

**Cyberphysical Systems**  
(Autonomous Systems)  
**in the Intersection of**  
**Controls, Learning and Formal Methods**

**Ufuk Topcu**

The University of Texas at Austin

[u-t-autonomous.info](http://u-t-autonomous.info)

**a**UT**onomous**  
SYSTEMS GROUP

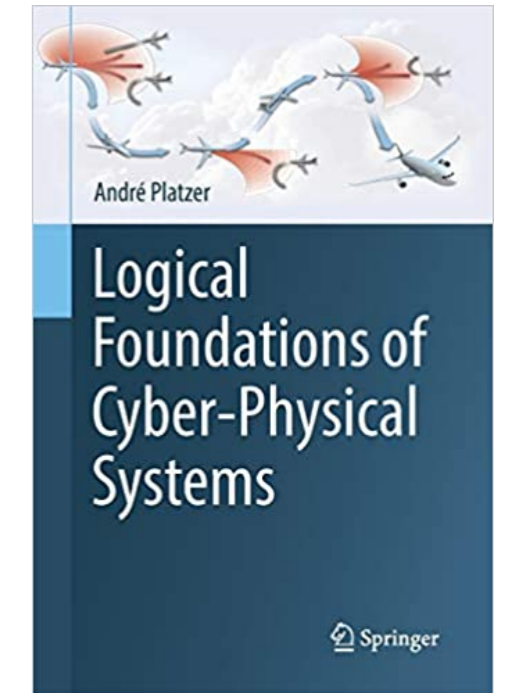
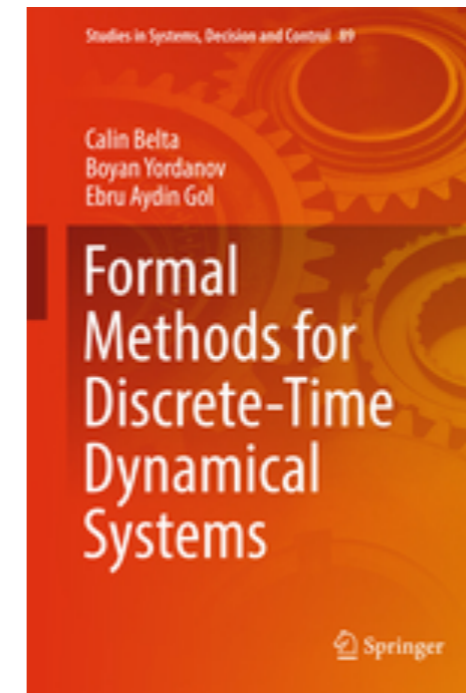
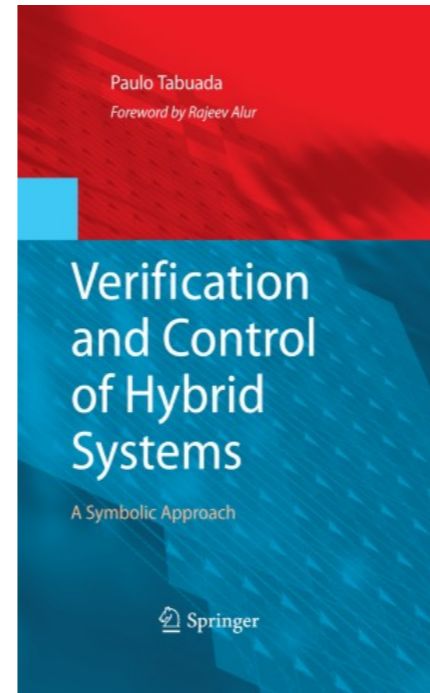
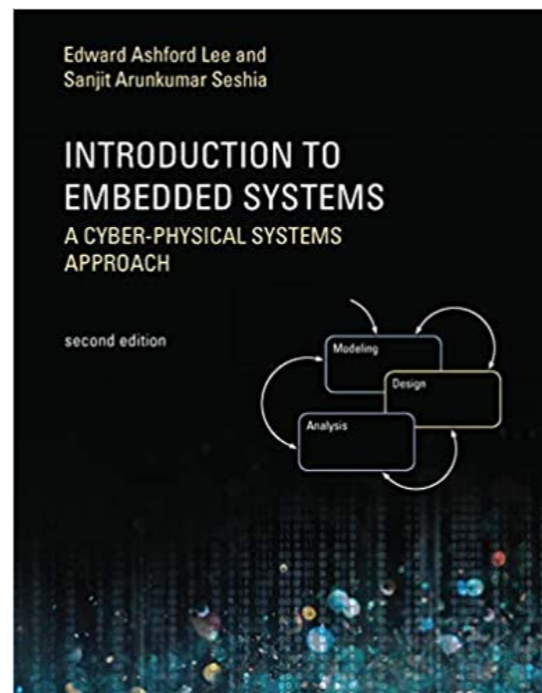
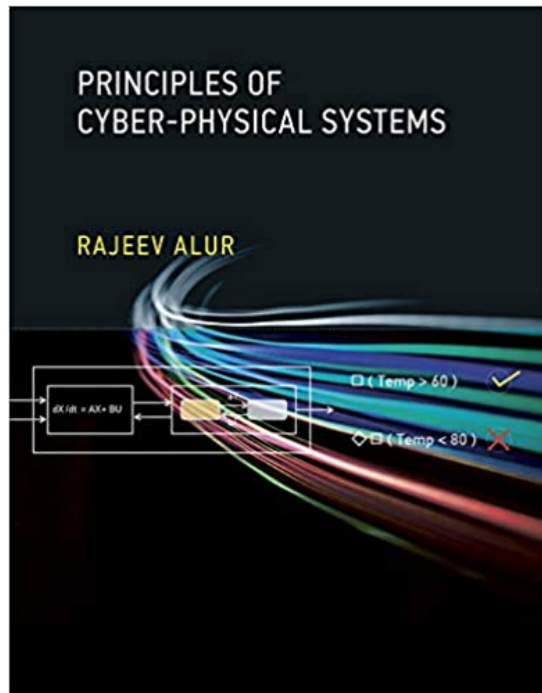
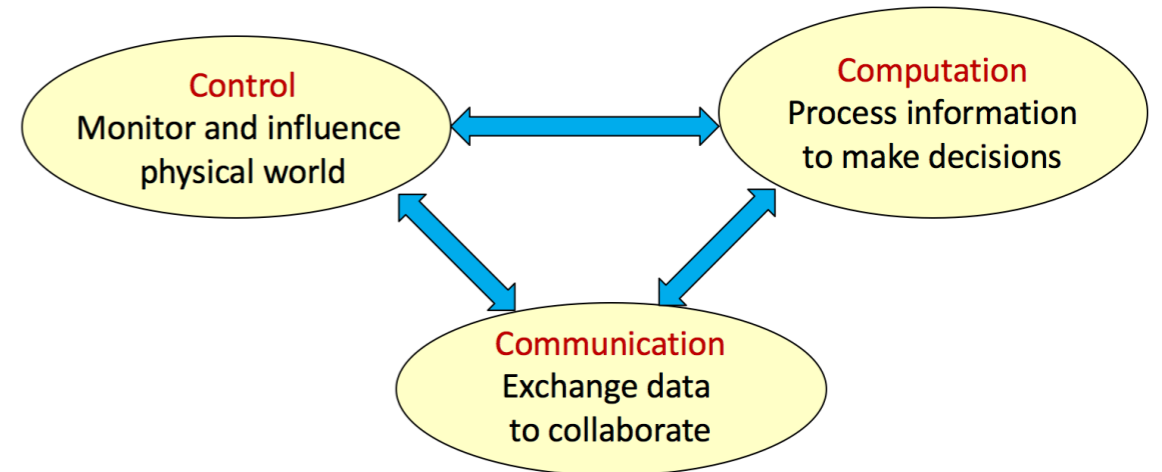
# What are cyberphysical systems (CPS)?

A cyberphysical system...

...consists of a collection of computing devices communicating with one another and...

...interacting with the physical world via sensors and actuators in a feedback loop.

Common among definitions: convergence





# Some properties of cyberphysical systems

## Reactivity and interaction with the physical world

A reactive system interacts with its environment in an ongoing manner via inputs (e.g., through sensors) and outputs (e.g., through actuators).

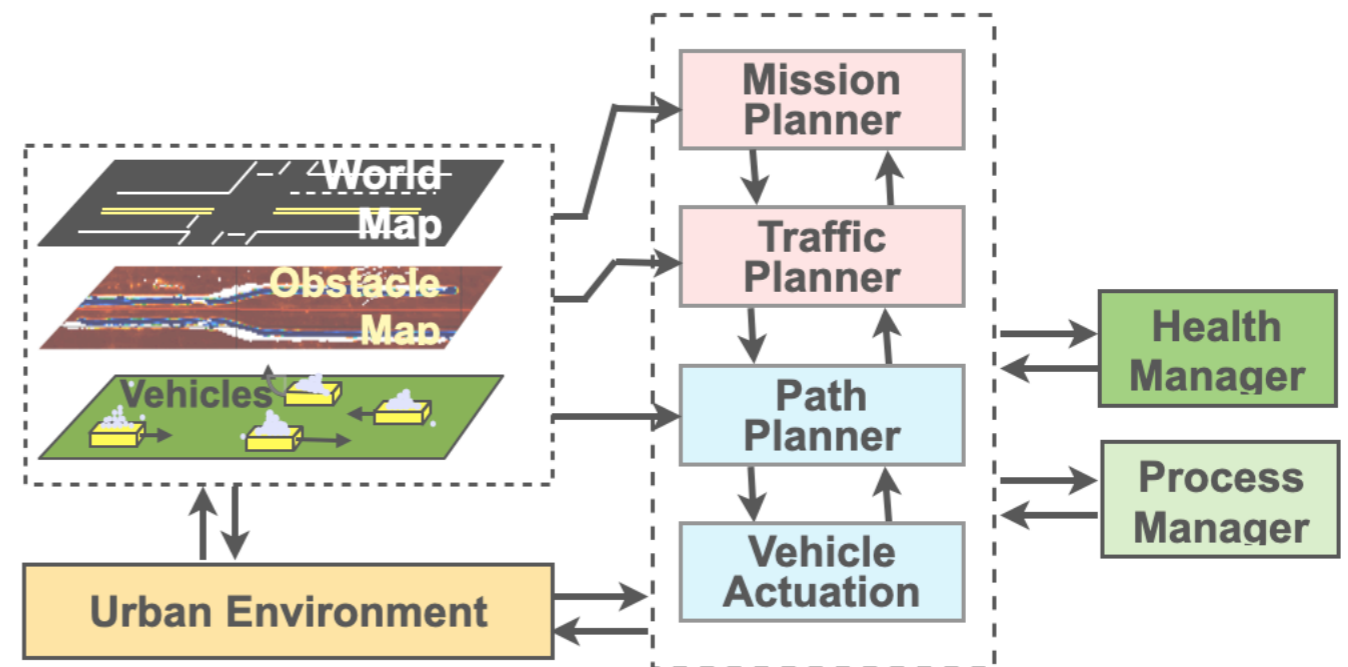


## Heterogeneity

Multiple, integrated functionality.

## Real-time decisions

Delays in computation and communication are critical.

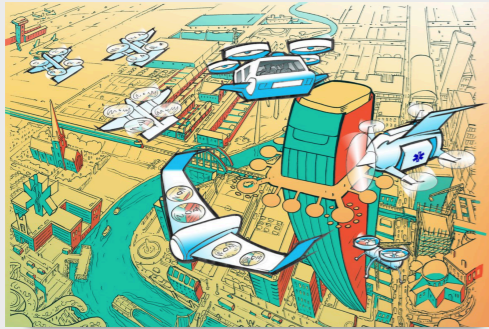


## Concurrency

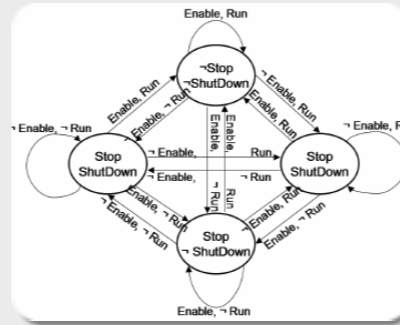
Multiple threads (components or processes) execute simultaneously, exchanging information to achieve a desired goal.



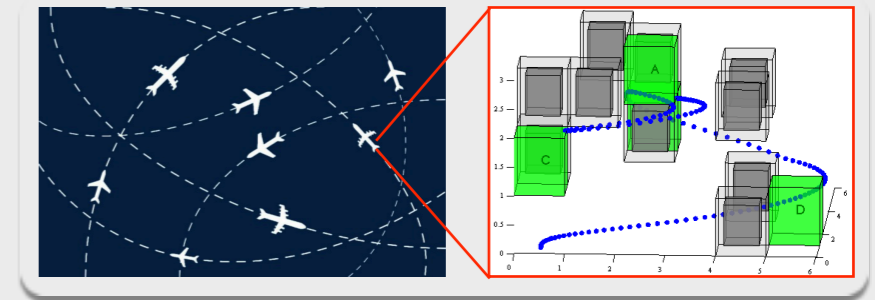
# What makes **autonomy** hard?



**dynamic environment**



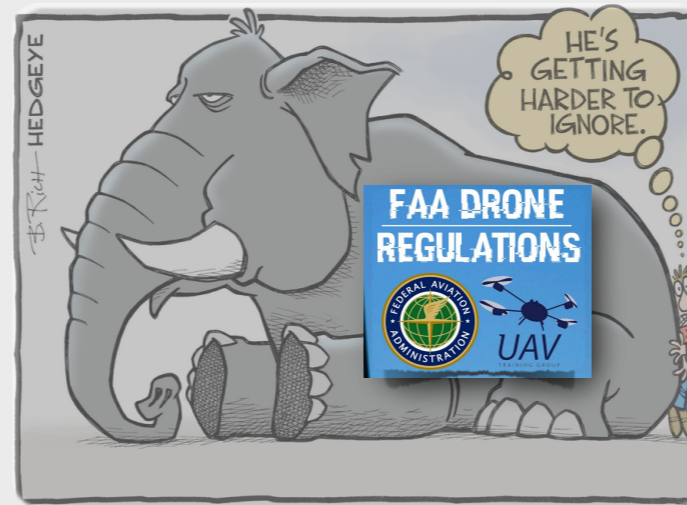
**complex missions**



**heterogenous decisions**



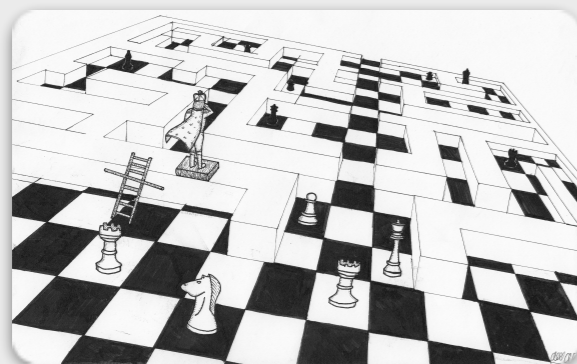
**run-time faults**



**verifiability**



**unknown environments**



**integrity of critical information**

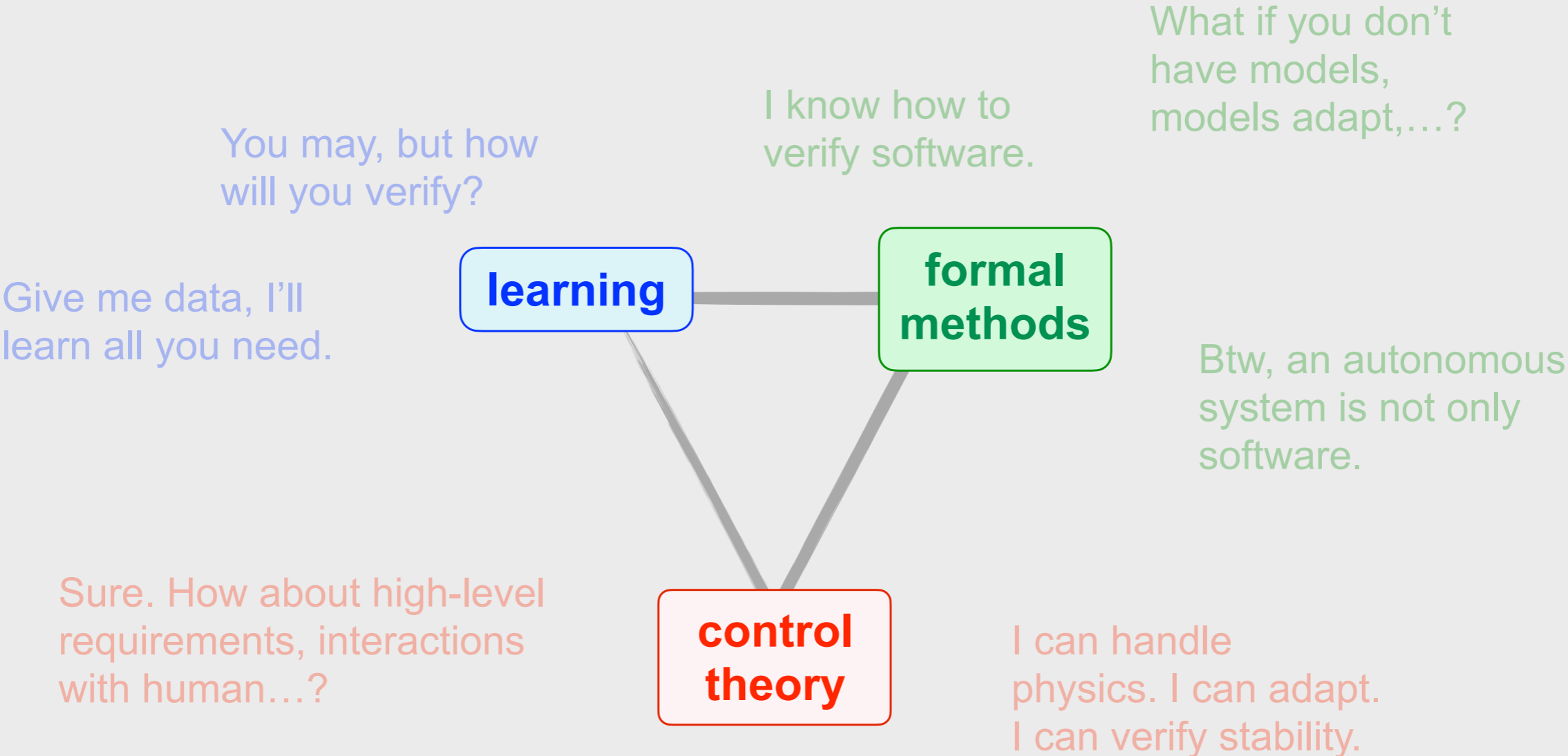


**imperfect perception**



**variations in user characteristics**

# Why are we not there yet? What is still missing?



**Autonomous systems are nobody's comfort zone:  
We need hybrid solutions.**

# Outline

## Correct-by-construction synthesis of hierarchical control protocols

- **formal methods** ↔ **controls**

What would it take to construct the control software of an autonomous system **in an hour**—**as opposed to days or even weeks**—to deliver a mission?

## Verifiable reinforcement learning

- **formal methods** ↔ **learning**

How can an autonomous system **learn** how to execute a new mission in an a priori unknown environment **efficiently** and **safely**?

## Planning in POMDPs

- **formal methods** ↔ **convex optimization**
- **formal methods** ↔ **learning**

How can we cope with imperfection and/or limitations in run-time information?



# A (sample) synthesis problem

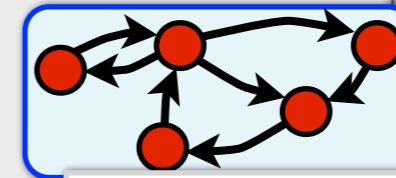
Given:

- **System model**

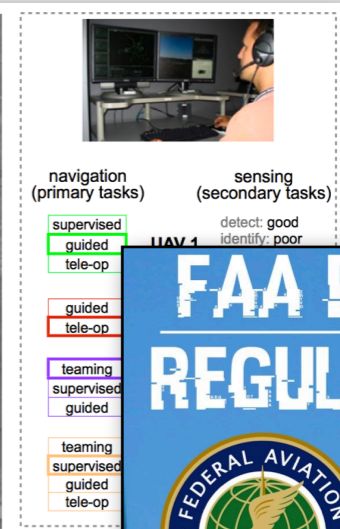
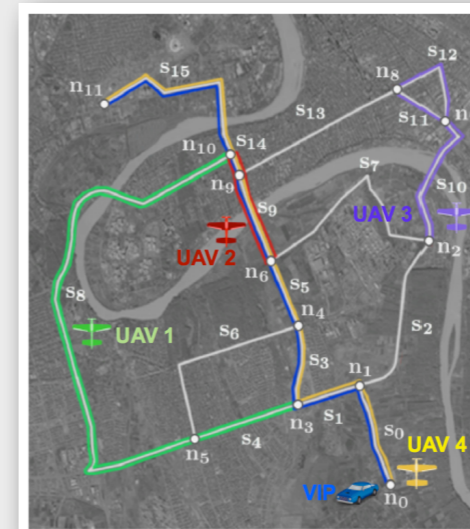
- both **continuous** & **discrete** evolution
- actuation limitations
- modeling uncertainties & disturbances

- **Specifications in “temporal logic”**

- high-level requirements & goals
- assumptions about the a priori unknown, dynamic environment



$$\dot{x} = f(x, u, \delta)$$
$$g(x, u) \geq 0$$



**Automatically synthesize a control protocol that**

- manages the system behavior;
- reacts to changes in allowable external environment; and
- is **provably correct** with respect to the specifications.

# De-tour: Specifying behavior with temporal logic

**Temporal Logic** = **Propositional Logic** + **Temporal Operators**

$\wedge$ (and)	$\diamond$ (eventually)
$\vee$ (or)	$\square$ (always)
$\rightarrow$ (implies)	$\mathcal{U}$ (until)
$\neg$ (not)	

- Reason about infinite sequences  $\sigma = s_0 s_1 s_2 \dots$  of states
- Many different dialects of temporal logic (with probabilistic and epistemic modalities)
- Specify safe, allowable, required, or desired behavior of system and/or environment.

## Coverage:

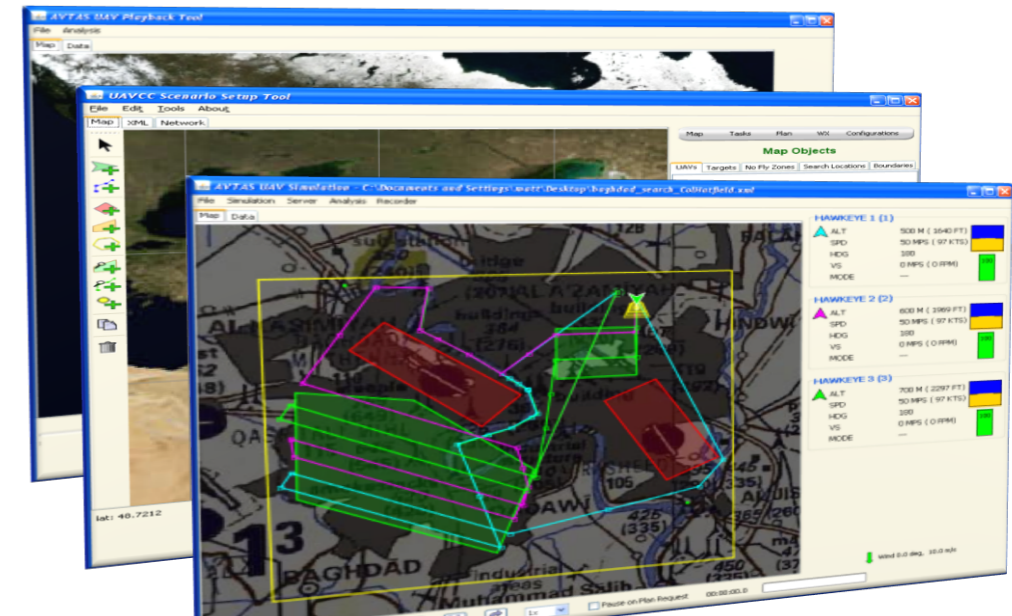
$$\square \diamond uav_1 = w_3 \wedge \square \diamond uav_1 = w_4 \wedge \square \diamond uav_1 = w_5 \wedge \diamond uav_2 = w_1 \wedge \diamond uav_2 = w_2 \wedge \diamond uav_2 = w_6$$

## Sequencing:

$$\diamond (uav = w_1 \wedge \diamond (uav = w_2 \wedge \diamond uav = w_3))$$

## Coordination with mutual exclusion:

$$\neg (uav_1 = uav_2 \vee uav_1 = uav_3 \vee uav_1 = uav_4 \vee uav_2 = uav_3 \vee uav_2 = uav_4 \vee uav_3 = uav_4) \mathcal{U} (uav_1 = w_1 \wedge uav_2 = w_1 \wedge uav_3 = w_1 \wedge uav_4 = w_1 \wedge \bigcirc^{-1} uav_1 = a_{1,1} \wedge \bigcirc^{-1} uav_2 = a_{1,3} \wedge \bigcirc^{-1} uav_3 = a_{1,5} \wedge \bigcirc^{-1} uav_4 = a_{1,7})$$



## Sequencing with avoidance:

$$\diamond uav = w_1 \wedge \diamond uav = w_2 \wedge \diamond uav = w_3 \wedge \neg w_2 \mathcal{U} w_1 \wedge \neg w_3 \mathcal{U} w_2$$

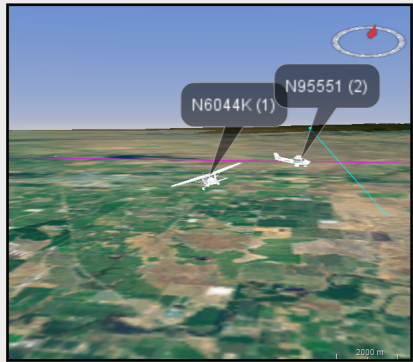
## Never after:

$$\square (uav = w_1 \rightarrow \bigcirc \square \neg uav = w_1)$$

# A Solution: Hierarchical Control Structure

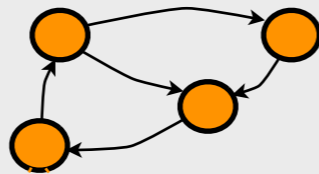
**Theorem: Correctness provably guaranteed by the construction of the abstractions and splitting of the specifications**

Different views



long-horizon specifications

Multi-scale models



Synthesis method

Iterative graph search

Control protocol



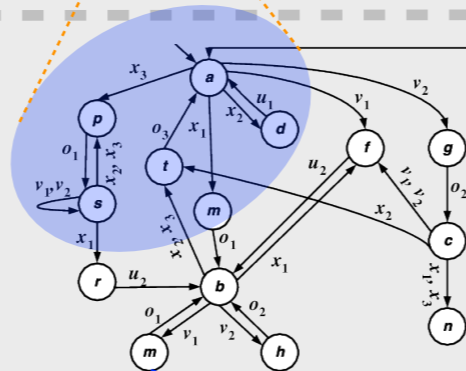
response

Goal Generator

$\mathcal{G}$

Reactive Controller

Two-player, turn-based graph game



short-horizon specifications

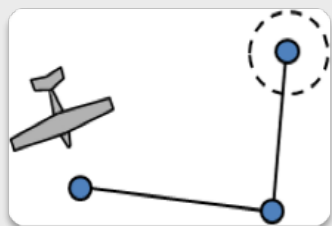
Constrained, finite-horizon optimal control

$$x_{t+1} = f(x_t, w_t, u_t)$$

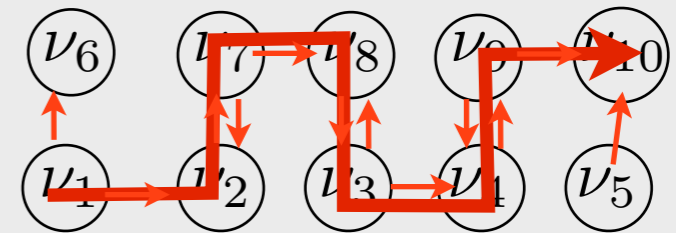
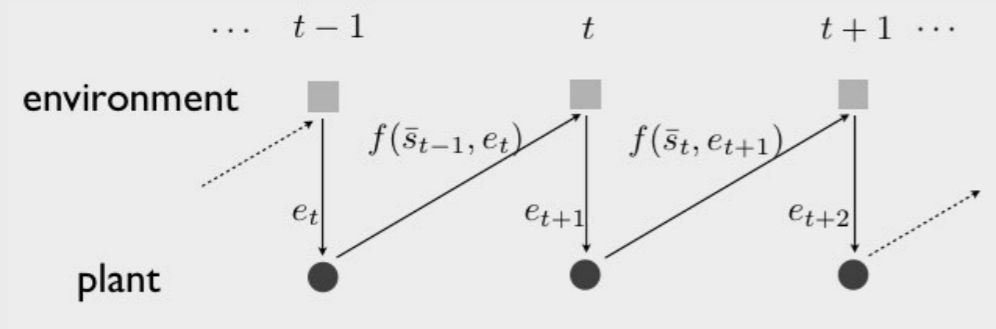
$$x \in \mathcal{X}, u \in \mathcal{U}, w \in \mathcal{W}$$

constraints on continuous state + input

Continuous Controller



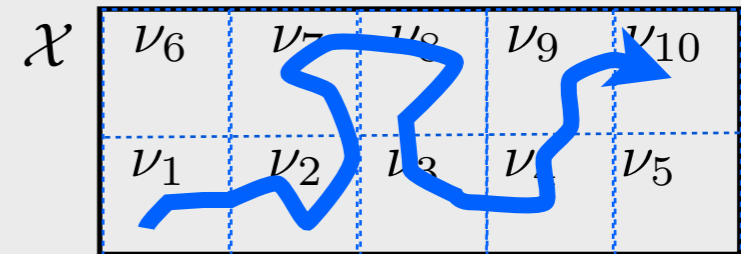
# Zoom in a bit...



$$x_{t+1} = f(x_t, w_t, u_t)$$

$$x \in \mathcal{X}, u \in \mathcal{U}, w \in \mathcal{W}$$

$\exists u \in \mathcal{U}$  and finite  $T > 0$  s.t.  
 $x_t \in \mathcal{X}_{\text{initial}} \cup \mathcal{X}_{\text{target}} \quad \forall t = 0, \dots, T,$   
 $x_T \in \mathcal{X}_{\text{target}}$   
 for all  $x_0 \in \mathcal{X}_{\text{initial}}$  and  $w \in \mathcal{W}$ ?



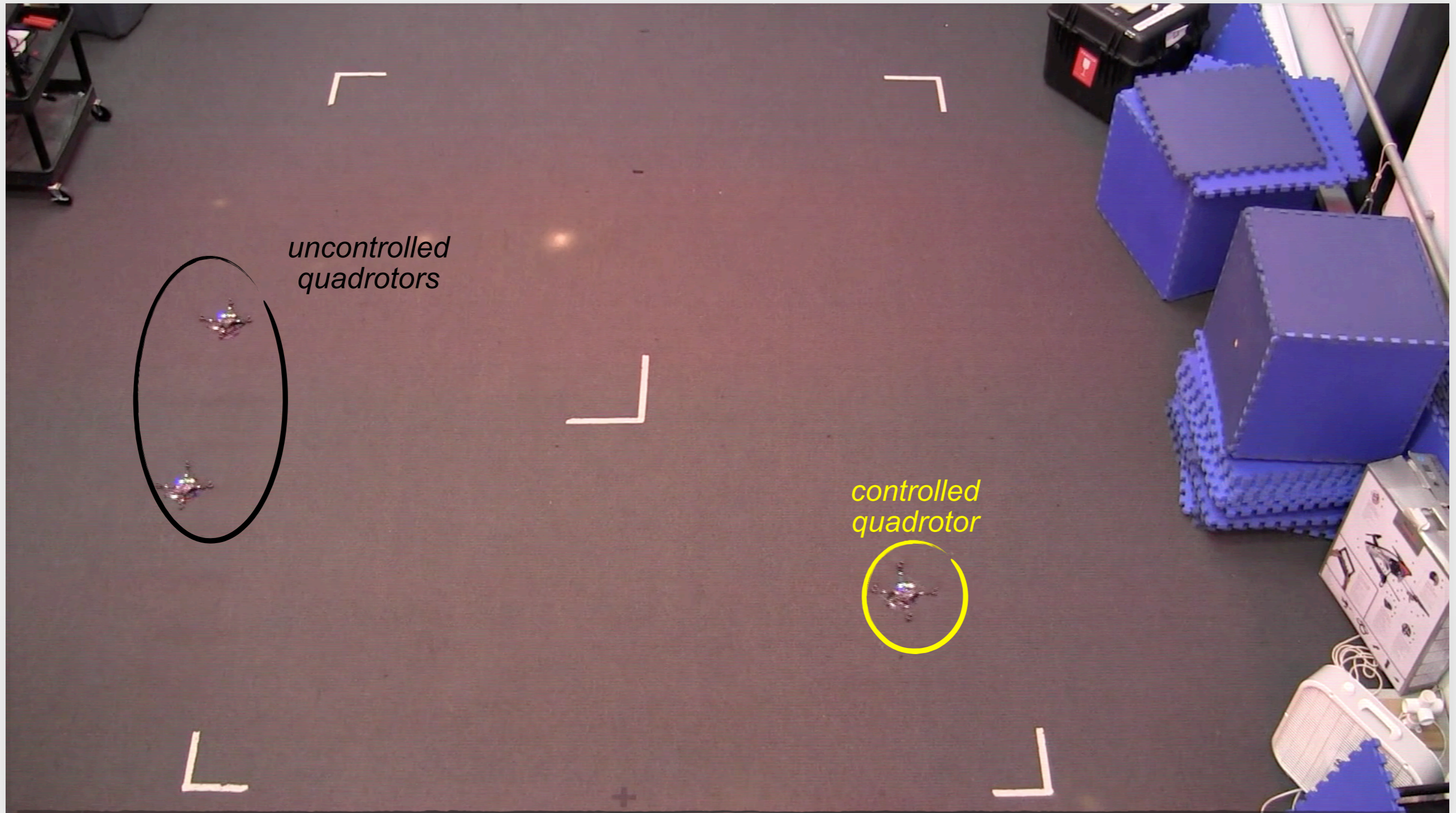
**Correctness:** For any discrete run satisfying the specification, there exists an admissible control signal leading to a continuous trajectory satisfying the specification.

**Proof — “correct by construction”:**

Constructive  $\rightarrow$  Finite-state model + Continuous control signals.



# Specify + Synthesize



# Outline

## Correct-by-construction synthesis of hierarchical control protocols

- **formal methods** ↔ **controls**

## Verifiable reinforcement learning

- **formal methods** ↔ **learning**

## Planning in POMDPs

- **formal methods** ↔ **convex optimization**
- **formal methods** ↔ **learning**

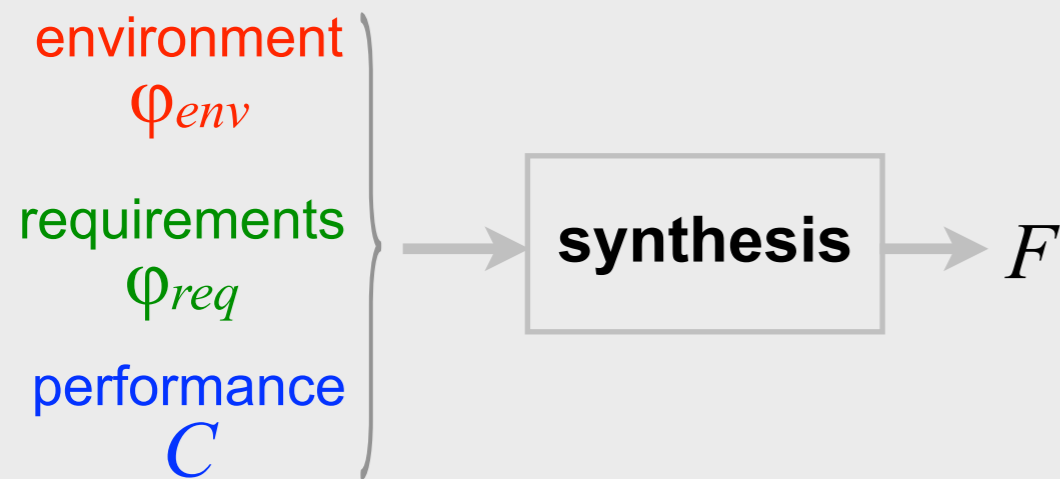
What would it take to construct the control software of an autonomous system **in an hour**—as opposed to days or even weeks—to deliver a mission?

How can an autonomous system **learn** how to execute a new mission in an a priori unknown environment **efficiently** and **safely**?

How can we cope with imperfection and/or limitations in run-time information?

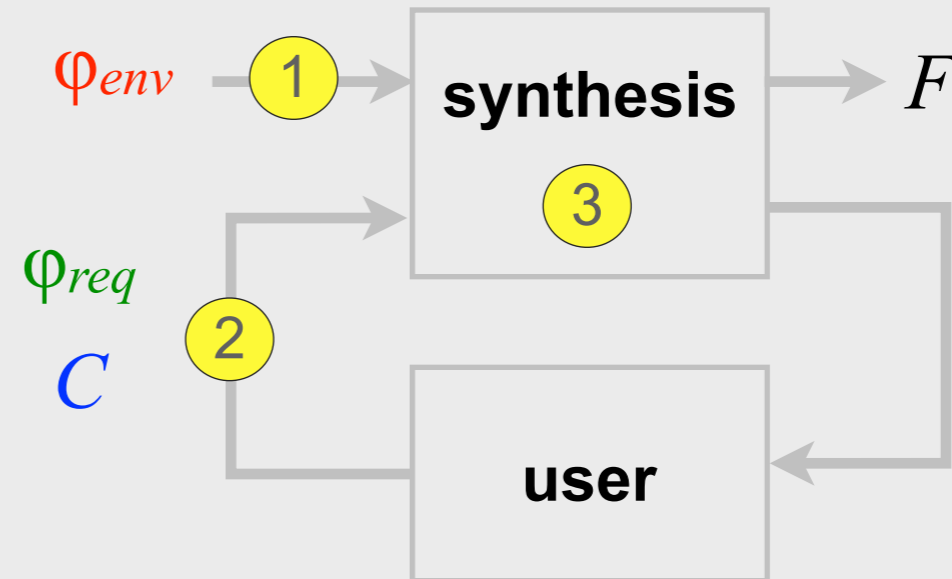
# Need for learning

## Conventional Synthesis



Known/fixed  
temporal logic constraints  
+ reward structure

## Synthesis + Learning

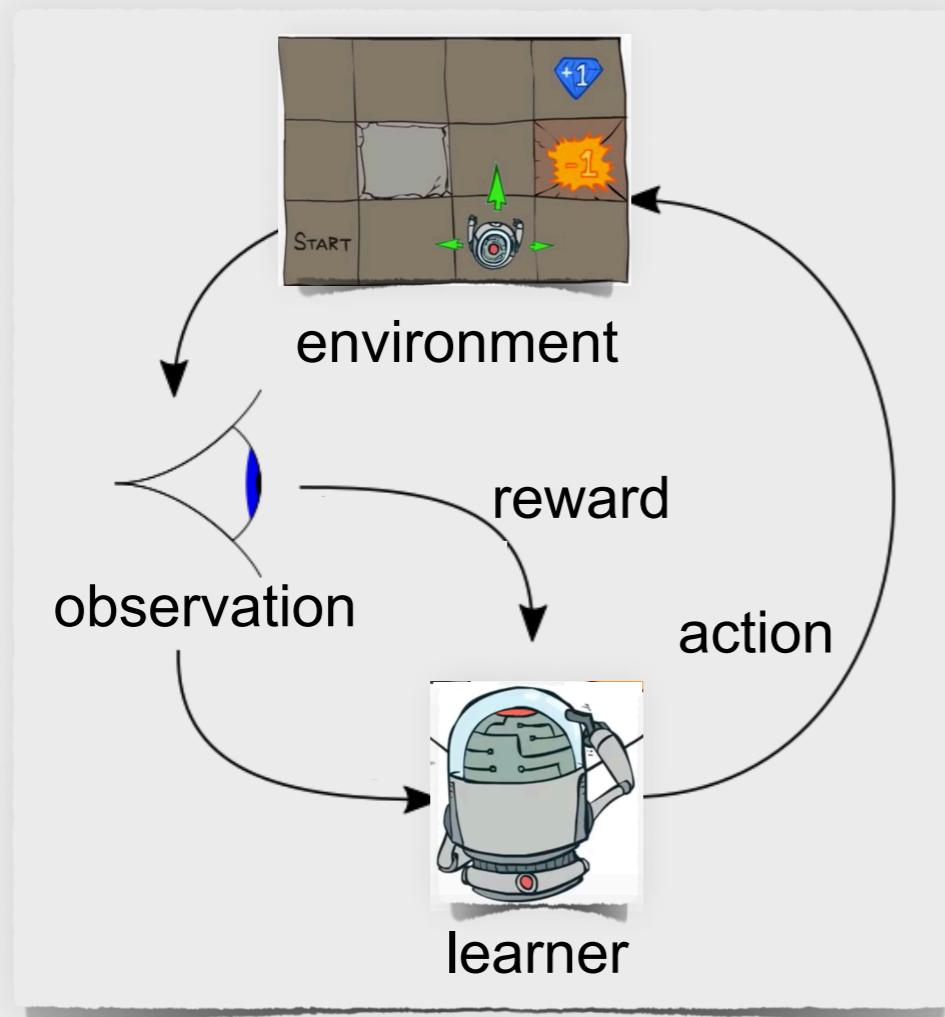


- 1 No a priori knowledge about environment. But, online (and offline) data available.
- 2 Operators do not “speak” temporal logic or cannot express complex reward structures.
- 3 Controller itself may be learned from data or examples.

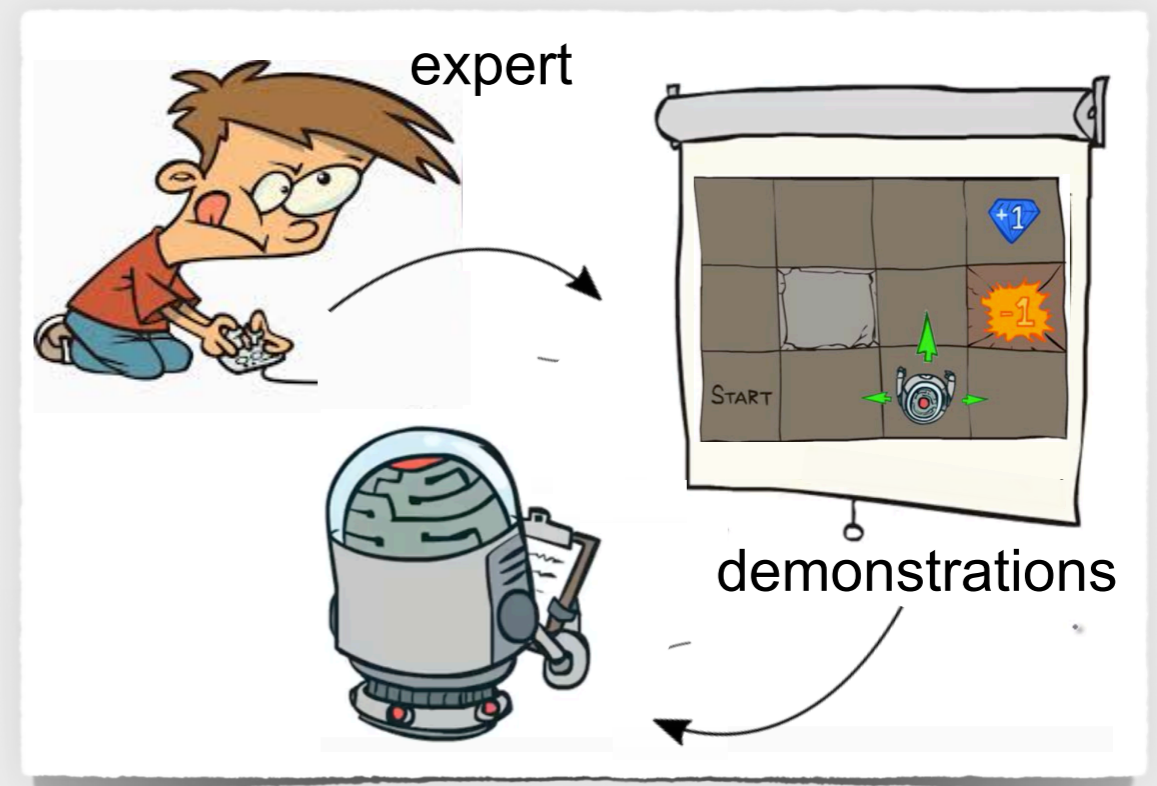


# Learning **subject to temporal logic specifications**

Reinforcement learning



Inverse reinforcement learning



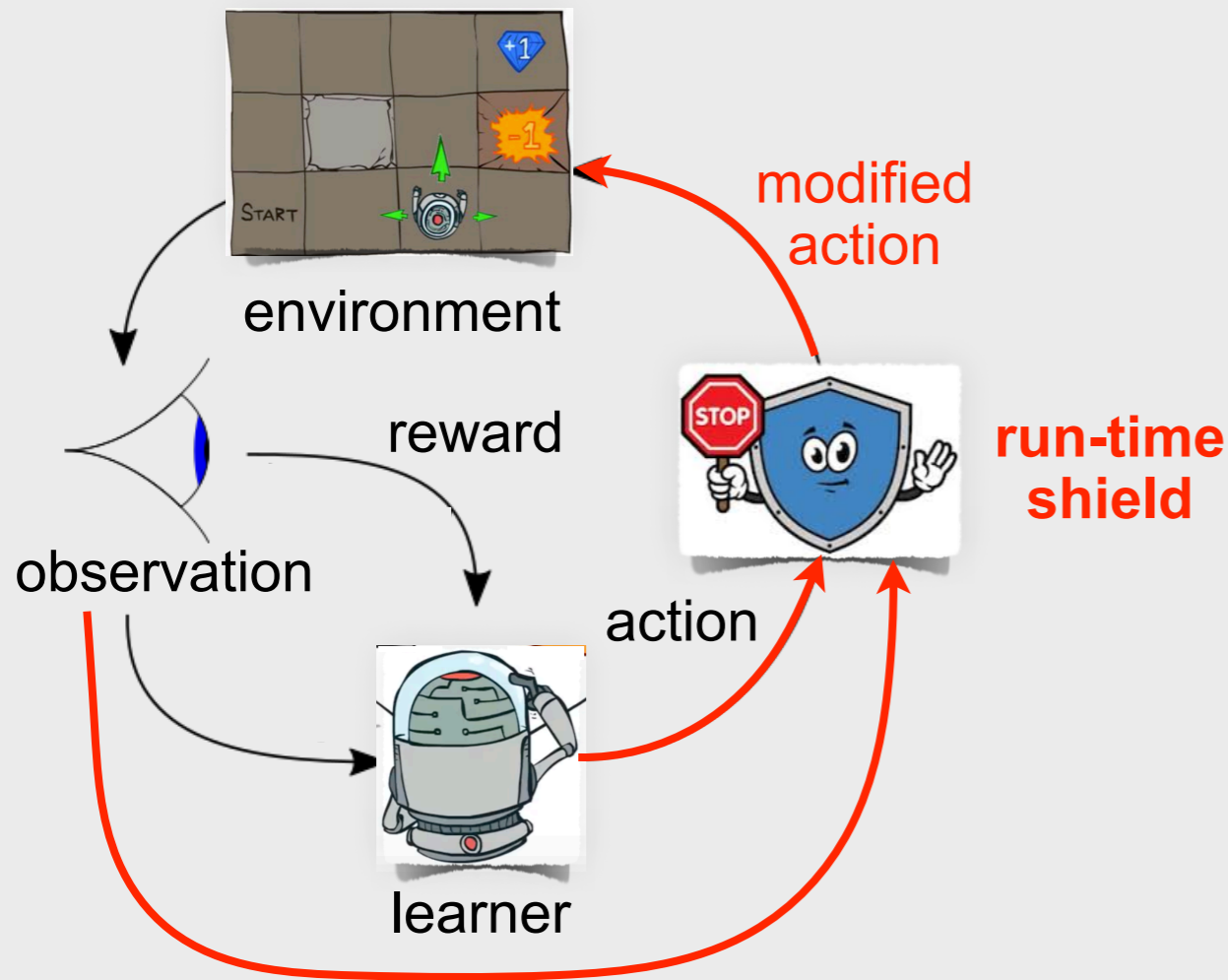
**new piece of problem data: temporal logic specification  $\varphi$**

**Does the learned strategy satisfy  $\varphi$ ?**

**Is  $\varphi$  violated during learning?**



# Safe reinforcement learning via shielding



## Run-time shielding...

- Corrective w.r.t. specifications in the safety fragment of linear temporal logic
- Minimally interfering
- Agnostic to the learner
- Can be used with function-approximation-based methods
- Preserves (constrained) convergence in some cases

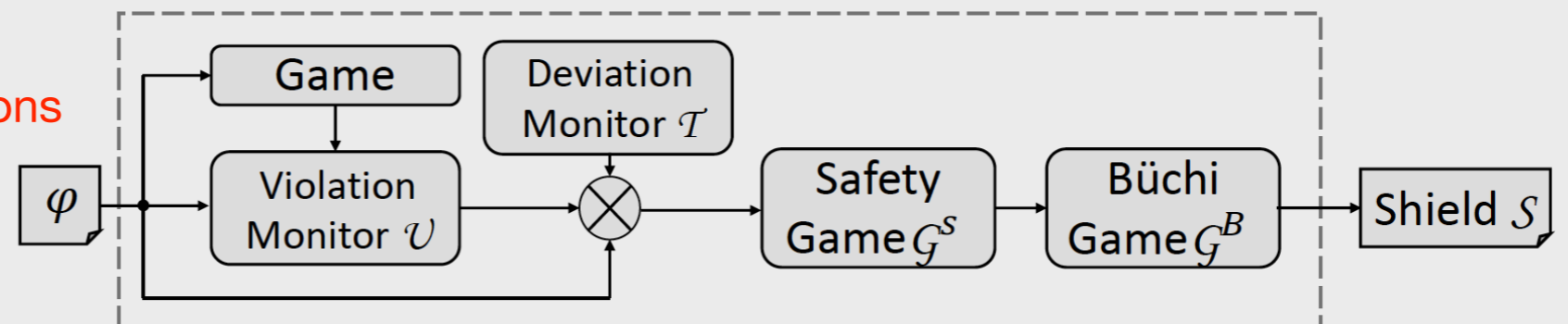
**Provable properties**

- Improves the data-efficiency in learning

**Empirically observed**

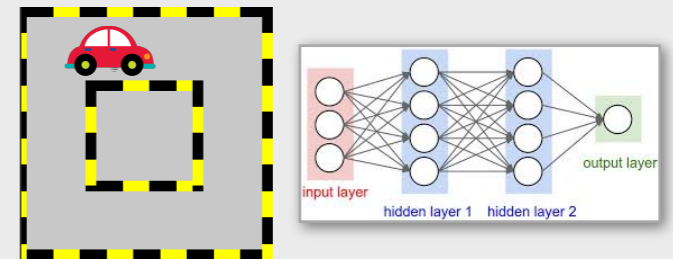
## Shield synthesis

specifications



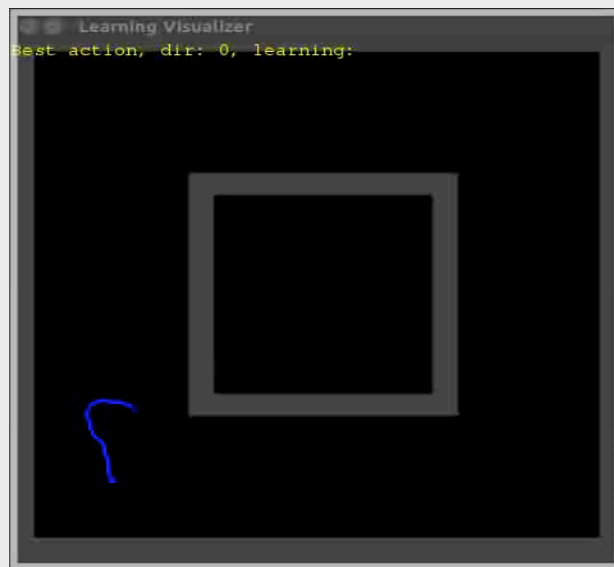
# Examples of shielded reinforcement learning

Deep reinforcement learning with a neural network of three layers and Boltzmann exploration



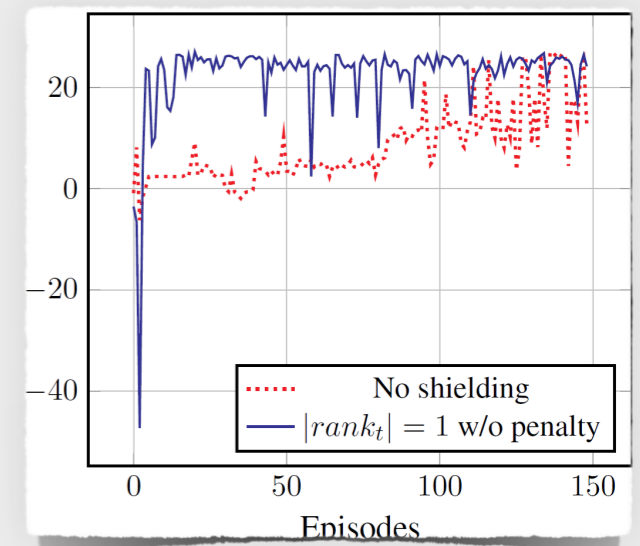
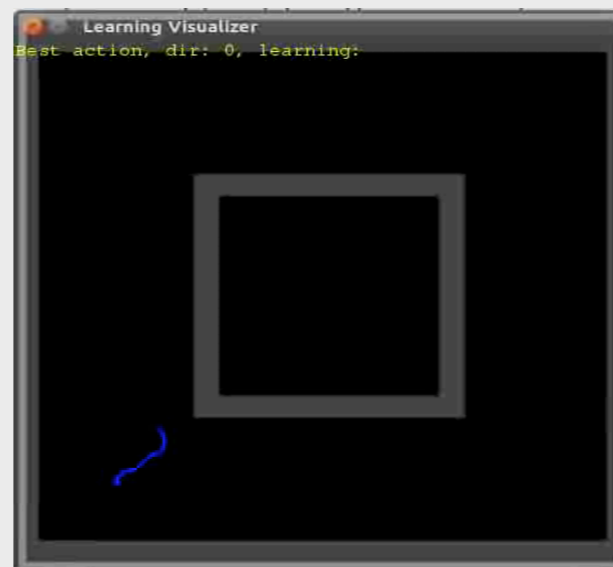
**without shielding**

(40 min training + 1600 crashes)

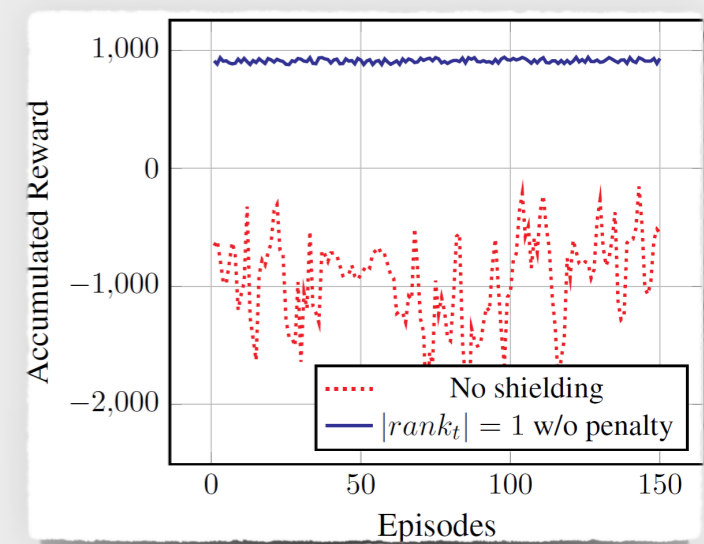
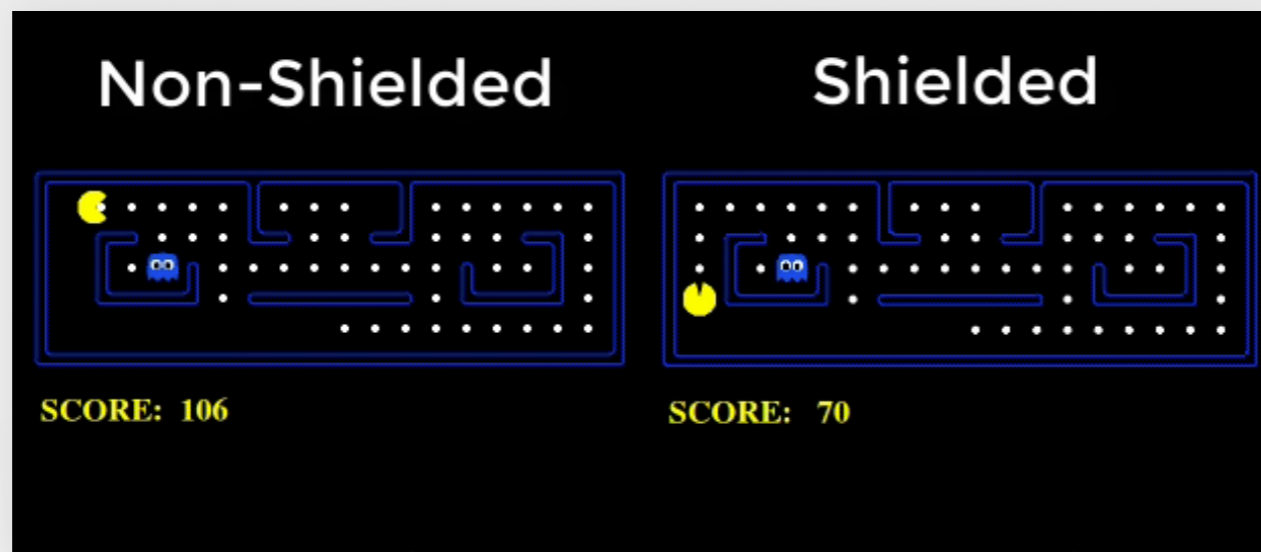


**with shielding**

(1 min training + no crashes)

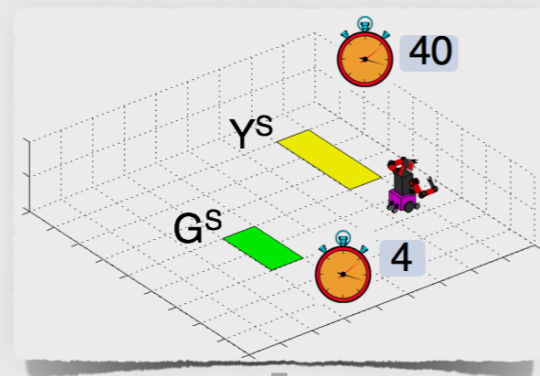


Reinforcement learning for PACMAN



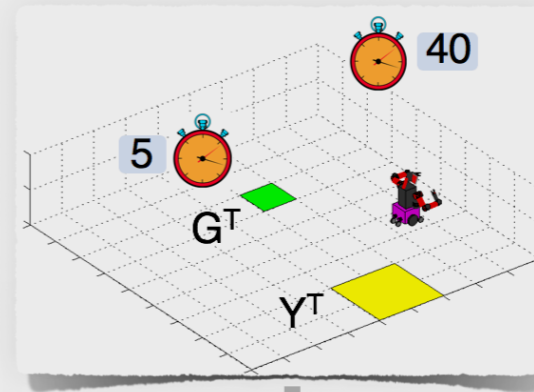
# Can **task** similarities help learn faster?

Source task



Data

Target task



Data

Inferred formulas

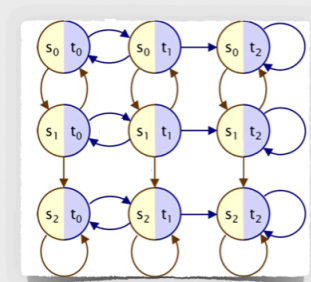
$$\phi^S = \diamond_{[1,15]} \square_{[0,3]} G^S \wedge \diamond_{[21,39]} Y^S$$

Logical similarity?

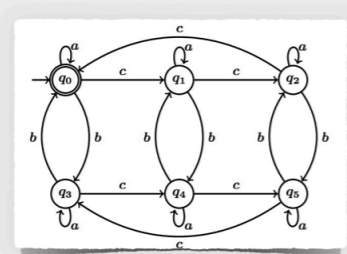
$$\phi^T = \diamond_{[5,18]} \square_{[0,5]} G^T \wedge \diamond_{[24,39]} \bar{Y}^T$$

Construct timed automaton

Construct timed automaton



x



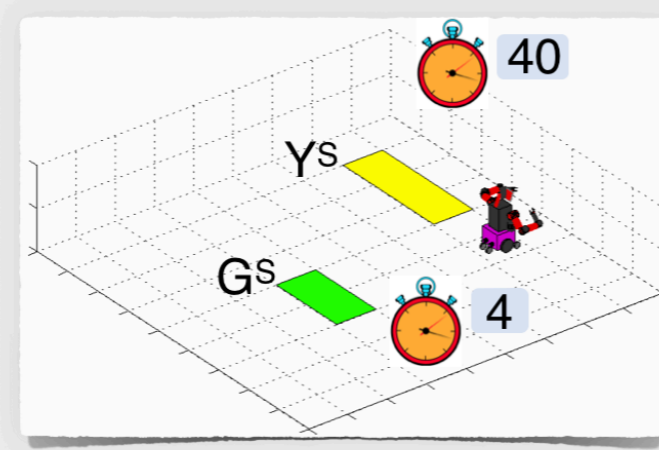
RL in extended state space

Extended Q-functions

RL in extended state space

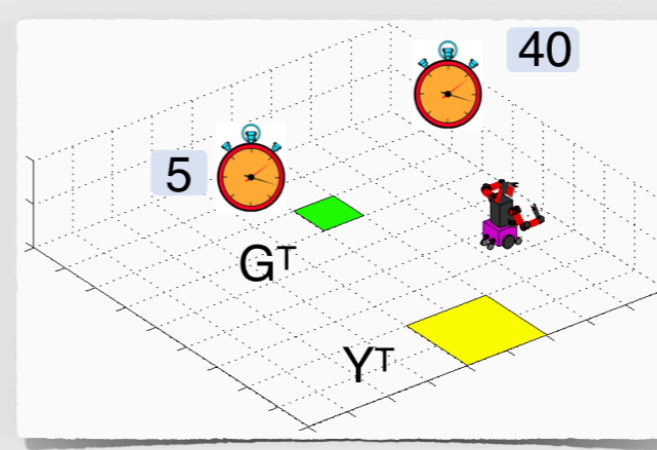
# Transfer of temporal logical similarities improves data efficiency by orders of magnitude.

Source task



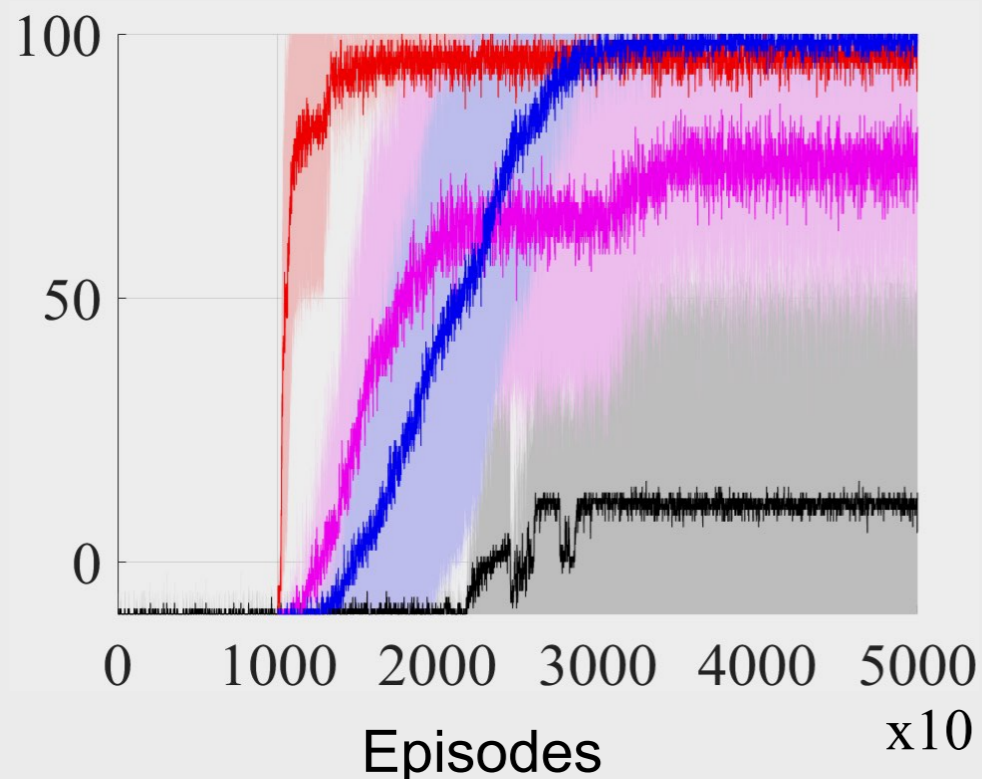
$$\diamond_{[1,15]} \square_{[0,3]} G^S \wedge \diamond_{[21,39]} Y^S$$

Target task



$$\diamond_{[5,18]} \square_{[0,5]} G^T \wedge \diamond_{[24,39]} \bar{Y}^T$$

Cumulative reward



+ transfer of extended Q-functions

+ “informative” MITL formulas

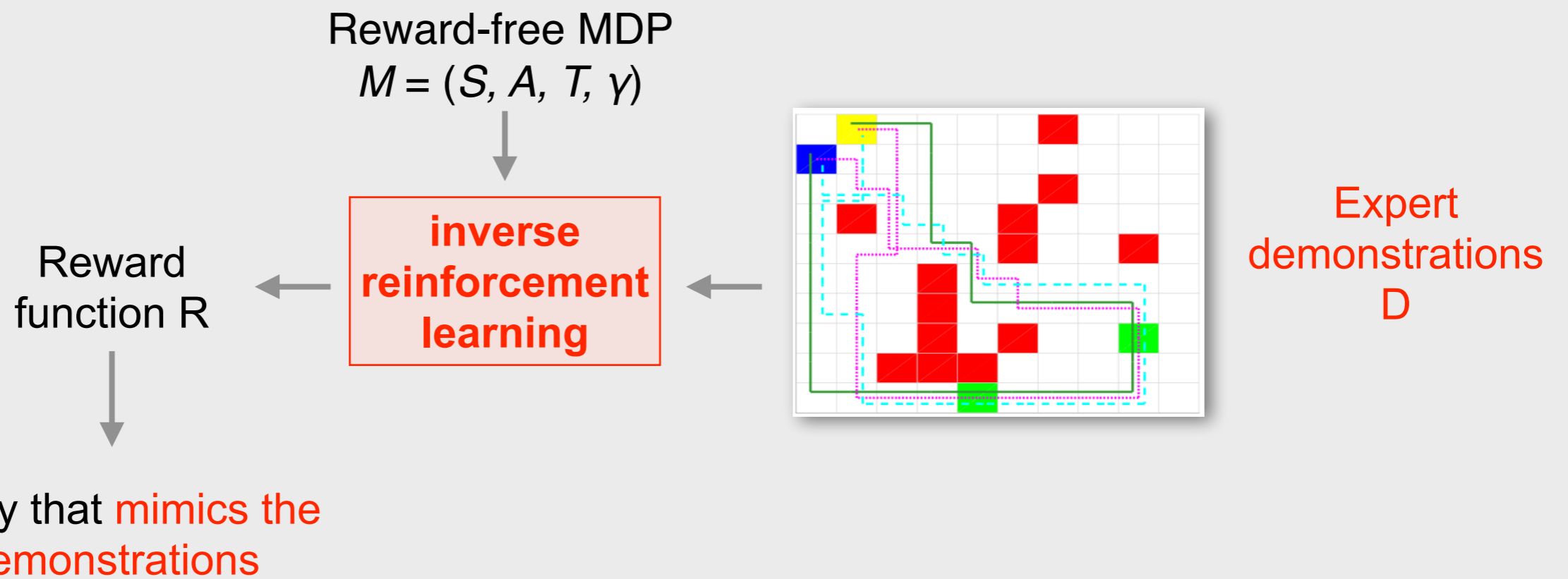
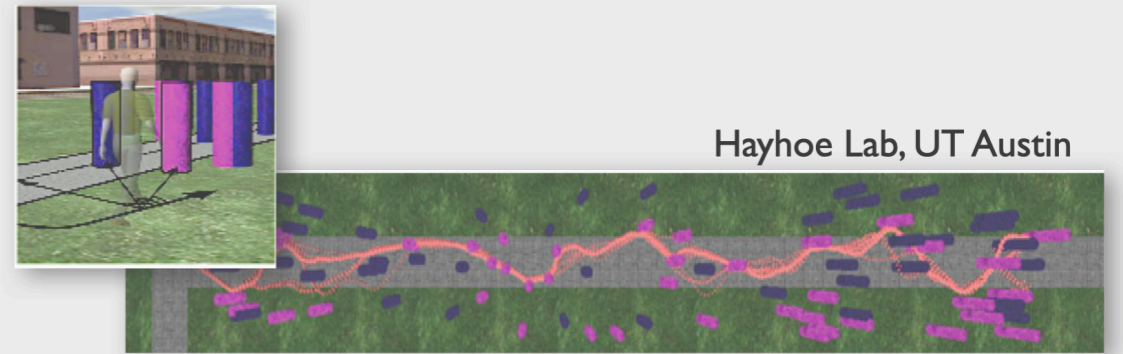
+ extended state space

+ maximize classification rate in MITL inference

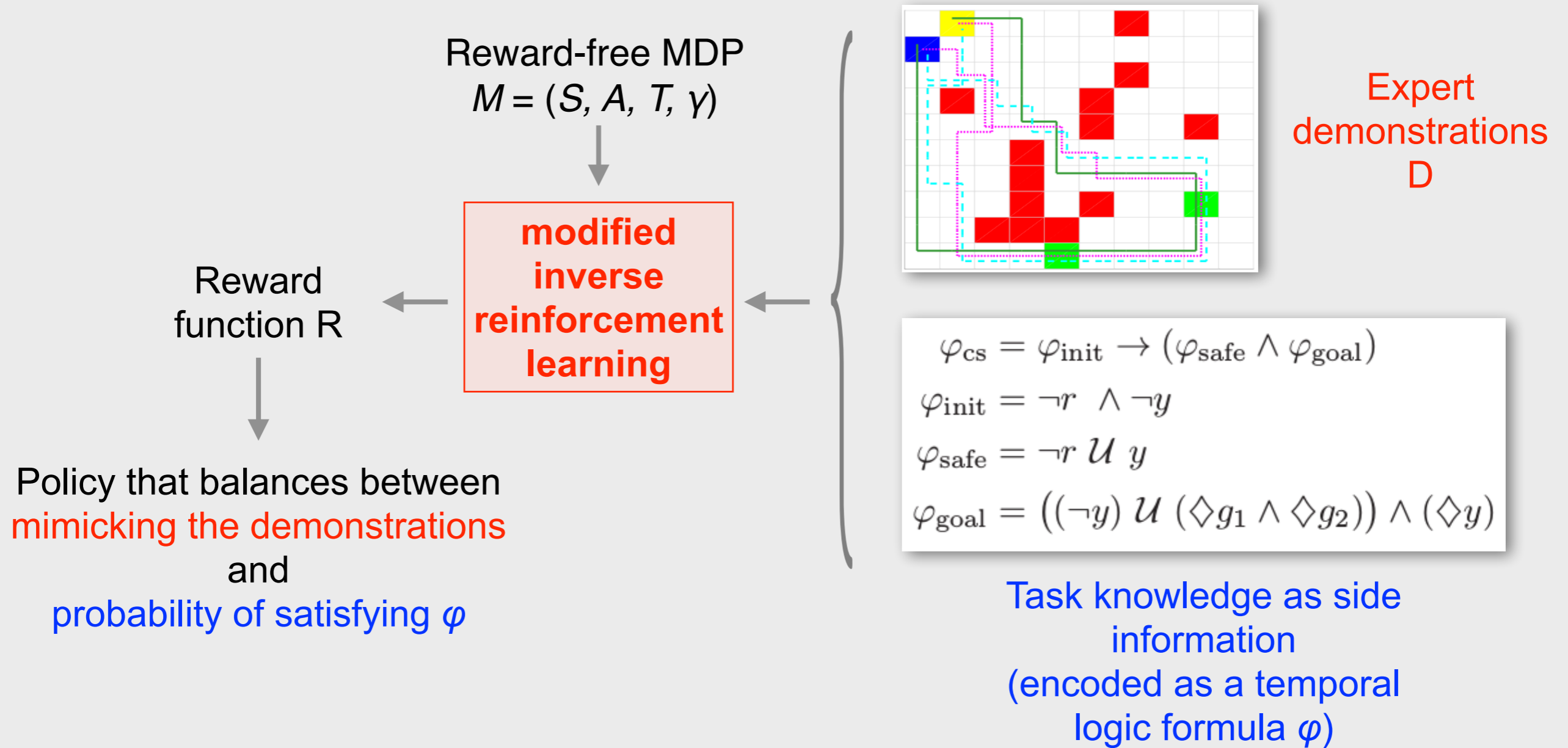
Q-learning



# Inverse reinforcement learning



# Inverse reinforcement learning with high-level task as side information



# Inverse reinforcement learning with high-level task as side information

$R = [f_1 \cdots f_k] \theta$       reward parameterized in  $\theta$  given the features  $f_1, \dots, f_k$

$Q_\pi(s, a) = \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R(s_k, a_k) \mid s_0, a_0 \right]$       expected discounted reward with policy  $\pi$  from state  $s$  and action  $a$

$\pi_Q(s, a) = \frac{\exp(Q(s, a))}{\sum_{a'} \exp(Q(s, a'))}$       softmax policy (randomized)

## Assumptions:

- Demonstrations are samples of a softmax, randomized policy.
- This unknown policy satisfies  $\varphi$  with probability higher than a threshold.

$$\begin{aligned} \min_{\theta} \quad & J^{\text{mle}}(\pi_\theta \mid M, D) \\ \text{s.t.} \quad & \text{Bellman equation}(\theta, M, D) \end{aligned}$$

negative log-likelihood

# Inverse reinforcement learning with high-level task as side information

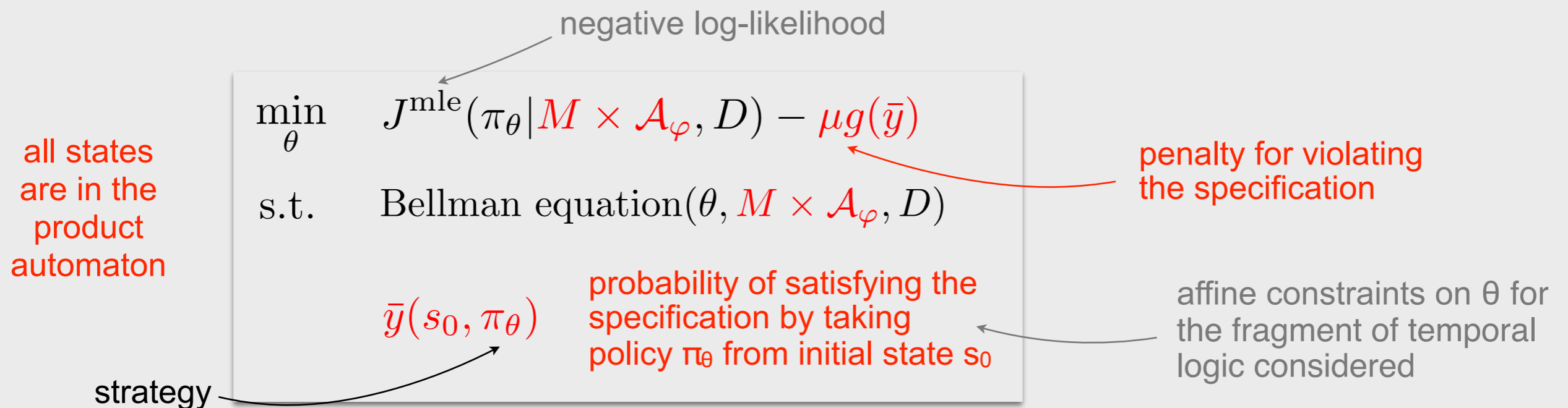
$R = [f_1 \cdots f_k] \theta$  reward parameterized in  $\theta$  given the features  $f_1, \dots, f_k$

$Q_\pi(s, a) = \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R(s_k, a_k) \mid s_0, a_0 \right]$  expected discounted reward with policy  $\pi$  from state  $s$  and action  $a$

$\pi_Q(s, a) = \frac{\exp(Q(s, a))}{\sum_{a'} \exp(Q(s, a'))}$  softmax policy (randomized)

## Assumptions:

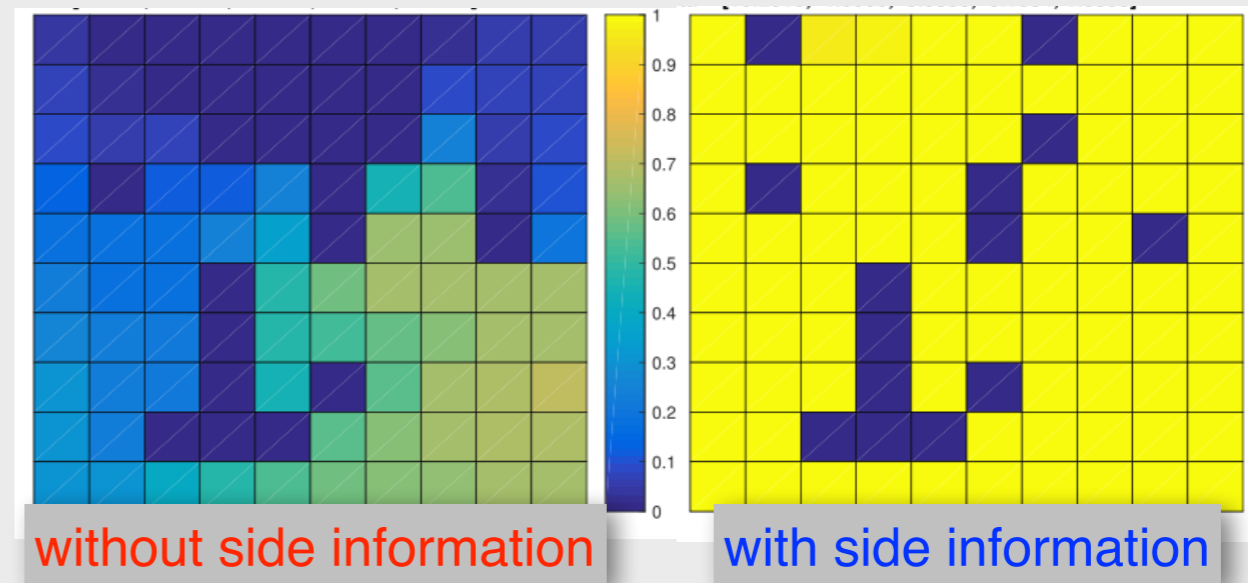
- Demonstrations are samples of a softmax, randomized policy.
- This unknown policy satisfies  $\varphi$  with probability higher than a threshold.



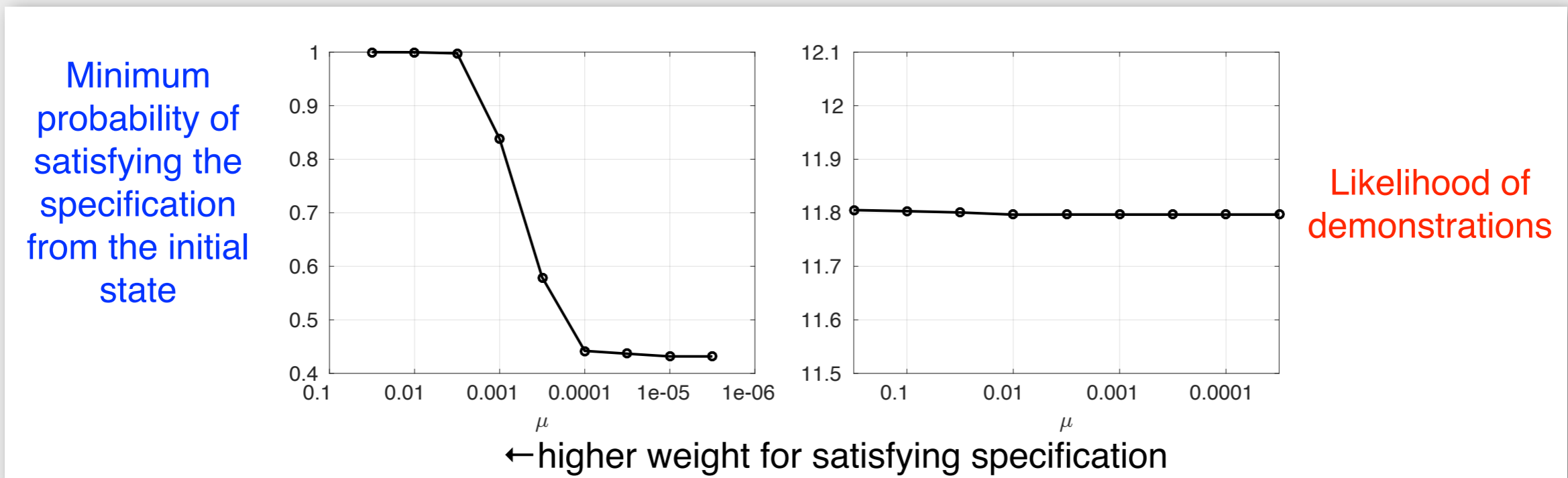


# Example

Probability of satisfying the specification for the given initial state:  
**Poor generalization without side information**



**Major increase in probability of satisfaction at negligible reduction in likelihood**



# Outline

## Correct-by-construction synthesis of hierarchical control protocols

- **formal methods** ↔ **controls**

What would it take to construct the control software of an autonomous system **in an hour**—as opposed to days or even weeks—to deliver a mission?

## Verifiable reinforcement learning

- **formal methods** ↔ **learning**

How can an autonomous system **learn** how to execute a new mission in an a priori unknown environment **efficiently** and **safely**?

## Planning in POMDPs

- **formal methods** ↔ **convex optimization**
- **formal methods** ↔ **learning**

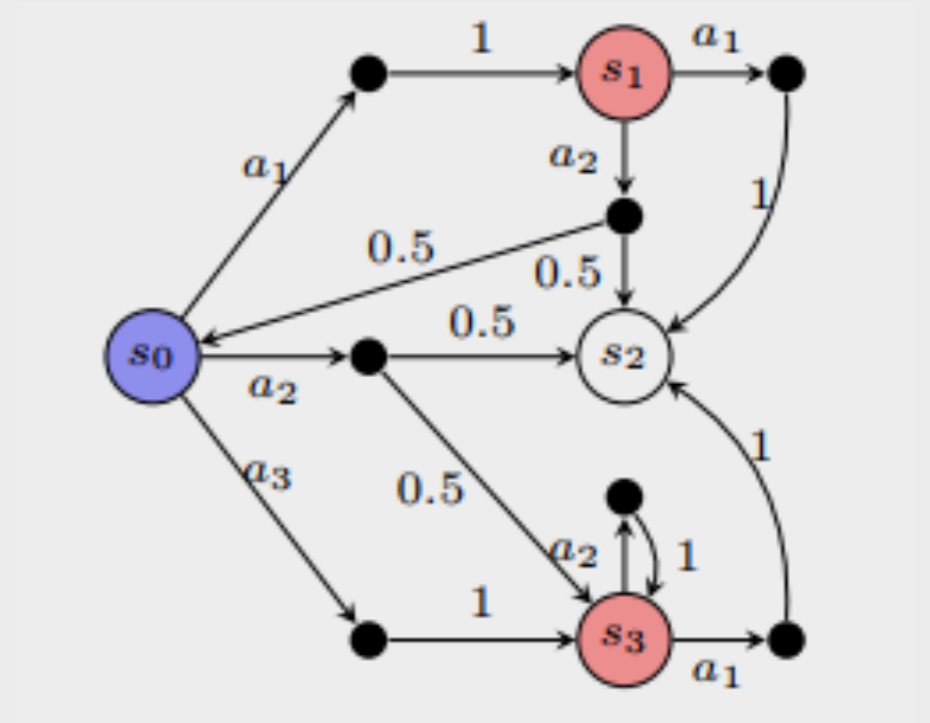
How can we cope with imperfection and/or limitations in run-time information?

# A few words on POMDPs

**POMDPs** — partially observable MDPs: Like MDPs but with limited information

$$\mathcal{M} = (S, A, P, Z, O, C)$$

set of states set of actions probabilistic transition function set of observations observation function cost function



Can distinguish only between the colors of the states

## Synthesis in POMDPs is hard.

- Undecidable as infinite memory may be necessary.
- Restriction to finite-memory strategies yields decidable yet still hard problems (and “suboptimality”).
- Finite-memory strategies: Randomized may be better than deterministic.

# Computing finite-state controllers for POMDPs by parameter synthesis

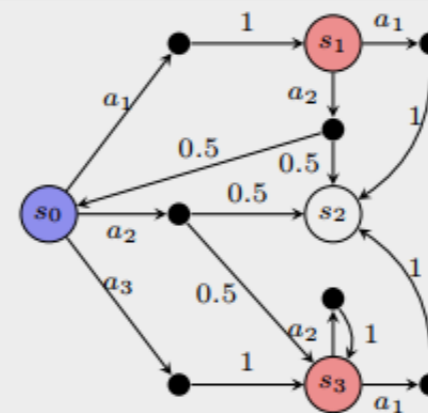
- The set of all finite-state controllers for a POMDP and a fixed memory bound can be represented by a parametric Markov chain (pMC).
- If POMDP states share an observation, the corresponding pMC states will share parameters at their transitions.
- A strategy in the POMDP, corresponds to a parameter instantiation in the pMC.

## Finite-state Controllers of POMDPs via Parameter Synthesis

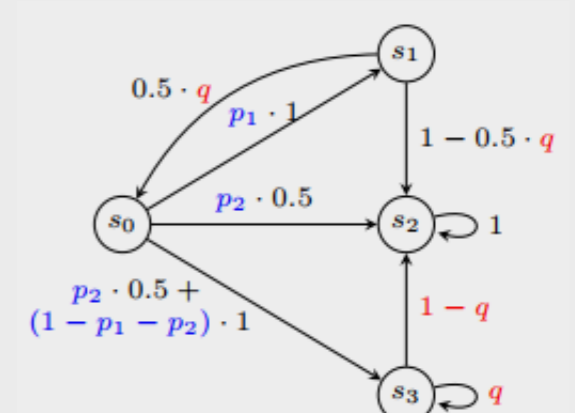
Sebastian Junges<sup>1</sup>, Nils Jansen<sup>2</sup>, Ralf Wimmer<sup>3</sup>, Tim Quatmann<sup>1</sup>,  
 Leonore Winterer<sup>3</sup>, Joost-Pieter Katoen<sup>1</sup>, and Bernd Becker<sup>3</sup>  
<sup>1</sup>RWTH Aachen University, Aachen, Germany  
<sup>2</sup>Radboud University, Nijmegen, The Netherlands  
<sup>3</sup>Albert-Ludwigs-Universität Freiburg, Freiburg im Breisgau, Germany

POMDP under FSC	pMC
states $\times$ memory	states
same observation	same parameter
strategy	parameter instantiation
permissive strategy	region of instantiations

Correspondence table between POMDP and pMCs



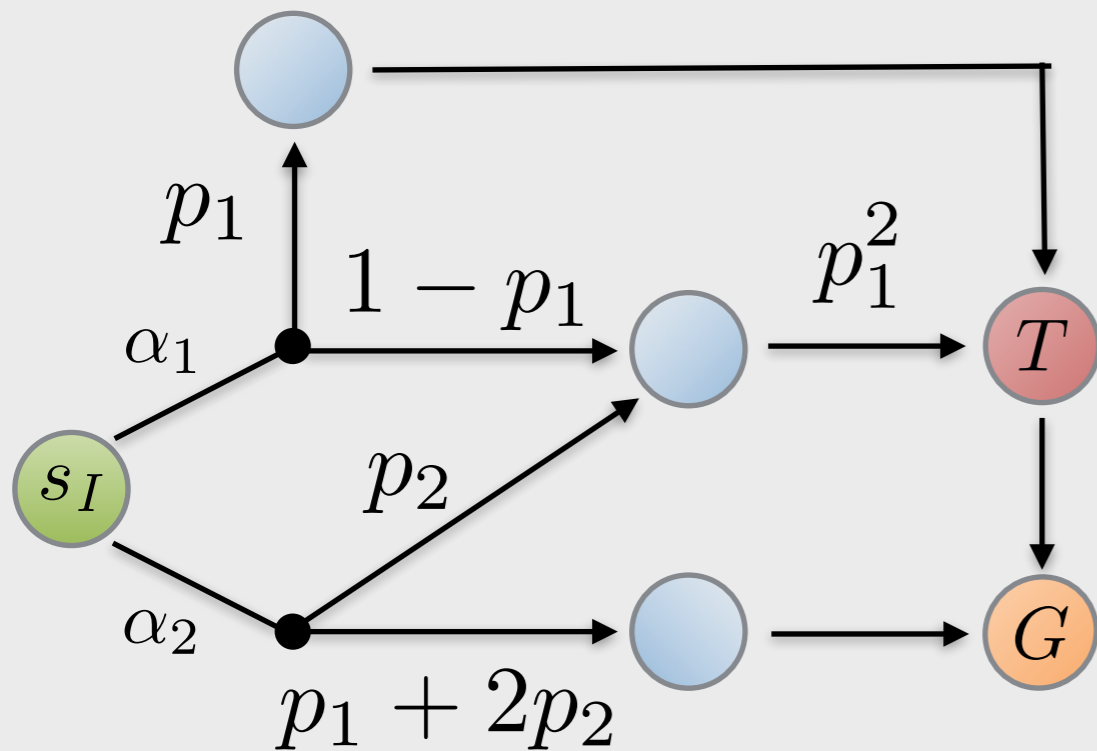
POMDP



Corresponding pMC with Parameters  $p_1, p_2$  and  $q$



# Synthesis in parametric MDPs (pMDPs)



Safety specification

$$\varphi = \mathbb{P}_{\leq \lambda}(\diamond T), \quad T \subseteq S$$

Performance specification

$$\psi = \mathbb{E}_{\leq \kappa}(\diamond G), \quad G \subseteq S$$

Objective function  $f: V \rightarrow \mathbb{R}$

Parameters  $p_1, p_2, \dots, p_n \in V$

Given pMDP  $\mathcal{M}$ , find a well-defined valuation of parameters and a scheduler  $\sigma \in \text{Sched}^{\mathcal{M}}$  such that

$$\mathcal{M}^{\sigma} \models \varphi \wedge \psi$$

and value for objective function  $f: V \rightarrow \mathbb{R}$  is minimal.

# Solution as a “nonlinear program”

minimize  $f$  **objective function  
over parameters**  
subject to

$$p_{s_I} \leq \lambda$$

$$c_{s_I} \leq \kappa$$

**safety and  
performance  
specifications**

$$\forall s \in S. \quad \sum_{\alpha \in Act(s)} \sigma^{s,\alpha} = 1$$

**well-defined  
schedulers and  
parameter  
instantiations**

$$\forall s \in S \quad \forall \alpha \in Act(s). \quad \sum_{s' \in S} \mathcal{P}(s, \alpha, s') = 1$$

$$\forall s \in T. \quad p_s = 1$$

$$\forall s \in S \setminus T. \quad p_s = \sum_{\alpha \in Act(s)} \sigma^{s,\alpha} \cdot \sum_{s' \in S} \mathcal{P}(s, \alpha, s') \cdot p_{s'}$$

**safety  
probability  
computation**

$$\forall s \in G. \quad c_s = 0$$

**expected performance computation**

$$\forall s \in S \setminus G. \quad c_s = \sum_{\alpha \in Act(s)} \sigma^{s,\alpha} \cdot \left( c(s, \alpha) + \sum_{s' \in S} \mathcal{P}(s, \alpha, s') \cdot c_{s'} \right)$$

## Problem variables:

Randomized scheduler:

$$\{\sigma^{s,\alpha} \mid s \in S, \alpha \in Act(s)\}$$

Probability of reaching T:

$$\{p_s \mid s \in S\}$$

Expected cost of reaching G:

$$\{c_s \mid s \in S\}$$

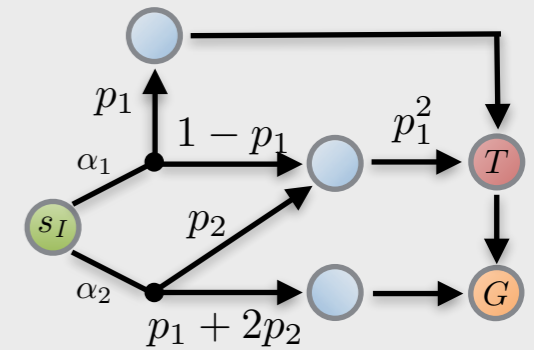
# A useful observation

non-negative-valued variable

signomials

non-negative-valued variable

$$p_s = \sum_{\alpha \in Act(s)} \sigma^{s, \alpha} \cdot \sum_{s' \in S} \mathcal{P}(s, \alpha, s') \cdot p_{s'}$$



$$f = \sum_{k=1}^K c_k \cdot \underbrace{x_1^{a_{1k}} \cdots x_n^{a_{nk}}}_{\text{monomial}}$$

strictly positive

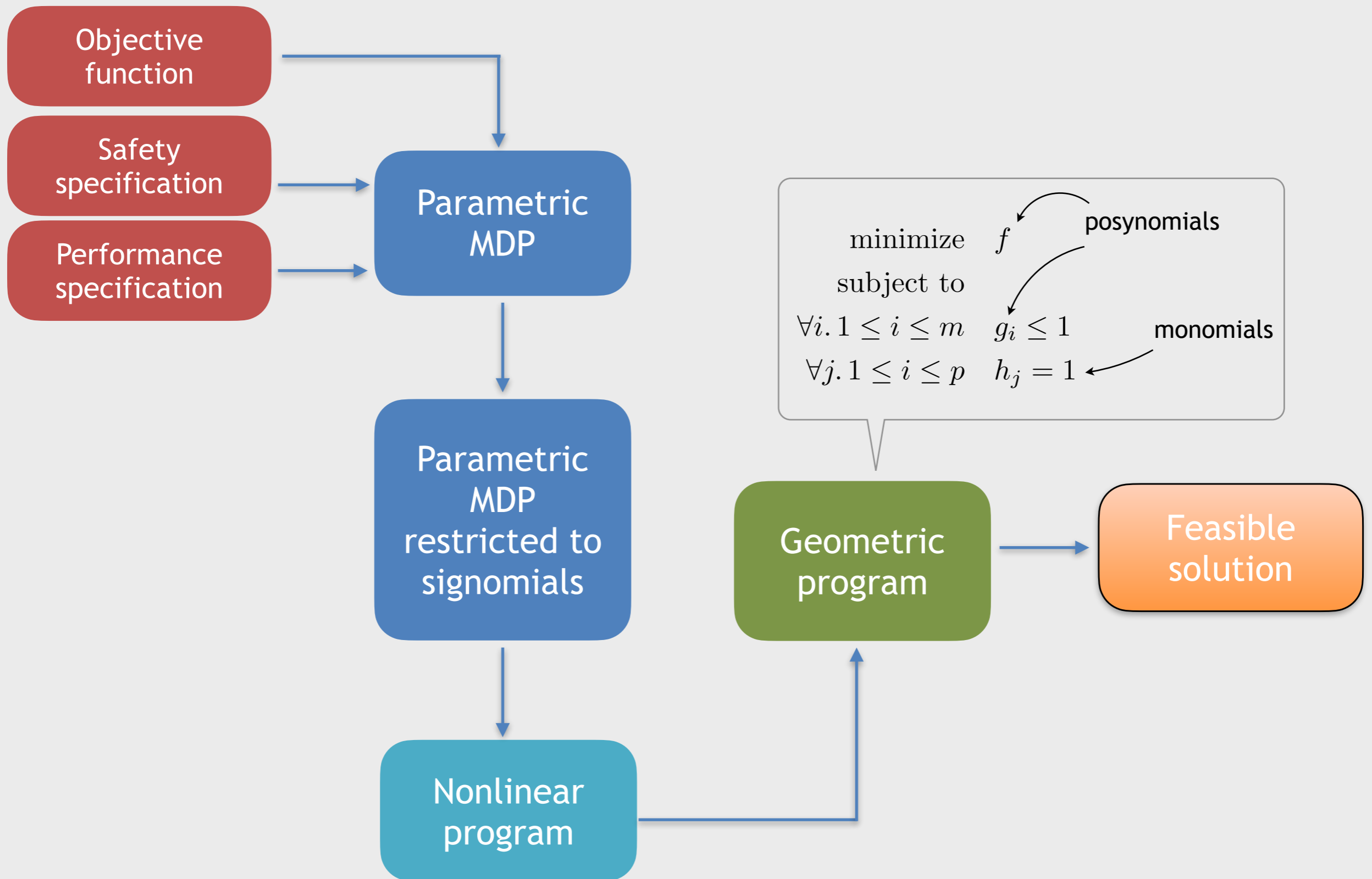
real

$> 0$ : posynomial

no restriction: signomial

**Question:** Can we somehow exploit this structure and solve the parameter synthesis problem as a convex optimization problem (maybe bunch of them)?

# Workflow





# Convexification: relaxation + lifting

$$p_s = \sum_{\alpha \in Act(s)} \sigma^{s,\alpha} \cdot \sum_{s' \in S} \mathcal{P}(s, \alpha, s') \cdot p_{s'}$$



$$\frac{\sum_{\alpha \in Act(s)} \sigma^{s,\alpha} \cdot \sum_{s' \in S} \mathcal{P}(s, \alpha, s') \cdot p_{s'}}{p_s} \leq 1$$

still signomials  
(not geometric  
program yet)

(assuming  $p_s \neq 0$ )

upper bound  
on actual probability

lifting: introduce  
new variables

$$\mathcal{P}(s, \alpha, \bar{s}) = 1 - \sum_{s' \in S \setminus \{\bar{s}\}} \mathcal{P}(s, \alpha, s')$$

signomial

posynomials

Turns out that all but one  
transition probability can  
be written as a  
posynomial.

# Geometric program (with relaxation tightening)

minimize  $\sum_{p \in V} \frac{1}{p} + \sum_{\bar{p} \in L} \frac{1}{\bar{p}} + \sum_{s \in S, \alpha \in Act(s)} \frac{1}{\sigma_{s,\alpha}}$  regularization

subject to

$$\frac{p_{s_I}}{\lambda} \leq 1$$

$$\frac{c_{s_I}}{\kappa} \leq 1$$

$$\forall s \in S. \quad \sum_{\alpha \in Act(s)} \sigma^{s,\alpha} \leq 1$$

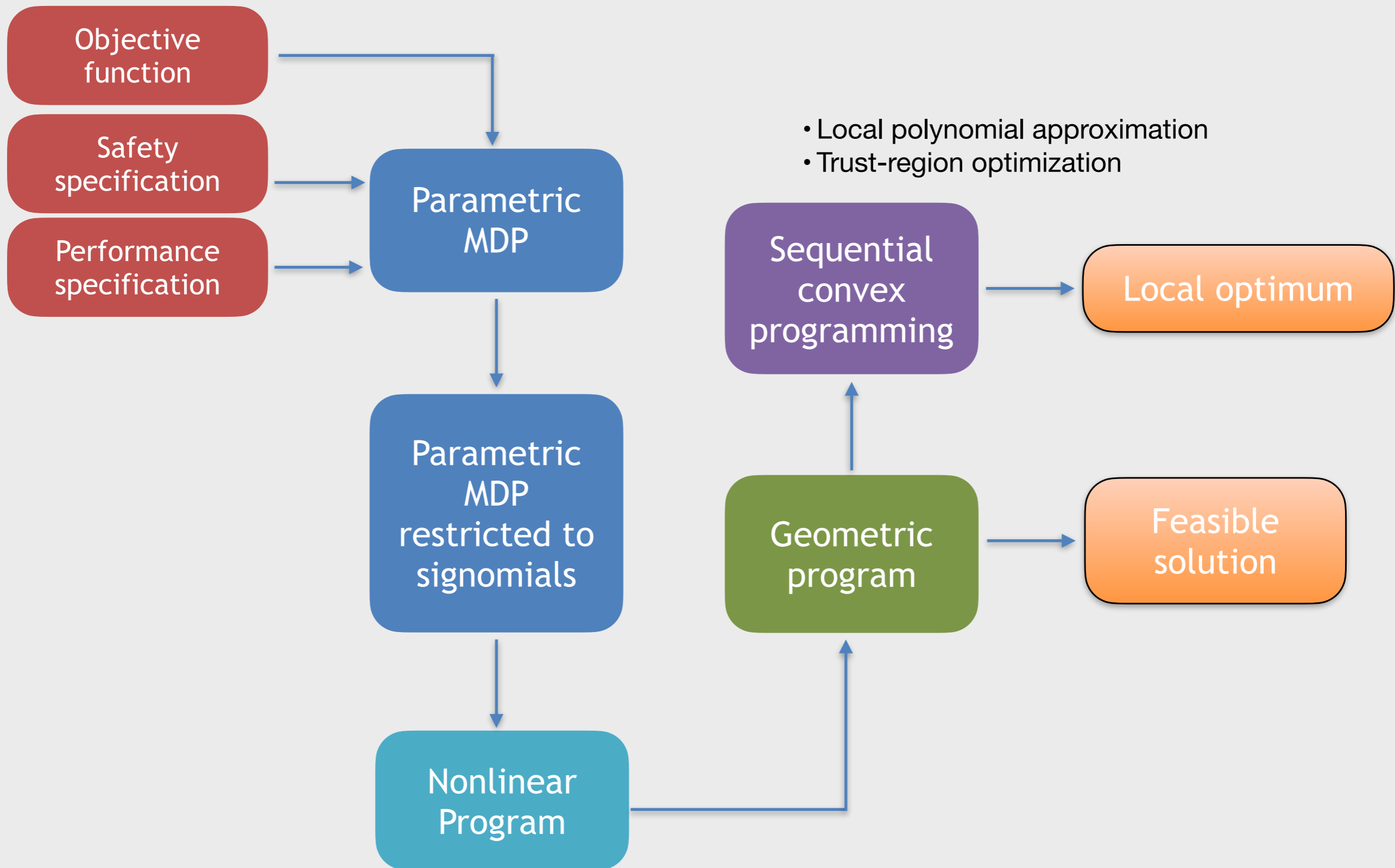
$$\forall s \in S \forall \alpha \in Act(s). \quad \sum_{s' \in S} \mathcal{P}(s, \alpha, s') \leq 1$$

$$\forall s \in S \setminus T. \quad \frac{\sum_{\alpha \in Act(s)} \sigma^{s,\alpha} \cdot \sum_{s' \in S} \mathcal{P}(s, \alpha, s') \cdot p_{s'}}{p_s} \leq 1$$

$$\forall s \in S \setminus G. \quad \frac{\sum_{\alpha \in Act(s)} \sigma^{s,\alpha} \cdot \left( c(s, \alpha) + \sum_{s' \in S} \mathcal{P}(s, \alpha, s') \cdot c_{s'} \right)}{c_s} \leq 1$$

**Theorem:** The solution to the geometric program gives a well-defined scheduler and parameter instantiation. But it may be sub-optimal.

# Workflow



# Compare against benchmarks

Alternative tools for optimization TO even in the smallest instances

Benchmark	#states	#par	specs	proposed method		only feasibility
				MOSEK (s)	Z3	
BRP (pMC)	5382	2	EC, $\mathbb{P}$ , *	23.17	(6.48)	—
	112646	2	EC, $\mathbb{P}$ , *	3541.59	(463.74)	—
	112646	4	EC, $\mathbb{P}$ , *	4173.33	(568.79)	—
	5382	2	EC, $\mathbb{P}$	3.61		904.11
	112646	2	EC, $\mathbb{P}$	479.08		TO
NAND (pMC)	4122	2	EC, $\mathbb{P}$ , *	14.67	(2.51)	—
	35122	2	EC, $\mathbb{P}$ , *	1182.41	(95.19)	—
	4122	2	EC, $\mathbb{P}$	1.25		1.14
	35122	2	EC, $\mathbb{P}$	106.40		11.49
BRP (pMDP)	5466	2	EC, $\mathbb{P}$ , *	31.04	(8.11)	—
	112846	2	EC, $\mathbb{P}$ , *	4319.16	(512.20)	—
	5466	2	EC, $\mathbb{P}$	4.93		1174.20
	112846	2	EC, $\mathbb{P}$	711.50		TO
CONS (pMDP)	4112	2	EC, $\mathbb{P}$ , *	102.93	(1.14)	—
	65552	2	EC, $\mathbb{P}$ , *	TO		—
	4112	2	EC, $\mathbb{P}$	6.13		TO
	65552	2	EC, $\mathbb{P}$	1361.96		TO



# More benchmarks

(using a related but different method)

parametric  
MCs

POMDPs

parametric  
MDPs

Set	Problem		Info			PSO			SMT	CCP		
	Inst	Spec	States	Trans.	Par.	tmin	tmax	tavg	t	t	solv	iter
Brp	16,2	$\mathbb{P}_{\leq 0.1}$	98	194	<b>2</b>	0	0	<b>0</b>	40	0	30%	3
	512,5	$\mathbb{P}_{\leq 0.1}$	6146	12290	<b>2</b>	24	36	<b>28</b>	TO	33	24%	3
	10,5	$\mathbb{P}_{\leq 0.1}$	42	82	<b>2</b>	4	5	5	8	4	2%	4
	5,10	$\mathbb{P}_{\leq 0.05}$	10492	20982	<b>2</b>	21	51	28	TO	<b>22</b>	21%	2
Zeroconf	10000	$\mathbb{E}_{\leq 10010}$	10003	20004	<b>2</b>	2	4	<b>3</b>	TO	57	81%	3
GridA	4	$\mathbb{P}_{\geq 0.84}$	1026	2098	<b>72</b>	11	11	<b>11</b>	TO	22	81%	11
GridB	8,5	$\mathbb{P}_{\geq 0.84}$	8653	17369	<b>700</b>	409	440	427	TO	<b>213</b>	84%	8
GridB	10,6	$\mathbb{P}_{\geq 0.84}$	16941	33958	<b>1290</b>	533	567	553	TO	<b>426</b>	84%	7
GridC	6	$\mathbb{E}_{\leq 4.8}$	1665	305	<b>168</b>	261	274	267	TO	<b>169</b>	90%	23
Maze	5	$\mathbb{E}_{\leq 14}$	1303	2658	<b>590</b>	213	230	219	TO	<b>67</b>	89%	8
Maze	5	$\mathbb{E}_{\leq 6}$	1303	2658	<b>590</b>	-	-	TO	TO	<b>422</b>	85%	97
...	7	$\mathbb{E}_{\leq 6}$	2580	5233	<b>1176</b>	-	-	TO	TO	<b>740</b>	90%	60
...	5,2	$\mathbb{E}_{\leq 11.5}$	21746	63158	<b>2420</b>	312	523	359	TO	<b>207</b>	39%	3
Netw	5,2	$\mathbb{E}_{\leq 10.5}$	21746	63158	<b>2420</b>	-	-	TO	TO	<b>210</b>	38%	4
Netw	4,3	$\mathbb{E}_{\leq 11.5}$	38055	97335	<b>4545</b>	-	-	TO	TO	MO	-	-
Repud	8,5	$\mathbb{P}_{\geq 0.1}$	1487	3002	<b>360</b>	16	22	18	TO	4	36%	2
Repud	8,5	$\mathbb{P}_{\leq 0.05}$	1487	3002	<b>360</b>	273	324	293	TO	<b>14</b>	72%	4
Repud	16,2	$\mathbb{P}_{\leq 0.01}$	790	1606	<b>96</b>	-	-	TO	TO	<b>15</b>	78%	9
Repud	16,2	$\mathbb{P}_{\geq 0.062}$	790	1606	<b>96</b>	-	-	TO	TO	TO	-	-

## Synthesis in pMDPs: A Tale of 1001 Parameters

Murat Cubuktepe<sup>1</sup>, Nils Jansen<sup>2</sup>, Sebastian Junges<sup>3</sup>,  
Joost-Pieter Katoen<sup>3</sup>, Ufuk Topcu<sup>1</sup>

<sup>1</sup> The University of Texas at Austin, Austin, TX, USA\*

<sup>2</sup> Radboud University, Nijmegen, The Netherlands

<sup>3</sup> RWTH Aachen University, Aachen, Germany\*\*

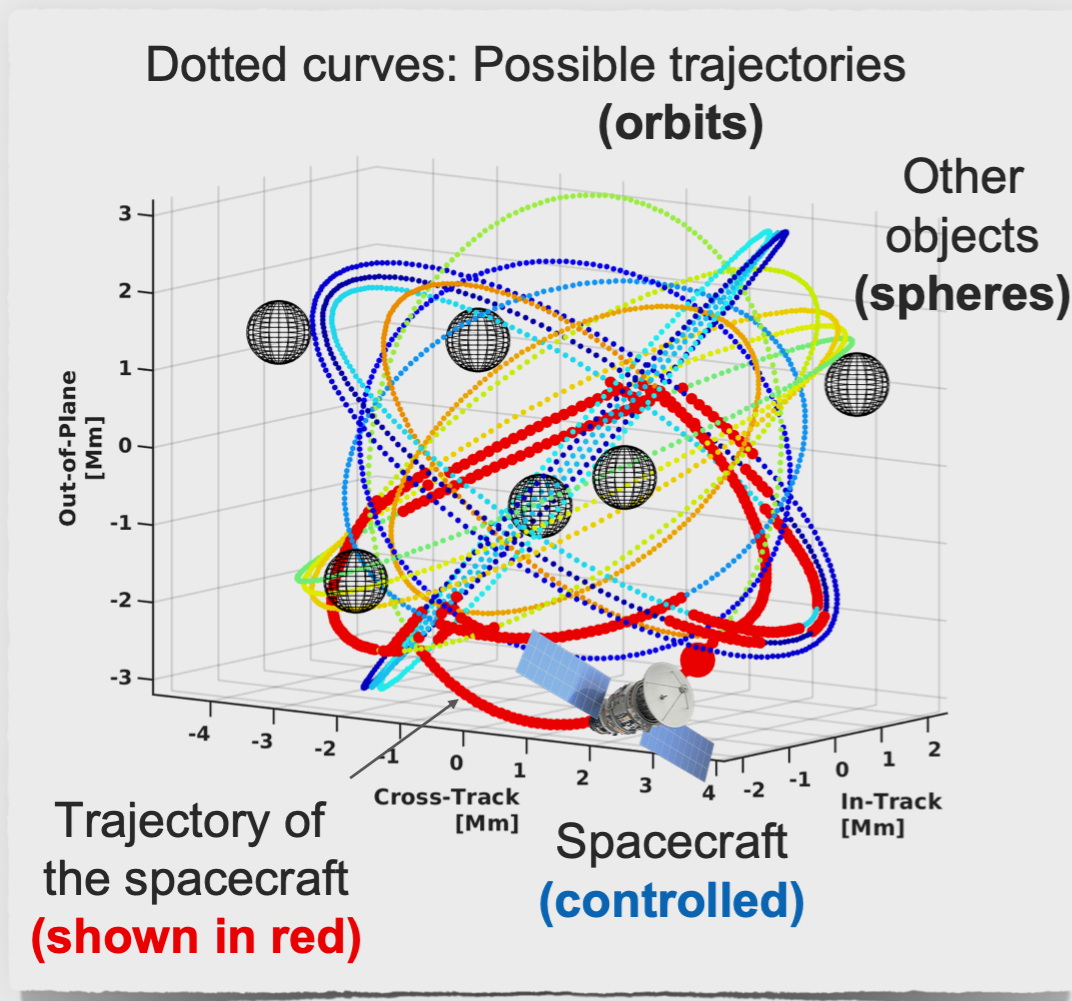
**Abstract.** This paper considers parametric Markov decision processes (pMDPs) whose transitions are equipped with affine functions over a finite set of parameters. The synthesis problem is to find a parameter valuation such that the instantiated pMDP satisfies a (temporal logic) specification under all strategies. We show that this problem can be formulated as a quadratically-constrained quadratic program (QCQP) and is non-convex in general. To deal with the NP-hardness of such problems, we exploit a convex-concave procedure (CCP) to iteratively obtain local optima. An appropriate interplay between CCP solvers and probabilistic model checkers creates a procedure — realized in the tool PROPheSY — that solves the synthesis problem for models with thousands of parameters.

## 1 Introduction

*The parameter synthesis problem.* Probabilistic model checking concerns the automatic verification of models such as Markov decision processes (MDPs). Unremitting improvements in algorithms and efficient tool implementations [14,22,26] have opened up a wide variety of applications, most notably in dependability, security, and performance analysis as well as systems biology. However, at early development stages, certain system quantities such as fault or reaction rates are often not fully known. This lack of information gives rise to parametric models where transitions are functions over real-valued parameters [12,21,27], forming symbolic descriptions of (uncountable) families of concrete MDPs. The *parameter synthesis problem* is: Given a finite-state parametric MDP, find a parameter

Set	Problem			Info			PSO			SMT	CCP	
	Inst	Spec	States	Act	Trans.	Par.	tmin	tmax	tavg	t	t	solv iter
BRP	4,128	$\mathbb{P}_{\leq 0.1}$	17131	17396	23094	<b>2</b>	45	47	46	TO	<b>39</b>	33% 4
Coin	32	$\mathbb{E}_{\leq 500}$	4112	6160	7692	<b>2</b>	117	119	<b>118</b>	TO	TO	- -
CoinX	32	$\mathbb{E}_{\leq 210}$	16448	24640	30768	<b>2</b>	1196	1222	1208	TO	<b>32</b>	78% 3
Zeroconf	1	$\mathbb{P}_{\geq 0.99}$	31402	55678	70643	<b>3</b>	18	19	<b>19</b>	TO	79	82% 2
CSMA	2,4	$\mathbb{E}_{\leq 69.3}$	7958	7988	10594	<b>26</b>	n.s.	n.s.	n.s.	TO	<b>79</b>	86% 10
Virus	-	$\mathbb{E}_{\leq 10}$	809	3371	6741	<b>18</b>	113	113	113	TO	<b>13</b>	76% 4
Wlan	0	$\mathbb{E}_{\leq 580}$	2954	3972	5202	<b>15</b>	n.s.	n.s.	n.s.	TO	<b>7</b>	72% 2

# Uncertain POMDPs, through a more visual example



Spacecraft motion planning:

Switching between orbits is possible if the orbits are close to each other, but **costs fuel**.

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

induced uncertain Markov chain

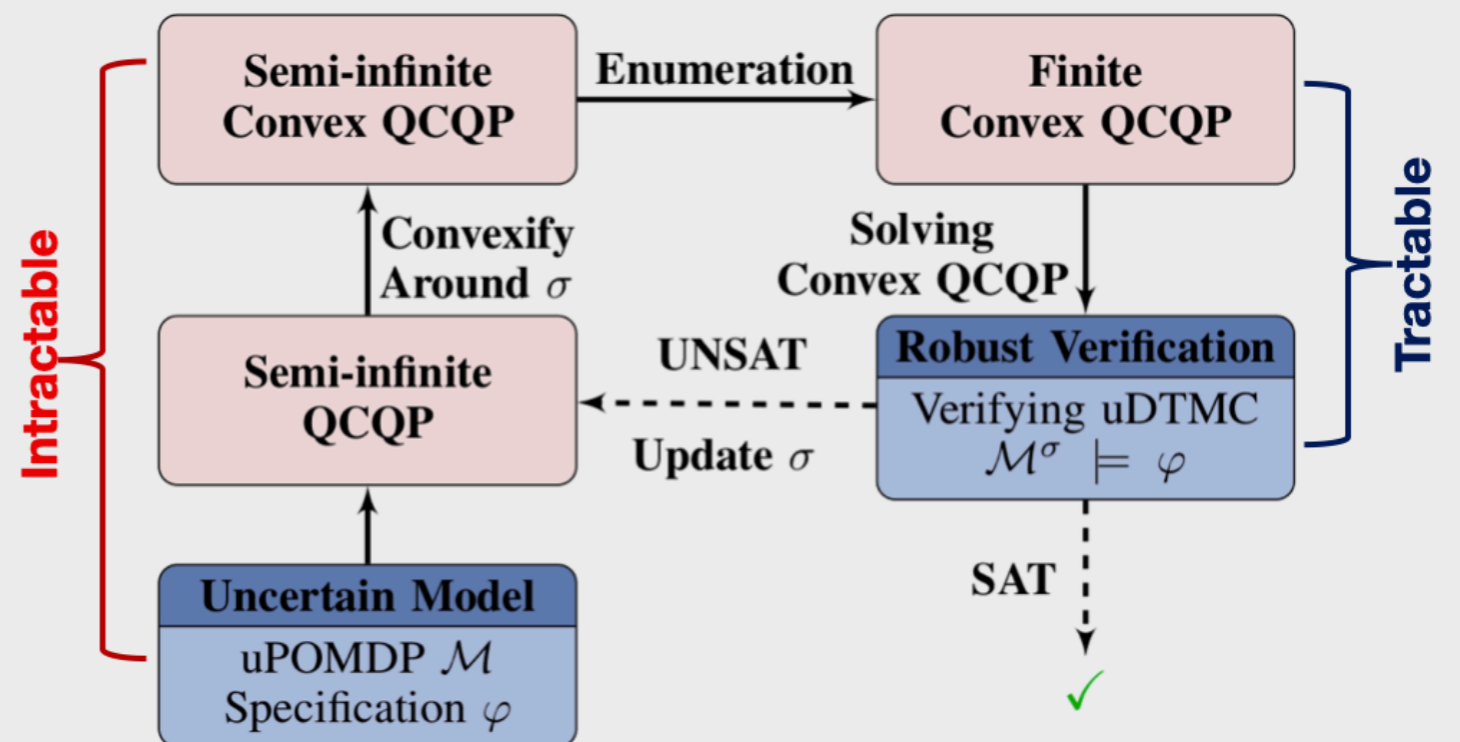
$$\mathcal{M}_\sigma^P \models \varphi$$

satisfies the specification

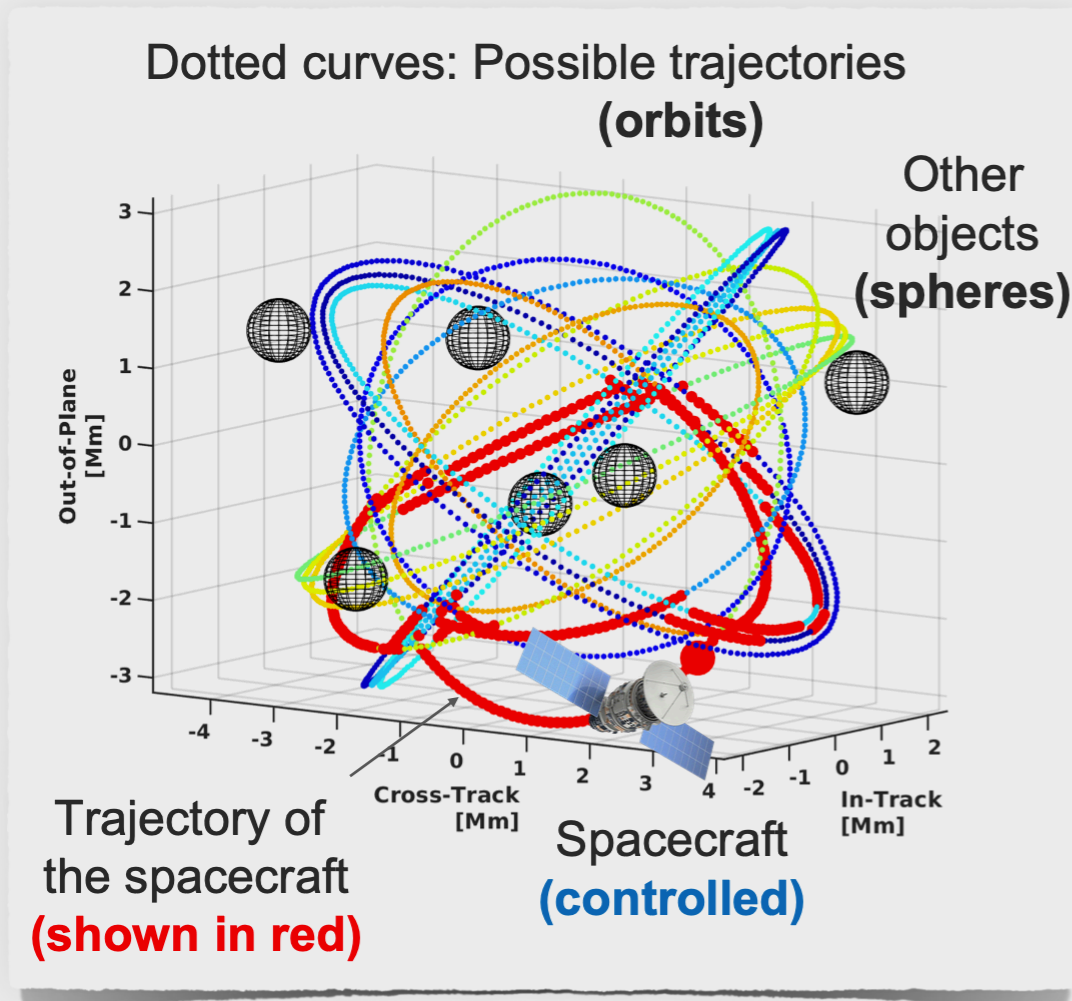
transition function

$$\text{for all } P \in \mathcal{P}$$

uncertainty set



# Uncertain POMDPs, through a more visual example



Spacecraft motion planning:

Switching between orbits is possible if the orbits are close to each other, but **costs fuel**.

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

induced uncertain Markov chain

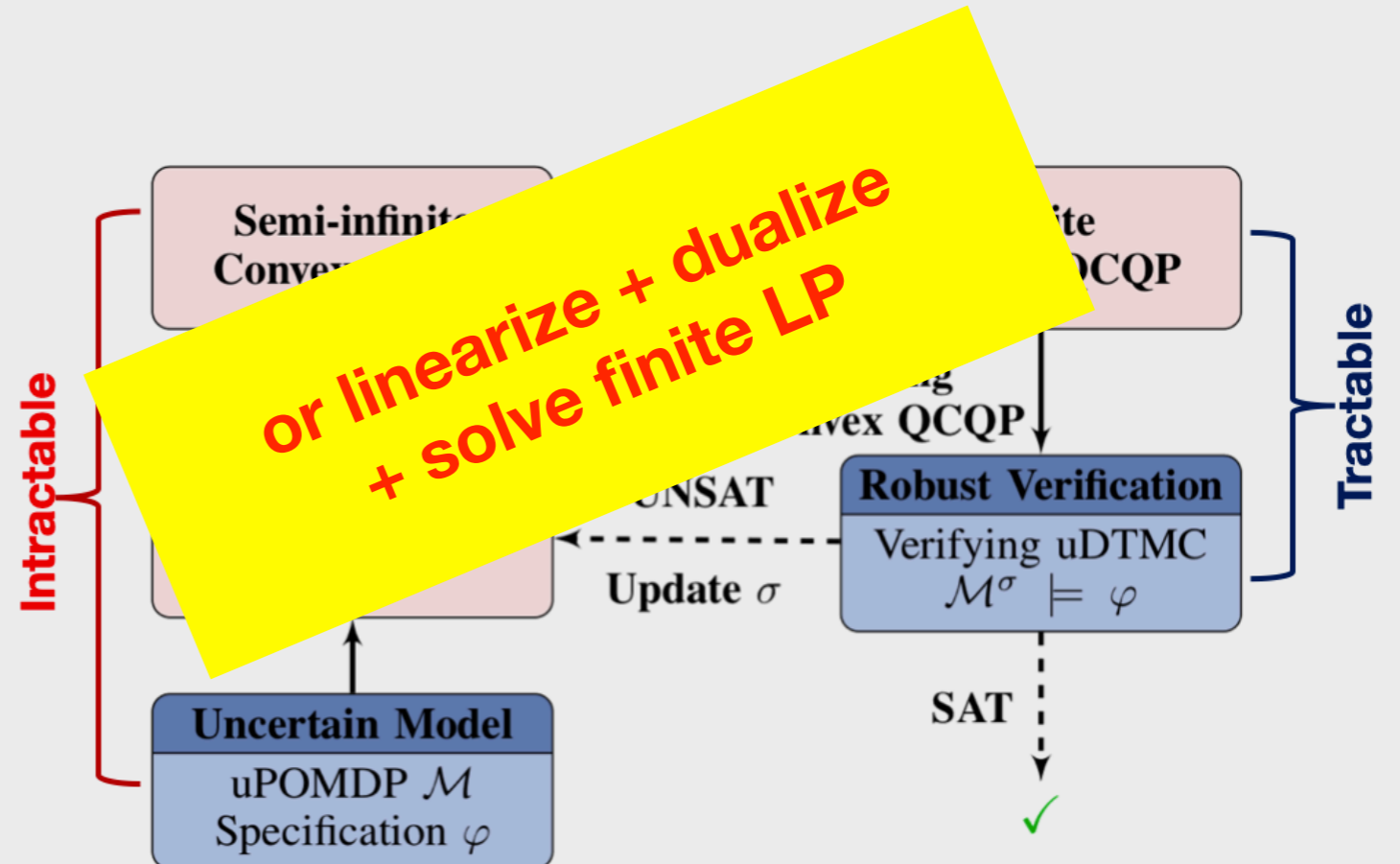
$$\mathcal{M}_\sigma^P \models \varphi$$

satisfies the specification

transition function

$$\text{for all } P \in \mathcal{P}$$

uncertainty set



# Outline

## Correct-by-construction synthesis of hierarchical control protocols

- **formal methods** ↔ **controls**

What would it take to construct the control software of an autonomous system **in an hour**—as opposed to days or even weeks—to deliver a mission?

## Verifiable reinforcement learning

- **formal methods** ↔ **learning**

How can an autonomous system **learn** how to execute a new mission in an a priori unknown environment **efficiently** and **safely**?

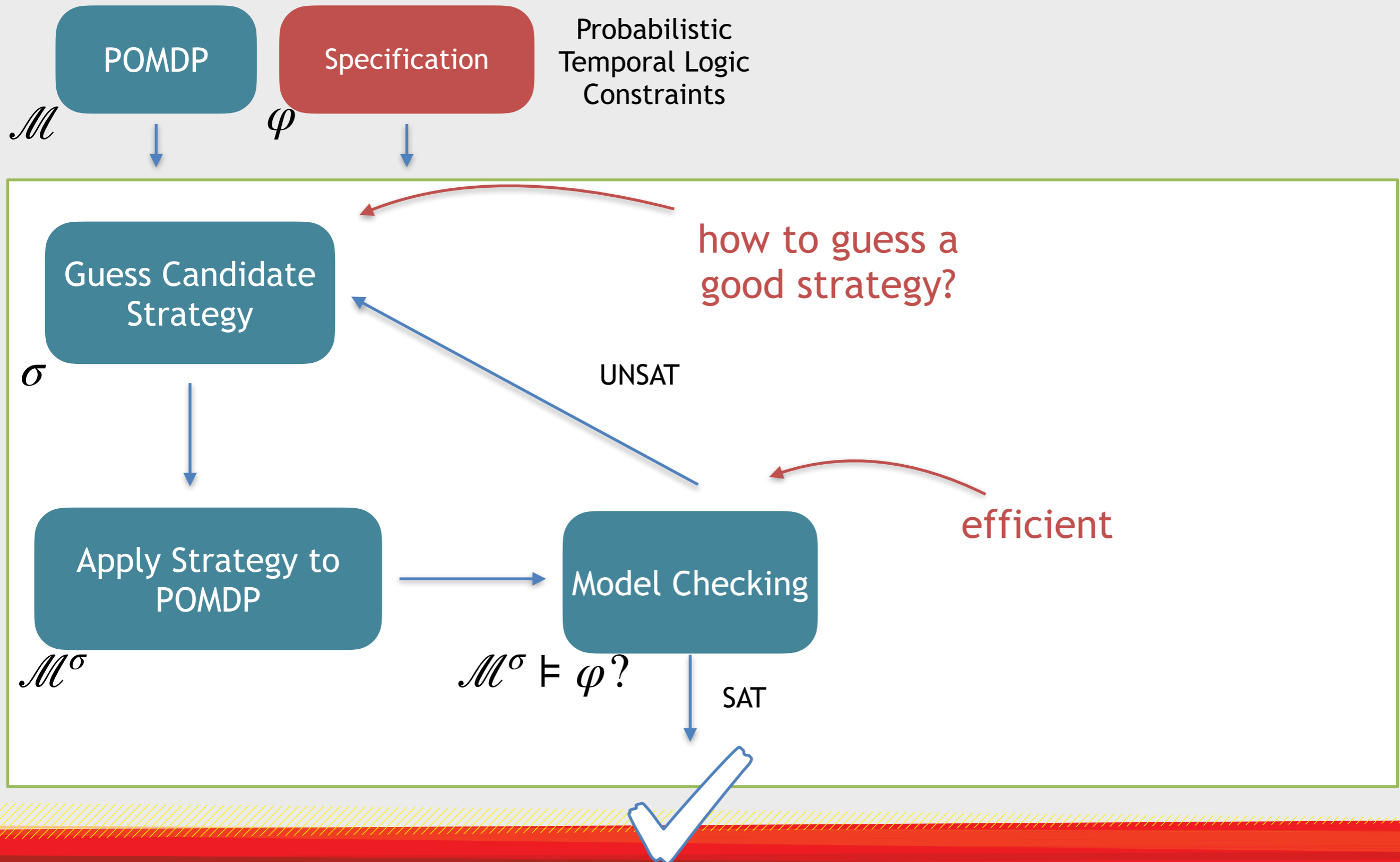
## Planning in POMDPs

- **formal methods** ↔ **convex optimization**
- **formal methods** ↔ **learning**

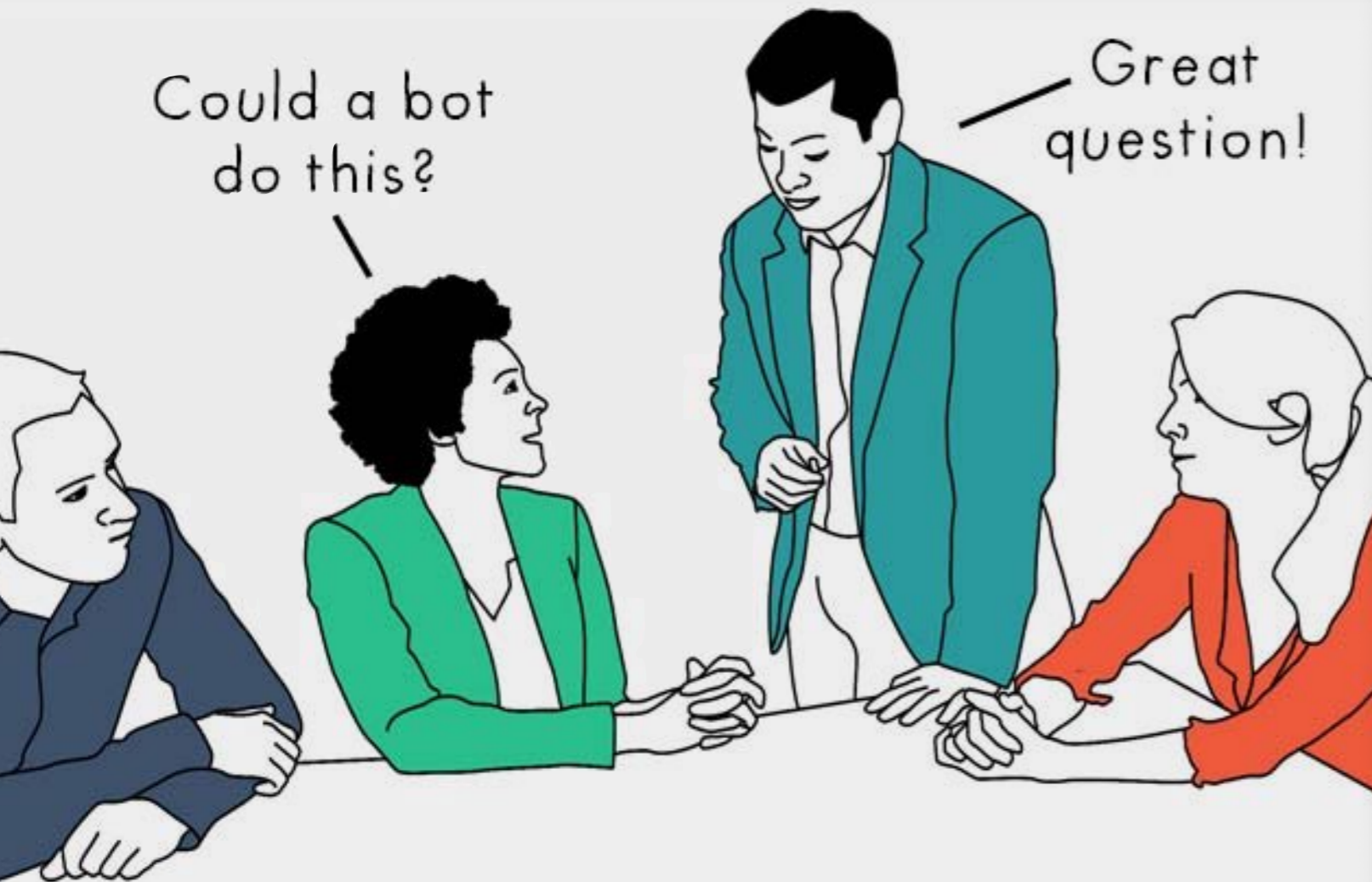
How can we cope with imperfection and/or limitations in run-time information?



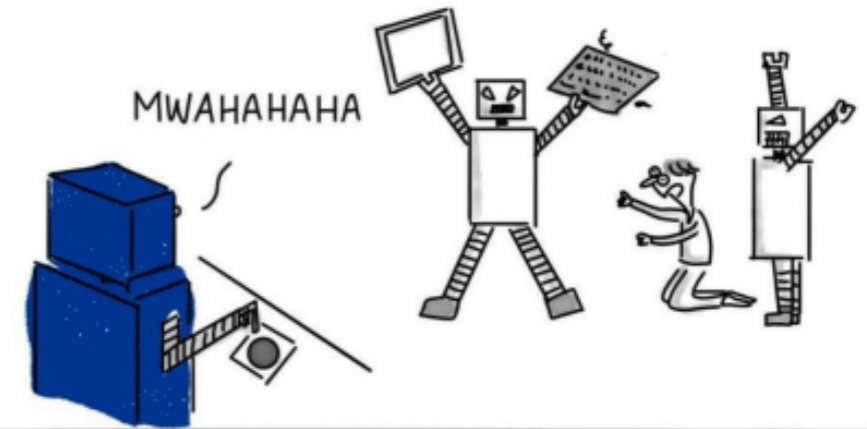
# Synthesis is hard. Guess a strategy and verify!



# How to guess a "good" strategy?



WHAT PEOPLE FEAR AI WILL DO



WHAT PEOPLE SAY AI CAN DO



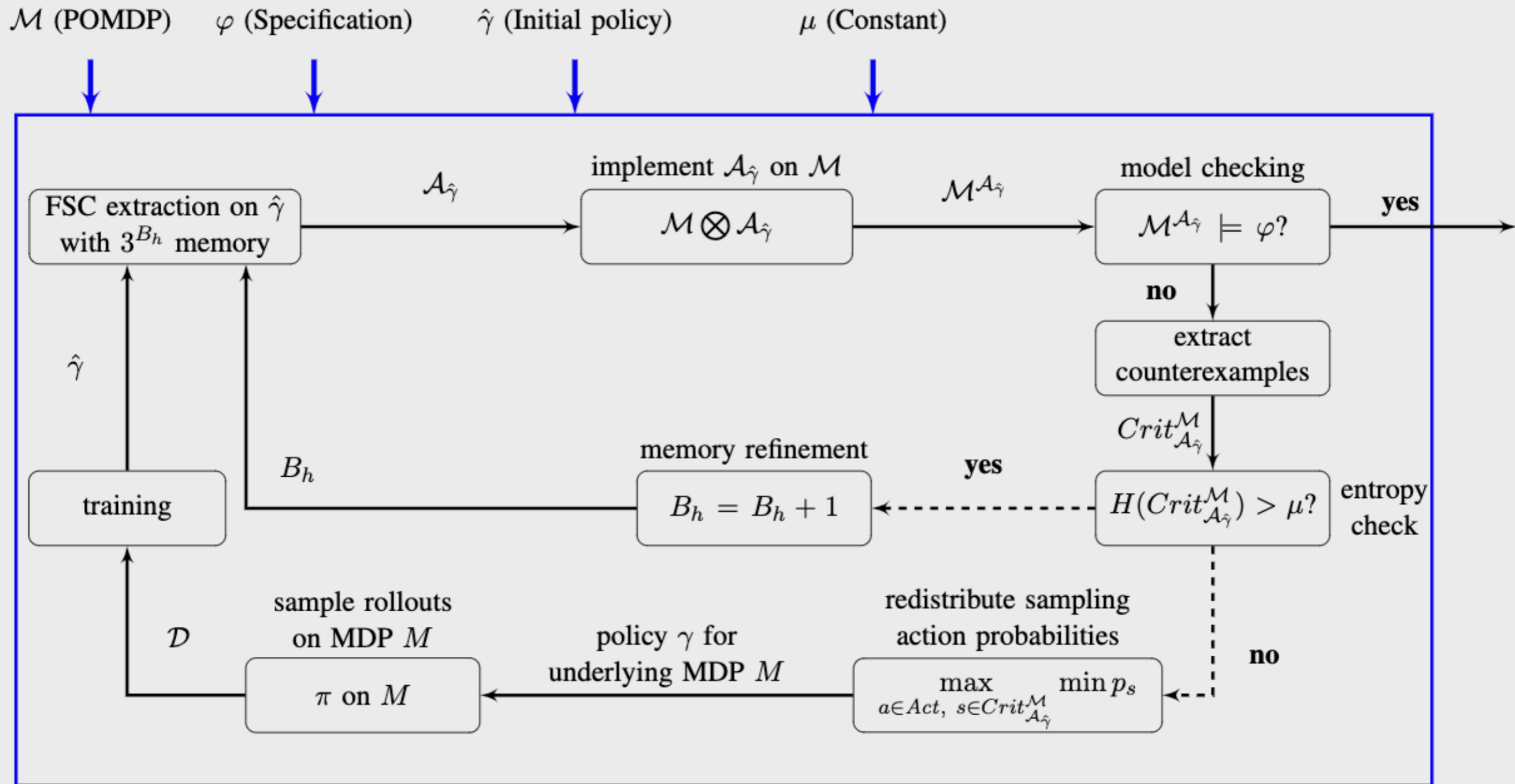
MEANWHILE...



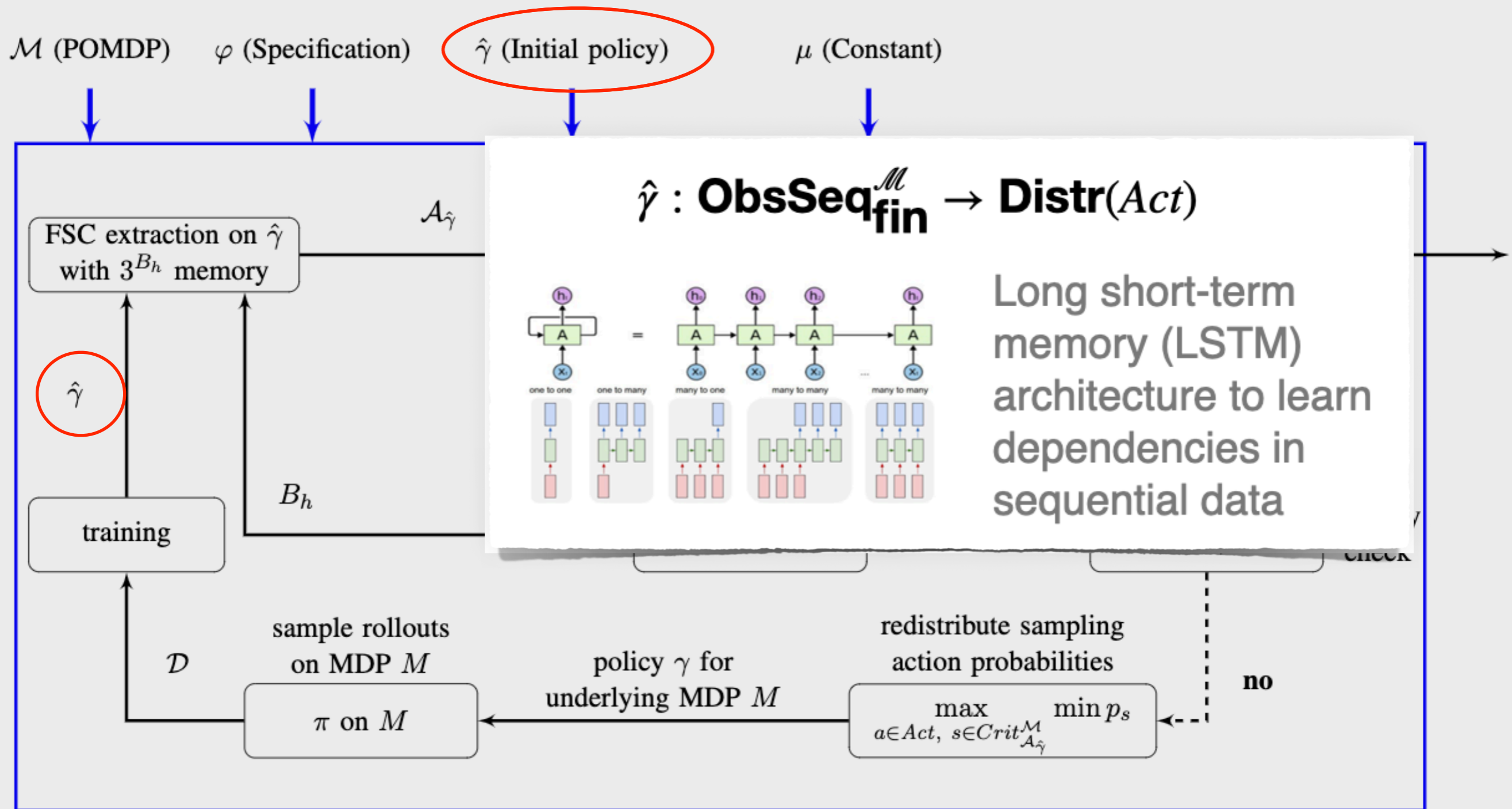
OURSKY · HK

FTFY

# Learn a guess and improve (with help from verification).



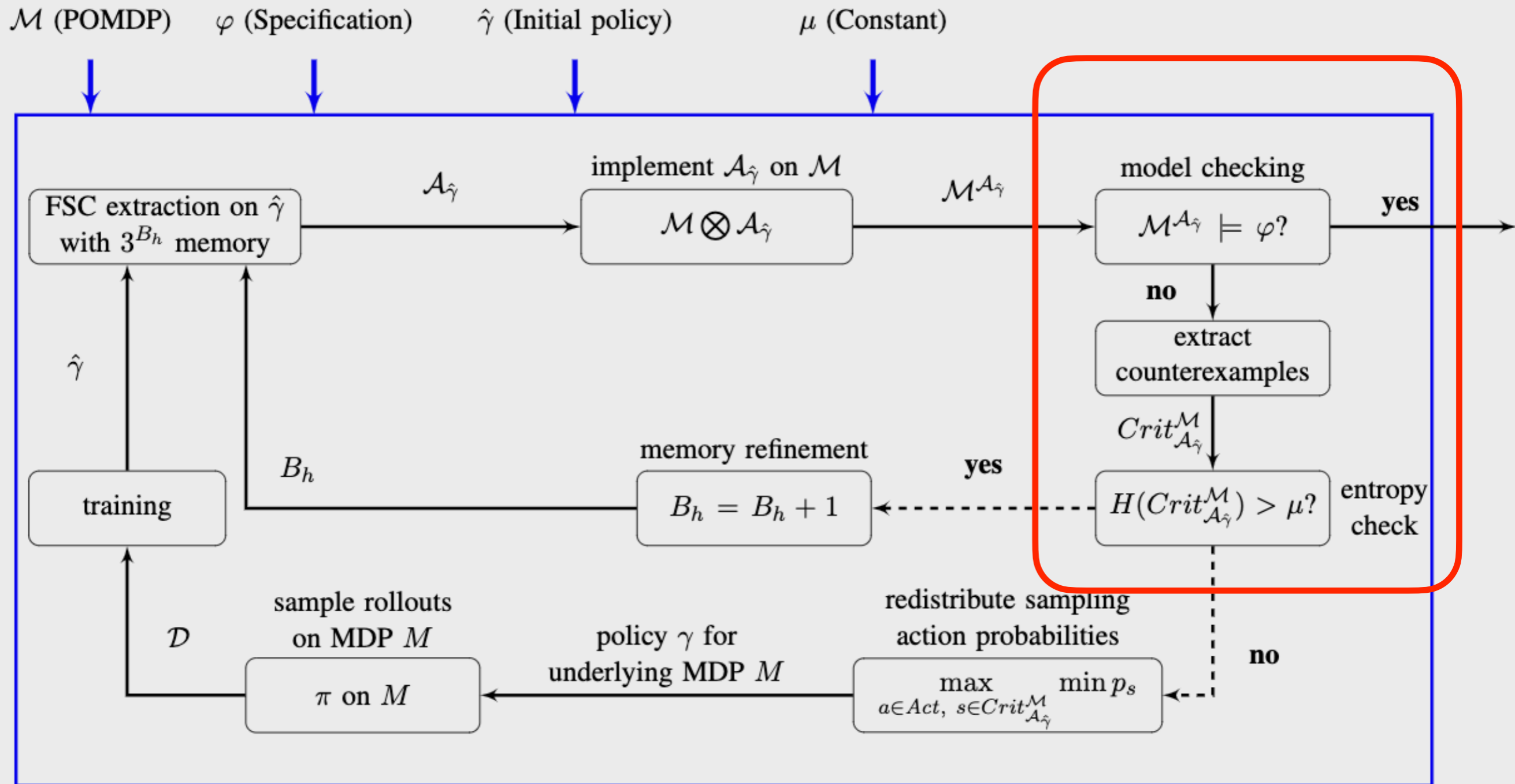
# Learn a guess and improve (with help from verification).



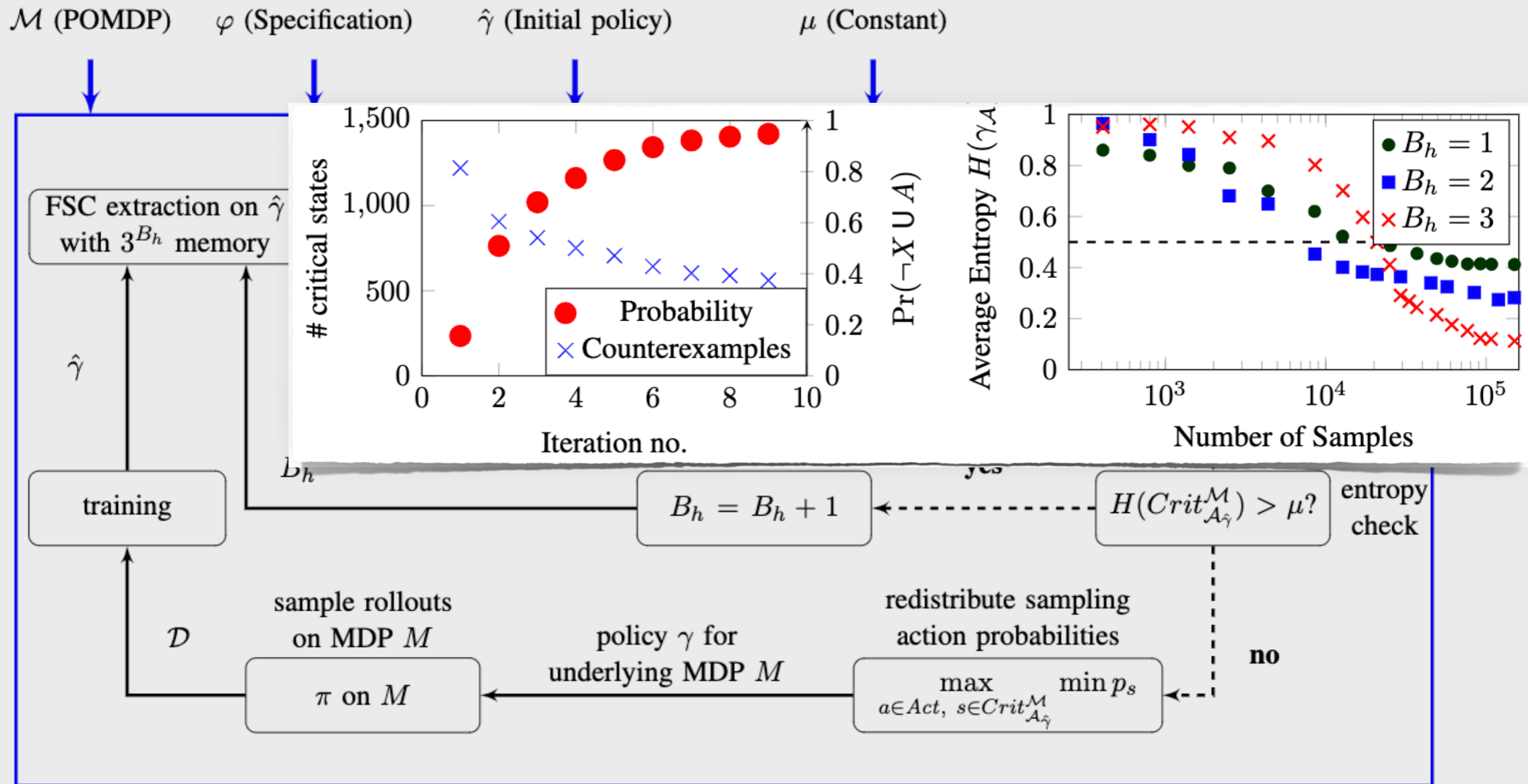




# Learn a guess and improve (with help from verification).



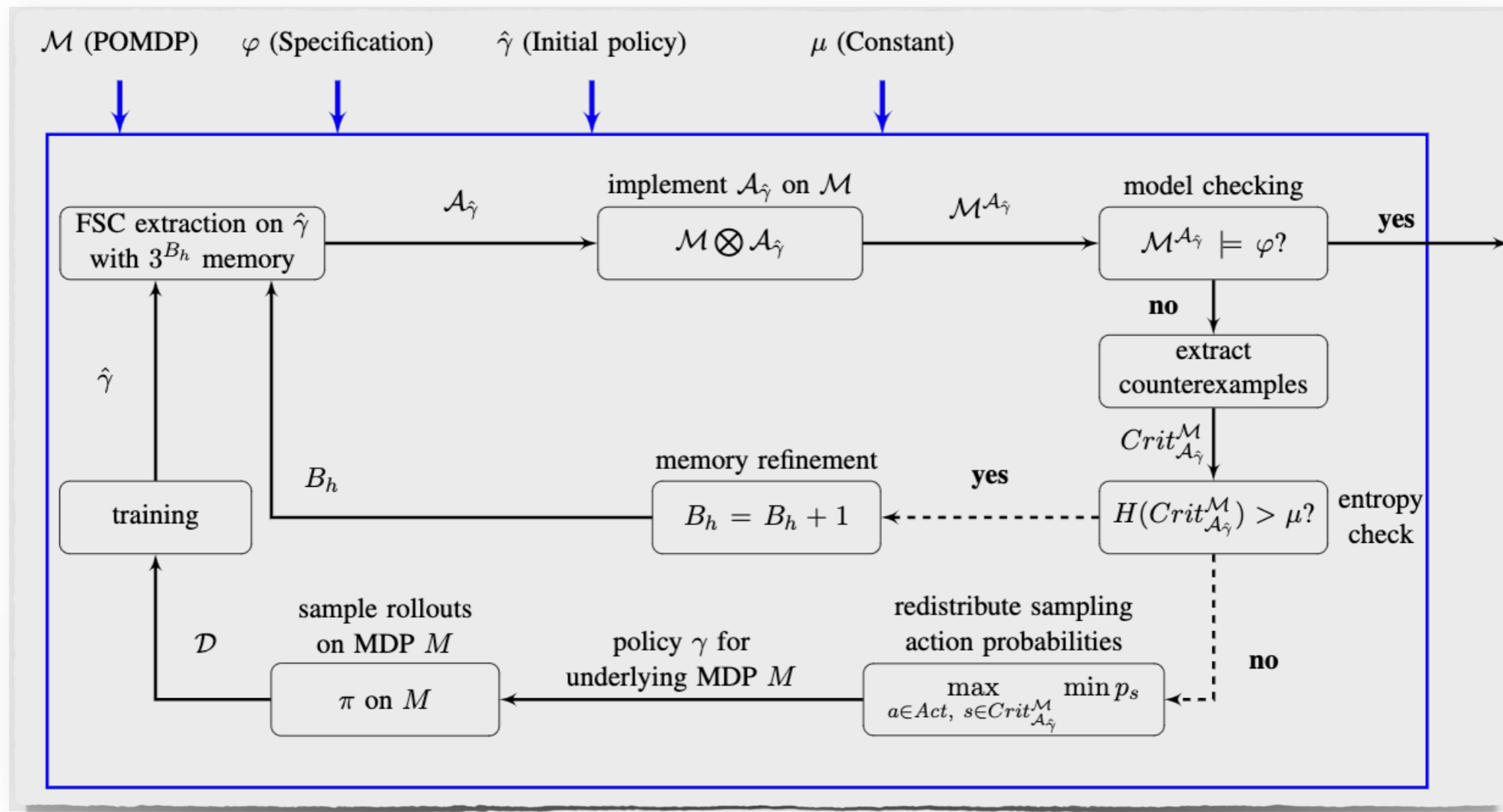
# Learn a guess and improve (with help from verification).



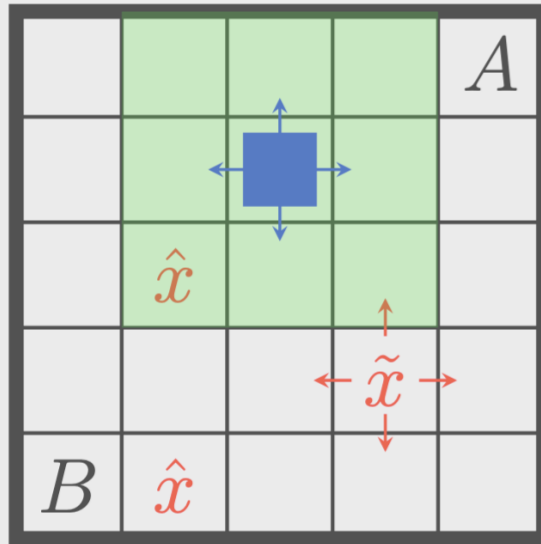
# What do we get at the end?

**Correct**, as each strategy prediction is evaluated using model checking.

**Not complete**, as we may never find a feasible strategy.  
Obviously, expected!



# Numerical examples with LTL constraints



Problem	$ S $	$ Act $	$ Z $
Navigation ( $c$ )	$c^4$	4	256
Delivery ( $c$ )	$c^2$	4	256
Slippery ( $c$ )	$c^2$	4	256
Maze( $c$ )	$3c + 8$	4	7
Grid( $c$ )	$c^2$	4	2
RockSample[4, 4]	257	9	2
RockSample[5, 5]	801	10	2
RockSample[7, 8]	12545	13	2

Problem	States	Type, $\varphi$	RNN-based Synthesis		PRISM-POMDP	
			Res.	Time (s)	Res.	Time (s)
Navigation (3)	333	$\mathbb{P}_{\max}^{\mathcal{M}}, \varphi_1$	0.74	<b>14.16</b>	<b>0.84</b>	73.88
Navigation (4)	1088	$\mathbb{P}_{\max}^{\mathcal{M}}, \varphi_1$	0.82	<b>22.67</b>	<b>0.93</b>	1034.64
Navigation (4) [2-FSC]	13373	$\mathbb{P}_{\max}^{\mathcal{M}}, \varphi_1$	0.91	47.26	–	–
Navigation (4) [4-FSC]	26741	$\mathbb{P}_{\max}^{\mathcal{M}}, \varphi_1$	0.92	59.42	–	–
Navigation (4) [8-FSC]	53477	$\mathbb{P}_{\max}^{\mathcal{M}}, \varphi_1$	<b>0.92</b>	85.26	–	–
Navigation (5)	2725	$\mathbb{P}_{\max}^{\mathcal{M}}, \varphi_1$	0.91	<b>34.34</b>	MO	MO
Navigation (5) [2-FSC]	33357	$\mathbb{P}_{\max}^{\mathcal{M}}, \varphi_1$	0.92	115.16	–	–
Navigation (5) [4-FSC]	66709	$\mathbb{P}_{\max}^{\mathcal{M}}, \varphi_1$	0.92	159.61	–	–
Navigation (5) [8-FSC]	133413	$\mathbb{P}_{\max}^{\mathcal{M}}, \varphi_1$	<b>0.92</b>	250.91	–	–
Navigation (10)	49060	$\mathbb{P}_{\max}^{\mathcal{M}}, \varphi_1$	0.79	<b>822.87</b>	MO	MO
Navigation (10) [2-FSC]	475053	$\mathbb{P}_{\max}^{\mathcal{M}}, \varphi_1$	0.83	1185.41	–	–
Navigation (10) [4-FSC]	950101	$\mathbb{P}_{\max}^{\mathcal{M}}, \varphi_1$	<b>0.85</b>	1488.77	–	–
Navigation (10) [8-FSC]	1900197	$\mathbb{P}_{\max}^{\mathcal{M}}, \varphi_1$	0.81	1805.22	–	–
Navigation (15)	251965	$\mathbb{P}_{\max}^{\mathcal{M}}, \varphi_1$	<b>0.91</b>	<b>1271.80*</b>	MO	MO
Navigation (20)	798040	$\mathbb{P}_{\max}^{\mathcal{M}}, \varphi_1$	<b>0.96</b>	<b>4712.25*</b>	MO	MO
Navigation (30)	4045840	$\mathbb{P}_{\max}^{\mathcal{M}}, \varphi_1$	<b>0.95</b>	<b>25191.05*</b>	MO	MO
Navigation (40)	–	$\mathbb{P}_{\max}^{\mathcal{M}}, \varphi_1$	TO	TO	MO	MO
Delivery (4) [2-FSC]	80	$\mathbb{E}_{\min}^{\mathcal{M}}, \varphi_2$	6.02	35.35	<b>6.0</b>	<b>28.53</b>
Delivery (5) [2-FSC]	125	$\mathbb{E}_{\min}^{\mathcal{M}}, \varphi_2$	8.11	<b>78.32</b>	<b>8.0</b>	102.41
Delivery (10) [2-FSC]	500	$\mathbb{E}_{\min}^{\mathcal{M}}, \varphi_2$	<b>18.13</b>	<b>120.34</b>	MO	MO
Slippery (4) [2-FSC]	460	$\mathbb{P}_{\max}^{\mathcal{M}}, \varphi_3$	0.78	67.51	<b>0.90</b>	<b>5.10</b>
Slippery (5) [2-FSC]	730	$\mathbb{P}_{\max}^{\mathcal{M}}, \varphi_3$	0.89	84.32	<b>0.93</b>	<b>83.24</b>
Slippery (10) [2-FSC]	2980	$\mathbb{P}_{\max}^{\mathcal{M}}, \varphi_3$	<b>0.98</b>	<b>119.14</b>	MO	MO
Slippery (20) [2-FSC]	11980	$\mathbb{P}_{\max}^{\mathcal{M}}, \varphi_3$	<b>0.99</b>	<b>1580.42</b>	MO	MO

# Numerical examples on standard POMDP benchmarks

Problem	Type	RNN-based Synthesis			PRISM-POMDP		pomdpSolve	
		States	Res	Time (s)	Res	Time (s)	Res	Time (s)
Maze (1)	$E_{\min}^{\mathcal{M}}$	68	4.31	31.70	<b>4.30</b>	<b>0.09</b>	4.30	0.30
Maze (2)	$E_{\min}^{\mathcal{M}}$	83	5.31	46.65	5.23	2.176	<b>5.23</b>	<b>0.67</b>
Maze (3)	$E_{\min}^{\mathcal{M}}$	98	8.10	58.75	7.13	38.82	<b>7.13</b>	<b>2.39</b>
Maze (4)	$E_{\min}^{\mathcal{M}}$	113	11.53	58.09	8.58	543.06	<b>8.58</b>	<b>7.15</b>
Maze (5)	$E_{\min}^{\mathcal{M}}$	128	14.40	<b>68.09</b>	13.00	4110.50	<b>12.04</b>	132.12
Maze (6)	$E_{\min}^{\mathcal{M}}$	143	22.34	<b>71.89</b>	MO	MO	<b>18.52</b>	1546.02
Maze (10)	$E_{\min}^{\mathcal{M}}$	203	100.21	<b>158.33</b>	MO	MO	MO	MO
Grid (3)	$E_{\min}^{\mathcal{M}}$	165	2.90	38.94	2.88	2.332	<b>2.88</b>	<b>0.07</b>
Grid (4)	$E_{\min}^{\mathcal{M}}$	381	4.32	79.99	4.13	1032.53	<b>4.13</b>	<b>0.77</b>
Grid (5)	$E_{\min}^{\mathcal{M}}$	727	6.623	91.42	MO	MO	<b>5.42</b>	<b>1.94</b>
Grid (10)	$E_{\min}^{\mathcal{M}}$	5457	<b>13.630</b>	<b>268.40</b>	MO	MO	MO	MO
RockSample[4, 4]	$E_{\max}^{\mathcal{M}}$	2432	17.71	35.35	N/A	N/A	<b>18.04</b>	<b>0.43</b>
RockSample[5, 5]	$E_{\max}^{\mathcal{M}}$	8320	18.40	<b>43.74</b>	N/A	N/A	<b>19.23</b>	621.28
RockSample[7, 8]	$E_{\max}^{\mathcal{M}}$	166656	20.32	<b>860.53</b>	N/A	N/A	<b>21.64</b>	20458.41



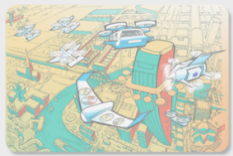
# Assured Autonomy: Path Toward Living With Autonomous Systems We Can Trust



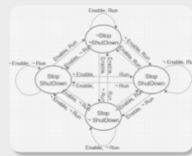
## Some common comments:

- No assurance = no useful autonomy
- Diverse set of vulnerabilities
- Open world
- Interdisciplinary approaches needed, **not** as an afterthought

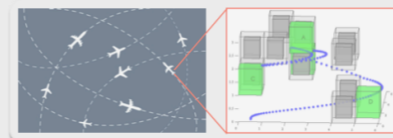
## What makes autonomy hard?



dynamic environment



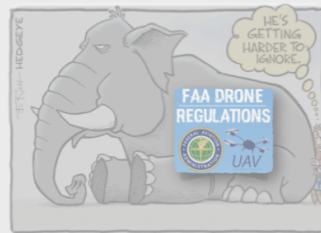
complex missions



heterogenous decisions



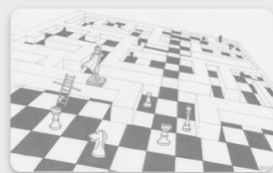
run-time faults



verifiability



unknown environments



integrity of critical information



imperfect perception



variations in user characteristics

1. Safety and verification
2. Security
3. Certification and regulation
4. Human-system integration and trust
5. Privacy
6. Ethics
7. Societal impacts
8. Governance and policy
9. Education