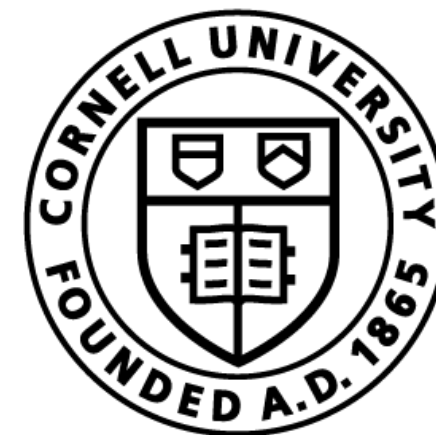


Dealing with Uncertainty: Part 2

Sanjiban Choudhury

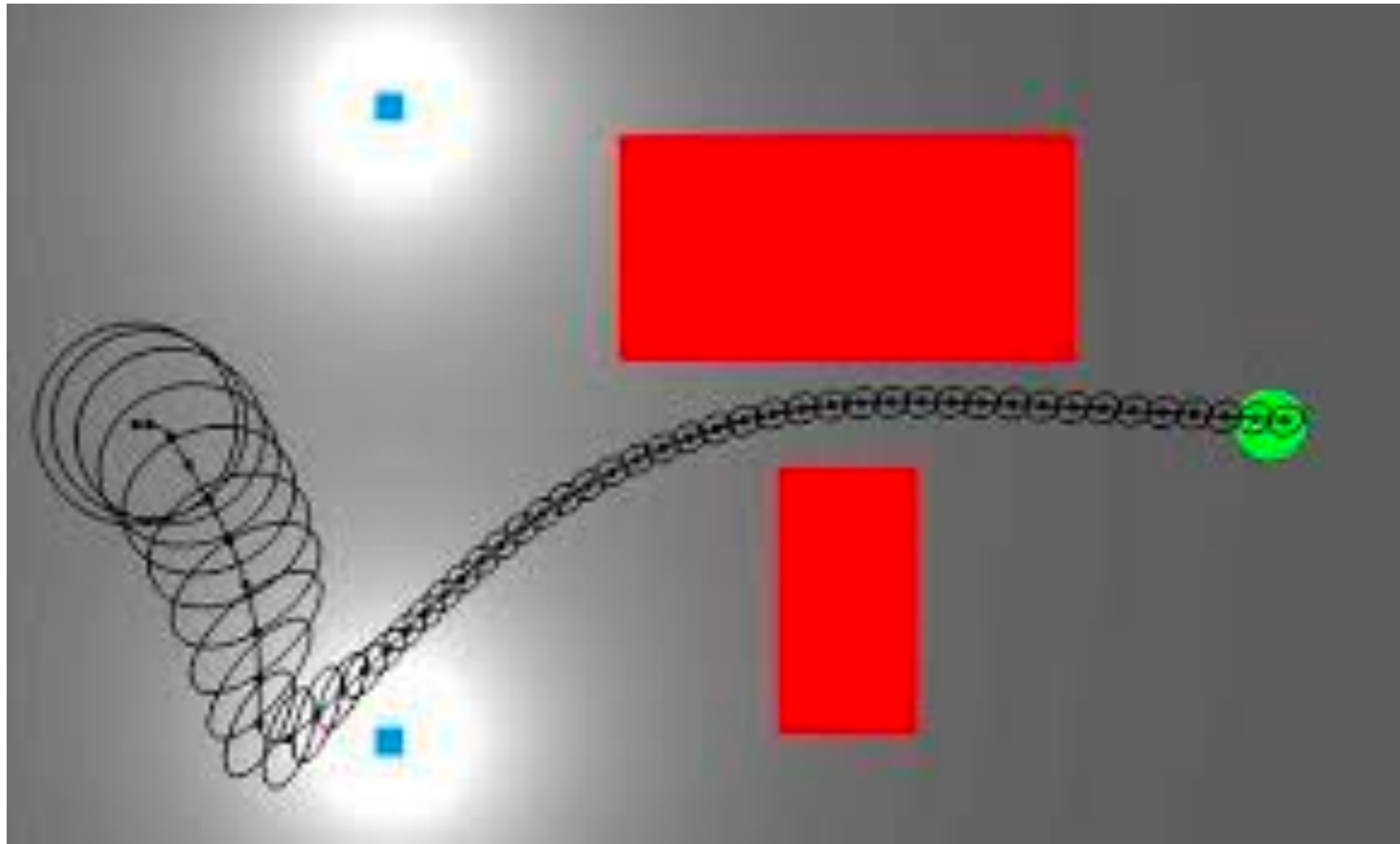


Cornell Bowers CIS
Computer Science

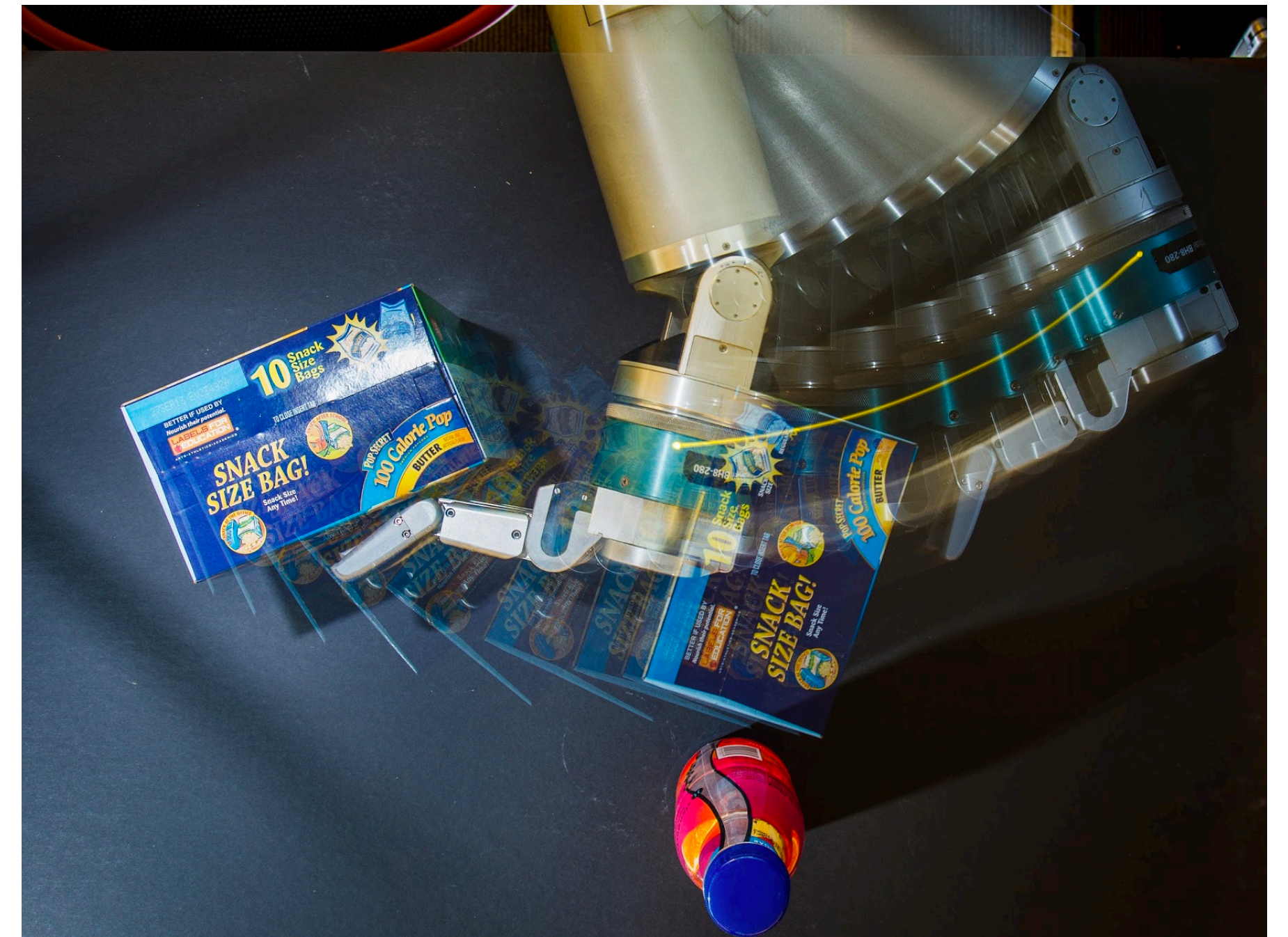
A grayscale, misty landscape of a forested valley. The scene is dominated by dense evergreen trees, with a thick layer of fog or mist filling the valley and obscuring the details of the forest in the distance. The lighting is soft and diffused, creating a sense of depth and atmosphere. The overall tone is somber and mysterious.

Uncertainty

Epistemic Uncertainty

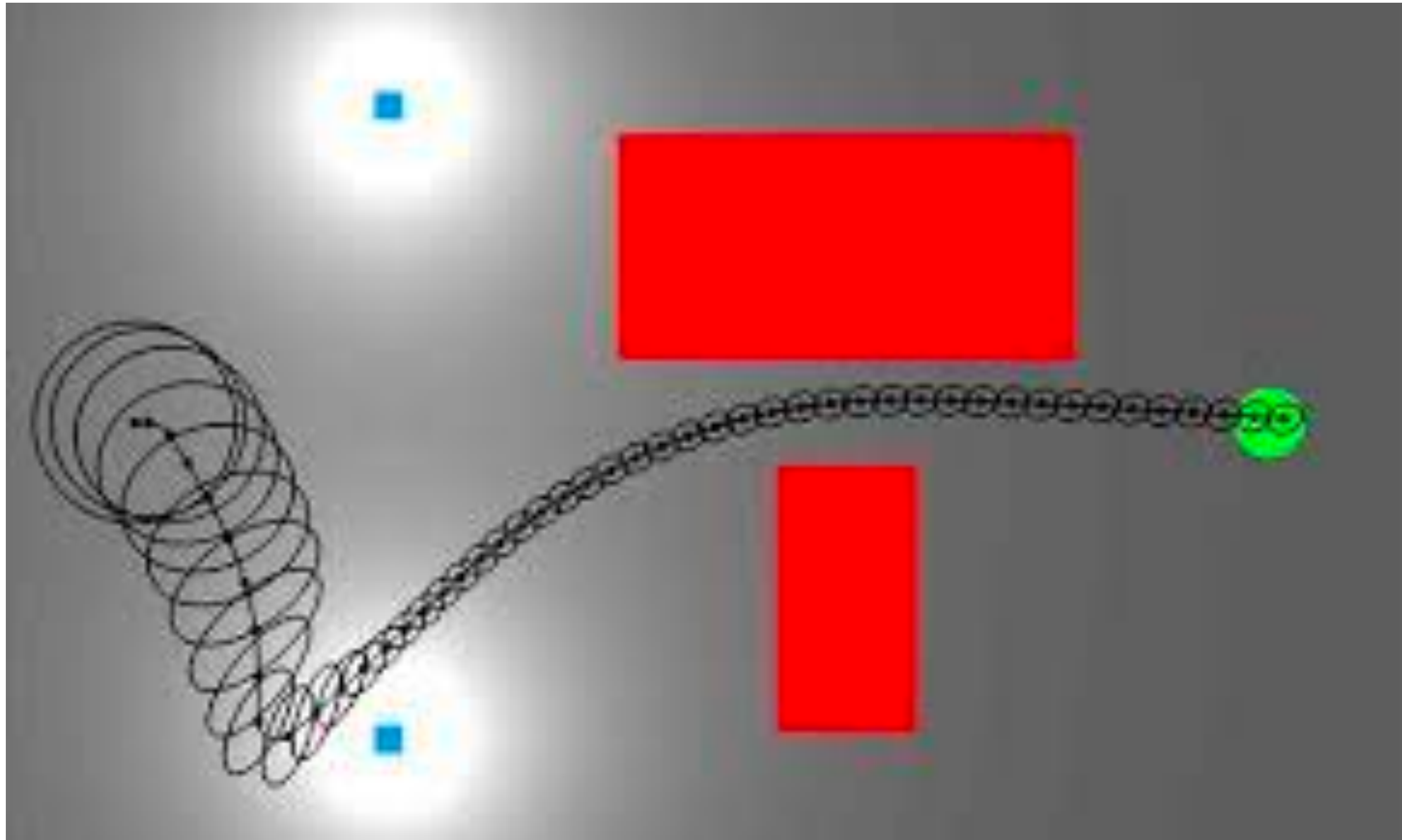


Uncertain about state

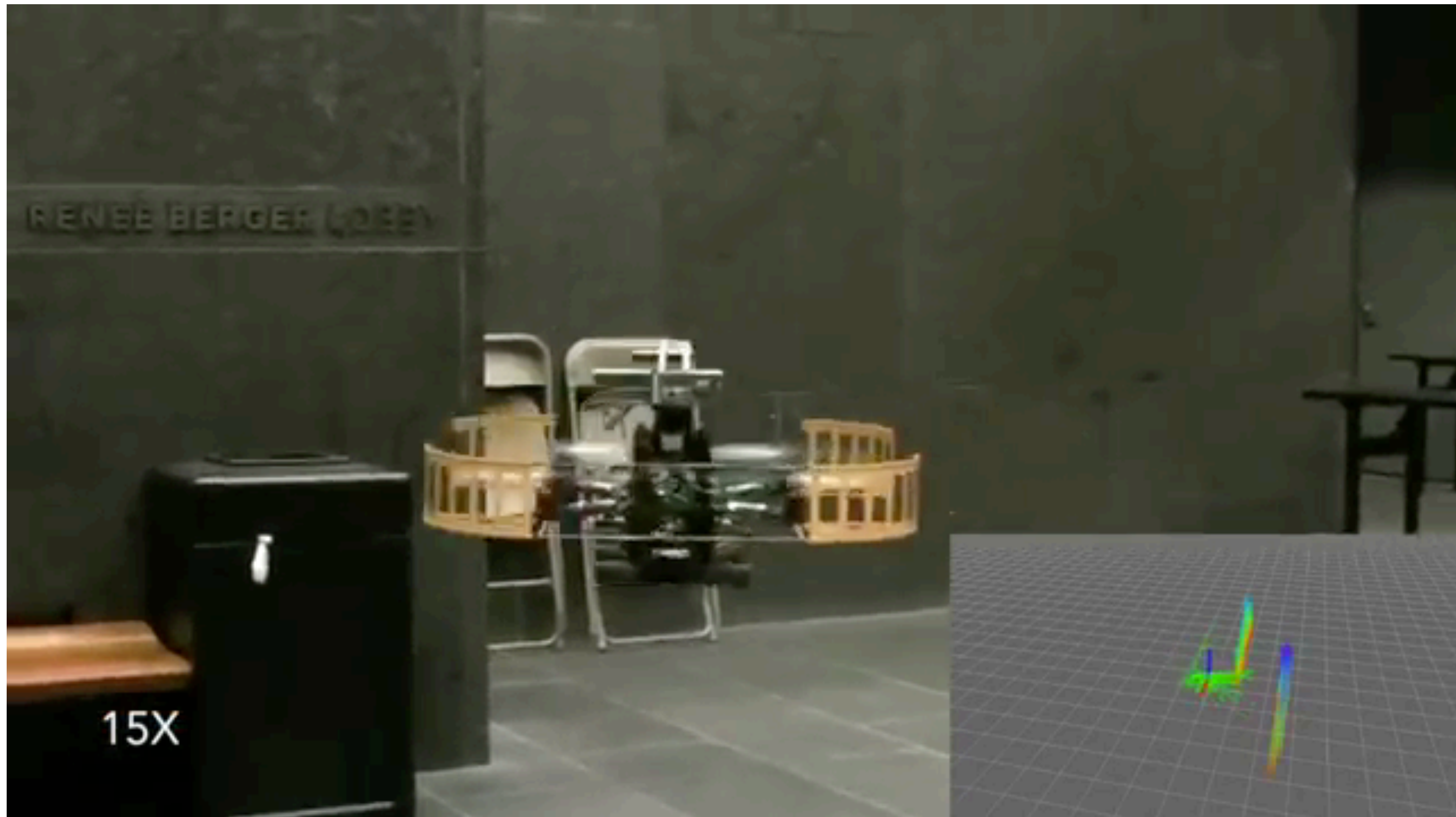


Uncertain about transitions

Uncertain about the robot pose

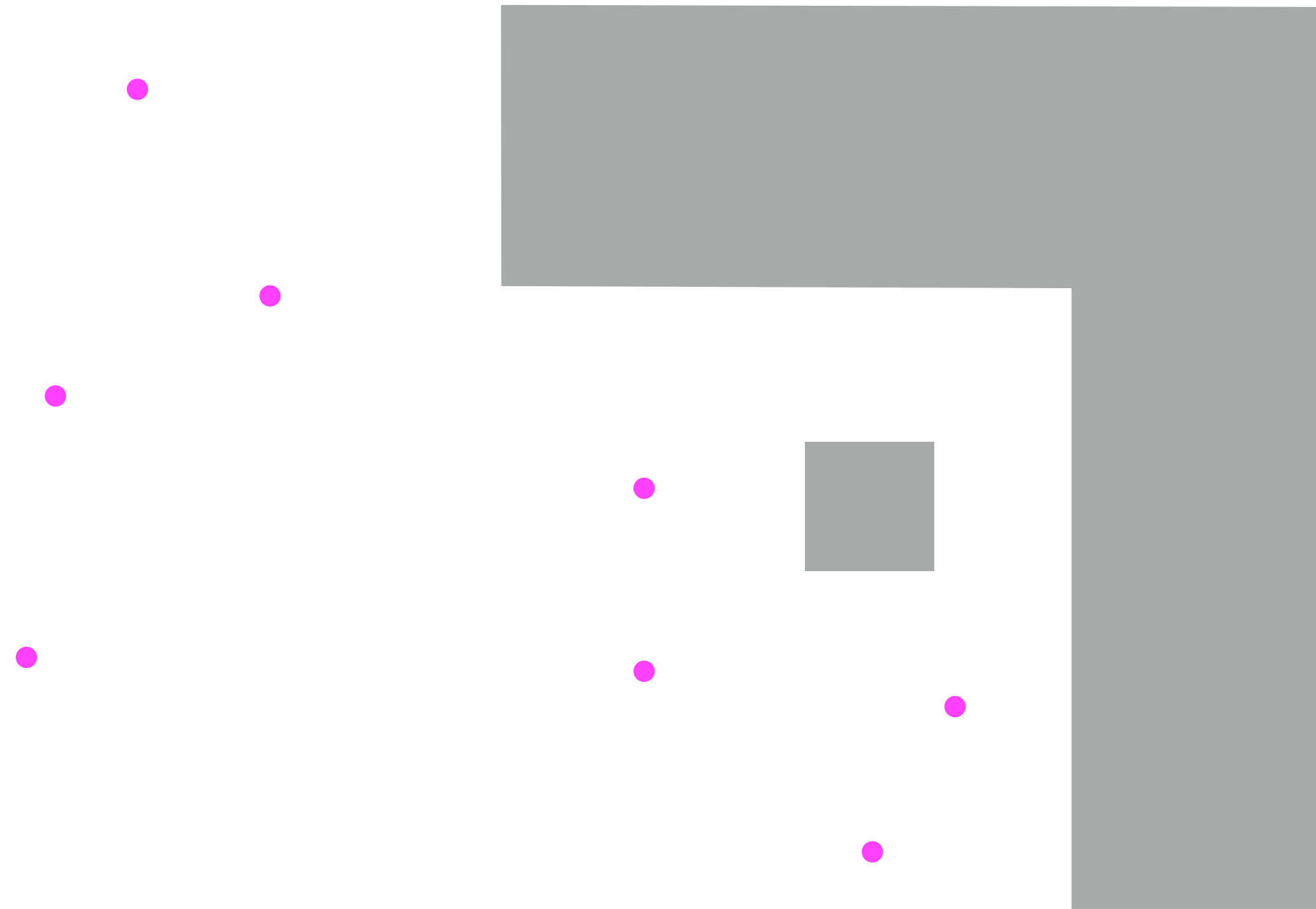


Uncertain about the world



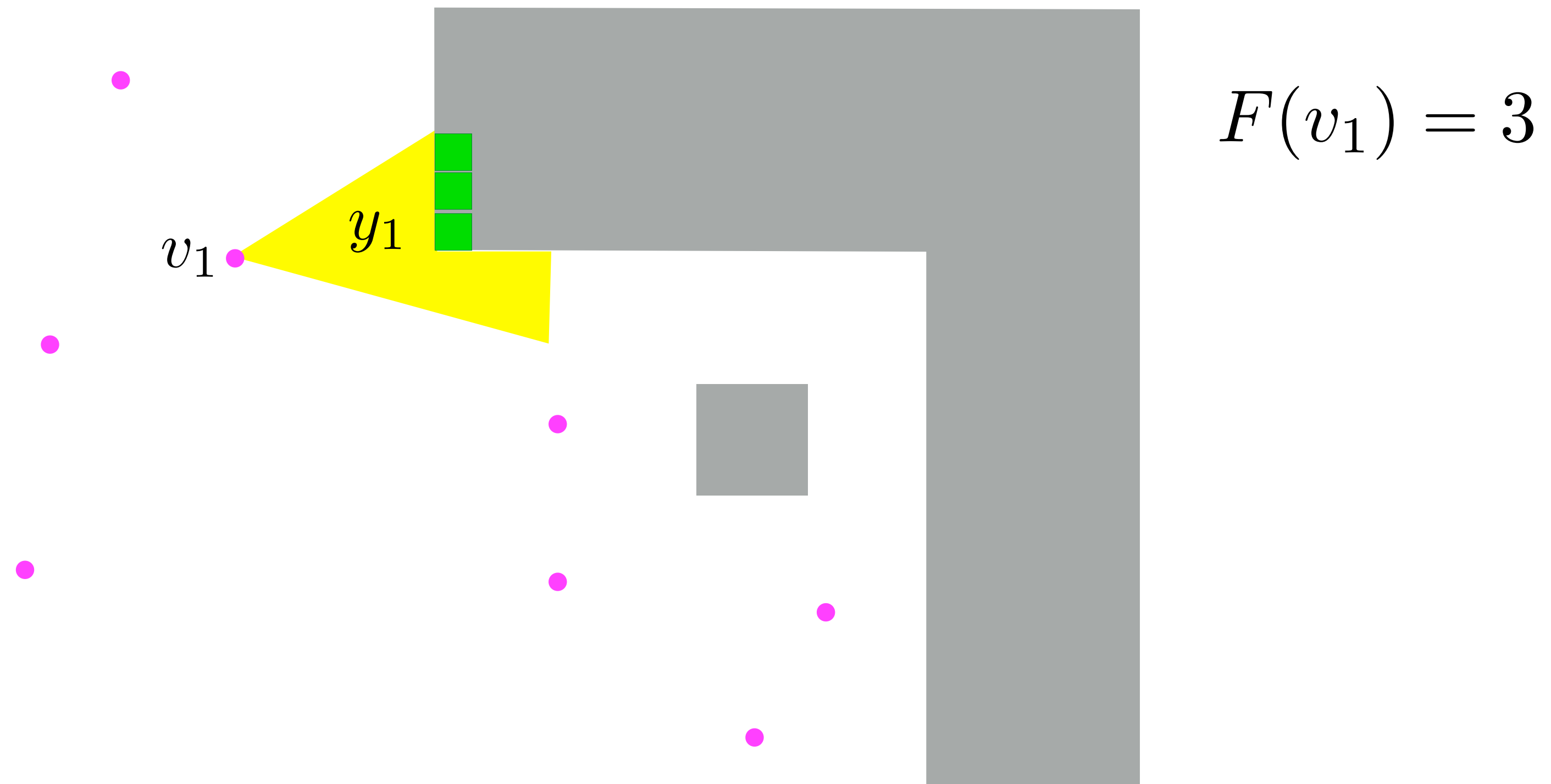
The coverage problem

$$\text{maximize}_{\{v_1, \dots, v_n\}} F(v_1, \dots, v_n)$$



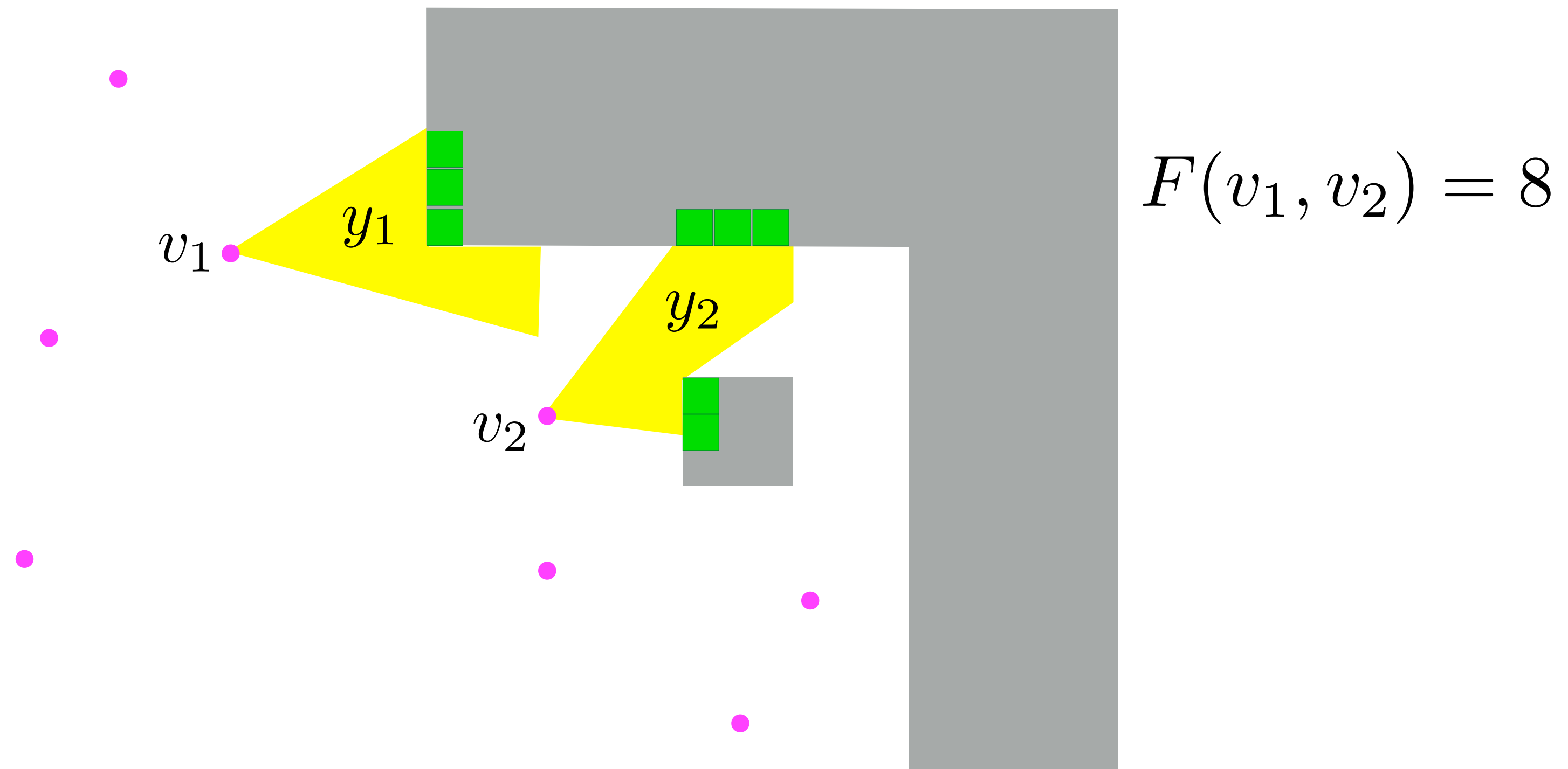
The coverage problem

$$\text{maximize } F(v_1, \dots, v_n) \\ \{v_1, \dots, v_n\}$$



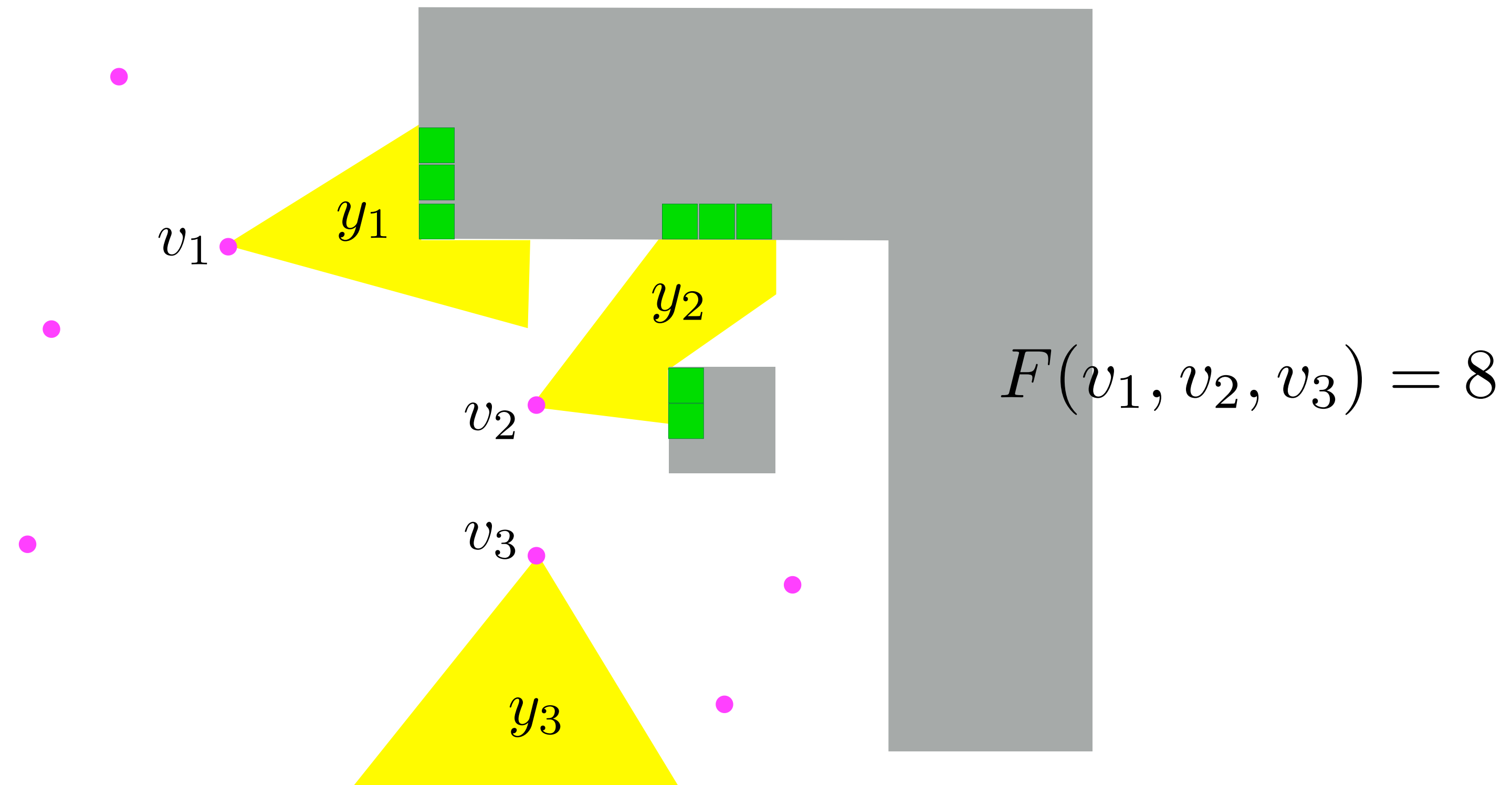
The coverage problem

$$\text{maximize } F(v_1, \dots, v_n) \\ \{v_1, \dots, v_n\}$$

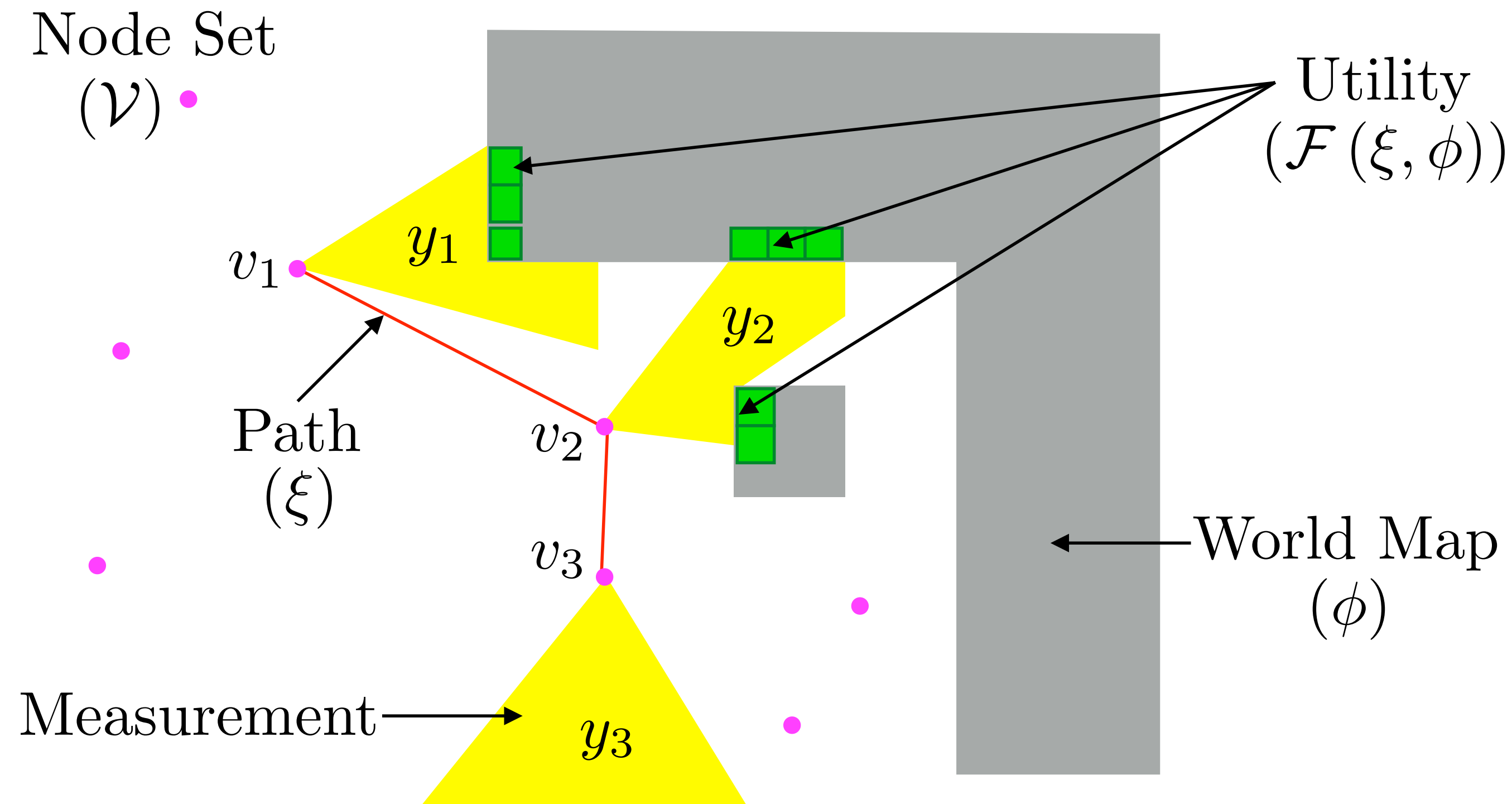


The coverage problem

$$\text{maximize } F(v_1, \dots, v_n) \\ \{v_1, \dots, v_n\}$$



The *budgeted* coverage problem



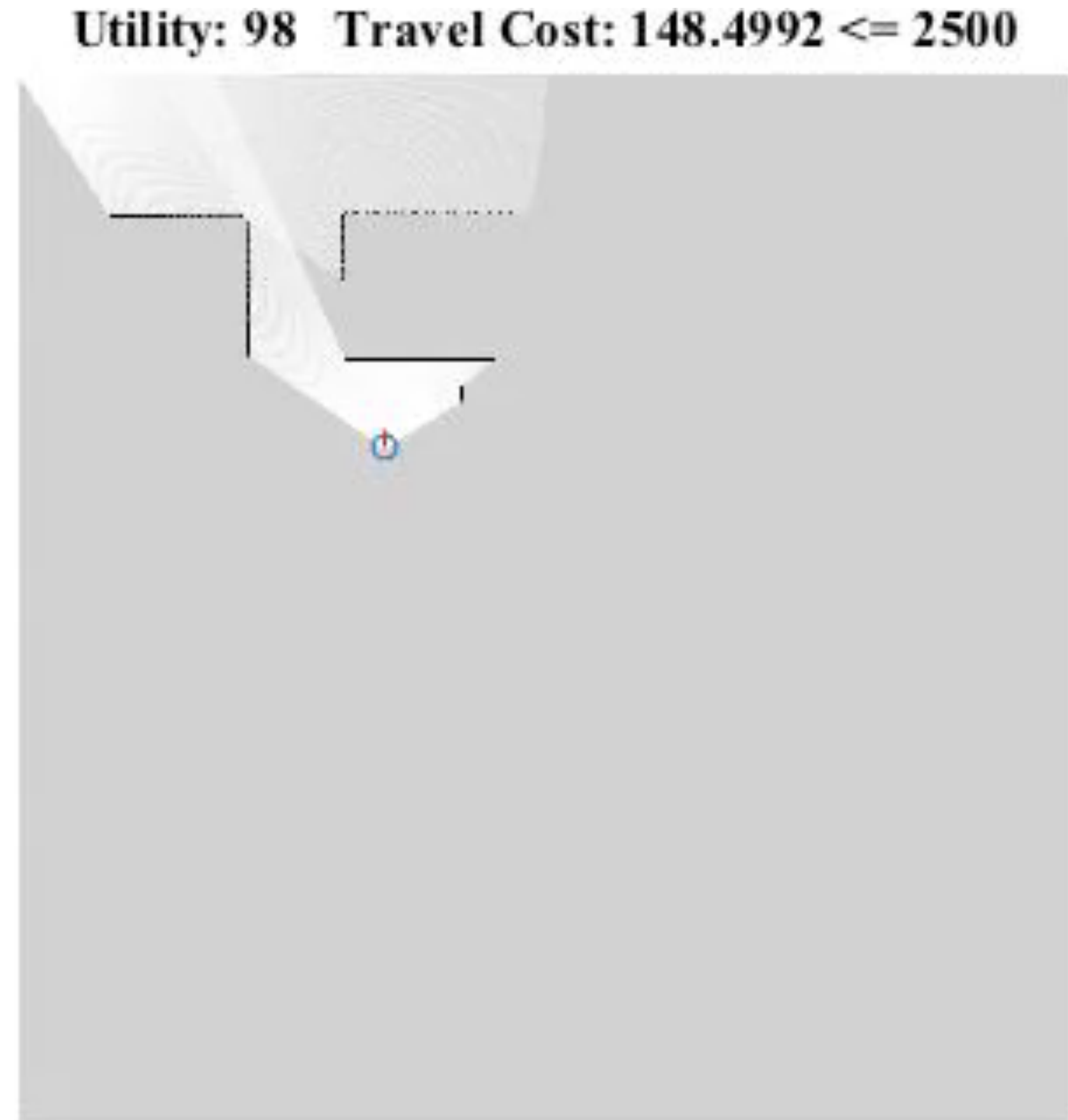
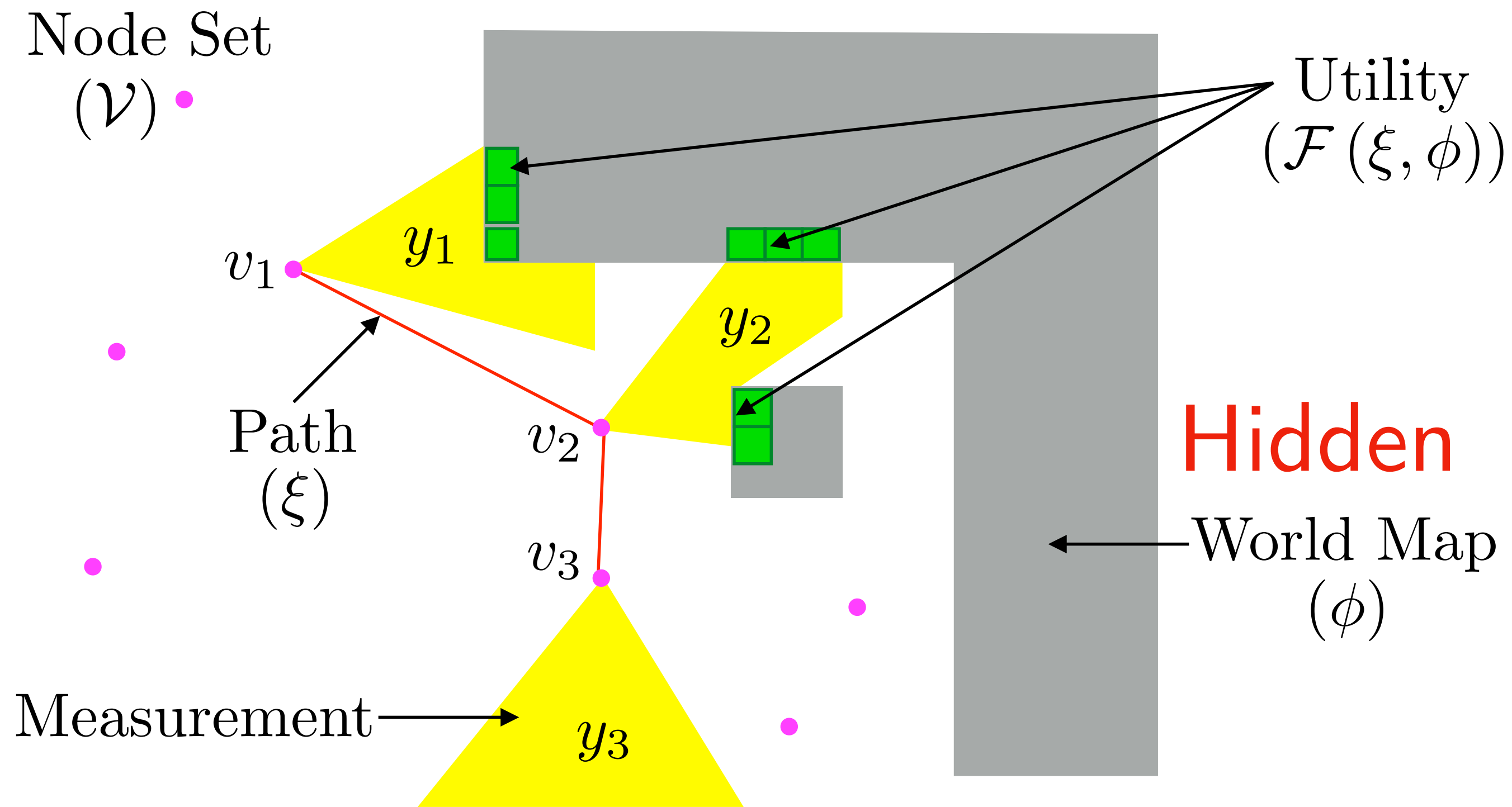
$$\arg \max_{\xi \in \Xi} \mathcal{F}(\xi, \phi)$$

$$s.t. \quad \mathcal{T}(\xi, \phi) \leq B$$

Cover as many cells

Subject to travel cost!

The *budgeted coverage* info-gathering problem



$$\arg \max_{\xi \in \Xi} \mathcal{F}(\xi, \phi)$$

$$s.t. \quad \mathcal{T}(\xi, \phi) \leq B$$

Activity!



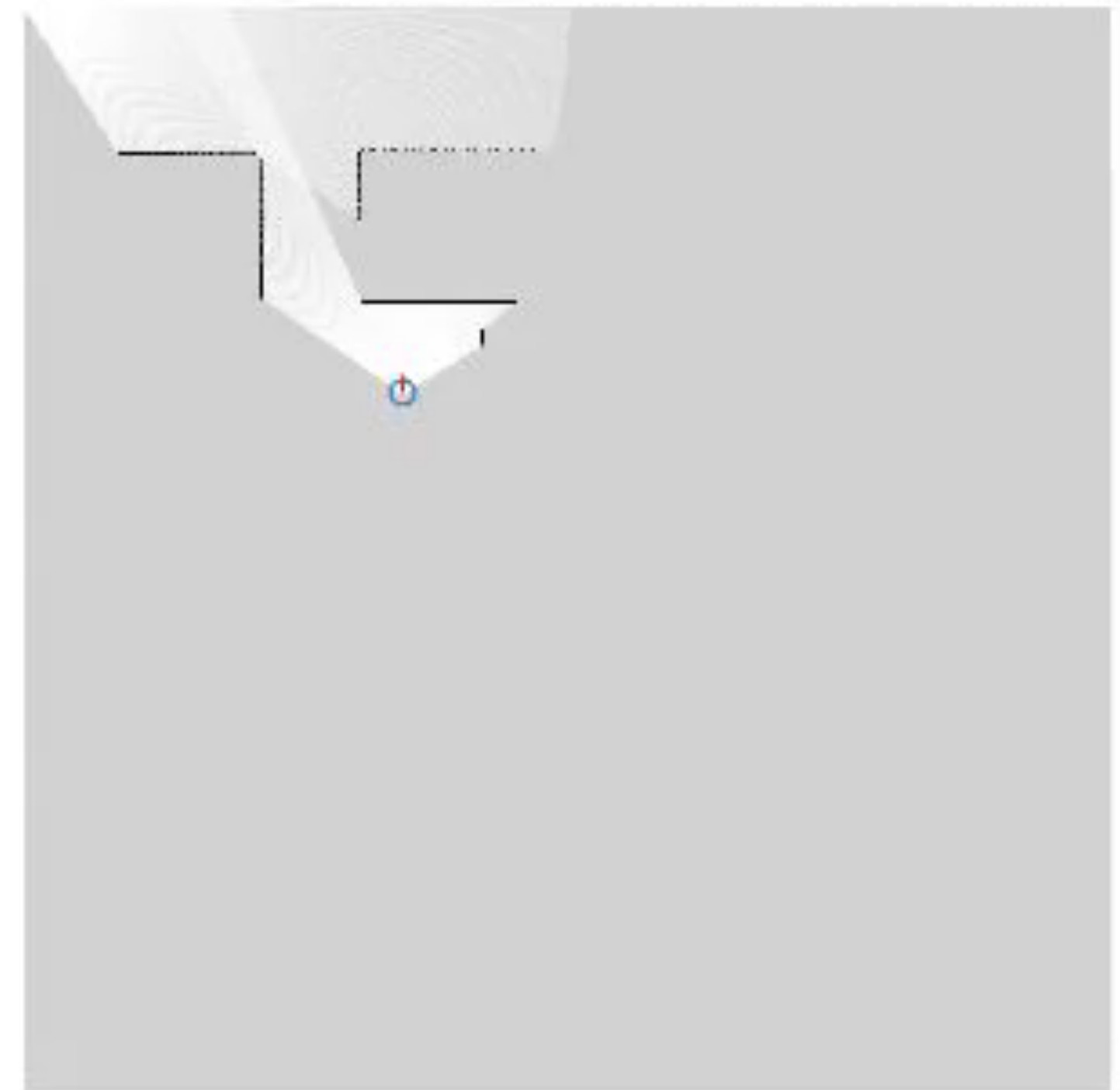
Think-Pair-Share

Think (30 sec): Can you think of some heuristics for solving the budgeted information gathering problem?

Pair: Find a partner

Share (45 sec): Partners exchange ideas

Utility: 98 Travel Cost: 148.4992 \leq 2500



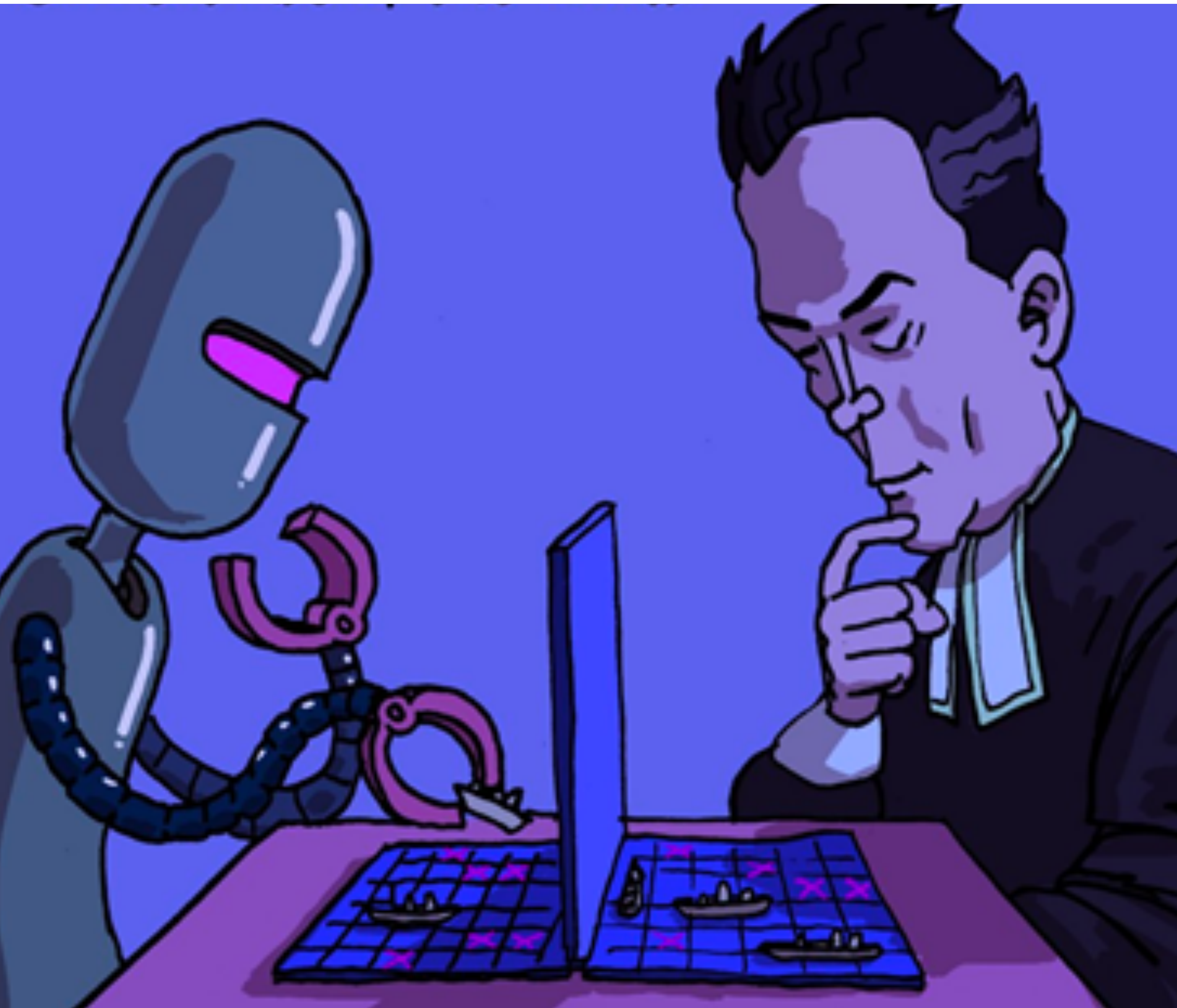


Belief Space Planning is NP-Hard
at best, undecidable at worst

Need to relax our problem!

What if we wanted to explore as optimally as possible using prior information?





Information Gain

20 Questions

Let's say you have a set of hypotheses

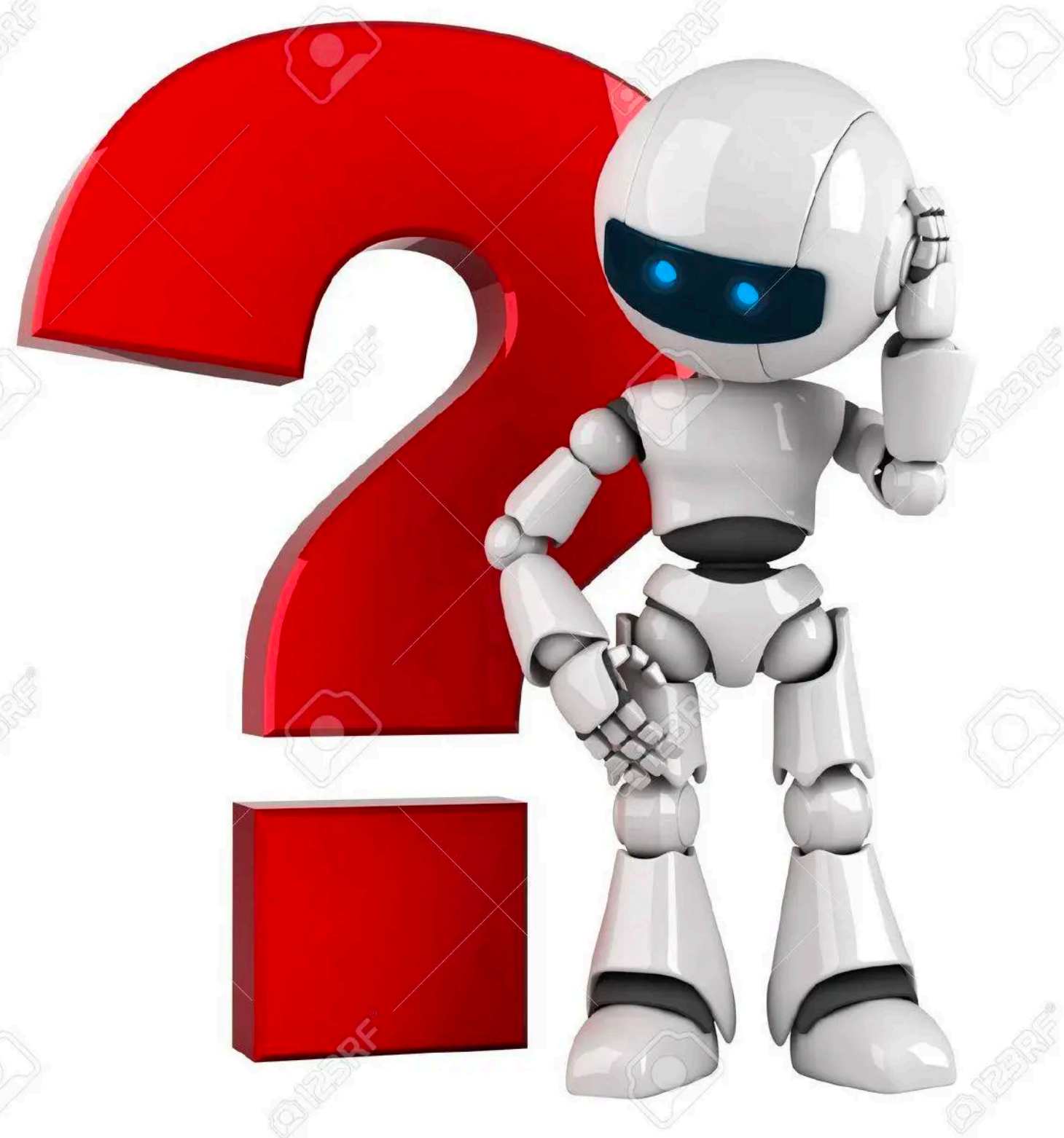
$$\{\theta_1, \theta_2, \dots, \theta_n\}$$

and a set of tests

$$\{t_1, t_2, \dots, t_n\}$$

Given a prior over hypotheses $P(\theta)$

Find the minimal number of tests to identify hypothesis



20 Questions

Let's say you have a set of hypotheses

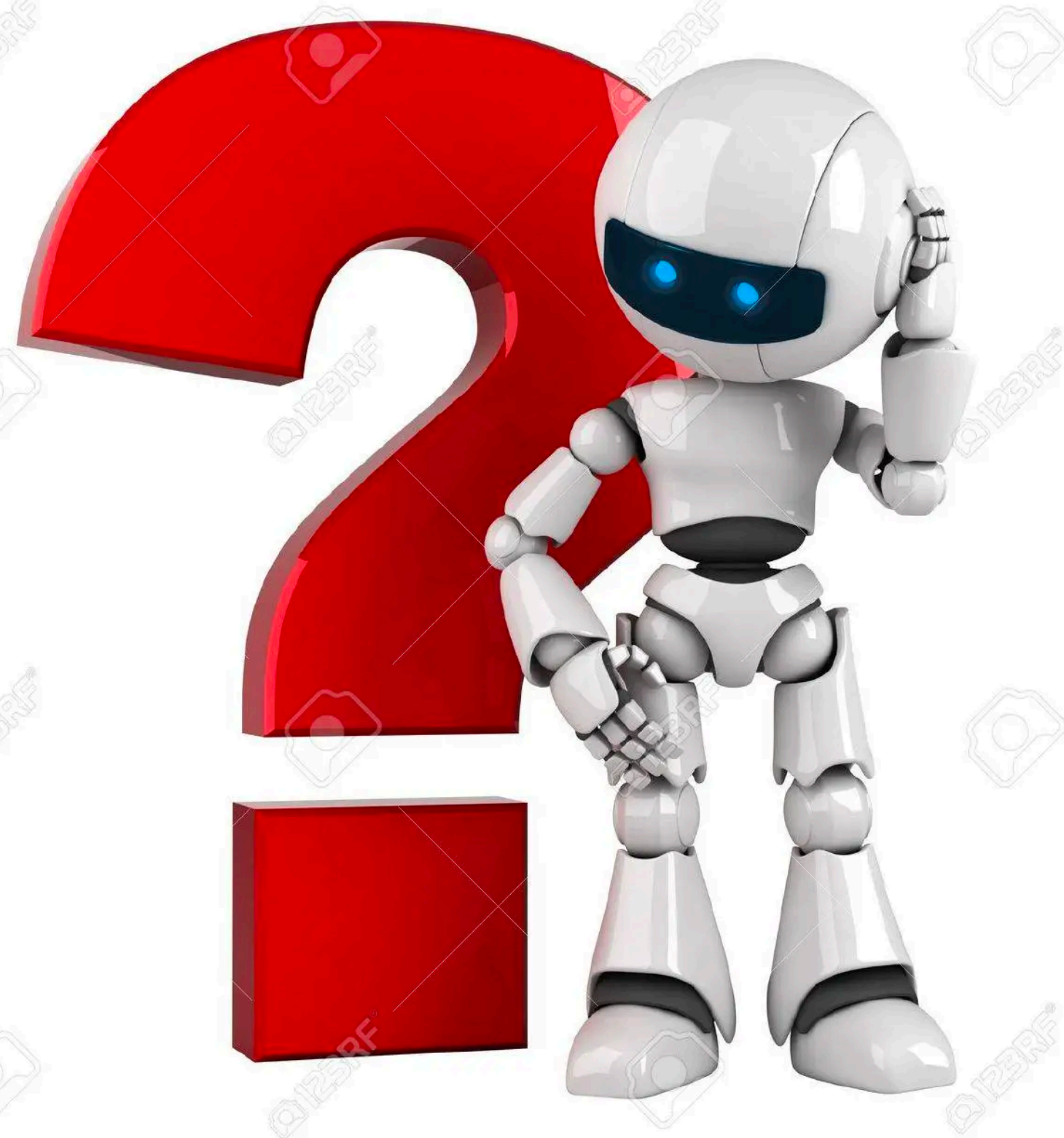
$$\{\theta_1, \theta_2, \dots, \theta_n\}$$

and a set of tests

$$\{t_1, t_2, \dots, t_n\}$$

Given a prior over hypotheses $P(\theta)$

Find the minimal number of tests to identify hypothesis



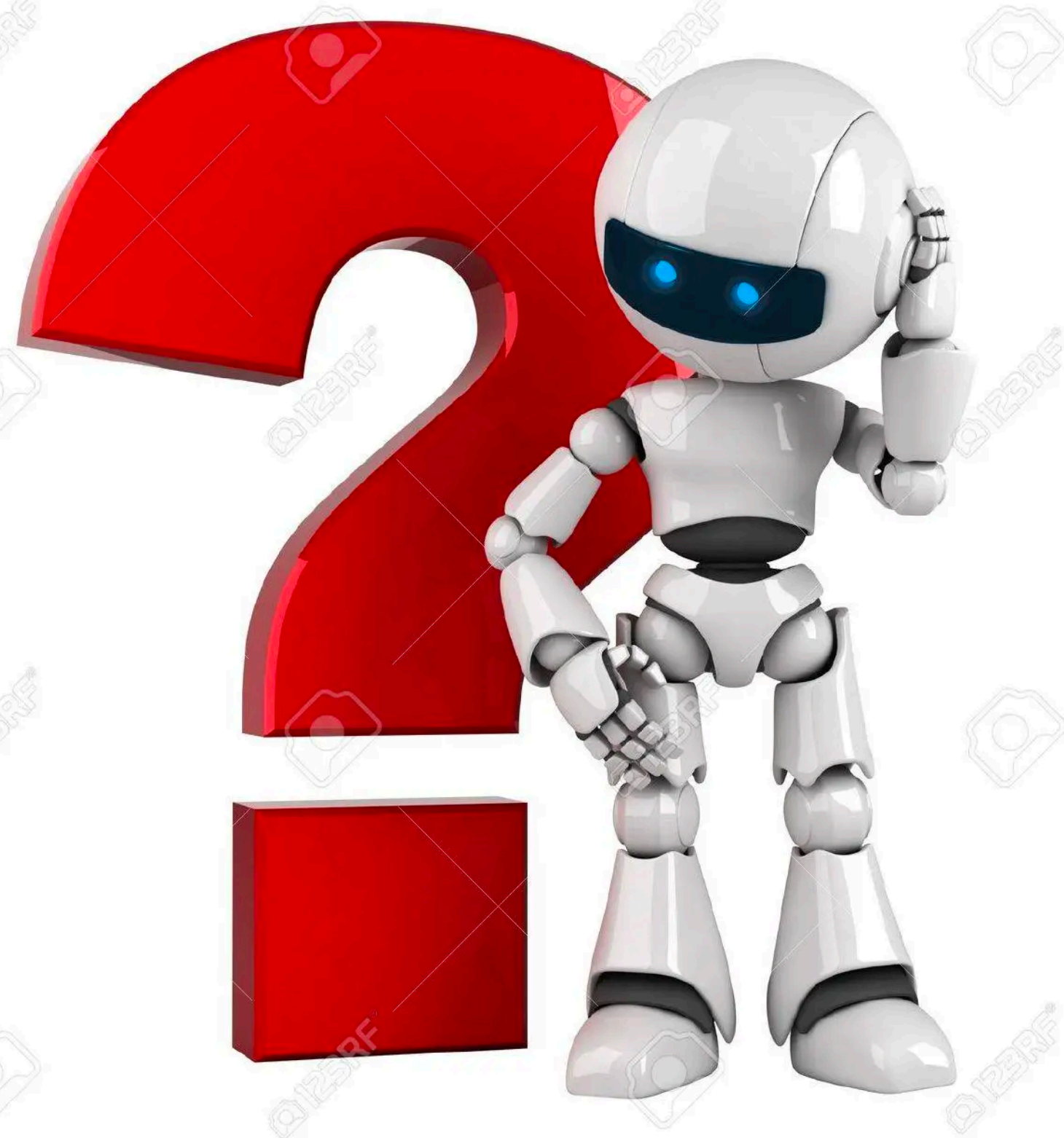
NP-HARD

A simple algorithm

Greedy pick the test that maximizes information gain

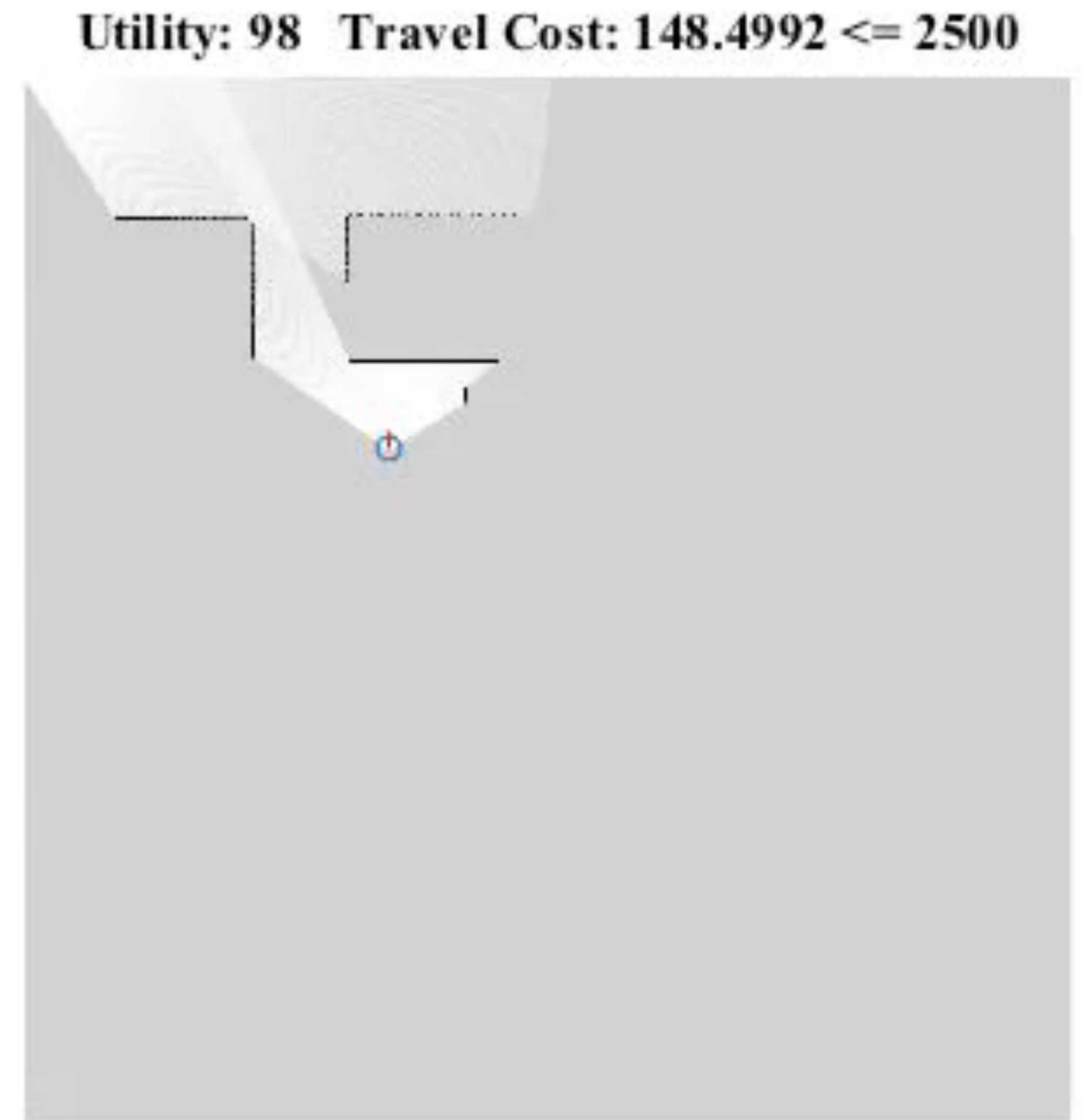
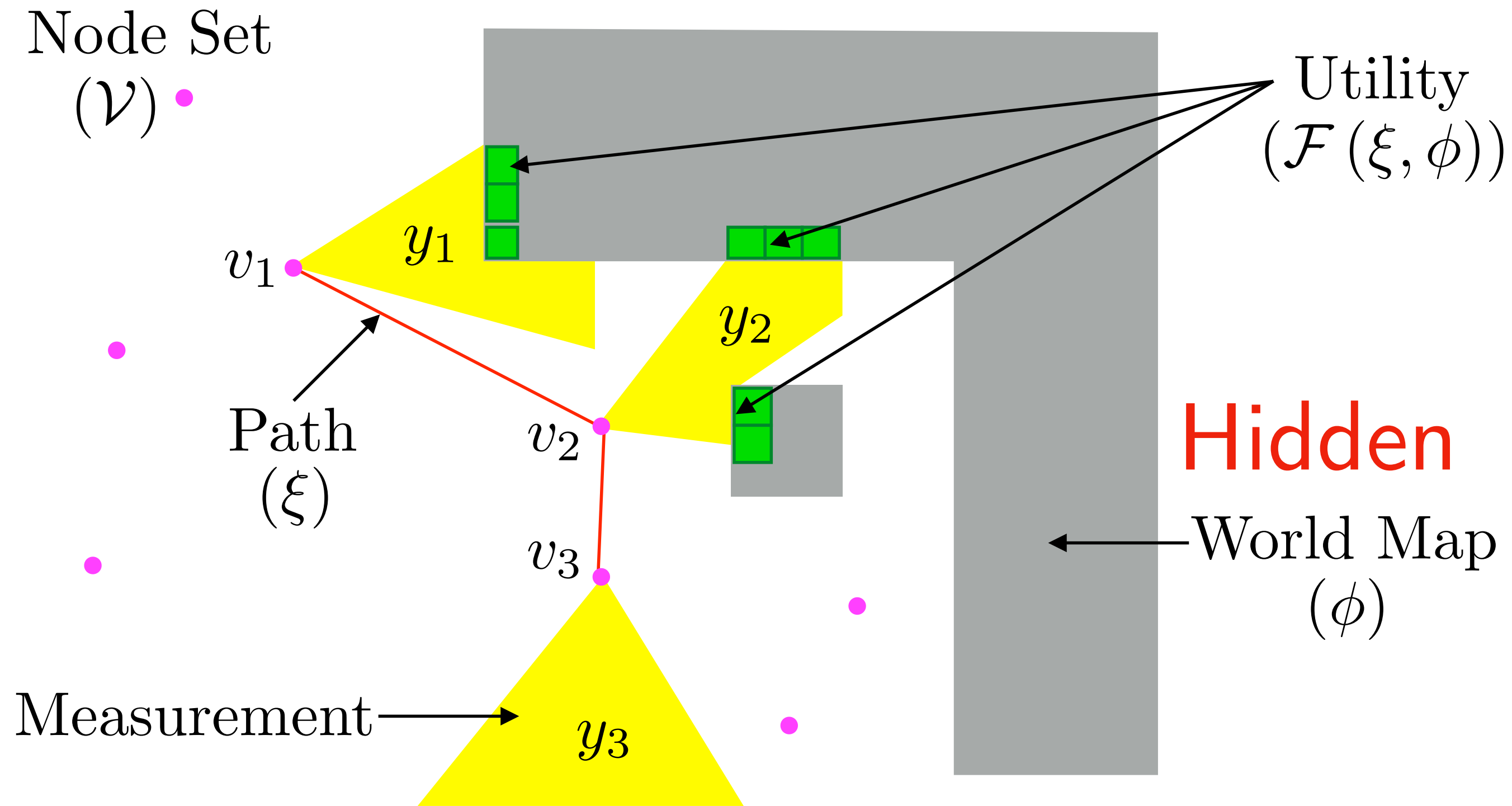
$$\max_t H(\theta) - \mathbb{E}_o H(\theta | t, o)$$

Entropy Posterior entropy



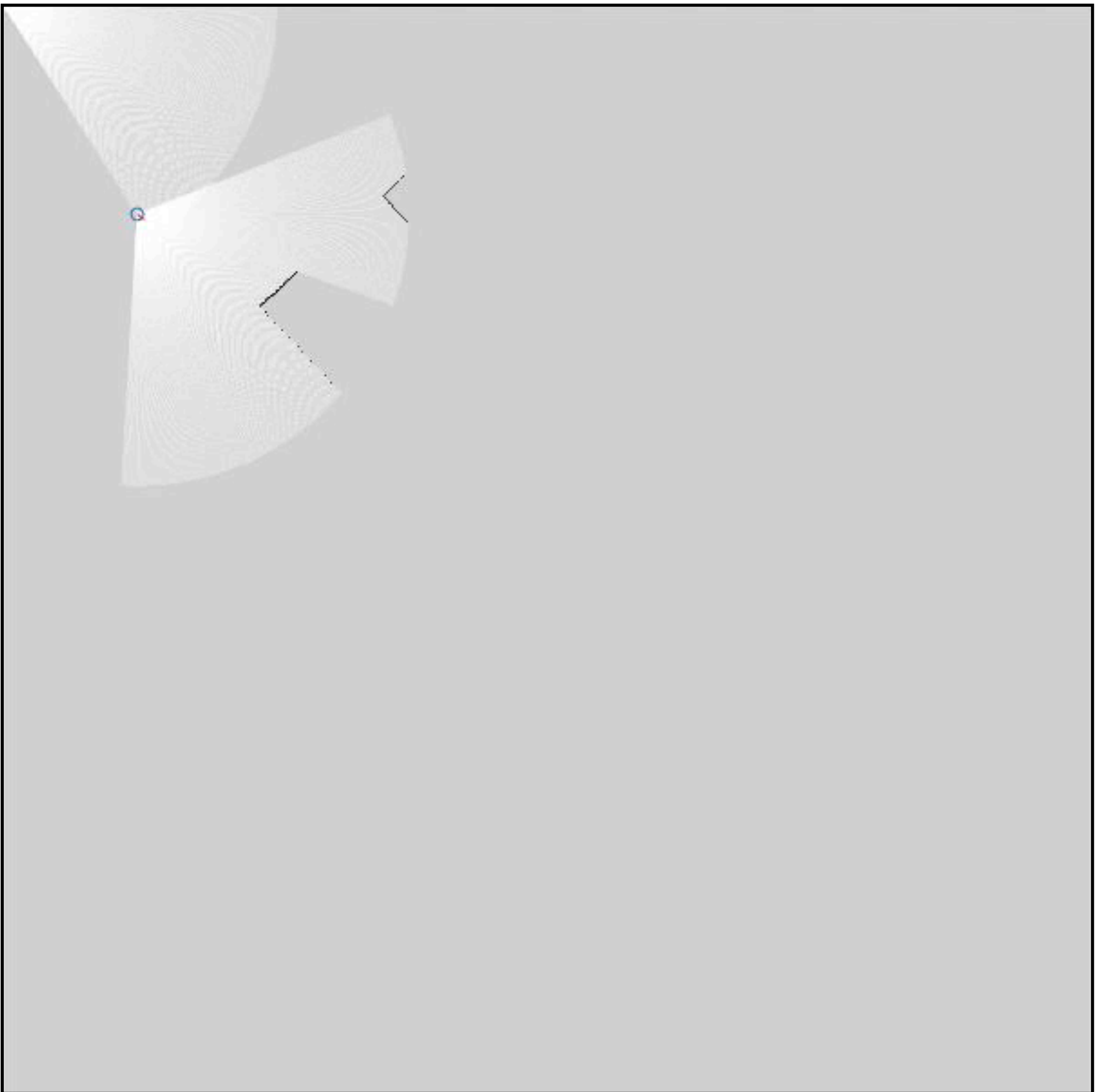
Entropy is adaptive sub modular \Rightarrow Greedy is near-optimal

So does information gain work for this problem?



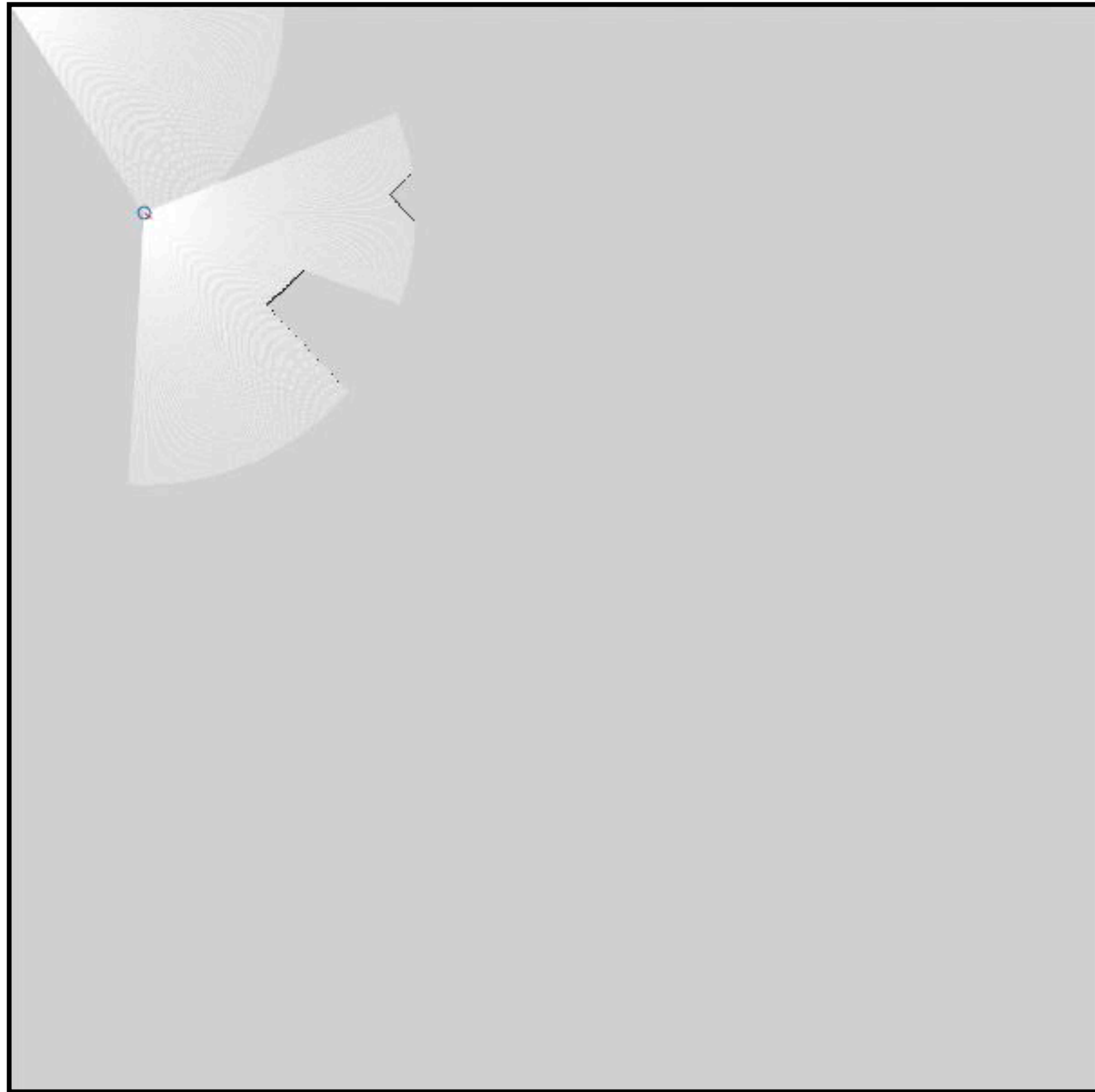
$$\arg \max_{\xi \in \Xi} \mathcal{F}(\xi, \phi)$$

$$s.t. \quad \mathcal{T}(\xi, \phi) \leq B$$



Information Gain
overexplores!

Why does this happen?

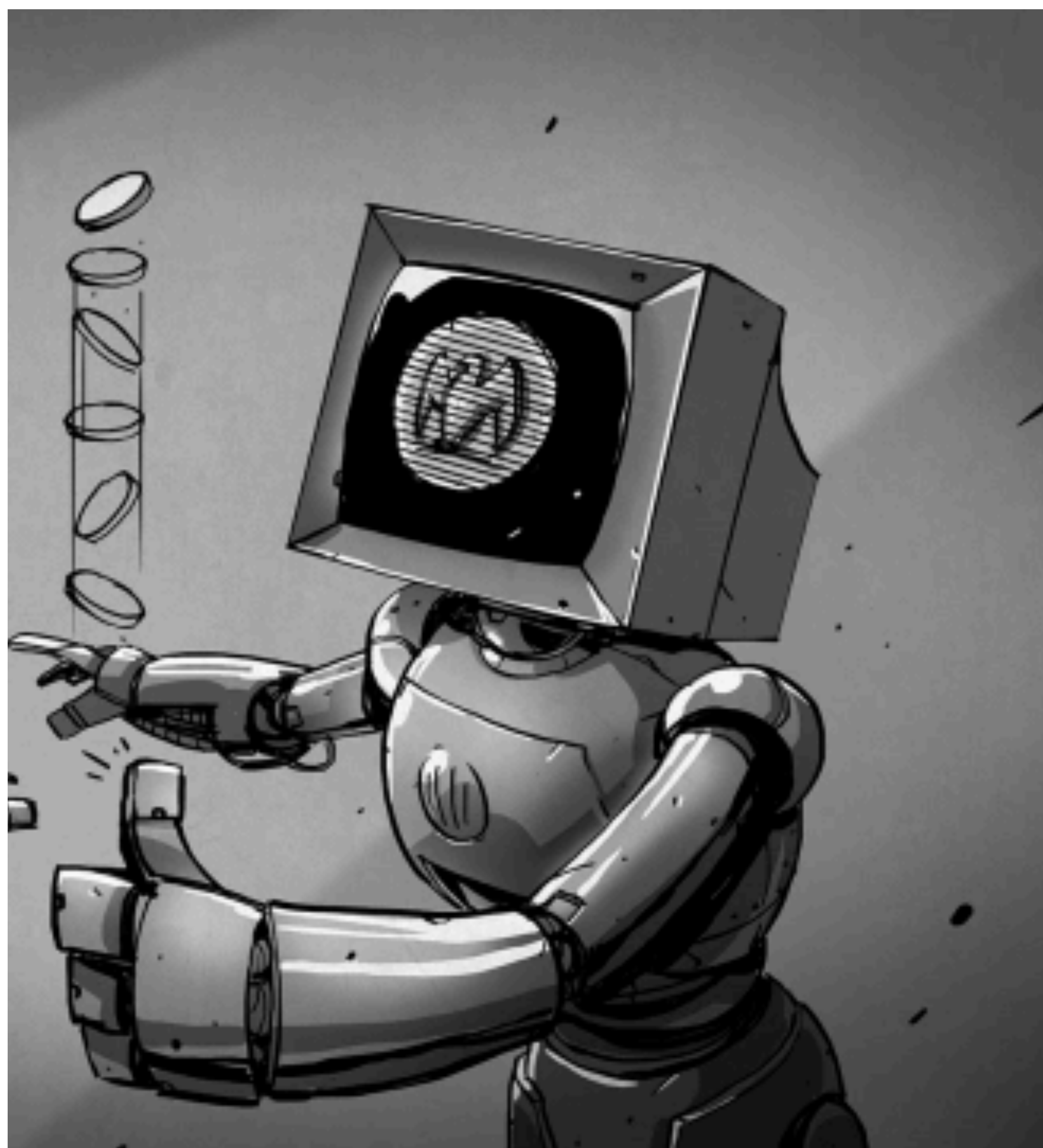


1. Information Gain does not take into account travel cost!

2. Uniform Bernoulli prior may not be the best prior!

Can we find a better
exploration / exploitation
algorithm?





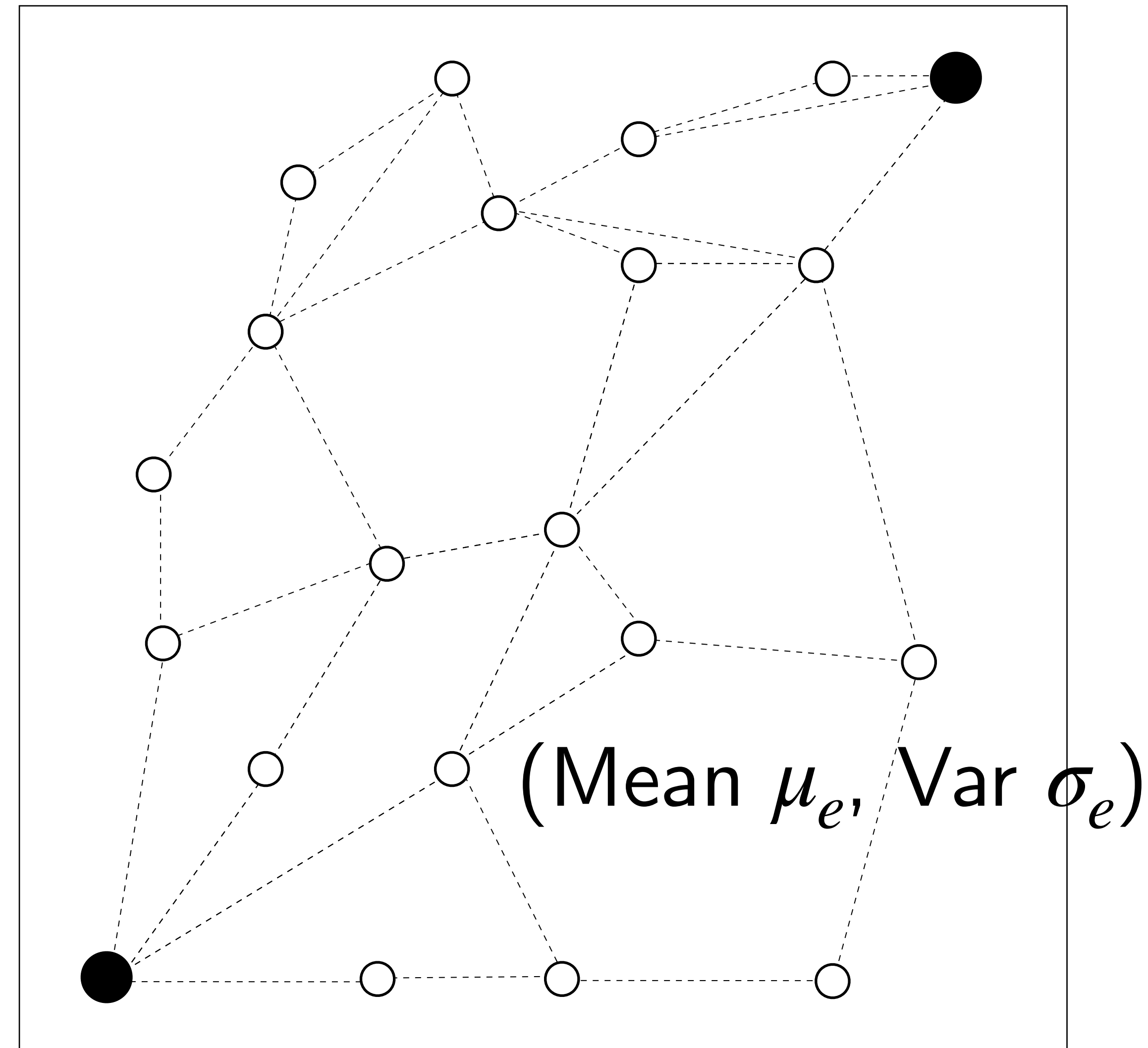
Posterior Sampling

The Online Shortest Path Problem

You just moved to Cornell and are traveling from office to home.

You would like to get home quickly but you are uncertain about travel times along each edge

Suppose we had a prior on travel time for each edge
(Mean μ_e , Var σ_e)



What if ...

... we just sampled travel times from our prior and solved the shortest path?



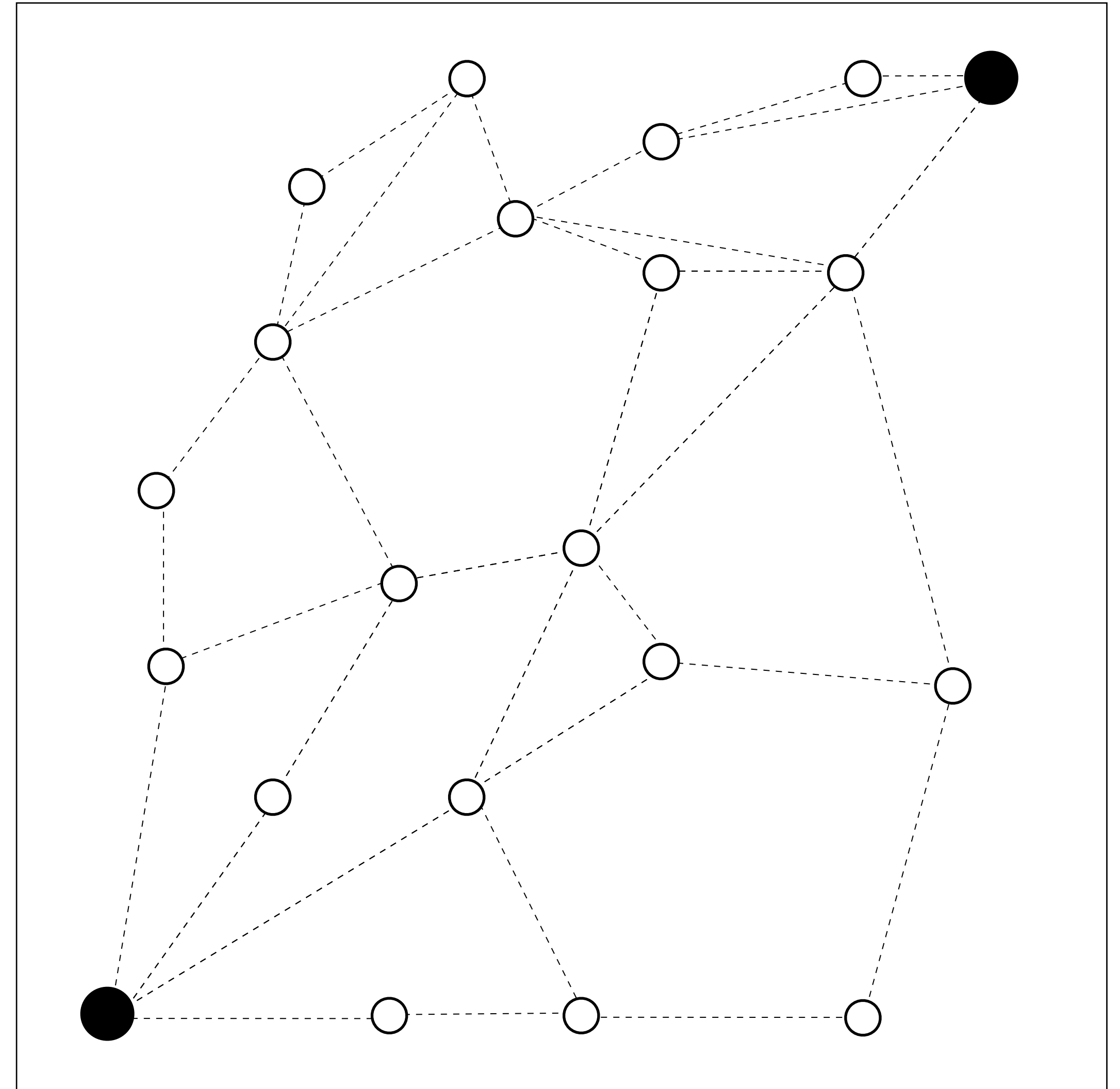
A suspiciously simple algorithm

Repeat forever:

Sample edge times from posterior

Compute shortest path

Travel along path, and update posterior



A suspiciously simple algorithm

Repeat forever:

Sample model from posterior

Compute optimal policy

Execute policy, observe s, a, s' ,
Update model

A Tutorial on Thompson Sampling

Daniel J. Russo¹, Benjamin Van Roy², Abbas Kazerouni², Ian Osband³ and Zheng Wen⁴

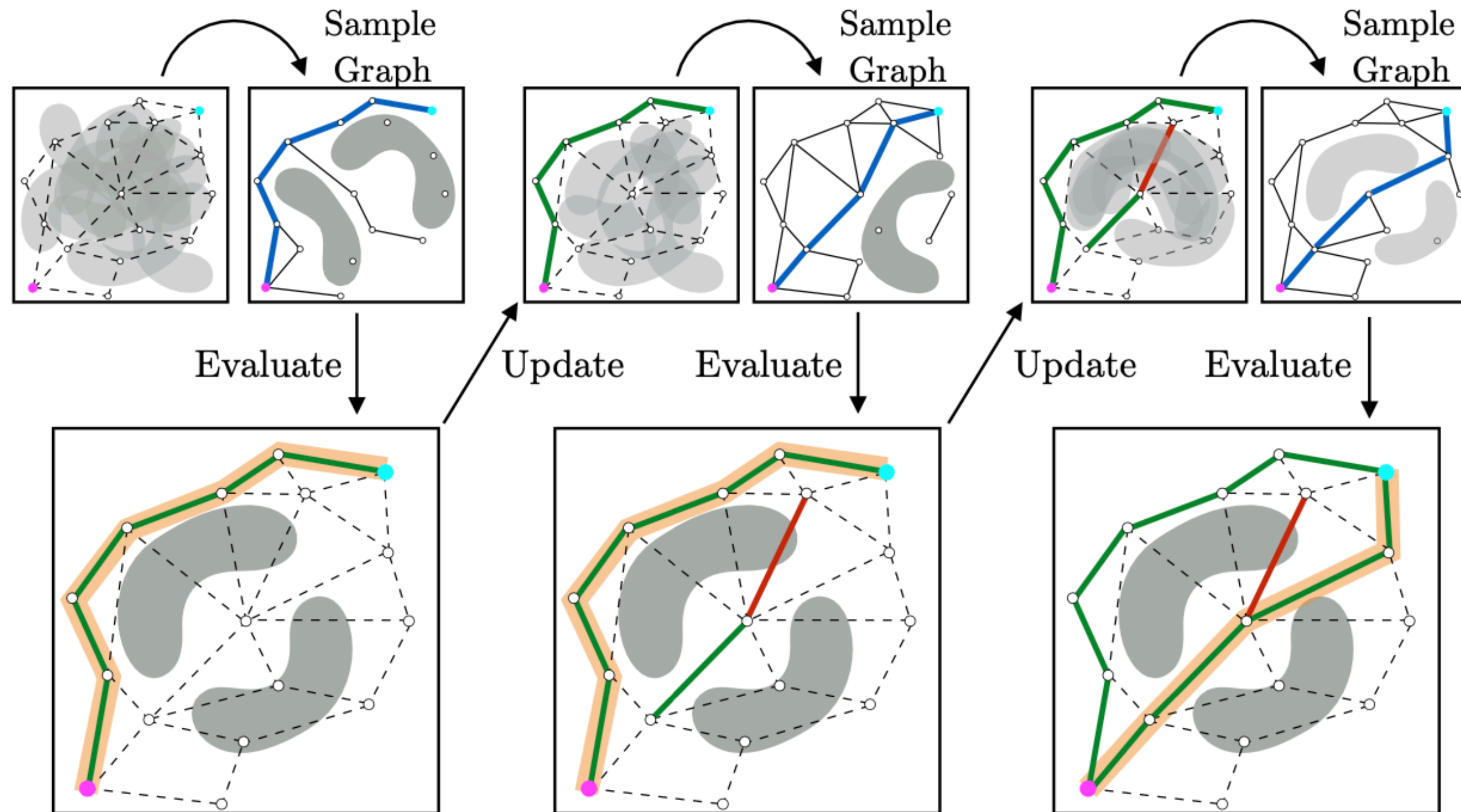
¹*Columbia University*

²*Stanford University*

³*Google DeepMind*

⁴*Adobe Research*

Posterior Sampling for Motion Planning

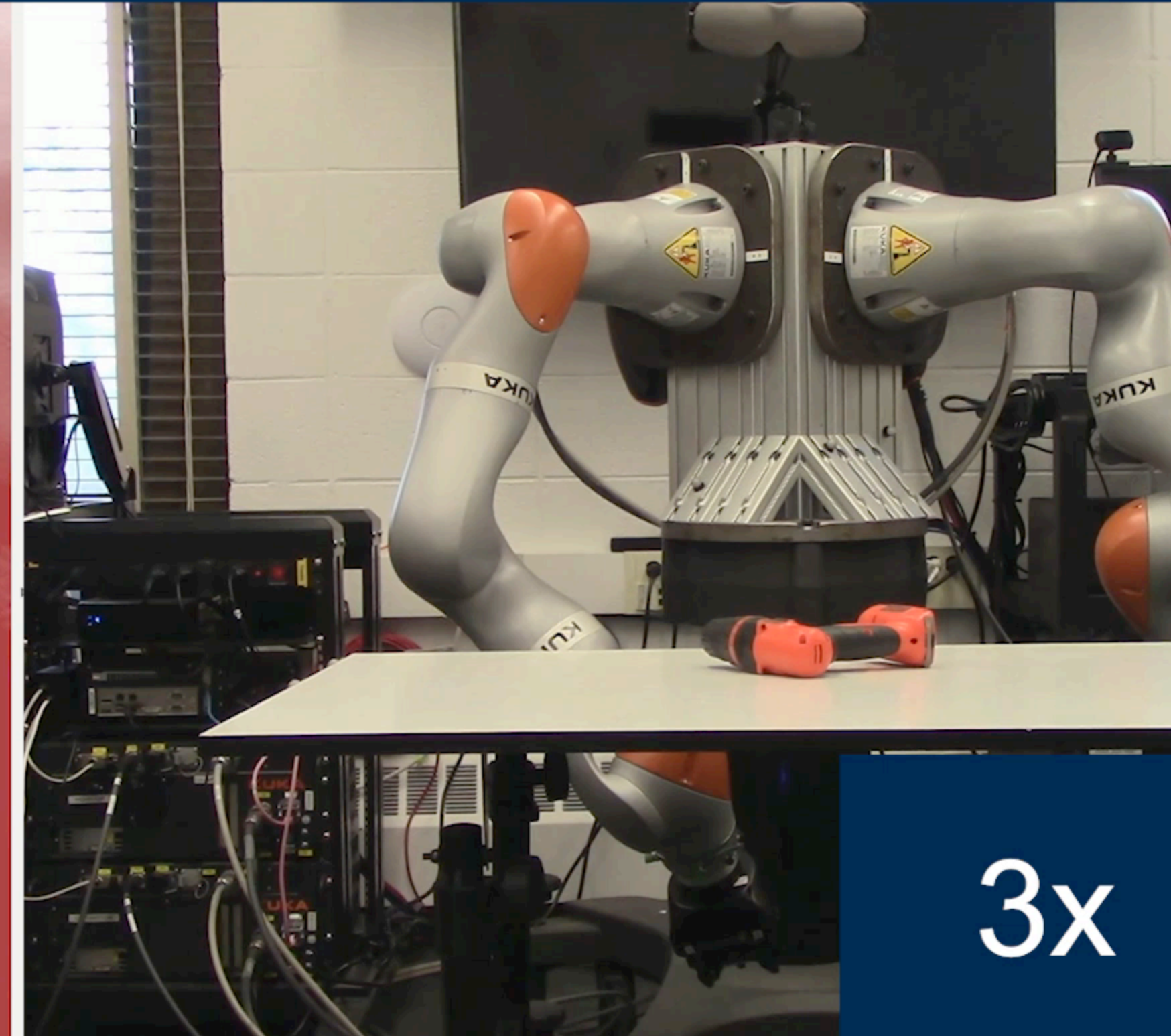
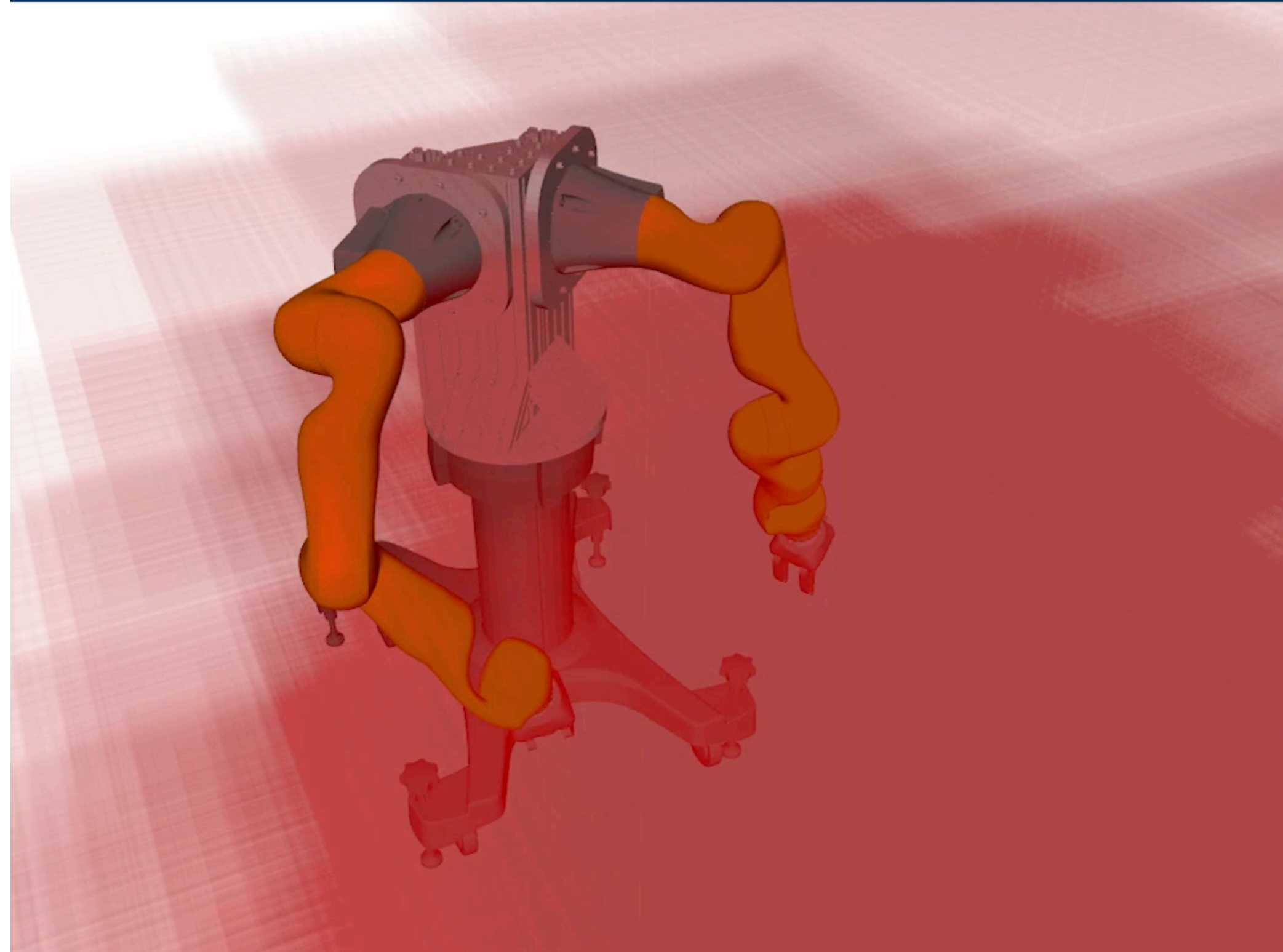


**Posterior Sampling for Anytime Motion Planning
on Graphs with Expensive-to-Evaluate Edges**


Real Robot Problems!

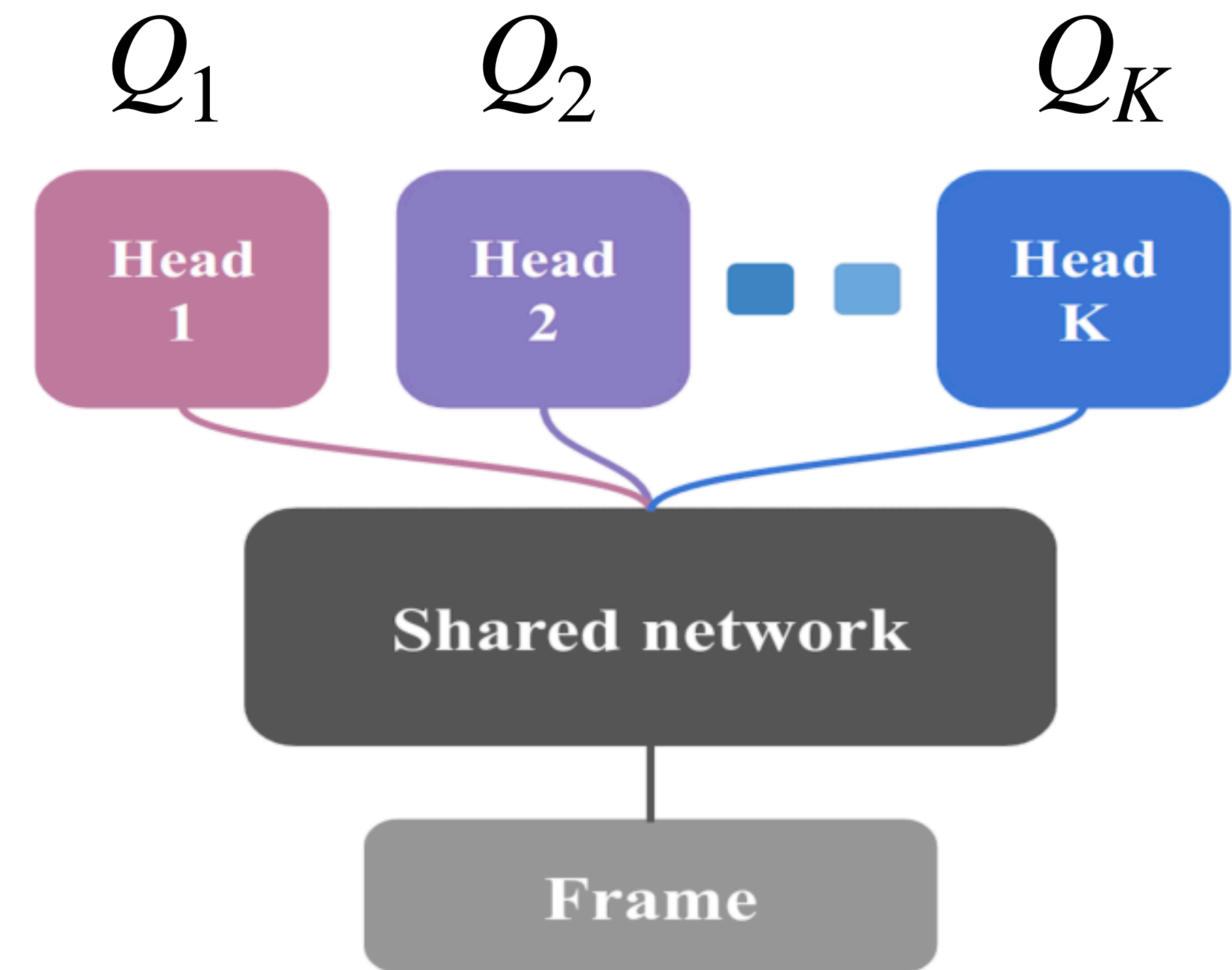


The Blindfolded Robot:
Bayesian Planning with Contact Feedback
[ISRR'19]



Posterior Sampling for Reinforcement Learning

- 
1. sample Q-function Q from $p(Q)$
 2. act according to Q for one episode
 3. update $p(Q)$




Bootstrapped Q Network

Deep Exploration via Bootstrapped DQN

Ian Osband^{1,2}, Charles Blundell², Alexander Pritzel², Benjamin Van Roy¹
¹Stanford University, ²Google DeepMind
{iosband, cblundell, apritzel}@google.com, bvr@stanford.edu

Posterior Sampling for Reinforcement Learning

Atari

- 
1. sample Q-function Q from $p(Q)$
 2. act according to Q for one episode
 3. update $p(Q)$

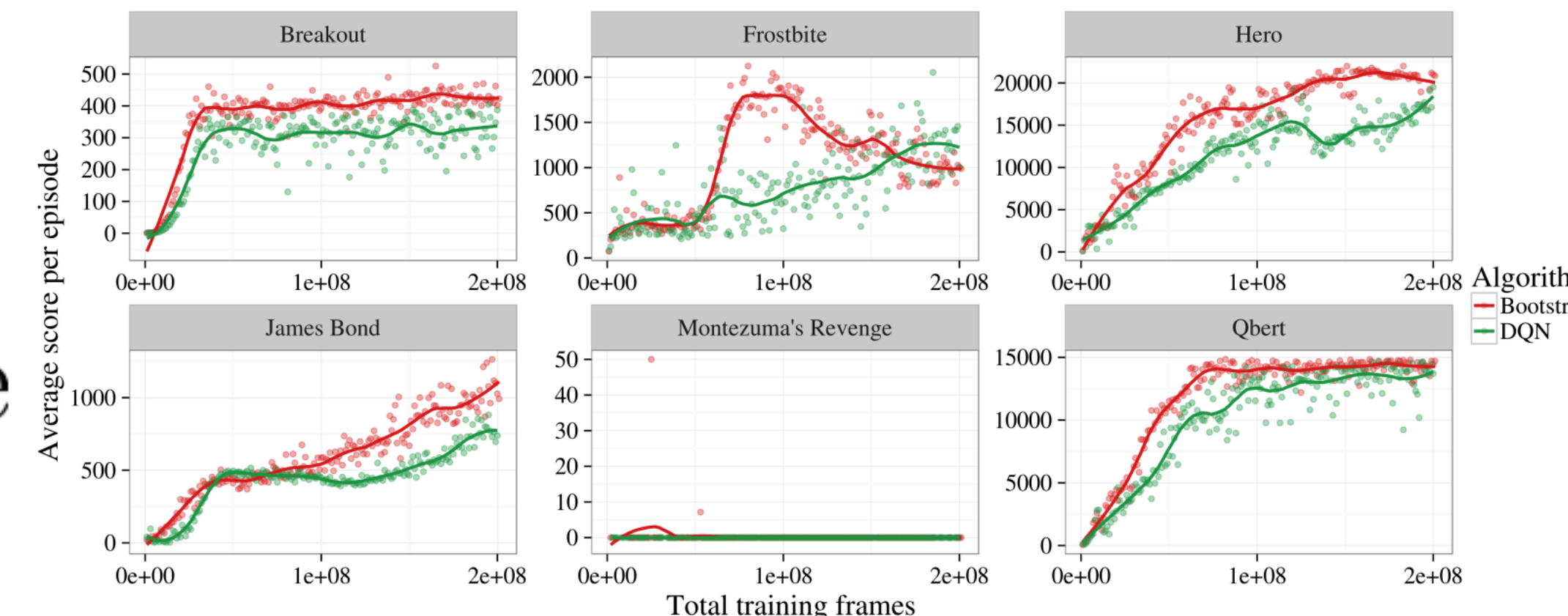


Figure 6: Bootstrapped DQN drives more efficient exploration.

Why does work better than taking random actions?

tl;dr

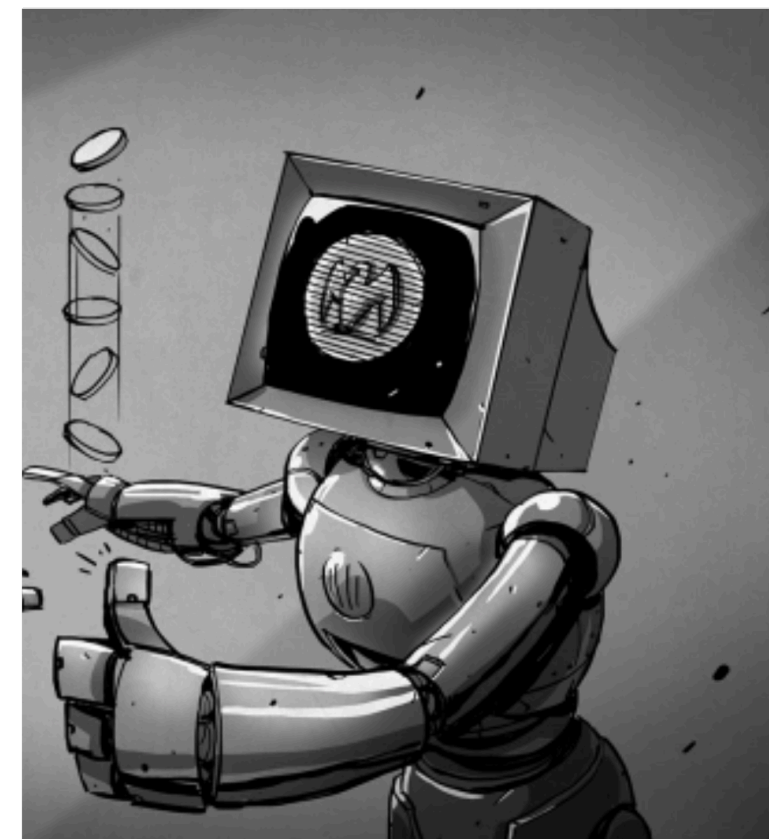


Belief Space Planning is NP-Hard
at best, undecidable at worst

Need to relax our problem!



Optimism
in the Face of
Uncertainty
(OFU)



Posterior
Sampling



Information
Gain

What is my prior is
intractable to represent
and sample from?

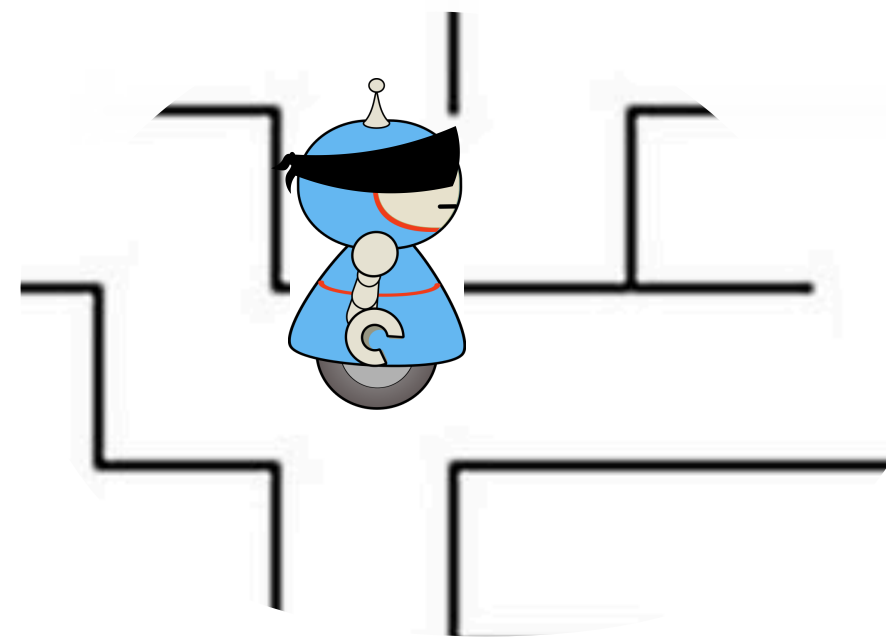




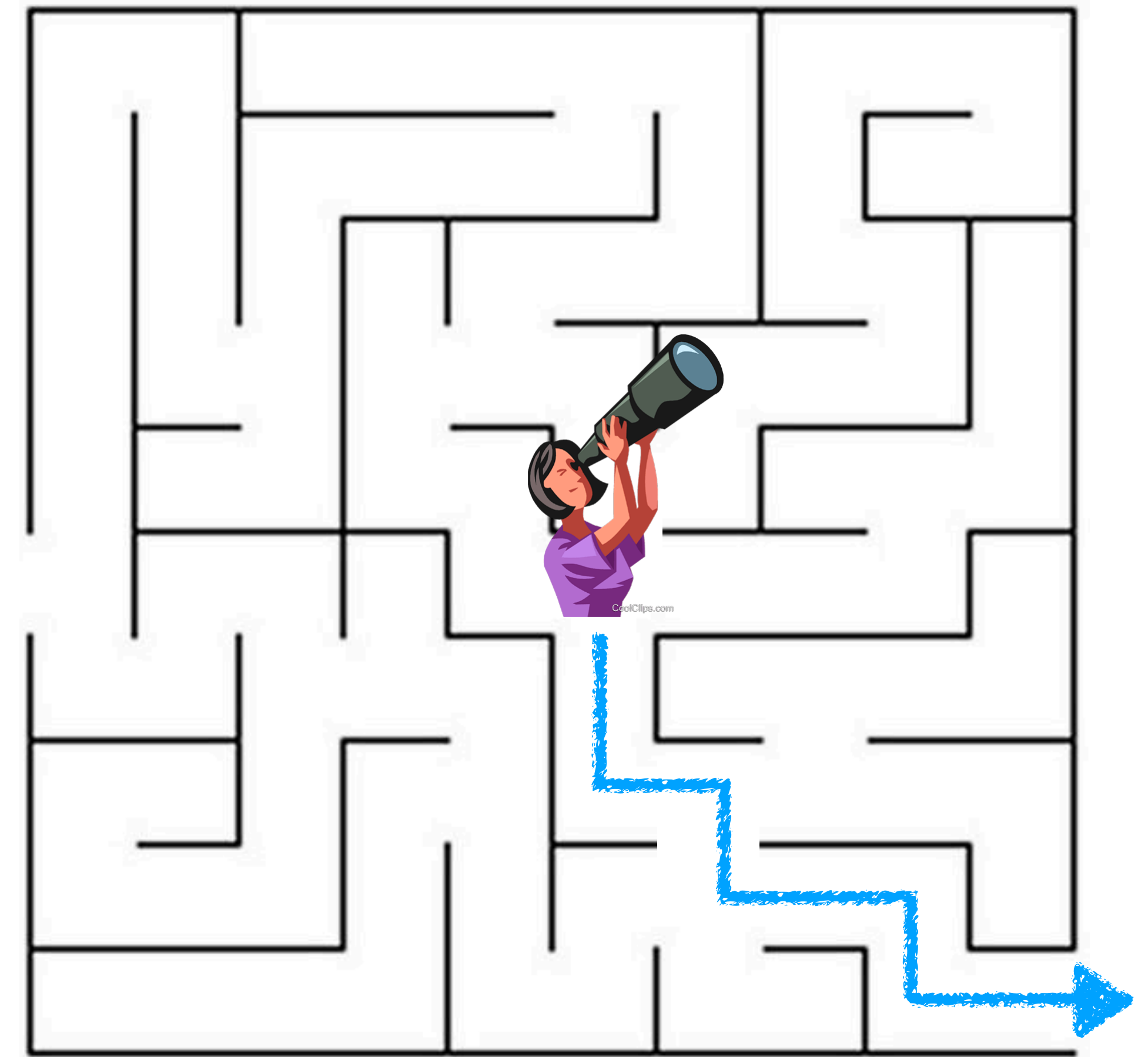
What if ...

... we trained a learner to
imitate a clairvoyant oracle?

Imitating Experts with Privileged Information



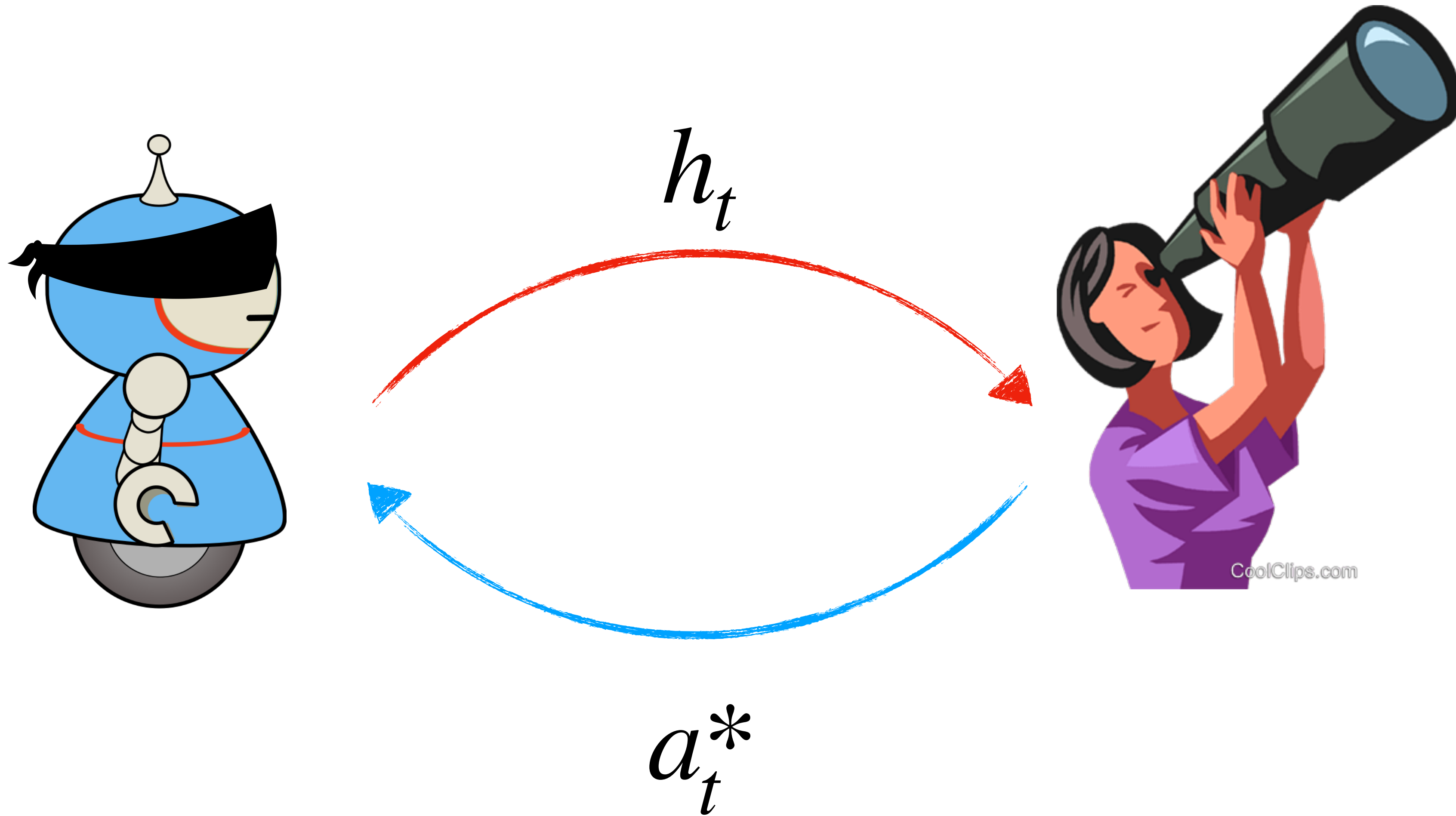
Imitate



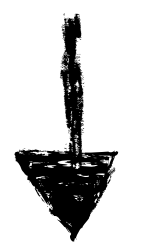
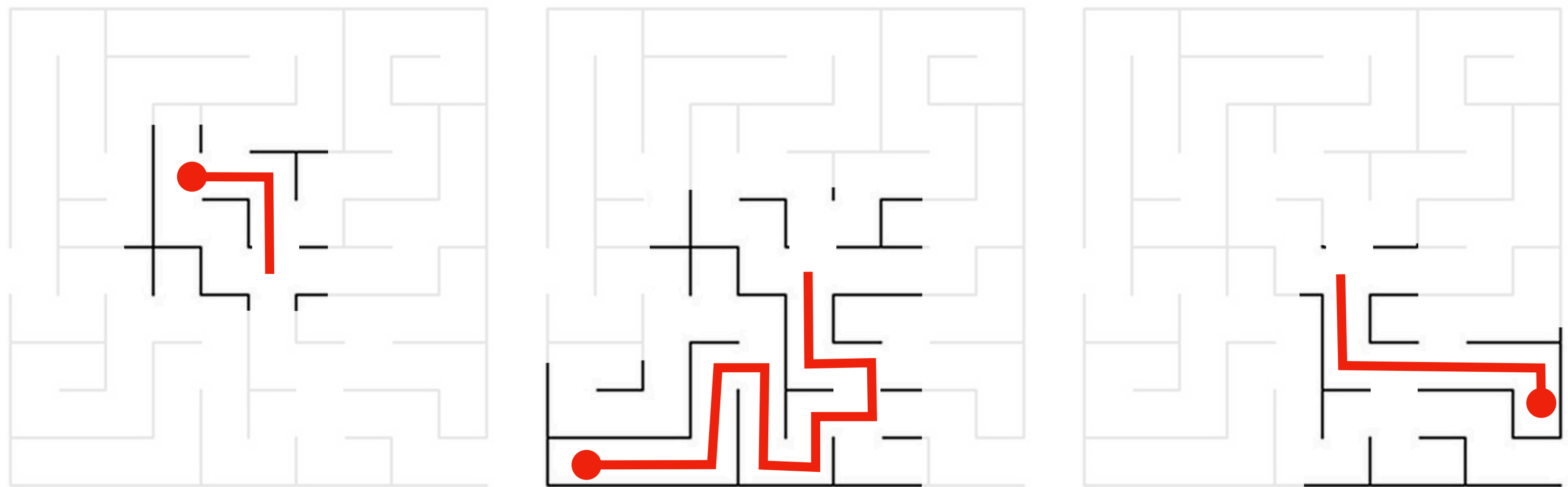
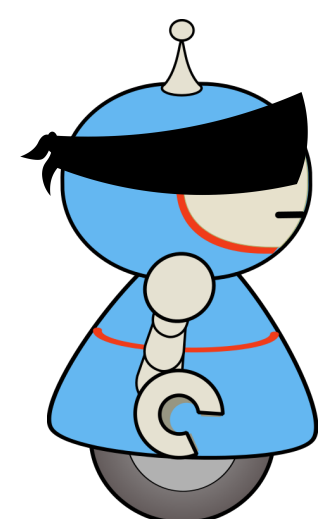
Learner
w/ limited sensing

Expert
can see further

Solution: **Interactively** query expert

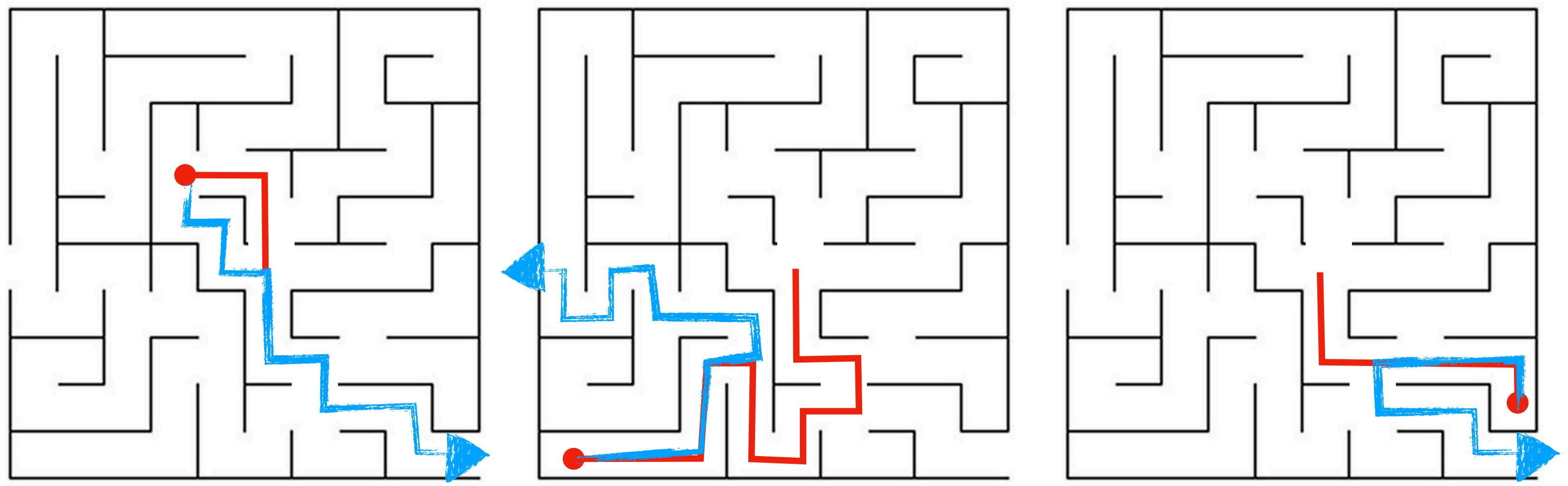


Solution: Interactively query expert

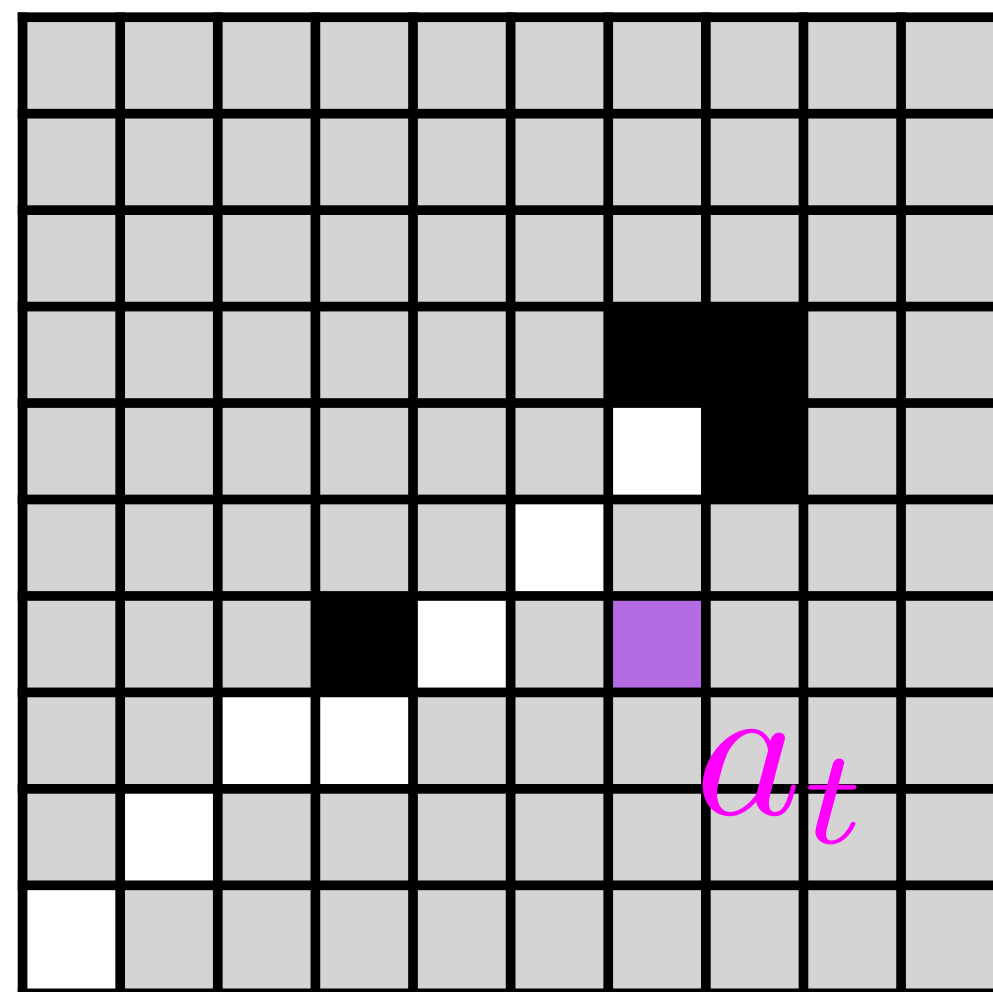


e.g DAGGER

- 1. Roll out learner
 - 2. Query Expert
 - 3. Aggregate Data
- and repeat!

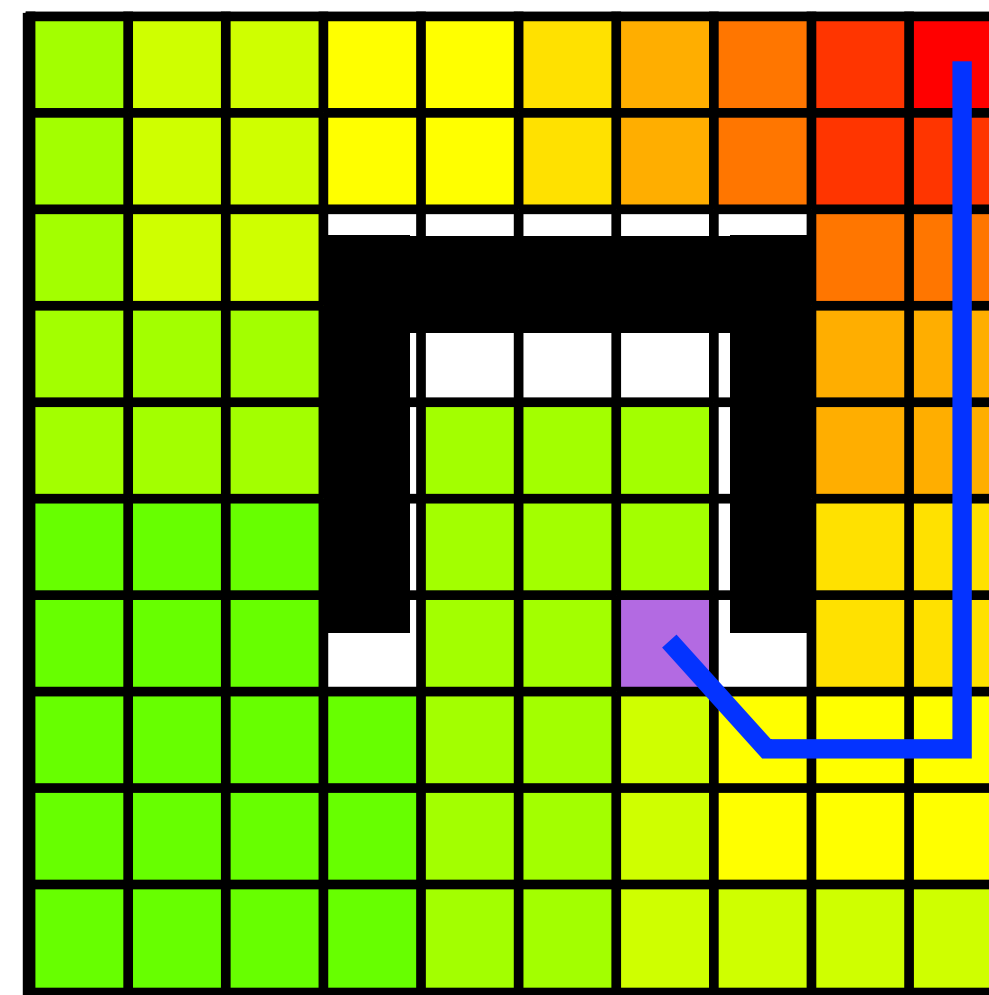


Privileged Information: Motion Planning

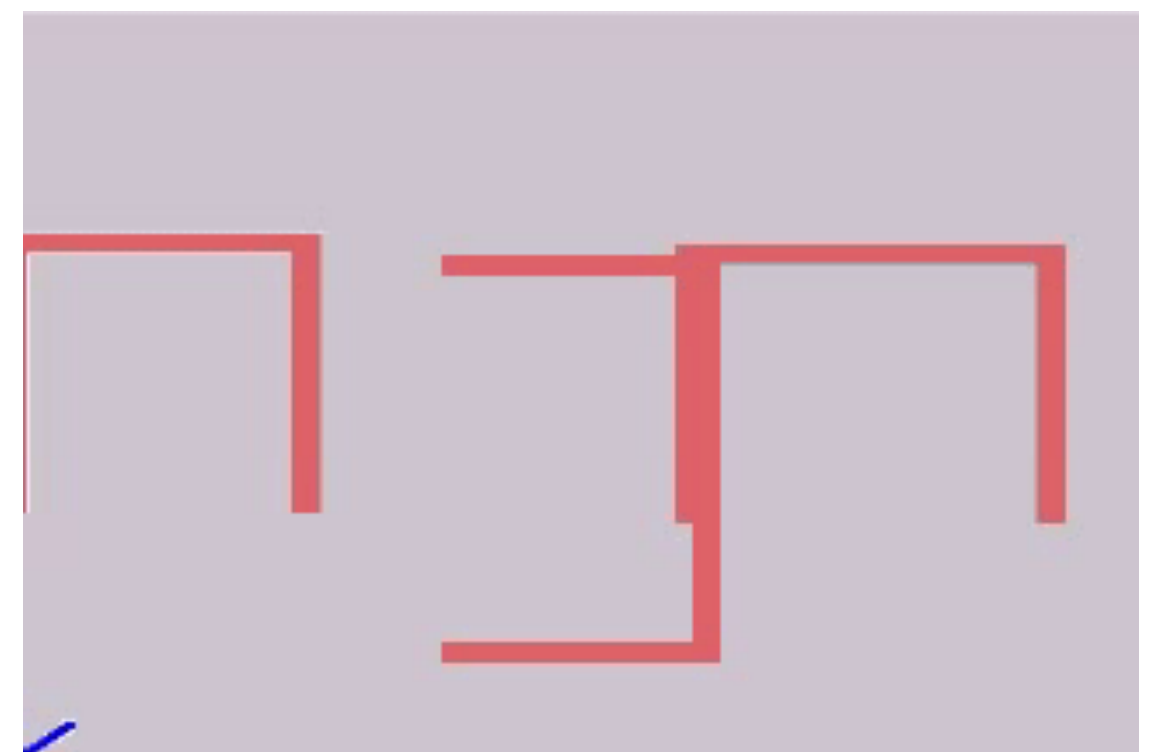
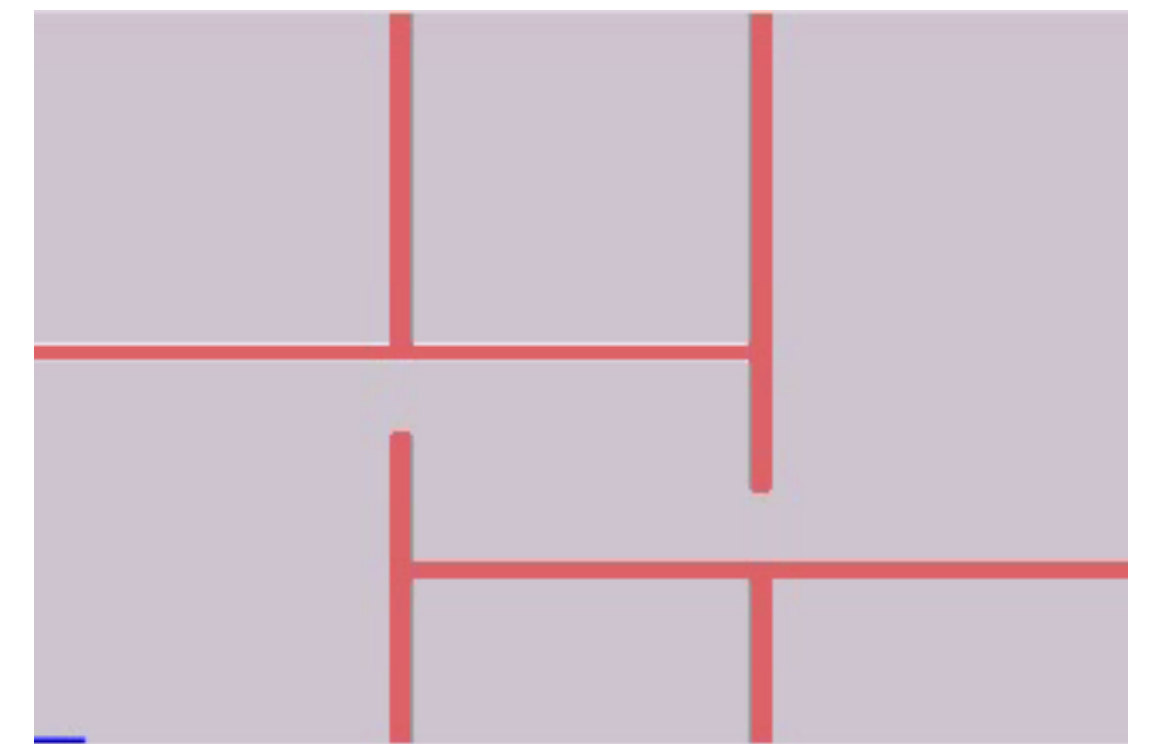
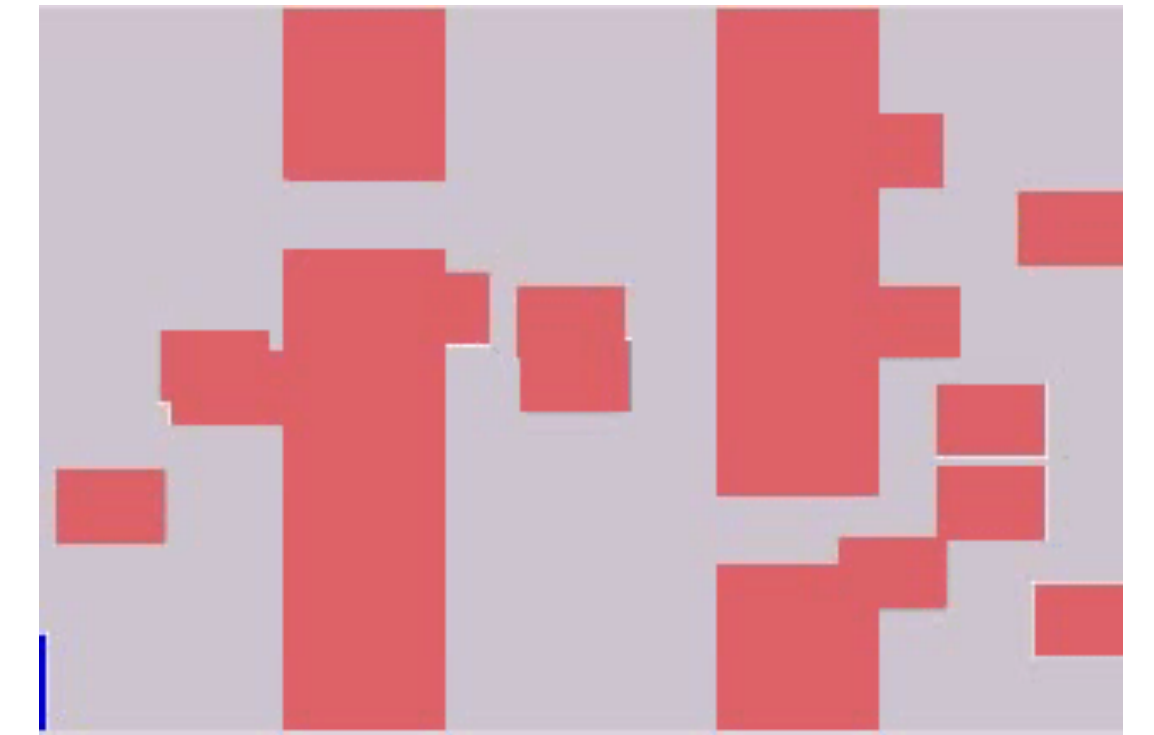


Learned
Search Heuristic

Imitate



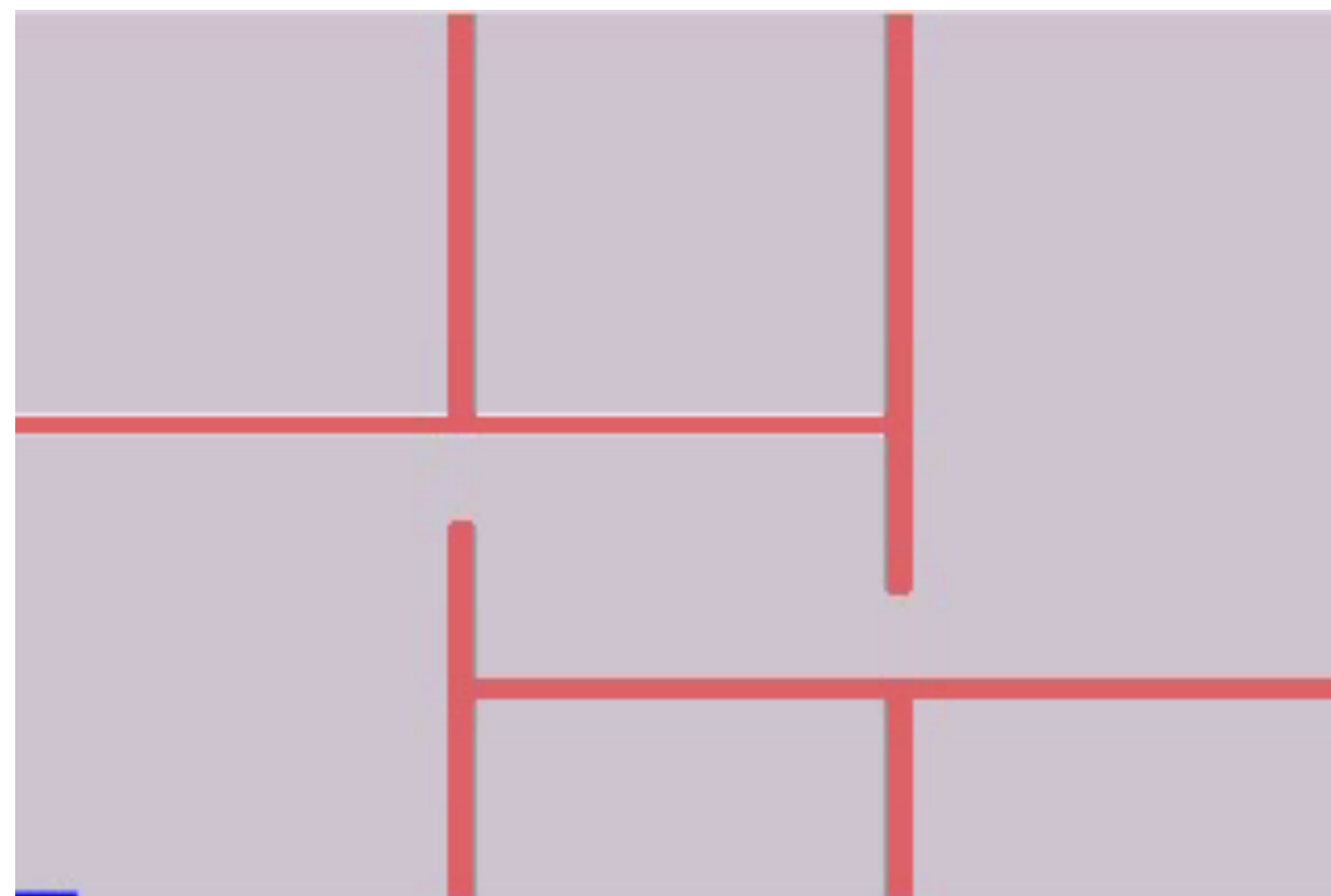
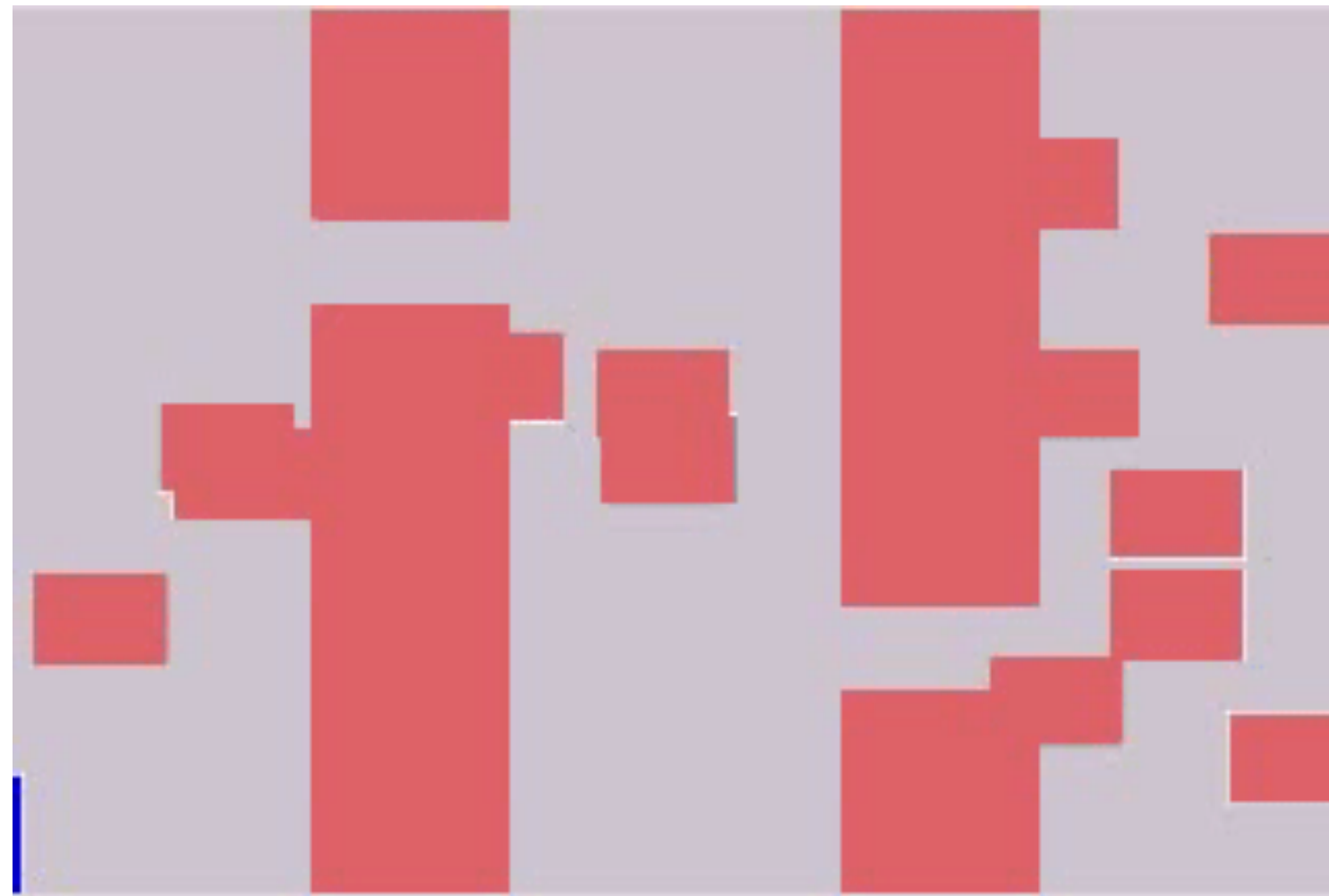
Optimal
Value Function



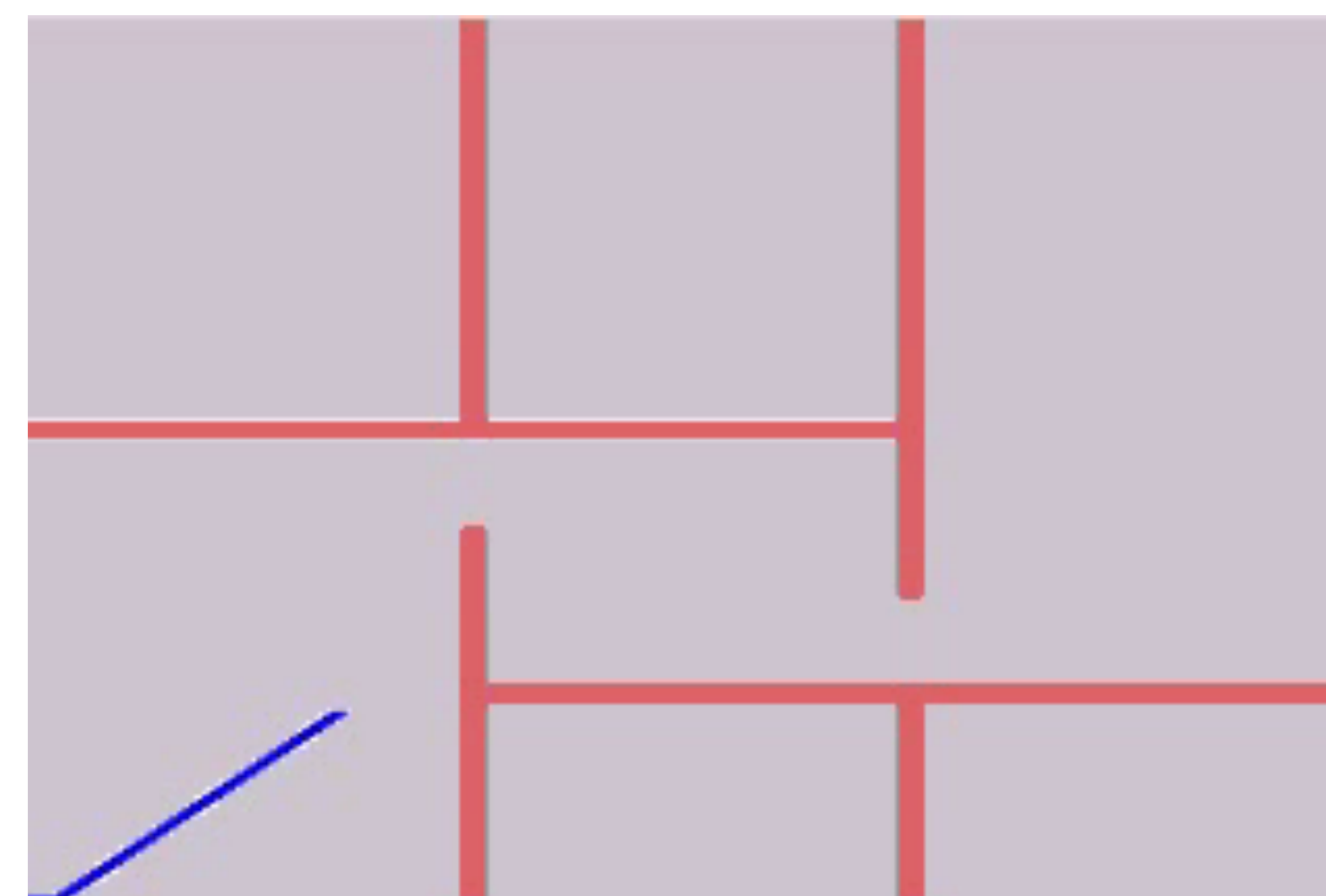
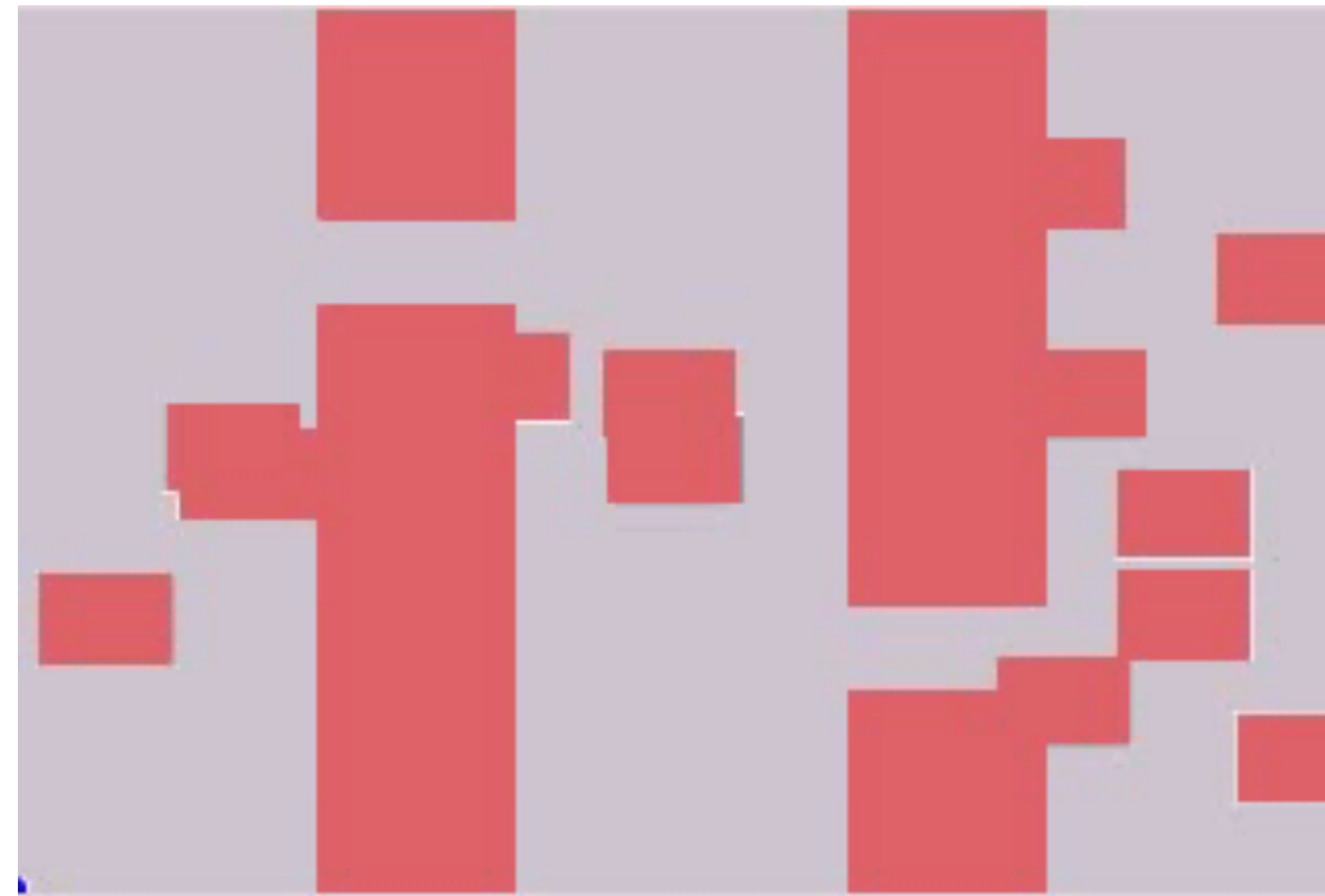
[Choudhury et al. '2018]

Example: Training search heuristics

On-policy (Aggrevate)



Behavior Cloning



Why / When does this work?

Proved that this approximates Hindsight Optimization / QMDP

Fails when you need to explicitly explore (i.e. asymptotic realizability not hold)