# Multi-task Imitation Learning for Linear Dynamical Systems

Thomas T. Zhang\*1
Katie Kang\*2
Bruce D. Lee\*1
Claire Tomlin2
Sergey Levine2
Stephen Tu3
Nikolai Matni1,3

TTZ2@SEAS.UPENN.EDU
KATIEKANG@EECS.BERKELEY.EDU
BRUCELE@SEAS.UPENN.EDU
TOMLIN@EECS.BERKELEY.EDU
SVLEVINE@EECS.BERKELEY.EDU
STEPHENTU@GOOGLE.COM
NMATNI@SEAS.UPENN.EDU

Editors: N. Matni, M. Morari, G. J. Pappas

### **Abstract**

We study representation learning for efficient imitation learning over linear systems. In particular, we consider a setting where learning is split into two phases: (a) a pre-training step where a shared k-dimensional representation is learned from H source policies, and (b) a target policy fine-tuning step where the learned representation is used to parameterize the policy class. We find that the imitation gap over trajectories generated by the learned target policy is bounded by  $\tilde{O}\left(\frac{kn_x}{HN_{\rm shared}} + \frac{kn_u}{N_{\rm target}}\right)$ , where  $n_x > k$  is the state dimension,  $n_u$  is the input dimension,  $N_{\rm shared}$  denotes the total amount of data collected for each policy during representation learning, and  $N_{\rm target}$  is the amount of target task data. This result formalizes the intuition that aggregating data across related tasks to learn a representation can significantly improve the sample efficiency of learning a target task. The trends suggested by this bound are corroborated in simulation.

Keywords: Imitation learning, transfer learning, multi-task learning, representation learning

#### 1. Introduction

Imitation learning (IL), which learns control policies by imitating expert demonstrations, has demonstrated success across a variety of domains including self-driving cars (Codevilla et al., 2018) and robotics (Schaal, 1999). However, using IL to learn a robust behavior policy may require a large amount of training data (Ross et al., 2011), and expert demonstrations are often expensive to collect. One remedy for this problem is multi-task learning: using data from other tasks (source tasks) in addition to from the task of interest (target task) to jointly learn a policy. We study the application of multi-task learning to IL over linear systems, and demonstrate improved sample efficiency when learning a controller via representation learning.

Our results expand on prior work that studies multi-task representation learning for supervised learning (Du et al., 2020; Tripuraneni et al., 2021), addressing the new challenges that arise in the imitation learning setting. First, the data for IL is temporally dependent, as it is generated from a dynamical system x[t+1] = f(x[t], u[t], w[t]). In contrast, the supervised learning setting assumes that both the train and test data are independent and identically distributed (i.i.d.) from

<sup>&</sup>lt;sup>1</sup> University of Pennsylvania, PA, USA

<sup>&</sup>lt;sup>2</sup> University of California, Berkeley, CA, USA

<sup>&</sup>lt;sup>3</sup> Google Research, Brain Team, NY, USA

<sup>\*</sup> Authors contributed equally to this work

the same underlying distribution. Furthermore, we are interested in the performance of the learned controller in closed-loop rather than its error on expert-controlled trajectories. Hence, bounds on excess risk, which corresponds to the one-step prediction error of the learned controller under the expert distribution, are not immediately informative for the closed-loop performance. We instead focus our analysis on the tracking error between the learned and expert policies, which requires us to account for the distribution shift between the learned and expert controllers.

We address these challenges in the setting of IL for linear systems. The following statement captures the benefits of multi-task representation learning on sample complexity:

**Theorem 1.1** (main result, informal) Suppose that the source task controllers are sufficiently related to the target task controller. Then, the tracking error between the learned target controller and the corresponding expert is bounded with high probability by:

$$tracking\ error \lesssim rac{rep.\ dimension\ imes state\ dimension}{\#\ source\ task\ datapoints} + rac{rep.\ dimension\ imes\ input\ dimension}{\#\ target\ task\ datapoints}.$$

The first term in this bound corresponds to the error from learning a common representation, and the second term the error in fitting the remaining weights of the target task controller. The key upshot of this result is that the numerator of the second term (rep. dimension  $\times$  input dimension) is smaller than the number of parameters (input dimension  $\times$  state dimension) in the target controller. This demonstrates an improvement in sample complexity of multi-task IL over direct IL, where the error scales as  $\frac{\# parameters}{\# datapoints}$ . Furthermore, we note that the error in learning the representation decays along all axes of the data: # of tasks  $\times$  # of trajs  $\times$  traj length for source tasks, and # of trajs  $\times$  traj length for the target task. It is non-trivial to demonstrate that the error decays with the trajectory length, and doing so requires tools that handle causally dependent data in our analysis.

The remainder of the paper formulates the multi-task IL problem, and the assumptions required to prove Theorem 1.1. The main contributions may be summarized as follows:

- We provide novel interpretable notions of source task overlap with the target task (§2 and §3).
- We bound the imitation gap achieved by multi-task IL as in Theorem 1.1 (§3).
- We empirically show the efficacy of multi-task IL when the assumptions are satisfied (§4).

#### 1.1. Related Work

Multi-task imitation and reinforcement learning: Multi-task RL and IL methods seek to represent policies solving different tasks with shared parameters, enabling the transfer of knowledge across related tasks (Teh et al., 2017; Espeholt et al., 2018; Hessel et al., 2018; Singh et al., 2020; Deisenroth et al., 2014), and rapid test-time adaptation to new tasks (Finn et al., 2017; Rakelly et al., 2019; Duan et al., 2016; Yu et al., 2021; Yang and Nachum, 2021). There also exists a body of work which theoretically analyses the sample complexity of representation learning in multi-task RL and IL (Lu et al., 2021; Cheng et al., 2022; Lu et al., 2022; Xu et al., 2020; Maurer et al., 2015; Arora et al., 2020). While this line of work considers a more general MDP setting compared with the linear dynamical systems we consider, the specific results are often stated with incompatible assumptions (such as bounded states/cost functions and discrete action spaces), and/or do not reflect how system-theoretic properties such as closed-loop task stability affect the final rates.

Multi-task system identification and adaptive control: Recent work has also considered applications of multi-task learning where the dynamics change between tasks, and the goal is to perform adaptive control (Harrison et al., 2018; Richards et al., 2021, 2022; Shi et al., 2021; Muthirayan

et al., 2022) or dynamics forecasting (Wang et al., 2021). Multi-task system identification (Modi et al., 2021) and stabilization using data from related systems (Li et al., 2022) have also been considered. Our work instead studies the problem of learning to imitate different expert controllers while the system remains the same, and demonstrates bounds on the tracking error between the learned controller and its corresponding expert.

Sample complexity of multi-task learning: Numerous works have studied the sample efficiency gains of multi-task learning for regression and classification under various task similarity assumptions (Baxter, 1995; Crammer et al., 2008; Maurer et al., 2016; Tripuraneni et al., 2020; Chua et al., 2021). Most closely related to our results are Du et al. (2020) and Tripuraneni et al. (2021), both of which show multi-task representation learning sample complexity bounds in the linear regression setting in which the error from learning the representation decays with the total number of source training samples. Our work leverages these results to tackle the setting of linear imitation learning, which has the additional challenges of non-i.i.d. data and test time distribution shift.

#### 2. Problem Formulation

# 2.1. Multi-Task Imitation Learning

Imitation learning uses state/action pairs  $(x,u) \in \mathbb{R}^{n_x} \times \mathbb{R}^{n_u}$  of expert demonstrations to learn a controller  $\hat{\pi}: \mathbb{R}^{n_x} \to \mathbb{R}^{n_u}$ , by matching the learned controller actions to the expert actions. In particular, if  $\mathcal{D}$  is the training set of expert state/action pairs, then  $\hat{\pi} \in \operatorname{argmin}_{\pi} \sum_{(x,u) \in \mathcal{D}} \|\pi(x) - u\|^2$ .

We are interested in the problem of *multi-task* imitation learning, where we consider H+1 different expert controllers. We call the first H controllers *source controllers* and the  $(H+1)^{\rm st}$  controller the *target controller*. We assume that we have access to  $N_1$  trajectories for each source task, and  $N_2 \leq N_1$  trajectories for the target task. For simplicity, we assume all trajectories are of the same length T. In particular, for each source task  $h \in \{1,\ldots,H\}$ , our source data consists of  $\left\{\left(x_i^{(h)}[t],u_i^{(h)}[t]\right)\right\}_{t=0}^{T-1}\right\}_{i=1}^{N_1}$ , while our target data consists of  $\left\{\left(x_i^{(H+1)}[t],u_i^{(H+1)}[t]\right)\right\}_{t=0}^{T-1}\right\}_{i=1}^{N_2}$ . Our goal is to learn a controller which effectively imitates the target controller. However, because we only have access to a small number  $(N_2)$  of target expert trajectories, we leverage the  $HN_1$  expert trajectories from the source controllers to accelerate the learning of the target controller.

To do so, we break our training into two stages: a pre-training stage which learns from the combined source task data, and a target training stage which only learns from the target task data. In the pre-training stage, we extract a common, low dimensional representation for the source controllers, which is used later in the target training stage. More specifically, we learn a common, low dimensional representation mapping  $\hat{\phi}: \mathbb{R}^{n_x} \to \mathbb{R}^k$ , where  $k < n_x$  is the dimension of the representation, and linear predictors  $\hat{F}^{(h)} \in \mathbb{R}^{n_u \times k}$  unique to each task:

$$\hat{\phi}, \hat{F}^{(1)}, \dots, \hat{F}^{(H)} \in \underset{\phi, F^{(1)}, \dots, F^{(H)}}{\operatorname{argmin}} \sum_{h=1}^{H} \sum_{i=1}^{N_1} \sum_{t=0}^{T-1} \left\| F^{(h)} \phi(x_i^{(h)}[t]) - u_i^{(h)}[t] \right\|^2. \tag{1}$$

We do not address the details of solving the empirical risk minimization problem, and instead perform our analysis assuming (1) can be solved to optimality. Note however that Tripuraneni et al. (2021) demonstrate in the linear regression setting that a method-of-moments-based algorithm can efficiently find approximate empirical risk minimizers.

Once a common representation  $\phi$  is obtained, we move on to target task training. During target task training, we use the common representation mapping  $\hat{\phi}$  learned from the pre-training step to

map the states into the lower dimensional representation, and learn an additional linear predictor  $\hat{F}^{(H+1)}$  unique to the target task to model the target controller:

$$\hat{F}^{(H+1)} = \underset{F}{\operatorname{argmin}} \sum_{i=1}^{N_2} \sum_{t=0}^{T-1} \left\| F \hat{\phi} \left( x_i^{(H+1)}[t] \right) - u_i^{(H+1)}[t] \right\|^2. \tag{2}$$

Since the representation  $\hat{\phi}$  is fixed from pre-training, (2) is an ordinary least squares problem.

# 2.2. System and Data Assumptions

We focus our analysis on a linear systems setting, with state  $x[t] \in \mathbb{R}^{n_x}$ , input  $u[t] \in \mathbb{R}^{n_u}$ , and Gaussian process noise  $w[t] \in \mathbb{R}^{n_x}$  obeying dynamics

$$x[t+1] = Ax[t] + Bu[t] + w[t]. (3)$$

Let each expert controller be of the form  $u[t] = K^{(h)}x[t] + z[t]$ , where  $z[t] \in \mathbb{R}^{n_u}$  is Gaussian actuator noise. We assume the system matrices (A,B) remain the same between tasks, but the process noise covariance and the controllers may change. In particular, we have  $w^{(h)}[t] \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \Sigma_w^{(h)})$  and  $z^{(h)}[t] \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma_z^2 I)$  with  $\Sigma_w^{(h)} \succ 0$  for all  $h \in [H+1]$  and  $\sigma_z^2 > 0$ .

We assume all of the expert controllers  $K^{(h)}$  are stabilizing, i.e., the spectral radii of  $A+BK^{(h)}$  are less than one. Note that this implies that (A,B) is stabilizable. No other assumptions on (A,B) are required. The state distribution of the system under each expert controller will converge to the stationary distribution  $\mathcal{N}(0,\Sigma_x^{(h)})$ , where  $\Sigma_x^{(h)}$  solves the following discrete Lyapunov equation:

$$\Sigma_x^{(h)} = (A + BK^{(h)})\Sigma_x^{(h)}(A + BK^{(h)})^{\top} + \sigma_z^2 BB^{\top} + \Sigma_w^{(h)}.$$

For simplicity, we assume that the initial states of the expert demonstrations in our datasets are sampled  $x_i^{(h)}[0] \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \Sigma_x^{(h)})$ . Thus at all times, the marginal state distributions of the expert demonstrations are equal to  $\mathcal{N}(0, \Sigma_x^{(h)})$ .<sup>3</sup>

Finally, we assume that the expert controllers share a low dimensional representation. Specifically, there exists some  $\Phi_{\star} \in \mathbb{R}^{k \times n_x}$  with  $2k \leq {n_x}^4$  and weights  $F_{\star}^{(1)}, F_{\star}^{(2)}, \dots, F_{\star}^{(H+1)} \in \mathbb{R}^{n_u \times k}$  such that for all  $h \in [H+1]$ ,  $K^{(h)} = F_{\star}^{(h)} \Phi_{\star}$ , and the action taken at time t for trajectory i is:<sup>5</sup>

$$u_i^{(h)}[t] = F_{\star}^{(h)} \Phi_{\star} x_i^{(h)}[t] + z_i^{(h)}[t].$$

Under this assumption, the learned common representation  $\hat{\phi}$  in Section 2.1 can be restricted to linear representations, i.e.,  $\hat{\phi}(x) = \hat{\Phi}x$ , where  $\hat{\Phi} \in \mathbb{R}^{k \times n_x}$ . Note that solving Problem (2) with  $\hat{\Phi}$ 

<sup>1.</sup> As the control actions are the labels in IL, actuator noise corresponds to label noise in supervised learning. In the absence of such noise, the controller is recovered by  $n_x$  linearly independent states and corresponding expert inputs.

<sup>2.</sup> This could be the case, for instance, if different controllers are designed for different levels of noise.

<sup>3.</sup> This assumption is not restrictive, as stable systems exponentially converge to stationarity.

<sup>4.</sup> While this assumption is more stringent than the intuitive  $k < n_x$  assumption, it arises from the fact that the residual of the stacked source controllers may be of rank 2k.

<sup>5.</sup> An example of a setting where expert controllers satisfy this assumption is when the system has high dimensional states which exhibit low dimensional structure, e.g. when A and B can be decomposed into  $A = \Phi_{\star}^{\dagger} \tilde{A} \Phi_{\star}$  and  $B = \Phi_{\star}^{\dagger} \tilde{B}$ , where  $\tilde{A} \in \mathbb{R}^{k \times k}$  and  $\tilde{B} \in \mathbb{R}^{k \times n_u}$ . Here, linear policies K which optimize some objective in terms of the low dimensional features of the system can be decomposed into  $K = \tilde{K} \Phi_{\star}$ , where  $\tilde{K} \in \mathbb{R}^{n_u \times k}$ , mirroring the assumptions of our expert controllers. We provide a concrete example in Section 4.

fixed involves solving for only  $kn_u$  parameters, which is smaller than the  $n_un_x$  unknown parameters when learning from scratch. In particular, by representing the controller as  $F^{(H+1)}\Phi$ , we have  $k(n_u+n_x)$  unknown parameters:  $kn_x$  of the parameters are, however, learned using the source task data, leaving only  $kn_u$  parameters to learn with target task data.

#### 2.3. Notation

The Euclidean norm of a vector x is denoted  $\|x\|$ . For a matrix A, the spectral norm is denoted  $\|A\|$ , and the Frobenius norm is denoted  $\|A\|_F$ . The spectral radius of a square matrix is denoted  $\rho(A)$ . We use  $\dagger$  to denote the Moore-Penrose pseudo-inverse. For a square matrix A with  $\rho(A) < 1$ , define  $\mathcal{J}(A) = \sum_{t \geq 0} \|A^t\| < \infty$ . A symmetric, positive semi-definite (psd) matrix  $A = A^{\top}$  is denoted  $A \succeq 0$ . Similarly  $A \succeq B$  denotes that A - B is positive semidefinite. The condition number of a positive definite matrix A is denoted  $\kappa(A) = \frac{\lambda_{\max}(A)}{\lambda_{\min}(A)}$ , where  $\lambda_{\max}$  and  $\lambda_{\min}$  denote the maximum and minimum eigenvalues, respectively. Similarly,  $\sigma_i(A)$  denote the singular values of A. We denote the normal distribution with mean  $\mu$  and covariance  $\Sigma$  by  $\mathcal{N}(\mu, \Sigma)$ . We use standard  $\mathcal{O}(\cdot)$ ,  $\Theta(\cdot)$  and  $\Omega(\cdot)$  to omit universal constant factors, and  $\tilde{\mathcal{O}}(\cdot)$ ,  $\tilde{\Theta}(\cdot)$  and  $\tilde{\Omega}(\cdot)$  to also omit polylog factors. We also use  $a \lesssim b$  to denote a = O(b). We use the indexing shorthand  $[K] := \{1, \dots, K\}$ . For a given task  $h \in [H+1]$ , the matrix of stacked states is defined as

$$\mathbf{X}^{(h)} = \begin{bmatrix} x_1^{(h)}[0] & \dots & x_1^{(h)}[T-1] & \dots & x_{N_1}^{(h)}[0] & \dots & x_{N_1}^{(h)}[T-1] \end{bmatrix}^{\top} \in \mathbb{R}^{N_1 T \times n_x}. \tag{4}$$

Lastly, let  $\bar{\lambda} = \max_{1 \le h \le H} \lambda_{\max}(\Sigma_x^{(h)})$  and  $\underline{\lambda} = \min_{1 \le h \le H} \lambda_{\min}(\Sigma_x^{(h)})$ .

# 3. Sample Complexity of Multi-Task Imitation Learning

In order to derive any useful information from source tasks for a downstream task, the source tasks must satisfy some notion of *task diversity* that sufficiently covers the downstream task. To that end, we introduce the following notions of source tasks covering the target task.

Definition 3.1 (target task covariance coverage (Du et al., 2020)) Define the constant c as:

$$c := \min_{h \in [H]} \lambda_{\min}((\Sigma_x^{(H+1)})^{-1/2} \Sigma_x^{(h)}(\Sigma_x^{(H+1)})^{-1/2}).$$
 (5)

Note that c is well-defined and positive by our assumption that  $\Sigma_w^{(h)} \succ 0$  for all  $h \in [H+1]$ .

Definition 3.1 captures the degree to which the closed-loop distribution of states for each source task aligns with that of the target task. We then introduce the following notion of task similarity between the source and target task weights, which generalizes the well-conditioning assumptions in Du et al. (2020) and Tripuraneni et al. (2021).

**Assumption 3.1 (diverse source controllers)** We assume the target task weights  $F_{\star}^{(H+1)}$  and the stacked source task weights  $\mathbf{F}_{\star} := \left[ \left( F_{\star}^{(1)} \right)^{\top} \quad \cdots \quad \left( F_{\star}^{(H)} \right)^{\top} \right] \in \mathbb{R}^{k \times n_u H}$  satisfy

$$\left\| \left( \mathbf{F}_{\star}^{\dagger} \right)^{\top} F_{\star}^{(H+1)} \right\|^{2} \leq \mathcal{O}\left(\frac{1}{H}\right). \tag{6}$$

Assumption 3.1 states that the alignment and loadings of the singular spaces between the stacked source task weights and target task weights closely match along the low-dimensional representation dimension. For example, if  $F_{\star}^{(h)} = F_{\star}^{(H+1)}$  for each  $h \in [H]$ , the RHS of (6) is 1/H. We note that this assumption subsumes and is more geometrically informative than a direct bound on the ratio of singular values, e.g.  $\sigma_{\max}^2(F_{\star}^{(H+1)})/\sigma_k^2(\mathbf{F}_{\star}) \leq \mathcal{O}(1/H)$ , which would follow by naively extending the well-conditioning assumptions in Du et al. (2020) and Tripuraneni et al. (2021). Notably, such a condition might not be satisfied even if  $F_{\star}^{(h)} = F_{\star}^{(H+1)}$ ,  $\forall h \in [H]$ , e.g., if  $F_{\star}^{(H+1)}$  is rank-deficient.

### 3.1. Excess Risk Bound: Generalization Along Expert Target Task Trajectories

First we show that learning controllers through multi-task representation learning leads to favorable generalization bounds on the excess risk of the learned controller inputs on the expert target task state distribution, analogous to the bounds on multi-task linear regression in Du et al. (2020); Tripuraneni et al. (2021). However, a key complicating factor in our setting is the fact that the input noise  $z^{(h)}[t]$  enters the process, and thus the data  $x^{(h)}[t]$  is causally dependent on the "label noise". In order to overcome this issue and preserve our statistical gains along time T, we leverage the theory of self-normalized martingales, in particular generalizing tools from Abbasi-Yadkori et al. (2011) to the matrix-valued setting. This culminates in the following target task excess risk bound.

**Theorem 3.1 (target task excess risk bound)** Given  $\delta \in (0,1)$ , suppose that

$$N_1 T \gtrsim \max_{h \in [H]} \mathcal{J}\left(A + BK^{(h)}\right)^2 \kappa\left(\Sigma_x^{(h)}\right) (n_x + \log(H/\delta)),$$
  
$$N_2 T \gtrsim \mathcal{J}\left(A + BK^{(H+1)}\right)^2 \kappa\left(\Sigma_x^{(H+1)}\right) (k + \log(1/\delta)).$$

Define  $\mathcal{P}^{(H+1)}_{0:T-1}$  as the distribution over target task trajectories  $(x^{(H+1)}[0], \cdots, x^{(H+1)}[T-1])$ . Then with probability at least  $1-\delta$ , the excess risk of the learned representation  $\hat{\Phi}$  and target task weights  $\hat{F}^{(H+1)}$  is bounded by

$$\operatorname{ER}(\hat{\Phi}, \hat{F}^{(H+1)}) := \frac{1}{2T} \mathbb{E}_{\mathcal{P}_{0:T-1}^{(H+1)}} \left[ \sum_{t=0}^{T-1} \left\| (F_{\star}^{(H+1)} \Phi_{\star} - \hat{F}^{(H+1)} \hat{\Phi}) x^{(H+1)}[t] \right\|^{2} \right]$$

$$\lesssim \sigma_{z}^{2} \left( \frac{k n_{x} \log \left( N_{1} T \frac{\bar{\lambda}}{\underline{\lambda}} \right)}{c N_{1} T H} + \frac{k n_{u} + \log(\frac{1}{\delta})}{N_{2} T} \right).$$
(7)

Note that when we are operating in the setting where we have much more source data than target data, the second term limits the excess risk bound in (7). The second term scales with  $kn_u$ , which is smaller than the number of total parameters in the controller  $n_un_x$ , or  $k(n_u+n_x)$  under the assumption of a low rank (rank-k) controller. Therefore, the benefit of multi-task learning exhibited by this bound is most clear in the setting of underactuation, i.e., when  $n_u \leq n_x$ . It should also be noted that the quantity  $kn_x$  in the numerator of the first term will only be smaller than the number of source controller parameters  $(n_xn_uH)$  if k is much smaller than  $n_uH$ . This is reasonable, because if  $k \geq n_uH$ , an optimal representation could simply contain all of the source task controllers.

### 3.2. Closed-Loop Guarantees: Tackling Distribution Shift

We show that using multi-task representation learning leads to favorable generalization bounds of the performance of the learned target controller in closed-loop. As we are studying the pure offline imitation learning ("behavioral cloning") setting, we do not assume that the expert controllers are optimizing any particular objective. Therefore, to quantify the performance of the controller, we bound the deviation of states generated by the learned and expert target controller run in closed-loop, i.e., the *tracking error*, which implies general expected-cost bounds.

In order to transfer a bound on the excess risk of the target task  $ER(\hat{\Phi}, \hat{F}^{(H+1)})$  into a bound on the tracking error, we must account for the fundamental distribution shift between the expert trajectories seen during training and the trajectories generated by running the learned controller in closed-loop. We leverage the recent framework of Pfrommer et al. (2022) to bound the tracking error, making the necessary modifications to handle stochasticity. Our bound formalizes the notion that "low training error implies low test error," even under the aforementioned distribution shift.

Let us define the following coupling of the states of the expert versus learned target task closed-loop systems: given a learned controller  $\hat{K} = \hat{F}^{(H+1)}\hat{\Phi}$  from solving the pre-training and fine-tuning optimization problems (1) and (2), for a realization of process randomness  $x[0] \sim \mathcal{N}(0, \Sigma_x^{(H+1)})$  and  $z[t] \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_z^2 I)$ ,  $w[t] \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \Sigma_w^{(H+1)})$  for  $t = 0, \dots, T-1$ , we write

$$x_{\star}[t+1] = (A + BK^{(H+1)})x_{\star}[t] + Bz[t] + w[t], \quad x_{\star}[0] = x[0],$$
  
 $\hat{x}[t+1] = (A + B\hat{K})\hat{x}[t] + Bz[t] + w[t], \quad \hat{x}[0] = x[0].$ 

Thus  $\hat{x}[t]$  and  $x_{\star}[t]$  are the states visited by the learned and expert target task systems with the *same* draw of process randomness. We show a high probability bound on the closed-loop tracking error  $||x_{\star}[t] - \hat{x}[t]||$  that scales with the excess risk of the learned controller. Denote by  $\mathcal{P}_{1:T}^{\star}$  and  $\hat{\mathcal{P}}_{1:T}$  the distributions of trajectories  $\{x_{\star}[t]\}_{t=1}^{T}$  and  $\{\hat{x}[t]\}_{t=1}^{T}$ .

**Theorem 3.2 (Target task tracking error bound)** Let  $(\hat{\Phi}, \hat{F}^{(H+1)})$  denote the learned representation and target task weights, and  $\mathrm{ER}(\hat{\Phi}, \hat{F}^{(H+1)})$  denote the corresponding excess risk. Define  $A_{\mathsf{cl}} := A + BK^{(H+1)}$ . Assume that the excess risk satisfies:

$$\operatorname{ER}(\hat{\Phi}, \hat{F}^{(H+1)}) \lesssim \frac{\lambda_{\min}\left(\Sigma_x^{(H+1)}\right)}{\mathcal{J}(A_{\mathsf{cl}})^2 \|B\|^2}.$$
 (8)

Then with probability greater than  $1-\delta$ , for a new target task trajectory sampled with process randomness  $x[0] \sim \mathcal{N}(0, \Sigma_x^{(H+1)})$  and  $z[t] \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_z^2 I)$ ,  $w[t] \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \Sigma_w^{(H+1)})$  for  $t=0,\ldots,T-1$ , the tracking error satisfies

$$\max_{1 \le t \le T} \|\hat{x}[t] - x_{\star}[t]\|^{2} \lesssim \mathcal{J}(A_{\mathsf{cl}})^{2} \|B\|^{2} \log\left(\frac{T}{\delta}\right) \operatorname{ER}(\hat{\Phi}, \hat{F}^{(H+1)}). \tag{9}$$

Furthermore, for any cost function  $h(\cdot)$  that is L-Lipschitz with respect to the trajectory-wise metric  $d(\boldsymbol{x}_{1:T}, \boldsymbol{y}_{1:T}) = \max_{1 \le t \le T} \|x[t] - y[t]\|$ , we have the following bound on the expected cost gap

$$\left| \mathbb{E}_{\hat{\mathcal{P}}_{1:T}}[h(\hat{\boldsymbol{x}}_{1:T})] - \mathbb{E}_{\mathcal{P}_{1:T}^{\star}}[h(\boldsymbol{x}_{1:T}^{\star})] \right| \lesssim L\mathcal{J}(A_{\mathsf{cl}}) \|B\| \sqrt{\log T} \sqrt{\mathrm{ER}(\hat{\Phi}, \hat{F}^{(H+1)})}$$
(10)

By invoking the bound on the excess risk from Theorem 3.1, condition (8) is satisfied with probability at least  $1 - \delta'$  if we have sufficiently many samples  $H, T, N_1, N_2$  such that

$$\sigma_z^2 \left( \frac{k n_x \log \left( N_1 T \frac{\bar{\lambda}}{\underline{\lambda}} \right)}{c N_1 T H} + \frac{k n_u + \log(\frac{1}{\delta'})}{N_2 T} \right) \lesssim \frac{\lambda_{\min} \left( \Sigma_x^{(H+1)} \right)}{\mathcal{J}(A_{\mathsf{cl}})^2 \left\| B \right\|^2}.$$

The bound on excess risk from Theorem 3.1 may also be substituted into the tracking error bound in (9) to find that with probability at least  $1 - \delta - \delta'$ , the tracking error satisfies

$$\max_{1 \le t \le T} \|\hat{x}[t] - x_{\star}[t]\|^2 \lesssim \mathcal{J}(A_{\mathsf{cl}})^2 \|B\|^2 \log \left(\frac{T}{\delta}\right) \sigma_z^2 \left(\frac{k n_x \log \left(N_1 T \frac{\bar{\lambda}}{\underline{\lambda}}\right)}{c N_1 T H} + \frac{k n_u + \log(\frac{1}{\delta'})}{N_2 T}\right).$$

The above inequality provides the informal statement of the main result in Theorem 1.1 by hiding log terms as well as the terms dependent on system parameters. A bound for the expected cost gap  $\left|\mathbb{E}_{\hat{\mathcal{P}}_{1:T}}[h(\hat{\boldsymbol{x}}_{1:T})] - \mathbb{E}_{\mathcal{P}_{1:T}^{\star}}[h(\boldsymbol{x}_{1:T}^{\star})]\right|$  can be similarly instantiated.

**Remark 3.1** The dependence of the tracking error bound in (9) on the stability of the target-task closed-loop system through  $\mathcal{J}(A_{cl})$  is tight (see Zhang et al. (2022)). Intuitively, less stable systems exacerbate the input errors from the learned controller.

Remark 3.2 Some immediate examples of  $h(\cdot)$  include LQR state costs  $h(\boldsymbol{x}_{1:T}) = \max_{t} \|Q^{1/2}x[t]\|$  and regularized tracking costs  $h(\boldsymbol{x}_{1:T}) = \max_{t} \|x[t] - x_{\text{goal}}[t]\| + \lambda \|Rx[t]\|$ . Since  $\frac{1}{T} \sum_{t=1}^{T} \|x[t] - y[t]\| \le \max_{1 \le t \le T} \|x[t] - y[t]\|$ , (10) holds with no modification for time-averaged costs  $h(\cdot)$ . Bounds on the full LQR cost  $h((\boldsymbol{x}_{1:T}, K)) := \max_{1 \le t \le T} \left\| \begin{bmatrix} Q^{1/2} \\ R^{1/2}K \end{bmatrix} x[t] \right\|$  can be similarly derived, and are detailed in Zhang et al. (2022).

### 4. Numerical Results

We consider a simple system with  $n_x = 4$  and  $n_u = 2$  from Hong et al. (2021). In particular, let

$$x[t+1] = \begin{bmatrix} .99 & .03 & -.02 & -.32 \\ .01 & .47 & 4.7 & .00 \\ .02 & -.06 & .40 & .00 \\ .01 & -.04 & .72 & .99 \end{bmatrix} x[t] + \begin{bmatrix} .01 & .99 \\ -3.44 & 1.66 \\ -.83 & .44 \\ -.47 & .25 \end{bmatrix} u[t] =: Ax[t] + Bu[t].$$

We generate a collection of stabilizing controllers  $K^{(1)}, K^{(2)}, \ldots, K^{(H+1)}$  as LQR controllers with different cost matrices. Specifically, let  $R=I_2$ , and  $Q^{(h)}=\alpha^{(h)}I_4$  for  $\alpha^{(h)}\in \text{logspace}(-2,2,H+1)$ , where H=9. The controllers  $K^{(h)}$  are then given by  $K^{(h)}=-(B^\top P^{(h)}B+R)^{-1}B^\top P^{(h)}A$ , where  $P^{(h)}$  solves the following Discrete Algebraic Riccati equation:  $P^{(h)}=A^\top P^{(h)}A+A^\top P^{(h)}B(B^\top P^{(h)}B+R)^{-1}B^\top P^{(h)}A+Q^{(h)}$ .

Next, assume that rather than directly observing the state, we obtain a high dimensional observation given by an injective linear function of the state: such an observation model can be viewed as a linear "perceptual sensor" or camera. In particular, we suppose that  $y_t = Gx_t$ , where  $G \in \mathbb{R}^{50 \times 4}$ . For simplicity, we select the elements of G i.i.d. from  $\mathcal{N}(0,1)$ , which

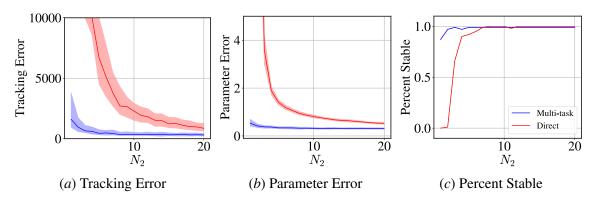


Figure 1: We plot the tracking error between trajectories from the expert and learned controllers,  $\max_{1 \leq t \leq T_{\text{test}}} \|\hat{y}[t] - y_{\star}[t]\|^2, \text{ the parameter error, } \left(\left\|\hat{F}^{(H+1)}\hat{\Phi} - \bar{K}^{(H+1)}\right\|_F\right), \text{ and the percent of stable closed-loop systems for varying amounts of target task data to compare multi-task IL to directly learning the controller from target task data only. All metrics are plotted with respect to the lifted system in Equation (11). Multi-task IL demonstrates a significant benefit over direct IL in all metrics, especially when there is limited target task data.$ 

ensures that G is injective almost surely. The dynamics of the observations may be written  $y[t+1] = GAx_t + GBu[t] = GAG^\dagger y[t] + GBu[t]$ , with the input  $u[t] = K^{(h)}x[t] = K^{(h)}G^\dagger y[t]$ . Define  $\bar{A} = GAG^\dagger$  and  $\bar{B} = GB$ , and  $\bar{K}^{(h)} = K^{(h)}G^\dagger$ . Consider the dynamics in the face of process noise  $w[t] \stackrel{i.i.d.}{\sim} \mathcal{N}(0, I_{50})$ , along with inputs corrupted by noise  $z[t] \stackrel{i.i.d.}{\sim} \mathcal{N}(0, I_{2})$ :

$$y[t+1] = (\bar{A} + \bar{B}\bar{K}^{(h)})y[t] + \bar{B}z[t] + w[t], \quad u[t] = \bar{K}^{(h)}y[t] + z[t]. \tag{11}$$

For the first H controllers, we collect  $N_1$  trajectories of length T=20 to get the pairs  $\left\{\left\{(y_i^h[t],u_i^h[t])\right\}_{t=0}^{T-1}\right\}_{i=1}^{N_1}\right\}_{h=1}^{H}$ . For the last controller, we collect  $N_2$  length T=20 trajectories to get the dataset  $\left\{\left\{(y_i^{H+1}[t],u_i^{H+1}[t])\right\}_{t=0}^{T-1}\right\}_{i=1}^{N_2}$ . Our goal is to learn the controller  $\bar{K}^{(H+1)}$  from the collected state measurements and inputs. We compare the following ways of doing so:

- Multi-task Imitation Learning: We observe that the data generating mechanism ensures the existence of a low dimensional representation. In particular, one possible  $\Phi_{\star}$  is  $G^{\dagger}$ . Therefore, the stage is set for the two step approach outlined in Section 2. In particular, we assume that the true underlying state dimension is known, and we set the low dimensional representation dimension to k=4, and jointly optimize over  $\Phi$ ,  $F^{(1)},\ldots,F^{(H)}$  in Problem (1). We approximately solve this problem with 10000 steps of alternating gradient descent using the adam optimizer (Kingma and Ba, 2014) in optax (Babuschkin et al., 2020) with a learning rate of 0.0001. The learned representation is then fixed, and the target training data is used to optimize  $F^{(H+1)}$ .
- Direct Imitation Learning: We compare multi-task learning to direct learning, which does not leverage the source data. In particular, given the target data, direct learning solves the problem  $\operatorname{minimize}_{F^{(H+1)}} \sum_{i=1}^{N_1} \sum_{t=0}^{T-1} \left\| F^{(H+1)} y_i^{(H+1)}[t] u_i^{(H+1)}[t] \right\|^2$ 6. In this setting, we let  $\hat{\Phi} = I_{50}$ .

<sup>6.</sup> Note that another baseline leverages the fact that a k dimensional representation of the state exists, and learns it using target task data only by solving  $\min_{F(H+1),\Phi} \sum_{i=1}^{N_1} \sum_{t=0}^{T-1} \left\| F^{(H+1)} \Phi y_i^{(H+1)}[t] - u_i^{(H+1)}[t] \right\|^2$ . For the current example, however,  $n_u < k$ , so this approach is less efficient than the direct learning approach proposed.

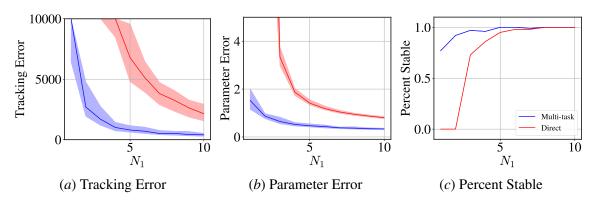


Figure 2: We plot the tracking error between trajectories from the expert and learned controllers,  $\max_{1 \leq t \leq T_{\text{test}}} \|\hat{y}[t] - y_{\star}[t]\|^2$ , the parameter error,  $\left(\left\|\hat{F}^{(H+1)}\hat{\Phi} - \bar{K}^{(H+1)}\right\|_F\right)$ , and the percent of stable closed-loop systems for varying amounts of target task data to compare multi-task IL to directly learning the controller from target task data only. Similar to the setting of multi-task IL for transfer to a new task, leveraging all source data to learn the controller for a single source task provides a significant benefit in all three metrics over direct IL.

Note that a  $2 \times 50$  controller has  $n_u \times n_x = 100$  parameters to learn from the target data. Meanwhile, multi-task imitation learning needs to learn a total of  $k \times n_u + k \times n_x = 208$  parameters for the target controller, but the  $k \times n_u$  parameters are learned using source task data. This leaves only 8 parameters to learn using target task data.

Figure 1 plots three metrics that provide insight into the efficacy of these approaches: the imitation gap given by  $\max_{1 \leq t \leq T_{\text{test}}} \|\hat{y}[t] - y_{\star}[t]\|^2$  for length  $T_{\text{test}} = 100$  observation trajectories  $\hat{y}$  and  $y_{\star}$  rolled out under the learned controller and expert controller, respectively, with the same noise realizations; the parameter error,  $\left\|\hat{F}^{(H+1)}\hat{\Phi} - \bar{K}^{(H+1)}\right\|_F$ ; and the percentage of trials where the learned controller is stabilizing. The trials are over ten realizations of G, as well as ten realizations of the noise, for a total of 100 trials. For each trial,  $N_1 = 10$ , while  $N_2$  sweeps values in [20]. The medians for the imitation gap and parameter error are shown, with the 20% - 80% quantiles shaded.

In Figure 2, we additionally plot these metrics on one of the H source training tasks (arbitrarily chosen as h=7) for varying amounts of training data to demonstrate the efficacy of the approach for multi-task learning. Here,  $N_1$  ranges from 1 to 10. We compare the controller  $\hat{F}^{(h)}\hat{\Phi}$  resulting from the shared training in Problem (1) with the controller from directly training a controller on this task without leveraging source task data. We note that our theoretical results, with mild modification, also support the efficacy of this simultaneous training of a representation and task weights.

#### 5. Conclusion and Future Work

We study the sample complexity of multi-task imitation learning for linear systems. We find that if the different sets of expert demonstrations share a low dimensional representation, and the demonstrations are sufficiently diverse, then doing multi-task representation learning will lead to a smaller tracking error when deploying the learned policy in closed-loop, compared to learning a policy only from target task data. Our results are a first step towards understanding how the performance of a controller trained on multi-task data relates to the characteristics of the multi-task training data and the system being controlled. Some exciting directions for future work would be to extend our analysis to nonlinear systems and nonlinear representation functions, as well as other types of learning algorithms such as model-based and model-free RL.

# Acknowledgments

Katie Kang is supported by a NSF Graduate Research Fellowship. Bruce D. Lee is supported by the DoD through a National Defense Science & Engineering Fellowship. Sergey Levine and Claire Tomlin are supported in part by the DARPA Assured Autonomy program. Nikolai Matni and Thomas Zhang are supported by NSF awards CPS-2038873, CAREER award ECCS-2045834, and a Google Research Scholar award.

### References

- Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Online least squares estimation with self-normalized processes: An application to bandit problems. *arXiv preprint arXiv:1102.2670*, 2011.
- Sanjeev Arora, Simon S. Du, Sham Kakade, Yuping Luo, and Nikunj Saunshi. Provable representation learning for imitation learning via bi-level optimization, 2020. URL https://arxiv.org/abs/2002.10544.
- Igor Babuschkin, Kate Baumli, Alison Bell, Surya Bhupatiraju, Jake Bruce, Peter Buchlovsky, David Budden, Trevor Cai, Aidan Clark, Ivo Danihelka, Claudio Fantacci, Jonathan Godwin, Chris Jones, Ross Hemsley, Tom Hennigan, Matteo Hessel, Shaobo Hou, Steven Kapturowski, Thomas Keck, Iurii Kemaev, Michael King, Markus Kunesch, Lena Martens, Hamza Merzic, Vladimir Mikulik, Tamara Norman, John Quan, George Papamakarios, Roman Ring, Francisco Ruiz, Alvaro Sanchez, Rosalia Schneider, Eren Sezener, Stephen Spencer, Srivatsan Srinivasan, Luyu Wang, Wojciech Stokowiec, and Fabio Viola. The DeepMind JAX Ecosystem, 2020. URL http://github.com/deepmind.
- Jonathan Baxter. Learning internal representations. In *Proceedings of the eighth annual conference on Computational learning theory*, pages 311–320, 1995.
- Yuan Cheng, Songtao Feng, Jing Yang, Hong Zhang, and Yingbin Liang. Provable benefit of multitask representation learning in reinforcement learning, 2022. URL https://arxiv.org/abs/2206.05900.
- Kurtland Chua, Qi Lei, and Jason D Lee. How fine-tuning allows for effective meta-learning. *Advances in Neural Information Processing Systems*, 34:8871–8884, 2021.
- Felipe Codevilla, Matthias Müller, Antonio López, Vladlen Koltun, and Alexey Dosovitskiy. Endto-end driving via conditional imitation learning. In *2018 IEEE international conference on robotics and automation (ICRA)*, pages 4693–4700. IEEE, 2018.
- Koby Crammer, Michael Kearns, and Jennifer Wortman. Learning from multiple sources. *Journal of Machine Learning Research*, 9(8), 2008.
- Marc Peter Deisenroth, Peter Englert, Jan Peters, and Dieter Fox. Multi-task policy search for robotics. In 2014 IEEE International Conference on Robotics and Automation (ICRA), pages 3876–3881, 2014. doi: 10.1109/ICRA.2014.6907421.
- Simon S Du, Wei Hu, Sham M Kakade, Jason D Lee, and Qi Lei. Few-shot learning via learning the representation, provably. *arXiv preprint arXiv:2002.09434*, 2020.

- Yan Duan, John Schulman, Xi Chen, Peter L. Bartlett, Ilya Sutskever, and Pieter Abbeel. Rl<sup>2</sup>: Fast reinforcement learning via slow reinforcement learning, 2016. URL https://arxiv.org/abs/1611.02779.
- Lasse Espeholt, Hubert Soyer, Remi Munos, Karen Simonyan, Volodymir Mnih, Tom Ward, Yotam Doron, Vlad Firoiu, Tim Harley, Iain Dunning, Shane Legg, and Koray Kavukcuoglu. Impala: Scalable distributed deep-rl with importance weighted actor-learner architectures, 2018. URL https://arxiv.org/abs/1802.01561.
- Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks, 2017. URL https://arxiv.org/abs/1703.03400.
- James Harrison, Apoorva Sharma, Roberto Calandra, and Marco Pavone. Control adaptation via meta-learning dynamics. In *Workshop on Meta-Learning at NeurIPS*, volume 2018, 2018.
- Matteo Hessel, Hubert Soyer, Lasse Espeholt, Wojciech Czarnecki, Simon Schmitt, and Hado van Hasselt. Multi-task deep reinforcement learning with popart, 2018. URL https://arxiv.org/abs/1809.04474.
- J Hong, N Moehle, and S Boyd. Lecture notes in "introduction to matrix methods", 2021.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint* arXiv:1412.6980, 2014.
- Lidong Li, Claudio De Persis, Pietro Tesi, and Nima Monshizadeh. Data-based transfer stabilization in linear systems. *arXiv preprint arXiv:2211.05536*, 2022.
- Rui Lu, Gao Huang, and Simon S. Du. On the power of multitask representation learning in linear mdp, 2021. URL https://arxiv.org/abs/2106.08053.
- Rui Lu, Andrew Zhao, Simon S. Du, and Gao Huang. Provable general function class representation learning in multitask bandits and mdps, 2022. URL https://arxiv.org/abs/2205.15701.
- Andreas Maurer, Massimiliano Pontil, and Bernardino Romera-Paredes. The benefit of multitask representation learning, 2015. URL https://arxiv.org/abs/1505.06279.
- Andreas Maurer, Massimiliano Pontil, and Bernardino Romera-Paredes. The benefit of multitask representation learning. *Journal of Machine Learning Research*, 17(81):1–32, 2016.
- Aditya Modi, Mohamad Kazem Shirani Faradonbeh, Ambuj Tewari, and George Michailidis. Joint learning of linear time-invariant dynamical systems. *arXiv preprint arXiv:2112.10955*, 2021.
- Deepan Muthirayan, Dileep Kalathil, and Pramod P Khargonekar. Meta-learning online control for linear dynamical systems. *arXiv preprint arXiv:2208.10259*, 2022.
- Daniel Pfrommer, Thomas TCK Zhang, Stephen Tu, and Nikolai Matni. Tasil: Taylor series imitation learning. *arXiv preprint arXiv:2205.14812*, 2022.

- Kate Rakelly, Aurick Zhou, Deirdre Quillen, Chelsea Finn, and Sergey Levine. Efficient off-policy meta-reinforcement learning via probabilistic context variables, 2019. URL https://arxiv.org/abs/1903.08254.
- Spencer M Richards, Navid Azizan, Jean-Jacques Slotine, and Marco Pavone. Adaptive-control-oriented meta-learning for nonlinear systems. *arXiv preprint arXiv:2103.04490*, 2021.
- Spencer M Richards, Navid Azizan, Jean-Jacques Slotine, and Marco Pavone. Control-oriented meta-learning. *arXiv preprint arXiv:2204.06716*, 2022.
- Stéphane Ross, Geoffrey Gordon, and Drew 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*, pages 627–635. JMLR Workshop and Conference Proceedings, 2011.
- Stefan Schaal. Is imitation learning the route to humanoid robots? *Trends in cognitive sciences*, 3 (6):233–242, 1999.
- Guanya Shi, Kamyar Azizzadenesheli, Michael O'Connell, Soon-Jo Chung, and Yisong Yue. Meta-adaptive nonlinear control: Theory and algorithms. *Advances in Neural Information Processing Systems*, 34:10013–10025, 2021.
- Avi Singh, Eric Jang, Alexander Irpan, Daniel Kappler, Murtaza Dalal, Sergey Levine, Mohi Khansari, and Chelsea Finn. Scalable multi-task imitation learning with autonomous improvement, 2020. URL https://arxiv.org/abs/2003.02636.
- Yee Whye Teh, Victor Bapst, Wojciech Marian Czarnecki, John Quan, James Kirkpatrick, Raia Hadsell, Nicolas Heess, and Razvan Pascanu. Distral: Robust multitask reinforcement learning, 2017. URL https://arxiv.org/abs/1707.04175.
- Nilesh Tripuraneni, Michael Jordan, and Chi Jin. On the theory of transfer learning: The importance of task diversity. *Advances in Neural Information Processing Systems*, 33:7852–7862, 2020.
- Nilesh Tripuraneni, Chi Jin, and Michael Jordan. Provable meta-learning of linear representations. In *International Conference on Machine Learning*, pages 10434–10443. PMLR, 2021.
- Rui Wang, Robin Walters, and Rose Yu. Meta-learning dynamics forecasting using task inference. *arXiv preprint arXiv:2102.10271*, 2021.
- Zhiyuan Xu, Kun Wu, Zhengping Che, Jian Tang, and Jieping Ye. Knowledge transfer in multi-task deep reinforcement learning for continuous control, 2020. URL https://arxiv.org/abs/2010.07494.
- Mengjiao Yang and Ofir Nachum. Representation matters: Offline pretraining for sequential decision making, 2021. URL https://arxiv.org/abs/2102.05815.
- Tianhe Yu, Aviral Kumar, Yevgen Chebotar, Karol Hausman, Sergey Levine, and Chelsea Finn. Conservative data sharing for multi-task offline reinforcement learning, 2021. URL https://arxiv.org/abs/2109.08128.

Thomas T Zhang, Katie Kang, Bruce D Lee, Claire Tomlin, Sergey Levine, Stephen Tu, and Nikolai Matni. Multi-task imitation learning for linear dynamical systems. *arXiv preprint arXiv:2212.00186*, 2022.