

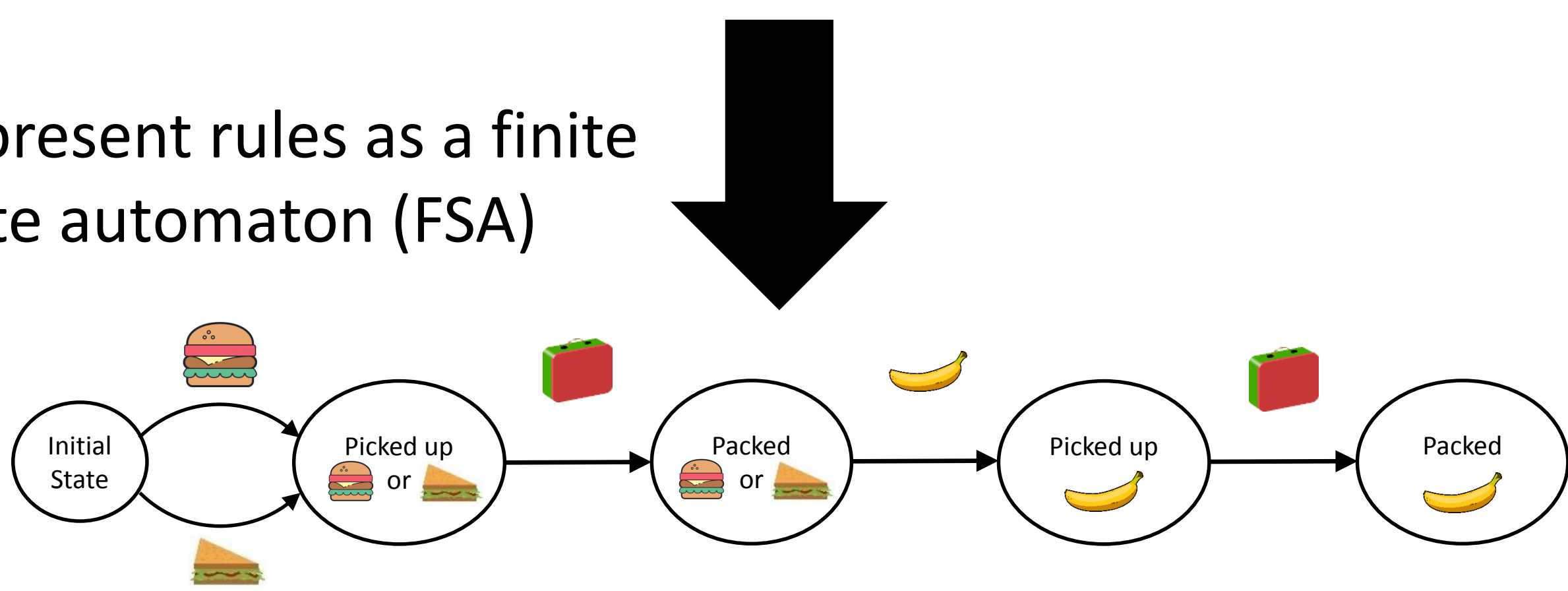


Goal: Learn rules and policies for rule-based environments in an **interpretable** and **manipulable** way

Interpretability

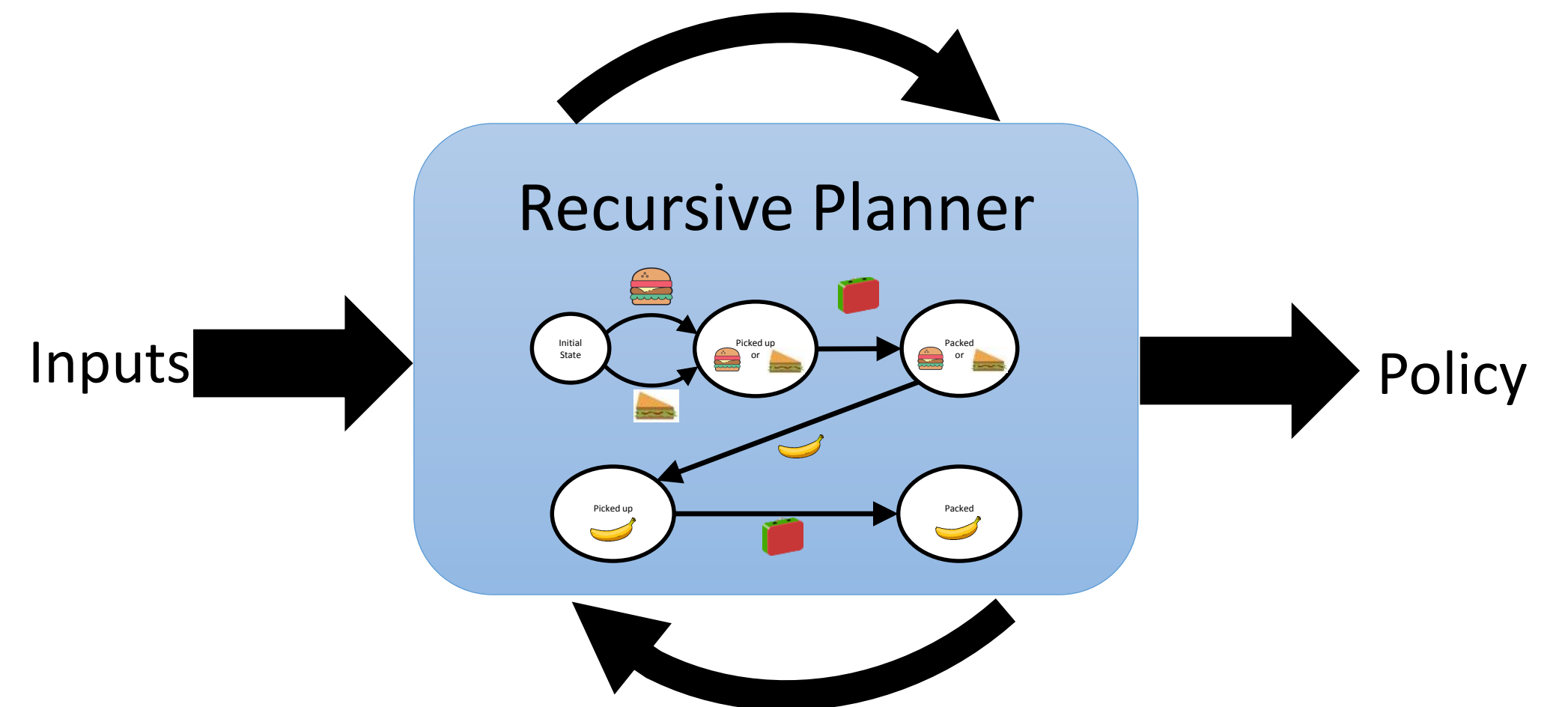
Pick up a sandwich or hamburger and pack it, then pick up a banana and pack it

Represent rules as a finite state automaton (FSA)



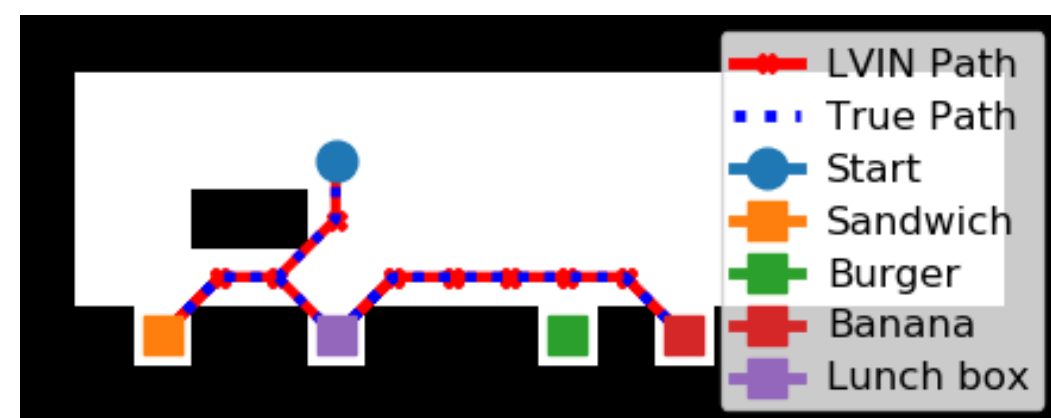
Manipulability

FSA is an input to a recursive planner, so changes to the FSA result in changes to the policy

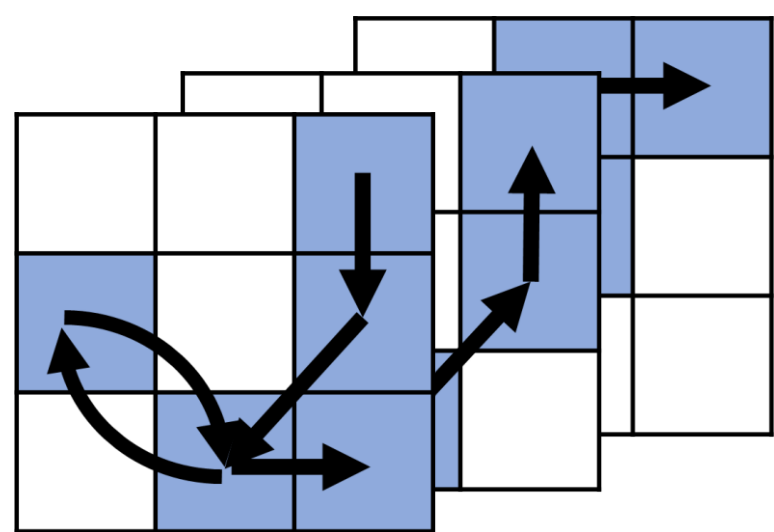


Data

A low-level environment

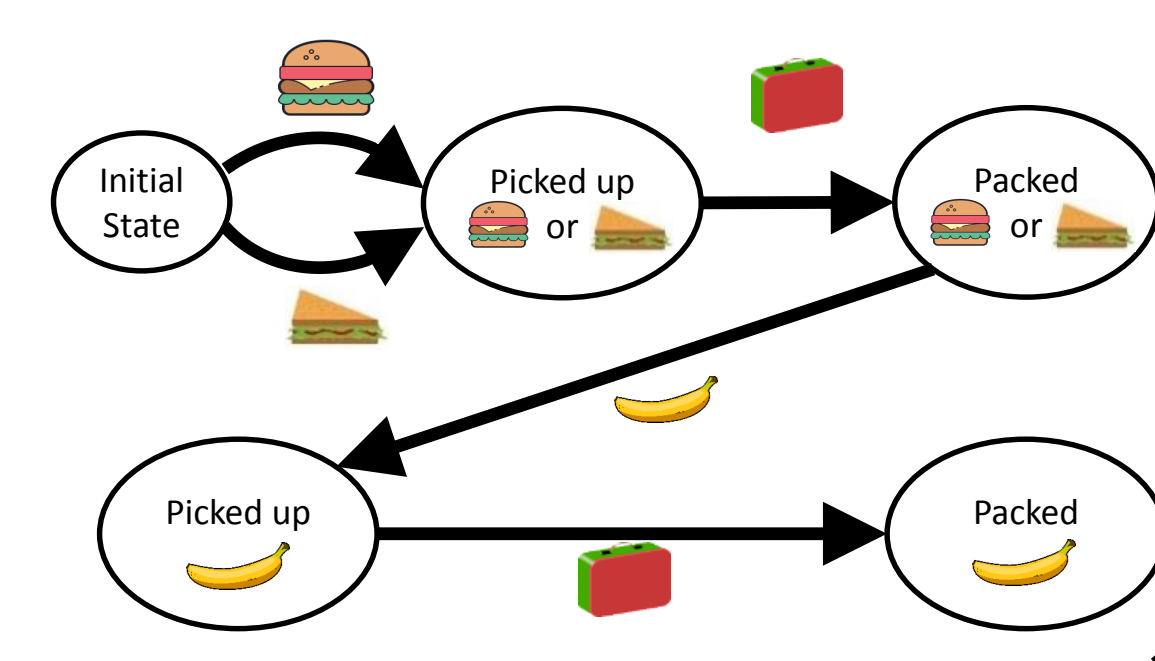


A dataset of trajectories



Learn

An FSA representation of the rules



A reward function

$\mathcal{R}(f) \rightarrow$ A policy

Summary of our Approach

1. Model the environment as a POMDP

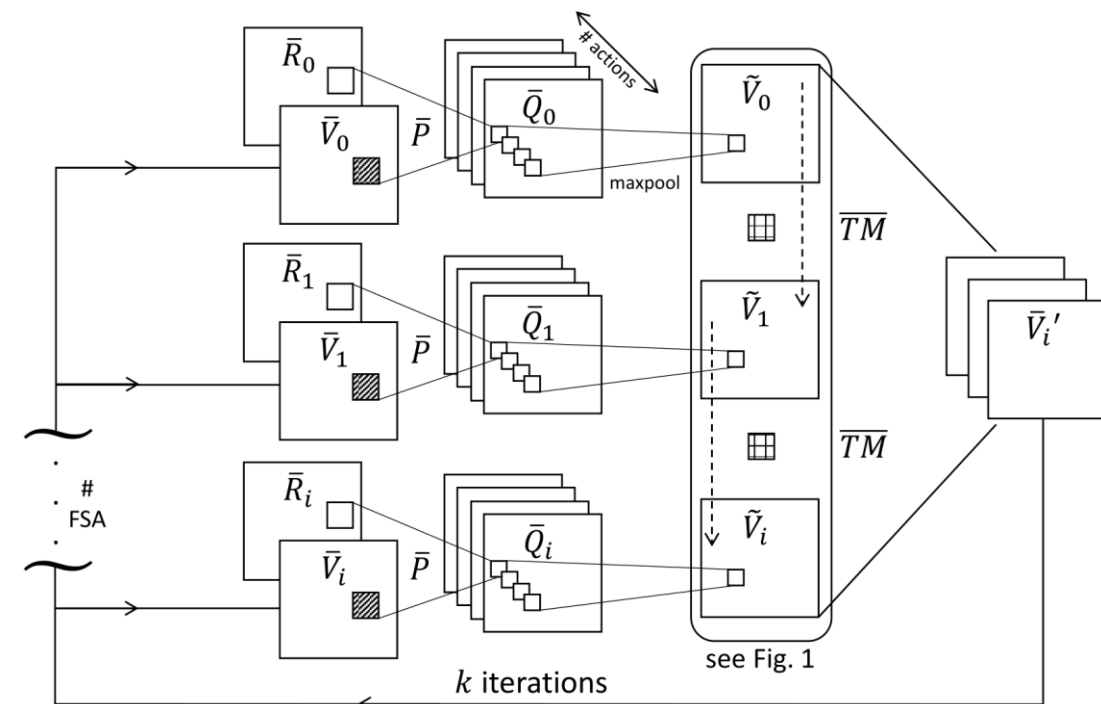
$$(\mathcal{S} \times \mathcal{P} \times \mathcal{F}, \mathcal{A}, T \times M \times \underline{TM}, \mathcal{R}, \mathcal{S} \times \mathcal{P}, \mathcal{O}, \gamma_d)$$

2. Parameterize the policy as a function of an FSA, a reward function, and a low-level environment

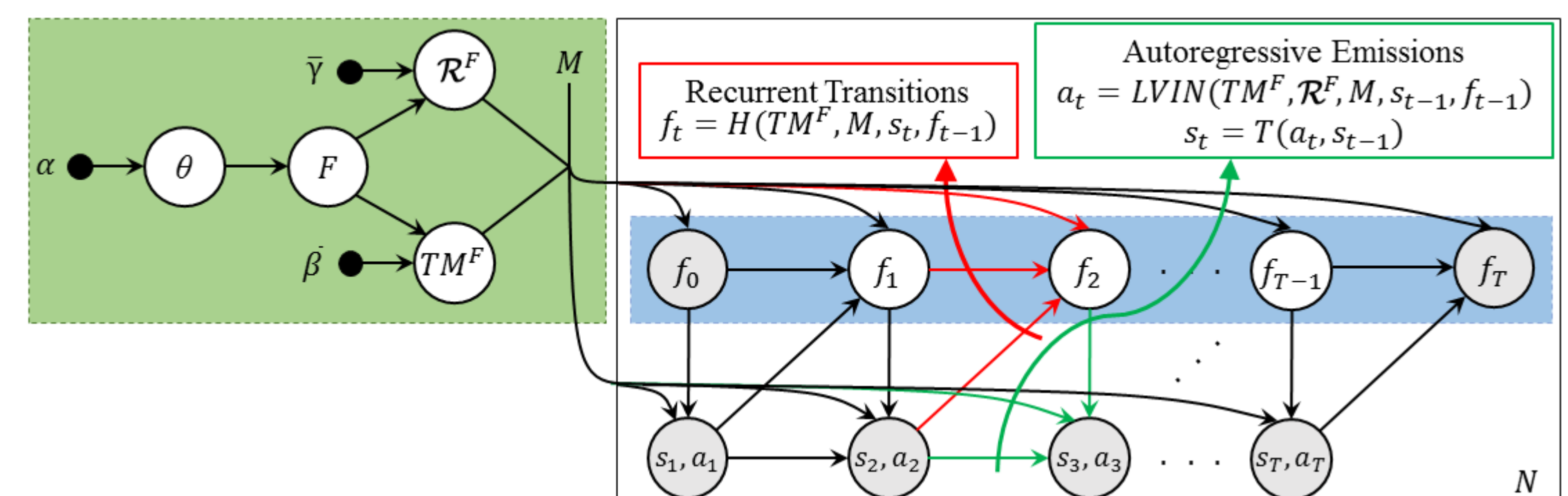
$$Q^{t+1}(s, f, a) \leftarrow R(s, f, a) + \gamma \sum_{s' \in \mathcal{S}} T(s'|s, a) V^t(s', f)$$

$$\hat{V}^{t+1}(s, f) \leftarrow \max_a Q^{t+1}(s, f, a)$$

$$V^{t+1}(s, f) \leftarrow \sum_{f' \in \mathcal{F}} TM(f'|f, M(s)) \hat{V}^t(s, f')$$



3. Model the policy rollout on the POMDP as an HMM

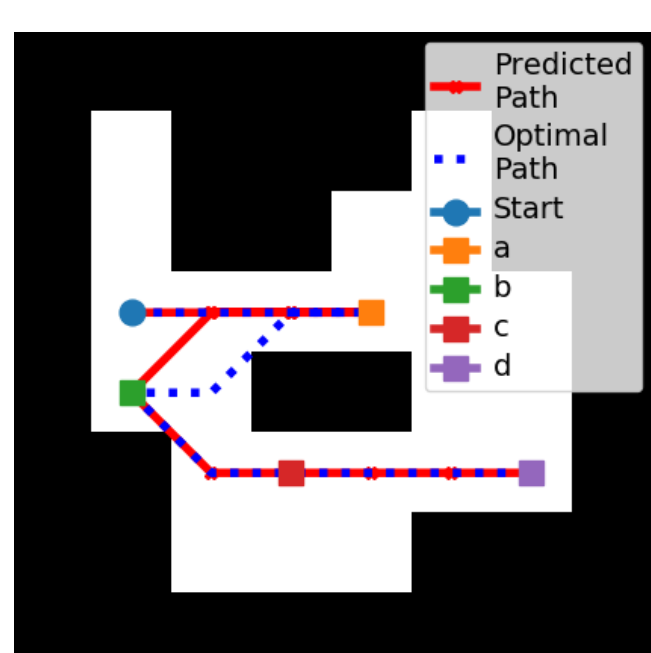


4. Use variational inference to infer the latent variables of the HMM including FSA and reward function

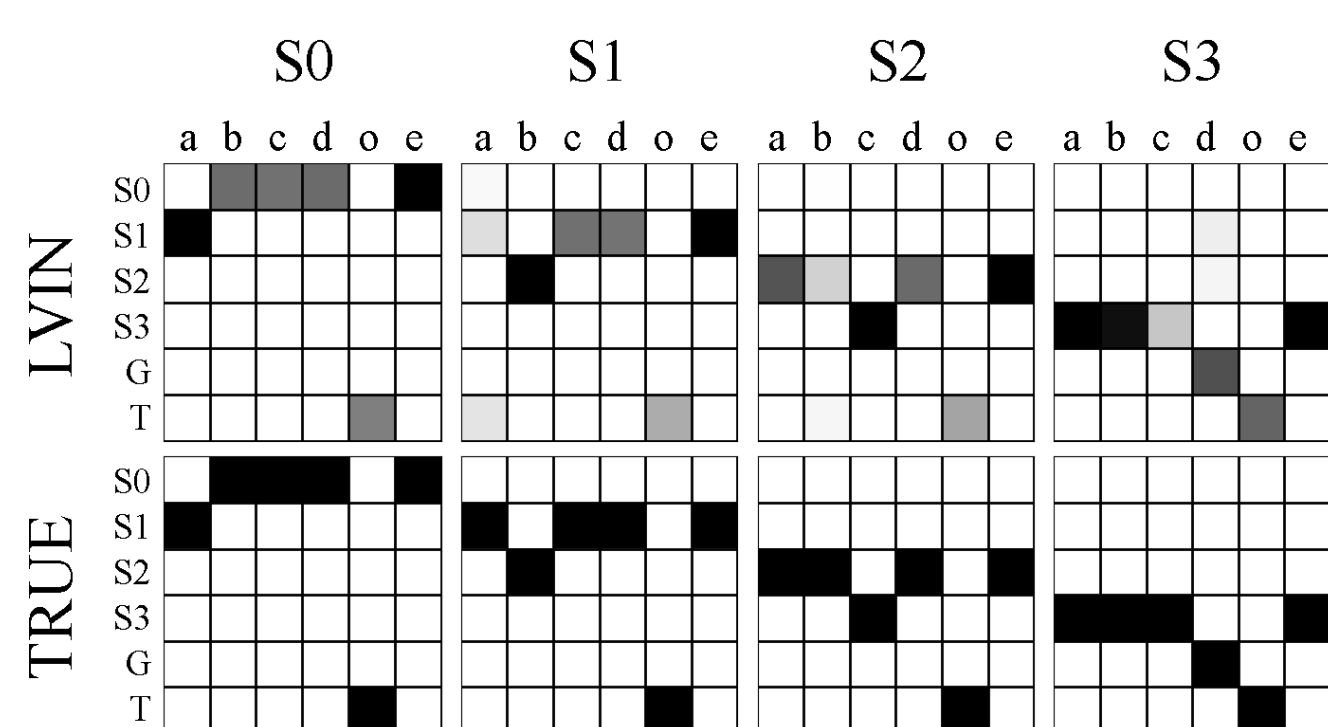
$$p(\overline{TM}, \overline{R}, \theta | \mathcal{D}, \alpha, \beta, \gamma) = \frac{p(\mathcal{D}, \overline{TM}, \overline{R}, \theta | \alpha, \beta, \gamma)}{p(\mathcal{D} | \alpha, \beta, \gamma)}$$

Experiments and Results

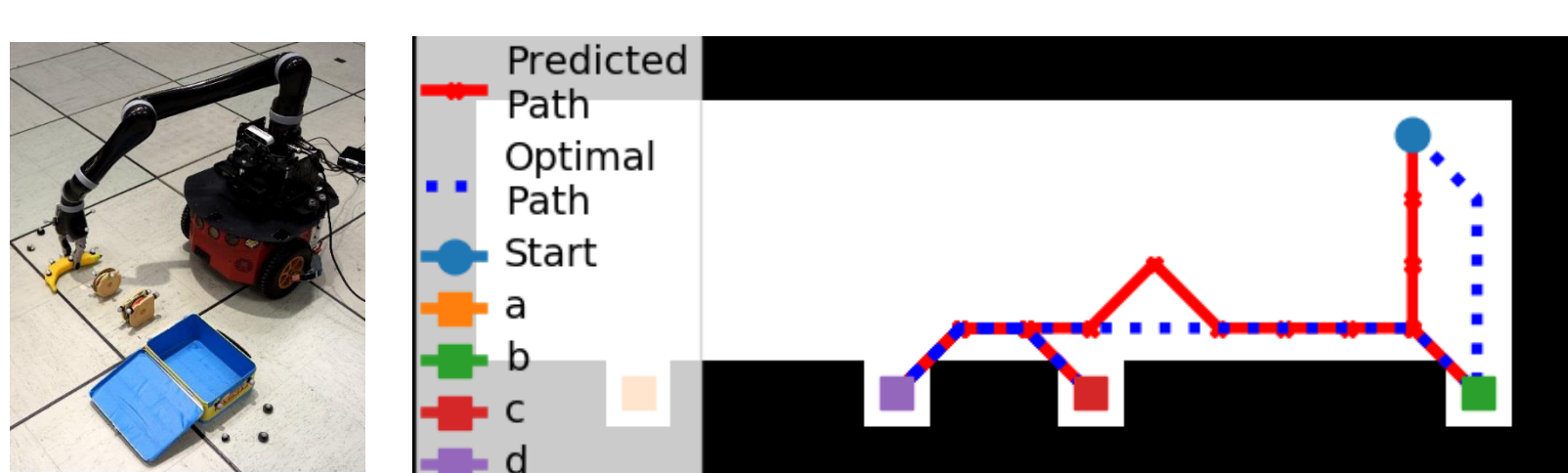
Gridworld



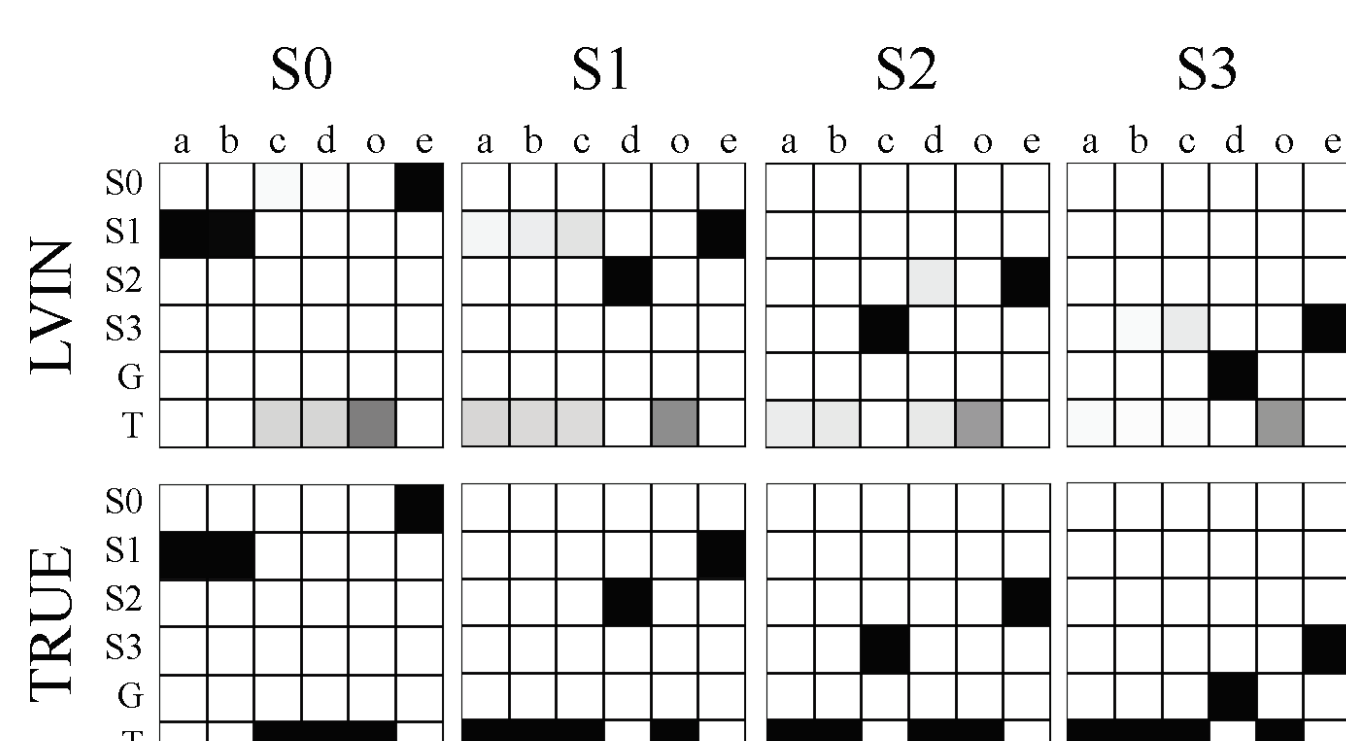
Go to a, then b, then c, then d



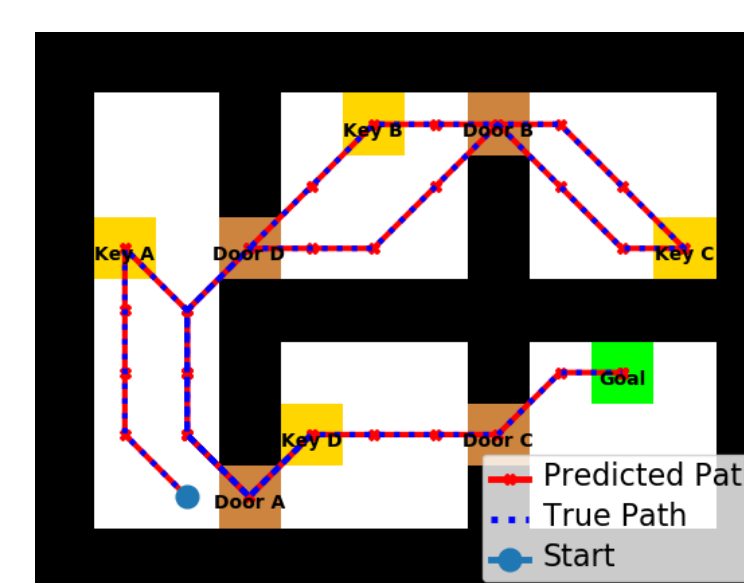
Lunchbox Packing



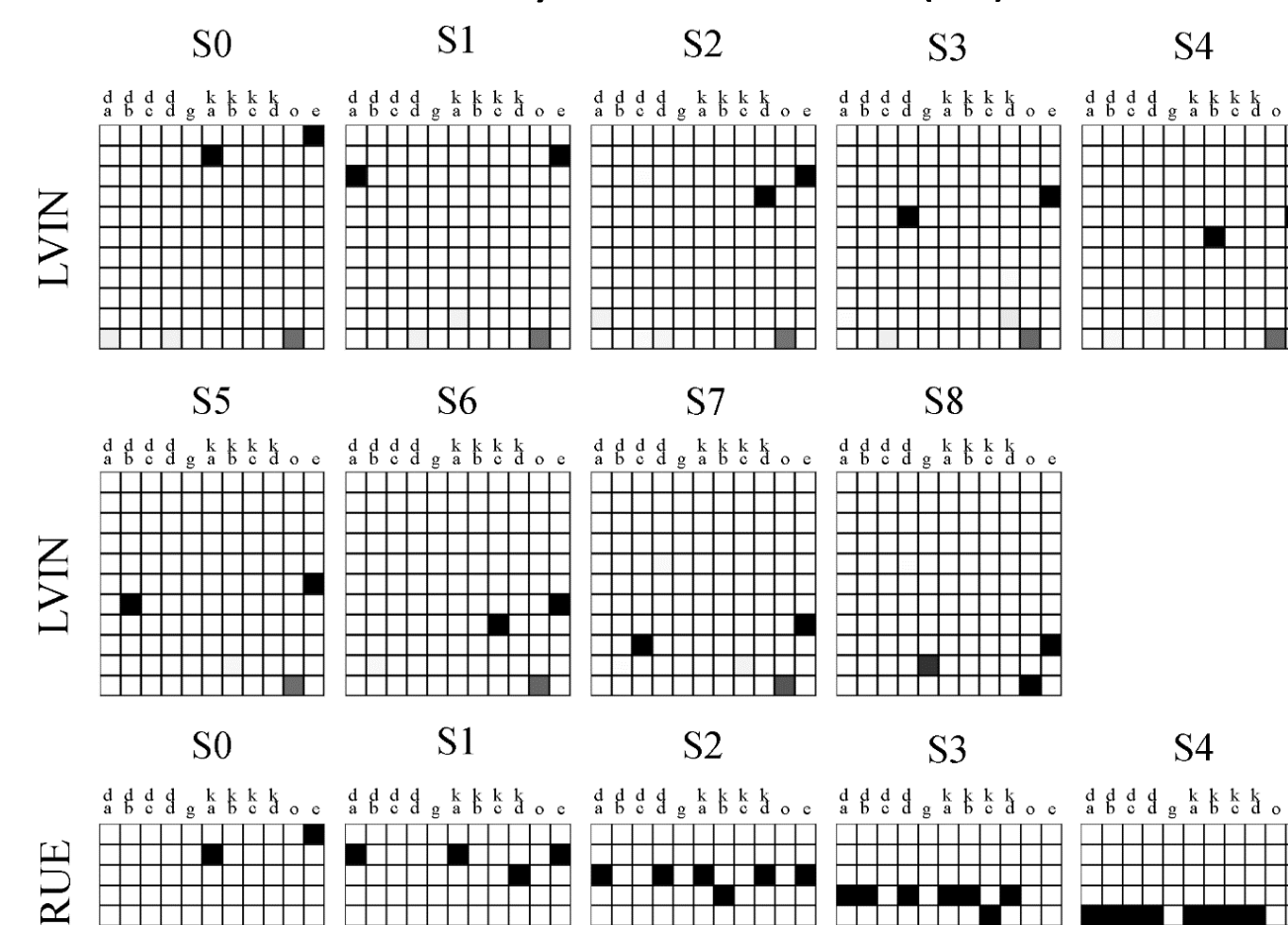
Pick up sandwich (a) or hamburger (b), put it in the lunchbox (d), then pick up banana (c) and put it in the lunchbox (d)



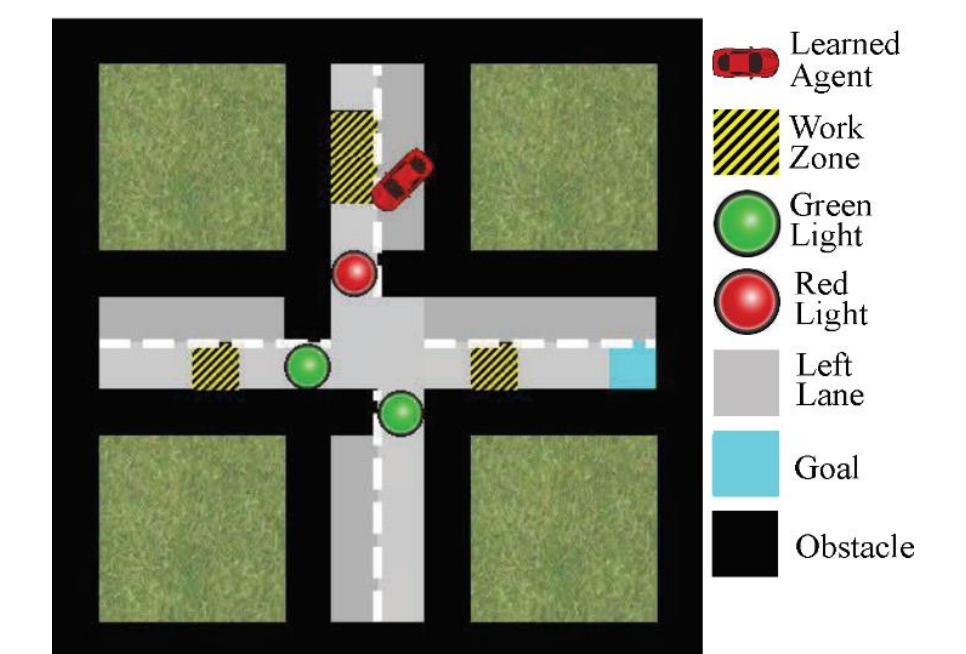
Dungeon



Go to the goal (g); can't pass through Door x (dx)



Driving Domain



- Go to the goal and avoid work zones and obstacles
- Stop if there's a red light in front of you and go if there's a green light
- Prefer the right lane to the left lane

