

Solving temporal logic tasks specified by regular languages without further learning by composing skill primitives



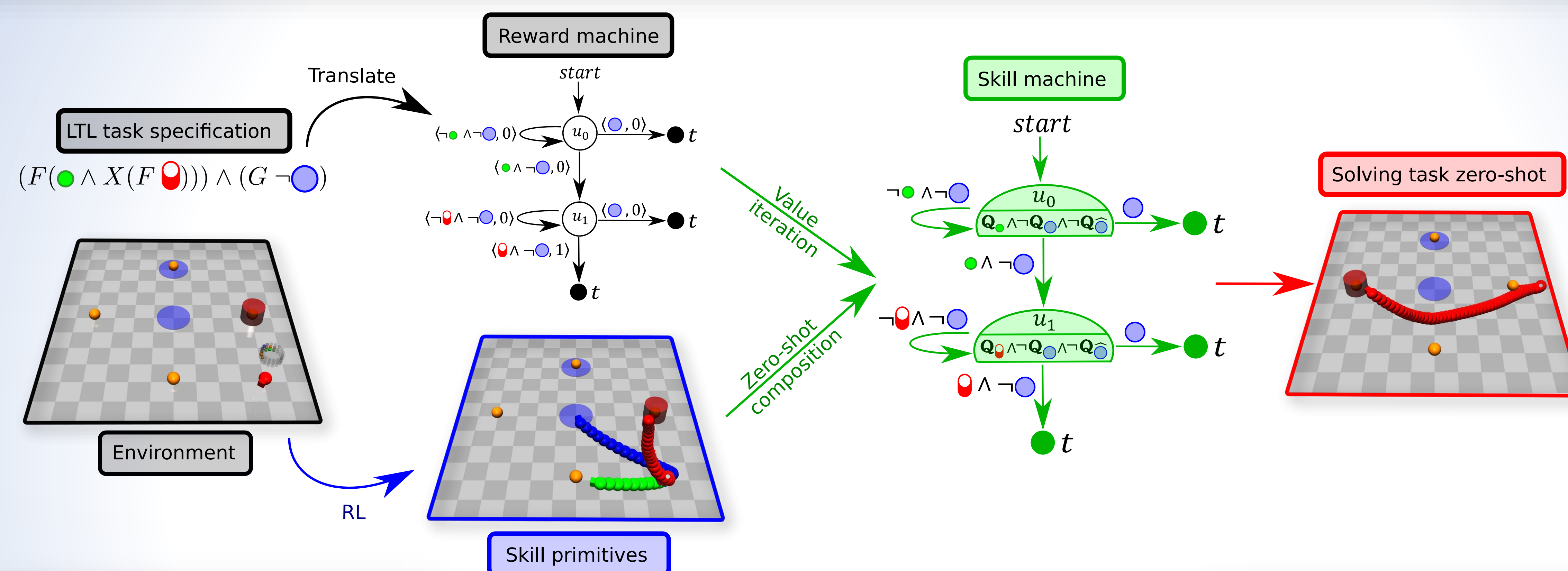
PAPER

What we want

We want agents (e.g robots) that can **solve multiple complex tasks** in an environment (e.g real world) **without further learning**

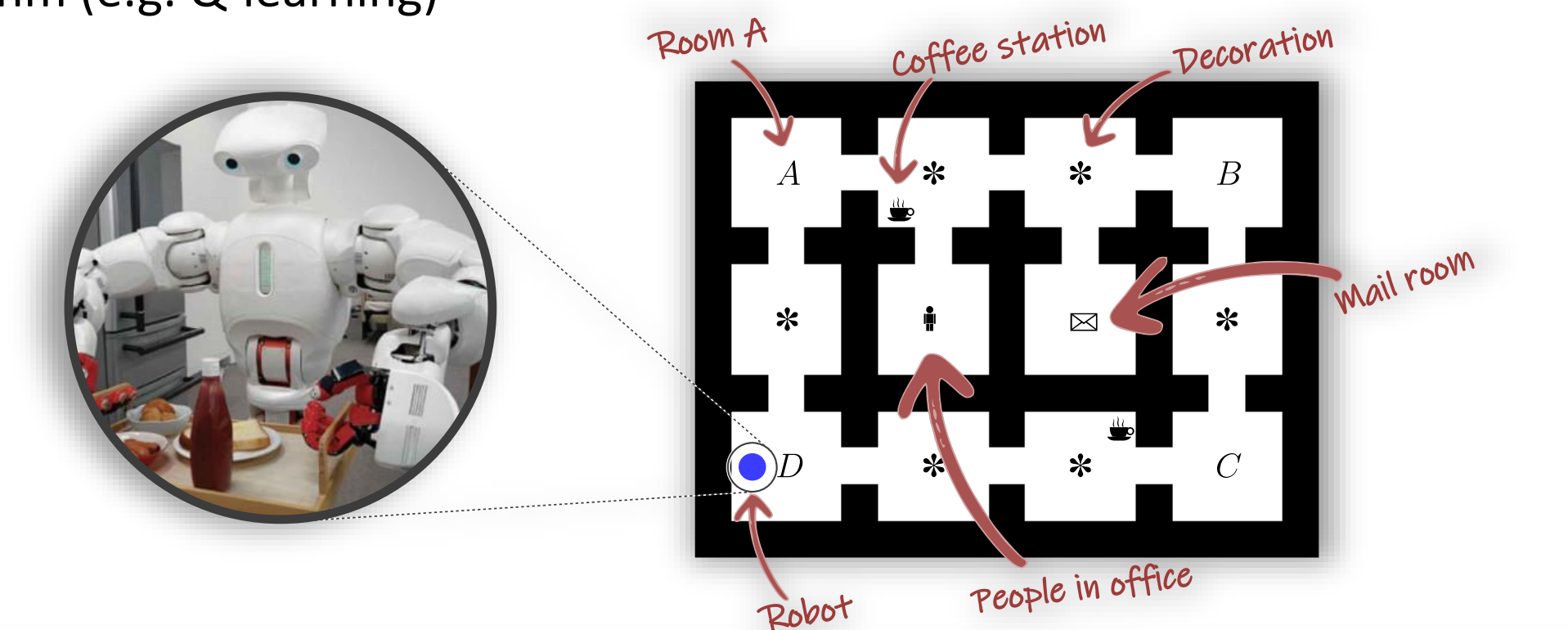
Defining tasks is hard!
Reward design requires expert knowledge

Learning is slow!
Sample efficiency in RL is terrible



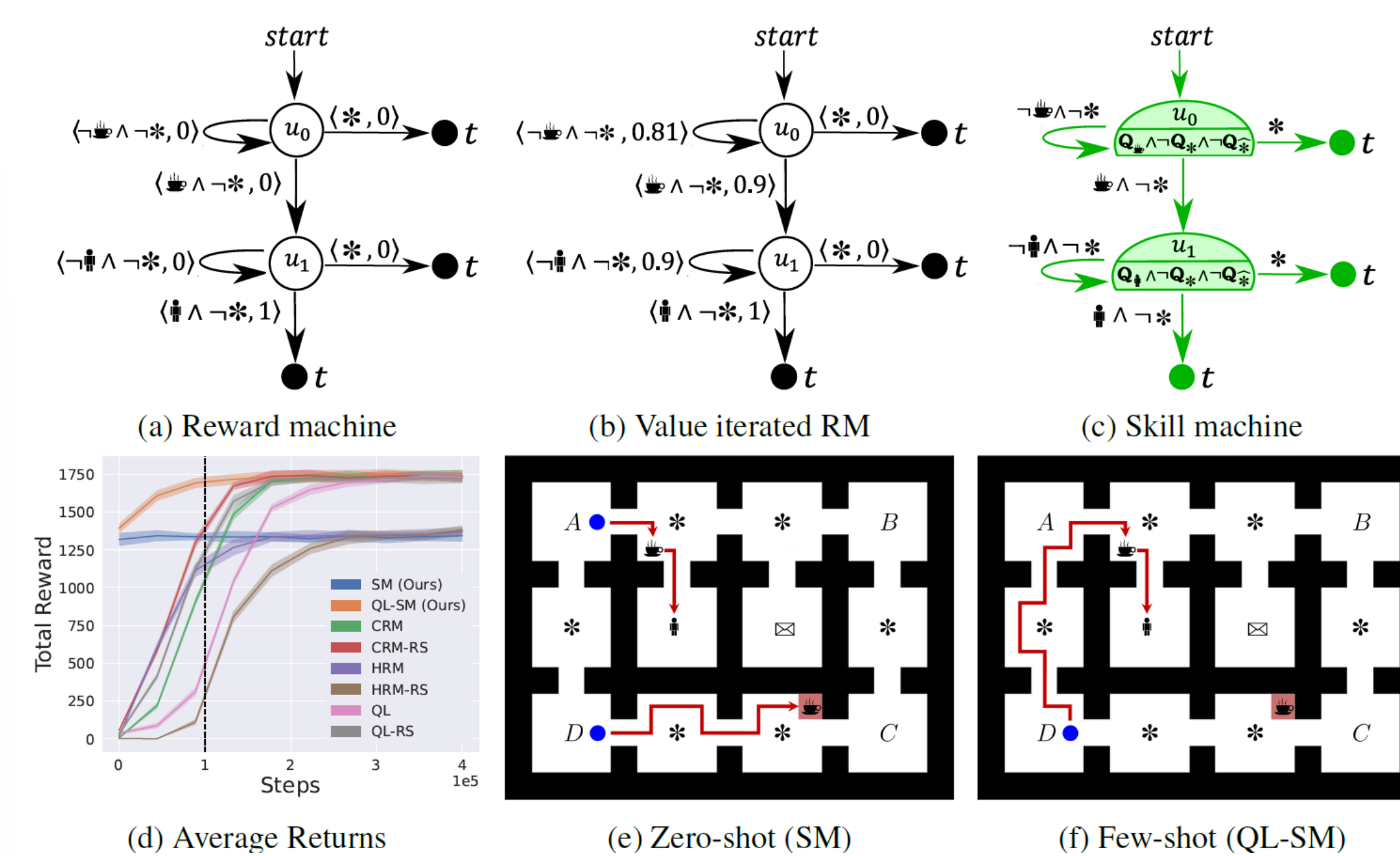
Few-shot transfer

We can improve the zero-shot policy to **optimality** by using any off-policy algorithm (e.g. Q-learning)



When reachability doesn't hold

"Deliver coffee to the office and never break a decoration"
 $(F(\text{coffee} \wedge X(F \text{ office}))) \wedge (G \neg \text{decoration})$



Curses of dimensionality

In general (for P environment propositions/sub-goals):

- Spatial curse:** Eventually 2^{2^P}
- Temporal curse:** Always 2^{2^P}

Task: "Navigate to a button and then to a cylinder while never entering blue regions"

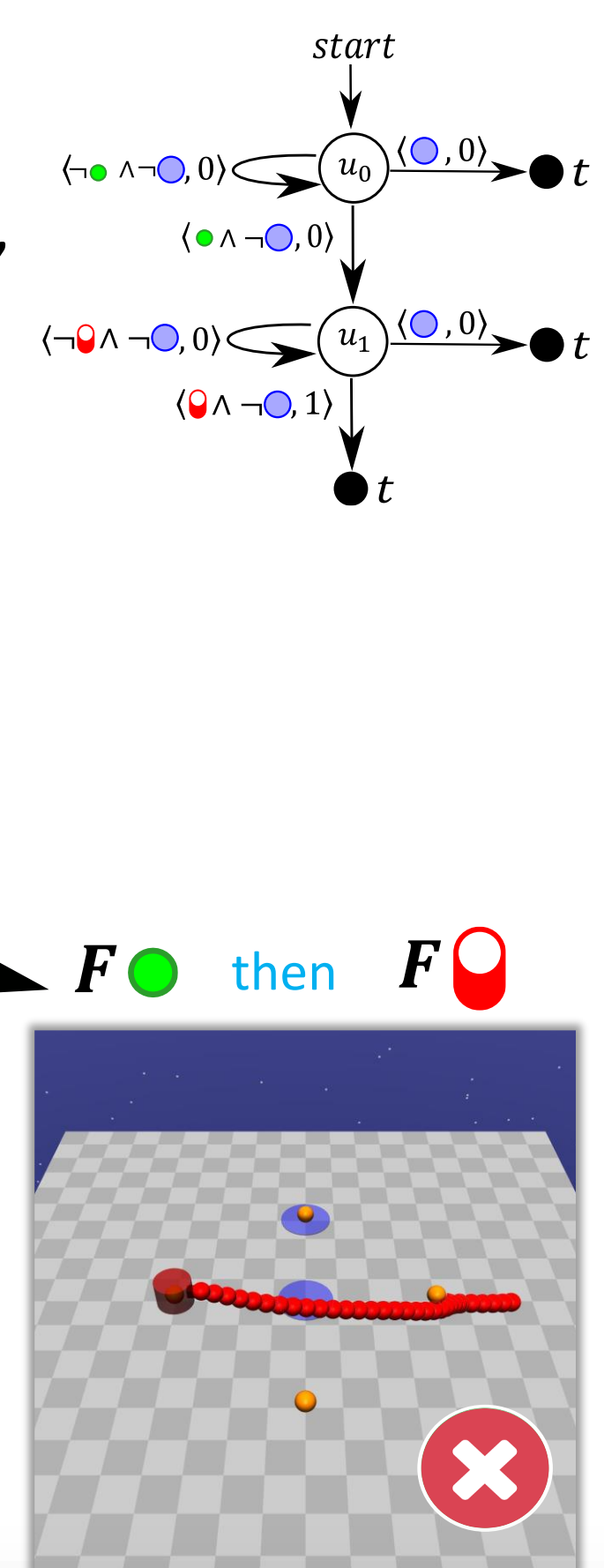
LTL: $(F(\text{button} \wedge X(F \text{ cylinder}))) \wedge (G \neg \text{blue})$

Agent needs to:

1. Eventually $\text{button} \wedge \neg \text{blue}$
2. Then eventually $\text{cylinder} \wedge \neg \text{blue}$
3. while always $\neg \text{blue}$

Learned skills (Options)

- "Go to a button" $F \text{ button}$
- "Go to a blue region" $F \text{ blue}$
- "Go to the cylinder" $F \text{ cylinder}$



Solution: primitives

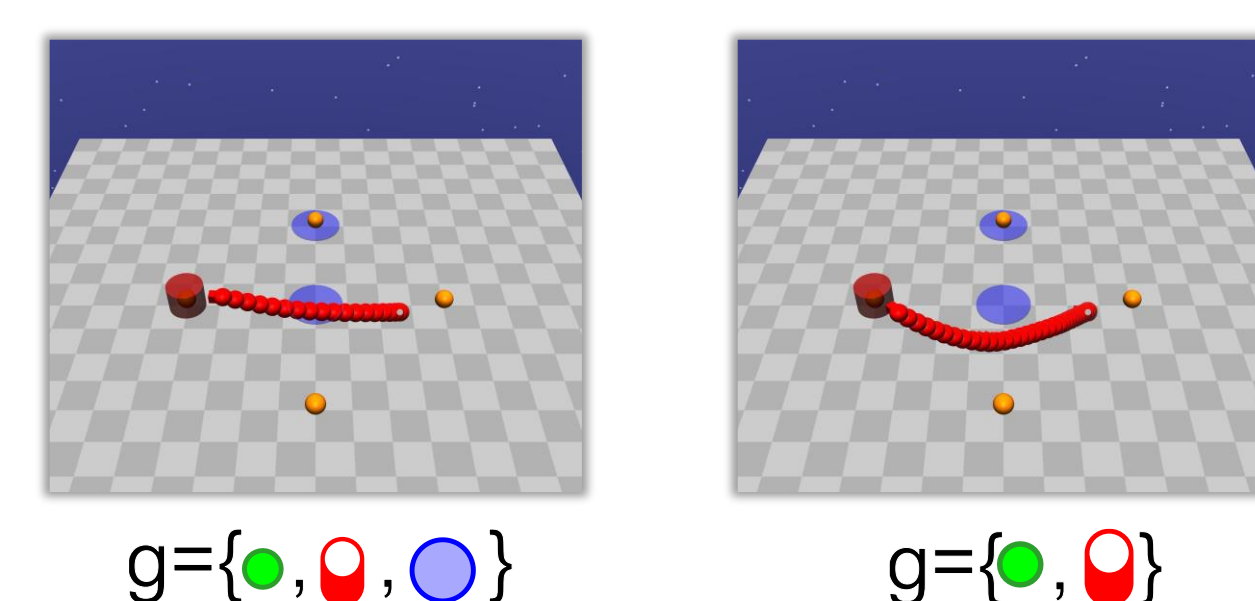
For **spatial curse**

- Task primitives (M_P): Composable MDPs for eventually satisfying each proposition
- Skill primitives (Q_P): WVF for each task primitive

For **temporal curse**

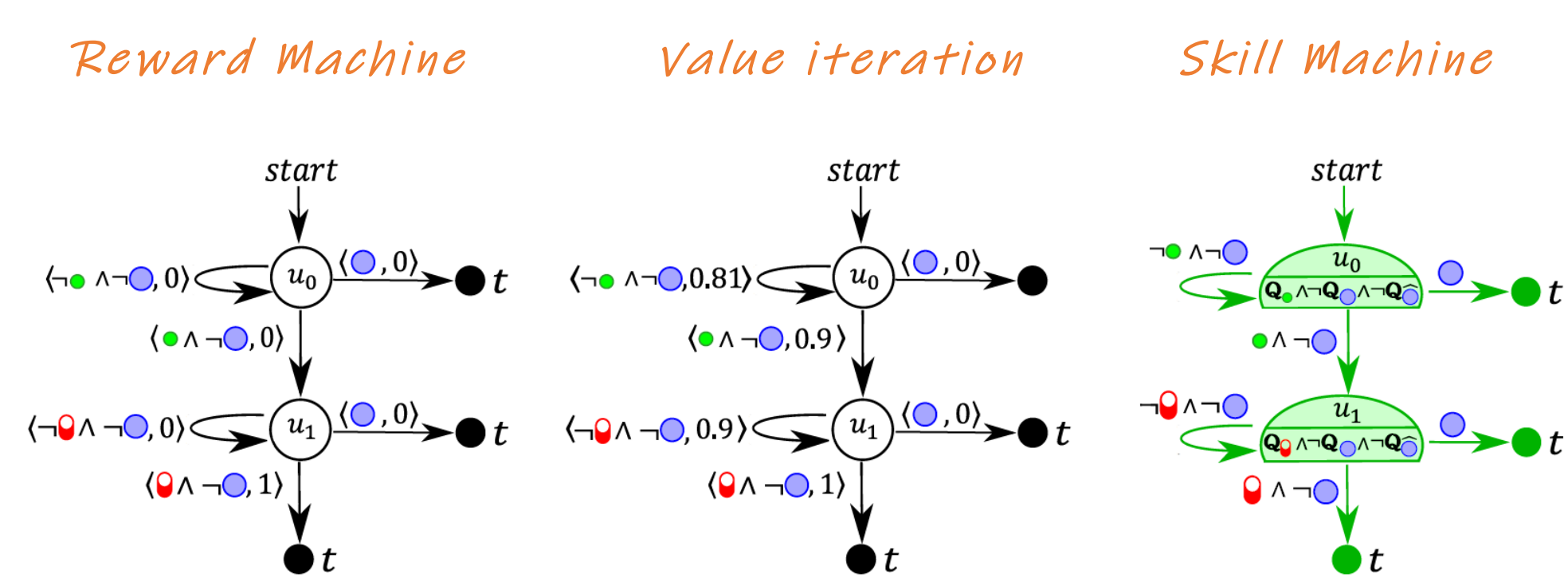
- Constraints (C): Propositions to always keep True or False in an episode
- Augments the state space of task primitives

E.g. M
where
 $C = \{\text{blue}\}$



Reward machine to skill machine

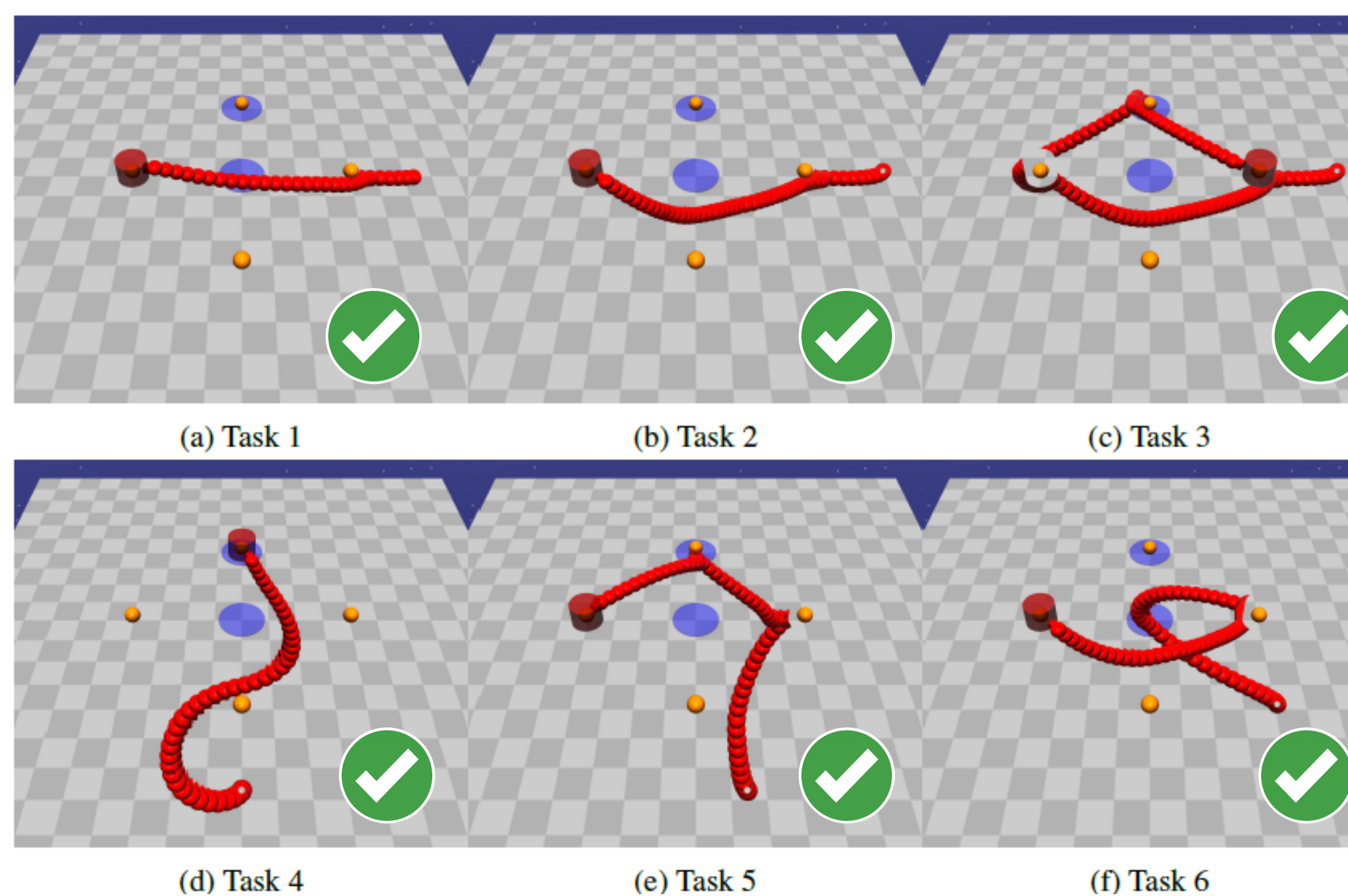
Plan over RM (e.g value iteration) then pick skills greedily



Zero-shot transfer

Task	Description — LTL
1	Navigate to a button and then to a cylinder. — $(F(\text{button} \wedge X(F \text{ cylinder})))$
2	Navigate to a button and then to a cylinder while never entering blue regions. — $(F(\text{button} \wedge X(F \text{ cylinder}))) \wedge (G \neg \text{blue})$
3	Navigate to a button, then to a cylinder without entering blue regions, then to a button inside a blue region, and finally to a cylinder again. — $F(\text{button} \wedge X(F(\text{button} \wedge X(F(\text{cylinder} \wedge \text{blue}))))$
4	Navigate to a button and then to a cylinder in a blue region. — $(F(\text{button} \wedge X(F \text{ cylinder}))) \wedge (F \text{ blue})$
5	Navigate to a cylinder, then to a button in a blue region, and finally to a cylinder again. — $(F \text{ cylinder} \wedge X(F(\text{button} \wedge X(F \text{ cylinder})))) \wedge (F \text{ blue})$
6	Navigate to a blue region, then to a button with a cylinder, and finally to a cylinder while avoiding blue regions. — $(F \text{ blue} \wedge X(F(\text{button} \wedge X(F \text{ cylinder})))) \wedge (G \neg \text{blue})$

Table A4: Tasks in the Safety Gym domains. The RMs are generated from the LTL expressions.



When choosing skills greedily is bad

"Go to the office but not until you have coffee"
 $(F \text{ office}) \wedge (\neg \text{office} \text{ until } F \text{ coffee})$

