

Behavior Sets for Adaptive Social Robotics

Chris Reinke

Inria Grenoble / Robotlearn



Socially Pertinent Robots in Gerontological Healthcare

Overall goal: Socially assistive robots



Goals of the team:

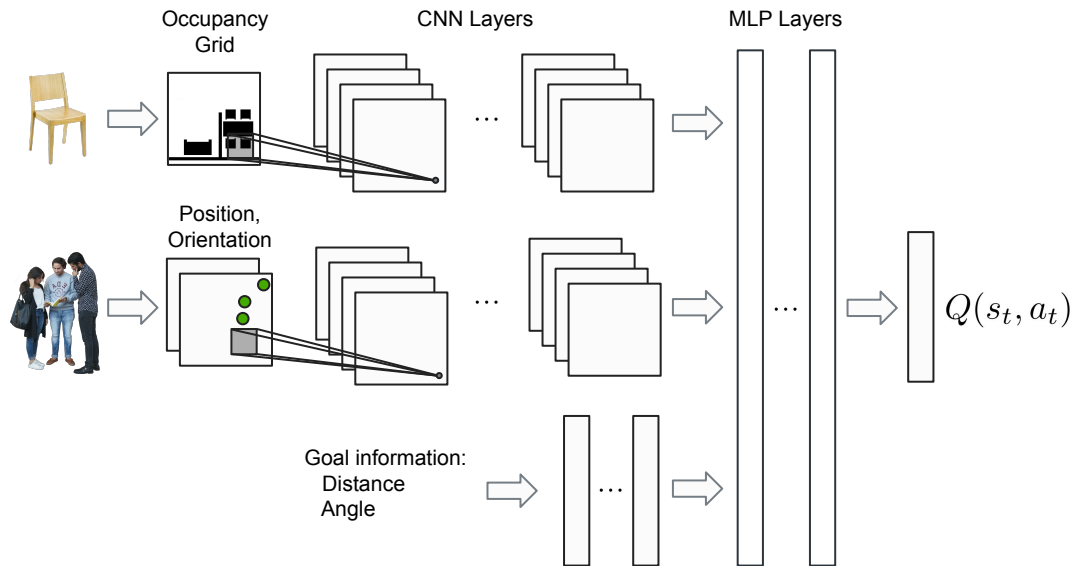
- ❑ Human-aware navigation
- ❑ Join persons and groups
- ❑ Optimal perception while navigating



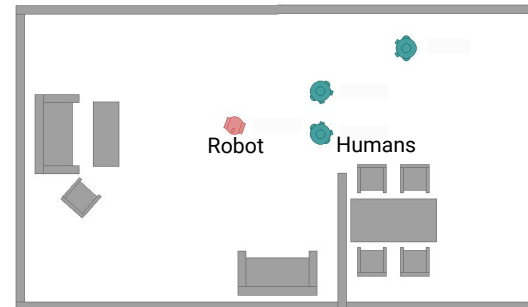
Reinforcement Learning for Social Robotics



Deep Reinforcement Learning



Training in simulation



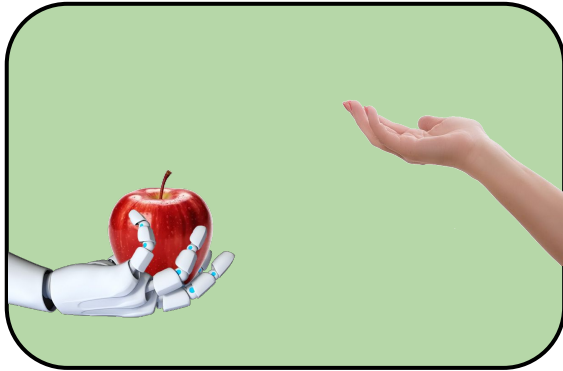
Sim2Real transfer



Challenge in Social Robotics

Social behaviors are highly user and context dependent

Task



User / Context



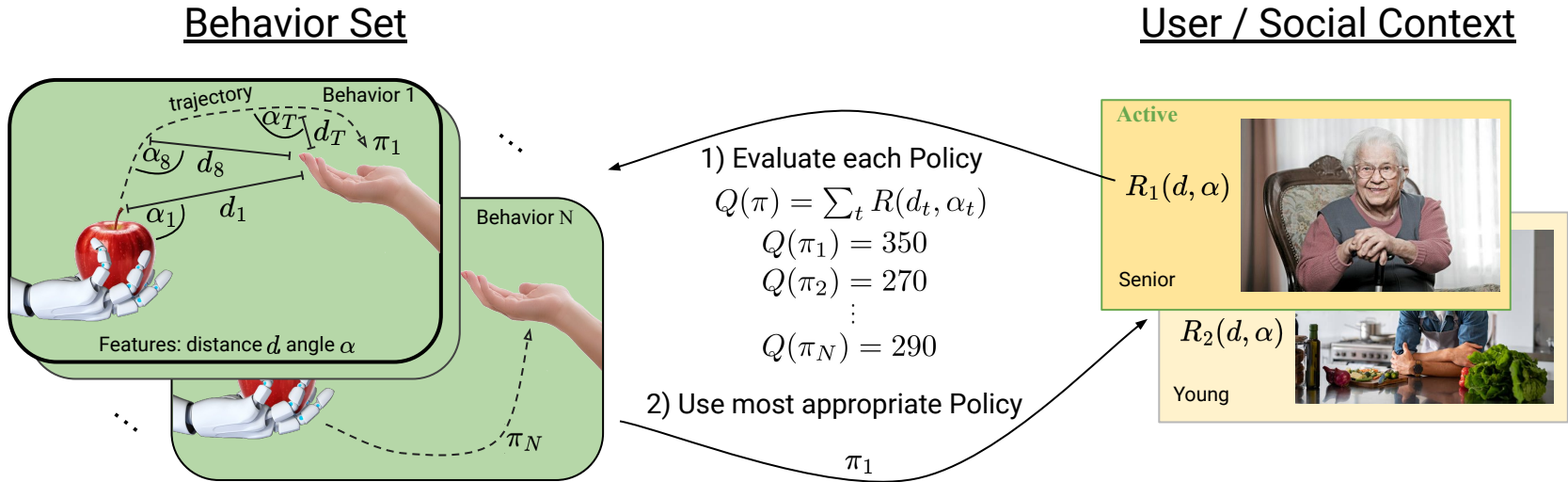
Senior



Younger adult

Research question: How to adapt behavior quickly?

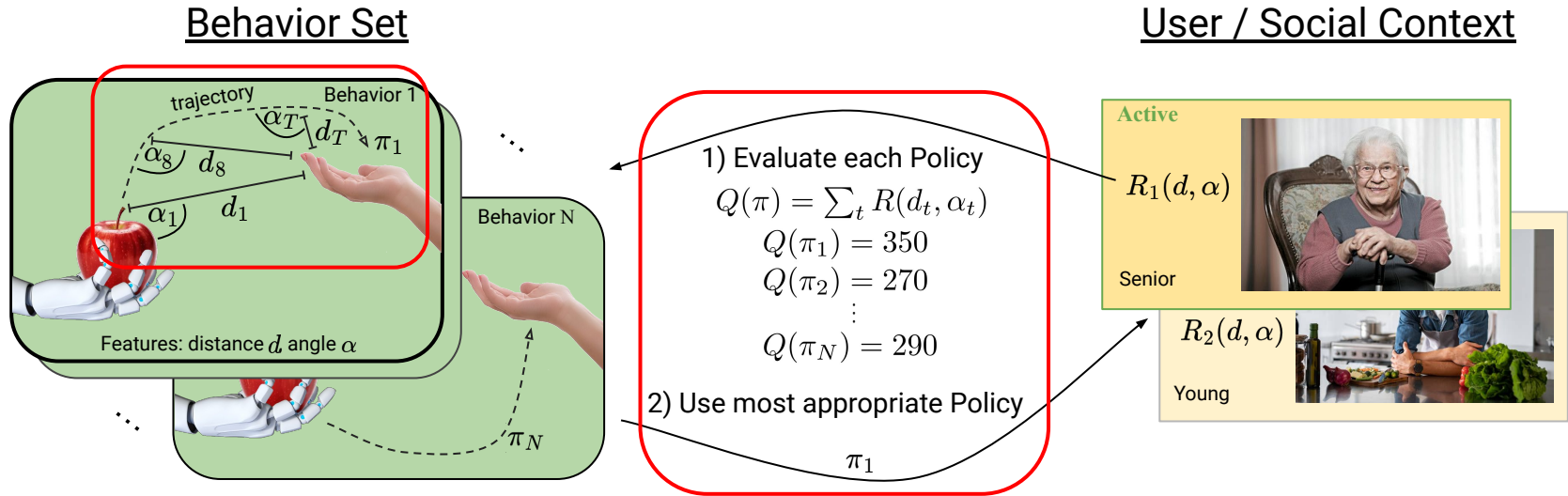
Proposal: Behavior Sets



Research Questions

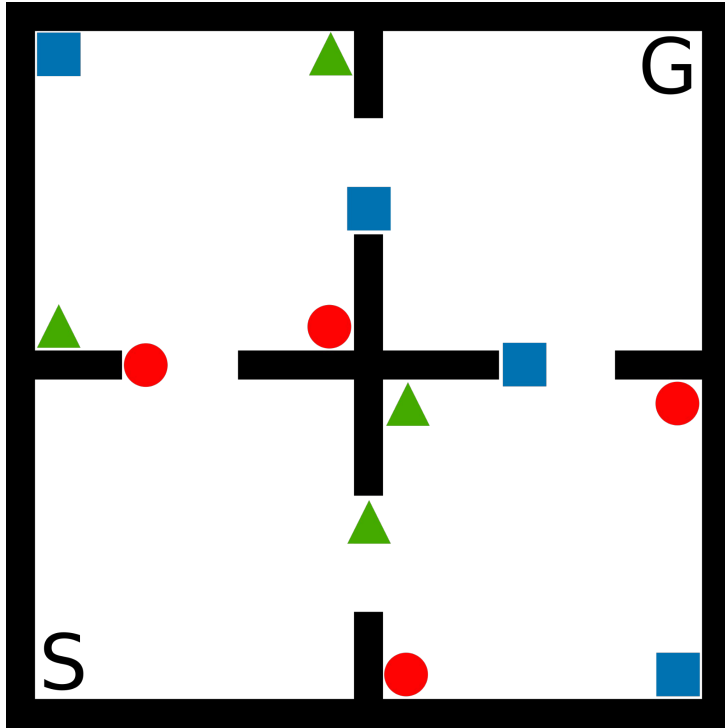
- 1) How to predict and evaluate behavior outcomes?
- 2) How to learn behavior sets?
- 3) How to model user preferences?

1) How to predict and evaluate behavior outcomes?



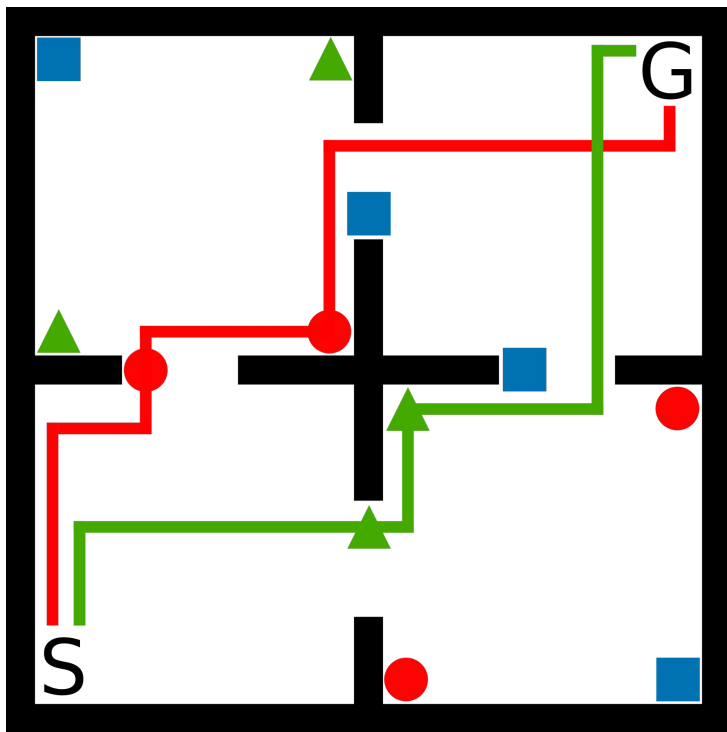
→ Successor Feature Representations

Simple Example

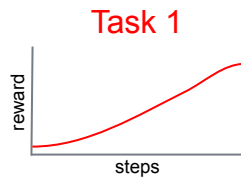


- ❑ 2D grid world
- ❑ Go from start (S) to goal (G)
- ❑ Get rewards for collecting objects
- ❑ Rewards differ between tasks

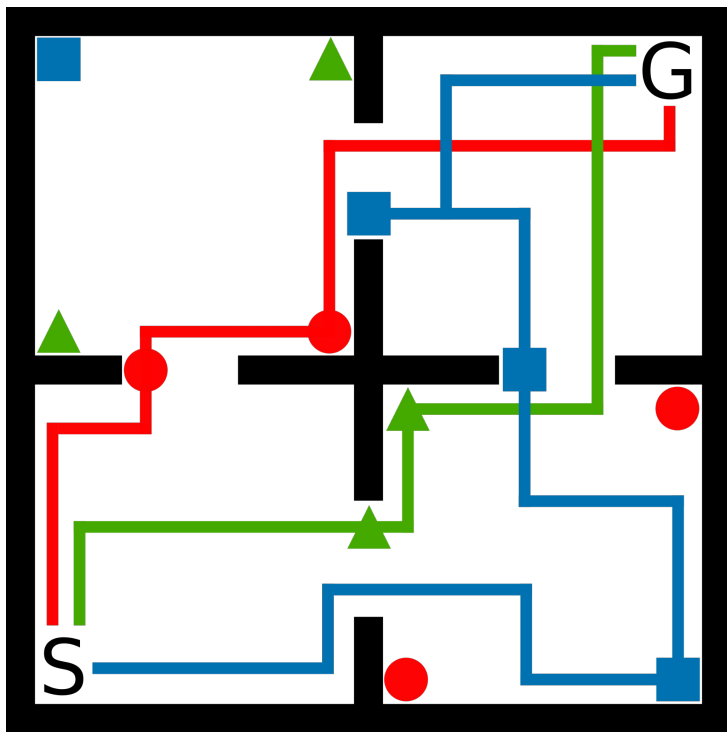
Simple Example



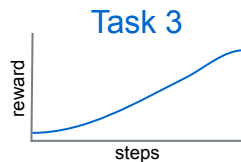
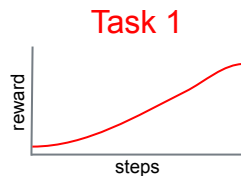
R	●	▲	■
Task 1	1.0	0.1	0.0
Task 2	0.1	1.0	0.0



Simple Example



R	●	▲	■
Task 1	1.0	0.1	0.0
Task 2	0.1	1.0	0.0
Task 3	0.2	0.0	1.0

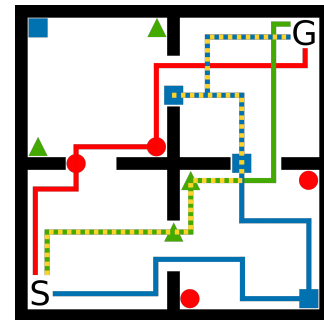


Successor Feature Representation (SFR)

Assumption: Low-dim features to encode rewards $r_t = R(\phi_t)$

Successor Feature Representation:

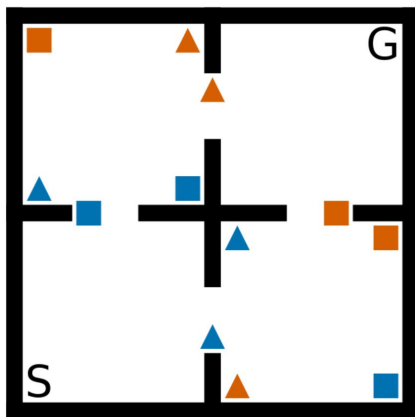
$$\begin{aligned}
 Q^\pi(s_t, a_t) &= \sum_{k=0}^{\infty} \gamma^k \mathbb{E}_{p(\phi_{t+k}|s_t, a_t; \pi)} \{R(\phi_{t+k})\} \\
 &= \sum_{k=0}^{\infty} \gamma^k \int_{\Phi} p(\phi_{t+k} = \phi | s_t, a_t; \pi) R(\phi) d\phi \\
 &= \int_{\Phi} R(\phi) \sum_{k=0}^{\infty} \gamma^k p(\phi_{t+k} = \phi | s_t, a_t; \pi) d\phi \\
 &= \int_{\Phi} R(\phi) \xi^\pi(s_t, a_t, \phi) d\phi = \sum_{\phi \in \Phi} R(\phi) \xi^\pi(s_t, a_t, \phi)
 \end{aligned}$$



$$\phi_t = \begin{pmatrix} \text{red circle} & ? \\ \text{green triangle} & ? \\ \text{blue square} & ? \end{pmatrix}$$

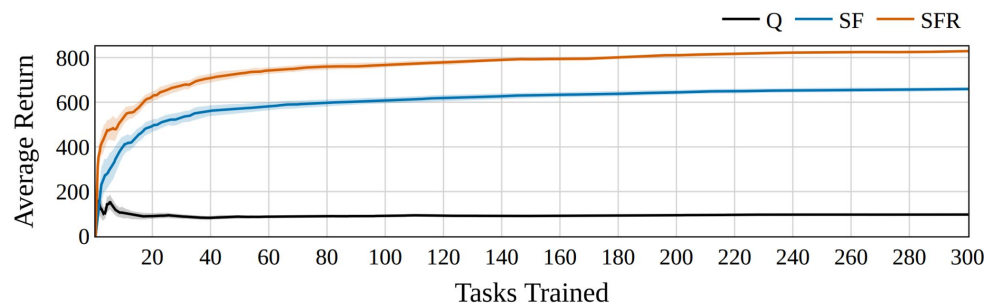
R	●	▲	■
Task 1	1.0	0.1	0.0
Task 2	0.1	1.0	0.0
Task 3	0.2	0.0	1.0
Task x	0.0	0.9	0.6

SFR Results - Discrete Features

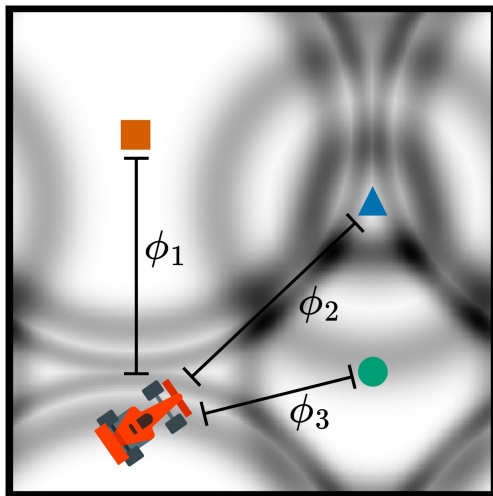


$$\phi_t = \begin{array}{l} \left. \begin{array}{l} \blacktriangle ? \\ \blacksquare ? \end{array} \right\} \text{Form} \\ \left. \begin{array}{l} \bullet ? \\ \bullet ? \end{array} \right\} \text{Color} \\ \mathbf{G} ? \quad \text{Goal} \end{array}$$

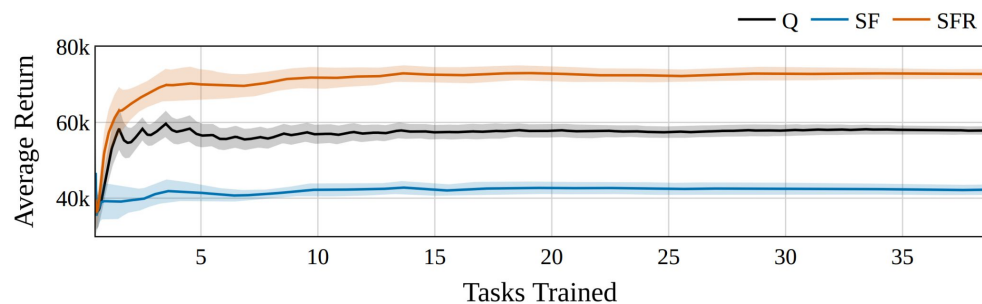
General Reward Functions: $r_t = R(\phi_t)$



SFR Results - Continuous Features

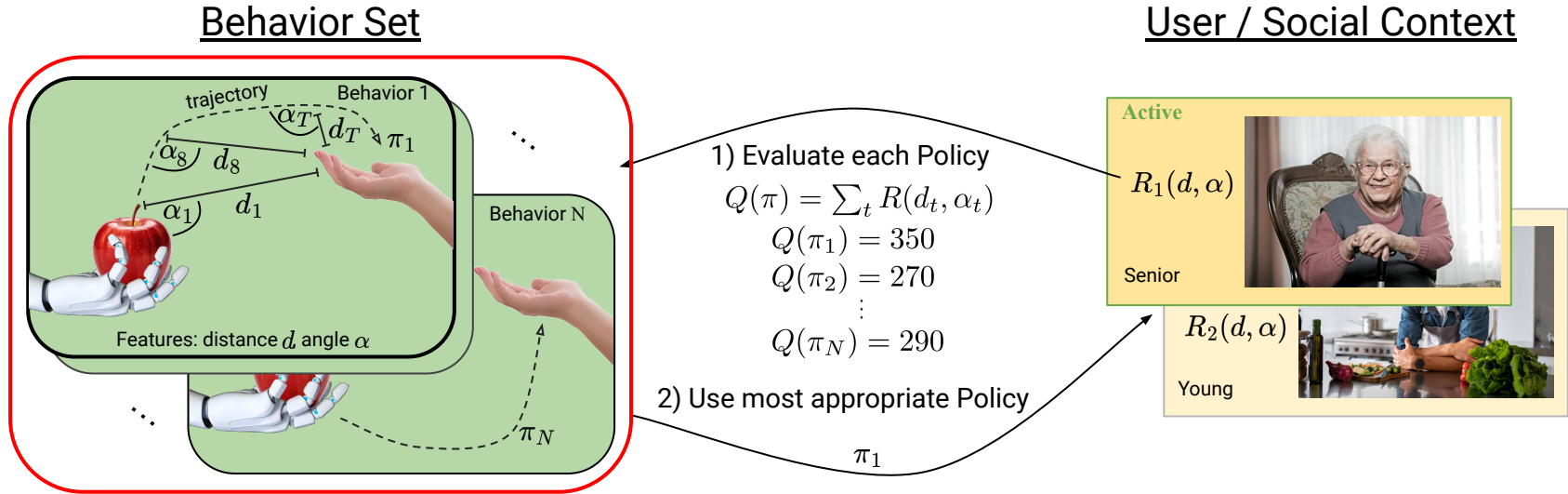


General Reward Functions: $r_t = R(\phi_t)$



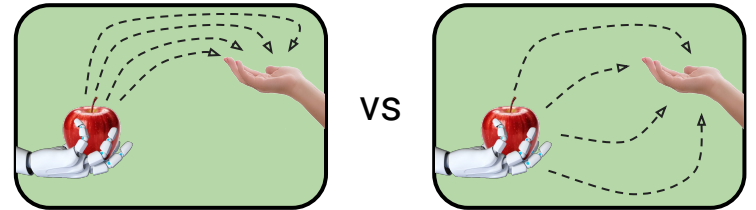
$$R(\phi) = \sum_{i=1}^3 r_i(\phi_i) \quad \text{with} \quad r_i(\phi_i) = \frac{1}{3} \max \left\{ \exp \left(-\frac{(\phi_i - \mu_i)^2}{\sigma_j} \right) \right\}_{j=1}^m$$

2) How to learn behavior sets?



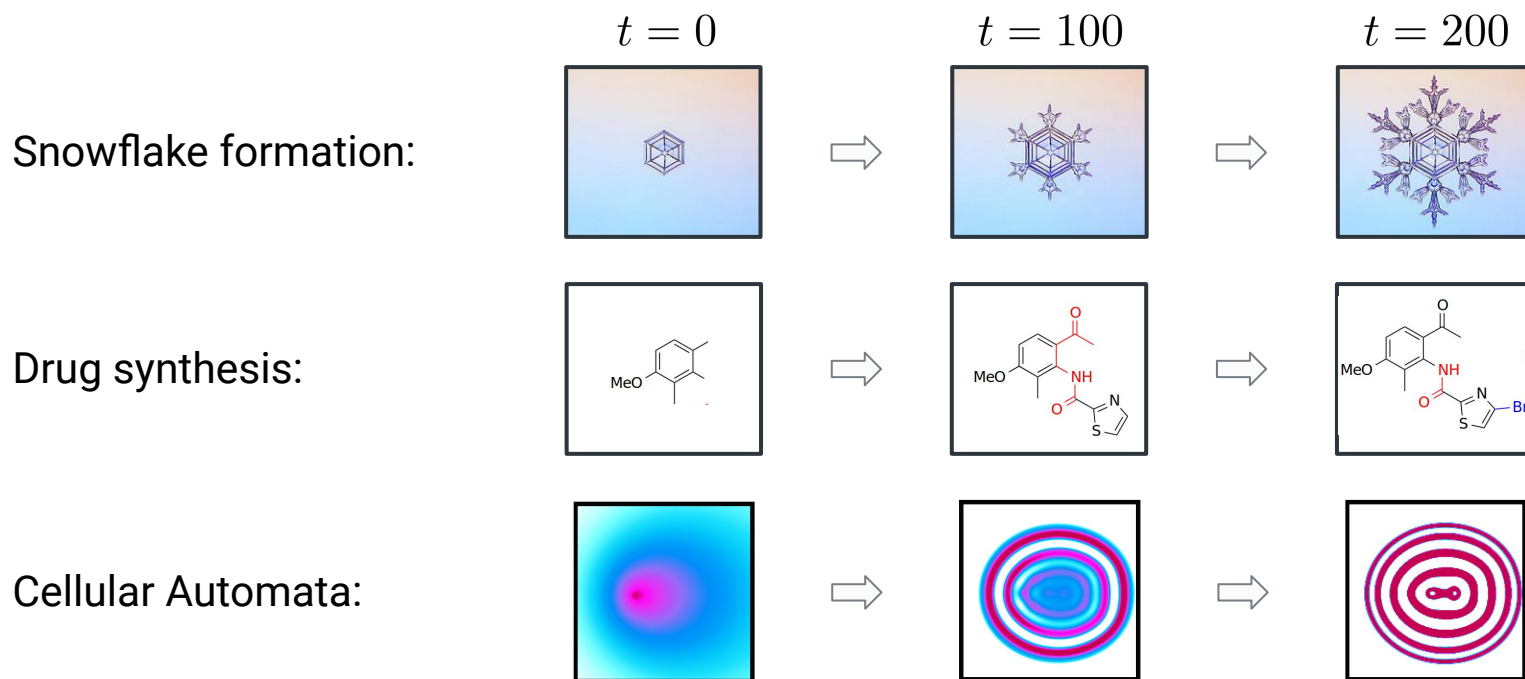
Goal: Learn a high diversity of behaviors

→ Diversity Exploration

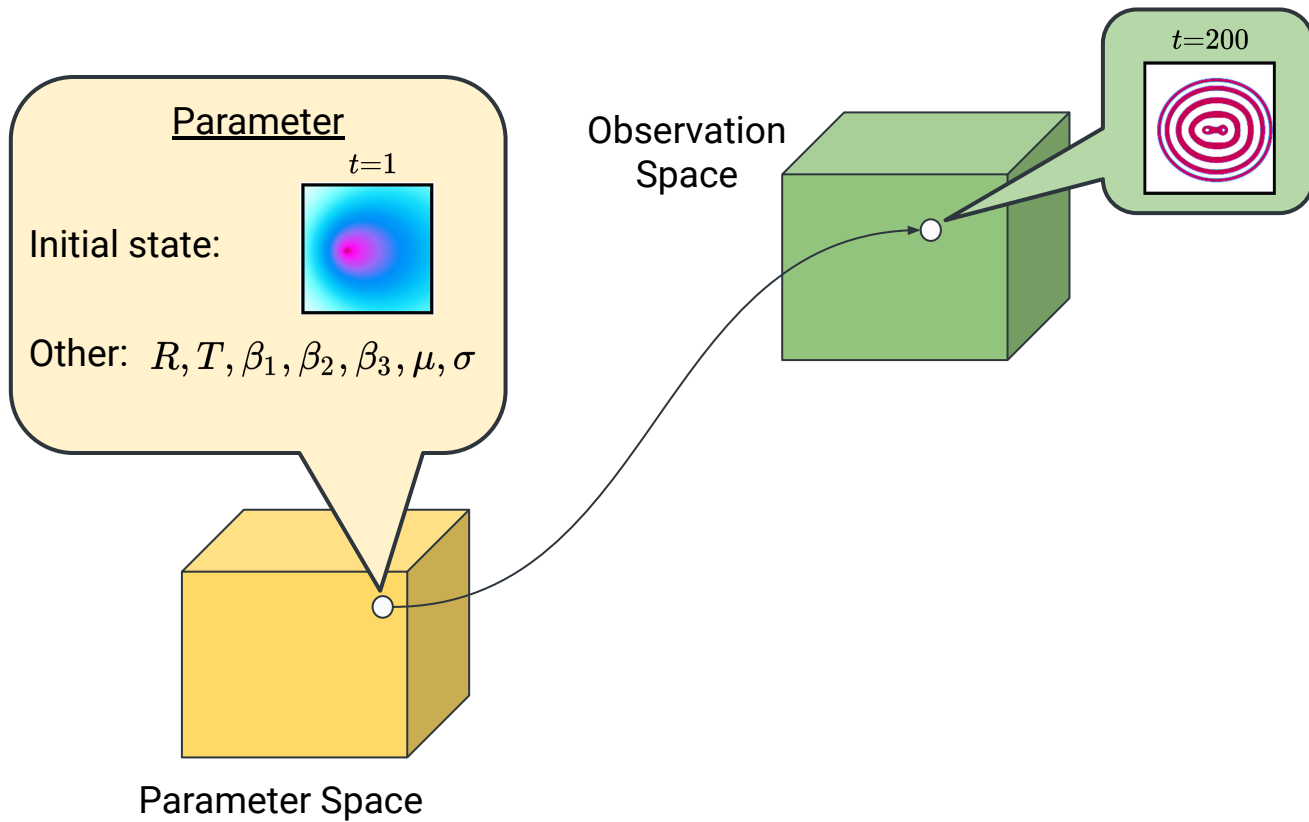


Diversity Exploration

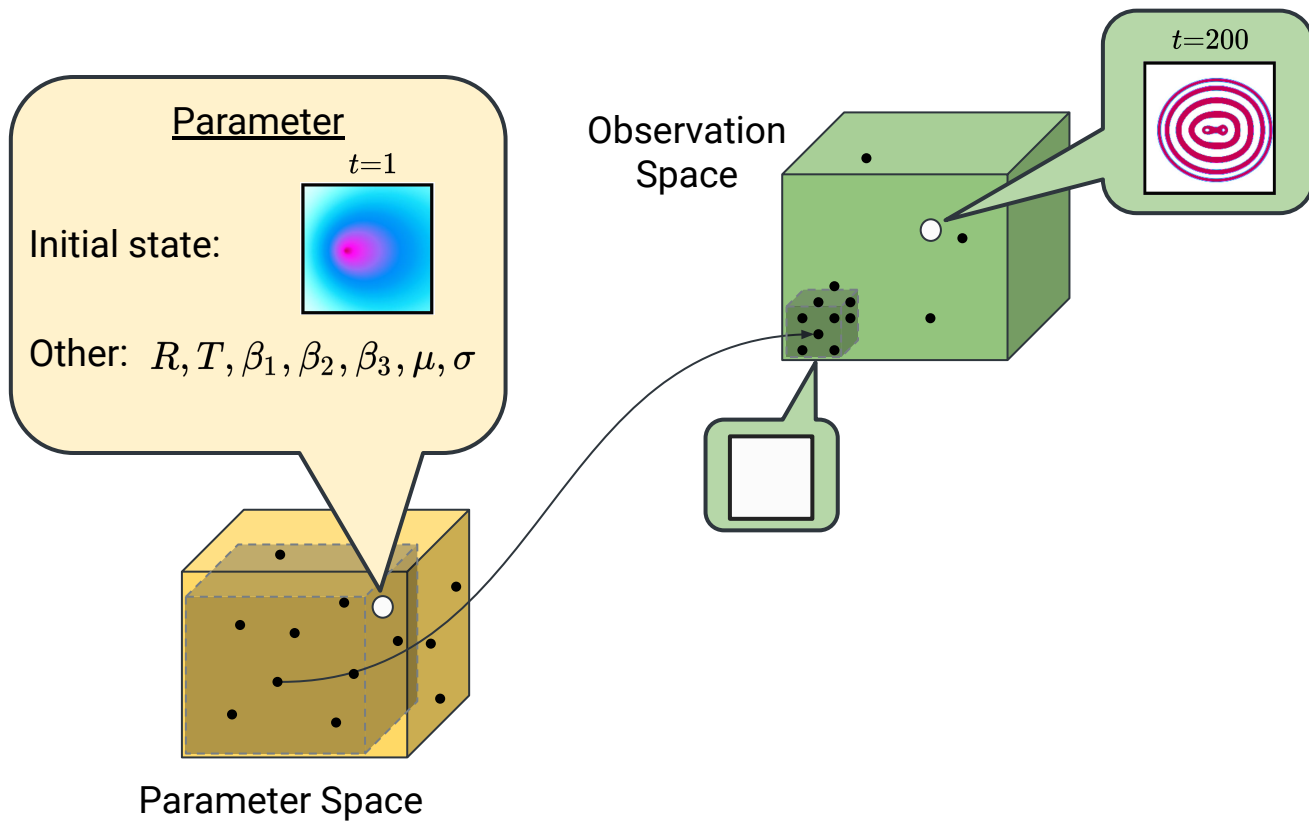
How to explore a high diversity of outcomes for high-dimensional dynamic (black box) systems?



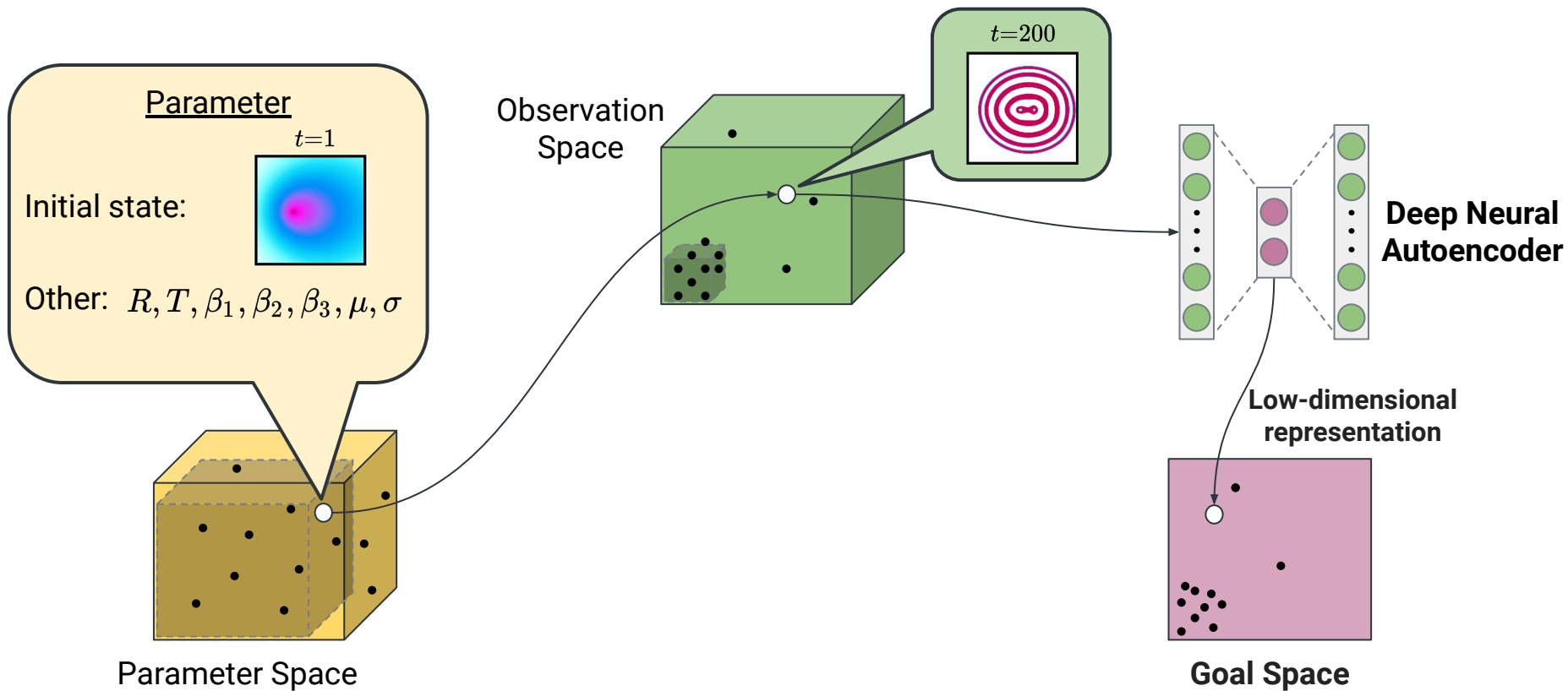
Diversity Exploration



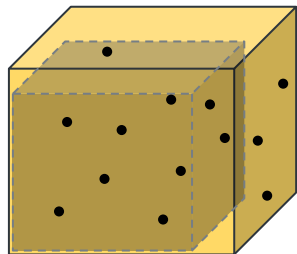
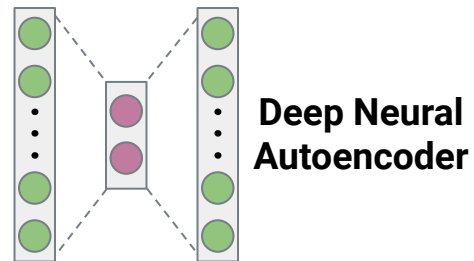
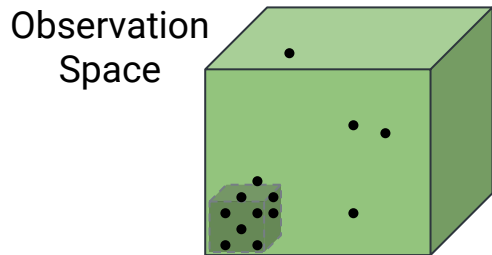
Diversity Exploration



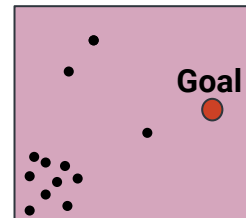
Diversity Exploration



Diversity Exploration

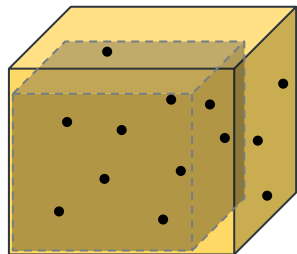
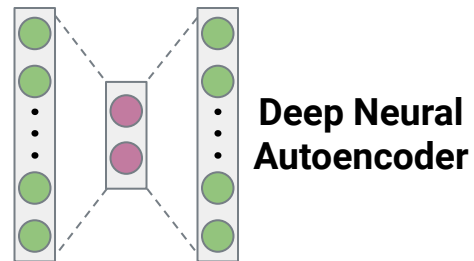
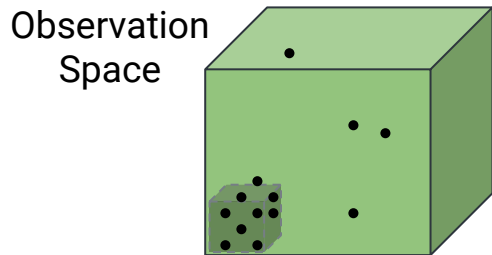


Parameter Space

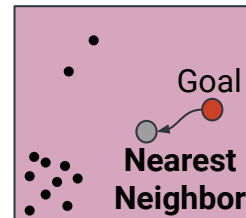


Goal Space

Diversity Exploration

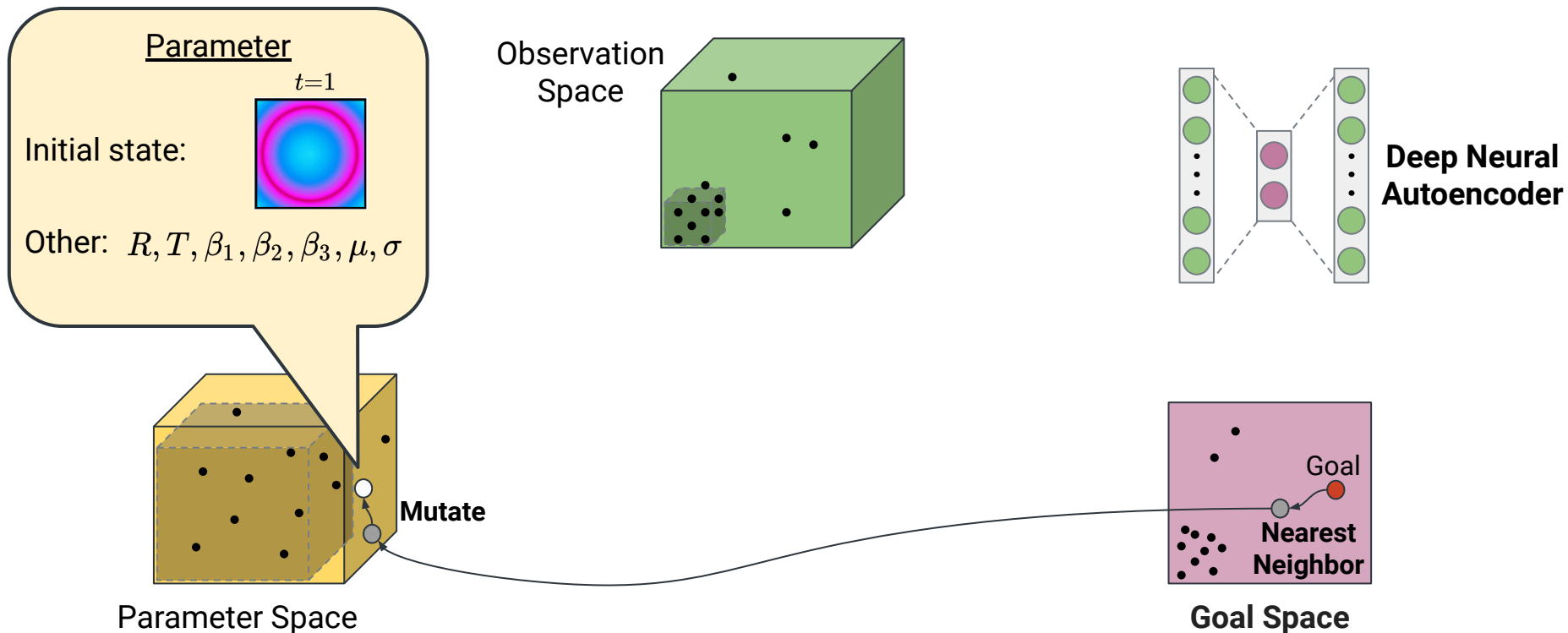


Parameter Space

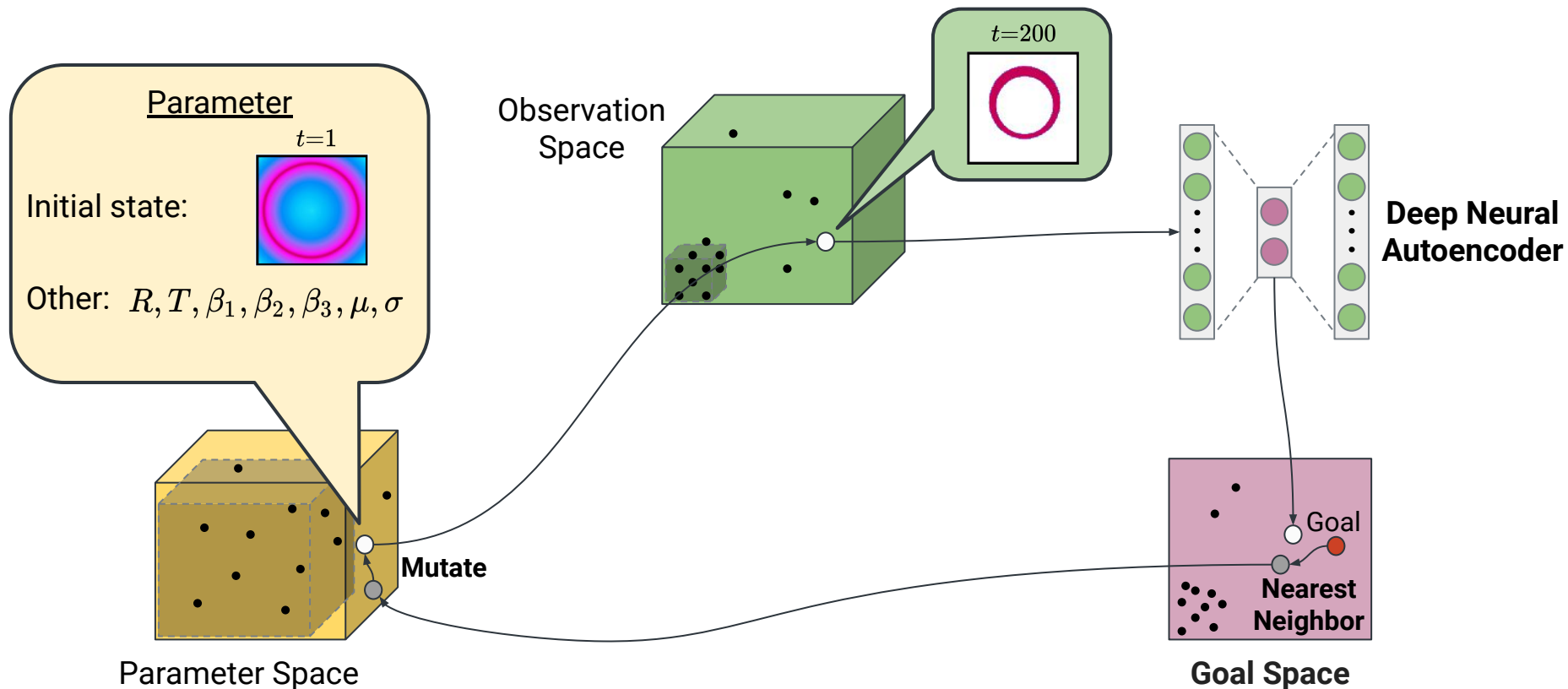


Goal Space

Diversity Exploration

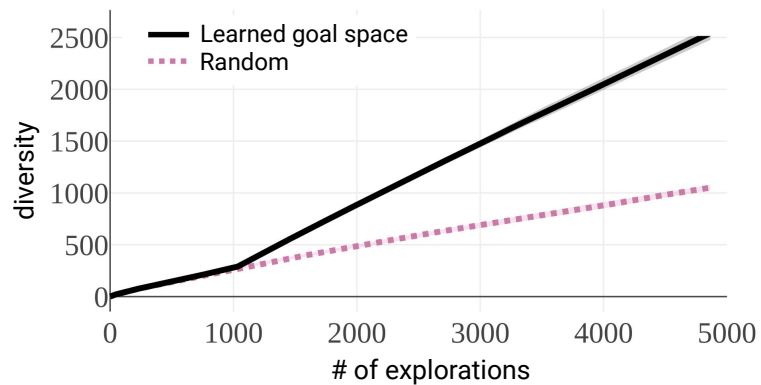
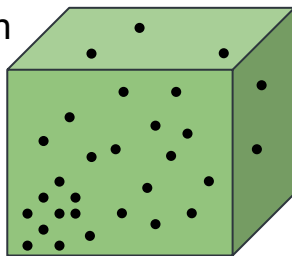


Diversity Exploration



Diversity Exploration

Observation Space



2) How to learn behavior sets?

Goal

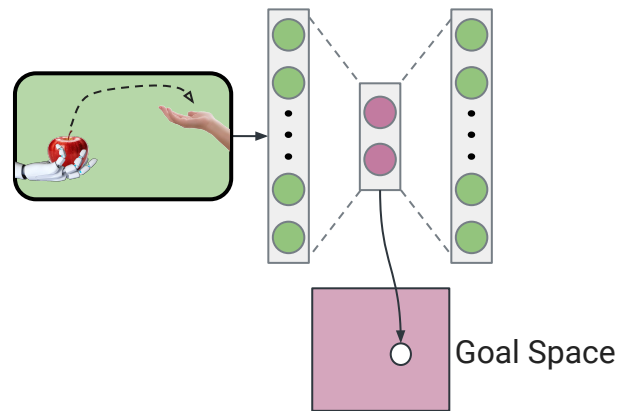
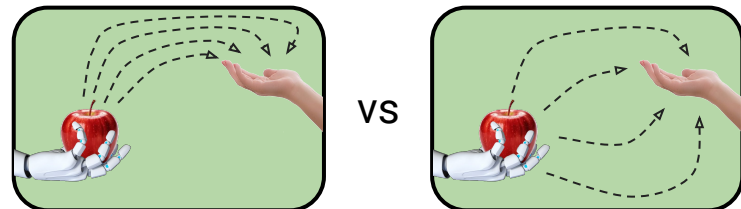
Learning a high diversity of behaviors

State of the art

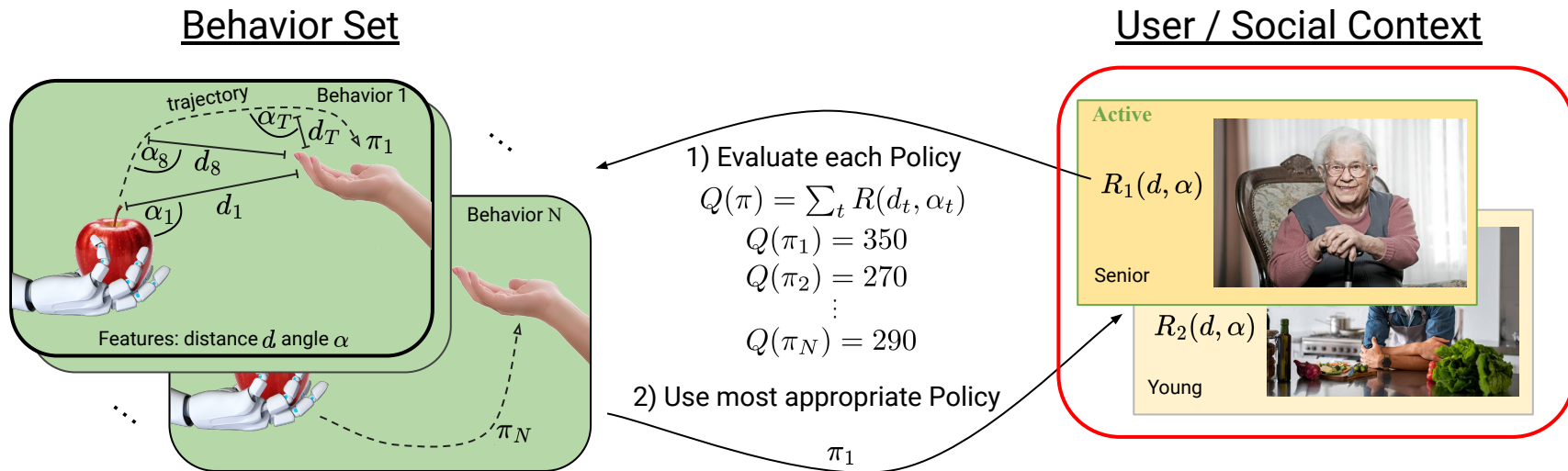
Diversity exploration focuses on final outcomes of systems and not dynamics (trajectories)

Proposed approach

Investigate unsupervised models that represent dynamics such as Dynamical VAEs



3) How to model user preferences?



3) How to model user preferences?

Goal

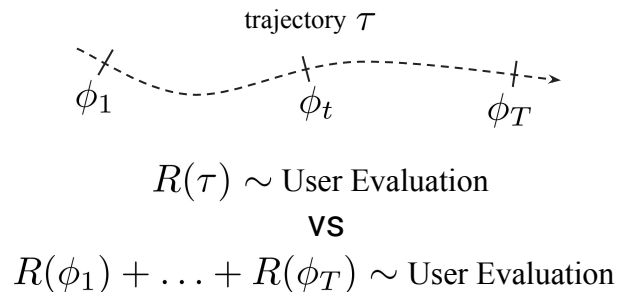
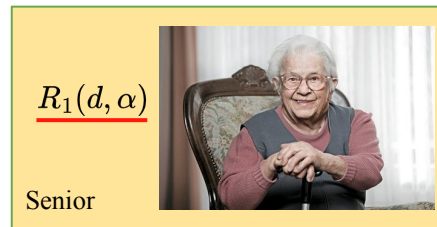
Correct and fast modeling of user preferences

State of the art

Model preferences over full trajectories

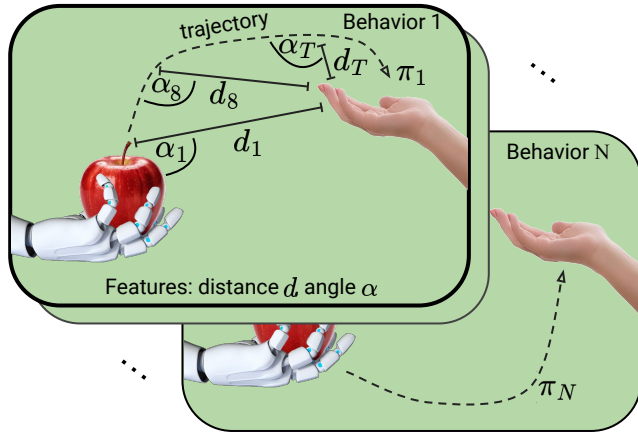
Proposed approaches

- ❑ Model preferences over features instead of trajectories to improve data efficiency
- ❑ Directly associating time dependent indirect preference signals (e.g. attention) with features

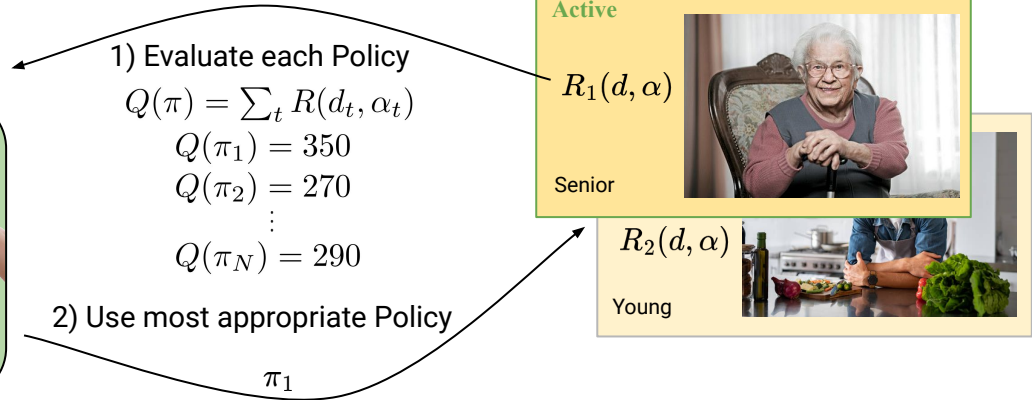


Thank You

Behavior Set



User / Social Context



More Information

❑ Myself:

www.scirei.net

❑ SFR:

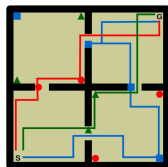
gitlab.inria.fr/robotlearn/sfr_learning

❑ Diversity Exploration:

automated-discovery.github.io

Background: Successor Features

- Assumption: Rewards are composed of features and weights $r_t = \phi_t^\top \mathbf{w}$



$$\phi_t = \begin{pmatrix} \text{red circle} & ? \\ \text{green triangle} & ? \\ \text{blue square} & ? \end{pmatrix}$$

R	●	▲	■	
Task 1	1.0	0.1	0.0	= \mathbf{w}_1
Task 2	0.1	1.0	0.0	= \mathbf{w}_2
Task 3	0.2	0.0	1.0	= \mathbf{w}_3
Task x	0.0	0.9	0.6	= \mathbf{w}_x

- Successor Features (SF) disassociate dynamics and rewards in Q-function

$$\begin{aligned} Q^\pi(s_t, a_t) &= \mathbb{E} [r_t + \gamma^1 r_{t+1} + \gamma^2 r_{t+2} + \dots] \\ &= \mathbb{E} [\phi_t^\top \mathbf{w} + \gamma^1 \phi_{t+1}^\top \mathbf{w} + \gamma^2 \phi_{t+2}^\top \mathbf{w} + \dots] \\ &= \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k \phi_{t+k} \right]^\top \mathbf{w} \\ &\equiv \boldsymbol{\psi}^\pi(s_t, a_t)^\top \mathbf{w} \end{aligned}$$

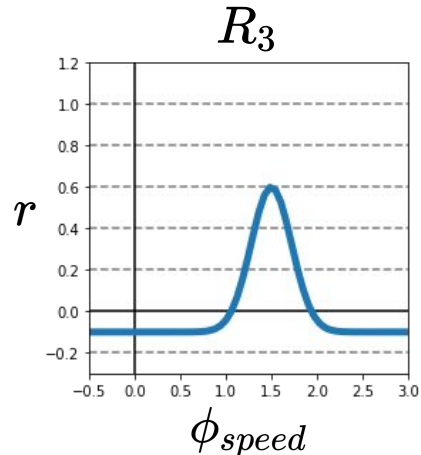
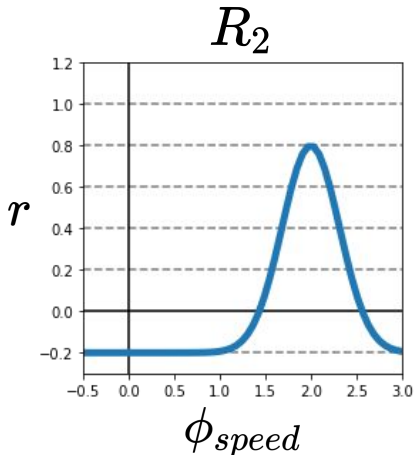
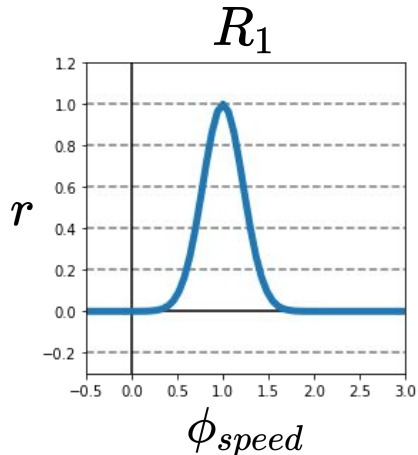
$$r_t \equiv R(s_t, a_t, s_{t+1}) \in \mathbb{R}$$

$$\phi_t \equiv \phi(s_t, a_t, s_{t+1}) \in \mathbb{R}^n, \quad \mathbf{w} \in \mathbb{R}^n$$

Background: Successor Features

- Problem: Assumption of linear relation of rewards $r_t = \phi_t^\top \mathbf{w}$

$$r(\phi_{speed}) = w \exp\left(-\frac{(\phi_{speed} - \mu)^2}{\sigma}\right) - d$$



Proposal: Successor Feature Representations

- Goal: Allow general reward functions $r_t = \phi_t^\top \mathbf{w} \longrightarrow r_t = R(\phi_t)$
- Successor Feature Representation (SFR):

$$\begin{aligned} Q^\pi(s_t, a_t) &= \sum_{k=0}^{\infty} \gamma^k \mathbb{E}_{p(\phi_{t+k}|s_t, a_t; \pi)} \{R(\phi_{t+k})\} \\ &= \sum_{k=0}^{\infty} \gamma^k \int_{\Phi} p(\phi_{t+k} = \phi | s_t, a_t; \pi) R(\phi) d\phi \\ &= \int_{\Phi} R(\phi) \sum_{k=0}^{\infty} \gamma^k p(\phi_{t+k} = \phi | s_t, a_t; \pi) d\phi \\ &= \int_{\Phi} R(\phi) \xi^\pi(s_t, a_t, \phi) d\phi \end{aligned}$$