# DeepLTL: Learning to Efficiently Satisfy Complex LTL Specifications for Multi-Task RL
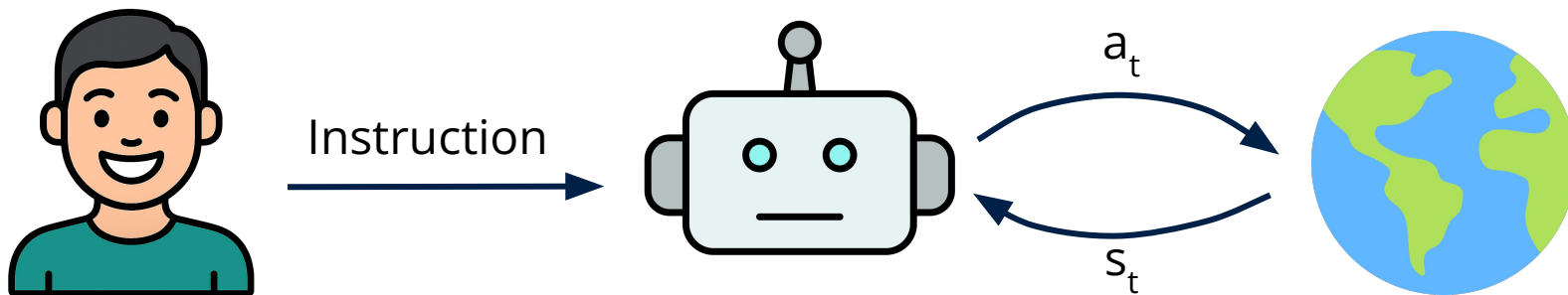
Mathias Jackermeier, Alessandro Abate

ICLR 2025, Singapore
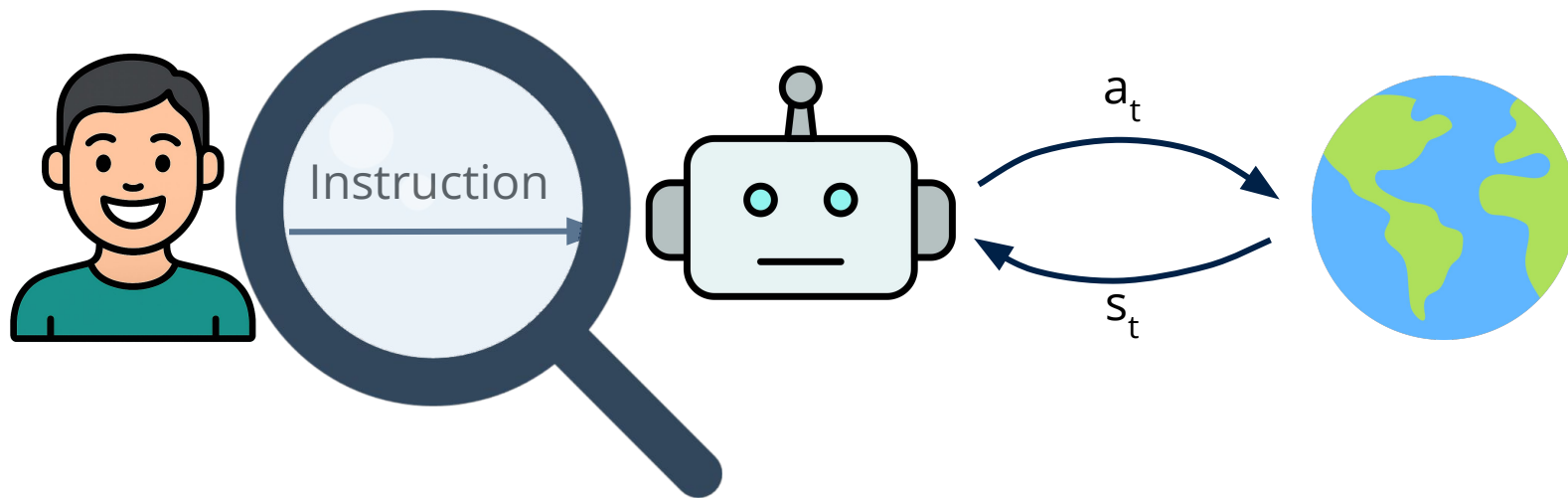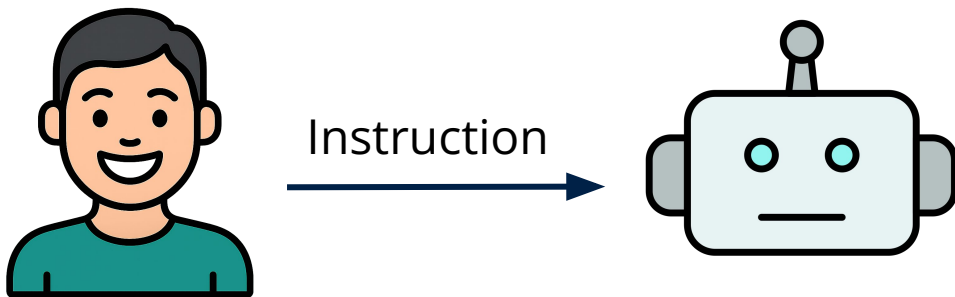
# Instruction-following RL agents

# Instruction-following RL agents

# Instruction-following RL agents



Instruction

**Natural language**:
✅ Intuitive
❌ Ambiguous
❌ Difficult to assess

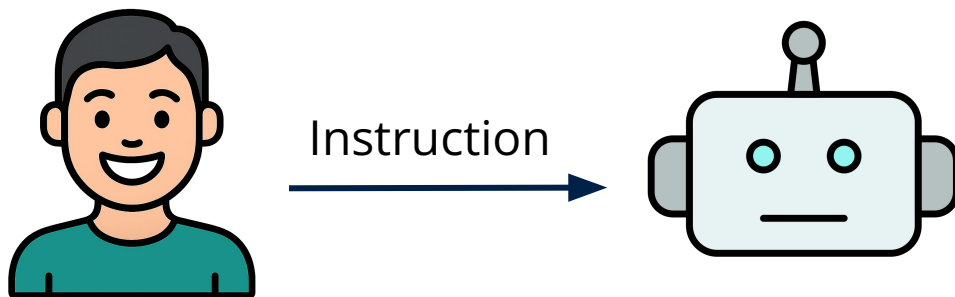# Instruction-following RL agents

Instruction

**Formal specifications**:
- ✅ Precise
- ✅ Easy to verify
- ✅ Explicit structure
- ❌ Difficult to formulate (?)

# Instruction-following RL agents

Instruction

**Formal specifications**:
- ✅ Precise
- ✅ Easy to verify
- ✅ Explicit structure
- ❌ Difficult to formulate (?)

→ Well suited when correctness is crucial, e.g. safety-critical settings

# Linear temporal logic (LTL)

Example specification:

$$(\neg \text{yellow } U \text{ purple}) \wedge G\,(\text{green} \Rightarrow F \text{ blue})$$

# Linear temporal logic (LTL)

Example specification:

$$(\neg\text{yellow U purple}) \wedge G\,(\text{green} \Rightarrow F\,\text{blue})$$

"Go to the purple zone while avoiding the yellow region,

# Linear temporal logic (LTL)

Example specification:

$$(\neg\text{yellow } U \text{ purple}) \wedge G\,(\text{green} \Rightarrow F\,\text{blue})$$

"Go to the purple zone while avoiding the yellow region, and

# Linear temporal logic (LTL)

Example specification:

$$(\neg \text{yellow } U \text{ purple}) \wedge G (\text{green} \Rightarrow F \text{ blue})$$

"Go to the purple zone while avoiding the yellow region, and
always,

# Linear temporal logic (LTL)

Example specification:

$$(\neg \text{yellow U purple}) \wedge \text{G} (\text{green} \Rightarrow \text{F blue})$$

"Go to the purple zone while avoiding the yellow region, and always, if you visit green you eventually have to go to blue."

# Linear temporal logic (LTL)

Example specification:

$$(\neg \text{yellow U purple}) \wedge \text{G (green} \Rightarrow \text{F blue)}$$

"Go to the purple zone while avoiding the yellow region, and always, if you visit green you eventually have to go to blue."

# Linear temporal logic (LTL)

Example specification:

$$(\neg\text{yellow U purple}) \land \text{G (green} \Rightarrow \text{F blue})$$

"Go to the purple zone while avoiding the yellow region, and always, if you visit green you eventually have to go to blue."

✅ Compositional     ✅ Temporally extended     ✅ Infinite horizon     ✅ Safety constraints

How can we train a **multi-task** policy to **zero-shot** execute **arbitrary** LTL specifications?
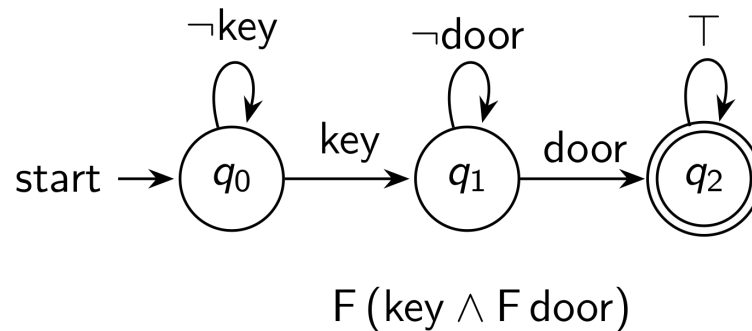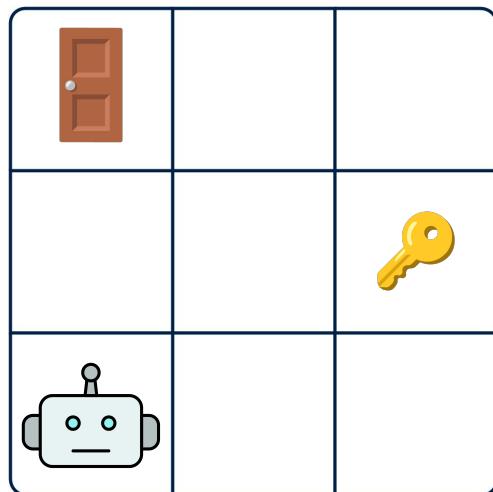
# From LTL specifications to automata

Any LTL specification can be converted to an equivalent (Büchi) **automaton**:

$$\text{G F blue} \wedge \text{G F green}$$

# Product MDP

Keeping track of the automaton state allows us to learn a **Markovian** policy

$$\pi : \mathcal{S} \times \mathcal{Q} \rightarrow \Delta(\mathcal{A})$$



$$\text{F}(\text{key} \wedge \text{F door})$$

# Product MDP

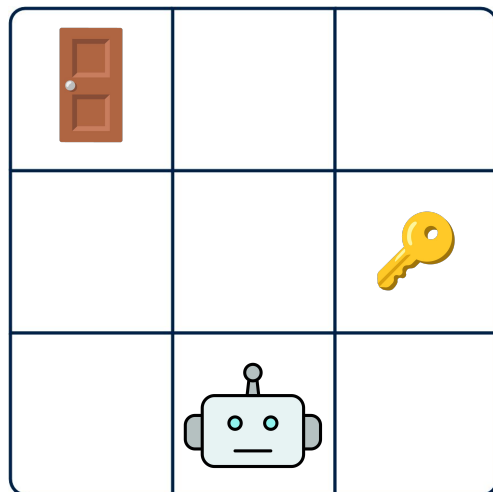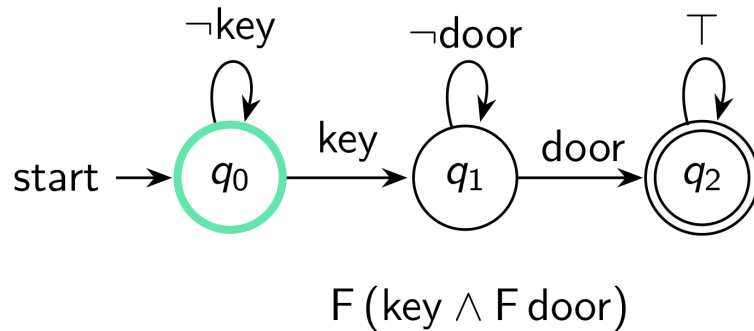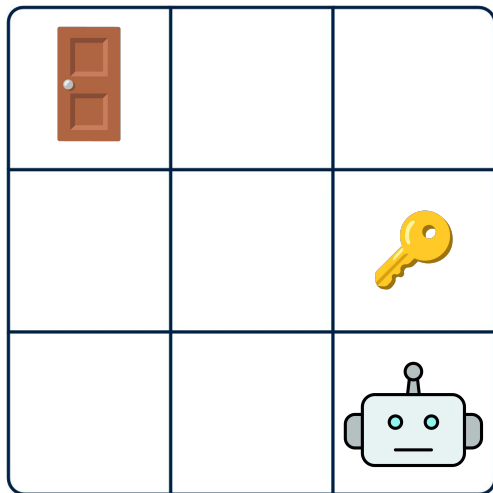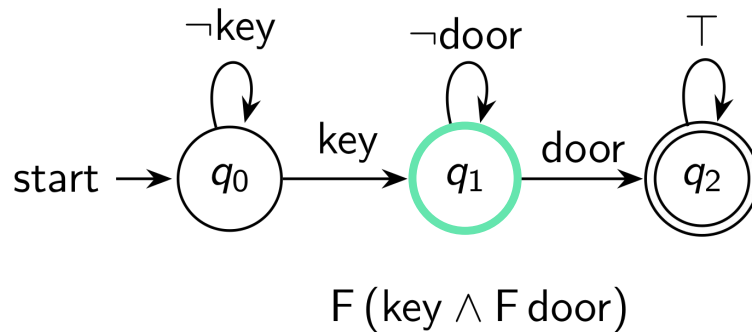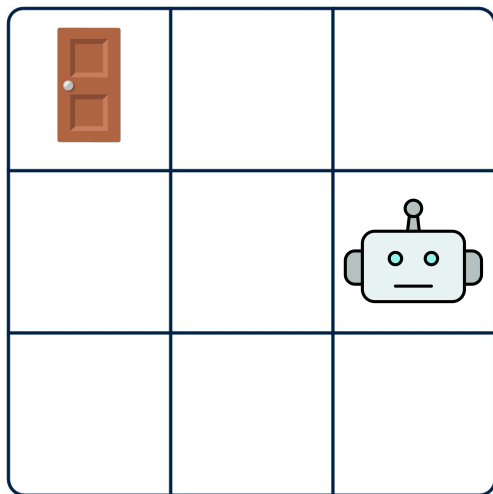Keeping track of the automaton state allows us to learn a **Markovian** policy
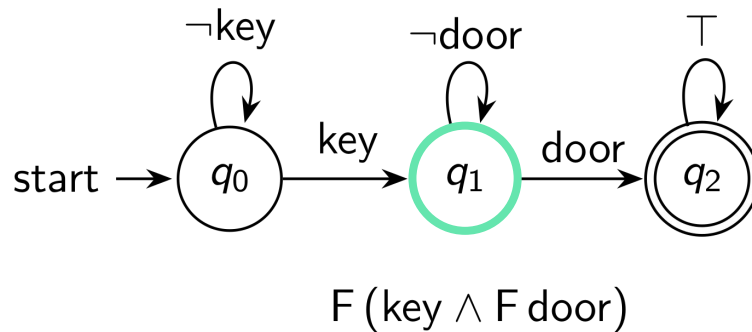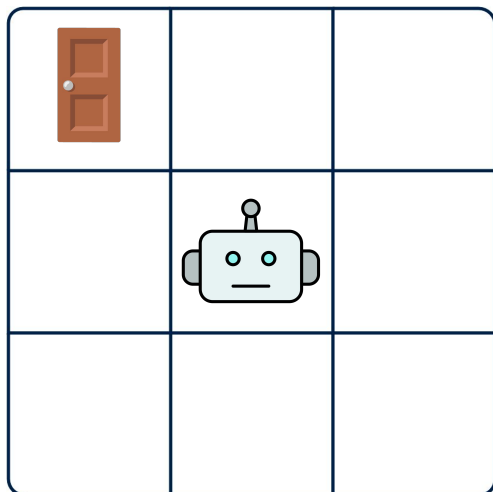
$$\pi : \mathcal{S} \times \mathcal{Q} \rightarrow \Delta(\mathcal{A})$$



F (key ∧ F door)

# Product MDP

Keeping track of the automaton state allows us to learn a **Markovian** policy

$$\pi : \mathcal{S} \times \mathcal{Q} \to \Delta(\mathcal{A})$$



$$\neg\text{key} \qquad \neg\text{door} \qquad \top$$

$$\text{start} \rightarrow q_0 \xrightarrow{\text{key}} q_1 \xrightarrow{\text{door}} q_2$$

$$F\,(\text{key} \wedge F\,\text{door})$$

# Product MDP

Keeping track of the automaton state allows us to learn a **Markovian** policy

$$\pi : \mathcal{S} \times \mathcal{Q} \rightarrow \Delta(\mathcal{A})$$



start $\rightarrow$ $q_0$ — key $\rightarrow$ $q_1$ — door $\rightarrow$ $q_2$

$\neg$key (self-loop on $q_0$), $\neg$door (self-loop on $q_1$), $\top$ (self-loop on $q_2$)

$\mathsf{F}\,(\mathsf{key} \wedge \mathsf{F}\,\mathsf{door})$

# Product MDP

Keeping track of the automaton state allows us to learn a **Markovian** policy
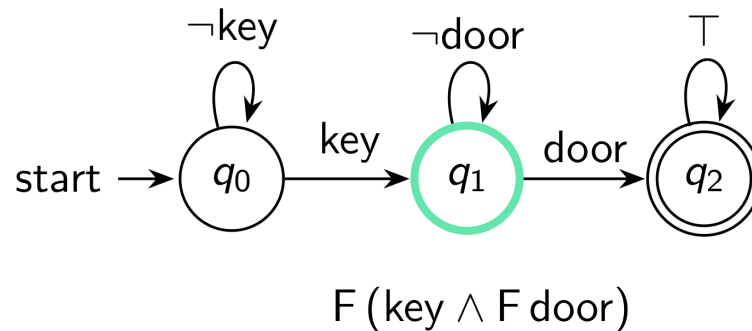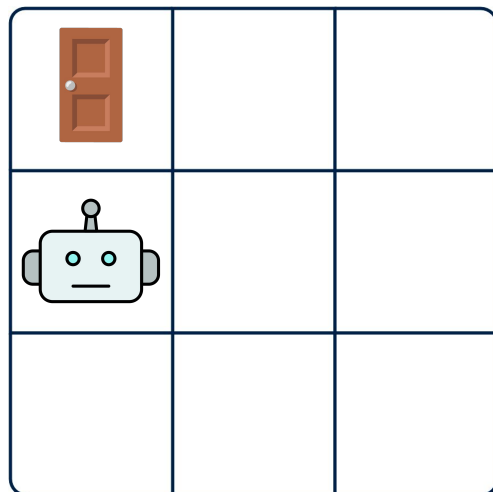
$$\pi : \mathcal{S} \times \mathcal{Q} \rightarrow \Delta(\mathcal{A})$$



F (key ∧ F door)

# Product MDP

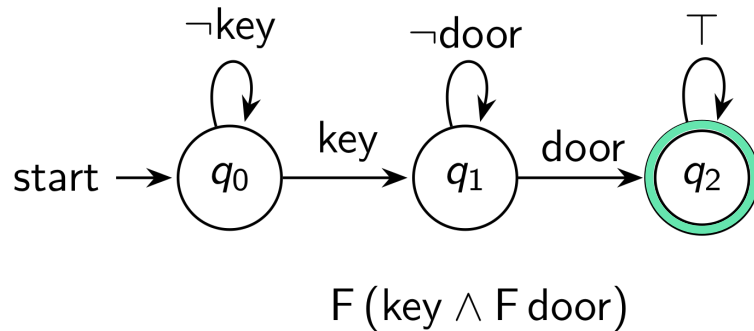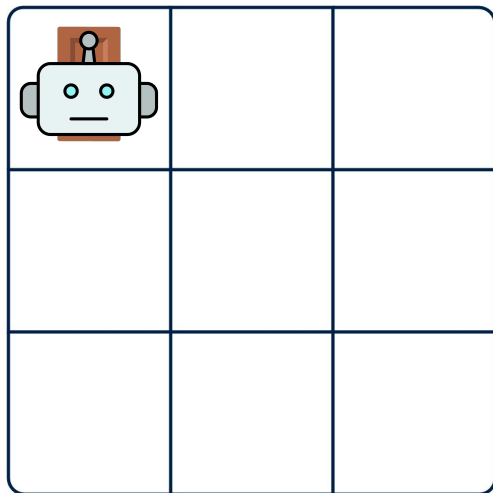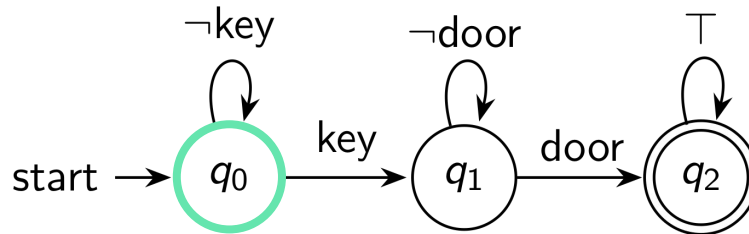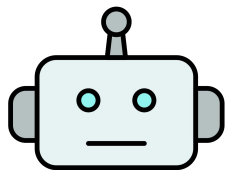Keeping track of the automaton state allows us to learn a **Markovian** policy

$$\pi : \mathcal{S} \times \mathcal{Q} \to \Delta(\mathcal{A})$$



$\neg key$

$\neg door$

$\top$

start $\longrightarrow$ $q_0$ $\xrightarrow{\text{key}}$ $q_1$ $\xrightarrow{\text{door}}$ $q_2$

F (key $\wedge$ F door)

# Product MDP

Keeping track of the automaton state allows us to learn a **Markovian** policy

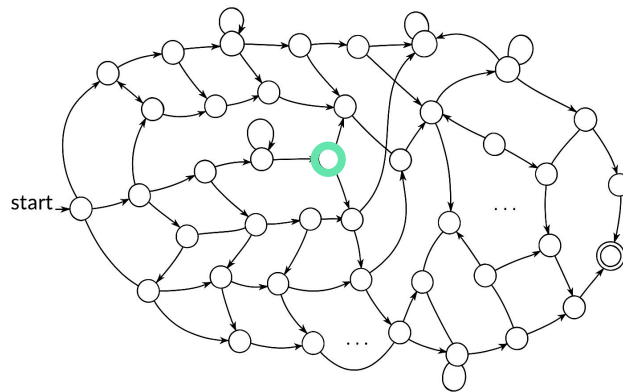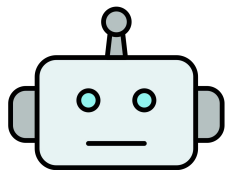$$\pi : \mathcal{S} \times \mathcal{Q} \rightarrow \Delta(\mathcal{A})$$



F (key ∧ F door)

# Product MDP

Keeping track of the automaton state allows us to learn a **Markovian** policy

$$\pi : \mathcal{S} \times \mathcal{Q} \rightarrow \Delta(\mathcal{A})$$



$$\neg\text{key} \qquad \neg\text{door} \qquad \top$$

start $\rightarrow$ $q_0$ $\xrightarrow{\text{key}}$ $q_1$ $\xrightarrow{\text{door}}$ $q_2$

$$\text{F}\,(\text{key} \wedge \text{F door})$$

# From single-task to multi-task

In a multi-task setting, we do not know the automaton beforehand

# From single-task to multi-task

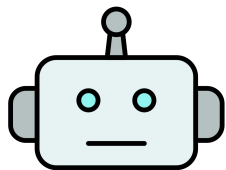In a multi-task setting, we do not know the automaton beforehand
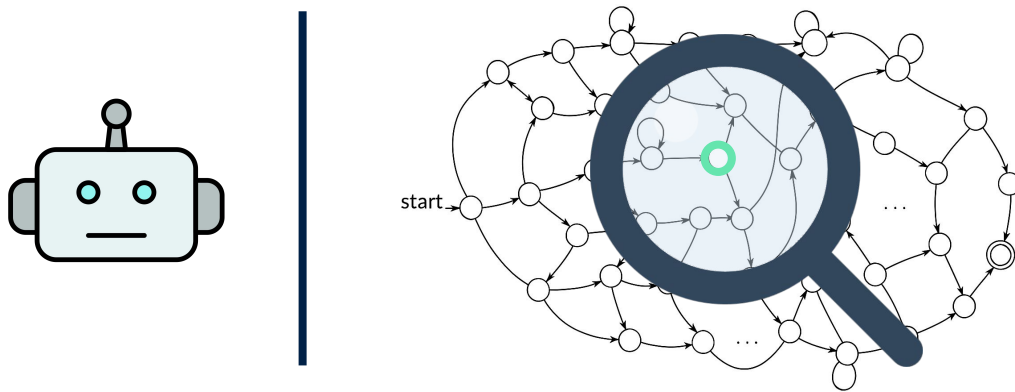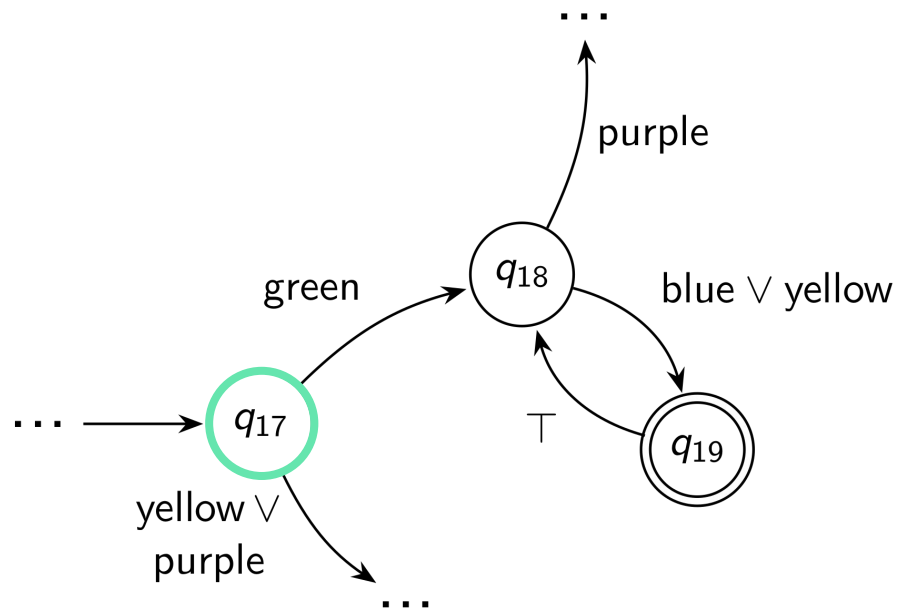
# From single-task to multi-task

In a multi-task setting, we do not know the automaton beforehand



What is a **general representation** of the automaton state that can be used to condition the policy?

# From single-task to multi-task

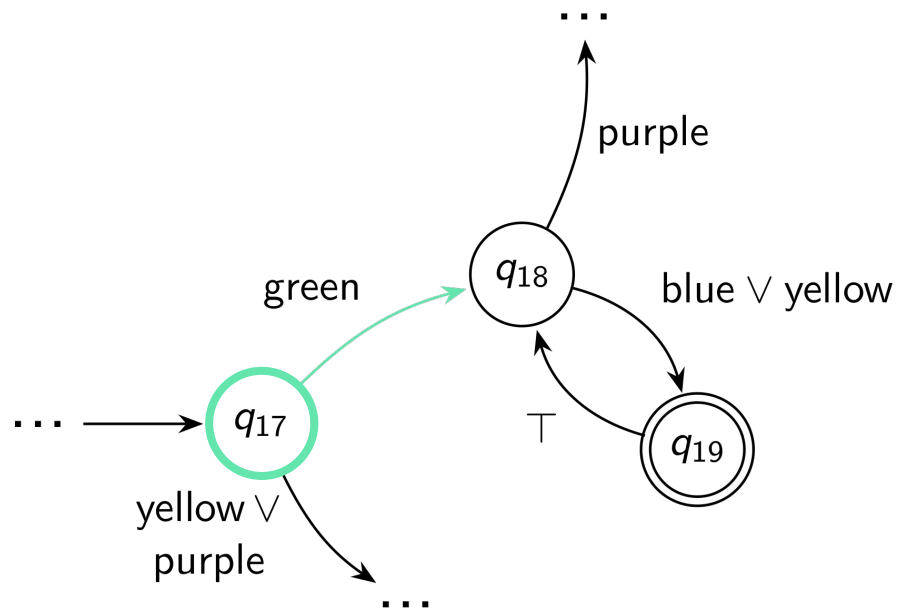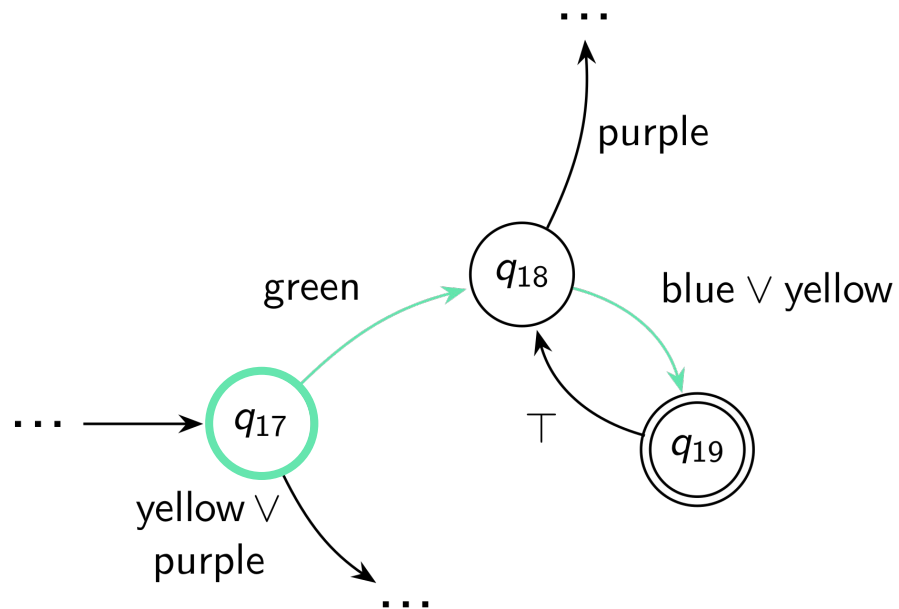In a multi-task setting, we do not know the automaton beforehand



What is a **general representation** of the automaton state that can be used to condition the policy?
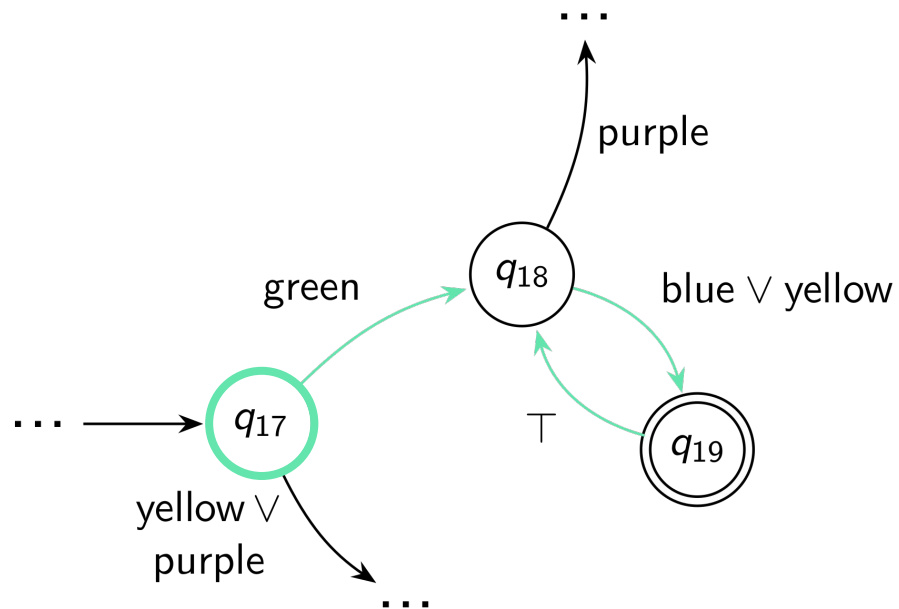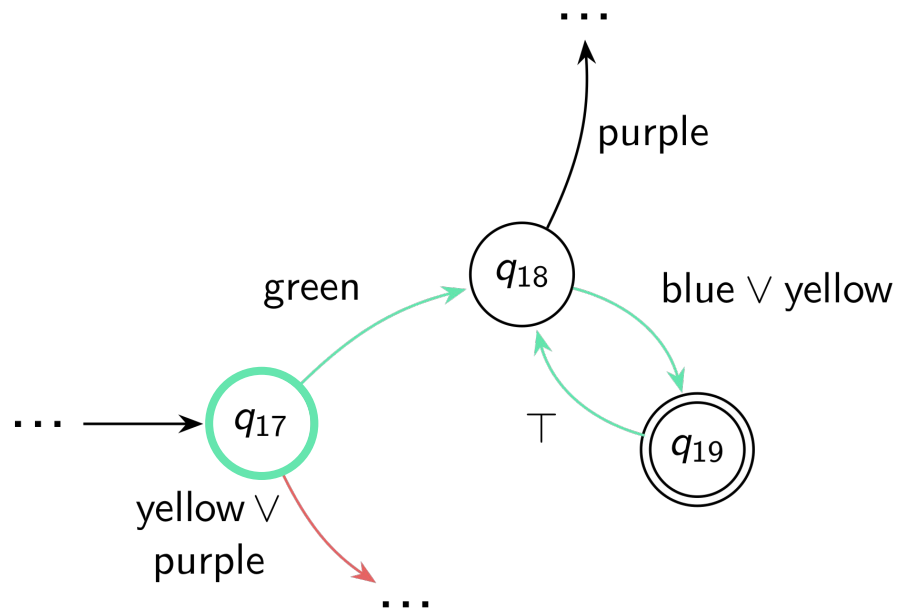
# Reach-avoid sequences

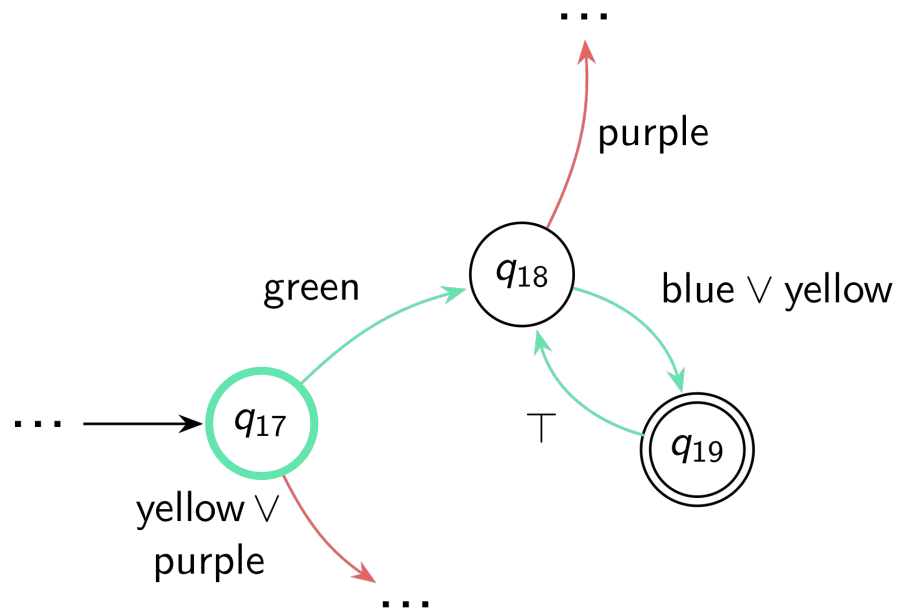# Reach-avoid sequences

# Reach-avoid sequences

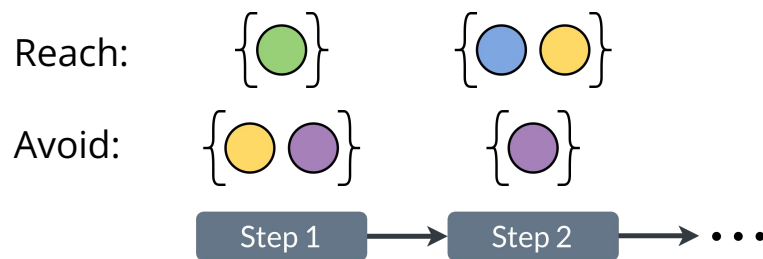# Reach-avoid sequences

# Reach-avoid sequences
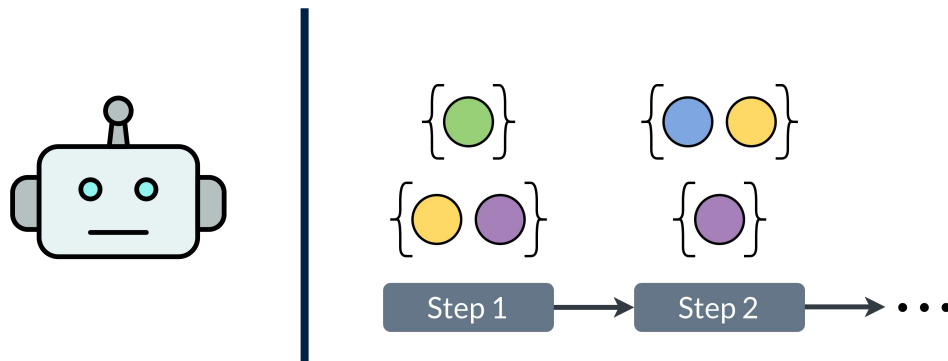
# Reach-avoid sequences

# Reach-avoid sequences

# Training a general policy

We use goal-conditioned RL to train a general policy:

# Training a general policy

We use goal-conditioned RL to train a general policy:

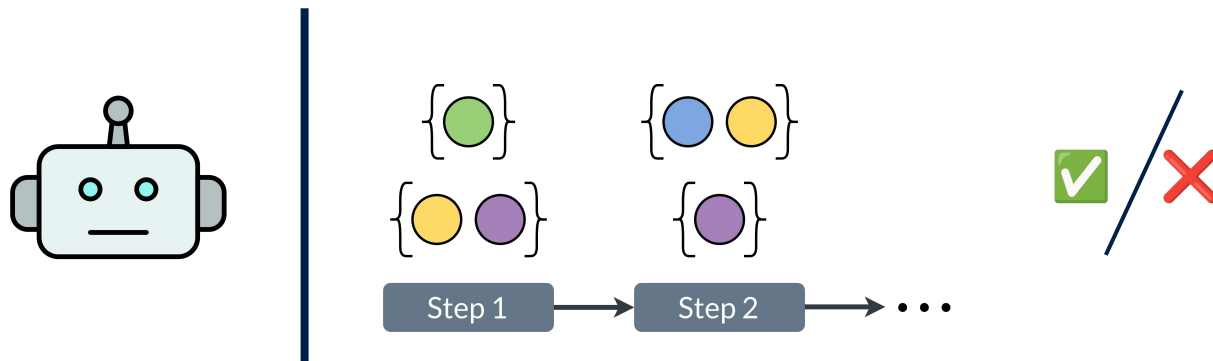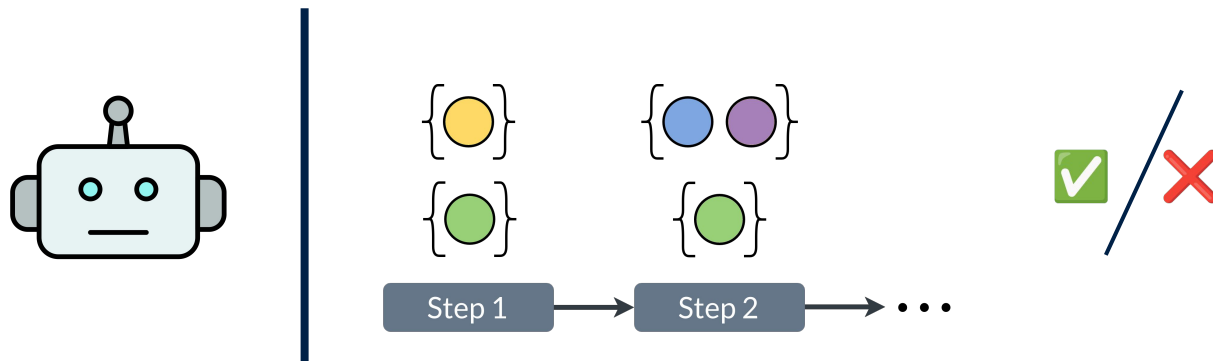# Training a general policy

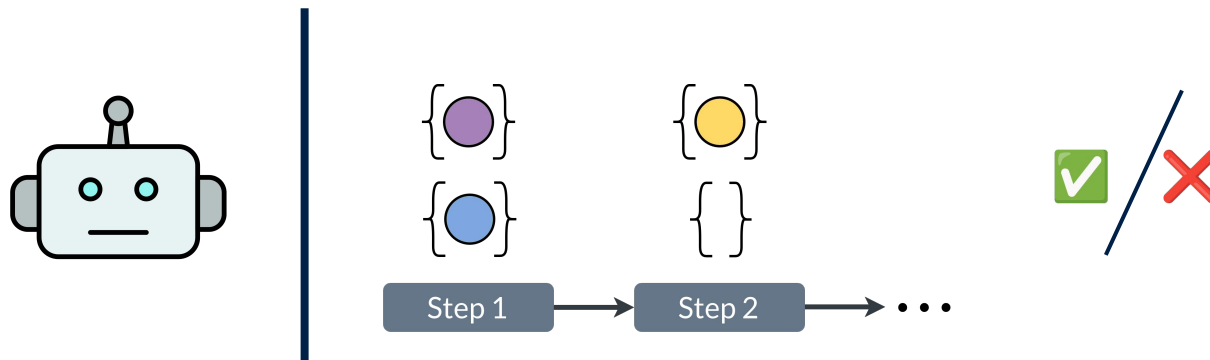We use goal-conditioned RL to train a general policy:

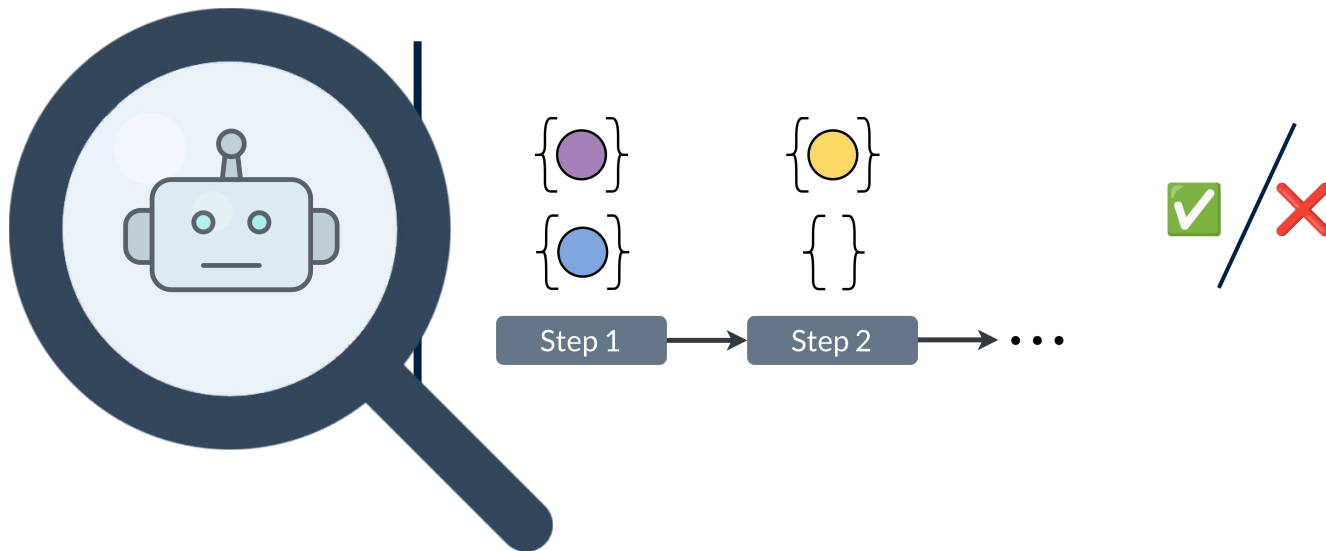# Training a general policy

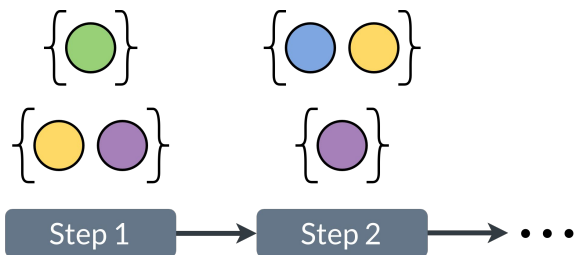We use goal-conditioned RL to train a general policy:
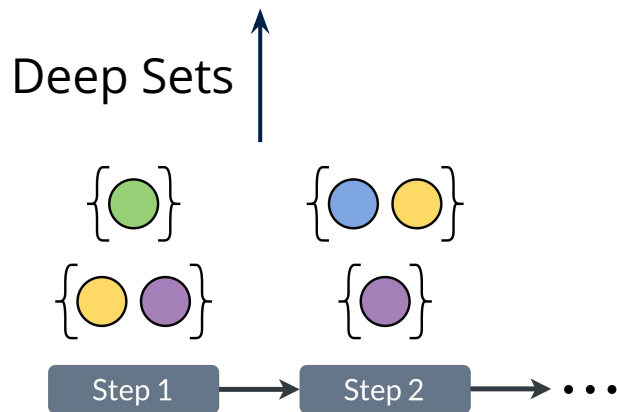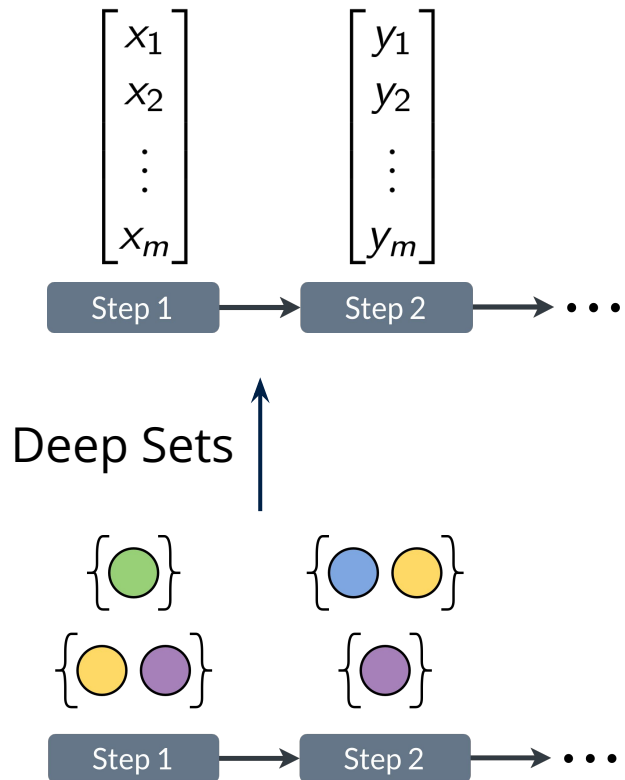
# Training a general policy

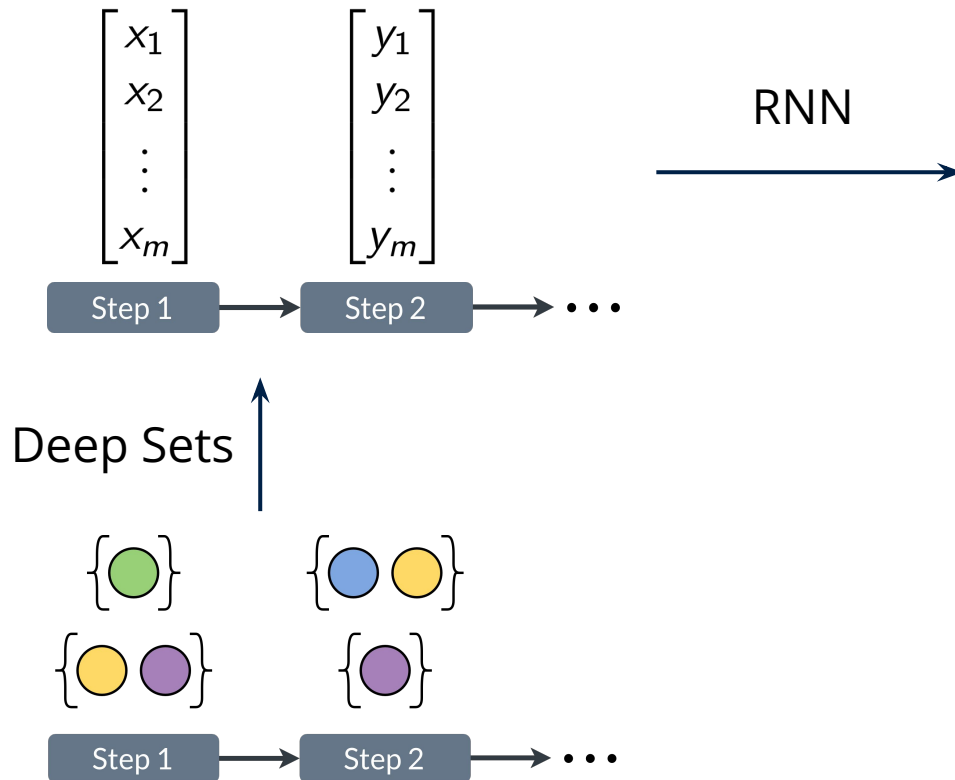We use goal-conditioned RL to train a general policy:

# Model architecture

# Model architecture



Deep Sets

{ 🟢 }    { 🔵 🟡 }

{ 🟡 🟣 }    { 🟣 }

Step 1 → Step 2 → • • •

# Model architecture

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix}$$

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$

Step 1 $\rightarrow$ Step 2 $\rightarrow$ $\bullet\bullet\bullet$

Deep Sets

$\{\ \bullet\ \}$  $\{\ \bullet\ \bullet\ \}$

$\{\ \bullet\ \bullet\ \}$  $\{\ \bullet\ \}$

Step 1 $\rightarrow$ Step 2 $\rightarrow$ $\bullet\bullet\bullet$
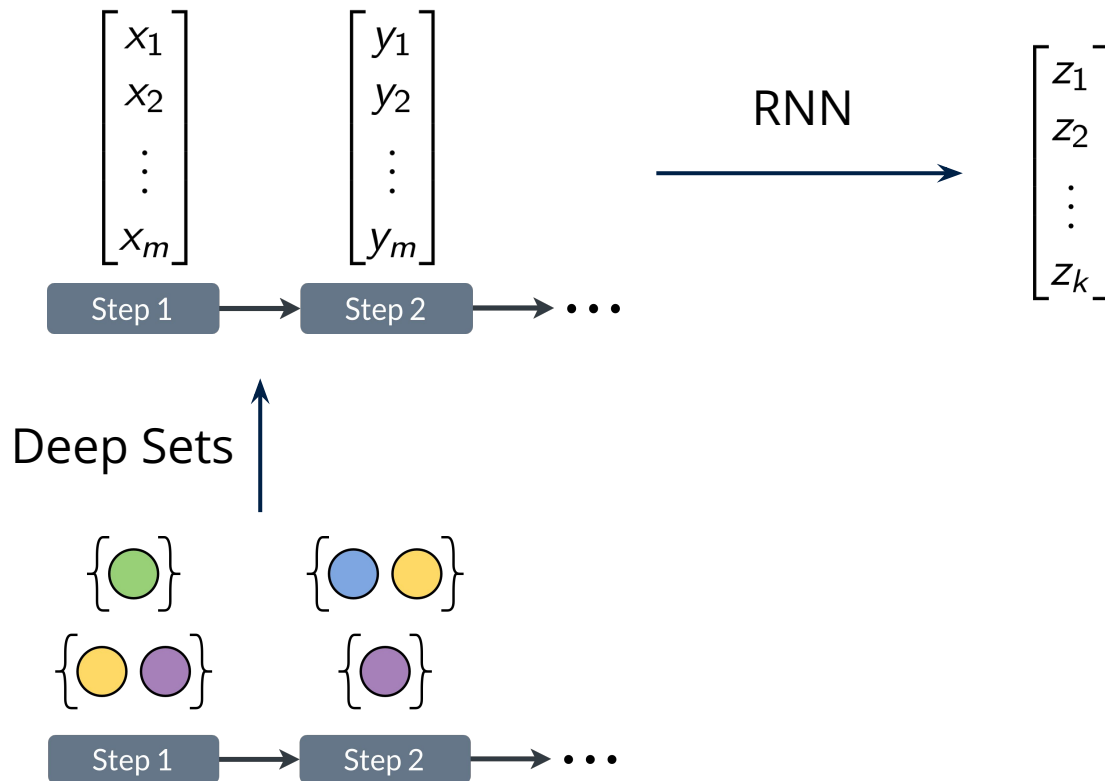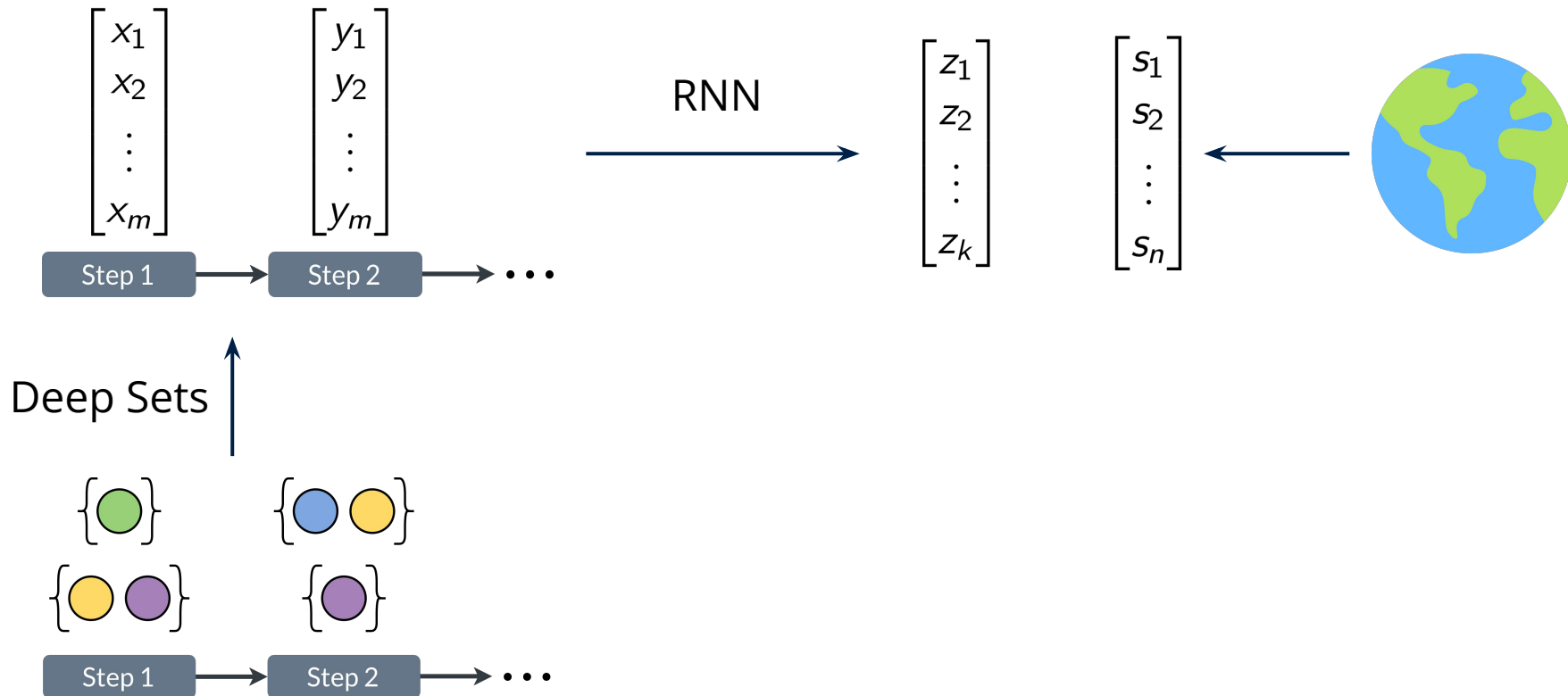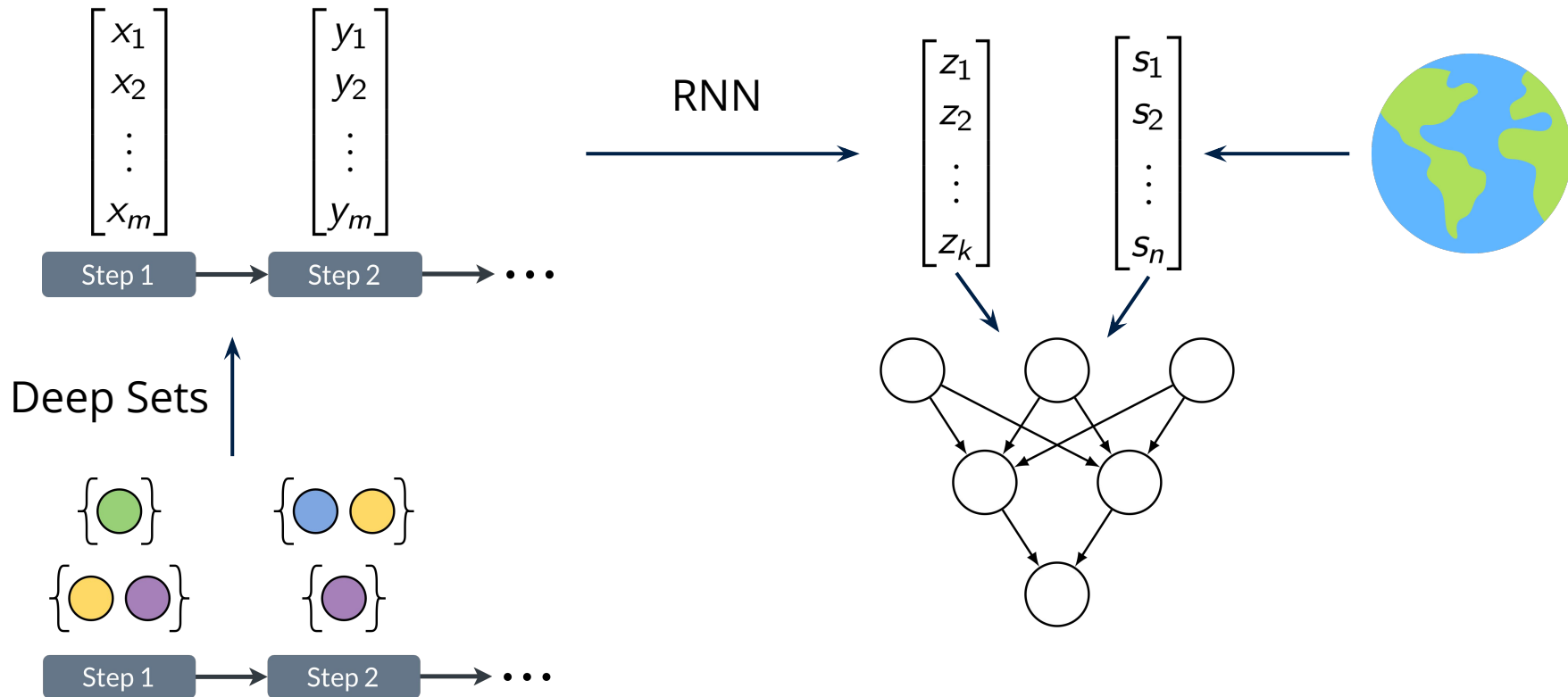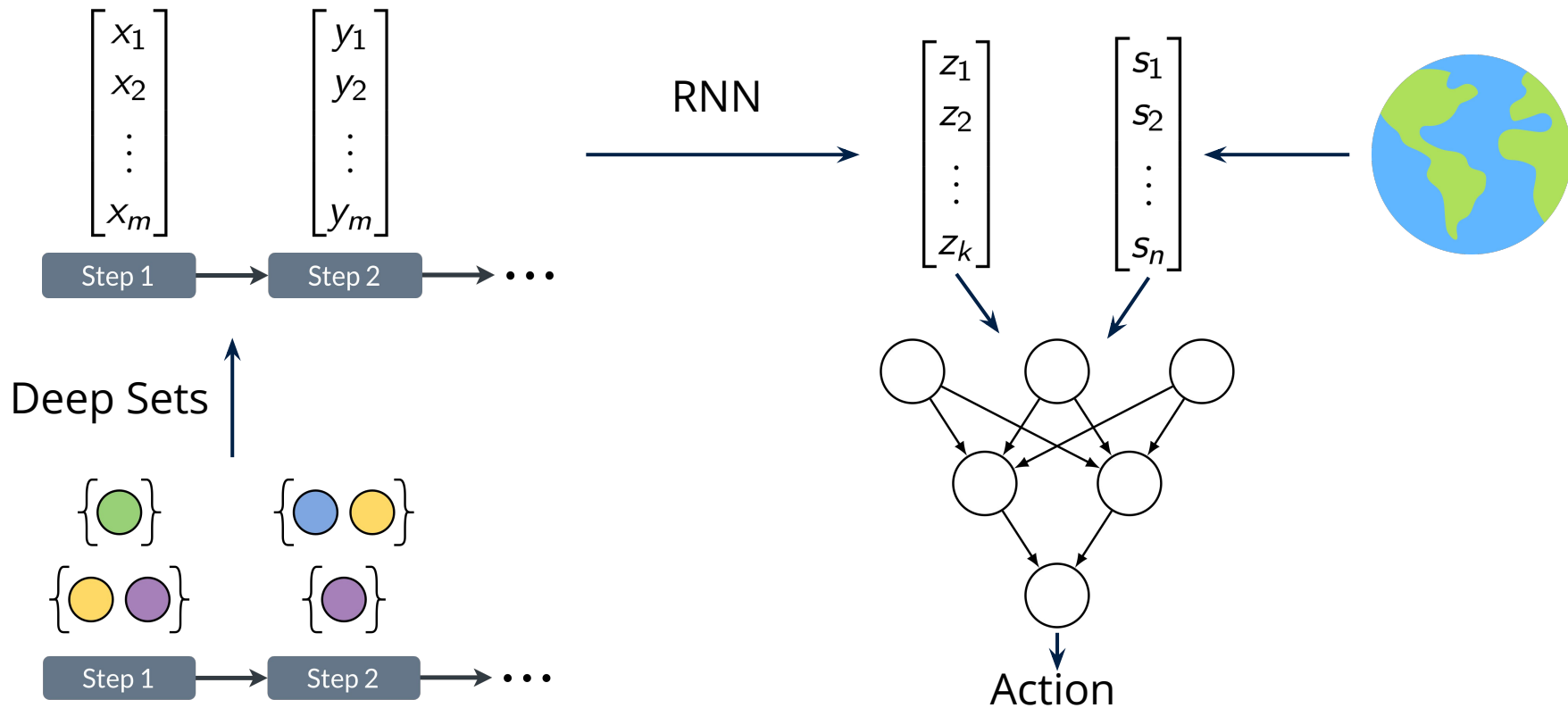
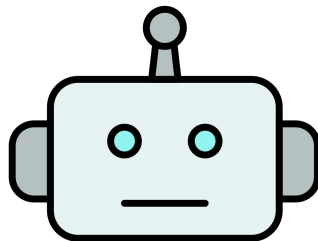# Model architecture

# Model architecture

# Model architecture

# Model architecture
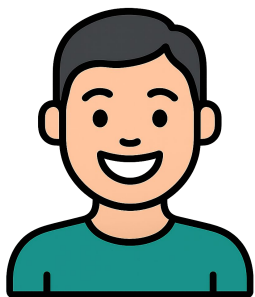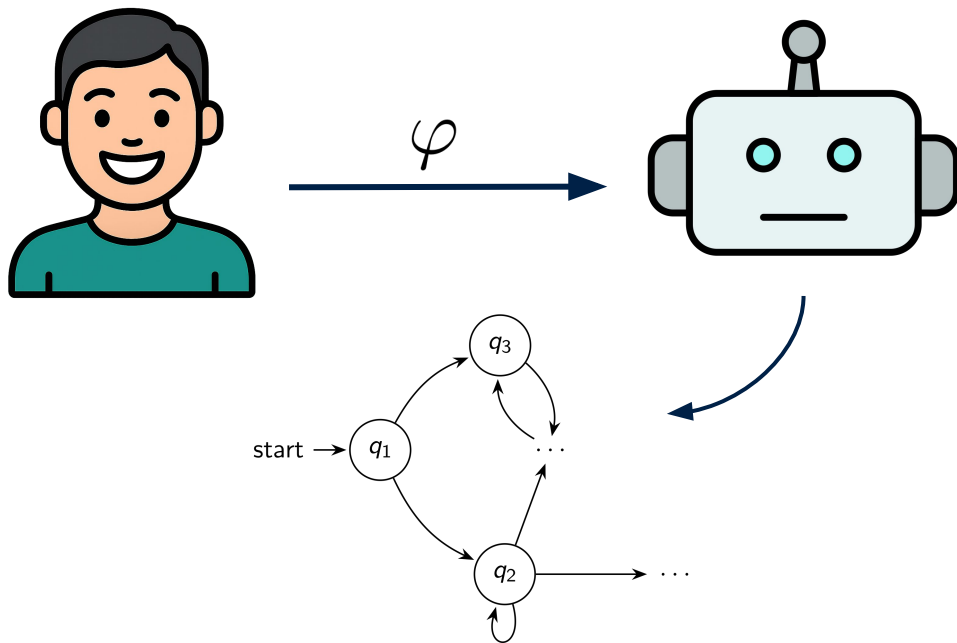
# Model architecture

# Test-time policy execution

# Test-time policy execution

# Test-time policy execution

# Test-time policy execution

# Test-time policy execution



$$\sigma^* = \arg\max_{\sigma} V^{\pi}(s, \sigma)$$

# Test-time policy execution



$\varphi$

$a \sim \pi(\cdot, \sigma^*)$

$\sigma^* = \arg\max_{\sigma} V^{\pi}(s, \sigma)$

start $\rightarrow q_1$

$q_2$

$q_3$

# Test-time policy execution



$$a \sim \pi(\cdot, \sigma^*)$$

$$s_t$$

$$\varphi$$

$$\sigma^* = \arg\max_{\sigma} V^{\pi}(s, \sigma)$$

start → $q_1$

$q_3$

$q_2$

...

...

# Test-time policy execution

# Discussion & Results

# Discussion



Infinite-horizon
tasks

Optimality

Safety

# Results



ZoneEnv          LetterEnv          FlatWorld

Discounted return

Number of steps

—— Ours (DeepLTL)     —— GCRL-LTL     —— LTL2Action

# Results



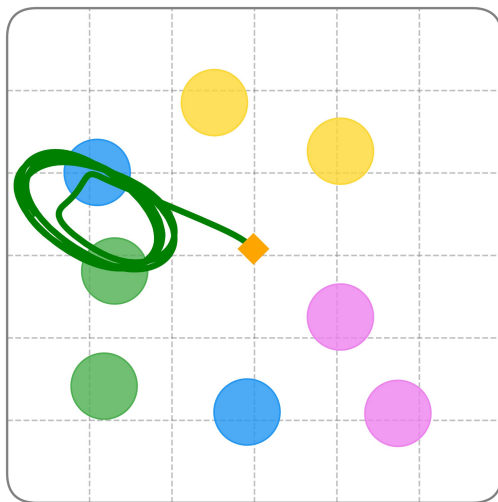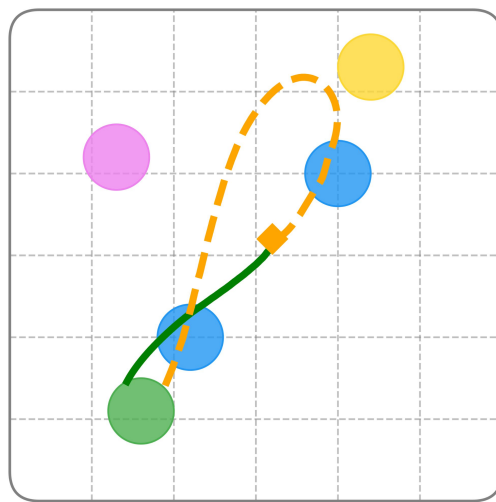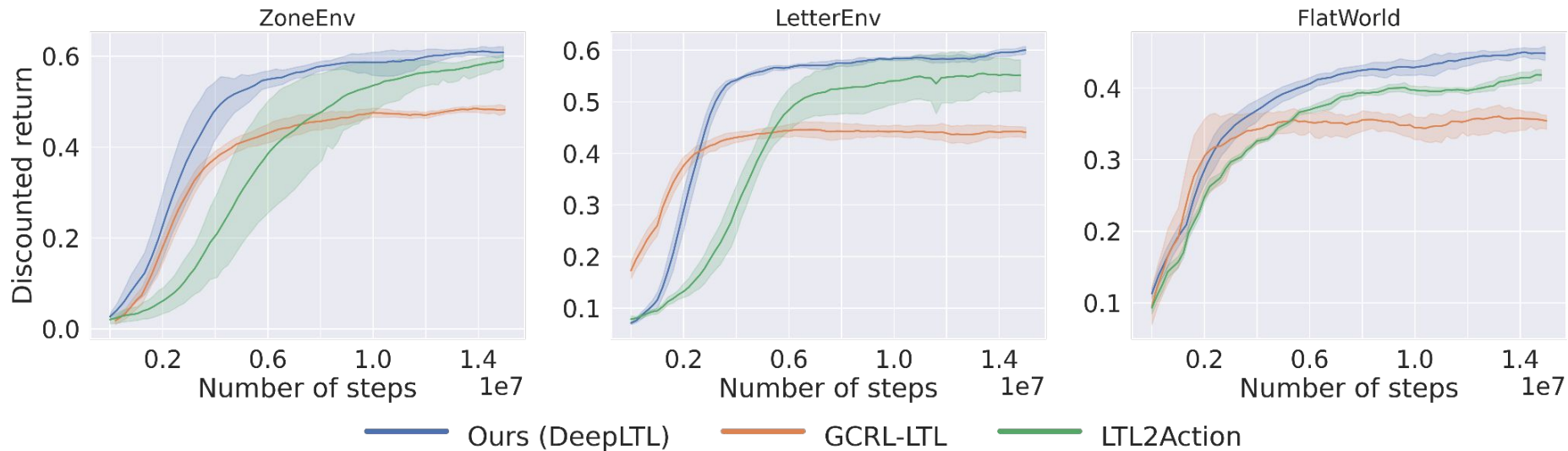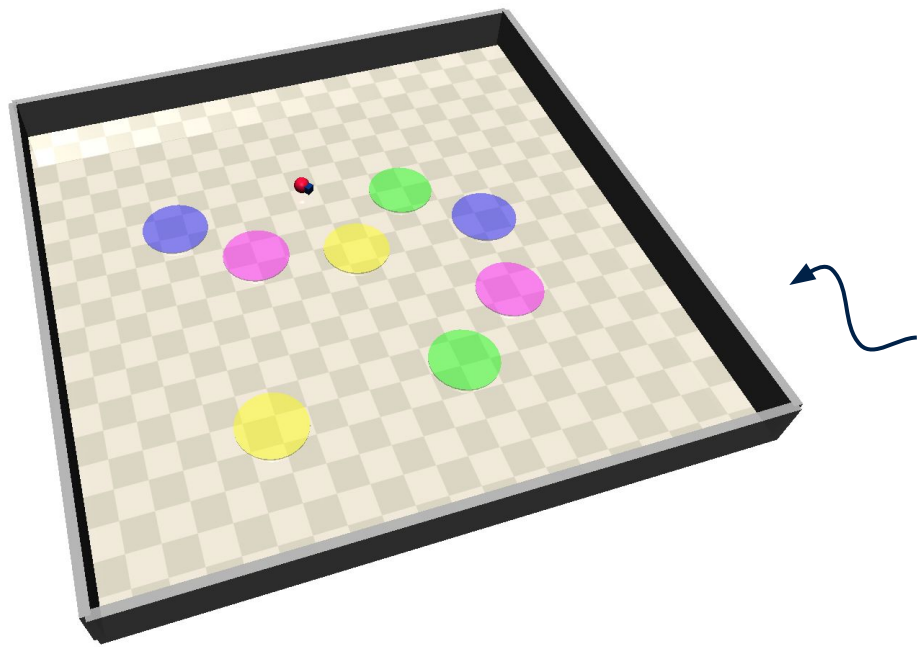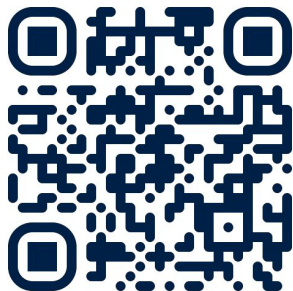|  |  | LTL2Action | GCRL-LTL | DeepLTL |
|---|---|---|---|---|
| LetterWorld | $\varphi_1$ | $0.75_{\pm 0.18}$ | $0.94_{\pm 0.05}$ | $\mathbf{1.00}_{\pm 0.00}$ |
|  | $\varphi_2$ | $0.79_{\pm 0.10}$ | $0.94_{\pm 0.03}$ | $\mathbf{0.98}_{\pm 0.00}$ |
|  | $\varphi_3$ | $0.41_{\pm 0.14}$ | $\mathbf{1.00}_{\pm 0.00}$ | $\mathbf{1.00}_{\pm 0.00}$ |
|  | $\varphi_4$ | $0.72_{\pm 0.17}$ | $0.82_{\pm 0.07}$ | $\mathbf{0.97}_{\pm 0.01}$ |
|  | $\varphi_5$ | $0.44_{\pm 0.26}$ | $\mathbf{1.00}_{\pm 0.00}$ | $\mathbf{1.00}_{\pm 0.00}$ |
| ZoneEnv | $\varphi_6$ | $0.60_{\pm 0.20}$ | $0.85_{\pm 0.03}$ | $\mathbf{0.92}_{\pm 0.06}$ |
|  | $\varphi_7$ | $0.14_{\pm 0.18}$ | $0.85_{\pm 0.05}$ | $\mathbf{0.91}_{\pm 0.03}$ |
|  | $\varphi_8$ | $0.67_{\pm 0.26}$ | $0.89_{\pm 0.04}$ | $\mathbf{0.96}_{\pm 0.04}$ |
|  | $\varphi_9$ | $0.69_{\pm 0.22}$ | $0.87_{\pm 0.02}$ | $\mathbf{0.90}_{\pm 0.03}$ |
|  | $\varphi_{10}$ | $0.66_{\pm 0.19}$ | $0.85_{\pm 0.02}$ | $\mathbf{0.91}_{\pm 0.02}$ |
|  | $\varphi_{11}$ | $0.93_{\pm 0.07}$ | $0.89_{\pm 0.01}$ | $\mathbf{0.98}_{\pm 0.01}$ |
| FlatWorld | $\varphi_{12}$ | $\mathbf{1.00}_{\pm 0.00}$ | $0.82_{\pm 0.41}$ | $\mathbf{1.00}_{\pm 0.00}$ |
|  | $\varphi_{13}$ | $0.63_{\pm 0.50}$ | $0.00_{\pm 0.00}$ | $\mathbf{1.00}_{\pm 0.00}$ |
|  | $\varphi_{14}$ | $0.71_{\pm 0.40}$ | $0.73_{\pm 0.41}$ | $\mathbf{0.98}_{\pm 0.01}$ |
|  | $\varphi_{15}$ | $0.07_{\pm 0.02}$ | $0.73_{\pm 0.03}$ | $\mathbf{0.86}_{\pm 0.01}$ |
|  | $\varphi_{16}$ | $0.56_{\pm 0.35}$ | $0.64_{\pm 0.08}$ | $\mathbf{1.00}_{\pm 0.01}$ |

# Further resources

**Website**: deep-ltl.github.io

**arXiv**: arxiv.org/abs/2410.04631

**GitHub**: mathiasj33/deep-ltl

💻 mathias-jackermeier.me

✉ mathias.jackermeier
@cs.ox.ac.uk

mathiasj33

@m_jackermeier