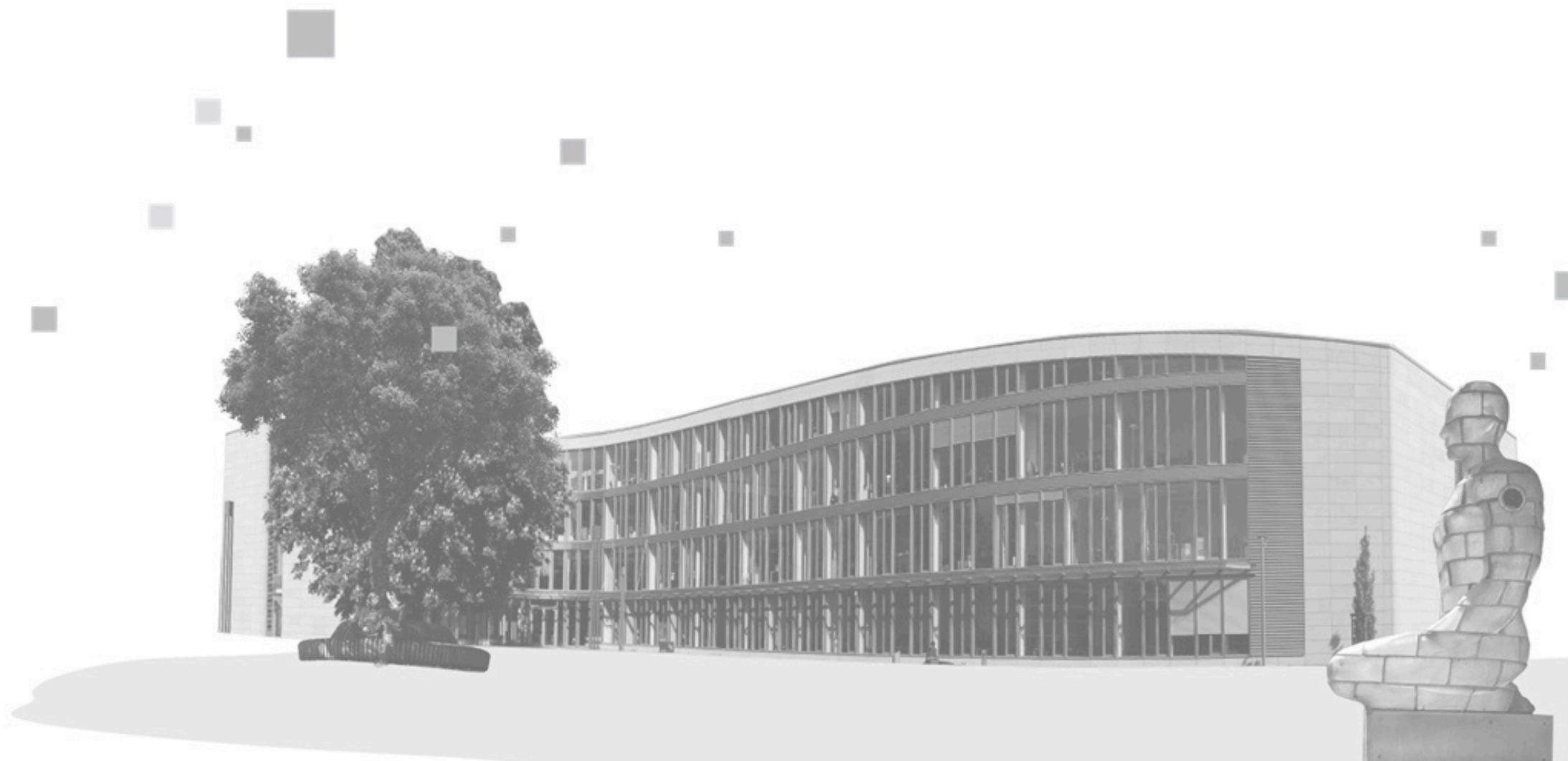
That's exactly why we're here today, discussing the importance of making sure AI systems are trustworthy and responsible.



# Reinforcement Learning I - How Agents Learn by Trial and Error

**Design IT.  
Create Knowledge.**

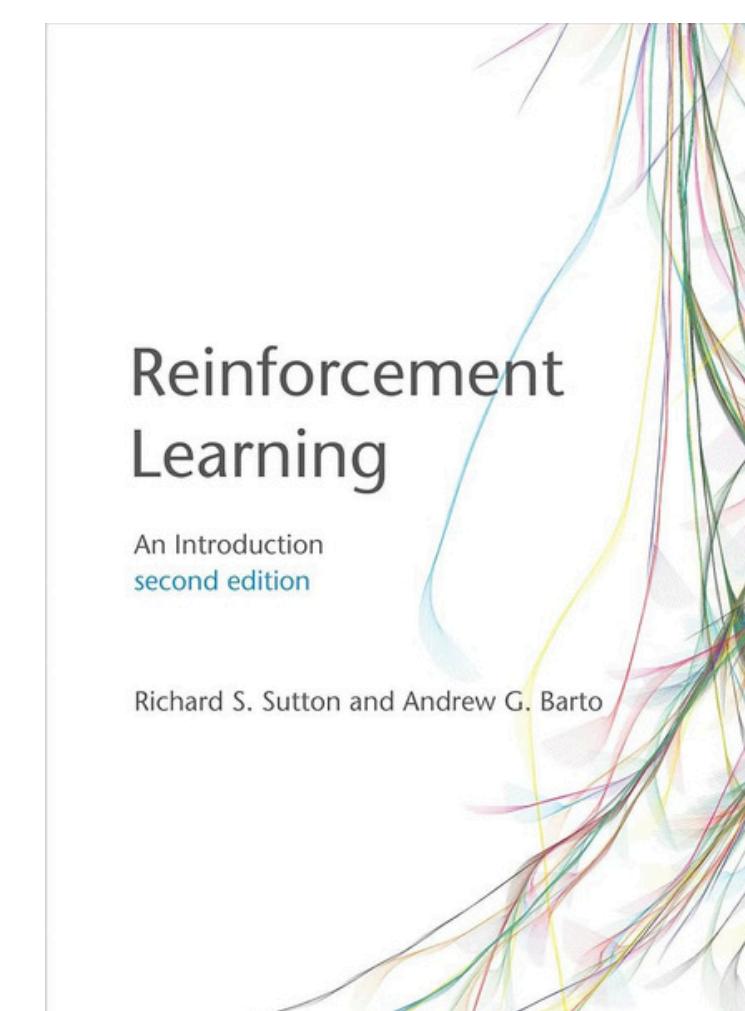
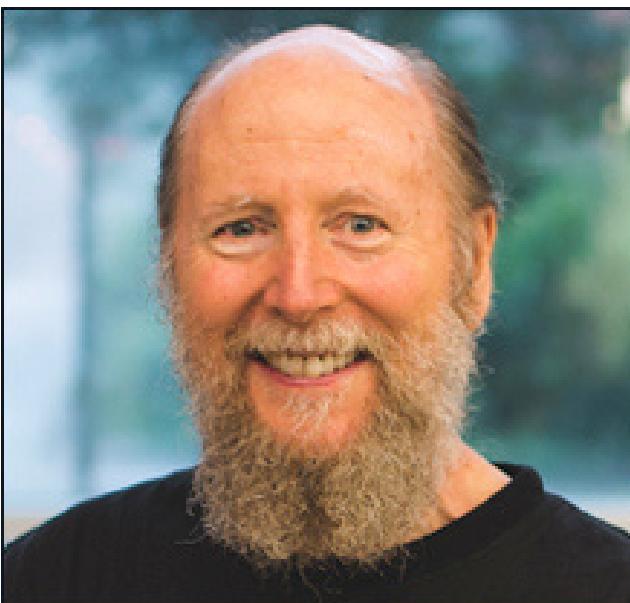
[www.hpi.de](http://www.hpi.de)



# Motivation and context I

“Of all the forms of machine learning, reinforcement learning is the closest to the kind of learning that humans and other animals do.” - Richard Sutton, Andrew Barto

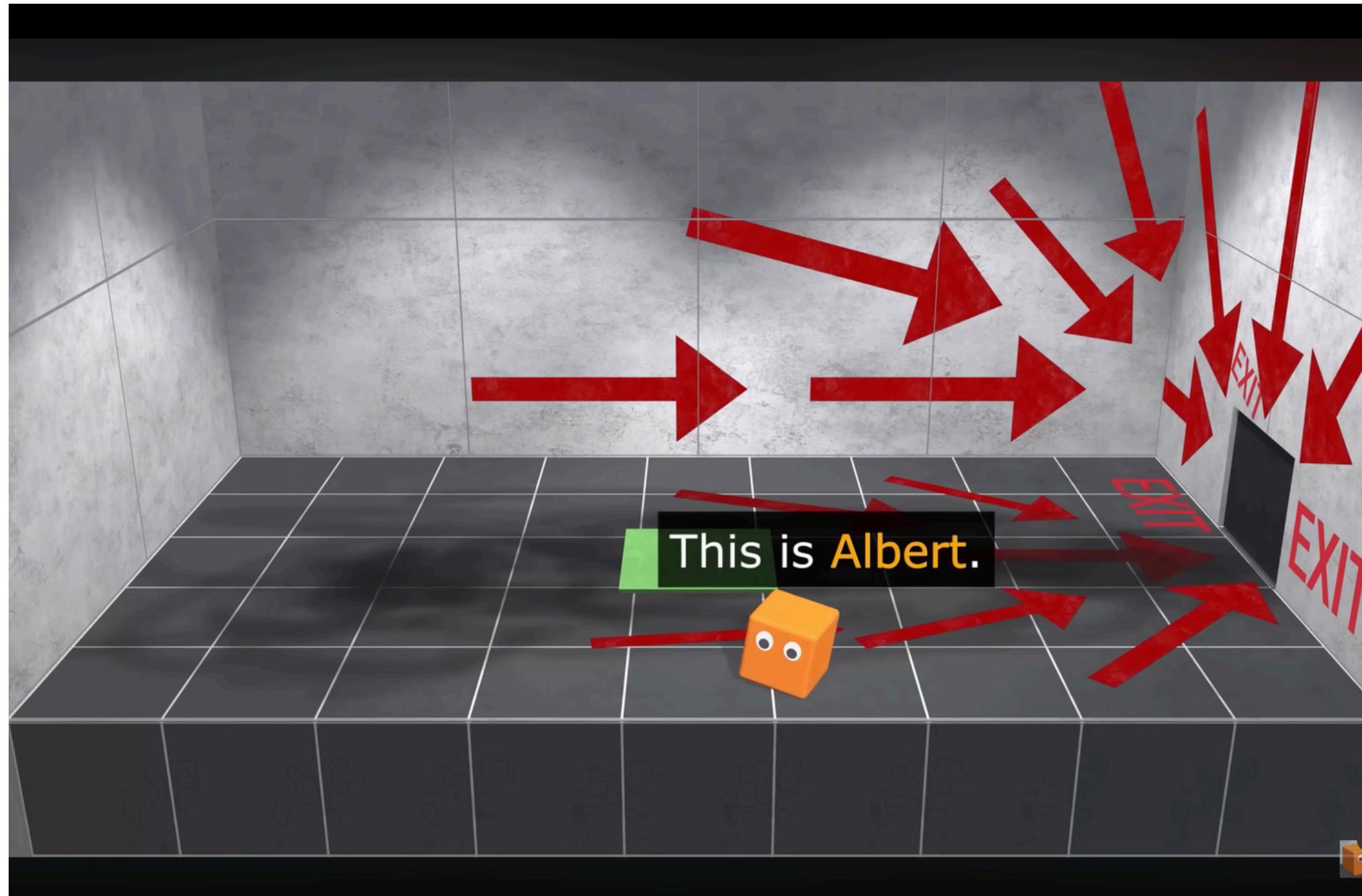
[kisz@hpi.de](mailto:kisz@hpi.de)  
[hpi.de/kisz](http://hpi.de/kisz)



# Motivation and context I

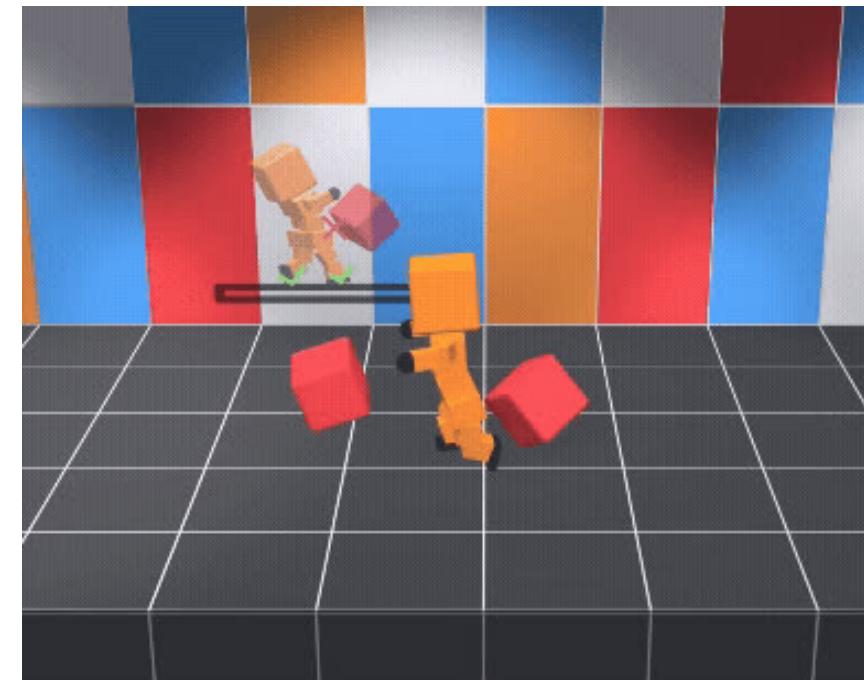
“Of all the forms of machine learning, reinforcement learning is the closest to the kind of learning that humans and other animals do.” -Richard S. Sutton

[kisz@hpi.de](mailto:kisz@hpi.de)  
[hpi.de/kisz](http://hpi.de/kisz)



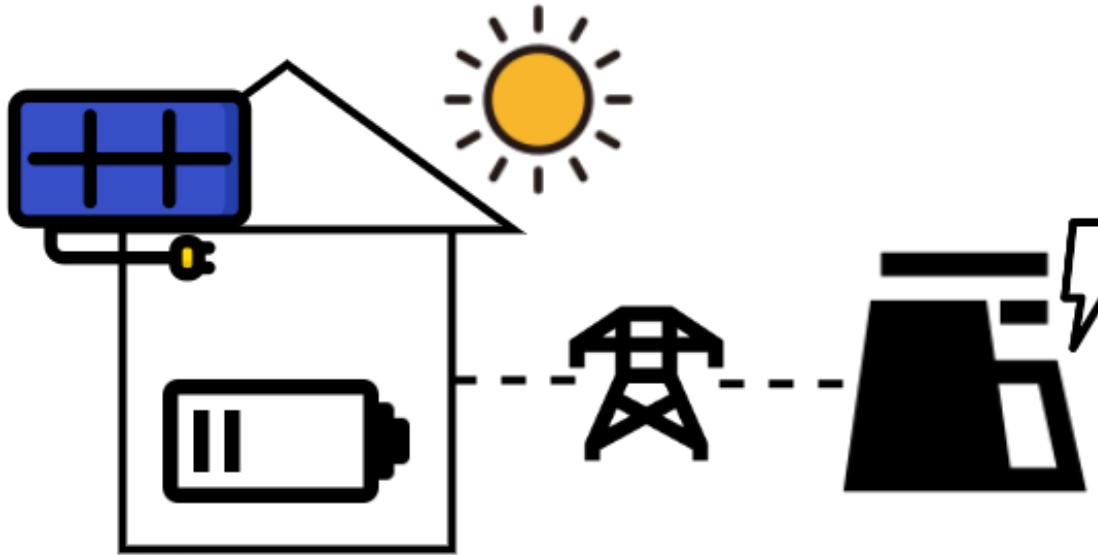
# Motivation and context II: Applications

## Robotics



[kisz@hpi.de](mailto:kisz@hpi.de)  
[hpi.de/kisz](http://hpi.de/kisz)

## Control Energy Systems



<https://github.com/tobirohrer/building-energy-storage-simulation>

## RLHF for Chatbots

What is the biggest animal in the world?

Wähle die beste Antwort

User-Auswahl	Modell-Präferenz	Beide
Score: 0.007 Tiger. seed: 11	Score: 0.016 The biggest animal in the world is a blue whale weighing up to 16,000 pounds (7,300 kg). seed: 22	Score: 0.007 Tiger. seed: 33

Keine Antwort ist gut

Danke für dein Feedback!

# Motivation and context III

[kisz@hpi.de](mailto:kisz@hpi.de)  
[hpi.de/kisz](http://hpi.de/kisz)

## Supervised Learning

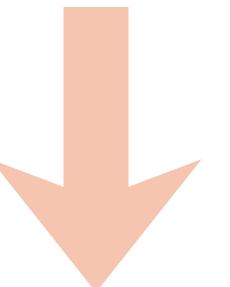
- Models are **trained on labeled data** to predict outputs

## Reinforcement Learning

- Learning optimal behavior by **interacting with an environment**

## Unsupervised Learning

- Learn to **identify previously unknown patterns** and relationships in unlabeled data



used when **optimal actions are too complex or unknown to define** but can be learned through experience

# What will we cover today?

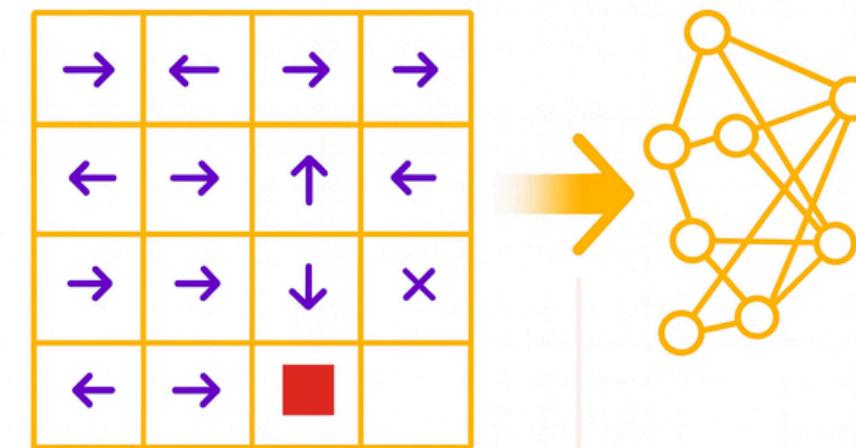
[kisz@hpi.de](mailto:kisz@hpi.de)  
[hpi.de/kisz](http://hpi.de/kisz)

## Today

- Basic concepts
- "Hello World" of RL:  
Q-learning on FrozenLake environment
- Build intuition about learning parameters
- First Coding

## Following workshops

- Deep RL
- When to apply which algorithm
- How to set up a RL problem
- Real-world applications



# What will we cover today?

[kisz@hpi.de](mailto:kisz@hpi.de)  
[hpi.de/kisz](http://hpi.de/kisz)

**RL Lab**  
Interactive Reinforcement Learning Visualization

### Configuration

Environment  
FrozenLake-v1-NoSlip

Algorithm  
Q-Learning

**START TRAINING**

**STOP PLAYBACK**

Learning Parameters

Num Episodes **100**

Training episodes. Must be an integer.

Exploration Rate **0.1**

0 ≤ ε ≤ 1 - probability of random action

Learning Rate **0.1**

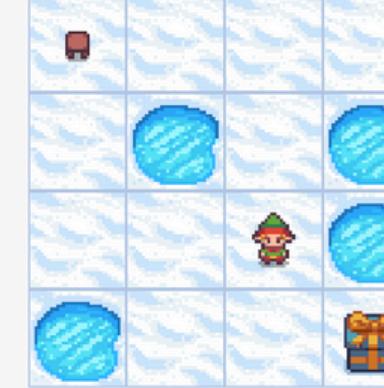
0 < α ≤ 1 - controls how much new information overrides old

Discount Factor **0.95**

0 ≤ γ < 1 - importance of future rewards

### Environment

Playing Policy



About FrozenLake (No Slip)

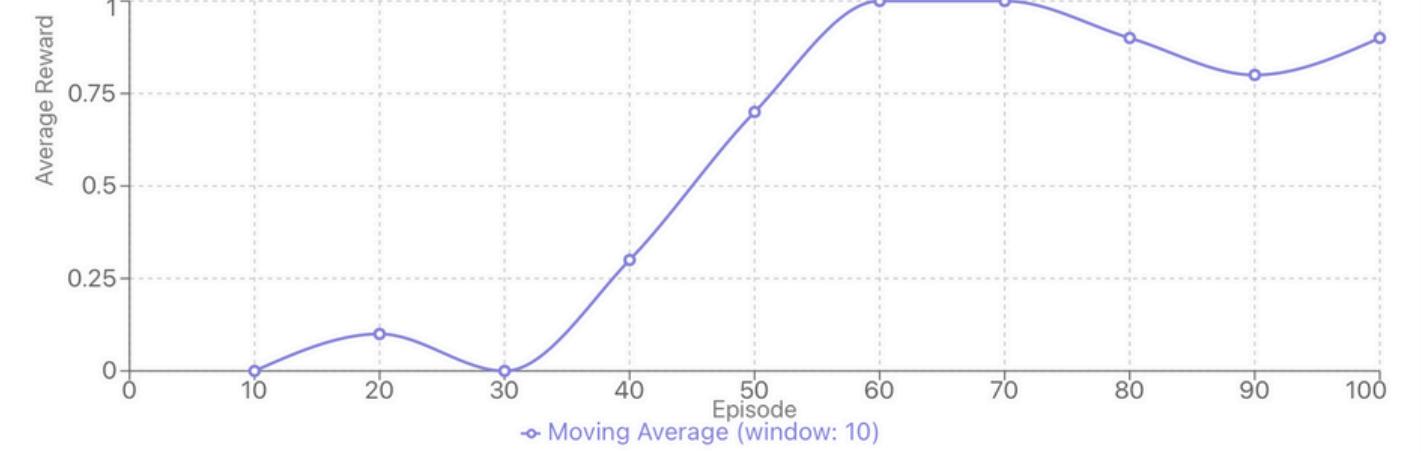
About Q-Learning

### Moving Average Reward

Episodes Trained: **100**

Current Average: **0.900**

Best Average: **1.000**



Average Reward

Episode

Moving Average (window: 10)

Min Q: **0.000**

Max Q: **0.998**

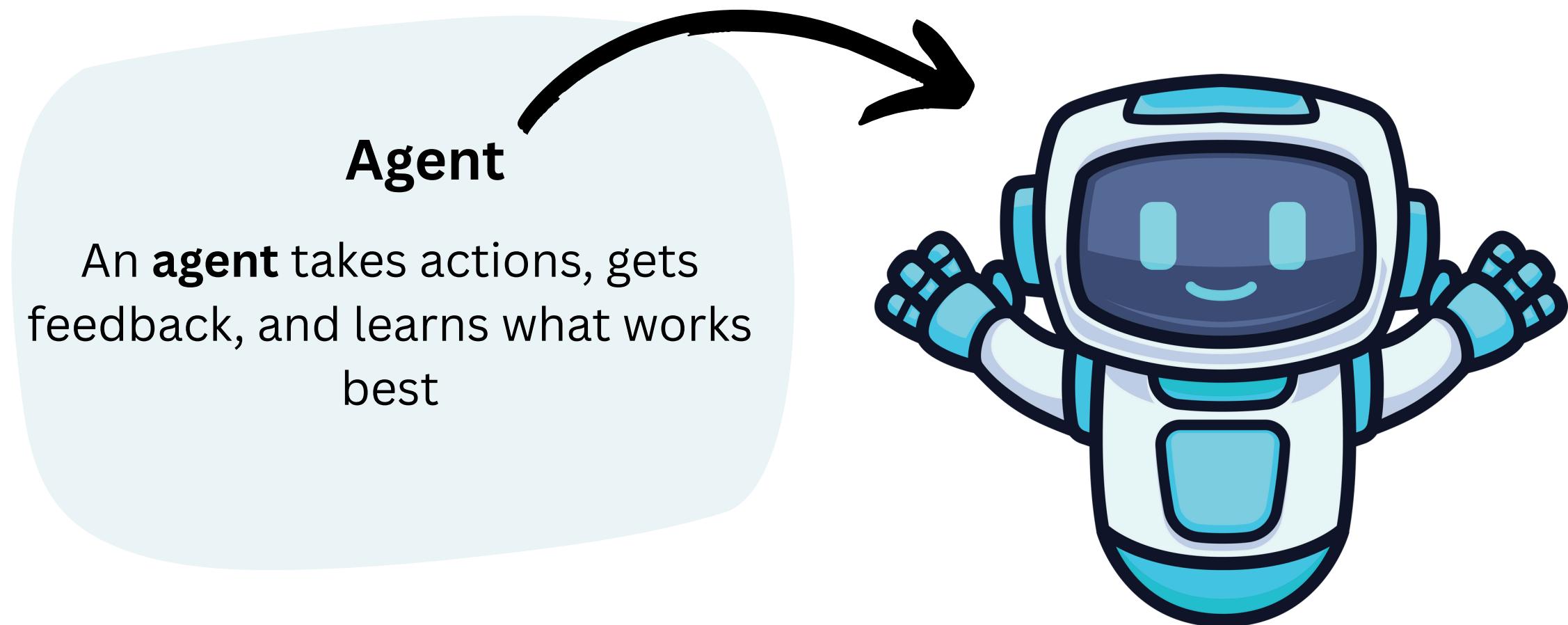
Avg Q: **0.085**

0	0.02	0.44	0.00	0.03	0.00	0.06	0.00
1	0.02	0.00	0.59	0.00	0.74	0.00	0.00
2	0.00	0.00	0.00	0.00	0.00	0.11	0.00
3	0.00	0.00	0.00	0.00	0.00	0.00	0.00
4	0.00	0.00	0.00	0.00	0.00	0.00	0.00
5	0.00	0.00	0.00	0.00	0.00	0.00	0.00
6	0.00	0.00	0.00	0.00	0.00	0.00	0.00
7	0.00	0.00	0.00	0.00	0.00	0.00	0.00

Q-Table Heatmap

# Key terms

[kisz@hpi.de](mailto:kisz@hpi.de)  
[hpi.de/kisz](http://hpi.de/kisz)



# Key terms

[kisz@hpi.de](mailto:kisz@hpi.de)  
[hpi.de/kisz](http://hpi.de/kisz)

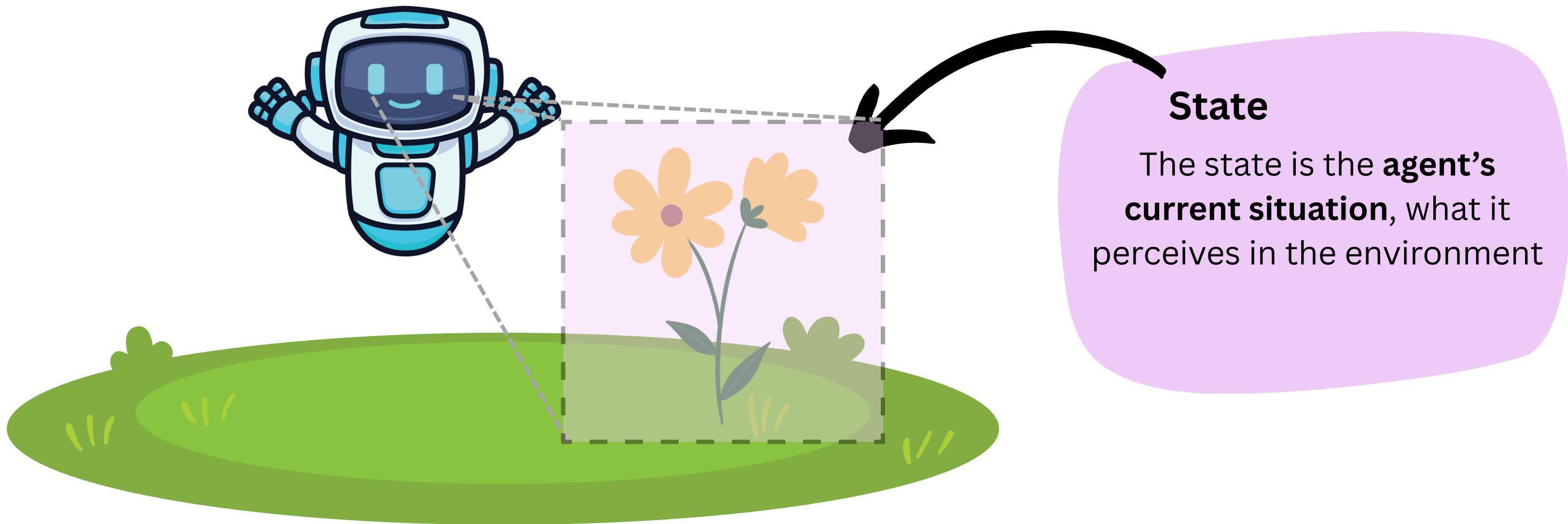
## Environment

The environment is **the world the agent interacts with**. It defines what happens when the agent takes an action.



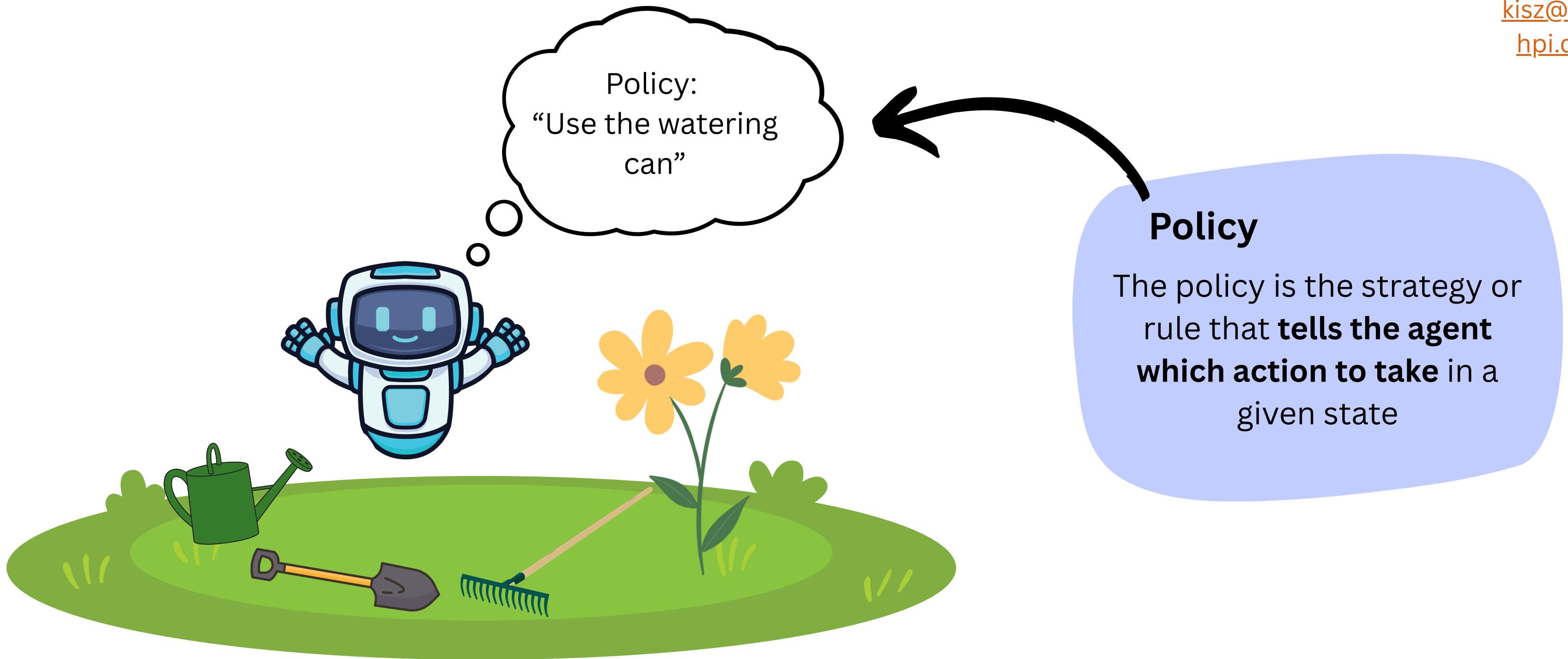
# Key terms

[kisz@hpi.de](mailto:kisz@hpi.de)  
[hpi.de/kisz](http://hpi.de/kisz)



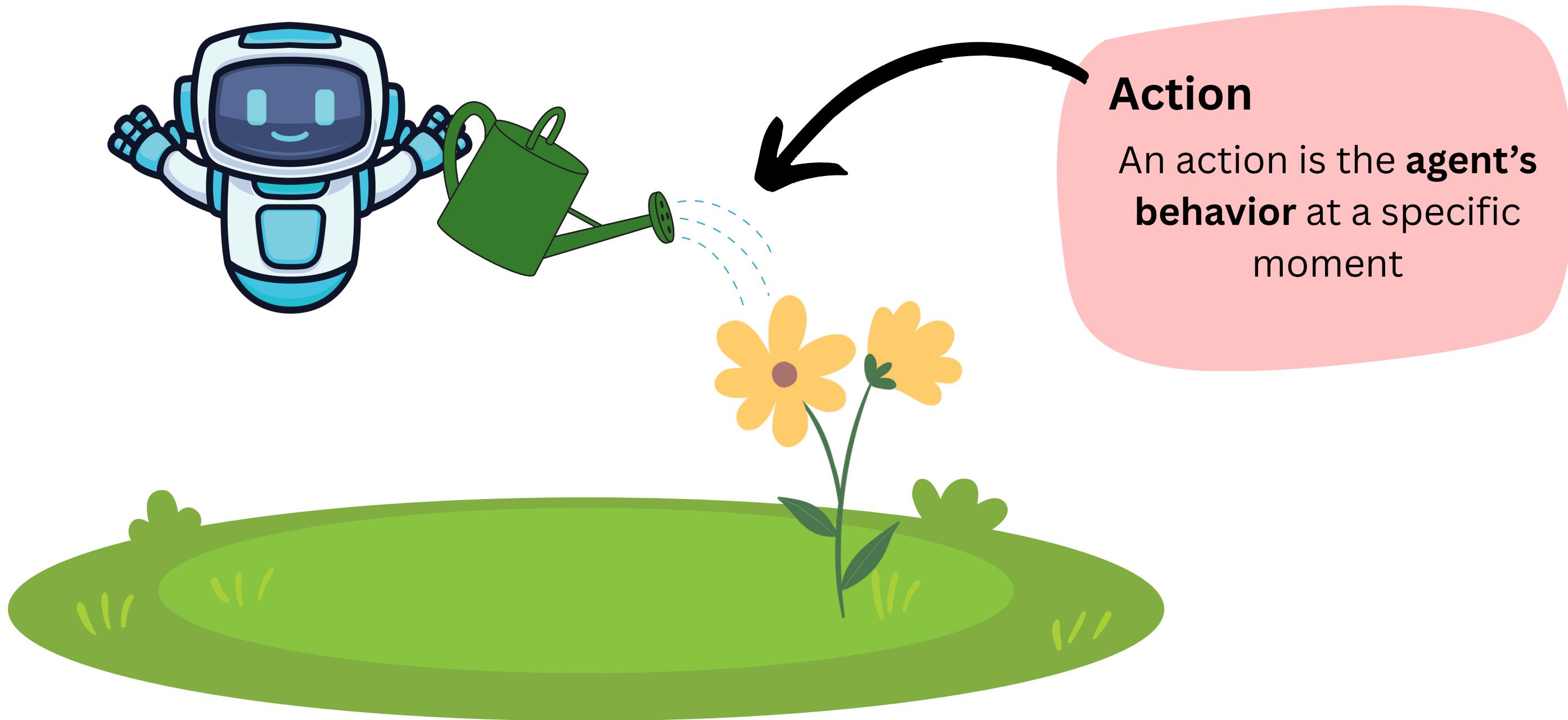
# Key terms

[kisz@hpi.de](mailto:kisz@hpi.de)  
[hpi.de/kisz](http://hpi.de/kisz)



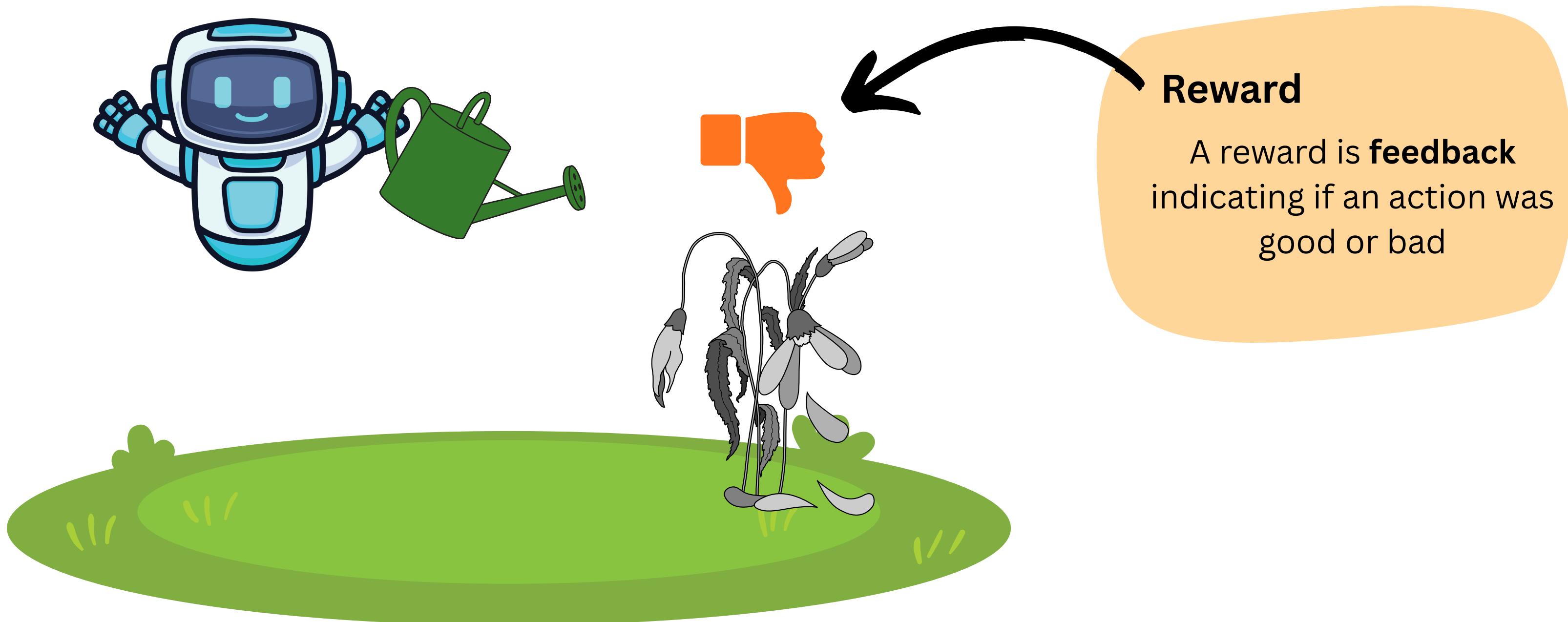
# Key terms

[kisz@hpi.de](mailto:kisz@hpi.de)  
[hpi.de/kisz](http://hpi.de/kisz)



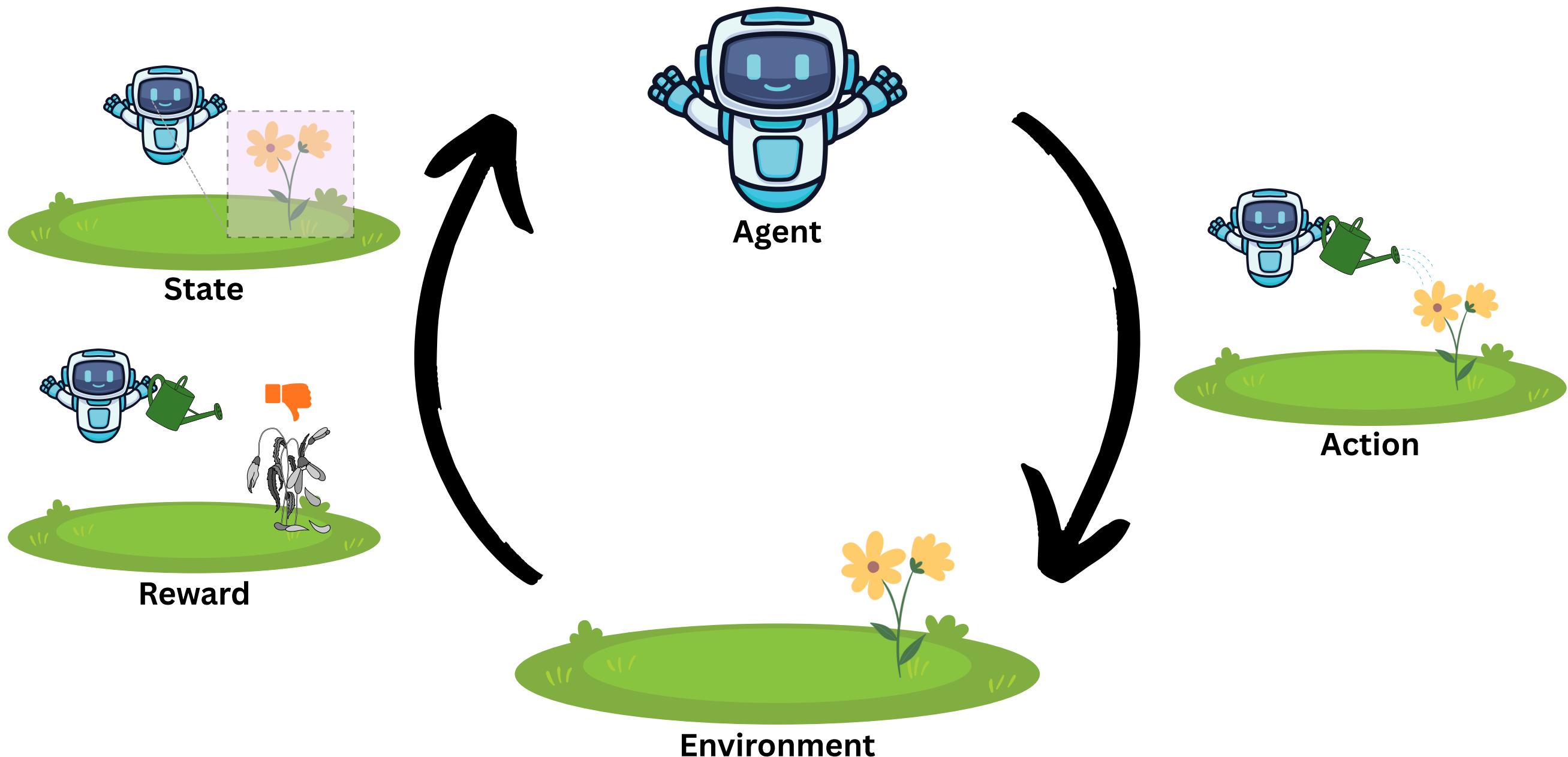
# Key terms

[kisz@hpi.de](mailto:kisz@hpi.de)  
[hpi.de/kisz](http://hpi.de/kisz)



# Key terms

[kisz@hpi.de](mailto:kisz@hpi.de)  
[hpi.de/kisz](http://hpi.de/kisz)



## Timestep

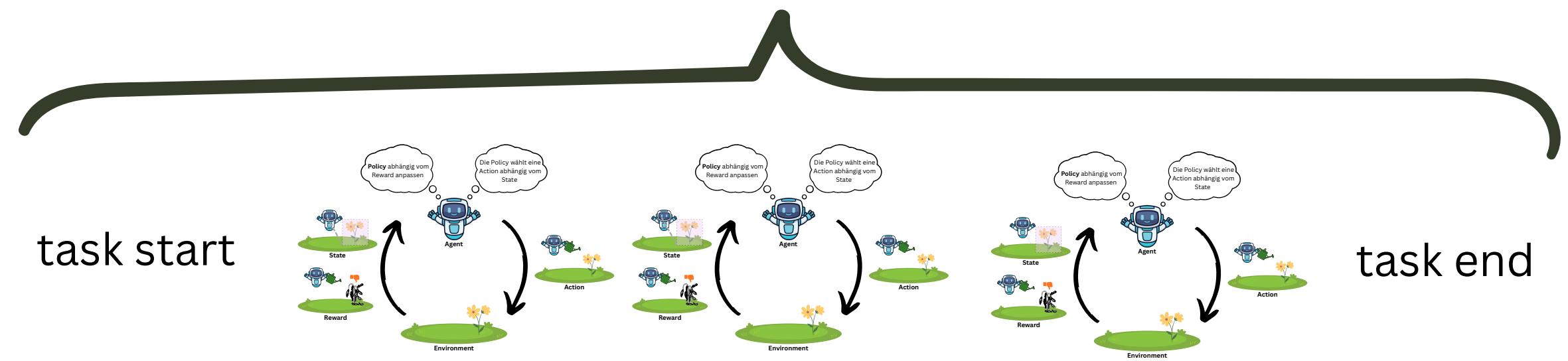
A timestep is one interaction with the environment

# Key terms

[kisz@hpi.de](mailto:kisz@hpi.de)  
[hpi.de/kisz](http://hpi.de/kisz)

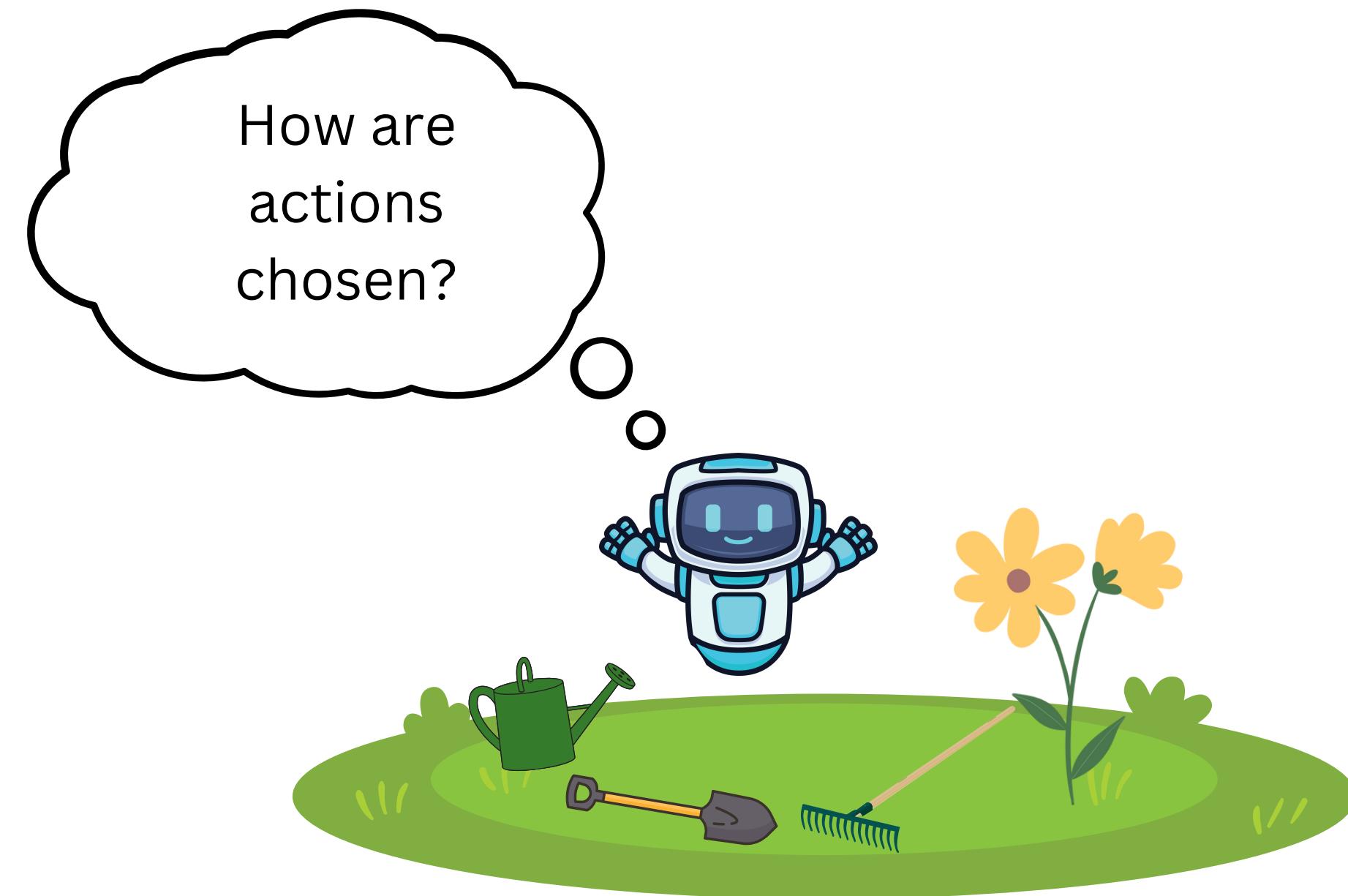
## Episode

An episode is a **finished sequence of interactions** forming one complete task run



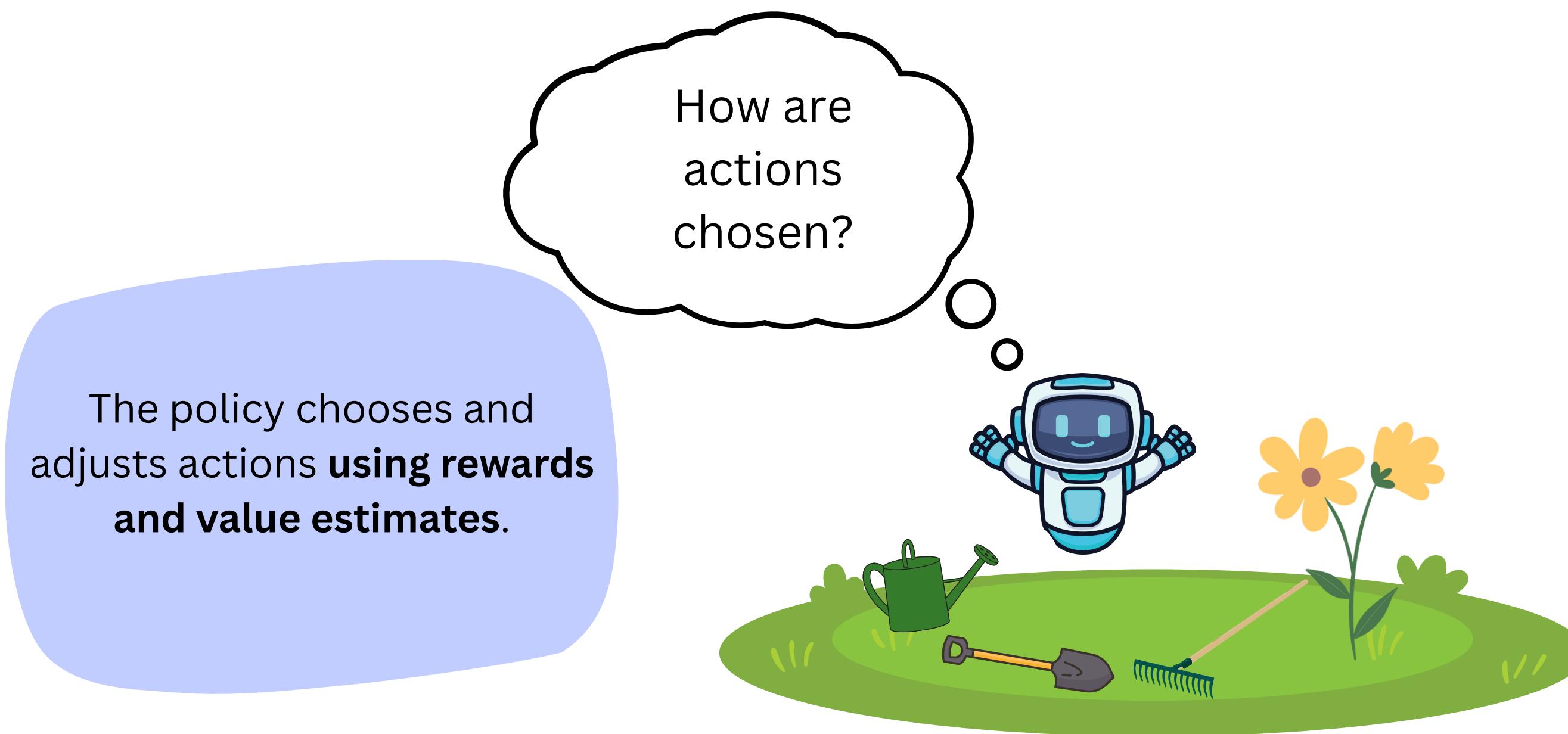
# How does the agent evaluate situations? I

[kisz@hpi.de](mailto:kisz@hpi.de)  
[hpi.de/kisz](http://hpi.de/kisz)



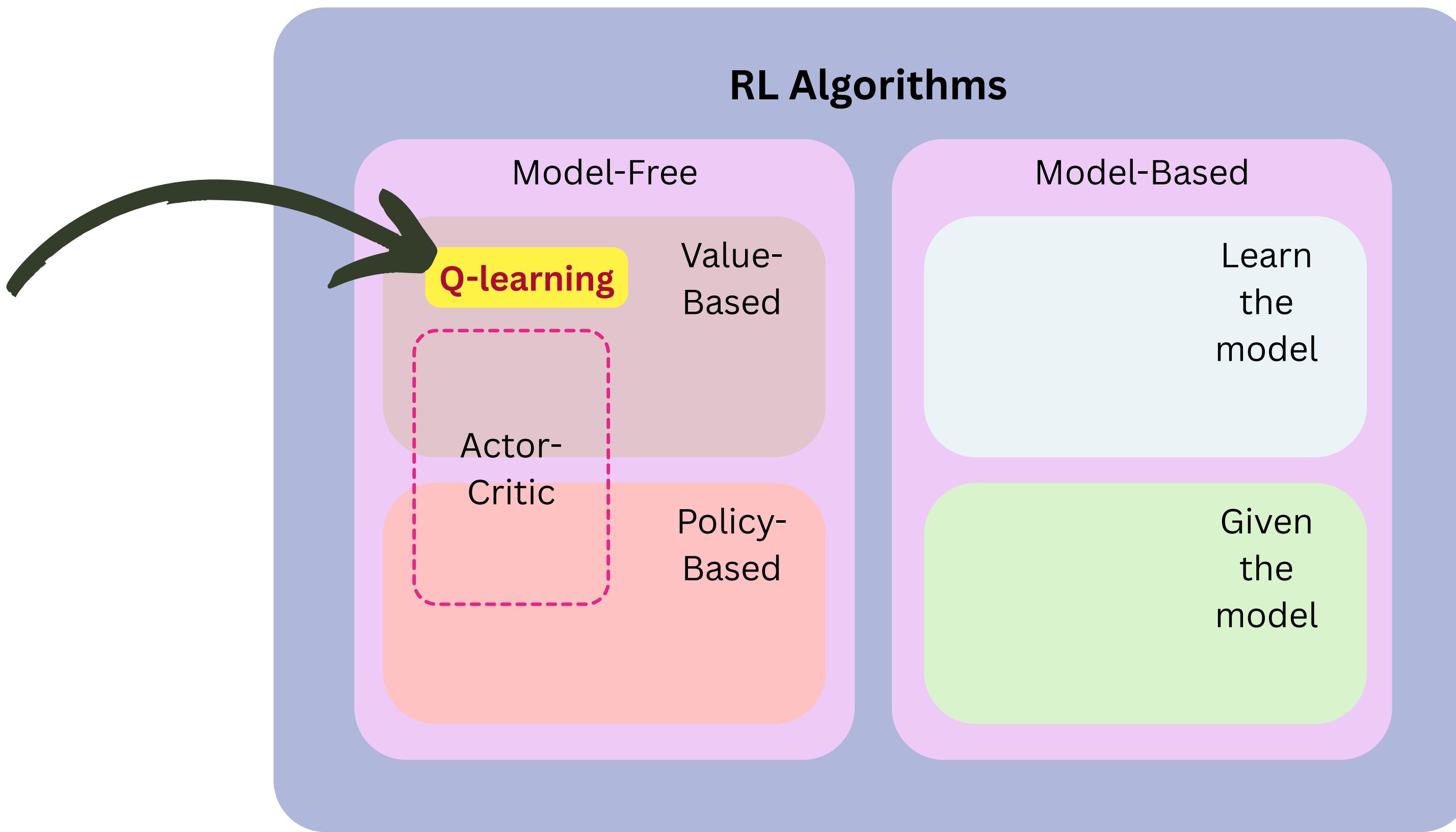
# How does the agent evaluate situations? II

[kisz@hpi.de](mailto:kisz@hpi.de)  
[hpi.de/kisz](http://hpi.de/kisz)



# Algorithm

[kisz@hpi.de](mailto:kisz@hpi.de)  
[hpi.de/kisz](http://hpi.de/kisz)

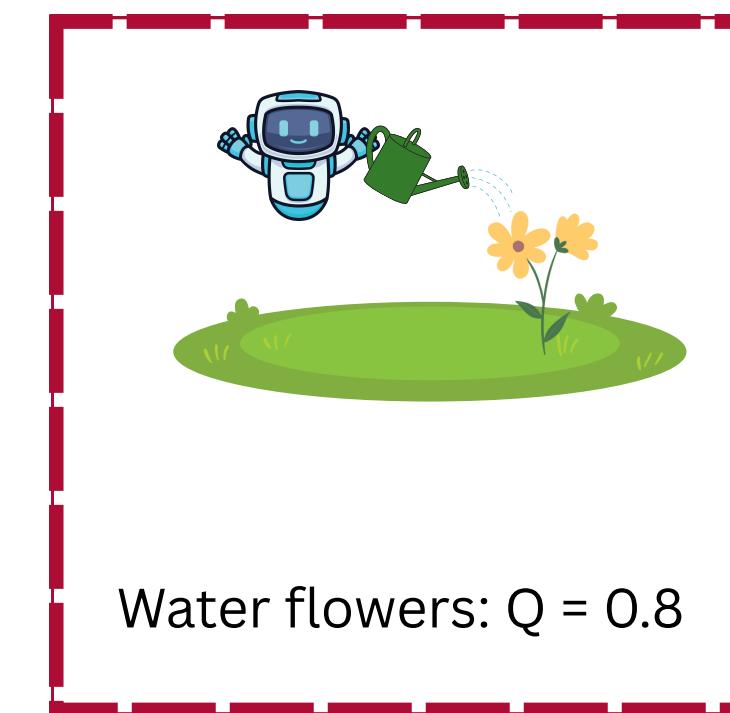
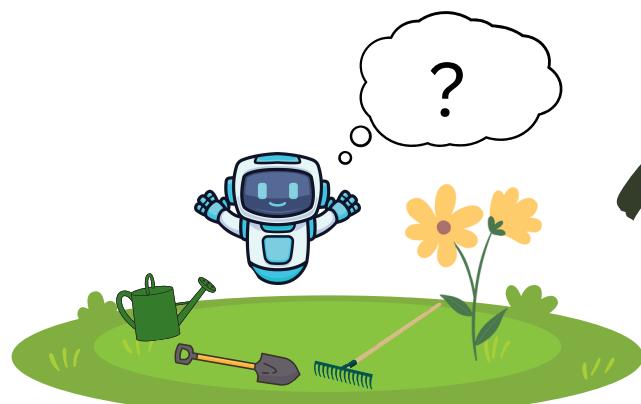


# What are Q-values?

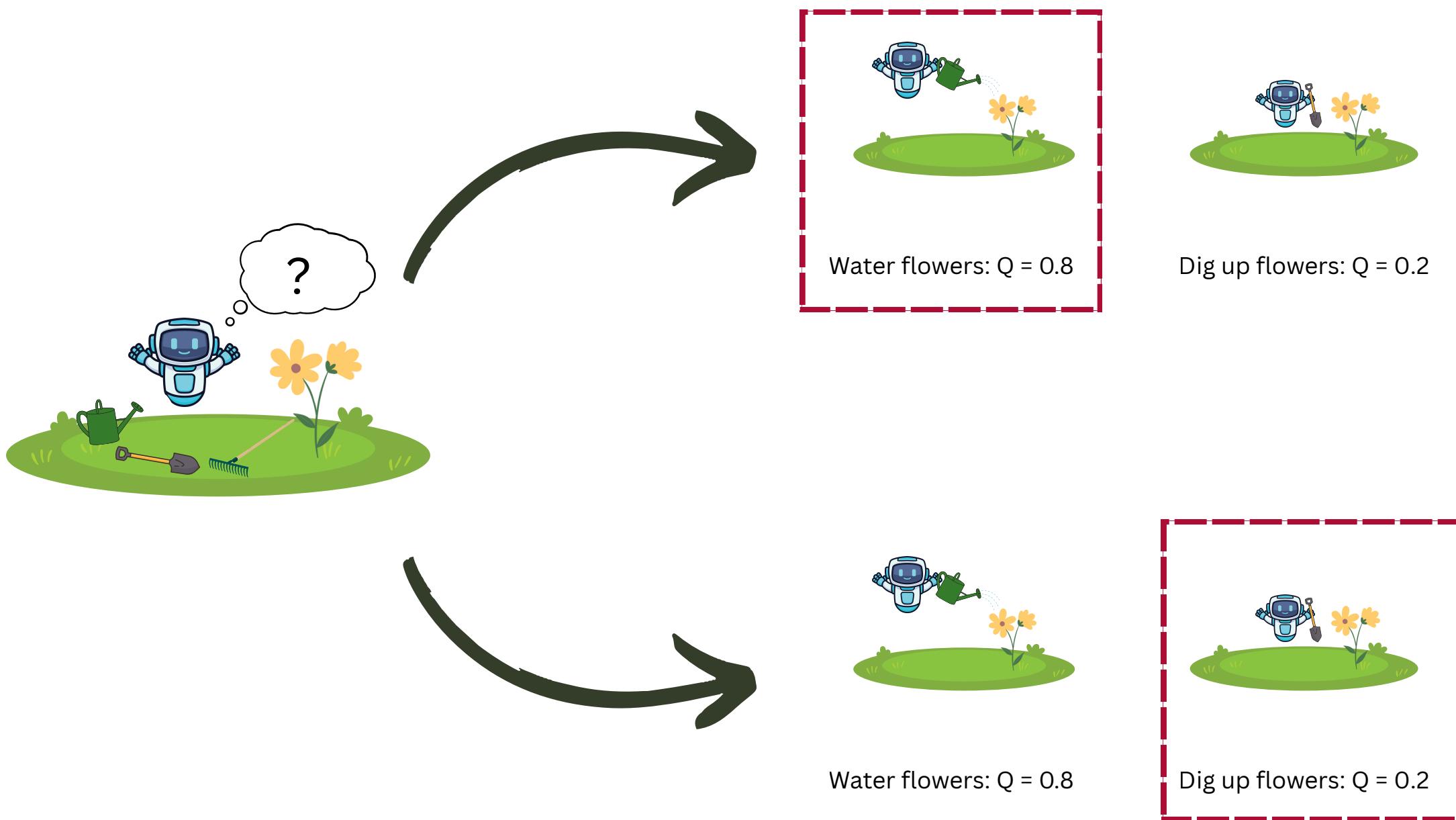
estimates **how good an action is in a given state**

used for **choosing the best action**

higher Q-values  
→ better actions



# Exploration vs. Exploitation



## Exploitation

Choosing the best-known action to get rewards

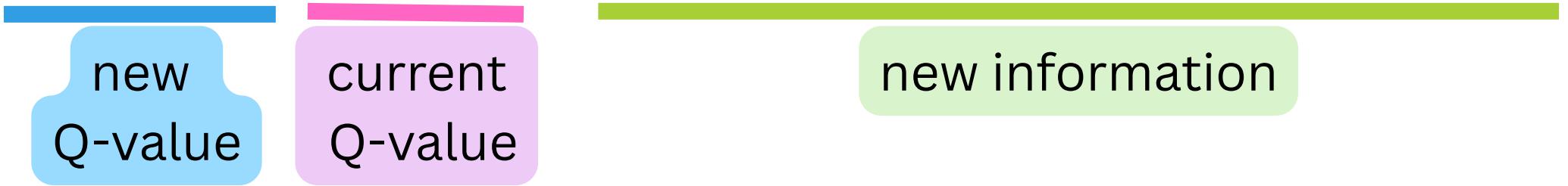
## Exploration

Trying new actions to discover better options

# How are Q-values learned?

[kisz@hpi.de](mailto:kisz@hpi.de)  
[hpi.de/kisz](http://hpi.de/kisz)

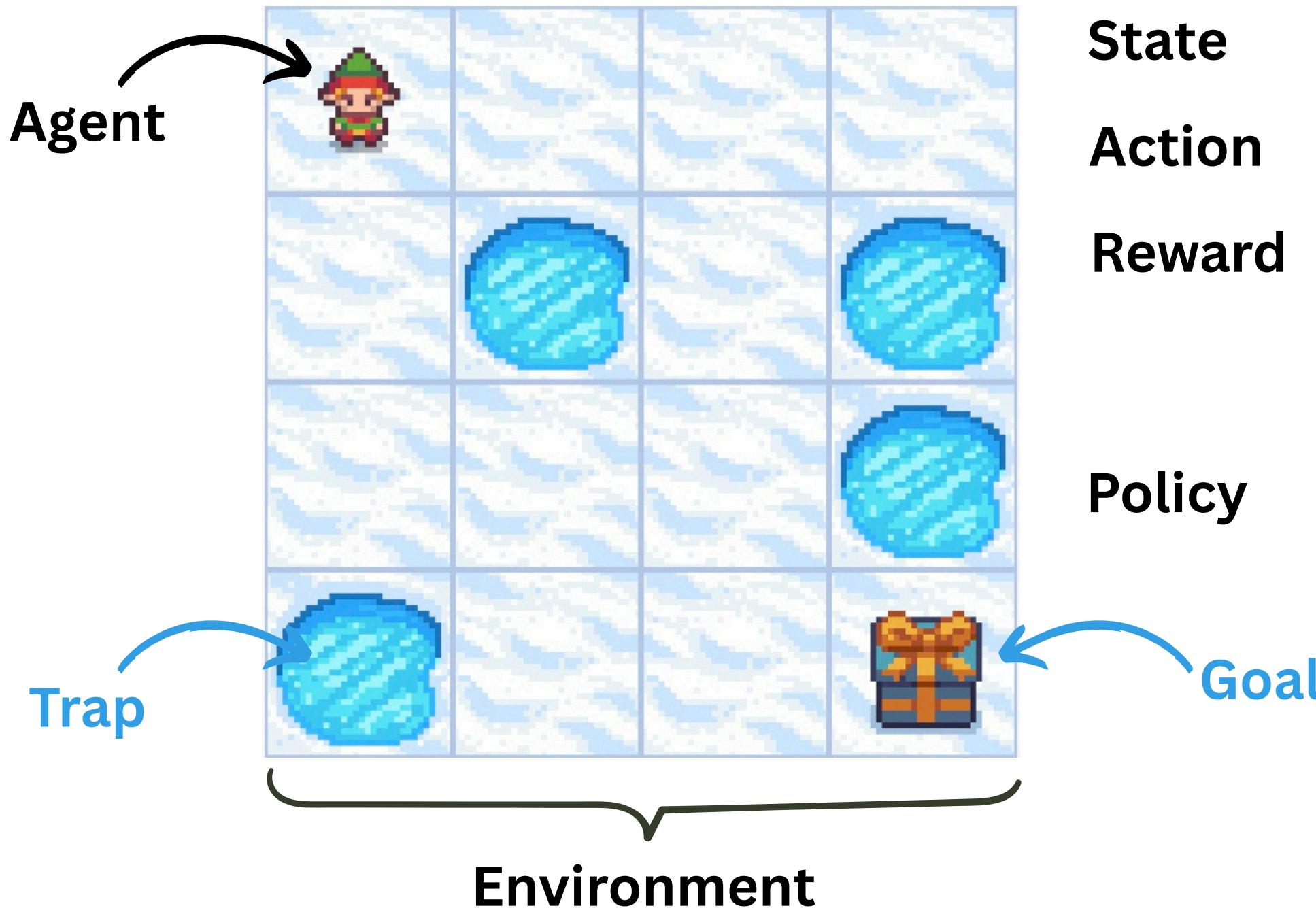
- Q-values are not known at the beginning

$$Q'(s, a) = Q(s, a) + \alpha \left[ r + \gamma \max_{a'} Q(s', a') - Q(s, a) \right]$$


Formula symbols	
$Q'(s, a)$	updated value
$Q(s, a)$	current value estimate
$\alpha$	learning rate
$r$	reward
$\gamma$	discount factor
$\max_{a'} Q(s', a')$	best next-state value
$s$	state
$a$	action
$s'$	next state

# Example: Frozen Lake

[kisz@hpi.de](mailto:kisz@hpi.de)  
[hpi.de/kisz](http://hpi.de/kisz)



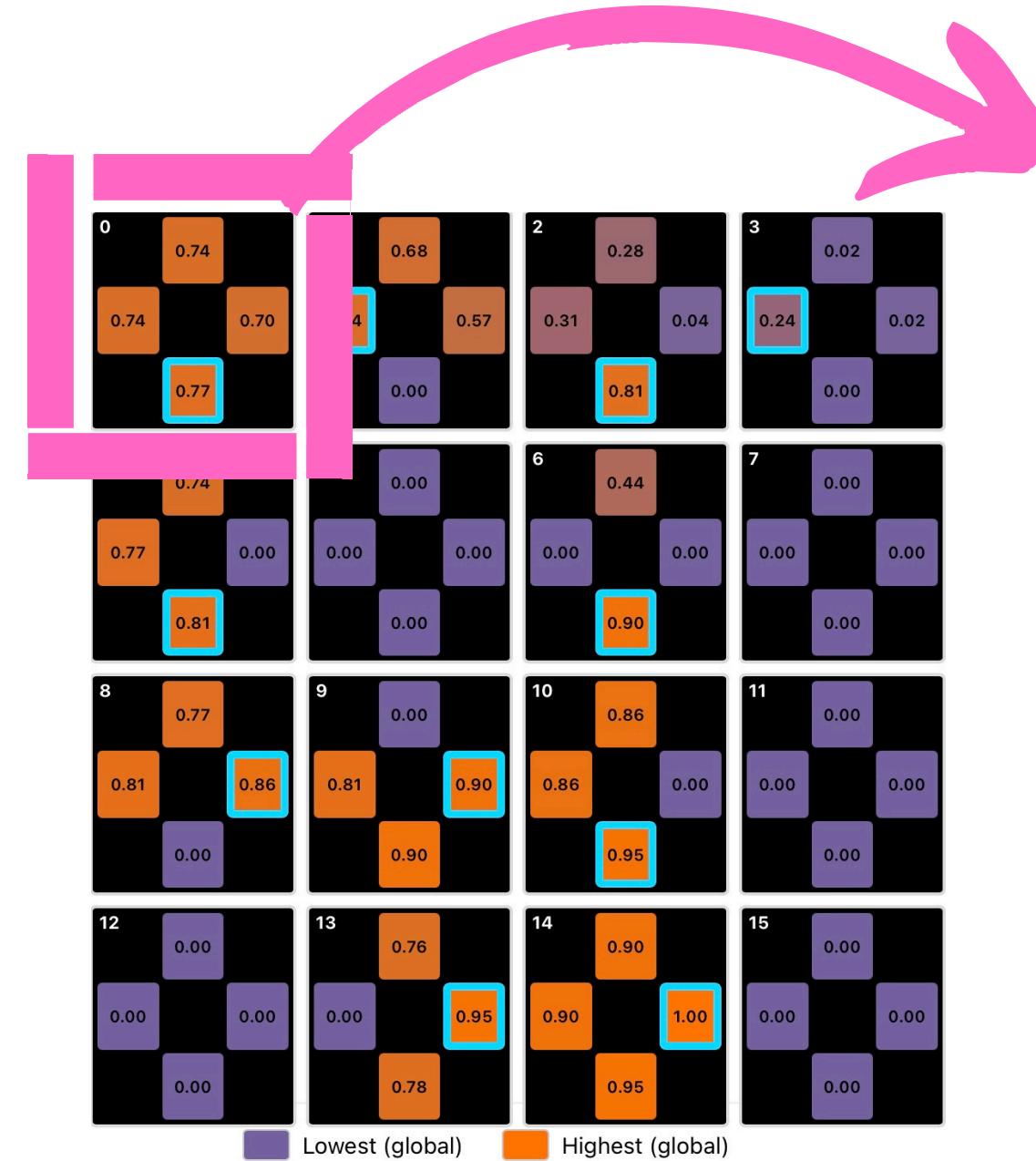
- State** = Position in the grid
- Action** = if possible, left, right, up, down
- Reward** = +1 for reaching the goal  
0 for every step on a frozen tile  
0 for falling into a hole
- Policy** = choose the action with the best Q-value

# Example: Frozen Lake

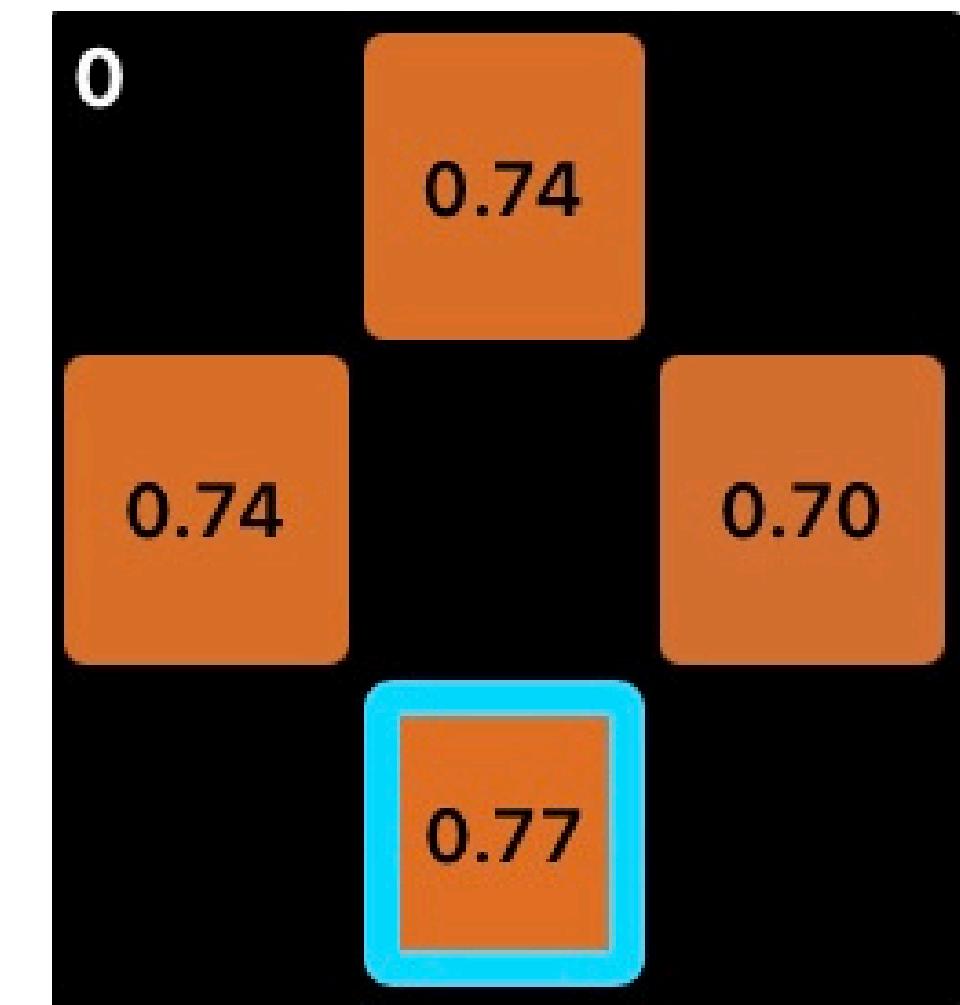
[kisz@hpi.de](mailto:kisz@hpi.de)  
[hpi.de/kisz](http://hpi.de/kisz)



Environment

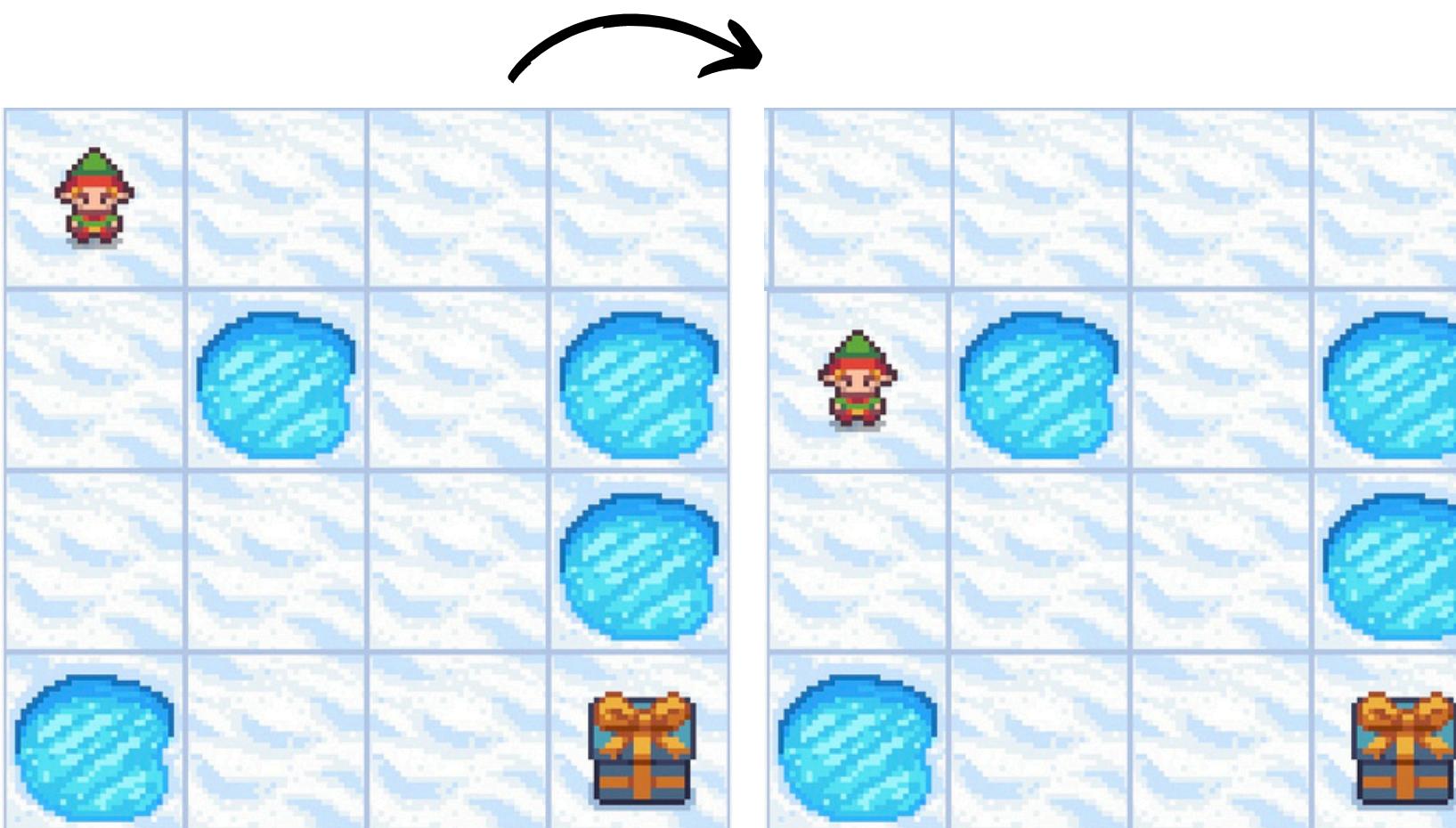


Q-Table



Q-Table for  
the starting position

# Example: Frozen Lake



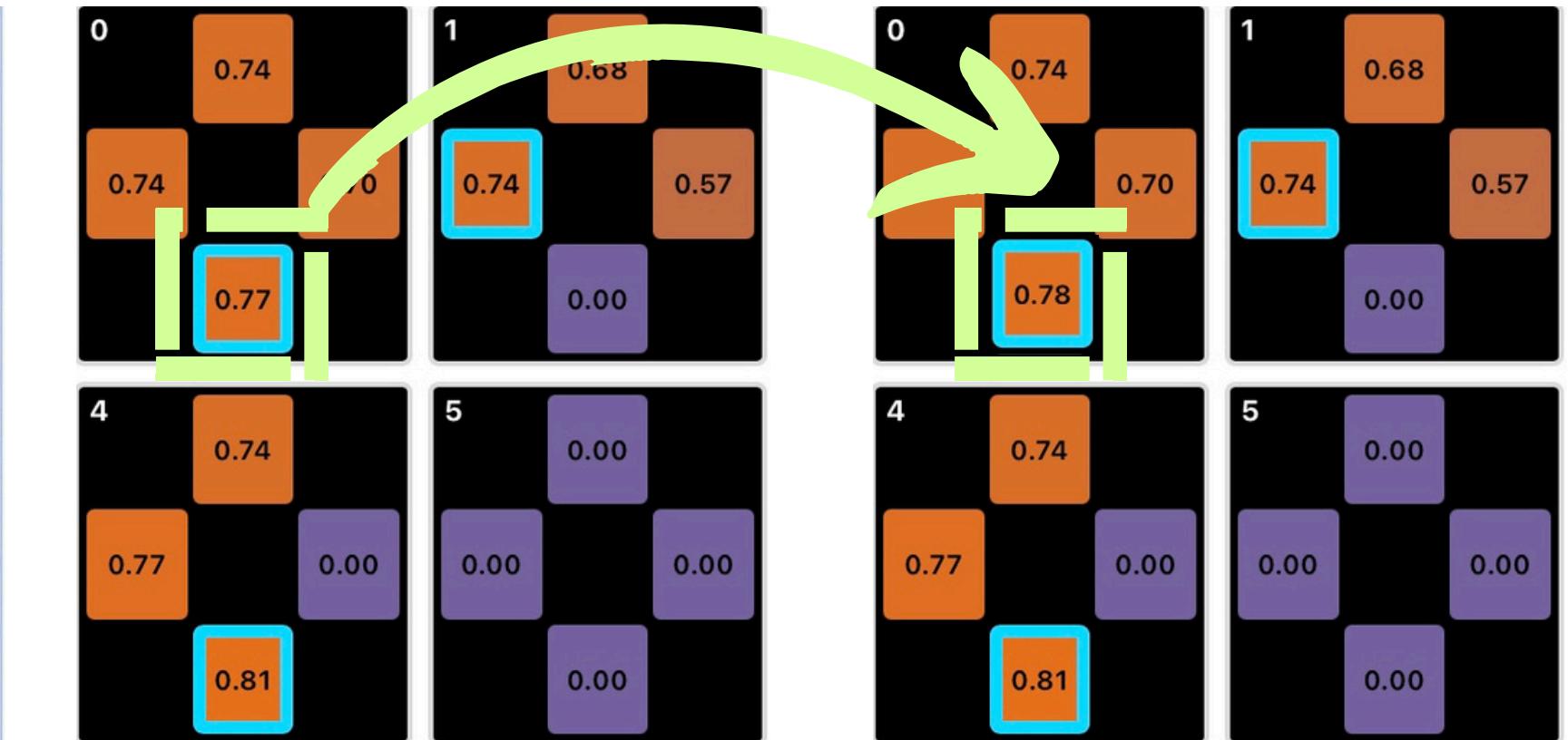
$$Q(s, a) = 0.77$$

$$R = 0$$

$$\max Q(s', a') = 0.81$$

$$\alpha = 0.25$$

$$\gamma = 1.0$$



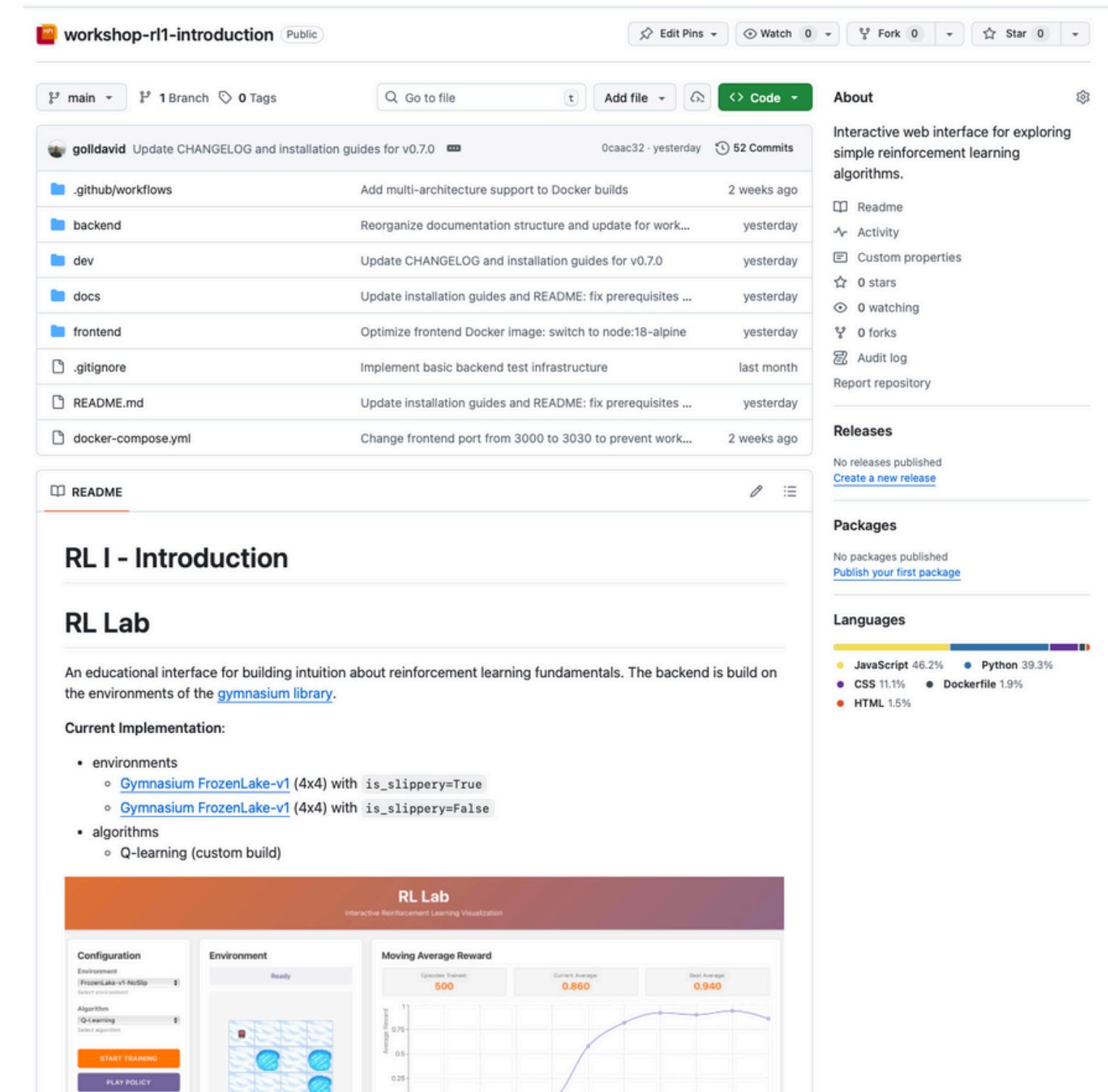
$$Q(s, a) \leftarrow Q(s, a) + \alpha(r + \gamma \max_{a'} Q(s', a') - Q(s, a))$$

$$Q'(s, a) = 0.77 + 0.25(0.81 - 0.77) = 0.78$$

# Try it out

<https://github.com/aihpi/workshop-rl1-introduction>

[kisz@hpi.de](mailto:kisz@hpi.de)  
[hpi.de/kisz](http://hpi.de/kisz)



The screenshot shows the GitHub repository page for `workshop-rl1-introduction`. The repository is public and has 1 branch and 0 tags. The main commit by `golldavid` is to update CHANGELOG and installation guides for v0.7.0. The repository contains several subfolders like `backend`, `dev`, `docs`, `frontend`, and `README.md`. The `README` file contains sections for **RL I - Introduction** and **RL Lab**. The **RL Lab** section describes it as an educational interface for building intuition about reinforcement learning fundamentals, built on the [gymnasium library](#). It lists current implementations: environments (FrozenLake-v1 variants) and algorithms (Q-learning). A preview window shows the **RL Lab** interface with configuration options for environment and algorithm, and a visualization of moving average reward over episodes trained.

# What learning parameters are there?

[kisz@hpi.de](mailto:kisz@hpi.de)  
[hpi.de/kisz](http://hpi.de/kisz)

$$Q'(s, a) = Q(s, a) + \alpha \left[ r + \gamma \max_{a'} Q(s', a') - Q(s, a) \right]$$

$$a = \begin{cases} \arg \max_{a'} Q(s, a') & \text{with probability } 1 - \epsilon, \\ \text{random action} & \text{with probability } \epsilon. \end{cases}$$

Parameter	Symbol	Meaning	Impact
learning rate	$\alpha$	determines how strongly new information influences the previous Q-value	high values → fast but unstable learning low values → slow but stable learning
discount factor	$\gamma$	weights future rewards relative to immediate rewards	high values → long-term focus low values → short-term focus
exploration rate	$\epsilon$	indicates how often random actions are chosen instead of the current best action	high values → explore many different actions low values → exploit the learned action more
number of episodes	-	number of learning cycles the model goes through	few episodes → poor learning many episodes → risk of overfitting, training time
initialization of Q-values	-	initial values of the Q-table	influences how quickly the model arrives at good decisions

# Challenges in RL

[kisz@hpi.de](mailto:kisz@hpi.de)  
[hpi.de/kisz](http://hpi.de/kisz)

## Exploration vs. Exploitation

- Balance between exploring new actions and exploiting known knowledge
- Too much exploration → inefficient
- Too little → gets stuck in poor strategies

## Reward Shaping

- Rewards must encourage the desired behavior
- Poor or overly simple reward designs can lead to unexpected strategies (CAREFUL!)

## Learning Stability & Convergence

- Updates can become unstable or diverge

## Sample Efficiency

- RL often requires many interactions
- In real scenarios (e.g., robotics) this is costly or slow

## Partial Observability & Uncertainty

- The agent often cannot observe the full state of the environment
- Must still make robust decisions

## Scalability & Complexity

- Large state or action spaces make learning computationally expensive/infeasible
- Often requires function approximation (e.g., neural networks)

# What have we learned so far?

[kisz@hpi.de](mailto:kisz@hpi.de)  
[hpi.de/kisz](http://hpi.de/kisz)

## Basics

- **Agent-Environment interaction**
- **Policy learns to maximize** long-term cumulative **reward**
- Learning through **trial and error**
- Parameters:  $\alpha$  (learning rate),  $\gamma$  (discount factor),  $\epsilon$  (exploration rate)

## Specifically

- **Q-Learning**
- **Q-Value function:** measures the expected future reward of an action in a specific state
- **Gym library** (FrozenLake)

## Potential and Limitations

- **Potential:** Provides a framework for continual learning via interaction (core idea of RL is closely related to biological learning)
- **Challenges:** Sample inefficiency, sim2real, reward design, instability, safety,

# Outlook

[kisz@hpi.de](mailto:kisz@hpi.de)  
[hpi.de/kisz](http://hpi.de/kisz)

**Deep Reinforcement Learning**

**Procedure in an RL project**

**Choice of algorithm**

**Current application areas**

[kisz@hpi.de](mailto:kisz@hpi.de)

[hpi.de/kisz](http://hpi.de/kisz)

Your opinion is  
relevant!



QR code to feedback form

Gefördert durch:



Bundesministerium  
für Forschung, Technologie  
und Raumfahrt

