

Temporal Logic Imitation

Yanwei Wang*, Nadia Figueroa†, Shen Li*, Ankit Shah‡, Julie Shah*

*MIT CSAIL, †University of Pennsylvania, ‡Brown University

Abstract—Learning from demonstration (LfD) methods have shown promise for solving multi-step tasks; however, these approaches do not guarantee successful reproduction of the task given perturbations. In this work, we identify the roots of such a challenge as the failure of the learned continuous policy to satisfy the discrete plan implicit in the demonstration. By utilizing modes (rather than subgoals) as the discrete abstraction and motion policies with both mode invariance and goal reachability properties, we show our learned continuous policy can simulate any given discrete plan. Consequently, the imitator is robust to both task- and motion-level perturbations and guaranteed to achieve task success. Project page: <https://yanweiw.github.io/tli/>

I. INTRODUCTION

In prior work, learning from demonstration (LfD) [1] has successfully enabled robots to accomplish multi-step tasks by segmenting demonstrations into subtasks/subgoals [2], phases [3], keyframes [4], and primitives [5]. Most of these abstractions assume reaching subgoals sequentially will deliver the desired outcomes. However, successful imitation of many manipulation tasks with spatial/temporal constraints cannot be reduced to imitation at the motion level unless the learned motion policy also satisfies these constraints. For example, transferring a spoonful of soup without restricting the orientation of the spoon can fail due to spilling even when the spoon reaches the target successfully.

We show that successful goal-reaching does not imply successful task execution in Fig. 1 with a 2D task, where task success depends on whether a continuous trajectory simulates a discrete plan of transitioning through the white, yellow, pink and green regions consecutively. Human demonstrations, shown in Fig. 1(a), are employed to learn a dynamical system (DS) policy [6] depicted by the streamlines in Fig. 1(b). Of all sampled trajectories, only the blue ones succeed in the task. The red ones fail as they simulate at least one discrete transition not physically realizable (e.g., white \Rightarrow pink). The issue is not mitigated by further segmenting the demonstrations into three subgoals and learning a DS for each subgoal, as seen in Fig. 1(c-f). While one can frame this problem as covariate shift and solve it by asking humans for more demonstrations [7], we frame it as the mismatch between a learned continuous policy and a discrete task plan and solve it by asking humans for a task specification. Specifically, the core challenges illustrated by this example are two-fold: 1) Subgoals only impose point constraints that are insufficient to represent the boundary of a discrete abstraction. 2) The continuous policy may deviate from a demonstrated discrete plan and in such cases cannot replan to ensure all discrete transitions are valid. Instead, our approach employs “modes” as the discrete abstraction. We define a *mode* as a set of robot and environment configurations that share the same

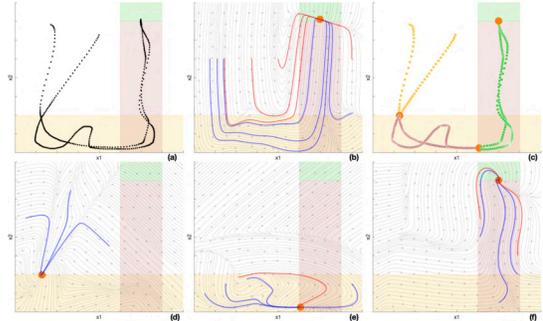


Fig. 1. A mode abstraction of a soup-scooping task. x_1 denotes spoon orientation and x_2 represents spoon distance from the soup. (a) Task: White region (all spoon configurations without soup on it) \Rightarrow yellow region (spoon in contact with soup) \Rightarrow pink region (spoon holding soup) \Rightarrow green region (soup at target). The black curves denote two successful demonstrations. (b) Learning a dynamical system (DS) policy [6] over unsegmented data can result in successful task replay (blue trajectories), but lacks a guarantee due to invalid transitions (red trajectories). (c) Trajectories segmented into three colored regions (modes), with orange circles denoting attractors. (d-f) Learning individual DS over segmented trajectories still results in *invariance failures*, i.e., red trajectories traveling outside of modes.

sensor reading [8]. Additionally, we use a task automaton as a receding horizon controller that replans when a perturbation causes the system to travel outside a mode boundary. For example, detecting a transition of yellow \Rightarrow white instead of the desired yellow \Rightarrow pink will result in a new plan: white \Rightarrow yellow \Rightarrow pink \Rightarrow green.

In this work, we assume the task automaton is given in the form of a Linear Temporal Logic formula. We denote the challenge of learning a continuous policy that can realize any discrete plan of valid mode transitions specified by the automaton as *temporal logic imitation* (TLI). In contrast to temporal logic planning (TLP) [9], where the workspace is typically partitioned into connected convex cells with known boundaries, TLI does not assume knowing mode boundaries. Consequently, the learned policy might prematurely exit a mode if the robot is perturbed to out-of-distribution states. To guarantee any discrete plan is feasible during execution at the continuous level, we show a learned policy with a global stability property can be refined to satisfy the bisimulation criteria [9] through human perturbations. By studying TLP in the setting of LfD, we are able to address covariate shift through learning motion policies that always obey discrete plans without additional online data collection.

II. TEMPORAL LOGIC IMITATION FORMULATION

A. Robot Model and Sensor Model

We use a first-order dynamical system $\dot{x} = f(x)$ to represent the desired evolution of a robot end-effector, where $x =$

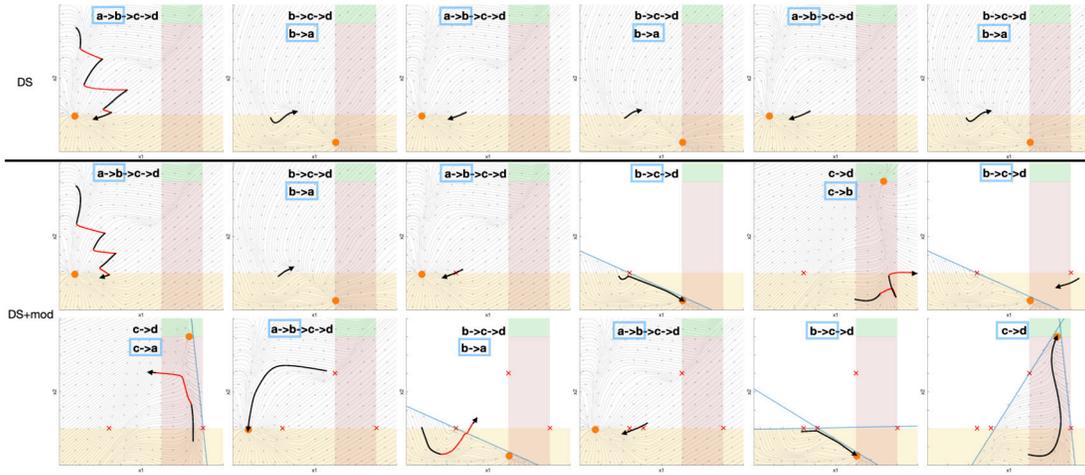


Fig. 2. Rollouts of multi-step scooping task under perturbations. The orange circle indicates the attractor for the current mode. We show the mode sequence planned by the automaton at the top of each sub-figure with the blue bounding box indicating the current mode transition actually detected. The trajectory in black is the rollout following the original policy, and the trajectory in red is driven by perturbations. The first row shows DS policies sequenced by an automaton but without boundary, estimation can lead to looping. The second and third rows show modulation can prevent looping and eventually allow the system to reach the goal mode despite repeated perturbations. Please check the video on the project page.

$[x_1, \dots, x_n]^T \in \mathbb{R}^n$ describes an n -dimensional continuous robot state. Let discrete sensor state $\alpha = [\alpha_1, \dots, \alpha_m]^T \in \{0, 1\}^m$ be an m -dimensional sensor variable. We define a system state as a tuple $s = (x, \alpha) \in \mathbb{R}^n \times \{0, 1\}^m$ and its corresponding mode $\sigma \in \Sigma$ as $\sigma = \mathcal{L}(\alpha)$. Overloading the notation, we use σ to represent the set of all system states in the same mode, i.e., $\sigma_i = \{s = (x, \alpha) \mid \mathcal{L}(\alpha) = \sigma_i\}$. In contrast, $\delta_i = \{x \mid s = (x, \alpha) \in \sigma_i\}$ represent the corresponding set of robot states.

B. Demonstrations and perturbations

Demonstrations are in the form $\{\{x^{t,k}, \dot{x}^{t,k}, \alpha^{t,k}\}_{t=1}^{T_k}\}_{k=1}^K$, where $x^{t,k}, \dot{x}^{t,k}, \alpha^{t,k}$ are robot state, velocity, and sensor state at time t in demonstration k . T_k is the length of each k -th trajectory. Given l number of unique sensor states, we can segment the K demonstrations into l clusters of modes. We learn a policy for each mode, which we stress test with either (1) motion-level perturbations that displace the continuous motion within the same mode, or (2) task-level perturbations that drive the system outside of the current mode.

C. Problem Statement

Given (1) a task automaton ϕ specifying valid mode transitions to achieve a task, and (2) successful demonstrations $\{\{x^{t,k}, \dot{x}^{t,k}, \alpha^{t,k}\}_{t=1}^{T_k}\}_{k=1}^K$, we want to imitate a continuous policy that can realize any discrete mode sequence planned by the automaton despite arbitrary perturbations.

D. Bisimulation between Discrete Plan and Continuous Policy

To realize any discrete plan, every mode's associated continuous policy must satisfy the bisimulation conditions: [9].

Condition 1 (Invariance). *Every continuous motion starting in a mode must remain within the same mode while following the current mode's policy; i.e., $\forall i \forall t (s^0 \in \sigma_i \rightarrow s^t \in \sigma_i)$*

Condition 2 (Reachability). *Every continuous motion starting in a mode must reach the next mode in the demonstration*

while following the current mode's policy; i.e., $\forall i \exists T (s^0 \in \sigma_i \rightarrow s^T \in \sigma_j)$

III. EXPERIMENTS

We now show empirically having a reactive discrete plan is insufficient to guarantee task success without mode invariance for tasks with multiple modes. Consider the task introduced in Fig. 1: scooping and transporting soup. Formally, we define four modes: (a) starting empty spoon, (b) sensing the spoon is in contact with the soup, (c) sensing the spoon has soup on it, and (d) sensing the spoon has arrived at a target location. During successful demonstrations, we observe the following discrete transitions a (reaching) $\Rightarrow b$ (scooping) $\Rightarrow c$ (transporting) $\Rightarrow d$ (done). Invariance of mode b enforces contact during scooping and invariance of mode c constrains the spoon orientation to avoid spilling. One might assume having a task automaton is sufficient to guarantee task success without modulation, as it only needs to replan a finite number of times assuming a finite number of perturbations; however, not enforcing mode invariance can lead to looping at the discrete level, and ultimately renders the goal unreachable, as depicted in the top row of Fig. 2. In contrast, looping is prevented when modulation is enabled, as the system experiences each invariance failure only once.

IV. CONCLUSION

In this paper, we introduce *temporal logic imitation* as the problem of learning plan-satisficing motion policies. We identify one challenge of applying LfD methods to multi-step tasks as being that the learned controllers do not necessarily satisfy the bisimulation criteria. To address this issue, we propose a DS-based approach that can iteratively estimate mode boundaries to ensure invariance and reachability. Combining the task-level reactivity of a task automaton and the motion-level reactivity of DS, we arrive at an imitation learning system that can robustly perform a multi-step scooping task under arbitrary perturbations given only a few demonstrations.

REFERENCES

- [1] H. Ravichandar, A. S. Polydoros, S. Chernova, and A. Billard, “Recent advances in robot learning from demonstration,” *Annual Review of Control, Robotics, and Autonomous Systems*, vol. 3, pp. 297–330, 2020.
- [2] A. Mandlekar, D. Xu, R. Martín-Martín, S. Savarese, and L. Fei-Fei, “Learning to generalize across long-horizon tasks from human demonstrations,” *arXiv preprint arXiv:2003.06085*, 2020.
- [3] O. Kroemer, C. Daniel, G. Neumann, H. Van Hoof, and J. Peters, “Towards learning hierarchical skills for multi-phase manipulation tasks,” in *2015 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2015, pp. 1503–1510.
- [4] B. Akgun, M. Cakmak, J. W. Yoo, and A. L. Thomaz, “Trajectories and keyframes for kinesthetic teaching: A human-robot interaction perspective,” in *Proceedings of the seventh annual ACM/IEEE international conference on Human-Robot Interaction*, 2012, pp. 391–398.
- [5] S. Niekum, S. Chitta, A. G. Barto, B. Marthi, and S. Osentoski, “Incremental semantically grounded learning from demonstration.” in *Robotics: Science and Systems*, vol. 9. Berlin, Germany, 2013, pp. 10–15 607.
- [6] N. Figueroa and A. Billard, “A physically-consistent bayesian non-parametric mixture model for dynamical system learning.” in *CoRL*, 2018, pp. 927–946.
- [7] S. Ross, G. Gordon, and D. Bagnell, “A reduction of imitation learning and structured prediction to no-regret online learning,” in *Proceedings of the fourteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings, 2011, pp. 627–635.
- [8] C. R. Garrett, R. Chitnis, R. Holladay, B. Kim, T. Silver, L. P. Kaelbling, and T. Lozano-Pérez, “Integrated task and motion planning,” *Annual review of control, robotics, and autonomous systems*, vol. 4, pp. 265–293, 2021.
- [9] H. Kress-Gazit, M. Lahijanian, and V. Raman, “Synthesis for robots: Guarantees and feedback for robot behavior,” *Annual Review of Control, Robotics, and Autonomous Systems*, vol. 1, pp. 211–236, 2018.