

"Scalable" Synthesis of Robust Strategies for Uncertain POMDPs

Murat Cubuktepe, Ufuk Topcu

The University of Texas at Austin



aUTonomous
SYSTEMS GROUP

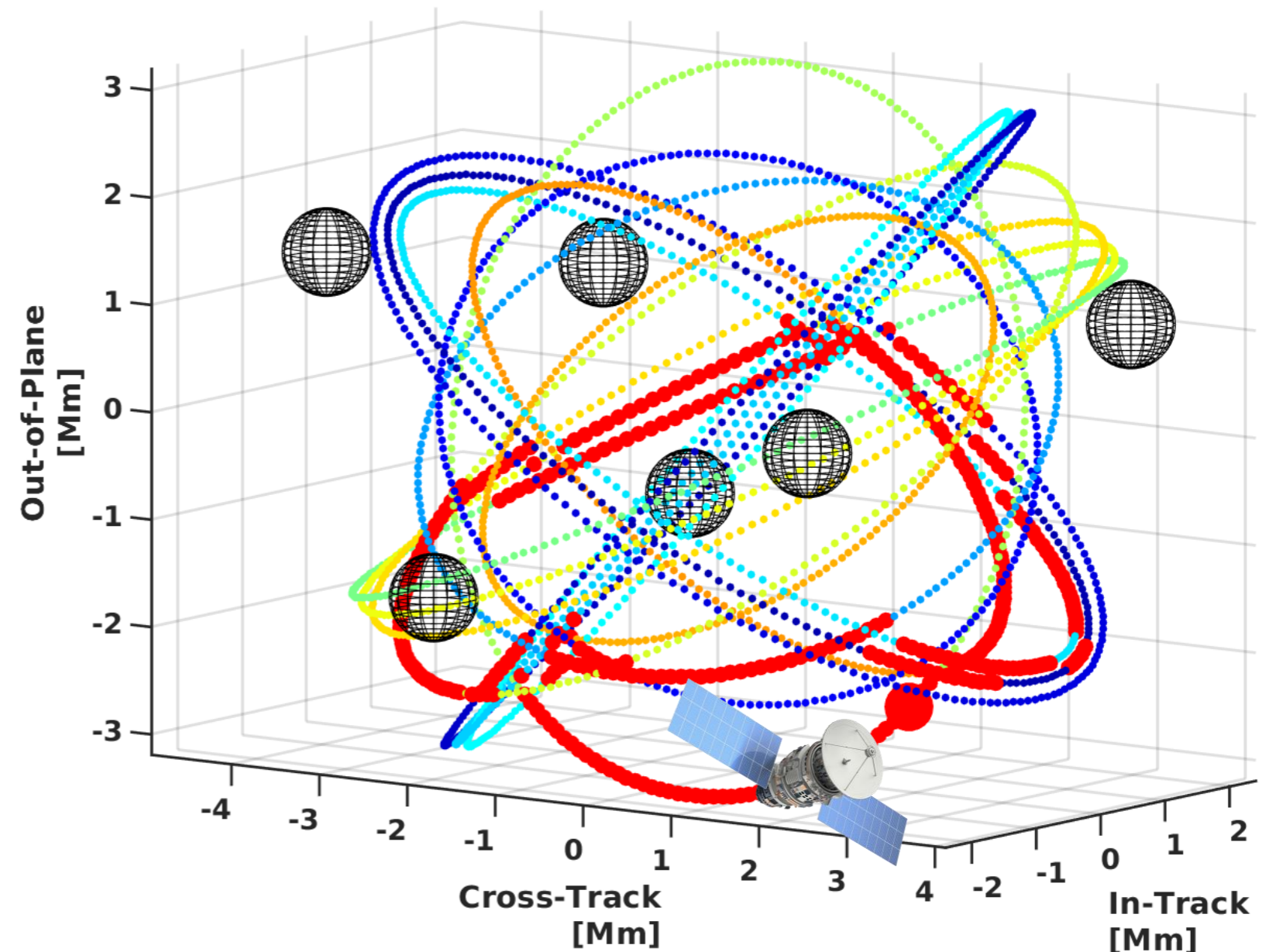
What is in this talk?

From **POMDPs** to
uncertain **POMDPs**.

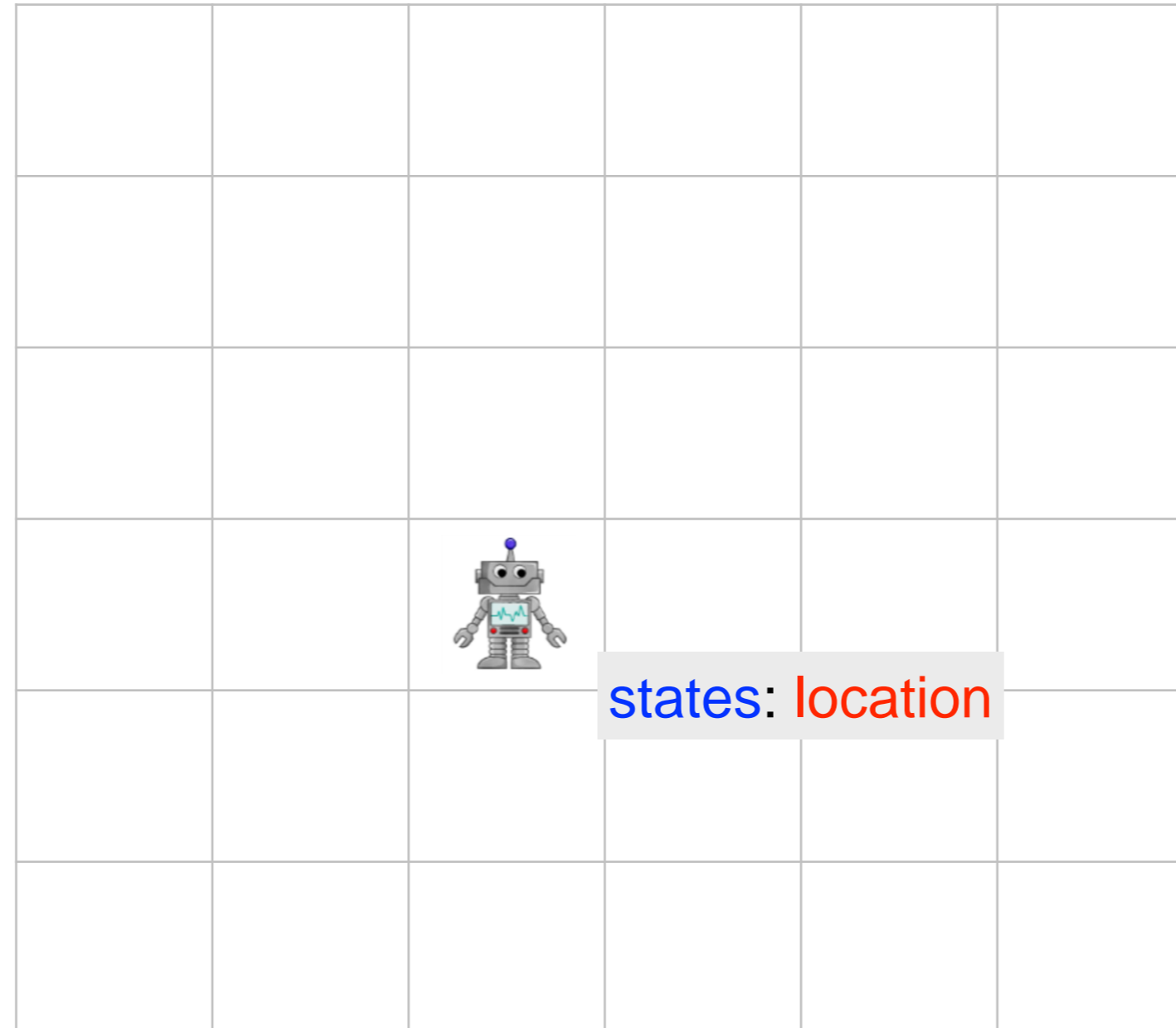
POMDPs are hard.
Uncertain POMDPs are harder.

We seem to have developed an
approach that scales.

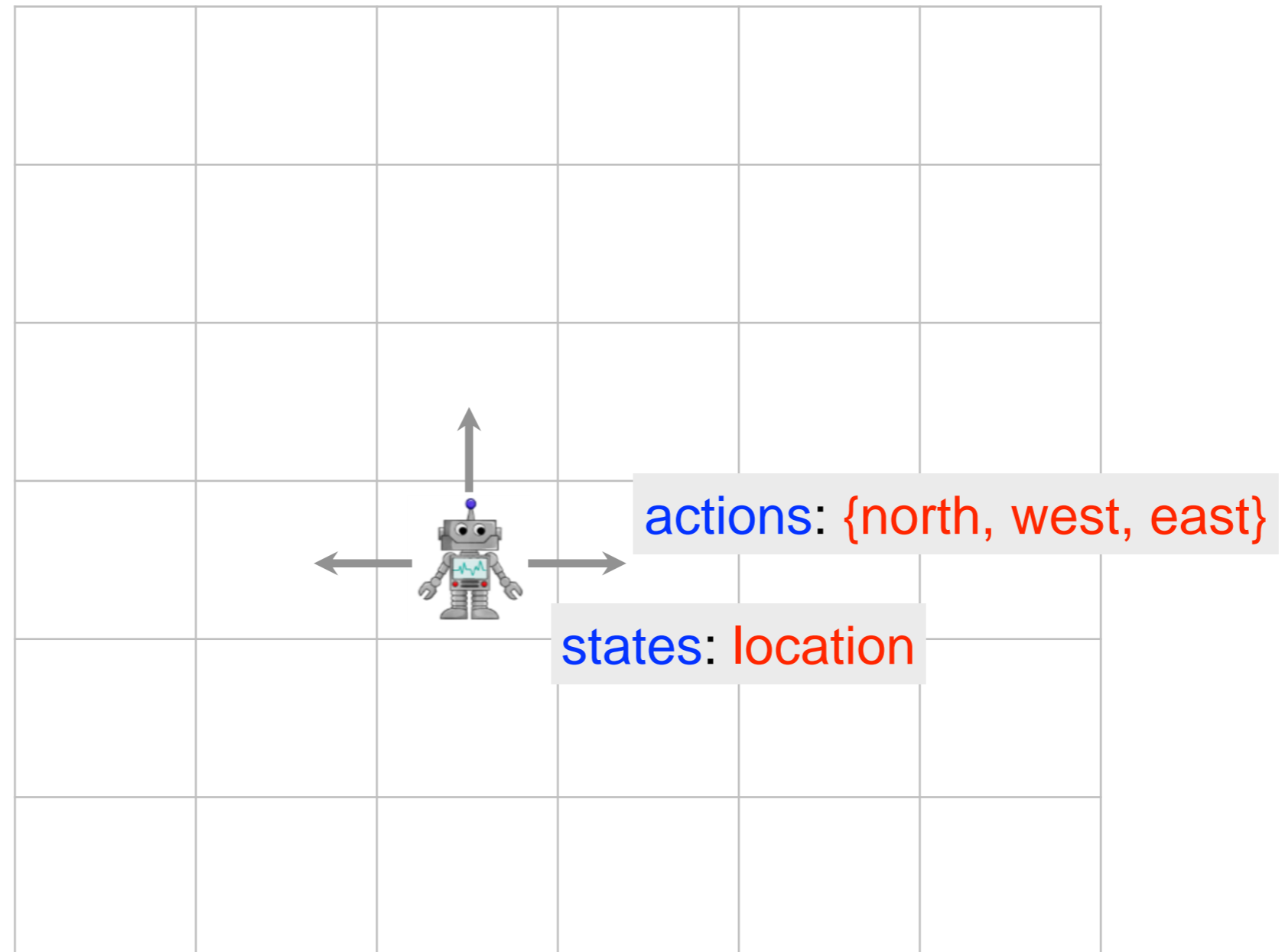
Demonstrated on spacecraft
motion planning.



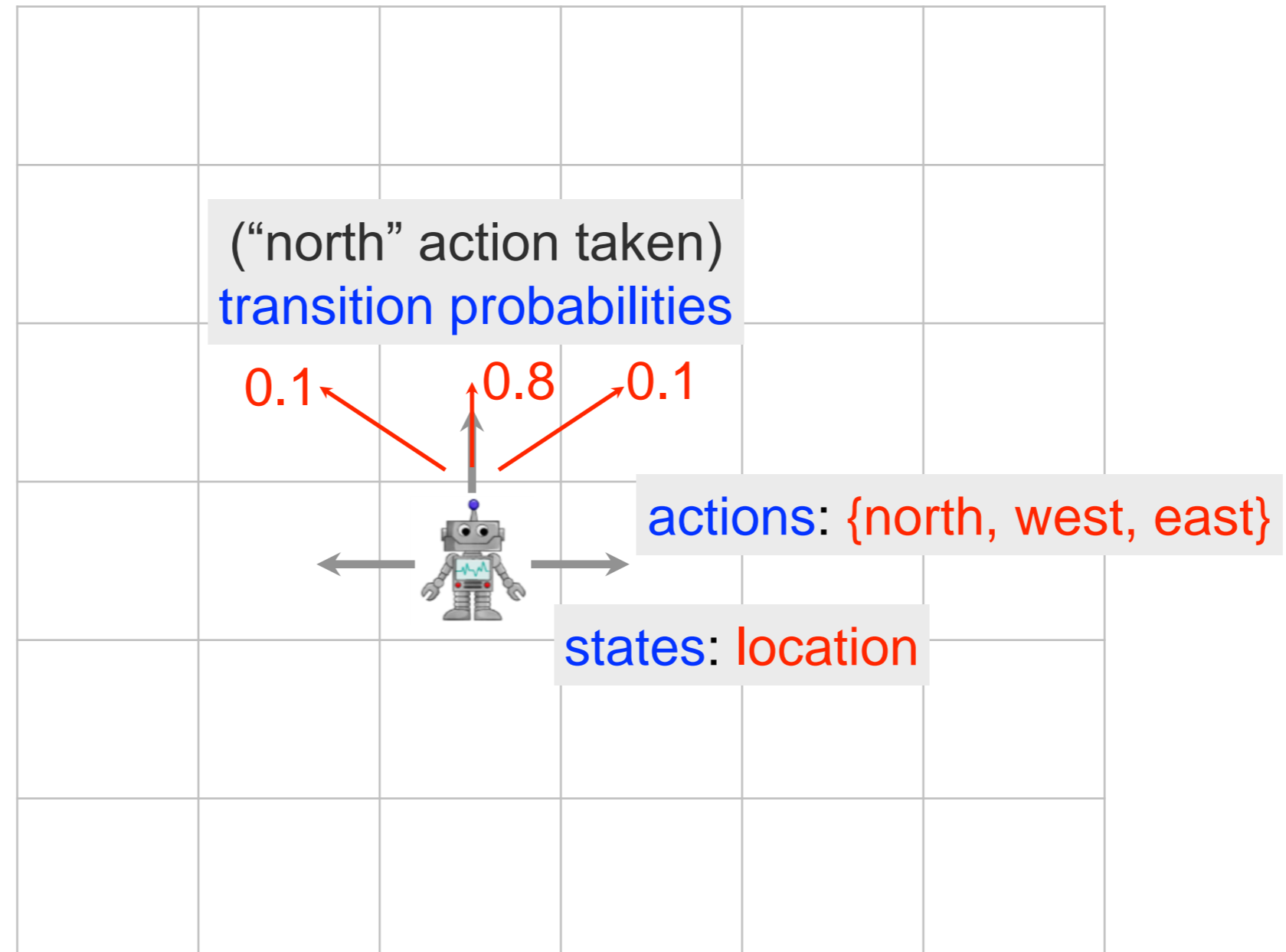
POMDPs through **an example**



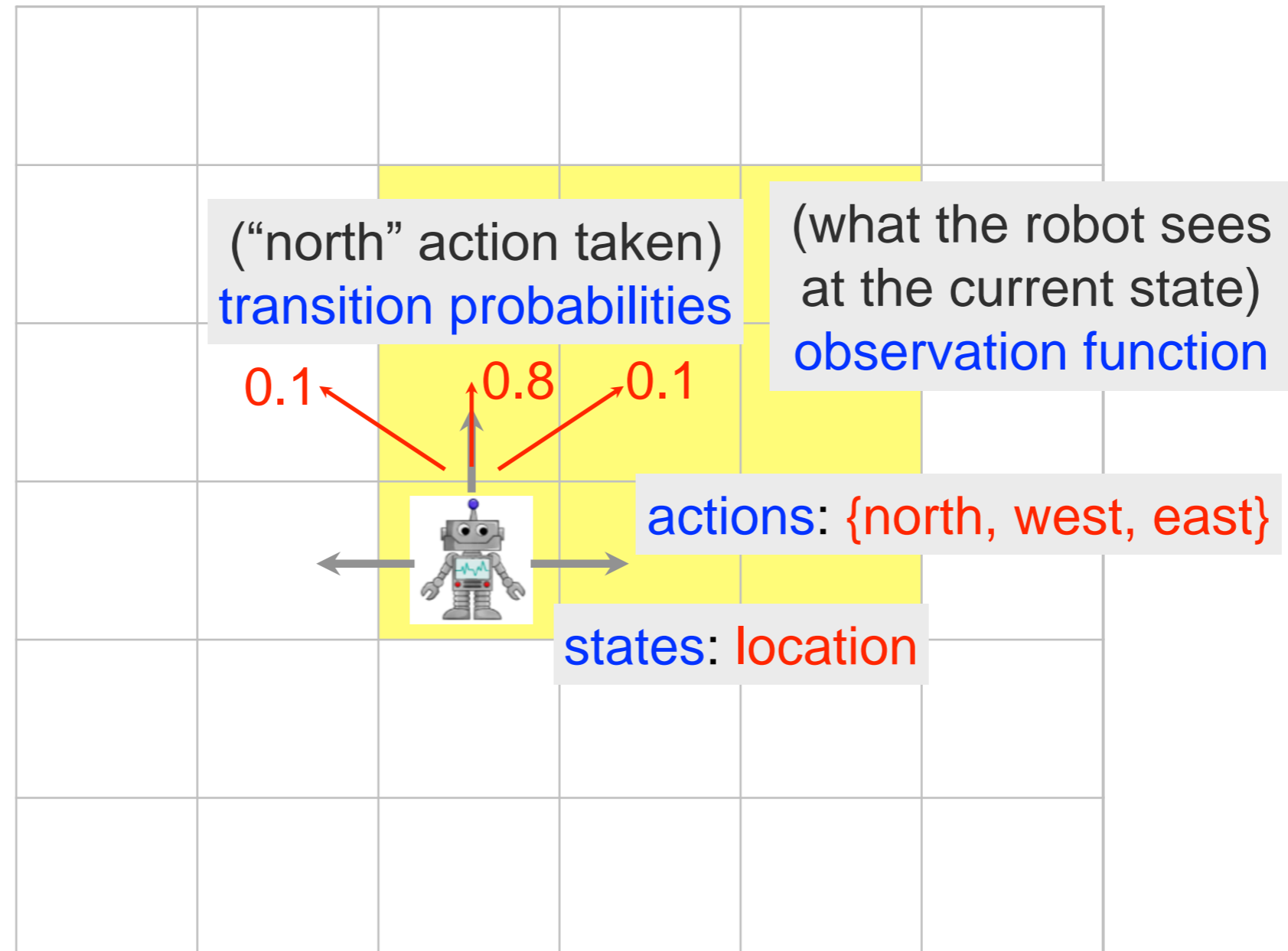
POMDPs through **an example**



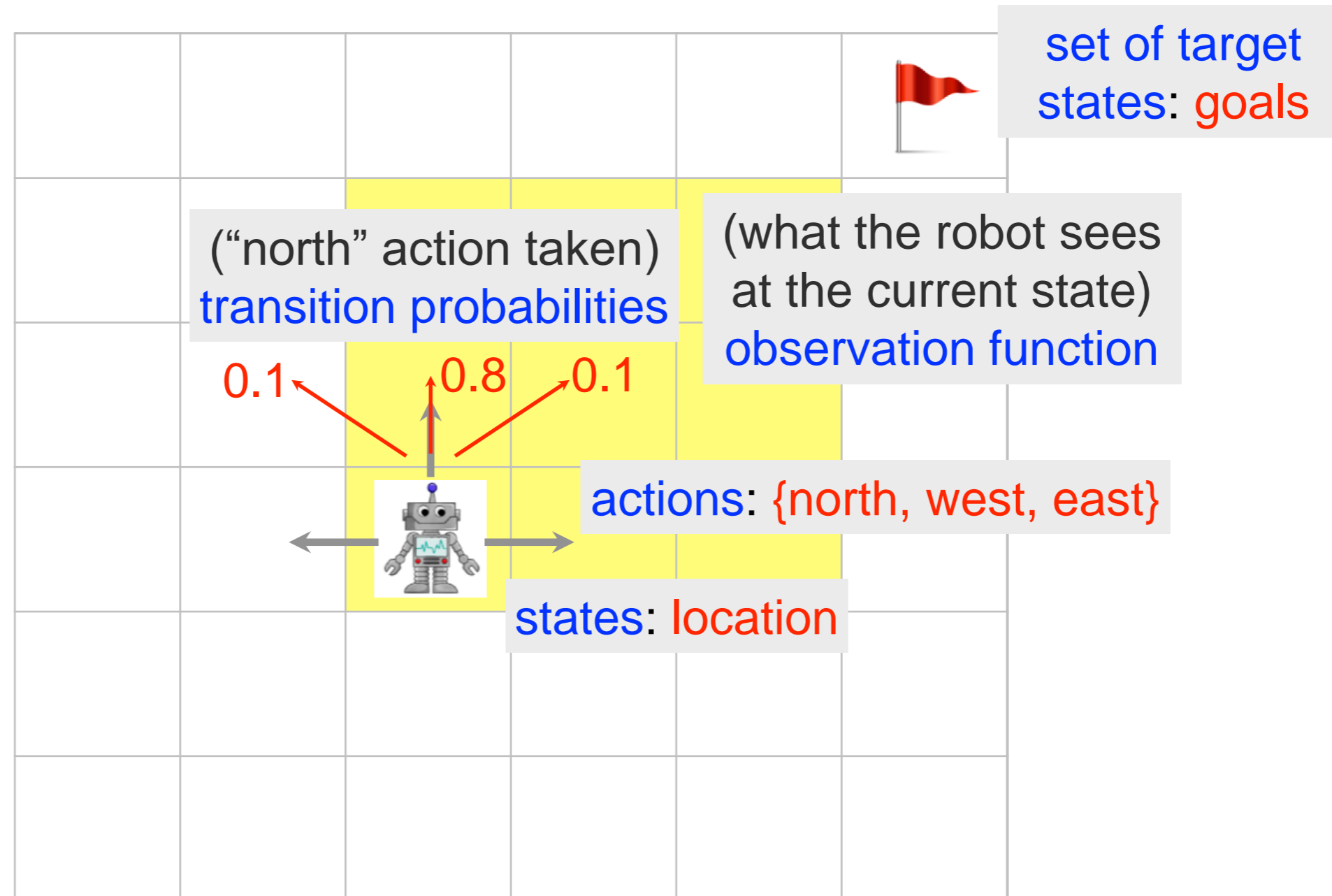
POMDPs through an example



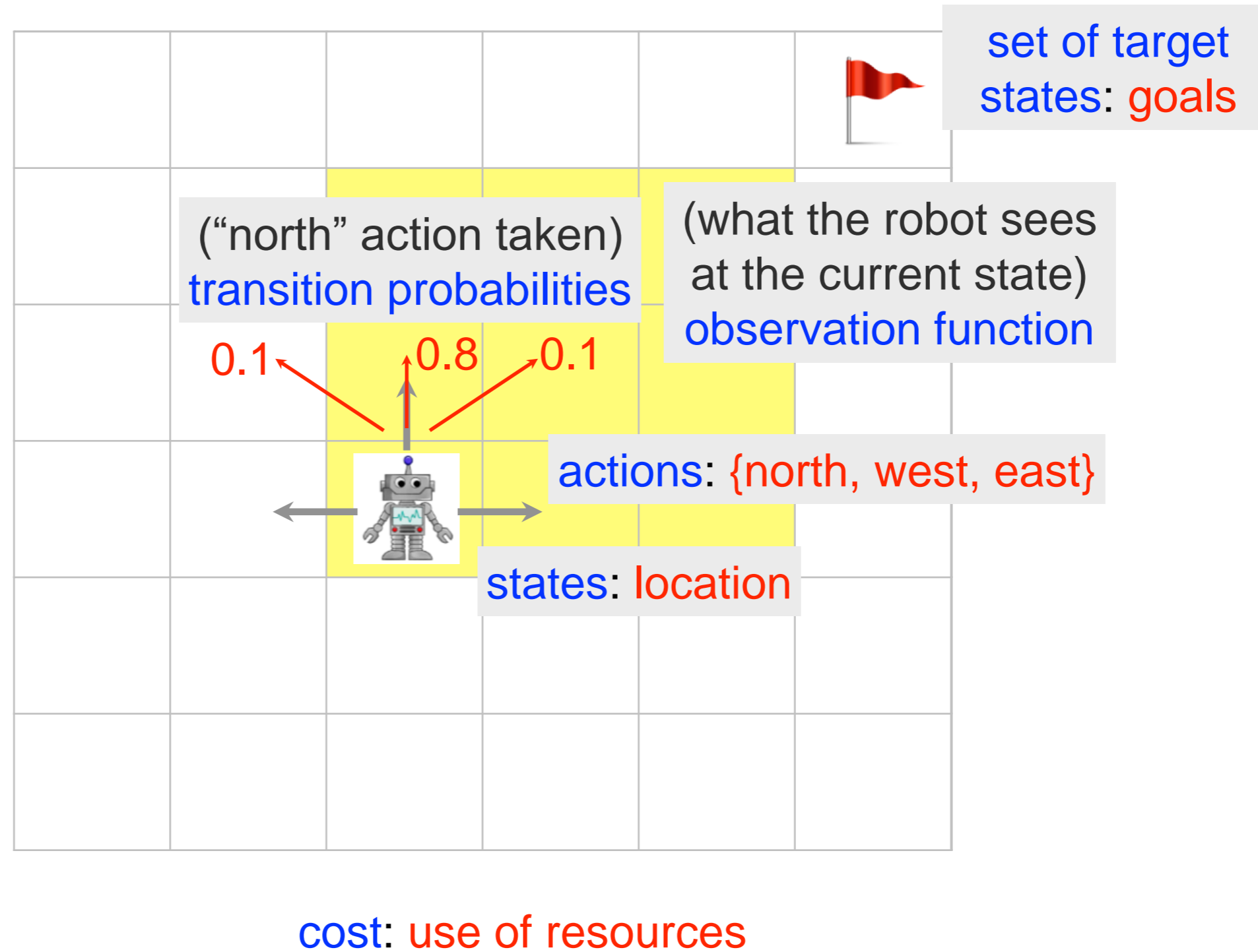
POMDPs through an example



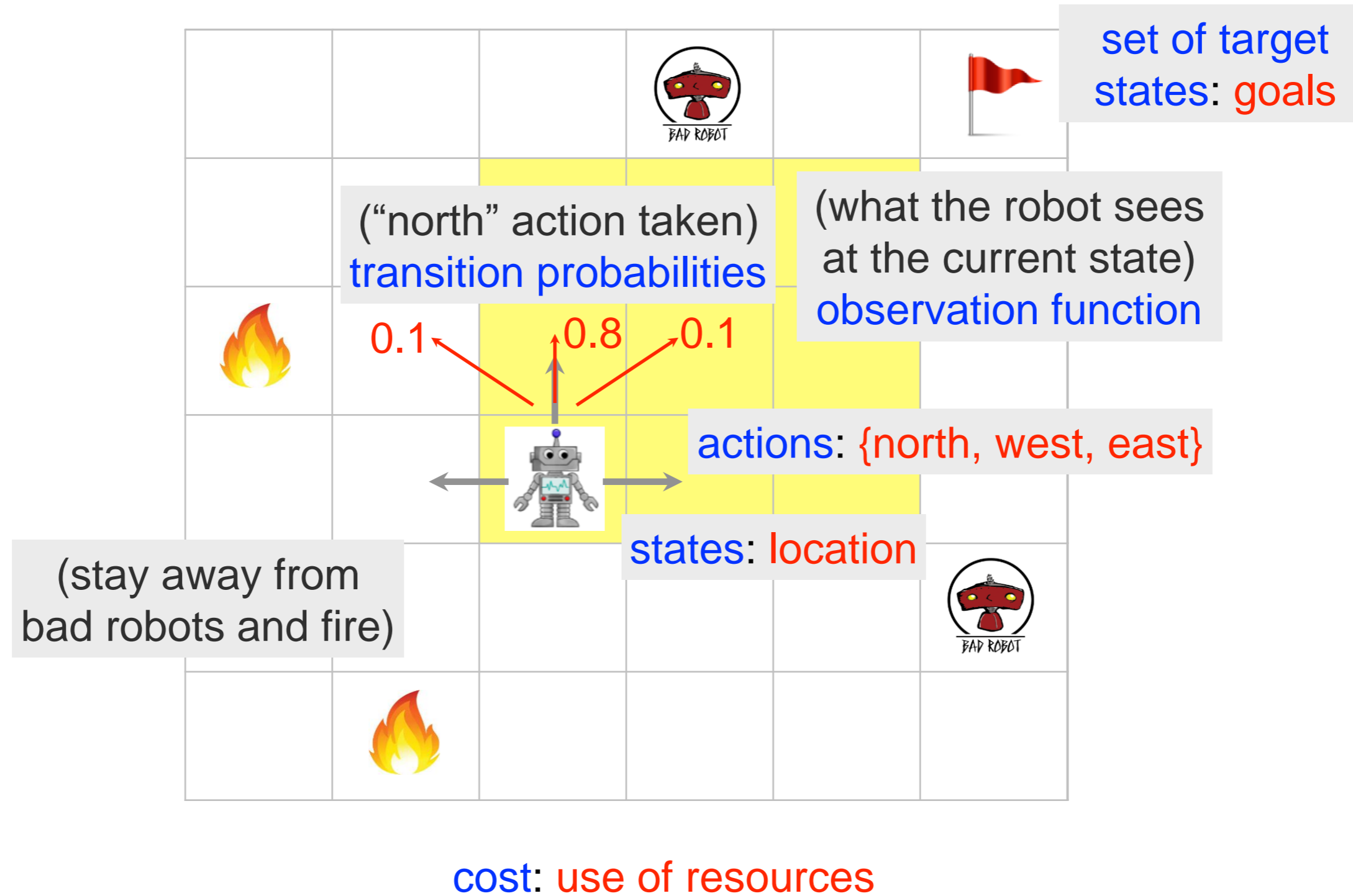
POMDPs through an example



POMDPs through an example



POMDPs through an example



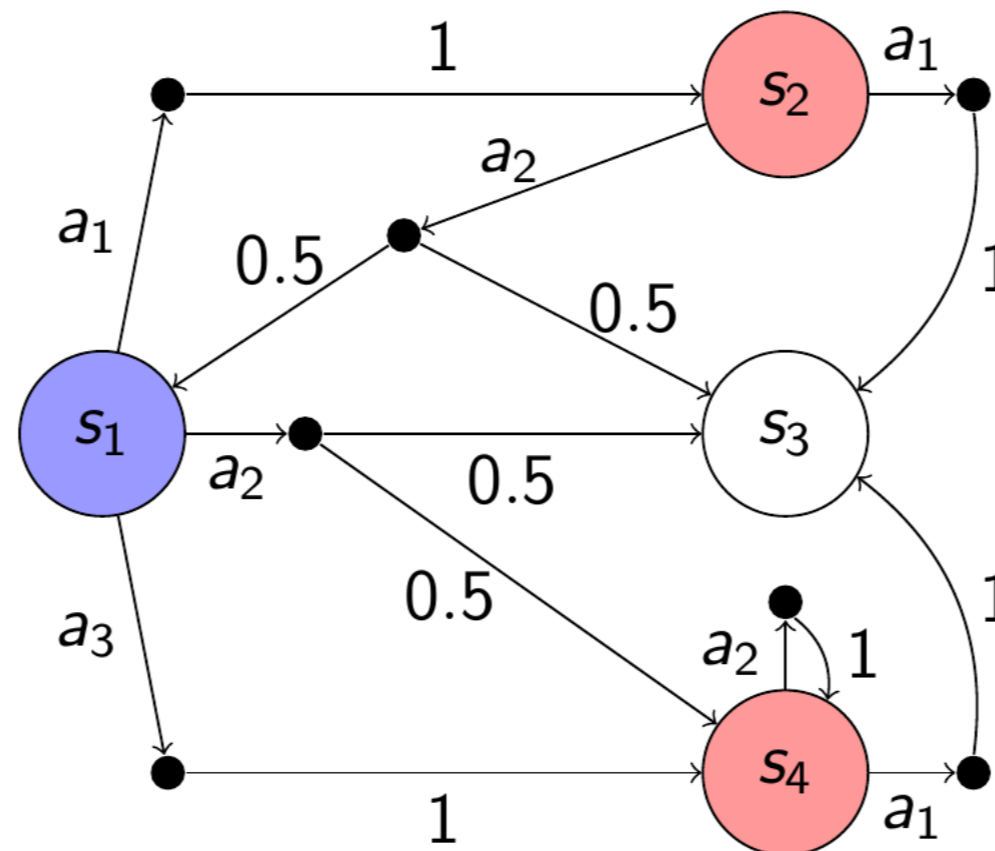
Partially Observable Markov Decision Processes

$$\mathcal{M} = (\underbrace{S}_{\text{set of states}}, \underbrace{A}_{\text{set of actions}}, \underbrace{P}_{\text{probabilistic transition function}}, \underbrace{Z}_{\text{set of observations}}, \underbrace{O}_{\text{observation function}}, \underbrace{C}_{\text{cost function}})$$

$$S = \{s_1, s_2, s_3, s_4\}$$

$$A = \{a_1, a_2, a_3\}$$

$$P(s_2, a_2, s_1) = 0.5$$



$$Z = \{ \text{red circle}, \text{blue circle}, \text{white circle} \}$$

$$O(s_4) = \text{red circle}$$

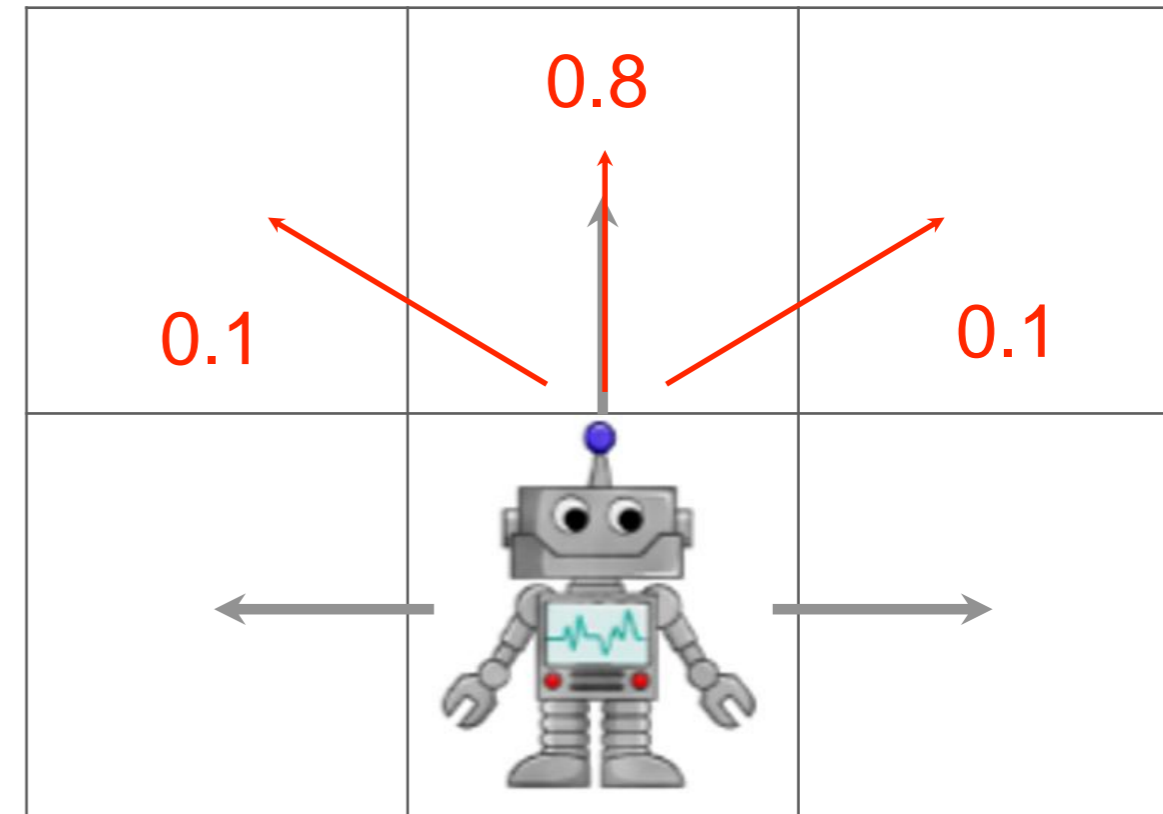
$$C(s_1, a_2) = 3$$

Uncertain POMDPs (uPOMDPs)

In a POMDP \mathcal{M} , transition function is assumed to be known

A uPOMDP $\mathcal{M}^{\mathcal{P}}$ extends a POMDP by **allowing for probability intervals**

transition function $P \in \mathcal{P}$ uncertainty set

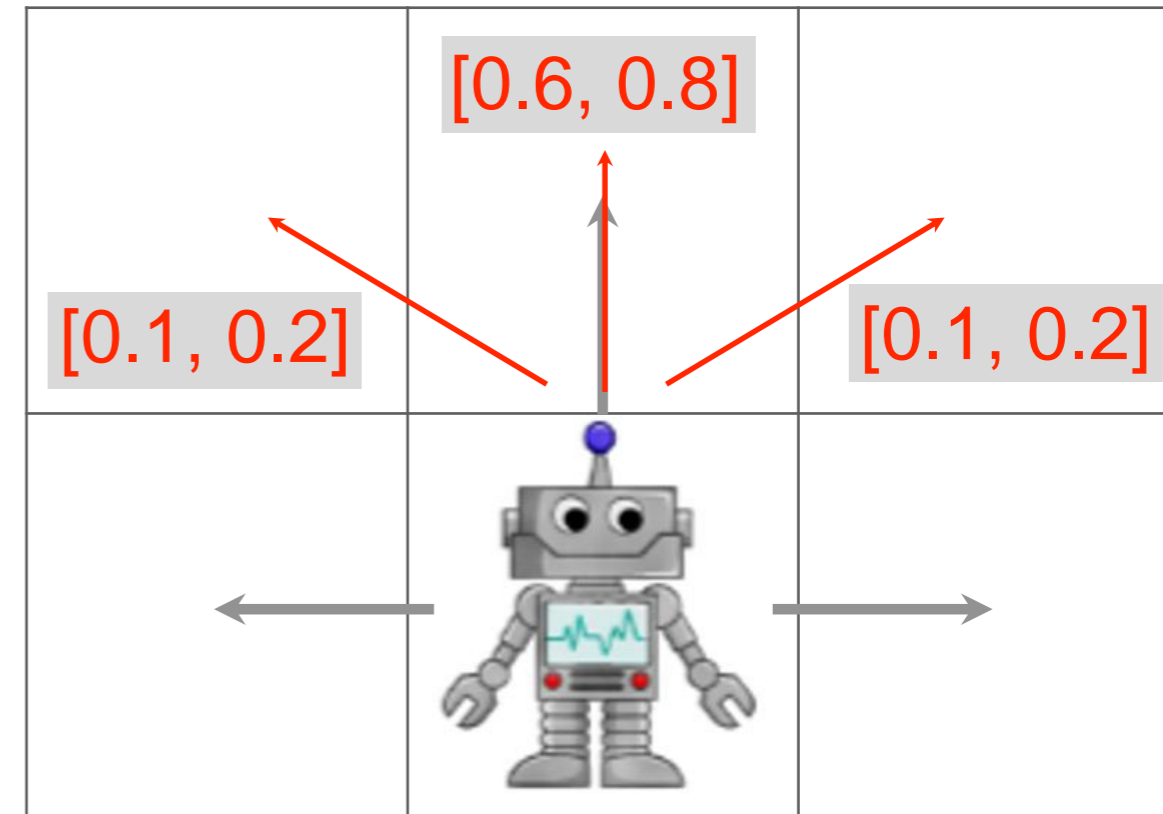


Uncertain POMDPs (uPOMDPs)

In a POMDP \mathcal{M} , transition function is assumed to be known

A uPOMDP $\mathcal{M}^{\mathcal{P}}$ extends a POMDP by **allowing for probability intervals**

transition function $P \in \mathcal{P}$ uncertainty set



Specifications

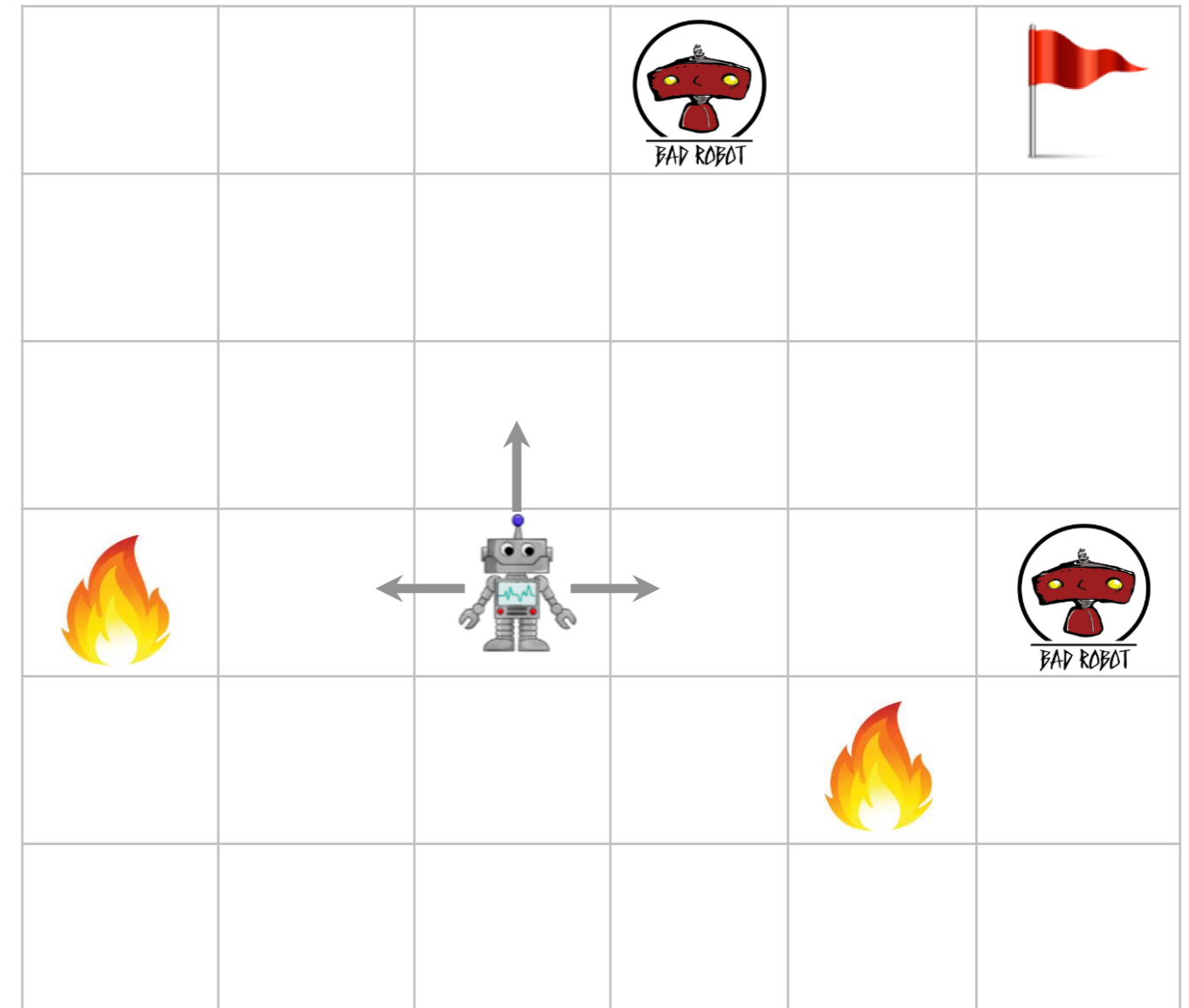
Reachability specification

$$\varphi_r = \mathbb{P}_{\geq \lambda}(\diamond T)$$

the probability of reaching a state in the target set T is greater than λ

reach the flag
while avoiding fire
and bad robots

$$T = \{ \text{flag} \}$$



Specifications

Reachability specification

$$\varphi_r = \mathbb{P}_{\geq \lambda}(\diamond T)$$

the probability of reaching a state in the target set T is greater than λ

reach the flag while avoiding fire and bad robots

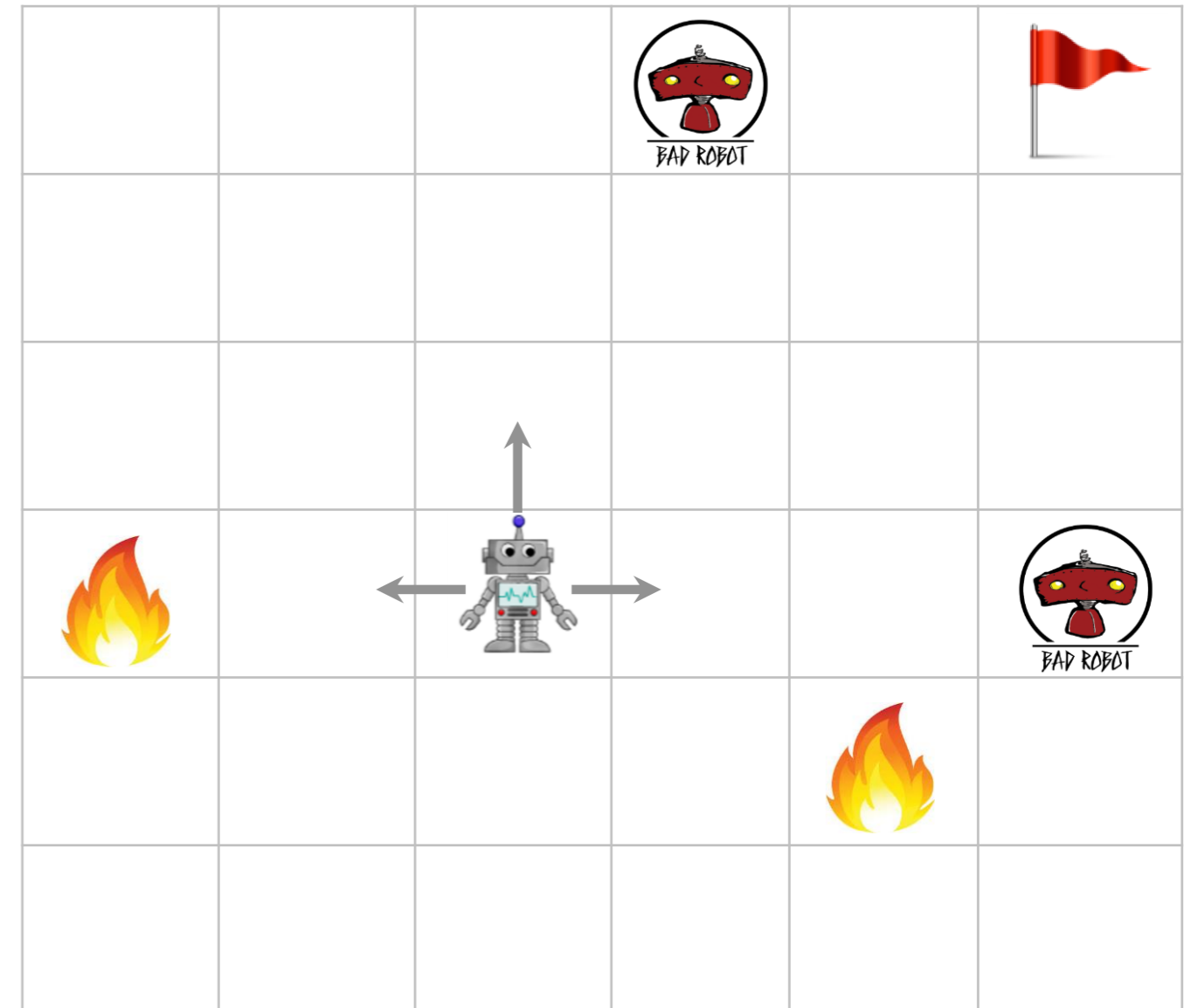
Performance specification

$$\varphi_p = \mathbb{E}_{\leq \kappa}(\diamond T)$$

the expected accumulated reward before reaching a state in T is less than κ

minimize use of resources before reaching the flag

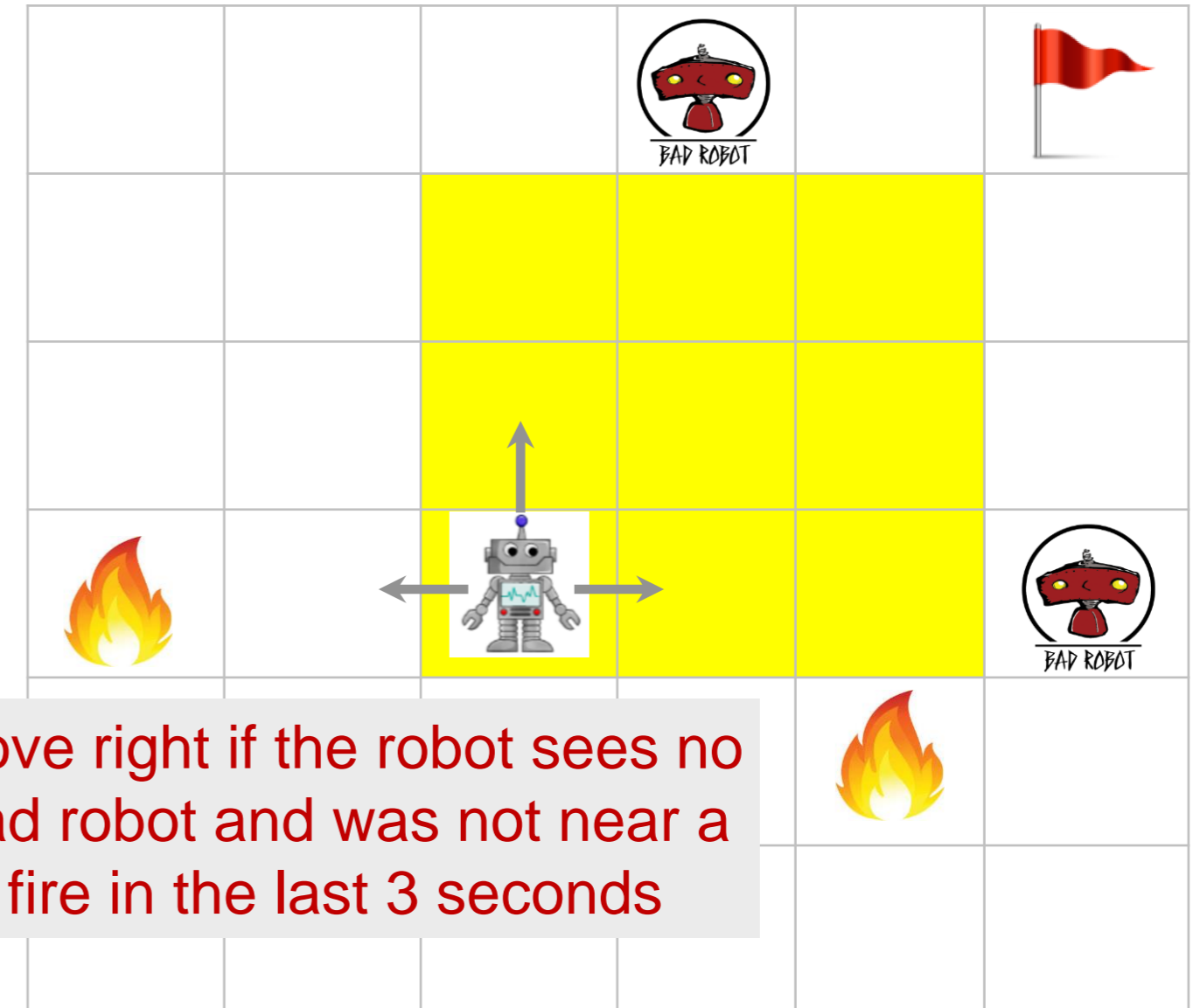
$$T = \{ \text{flag} \}$$



Policies

A policy σ for a uPOMDP maps sequences of observations and actions, to a distribution over actions

$$\sigma : \underbrace{(Z \times A)^*}_{\text{memory}} \times \underbrace{Z}_{\text{current observation}} \rightarrow \underbrace{\text{Distr}(A)}_{\text{distribution over actions}}$$



Synthesis of Robust Policies in uPOMDPs

Given a uPOMDP $\mathcal{M}^{\mathcal{P}}$, compute a policy σ such that

$$\begin{array}{c} \text{induced uncertain} \\ \text{Markov chain (uMC)} \end{array} \mathcal{M}_{\sigma}^{\mathcal{P}} \models \varphi \quad \begin{array}{c} \text{transition} \\ \text{function} \end{array} \quad \text{for all } P \in \mathcal{P} \quad \begin{array}{c} \text{uncertainty} \\ \text{set} \end{array}$$

satisfies the specification

Policy synthesis in uPOMDPs is hard(er)

Partial observability over the state of the agent makes policy synthesis in uPOMDPs computationally hard

- **Exponential** in the number of states, actions, and observations

Policy synthesis in uPOMDPs is hard(er)

Partial observability over the state of the agent makes policy synthesis in uPOMDPs computationally hard

- **Exponential** in the number of states, actions, and observations

Undecidable (policy requires infinite memory of observations)

Undecidable even when the transition function is known

The Main Ideas

Restriction to finite-memory policies
yields a decidable problem, **NP-hard** though

Dualization (of a **semi-infinite** optimization problem)
yields a **finite** (yet still **nonconvex**) problem

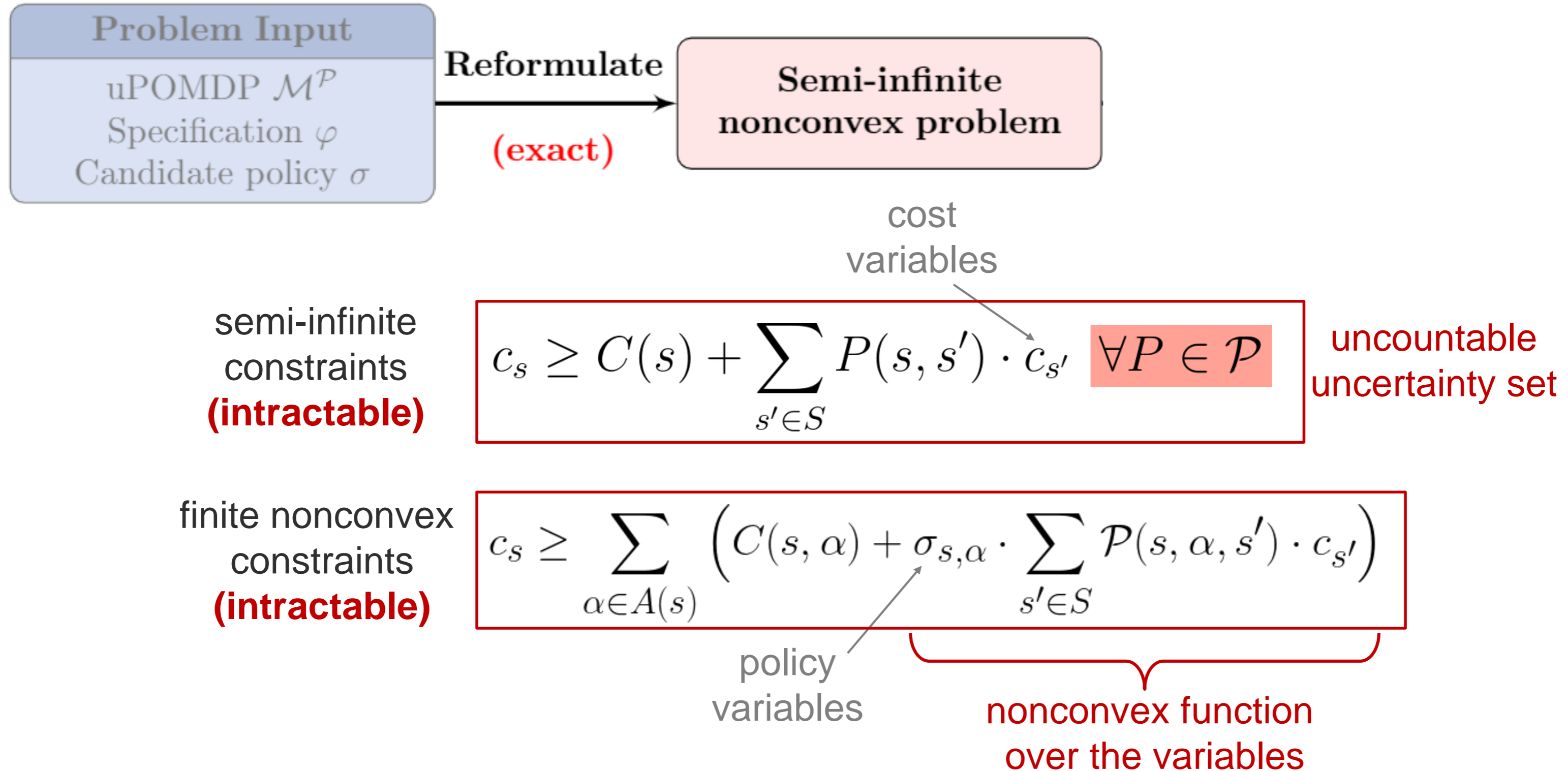
Linearization + verification
yields a **finite, convex** problem

Overview of the Algorithm

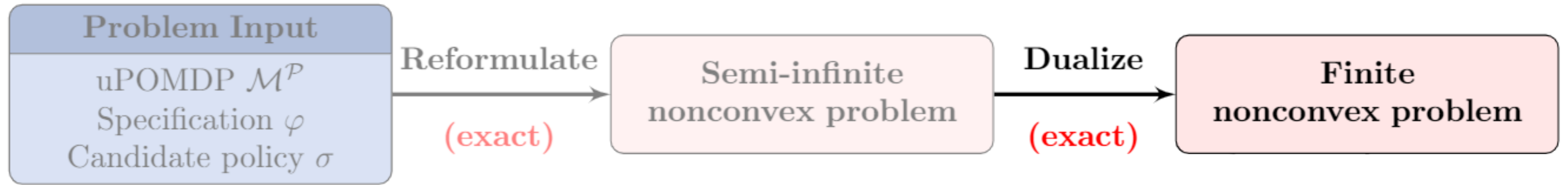
Problem Input

uPOMDP $\mathcal{M}^{\mathcal{P}}$
Specification φ
Candidate policy σ

Overview of the Algorithm



Overview of the Algorithm



finite affine
constraints
(tractable)

$$c_s \geq C(s) + \mu_s^\top g_s,$$

$$D_s^\top \mu_s + q = 0, \quad \mu_s \geq 0$$

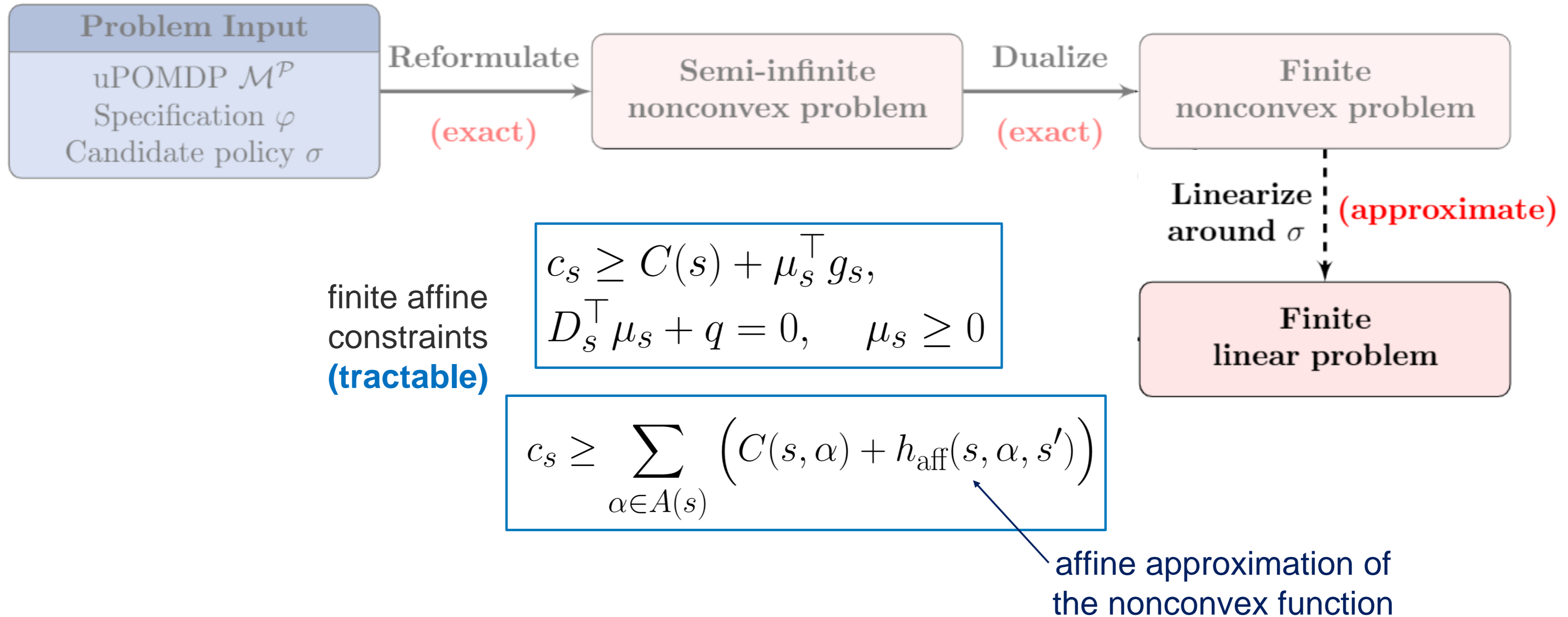
finite nonconvex
constraints
(intractable)

$$c_s \geq \sum_{\alpha \in A(s)} \left(C(s, \alpha) + \sigma_{s, \alpha} \cdot \sum_{s' \in S} \mathcal{P}(s, \alpha, s') \cdot c_{s'} \right)$$

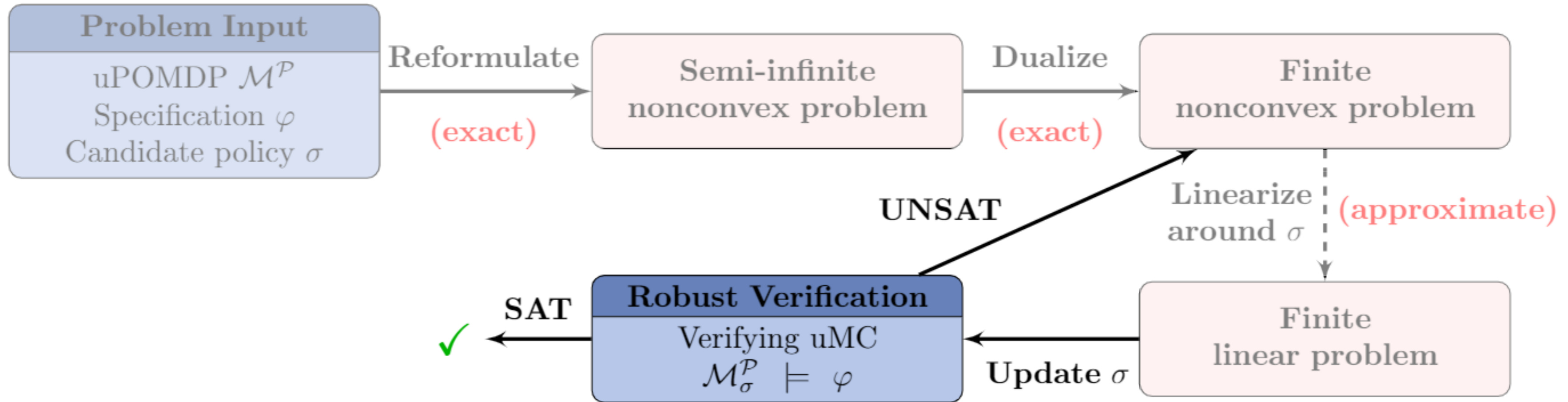
policy
variables

nonconvex function
over the variables

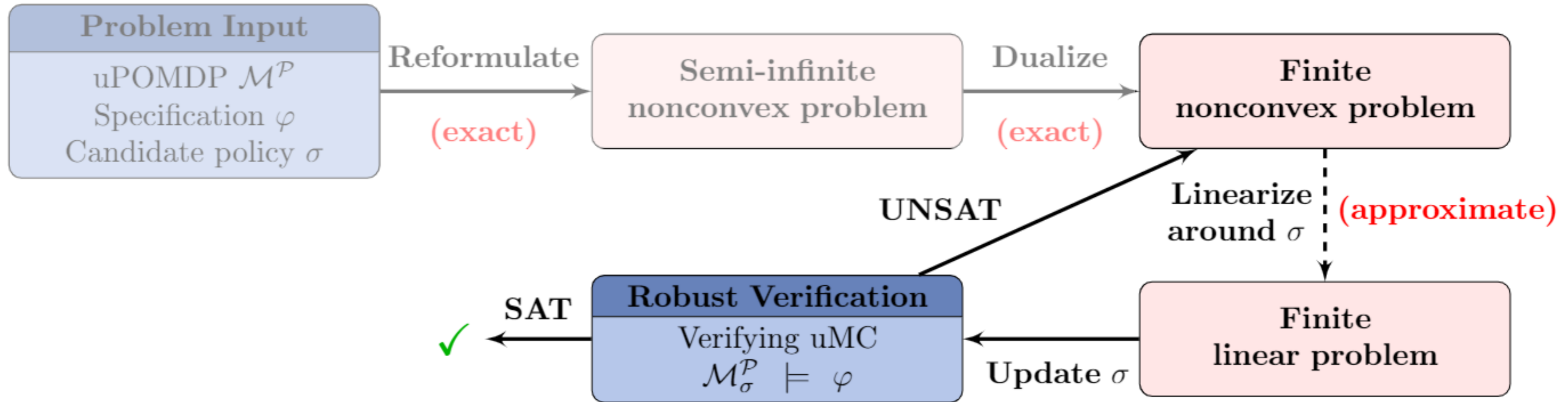
Overview of the Algorithm



Overview of the Algorithm



Overview of the Algorithm



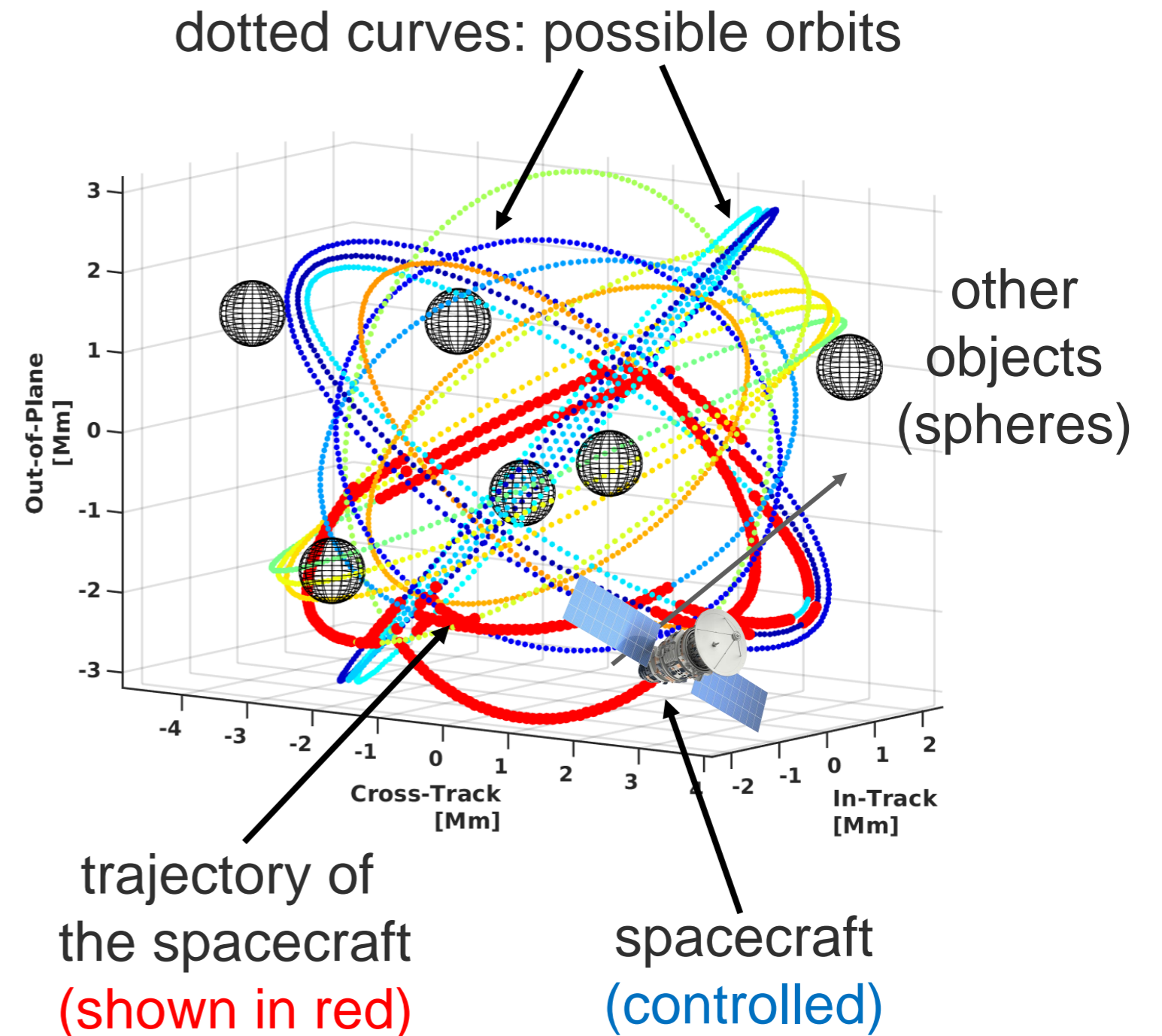
Spacecraft Motion Planning (with operator in the loop)



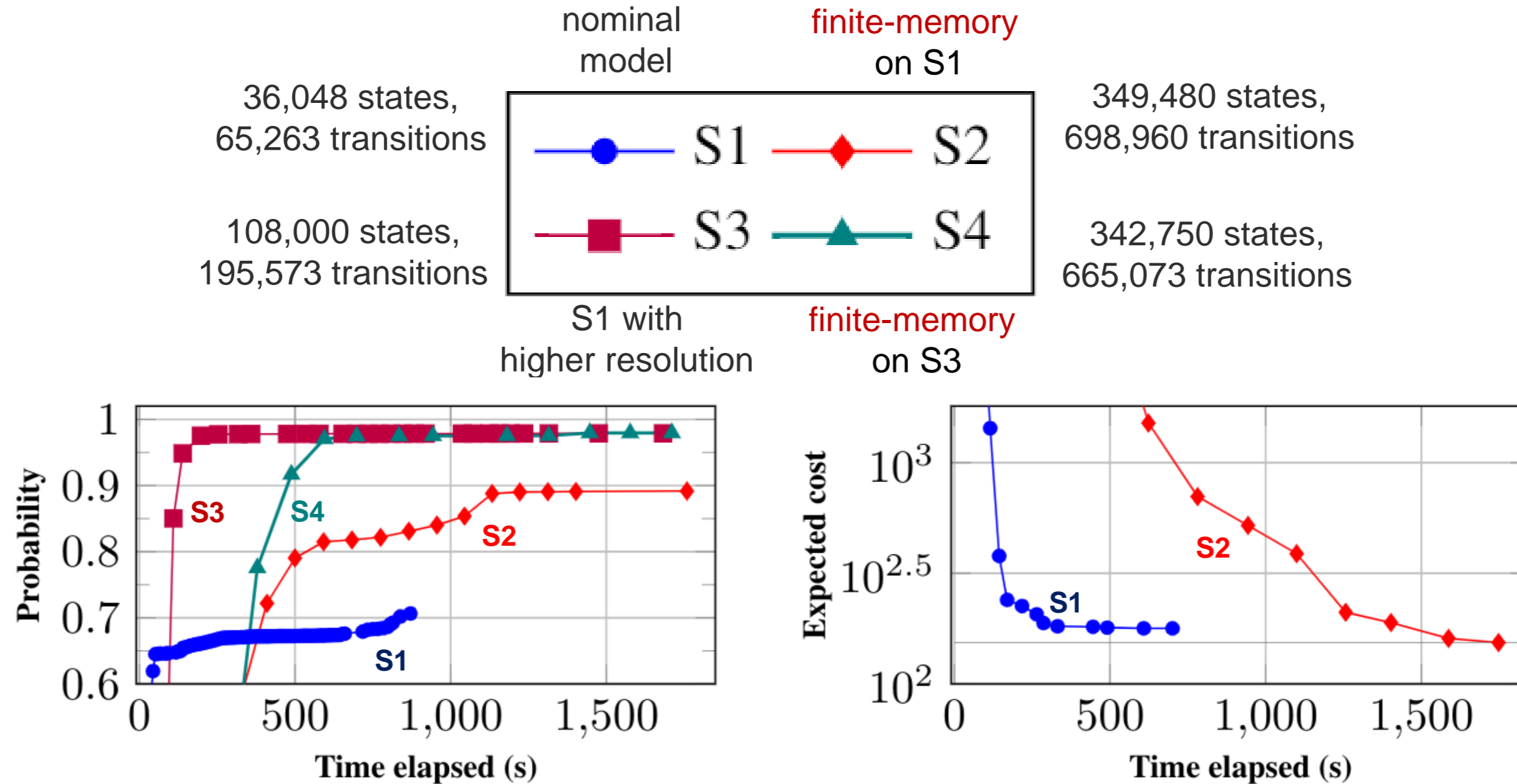
Problem:

Switching between orbits is possible if the orbits are close to each other

Partial observability over the current position of spacecraft, **uncertainty** on the location of other objects and operator



Results on Spacecraft Motion Planning

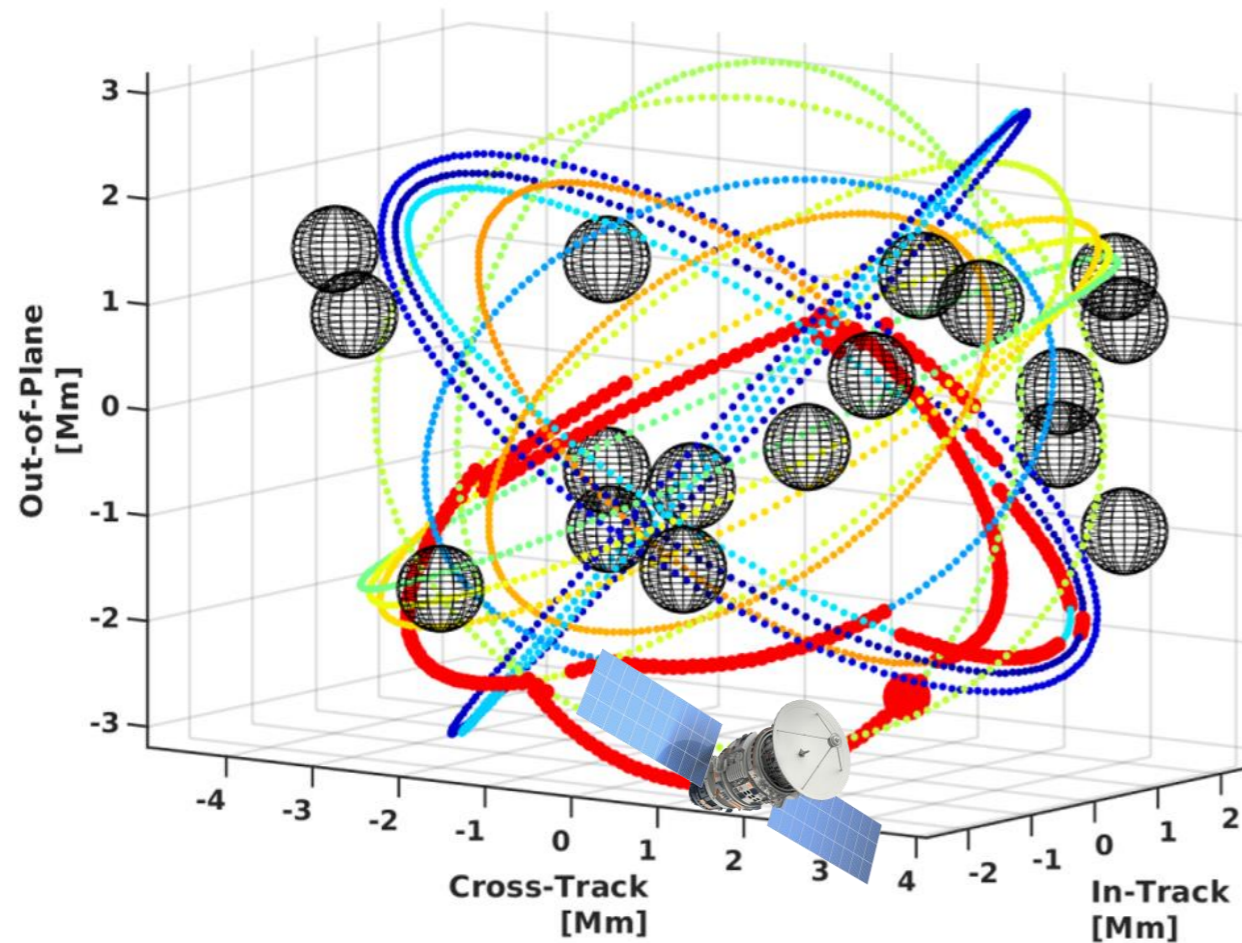


Can solve uPOMDPs with **hundreds of thousands states in minutes**

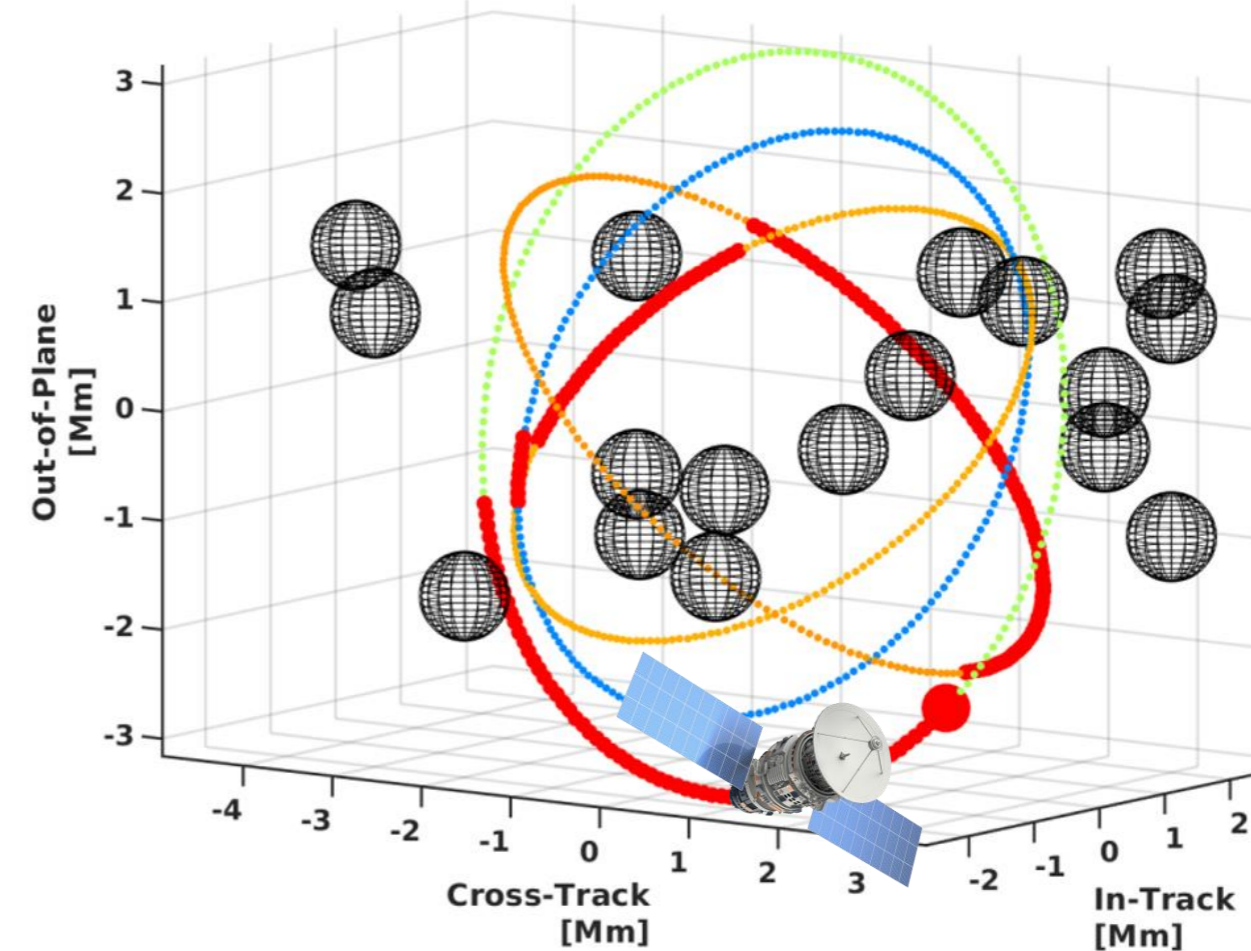
Models with memory take a **longer time to solve** but **yield better performance**

Results with Finite-Memory Policies

The memoryless policy (S1) **makes more orbit changes** than the policy with 5 memory nodes (S2)



memoryless policy (S1)



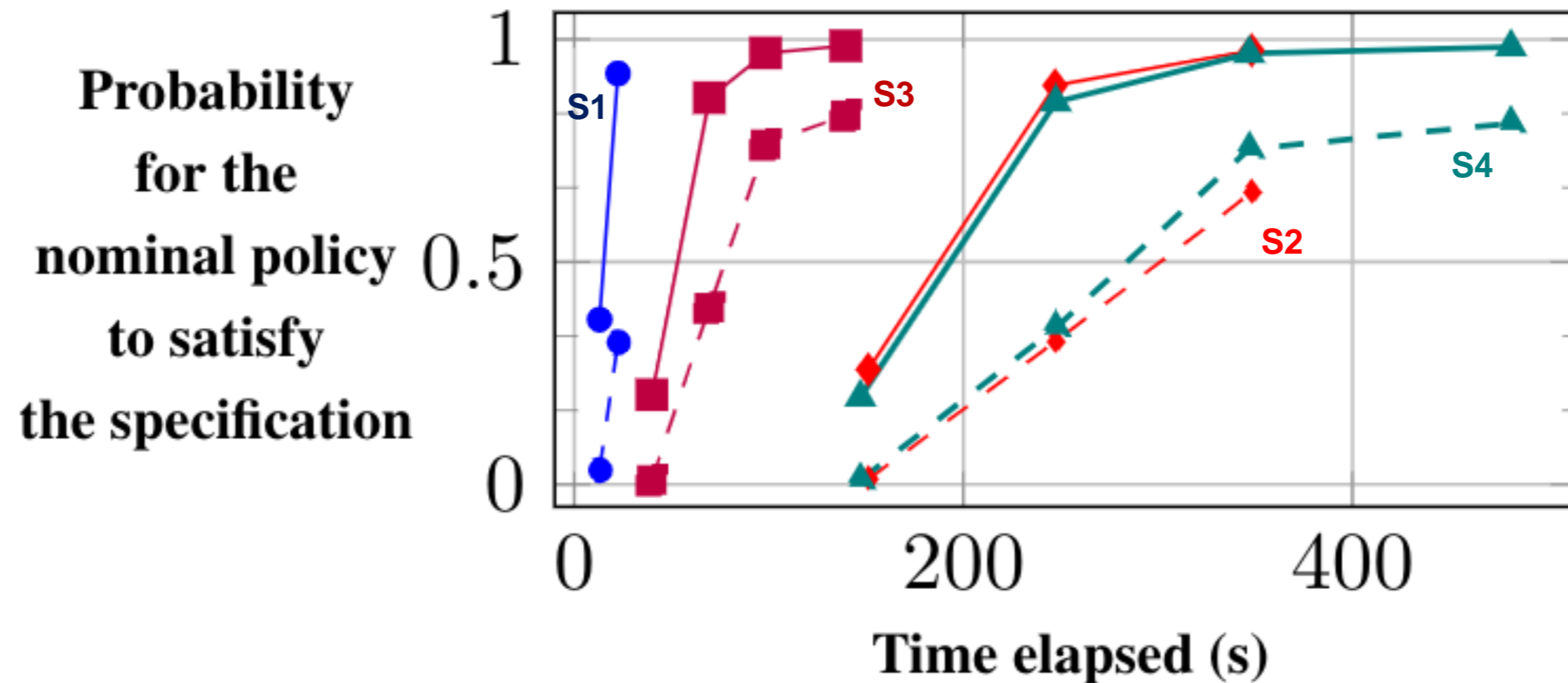
finite-memory policy (S2)

Robust Policies are Indeed More Robust

compute a policy
on the nominal POMDP
(solid lines)

AND

apply to the
uPOMDP
(dashed lines)



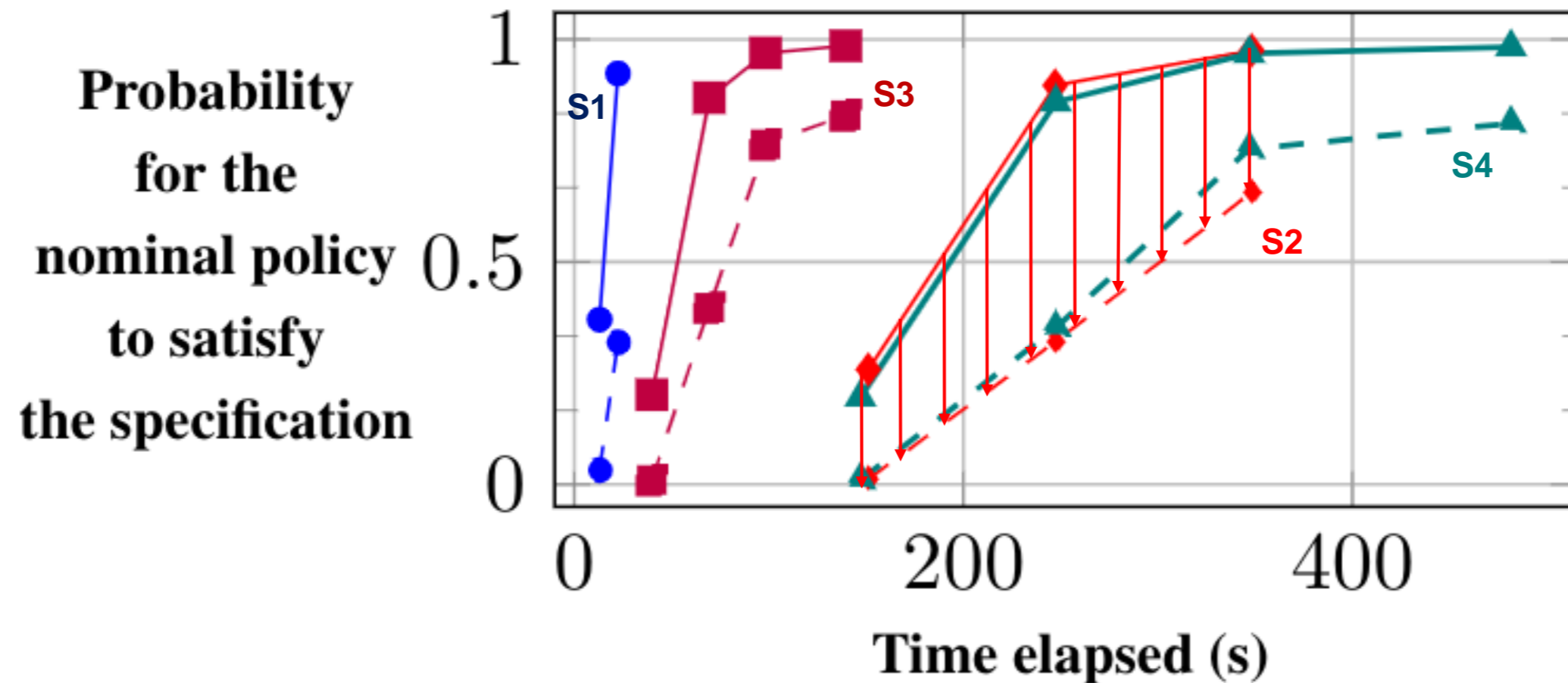
The performance of the nominal policies can **reduce up to 60 percentage points under uncertainty**

Robust Policies are Indeed More Robust

compute a policy
on the nominal POMDP
(solid lines)

AND

apply to the
uPOMDP
(dashed lines)



The performance of the nominal policies can **reduce up to 60 percentage points under uncertainty**

Conclusions and Future Work

Developed algorithms that scale to uPOMDPs that are **3 orders of magnitude larger** than previous approaches

Future work:

Uncertainty sets with correlations between different states

Incorporate these algorithms for safety in reinforcement learning