
On the Imitation of Non-Markovian Demonstrations: From Low-Level Stability to High-Level Planning

Anonymous Authors¹

Abstract

We propose a theoretical framework for studying the imitation of stochastic, non-Markovian, potentially multi-modal expert demonstrations in nonlinear dynamical systems. Our framework invokes low-level controllers - either learned or implicit in position-command control - to stabilize imitation policies around expert demonstrations. We show that with (a) a suitable low-level stability guarantee and (b) a stochastic continuity property of the learned policy we call “total variation continuity” (TVC), an imitator that accurately estimates actions on the demonstrator’s state distribution closely matches the demonstrator’s distribution over entire trajectories. We then show that TVC can be ensured with minimal degradation of accuracy by combining a popular data-augmentation regimen with a novel algorithmic trick: adding augmentation noise at execution time. We instantiate our guarantees for policies parameterized by diffusion models and prove that if the learner accurately estimates the score of the (noise-augmented) expert policy, then the distribution of imitator trajectories is close to the demonstrator distribution in a natural optimal transport distance. Our analysis constructs intricate couplings between noise-augmented trajectories, a technique that may be of independent interest. We conclude by empirically validating our algorithmic recommendations.

1. Introduction

Training dynamic agents from datasets of expert examples, known as *imitation learning*, promises to take advantage of the plentiful demonstrations available in the modern data

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the Workshop on New Frontiers in Learning, Control, and Dynamical Systems at the International Conference on Machine Learning (ICML). Do not distribute.

environment, in an analogous manner to the recent successes of language models conducting unsupervised learning on enormous corpora of text (Thoppilan et al., 2022; Vaswani et al., 2017). Imitation learning is especially exciting in robotics, where mass stores of pre-recorded demonstrations on Youtube (Abu-El-Haija et al., 2016) or cheaply collected simulated trajectories (Mandlekar et al., 2021; Dasari et al., 2019) can be converted into learned robotic policies.

An outstanding challenge for imitation learning is that demonstrator policies correlate with past actions in sophisticated ways. Multi-modal trajectories present a key example. Consider a robot navigating around an obstacle; because there is no difference between navigating around the object to the right and around to the left, the dataset of expert trajectories may include examples of both options. This bifurcation of good trajectories can make it difficult for the agent to effectively choose which direction to go, possibly even causing the robot to oscillate between directions and run into the object instead of going around it (Chi et al., 2023). Crucially, human demonstrators correlate current actions with the past in order to *commit* to either a right or left path, which makes even formulating the idea of an “expert *policy*” a conceptually challenging one.

In this paper, we develop a theory of imitation learning flexible enough to imitate non-Markovian (e.g. multi-modal or bifurcated as in the above example) demonstrations in smooth, nonlinear control systems. As in previous work, we formalize imitation learning in two stages: at *train-time*, we learn a map from observations to distributions over actions, supervised by (state, action)-pairs from expert demonstrations, while at *test-time*, the learned map, or *policy*, is executed on random initial states (distributed identically to initial training states). What makes imitation learning more challenging than supervised learning is the problem of compounding errors, which may bring the agent to regions of state space not seen during training. Unless one is permitted to collect data adaptively (Laskey et al., 2017; Ross et al., 2011), it is understood that some form of “stability” is required so that the agent navigates back from deviations (Tu et al., 2022; Havens & Hu, 2021).

Contributions. We propose a hierarchical formulation of stability to analyze imitation learning. During training, the

055 learner synthesizes sequences of *primitive controllers* - time-
 056 varying affine control policies which locally stabilize around
 057 each demonstration trajectory. We break these {demonstrator
 058 trajectory, primitive controller} pairs into sub-trajectories
 059 we call “chunks.” Building on (Chi et al., 2023), we use
 060 DDPMs to estimate the conditional distribution of primitive
 061 controller chunks conditioned on recent states from the pre-
 062 vious chunk. We also adopt a popular data-augmentation
 063 technique that corrupts trajectories (but not supervising ac-
 064 tions) with a small amount of Gaussian noise (Ke et al.,
 065 2021; Laskey et al., 2017; Ross et al., 2011). Unlike prior
 066 work, we propose adding noise *back into the policies at*
 067 *inference time*, a technique which is both both provably
 068 indispensable in our analysis, and which our simulations
 069 suggest yields considerable benefit over the conventional
 070 approach of not adding noise at inference time.

071 We prove that the learner can approximate the expert’s trajec-
 072 tory distribution provided three conditions hold: along each
 073 expert trajectory, (a) the dynamics are sufficiently smooth;
 074 (b) one can synthesize primitive controllers that stabilize
 075 the Jacobian-linearized dynamics; and (c) one can approxi-
 076 mately sample from conditional distributions over sequences
 077 of primitive controllers. For concreteness, we formulate
 078 part (c) in the language of Denoising Diffusion Probabilistic
 079 Models (DDPMs), although our results hold for arbitrary
 080 generative models. Our notion of trajectory approximation
 081 is a natural optimal transport metric, which considers a
 082 Wasserstein-like distance between the *marginal* distribu-
 083 tions of visited states, which is strong enough to ensure
 084 closeness of Lipschitz trajectory costs which decompose
 085 across time-steps.
 086

087 Our analysis reformulates our setting as imitation in a com-
 088 posite MDP, where composite states s_h corresponds to tra-
 089 jectory chunks, and composite-actions a_h correspond to sub-
 090 sequences of primitive controllers. A learner’s policy maps
 091 composite-states to distributions over composite-actions,
 092 and a marginalization trick lets us represent non-Markovian
 093 demonstrator trajectories in the same format. The primi-
 094 tive controller sequences a_h provide the requisite stability,
 095 and we show that noising the learner policy at inference
 096 time ensures continuity in the total variation distance (TVC).
 097 Our proof is inspired by the notion of replica symmetry in
 098 statistical physics (Mezard & Montanari, 2009): we show
 099 that by noising at inference time, we consistently estimate
 100 a “replica” policy, which, up to the stability of controllers,
 101 has marginals over states and actions close to those of the
 102 expert policy. The proof constructs a sophisticated coupling
 103 between the learned policy, replica policy, and other inter-
 104 polating sequences; this construction is enabled by subtle
 105 measure-theoretic arguments demonstrating consistency of
 106 our couplings. We also establish stability guarantees for
 107 sequences of primitive controllers in non-linear control sys-
 108 tems, which may be of independent interest. Finally, we
 109

empirically validate the benefits of our proposed augmenta-
 tion strategy in simulated robotic manipulation tasks.

Abridged Related Work. Due to space, we defer a full
 comparison to past work to Appendix B. DDPMs, proposed
 in (Ho et al., 2020; Sohl-Dickstein et al., 2015), along with
 their relatives have seen success in image generation (Song
 & Ermon, 2019; Ramesh et al., 2022), along with imitation
 learning (without data augmentation) (Janner et al., 2022;
 Chi et al., 2023; Pearce et al., 2023), which is the start-
 ing point of our work. Data augmentation is ubiquitous
 in modern imitation learning (Laskey et al., 2017) and our
 approach corresponds to that of (Ke et al., 2021) but with
 noise added at inference time. Despite the benefits of adap-
 tive data collection (Ross et al., 2011; Laskey et al., 2017),
 adaptive demonstrations are more expensive to collect. Pre-
 vious analyses of imitation learning without adaptive data
 collection have focused on classical control-theoretic no-
 tions of stability, notably incremental stability, (Tu et al.,
 2022; Havens & Hu, 2021; Pfrommer et al., 2022), which
 require continuity, Markovianity, and often determinism,
 and preclude the bifurcations permitted in our setting.

Organization. In Section 2 we formally introduce our set-
 ting as well as some preliminary notation and our main
 desideratum. We then state our assumptions and our pro-
 posed algorithm, TODA before giving our main guarantee
 (Theorem 1) in Section 3. In Section 4 we describe our
 proof techniques and provide a high level overview before
 concluding with some experiments in Section 5. The orga-
 nization of our many appendices is given in Appendix A.

2. Setting

Notation and Preliminaries. Appendix A gives a full re-
 view of notation. Bold lower-case (resp. upper-case) denote
 vectors (resp. matrices). We abbreviate the concatenation
 of sequences via $\mathbf{z}_{1:n} = (\mathbf{z}_1, \dots, \mathbf{z}_n)$. Norms $\|\cdot\|$ are Eu-
 clidean for vectors and operator norms for matrices unless
 otherwise noted. Rigorous probability-theoretic prelimi-
 naries are provided in Appendix C. In short, all random
 variables take values in Polish spaces \mathcal{X} (which include
 real vector spaces), the space of Borel distributions on \mathcal{X}
 is denoted $\Delta(\mathcal{X})$. We rely heavily on *couplings* from optimal
 transport theory: given measures $X \sim P$ and $X' \sim P'$
 on \mathcal{X} and \mathcal{X}' respectively, $\mathcal{C}(P, P')$ denotes the space of
 joint distributions $\mu \in \Delta(\mathcal{X} \times \mathcal{X}')$ called “couplings” such
 that $(X, X') \sim \mu$ has marginals $X \sim P$ and $X' \sim P$.
 $\Delta(\mathcal{X} | \mathcal{Y})$ denotes the space of *kernels* $Q : \mathcal{Y} \rightarrow \Delta(\mathcal{X})$
 ; Appendix C rigorously justifies that, in our setting, all
 conditional distributions can be expressed as kernels (which
 we do throughout the paper without comment).

Dynamics and Demonstrations. We consider a discrete-
 time, control system with states $\mathbf{x}_t \in \mathcal{X} := \mathbb{R}^{d_x}$, and inputs

$\mathbf{u}_t \in \mathcal{U} := \mathbb{R}^{d_u}$, obeying the following nonlinear dynamics

$$\mathbf{x}_{t+1} = f(\mathbf{x}_t, \mathbf{u}_t), \quad t \geq 1. \quad (2.1)$$

Given length $T \in \mathbb{N}$, we call sequences $\boldsymbol{\rho}_T = (\mathbf{x}_{1:T+1}, \mathbf{u}_{1:T}) \in \mathcal{P}_T := \mathcal{X}^{T+1} \times \mathcal{U}^T$ **trajectories**. For simplicity, we assume that (2.1) deterministic and address stochastic dynamics in Appendix J. Though the dynamics are Markov and deterministic, we consider a stochastic and possibly *non-Markovian* demonstrator, which allows for the multi-modality described in the Section 1.

Definition 2.1 (Expert Distribution). Let $\mathcal{D}_{\text{exp}} \in \Delta(\mathcal{P}_T)$ denote an **expert distribution** over trajectories to be imitated. $\mathcal{D}_{\mathbf{x}_1}$ denotes the distribution of \mathbf{x}_1 under $\boldsymbol{\rho}_T = (\mathbf{x}_{1:T+1}, \mathbf{u}_{1:T}) \sim \mathcal{D}_{\text{exp}}$.

Primitive Controllers and Synthesis Oracle. Let \mathcal{K} denote the space of affine mappings $\mathcal{X} \rightarrow \mathcal{U}$ (redundantly) parameterized as $\mathbf{x} \mapsto \bar{\mathbf{u}} + \mathbf{K}(\mathbf{x} - \bar{\mathbf{x}})$; we call these **primitive controllers**. We say $\kappa_{1:T} \in \mathcal{K}^T$ is *consistent with* a trajectory $\boldsymbol{\rho} = (\mathbf{x}_{1:T+1}, \mathbf{u}_{1:T}) \in \mathcal{P}_T$ if $\bar{\mathbf{x}}_t = \mathbf{x}_t$ and $\bar{\mathbf{u}}_t = \mathbf{u}_t$ for all $t \in [T]$; note that this implies that $\kappa_t(\mathbf{x}_t) = \mathbf{u}_t$ for all t . A **synthesis oracle** synth maps $\mathcal{P}_T \rightarrow \mathcal{K}^T$ such that, for all $\boldsymbol{\rho}_T \in \mathcal{P}_T$, $\kappa_{1:T} = \text{synth}(\boldsymbol{\rho}_T)$ is consistent with $\boldsymbol{\rho}_T$. For our theory, we assume access to a synthesis oracle at training time, and assume the ability to estimate conditional distributions over joint sequences of primitive controllers; Appendix G explains how this can be implemented by solving Riccati equations if dynamics are known (e.g. in a simulator), smooth, and stabilizable. In our experimental environment, control inputs are desired robot configurations, which the simulated robot executes by applying feedback gains.

Chunking Policies and Indices. The expert distribution \mathcal{D}_{exp} may involve non-Markovian sequences of actions. We imitate these sequences via **chunking policies**. Fix a **chunk length** $\tau_c \in \mathbb{N}$ and **memory length** $\tau_m \leq \tau_c$, and define time indices $t_h = (h-1)\tau_c + 1$. For simplicity, we assume τ_c divides T , and set $H = T/\tau_c$. Given a $\boldsymbol{\rho}_T \in \mathcal{P}_T$, define the **trajectory-chunks** $\boldsymbol{\rho}_{c,h} := (\mathbf{x}_{t_{h-1}:t_h}, \mathbf{u}_{t_{h-1}:t_h-1}) \in \mathcal{P}_{\tau_c}$ and **memory-chunks** $\boldsymbol{\rho}_{m,h} := (\mathbf{x}_{t_h-\tau_m+1:t_h}, \mathbf{u}_{t_h-\tau_m+1:t_h-1}) \in \mathcal{P}_{\tau_m-1}$ for $h > 1$, and $\boldsymbol{\rho}_{c,1} = \boldsymbol{\rho}_{m,1} = \mathbf{x}_1$. We call τ_c -length sequences of primitive controllers **composite actions** $\mathbf{a}_h = \kappa_{t_h:t_{h-1}} \in \mathcal{A} := \mathcal{K}^{\tau_c}$. A **chunking policy** $\pi = (\pi_h)$ consists of functions π_h mapping memory chunks $\boldsymbol{\rho}_{m,h}$ to distributions $\Delta(\mathcal{A})$ over composite actions and interacting with the dynamics (2.1) by $\mathbf{a}_h = \kappa_{t_h:t_{h-1}} \sim \pi_h(\boldsymbol{\rho}_{m,h})$, and executing $\mathbf{u}_t = \kappa_t(\mathbf{x}_t)$. The chunking scheme is represented in Figure 1 in Section 4, alongside the abstraction we use in our analysis.

Desideratum. The quality of imitation of a deterministic policy is naturally measured in terms of step-wise closeness of state and action (Tu et al., 2022; Pfrommer et al., 2022).

In stochastic settings, however, two rollouts of even the same policy can visit different states. We propose measuring *distributional closeness* via *couplings* introduced in the preliminaries above. We define the following losses:

Definition 2.2. Given $\varepsilon > 0$ and a (chunking) policy π , the imitation loss $\mathcal{L}_{\text{marg},\varepsilon}(\pi)$ is defined to be

$$\max_{t \in [T]} \inf_{\mu} \left\{ \mathbb{P}_{\mu} [\|\mathbf{x}_{t+1}^{\text{exp}} - \mathbf{x}_{t+1}^{\pi}\| > \varepsilon], \mathbb{P}_{\mu} [\|\mathbf{u}_t^{\text{exp}} - \mathbf{u}_t^{\pi}\| > \varepsilon] \right\}$$

where the infimum is over all couplings μ between the distribution of $\boldsymbol{\rho}_T$ under \mathcal{D}_{exp} and that induced by the policy π as described above, such that $\mathbb{P}_{\mu}[\mathbf{x}_1^{\text{exp}} = \mathbf{x}_1^{\pi}] = 1$. Also define $\mathcal{L}_{\text{fin},\varepsilon}(\pi) := \inf_{\mu} \mathbb{P}_{\mu} [\|\mathbf{x}_{T+1}^{\text{exp}} - \mathbf{x}_{T+1}^{\pi}\| > \varepsilon]$, the loss restricted to the final states under each distribution.

Under stronger conditions (whose necessity we establish), we can also imitate joint distributions over actions (Appendix I). Observe that $\mathcal{L}_{\text{fin},\varepsilon} \leq \mathcal{L}_{\text{marg},\varepsilon}$, and that both losses are equivalent to Wasserstein-type metrics on bounded domains (and correspond to total variation analogues of shifted Renyi divergences (Altschuler & Talwar, 2022; Altschuler & Chewi, 2023)). While empirically evaluating these infima over couplings is challenging, $\mathcal{L}_{\text{marg},\varepsilon}$ upper bounds the difference in expectation between any bounded and Lipschitz control cost decomposing across time steps, states and inputs, and $\mathcal{L}_{\text{fin},\varepsilon}$ upper bounds differences in bounded, Lipschitz final-state costs; see Appendix I for further discussion.

Diffusion Models. Our analysis provides imitation guarantees when chunking policies π_h select \mathbf{a}_h via a sufficiently accurate generative model. Given their recent success, we adopt the popular Denoising Diffusion Probabilistic Models (DDPM) framework (Chen et al., 2022; Lee et al., 2023) that allows the learner to sample from a density $q \in \Delta(\mathbb{R}^d)$ assuming that the *score* $\nabla \log q$ is known to the learner. More precisely, suppose the learner is given an observation $\boldsymbol{\rho}_{m,h}$ and wishes to sample $\mathbf{a}_h \sim q(\cdot | \boldsymbol{\rho}_{m,h})$ for some family of probability kernels $q(\cdot | \cdot)$. A DDPM starts with some \mathbf{a}_h^0 sampled from a standard Gaussian noise and iteratively “denoises” for each DDPM-time step $0 \leq j < J$:

$$\mathbf{a}_h^j = \mathbf{a}_h^{j-1} - \alpha \cdot \mathbf{s}_{\theta,h}(\mathbf{a}_h^{j-1}, \boldsymbol{\rho}_{m,h}, j) + 2 \cdot \mathcal{N}(0, \alpha^2 \mathbf{I}), \quad (2.2)$$

where $\mathbf{s}_{\theta,h}(\mathbf{a}_h^j, \boldsymbol{\rho}_{m,h}, j)$ estimates the true score $\mathbf{s}_{*,h}(\mathbf{a}_h, \boldsymbol{\rho}_{m,h}, \alpha j)$, formally defined for any continuous argument $t \leq J\alpha$ to be $\mathbf{s}_{*,h}(\mathbf{a}, \boldsymbol{\rho}_{m,h}, t) := \nabla_{\mathbf{a}} \log q_{[t]}(\mathbf{a} | \boldsymbol{\rho}_{m,h})$, where $q_{[t]}(\cdot | \boldsymbol{\rho}_{m,h})$ is the distribution of $e^{-t} \mathbf{a}_h^{(0)} + \sqrt{1 - e^{-2t}} \boldsymbol{\gamma}$ with $\mathbf{a}_h^{(0)} \sim q(\cdot | \boldsymbol{\rho}_{m,h})$ and $\boldsymbol{\gamma} \sim \mathcal{N}(0, \mathbf{I})$ is a standard Gaussian. We will denote by DDPM($\mathbf{s}_{\theta}, \boldsymbol{\rho}_{m,h}$) the law of \mathbf{a}_h^J sampled according to the DDPM using $\mathbf{s}_{\theta}(\cdot, \boldsymbol{\rho}_{m,h}, \cdot)$ as a score estimator. Preliminaries on DPPMs are detailed in Appendix H.

3. Algorithm and Results

We show that trajectories of the form given in [Definition 2.1](#) can be efficiently imitated if (a) we are given a synthesis oracle that locally stabilizes chunks of the trajectory with primitive controllers and (b) the score of the following conditional distributions (whose existence is guaranteed by [Appendix C](#)) lies in a class Θ of bounded statistical complexity.

Formal Assumptions. We say trajectory $\rho_\tau = (\mathbf{x}_{1:\tau+1}, \mathbf{u}_{1:\tau}) \in \mathcal{P}_\tau$ is *feasible* if it obeys the dynamics in [\(2.1\)](#). We assume that the transition map f takes the form of an Euler-like discretization

$$f(\mathbf{x}_t, \mathbf{u}_t) = \mathbf{x}_t + \eta f_\eta(\mathbf{x}_t, \mathbf{u}_t)$$

for a small step size $\eta > 0$ and say ρ_τ is $(R_{\text{dyn}}, L_{\text{dyn}}, M_{\text{dyn}})$ -regular if, for any $t \in [\tau]$ and $(\mathbf{x}'_t, \mathbf{u}'_t) \in \mathbb{R}^{d_x} \times \mathbb{R}^{d_u}$ such that $\|\mathbf{x}'_t - \mathbf{x}_t\| \vee \|\mathbf{u}'_t - \mathbf{u}_t\| \leq R_{\text{dyn}}$, it holds that $\|\nabla f_\eta(\mathbf{x}'_t, \mathbf{u}'_t)\|_{\text{op}} \leq L_{\text{dyn}}$ and $\|\nabla^2 f_\eta(\mathbf{x}'_t, \mathbf{u}'_t)\|_{\text{op}} \leq M_{\text{dyn}}$.¹ The **Jacobian linearizations** along a path $\rho_\tau = (\mathbf{x}_{1:\tau+1}, \mathbf{u}_{1:\tau}) \in \mathcal{P}_\tau$ are matrices $\mathbf{A}_t(\rho_\tau) := \frac{\partial}{\partial \mathbf{x}} f(x_t, u_t)$ and $\mathbf{B}_t(\rho_\tau) := \frac{\partial}{\partial \mathbf{u}} f(x_t, u_t)$ for $t \in [\tau]$. Given $\rho_\tau \in \mathcal{P}_\tau$ and primitive controllers $\kappa_{1:\tau}$, expressed as $\kappa_t(\mathbf{x}) = \bar{\mathbf{K}}_t(\mathbf{x} - \bar{\mathbf{x}}_t) + \bar{\mathbf{u}}_t(x)$, we say $(\rho_\tau, \kappa_{1:\tau})$ are $(R_{\text{stab}}, B_{\text{stab}}, L_{\text{stab}})$ -**Jacobian stable** if (a) $\kappa_{1:\tau}$ is consistent with ρ_τ (b) $\max_{t \in [\tau]} \|\bar{\mathbf{K}}_t\| \vee \|\bar{\mathbf{x}}_t\| \vee \|\bar{\mathbf{u}}_t\| \leq R_{\text{stab}}$, and (c) the linearized closed-loop transition operator has exponential decay:

$$\|\Phi_{\text{cl},k,j}\|_{\text{op}} \leq B_{\text{stab}} \left(1 - \frac{\eta}{L_{\text{stab}}}\right)^{k-j}$$

$$\Phi_{\text{cl},k,j} := (\mathbf{I} + \eta \mathbf{A}_{\text{cl},k-1}) \cdot (\mathbf{I} + \eta \mathbf{A}_{\text{cl},k-2}) \cdots (\mathbf{I} + \eta \mathbf{A}_{\text{cl},j}),$$

where above $\mathbf{A}_{\text{cl},k} = \mathbf{A}_k(\rho_\tau) + \mathbf{B}_{k-1}(\rho_\tau) \bar{\mathbf{K}}_{k-1}$. Our first two assumptions are as follows.

Assumption 3.1. The $\rho_T \sim \mathcal{D}_{\text{exp}}$ is feasible and $(R_{\text{dyn}}, L_{\text{dyn}}, M_{\text{dyn}})$ -regular with probability 1.

Assumption 3.2. With probability 1 over $\rho_T \sim \mathcal{D}_{\text{exp}}$ and $\kappa_{1:T} = \text{synth}(\rho_T)$, the chunk-action pairs $(\rho_{c,h+1}, \mathbf{a}_h)$ are $(R_{\text{stab}}, B_{\text{stab}}, L_{\text{stab}})$ -Jacobian Stable for $1 \leq h \leq H$.

[Assumption 3.1](#) enforces smoothness of the dynamics, but *not* smoothness or continuity of the underlying policy. [Assumption 3.2](#) generalizes popular quantifications of stability (e.g. strong stability ([Cohen et al., 2019](#))), and is satisfied when primitive controllers are synthesized via Riccati equations of dynamics with stabilizable linearizations ([Appendix G](#)). Finally, we require access to a class of score functions rich enough to represent the deconvolution conditionals, defined as follows.

Definition 3.1 (Deconvolution Conditionals). For $h \in [H]$, let $\pi_{\text{dec},h}^* \in \Delta(\mathcal{A} | \mathcal{P}_{\tau_m-1})$ denote a conditional distribution

¹Here, $\|\nabla^2 f_\eta(\mathbf{x}'_t, \mathbf{u}'_t)\|_{\text{op}}$ denotes the operator-norm of a three-tensor.

of $\mathbf{a}_h = \kappa_{t_h:t_{h+1}-1} | \tilde{\rho}_{m,h}$, where $\rho_T \sim \mathcal{D}_{\text{exp}}$, $\kappa_{1:T} = \text{synth}(\rho_T)$, and $\rho_{m,h}$ is the memory chunk of ρ_T at step h , and $\tilde{\rho}_{m,h} \sim \mathcal{N}(\rho_{m,h}, \sigma^2 \mathbf{I})$ augments $\rho_{m,h}$ with noise.

Assumption 3.3. For $h \in [H]$ let $\pi_{\text{dec},h,[t]}^* \in \Delta(\mathcal{A} | \mathcal{P}_{\tau_m-1})$ denote $q_{[t]}$ as defined below [\(2.2\)](#) for $q = \pi_{\text{dec},h}^*$ the deconvolution policy defined above. For fixed $\sigma, \alpha > 0$ and $j \in \mathbb{N}$, let $\mathbf{s}_{*,h,\sigma,[j]}$ denote the score function of $\pi_{\text{dec},h,[\alpha j]}^*$. We suppose that for any $J \in \mathbb{N}$ and $\alpha, \sigma > 0$, we are given a class of scores $\Theta = \Theta(\tau_c, \tau_m, \sigma) = \{\mathbf{s}_{\theta,1:H}\} = \bigcup_{j \in [J]} \Theta_j$ such that (a) for all $1 \leq j \leq J$, $\mathbf{s}_{*,h,\sigma,[\alpha j]} \in \Theta_j$ and (b) a Rademacher-like complexity of Θ_j , $\mathcal{R}_n(\Theta_j)$ (defined in [Appendix H](#)) has polynomial decay $\mathcal{R}_n(\Theta_j) \leq C_\Theta (1/\alpha)^\nu n^{-1/\nu}$ for some $\nu \geq 1$ and $C_\Theta = C_\Theta(\sigma, \tau_c, \tau_m)$.

As justified in [Appendix H](#), the above assumption is a natural for statistical learning, the decay condition on $\mathcal{R}_n(\Theta)$ holds for most common function classes (often with $\nu \leq 2$ and even more benign dependence on J, α), and our results extend to approximate realizability. $\mathcal{R}_n(\Theta)$ depends implicitly on chunk and memory lengths $\tau_c, \tau_m > 0$ and problem dimension through the specification of $\mathbf{s}_{*,h,\sigma,[\alpha j]}$. Realizability is motivated by the approximation power of deep neural networks ([Bartlett et al., 2021](#)).

Algorithm. Our proposed algorithm, TODA ([Algorithm 1](#)) combines DDPM-learning of chunked policies as in ([Chi et al., 2023](#)) with a popular form of data-augmentation ([Ke et al., 2021](#)). We collect N_{exp} expert trajectories, synthesis gains, and segment trajectories into memory chunks $\rho_{m,h}$ and composite actions \mathbf{a}_h as described in [Section 2](#). We perturb each $\rho_{m,h}$ to form N_{aug} chunks $\tilde{\rho}_{m,h}$, as well as horizon indices $j \in [J]$ and inference noises $\gamma \sim \mathcal{N}(\mathbf{a}_h, (\alpha j h)^2 \mathbf{I})$, and add these tuples $(\mathbf{a}_h, \tilde{\rho}_{m,h}, j_h, \gamma_h, h)$ to our data \mathcal{D} . We end the training phase by minimizing the standard DDPM loss ([Song & Ermon, 2019](#)) $\mathcal{L}_{\text{DDPM}}(\theta, \mathcal{D})$:

$$\sum \left\| \gamma_h - \mathbf{s}_{\theta,h} \left(e^{-\alpha j} \mathbf{a}_h + \sqrt{1 - e^{-2\alpha j}} \gamma_h, \tilde{\rho}_{m,h}, j_h \right) \right\|^2, \quad (3.1)$$

where the sum is over $(\mathbf{a}_h, \tilde{\rho}_{m,h}, j_h, \gamma_h, h) \in \mathcal{D}$. Our algorithm differs subtly from past work in [Line 8](#): we add augmentation noise *back in* at test time. Here, the notation $\text{DDPM}(\mathbf{s}_{\theta,h}, \cdot) \circ \mathcal{N}(\rho_{m,h}, \sigma^2 \mathbf{I})$ means, given $\rho_{m,h}$, we perturb it to $\tilde{\rho}_{m,h} \sim \mathcal{N}(\rho_{m,h}, \sigma^2 \mathbf{I})$, and sample $\mathbf{a}_h \sim \text{DDPM}(\mathbf{s}_{\theta,h}, \tilde{\rho}_{m,h})$. The motivation for this is that adding noise at inference time *removes distribution shift* coming from training on augmented data; this simple observation is crucial for our theoretical guarantees.

Theoretical Guarantee. We now state our main theorem, which bounds the imitation losses of TODA trained on expert demonstrations. Let $d = \tau_c(d_x + d_u + d_x d_u)$, and let c_1, \dots, c_5 denote terms given in [Appendix G](#) that are polynomial in the parameters in [Assumptions 3.1](#) and [3.2](#).

Algorithm 1 Trajectory Optimization with Data Augmentation (TODA)

```

1: Initialize Synthesis oracle  $\text{synth}$ , sample sizes
    $N_{\text{exp}}, N_{\text{aug}} \in \mathbb{N}$ ,  $\sigma \geq 0$ , DDPM step size  $\alpha > 0$ ,
   DDPM horizon  $J$ , function class  $\{\mathbf{s}_\theta\}_{\theta \in \Theta}$ , gain magni-
   tude  $R > 0$ , empty data buffer  $\mathcal{D} \leftarrow \emptyset$ .
2: For no augmentation, set  $\sigma = 0$  and
    $N_{\text{aug}} = 1$ 
3: for  $n = 1, 2, \dots, N_{\text{exp}}$  do
4:   Sample  $\boldsymbol{\rho}_T = (x_{1:T+1}, u_{1:T}) \sim \mathcal{D}_{\text{exp}}$  and set
    $\kappa_{1:T} = \text{synth}(\boldsymbol{\rho})$ 
5:   % Segment  $\boldsymbol{\rho}_{m,1:H}$  from  $\boldsymbol{\rho}_T$  and  $\mathbf{a}_{1:H}$  from
    $\kappa_{1:T}$ 
6:   for  $i = 1, 2, \dots, N_{\text{aug}}$  and  $h = 1, 2, \dots, H$  do
7:     Sample  $\tilde{\boldsymbol{\rho}}_{m,h} \sim \mathcal{N}(\boldsymbol{\rho}_{m,h}, \sigma^2 \mathbf{I})$ ,  $j_h \sim \text{Unif}([J])$ 
   and  $\gamma_h \sim \mathcal{N}(\mathbf{a}_h, (j_h \alpha)^2 \mathbf{I})$ .
8:      $\mathcal{D} \leftarrow \mathcal{D}.\text{append}(\{(\mathbf{a}_h, \tilde{\boldsymbol{\rho}}_{c,h}, j_h, \gamma_h, h)\})$ 
9:   Fit  $\theta \in \arg \min_{\theta \in \Theta} \mathcal{L}_{\text{DDPM}}(\theta, \mathcal{D})$ 
10:  return  $\hat{\pi}_\sigma = (\hat{\pi}_{1:H})$ , where  $\hat{\pi}_{h,\sigma}(\boldsymbol{\rho}_{m,h}) =$ 
   DDPM( $\mathbf{s}_{\theta,h}, \cdot$ )  $\circ \mathcal{N}(\boldsymbol{\rho}_{m,h}, \sigma^2 \mathbf{I})$ .

```

Theorem 1. Consider running TODA for $\sigma > 0$ with parameters J, α polynomial in the parameters given in Assumptions 3.1 and 3.2 specified in Appendix H. Suppose that Assumptions 3.1 to 3.3 hold and further suppose the chunk length satisfies $\tau_c \geq c_3/\eta$. Given $\sigma, \delta > 0$, select any $\varepsilon > 0$ for which $5d_x + 2 \log\left(\frac{4\sigma}{\varepsilon}\right) \leq c_4^2/(16\sigma^2)$. If $N_{\text{exp}} \geq \text{poly}(C_\Theta, \varepsilon/\sigma, R_{\text{stab}}, d, \log(H/\delta))^\nu$, then for $\hat{\pi}_\sigma$ the policy output by TODA, it holds with probability $1 - \delta$ over the training data that both $\mathcal{L}_{\text{marg},\varepsilon_1}(\hat{\pi}_\sigma)$ and $\mathcal{L}_{\text{fin},\varepsilon_2}(\hat{\pi}_\sigma)$ are upper bounded by

$$H \left(\frac{3\varepsilon}{\sigma} + 6c_5 \sqrt{5d_x + 2 \log\left(\frac{4\sigma}{\varepsilon}\right)} e^{-\frac{\eta(\tau_c - \tau_m)}{L_{\text{stab}}}} \right) \quad (3.2)$$

where $\varepsilon_1 = \varepsilon + 4c_5\sigma \cdot (5d_x + 2 \log\left(\frac{4\sigma}{\varepsilon}\right))^{1/2}$ and $\varepsilon_2 = \varepsilon + 4c_5 e^{-\eta\tau_c/L_{\text{stab}}} \sigma \cdot (5d_x + 2 \log\left(\frac{4\sigma}{\varepsilon}\right))^{1/2}$.

Theorem 1 guarantees imitation of the distribution of marginals and final states of \mathcal{D}_{exp} . Each term in (3.2) can be made small by decreasing the amount of noise σ in the augmentation, increasing the number of trajectories, and increasing the chunk length τ_c . Increasing τ_c comes at the (implicit) expense of requiring a more expressive score class Θ (requiring greater N_{exp}); similarly, as expressed in Appendix H, the scores $\mathbf{s}_{*,h,\sigma,[\alpha j]}$ may become harder to learn σ decreases. Note that the contribution of the additive σ -term in ε_2 , used for the final-state loss $\mathcal{L}_{\text{fin},\varepsilon}$, is exponentially-in- τ_c smaller than that in ε_1 . Interestingly, our theory suggest no benefit to increasing τ_m (corroborated empirically in (Chi et al., 2023)). Appendix I gives guarantees for imitating joint trajectories under the further assumptions that (a) the demonstrator has memory (or, more

generally, a mixing time) of at most τ_m , and (b) either the demonstrator distribution happens to satisfy a certain continuity property, or $\sigma = 0$ and instead the learned $\hat{\pi}$ satisfies that same property.

Theorem 1 leverages statistical learning guarantees for DPPMs to show our learned policy approximately samples from $\pi_{\text{dec},h}^*$ in a truncated Wasserstein distance (Appendix H). These steps are combined with a general template for imitation learning developed in Section 4, with the final proof deferred to Appendix I. In Appendix F we show that this framework is essentially tight and thus suboptimality in Theorem 1 comes from looseness in conditional sampling guarantees. If we were able to approximately sample from $\pi_{\text{dec},h}^*$ in total variation, rather than a truncated Wasserstein distance, the imitation learning problem would be trivialized (Appendix I). Appendix H explains that the needed assumptions for this stronger sense of approximate sampling do not hold in our setting, because expert distributions over actions typically lie on low-dimensional manifolds.

Stability, limitations, and future work. We never explicitly model bifurcations; rather, we allow expert demonstrations to be sufficiently rich as to permit them. Eschewing global stability, τ_c ensures that trajectories are long enough for the strictly local stability assumptions in Assumption 3.2 to provide benefit. Thus, non-Markovianity is challenging only insofar as it relates to the difficulty of local stabilization. A key limitation of our work is that, to take advantage of local stability, we rely on either synthesized primitive controllers (in our analysis) or low-level stabilizing controllers built into problem environments (in our experiments). Developing a more comprehensive approach to stability (perhaps one that does not require explicit gain synthesis, and extends to non-smooth systems) is an exciting direction for future work. Appendix B compares our hierarchical approach to stability to more standard notions, which we show rule out the possibility for bifurcated demonstrations.

4. Analysis

Our analysis abstracts away the vector-valued dynamics into a deterministic MDP with *composite-states* $\mathbf{s} \in \mathcal{S}$ and *composite-actions* $\mathbf{a} \in \mathcal{A}$, with dynamics

$$\mathbf{s}_{h+1} = F_h(\mathbf{s}_h, \mathbf{a}_h), \quad h \in \{1, 2, \dots, H\} \quad (4.1)$$

A *composite-policy* π is a sequence of kernels $\pi_1, \pi_2, \dots, \pi_H : \mathcal{S} \rightarrow \Delta(\mathcal{A})$. We let P_{init} denote the distribution of initial state \mathbf{s}_1 , and D_π denote the distribution of $(\mathbf{s}_{1:H+1}, \mathbf{a}_{1:H})$ subject to $\mathbf{s}_1 \sim P_{\text{init}}$, $\mathbf{a}_h | \mathbf{s}_{1:h}, \mathbf{a}_{1:h-1} \sim \pi_h(\mathbf{s}_h)$, and the composite-dynamics (4.1). We assume that we have an optimal policy π^* to be imitated, and define P_h^* as the marginal distribution of \mathbf{s}_h under D_{π^*} .

Structure of the proof. We begin by explaining key objects,

275 stability and continuity properties required in the composite
 276 MDP. Then, [Section 4.1](#) relates the composite MDP to our
 277 original setting by taking composite-states $s_h = \rho_{c,h}$ as
 278 chunks, and taking composite actions as sequences of primi-
 279 tive controllers $a_h = \kappa_{t_h:t_{h+1}-1}$ as in [Section 2](#). We also
 280 explain why relevant stability and continuity conditions are
 281 met. Finally, we derive [Theorem 1](#) from a generic guaran-
 282 tee for smoothed imitation learning in the composite MDP,
 283 [Theorem 2](#), and from sampling guarantees in [Appendix H](#).

284 We consider two pseudometrics on the space \mathcal{S} : $d_S, d_{\text{TVC}} :$
 285 $\mathcal{S}^2 \rightarrow \mathbb{R}_{\geq 0}$, and a function $d_A : \mathcal{A}^2 \rightarrow \mathbb{R}_{\geq 0}$. For conven-
 286 ience, *do not require* d_A to satisfy the axioms of a pseudo-
 287 metric. We use d_S and d_A to measure error between states
 288 and actions, respectively, and $d_{\text{TVC}}(\cdot, \cdot)$ for a probabilistic
 289 continuity property described below. In terms of d_S and d_A ,
 290 we consider three measures of imitation error: error on the
 291 (i) joint distribution of trajectories ($\Gamma_{\text{joint},\varepsilon}$) (ii) marginal
 292 distribution of trajectories ($\Gamma_{\text{marg},\varepsilon}$) and (iii) one-step error
 293 in actions ($d_{\text{os},\varepsilon}$). Formally:

294 **Definition 4.1** (Imitation Errors). Given an error param-
 295 eter $\varepsilon > 0$, define the **joint-error** $\Gamma_{\text{joint},\varepsilon}(\hat{\pi} \parallel \pi^*) :=$
 296 $\inf_{\mu_1} \mathbb{P}_{\mu_1} [\max_{h \in [H]} \max\{d_S(s_{h+1}^*, \hat{s}_{h+1}), d_A(a_h^*, \hat{a}_h)\} > \varepsilon]$,
 297 where the first infimum is over trajectory cou-
 298 plings $((\hat{s}_{1:H+1}, \hat{a}_{1:H}), (s_{1:H+1}^*, a_{1:H}^*)) \sim \mu_1 \in$
 299 $\mathcal{C}(\mathcal{D}_{\hat{\pi}}, \mathcal{D}_{\pi^*})$ satisfying $\mathbb{P}_{\mu_1}[\hat{s}_1 = s_1^*] = 1$.
 300 Define the **marginal error** $\Gamma_{\text{marg},\varepsilon}(\hat{\pi} \parallel$
 301 $\pi^*) := \max_{h \in [H]} \{\inf_{\mu_1} \mathbb{P}_{\mu_1}[d_S(s_{h+1}^*, \hat{s}_{h+1}) >$
 302 $\varepsilon], \inf_{\mu_1} \mathbb{P}_{\mu_1}[d_A(a_h^*, \hat{a}_h) > \varepsilon]\}$ to be the same as the
 303 to joint-gap, with the “max” outside the probability
 304 and inf over couplings. Lastly, define the **one-step error**
 305 $d_{\text{os},\varepsilon}(\hat{\pi}_h(s) \parallel \pi_h^*(s)) := \inf_{\mu_2} \mathbb{P}_{\mu_2}[d_A(\hat{a}_h, a_h^*) \leq \varepsilon]$, where
 306 the infimum is over $(a_h^*, \hat{a}_h) \sim \mu_2 \in \mathcal{C}(\hat{\pi}_h(s), \pi_h^*(s))$.

307 **Stability.** Our hierarchical approach offloads stability of
 308 stochastic π^* onto that of its composite-actions a_h , instanta-
 309 tated as *primitive controllers* (not raw inputs!). This allows
 310 us to circumvent more challenging incremental senses of
 311 stability (see [Appendix B](#) for further discussion).

312 **Definition 4.2** (Input-Stability). A trajectory $(s_{1:H+1}, a_{1:H})$
 313 is **input-stable** if all sequences $s'_1 = s_1$ and $s'_{h+1} =$
 314 $F_h(s'_h, a'_h)$ satisfy $d_S(s'_{h+1}, s_{h+1}) \vee d_{\text{TVC}}(s'_{h+1}, s_{h+1}) \leq$
 315 $\max_{1 \leq j \leq h} d_A(a'_j, a_j)$, $\forall h \in [H]$. A policy π is **input-**
 316 **stable** if $(s_{1:H}, a_{1:H}) \sim \mathcal{D}_{\pi}$ is **input-stable** almost surely.

317 **TVC.** Continuity of probability kernels and policies in TV
 318 distance are measured in terms of d_{TVC} .

319 **Definition 4.3.** For a measure-space \mathcal{X} and non-decreasing
 320 $\gamma : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}_{\geq 0}$, we call a probability kernel $W : \mathcal{S} \rightarrow$
 321 $\Delta(\mathcal{X})$ **γ -total variation continuous (γ -TVC)** if, for all
 322 $s, s' \in \mathcal{S}$, $\text{TV}(W(s), W(s')) \leq \gamma(d_{\text{TVC}}(s, s'))$. A policy
 323 π is **γ -TVC** if $\pi_h : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ is γ -TVC $\forall h \in [H]$.

324 **Smoothing.** In [Appendix D](#), we show that under the strong
 325 condition that the learned policy $\hat{\pi}$ is γ -TVC, then TODA
 326

with no data augmentation ($\sigma = 0$) learns the distribution.
 Frequently, however, $\hat{\pi}$ may not satisfy this condition, such
 as when the ground truth π^* is not also TVC. We circumvent
 this by introducing a *smoothing kernel* $W_\sigma : \mathcal{S} \rightarrow \Delta(\mathcal{S})$
 that corresponds to the data augmentation; in TODA we let
 the kernel be a Gaussian, sending $\rho_{m,h}$ to $\mathcal{N}(\rho_{m,h}, \sigma^2 \mathbf{I}) \in$
 $\Delta(\mathcal{P}_{\rho_{m,h}})$. We will thus be able to replace TVC of $\hat{\pi}$ with
 TVC of W_σ . We now introduce a few key objects.

Definition 4.4. Given a policy π , we define its **smoothed**
policy $\pi \circ W_\sigma$ via components $(\pi \circ W_\sigma)_h = \pi_h \circ W_\sigma :$
 $\mathcal{S} \rightarrow \Delta(\mathcal{A})$. For π^* fixed, define the **augmented distribu-**
tion $\mathbb{P}_{\text{aug},h}^*$ as the joint distribution over $(s_h^* \sim \mathbb{P}_{\pi^*}^*, a_h^* \sim$
 $\pi_h^*(s_h^*), \tilde{s}_h^* \sim W_\sigma(s_h^*))$, with $a_h^* \perp \tilde{s}_h^* \mid s_h^*$. The **deconvolu-**
tion policy π_{dec}^* is defined by letting $\pi_{\text{dec},h}^*(s)$ denote the
 distribution of $a_h^* \mid \tilde{s}_h^* = s$, where a_h^*, \tilde{s}_h^* are drawn from
 $\mathbb{P}_{\text{aug},h}^*$. Finally, the **replica policy** is $\pi_{\text{O}}^* = \pi_{\text{dec}}^* \circ W_\sigma$.

The operator $\pi \circ W_\sigma$ composes π with the smoothing ker-
 nel. The deconvolution policy π_{dec}^* captures the distribution
 of actions under π^* given an augmented state, and corre-
 sponds to $\pi_{\text{dec}}^* = (\pi_{\text{dec},h}^*)_{h=1}^H$. We argue that if a policy
 $\hat{\pi}$ approximates π_{dec}^* at each step, then $\hat{\pi} \circ W_\sigma$ imitates
 $\pi_{\text{O}}^* = \pi_{\text{dec}}^* \circ W_\sigma$. We explain the “replica policy”, and
 importance of imitating it, after we state our main theorem.
 First, we define a notion of stability to smoothing, taking
 d_{TVC}, d_S, d_A as given.

Definition 4.5. For a non-decreasing maps $\gamma_{\text{IPS},1}, \gamma_{\text{IPS},2} :$
 $\mathbb{R}_{\geq 0} \rightarrow \mathbb{R}_{\geq 0}$ a pseudometric $d_{\text{IPS}} : \mathcal{S} \times \mathcal{S} \rightarrow \mathbb{R}$ (pos-
 sibly other than d_S or d_{TVC}), and $r_{\text{IPS}} > 0$, we say
 a policy π is $(\gamma_{\text{IPS},1}, \gamma_{\text{IPS},2}, d_{\text{IPS}}, r_{\text{IPS}})$ -**input-&-process sta-**
ble (IPS) if the following holds for any $r \in [0, r_{\text{IPS}}]$.
 Consider any sequence of kernels $W_1, \dots, W_H : \mathcal{S} \rightarrow$
 $\Delta(\mathcal{S})$ satisfying $\max_{h,s \in \mathcal{S}} \mathbb{P}_{\tilde{s} \sim W_h(s)}[d_{\text{IPS}}(\tilde{s}, s) \leq r] = 1$,
 and define a process $s_1 \sim P_{\text{init}}, \tilde{s}_h \sim W_h(s_h), a_h \sim$
 $\pi_h(\tilde{s}_h)$, and $s_{h+1} := F_h(s_h, a_h)$. Then, almost surely, (a)
 the sequence $(s_{1:H+1}, a_{1:H})$ is input-stable w.r.t (d_S, d_A)
 (b) $\max_{h \in [H]} d_{\text{TVC}}(F_h(\tilde{s}_h, a_h), s_{h+1}) \leq \gamma_{\text{IPS},1}(r)$ and (c)
 $\max_{h \in [H]} d_S(F_h(\tilde{s}_h, a_h), s_{h+1}) \leq \gamma_{\text{IPS},2}(r)$.

Condition (a) means that the policy $\tilde{\pi}$ defined by $\tilde{\pi}_h = \pi_h \circ$
 W_h is input-stable. In the appendix, we instantiate $W_{1:H}$
 not as W_σ , but as (a truncation of) *replica kernels* $W_{\text{O},h}^*$
 for which $\pi_{\text{O},h}^* = \pi_h^* \circ W_{\text{O},h}^*$. We show that the replica kernel
 inherits any concentration satisfied by W_σ , ensuring (via
 truncation) that $\mathbb{P}_{\tilde{s} \sim W_h(s)}[d_{\text{IPS}}(\tilde{s}, s)] \leq r$. Conditions (b &
 c) merely require that one-step dynamics are robust to small
 changes in state, measured in terms of both d_{TVC} and d_S .

4.1. Instantiation for control

Here we explain the mapping from the control set-
 ting of interest to the composite MDP; in so doing
 we distinguish between the case $h > 1$ and $h =$
 1 with reference to composite-states. In the former
 case, $s_h = (\mathbf{x}_{t_{h-1}:t_h}, \mathbf{u}_{t_{h-1}:t_h-1}) \in \mathcal{P}_{\tau_m}$, and $a_h =$

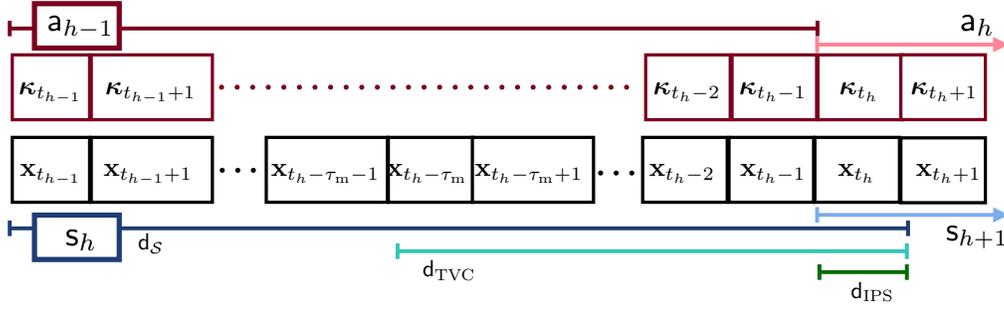


Figure 1. Schematic depicting the composite MDP. States \mathbf{x} and stabilizing gains κ are chunked into composite states \mathbf{s} and composite actions \mathbf{a} (control inputs \mathbf{u} not depicted). The distance labels correspond to the domain over which each distance is evaluated. Note that \mathbf{a}_h begins at the same time that \mathbf{s}_{h+1} does, an indexing convention that we adopt to make the notation in the proofs simpler.

$\kappa_{t_h:t_{h+1}-1}$ (as in Section 2). Importantly, \mathbf{a}_h are **primitive controllers, which allows us to meet the strong stability condition in Definition 4.2**. Figure 1 provides a visual aid for the subtle indexing. For $\mathbf{s}_h, \mathbf{s}'_h$, we define $d_S(\mathbf{s}_h, \mathbf{s}'_h) = \max_{t \in [t_h-1:t_h]} \|\mathbf{x}_t - \mathbf{x}'_t\| \vee \max_{t \in [t_{h-1}:t_h-1]} \|\mathbf{u}_t - \mathbf{u}'_t\|$, which measures distance on the full subtrajectory, $d_{\text{TVC}}(\mathbf{s}_h, \mathbf{s}'_h) = \max_{t \in [t_h-\tau_m:t_h]} \|\mathbf{x}_t - \mathbf{x}'_t\| \vee \max_{t \in [t_h-\tau_m:t_h-1]} \|\mathbf{u}_t - \mathbf{u}'_t\|$, which measures distance on the last τ_m steps, and $d_{\text{IPS}}(\mathbf{s}_h, \mathbf{s}'_h) = \|\mathbf{x}_{t_h} - \mathbf{x}'_{t_h}\|$, which is only on the last step. In the latter case, when $h = 1$, we let $\mathbf{s}_1 = \mathbf{x}_1 \in \mathcal{X}$, and we let $d_S, d_{\text{TVC}}, d_{\text{IPS}}$ all denote the Euclidean distance on \mathcal{X} . In all cases, the transition dynamics F_h are induced by the dynamics (2.1) with $\mathbf{u}_t = \kappa_t(\mathbf{x}_t)$. Finally, for $\mathbf{a} = (\bar{\mathbf{u}}_{1:\tau_c}, \bar{\mathbf{x}}_{1:\tau_c}, \bar{\mathbf{K}}_{1:\tau_c})$ and $\mathbf{a}' = (\bar{\mathbf{u}}'_{1:\tau_c}, \bar{\mathbf{x}}'_{1:\tau_c}, \bar{\mathbf{K}}'_{1:\tau_c})$, we choose a $d_{\mathcal{A}}$ that takes value ∞ when primitive controllers are too far apart as $d_{\mathcal{A}}(\mathbf{a}, \mathbf{a}')$ defined to be

$$c_1 \max_{k \in [\tau_c]} (\|\bar{\mathbf{u}}_k - \bar{\mathbf{u}}'_k\| + \|\bar{\mathbf{x}}_k - \bar{\mathbf{x}}'_k\| + \|\bar{\mathbf{K}}_k - \bar{\mathbf{K}}'_k\|) \quad (4.2)$$

$$+ \mathbf{I}_{0,\infty}\{\mathcal{E}\},$$

where we define $\mathcal{E} := \{\max_{1 \leq k \leq \tau_c} \max\{\|\bar{\mathbf{u}}_k - \bar{\mathbf{u}}'_k\|, \|\bar{\mathbf{x}}_k - \bar{\mathbf{x}}'_k\|, \|\bar{\mathbf{K}}_k - \bar{\mathbf{K}}'_k\|\} \leq c_2\}$, $\mathbf{I}_{0,\infty}$ is the indicator taking infinite value when the event fails to hold, and c_1 and c_2 are constants depending polynomially on all problem parameters, given in Appendix G.

We let the expert policy π^* be the concatenation of policies π_h^* , each of which is defined to be the distribution of \mathbf{a}_h conditioned on $\rho_{\mathbf{m},h}$ under \mathcal{D}_{exp} (see Appendix I for a rigorous definition). As noted above, we take the smoothing kernel W_σ to map $\rho_{\mathbf{m},h}$ to a $\mathcal{N}(\rho_{\mathbf{m},h}, \sigma^2 \mathbf{I}) \in \Delta(\mathcal{P}_{\rho_{\mathbf{m},h}})$, which that same appendix shows is $\frac{1}{2\sigma}$ -TVC (w.r.t. d_{TVC} defined above). We note that under these substitutions, the deconvolution policy $\pi_{\text{dec}}^* = (\pi_{\text{dec},h}^*)_{h=1}^H$ is **precisely as defined in Definition 3.1**.

Appendix G shows that Assumptions 3.1 and 3.2 imply

that π^* enjoys the IPS property in the composite MDP thus instantiated, along with many more granular stability guarantees for time-varying affine feedback in nonlinear control systems, which may be of independent interest.

Proposition 4.1. *Let $c_3, c_4, c_5 > 0$ be as given in Appendix G (and polynomial in relevant quantities). Suppose $\tau_c \geq c_3/\eta$, and let $r_{\text{IPS}} = c_4$, $\gamma_{\text{IPS},1}(u) = c_5 u \exp(-\eta(\tau_c - \tau_m)/L_{\text{stab}})$, $\gamma_{\text{IPS},2}(u) = c_5 u$. Then, for $d_S, d_{\text{TVC}}, d_{\text{IPS}}$ as above, we have that π^* is $(\gamma_{\text{IPS},1}, \gamma_{\text{IPS},2}, d_{\text{IPS}}, r_{\text{IPS}})$ -IPS.*

4.2. A Guarantee in the Composite MDP Stability, and the derivation of Theorem 1

With the substitutions in Section 4.1, it suffices to prove an imitation guarantee in the composite MDP, assuming π^* is IPS, and $\hat{\pi}$ is close to π_{dec}^* in the appropriate sense.

Theorem 2. *Suppose π^* is $(\gamma_{\text{IPS},1}, \gamma_{\text{IPS},2}, d_{\text{IPS}}, r_{\text{IPS}})$ -IPS and W_σ is γ_σ -TVC. Let $\varepsilon > 0$, $r \in (0, \frac{1}{2}r_{\text{IPS}}]$; define $p_r := \sup_{\mathbf{s}' \sim W_\sigma(\mathbf{s})} [d_{\text{IPS}}(\mathbf{s}', \mathbf{s}) > r]$ and $\varepsilon^r := \varepsilon + \gamma_{\text{IPS},2}(2r)$. Then, for any policy $\hat{\pi}$, both $\Gamma_{\text{joint},\varepsilon}(\hat{\pi} \circ W_\sigma \parallel \pi_{\text{dec}}^*)$ and $\Gamma_{\text{marg},\varepsilon^r}(\hat{\pi} \circ W_\sigma \parallel \pi^*)$ are upper bounded by*

$$H(2p_r + 3\gamma_\sigma(\max\{\varepsilon, \gamma_{\text{IPS},1}(2r)\}))$$

$$+ \sum_{h=1}^H \mathbb{E}_{\mathbf{s}_h^* \sim P_h^*} \mathbb{E}_{\mathbf{s}_h \sim W_\sigma(\mathbf{s}_h^*)} d_{\text{os},\varepsilon}(\hat{\pi}_h(\mathbf{s}_h^*) \parallel \pi_{\text{dec}}^*(\mathbf{s}_h^*)).$$

Deriving Theorem 1 from Theorem 2. A full proof is given in Appendix I, using the subtlety that π^* as described above yields trajectories with the same marginals (but possibly different joint distributions) as $\rho_T \sim \mathcal{D}_{\text{exp}}$; thus, to bound losses in Definition 2.2, it suffices to bound the imitation gaps in Definition 4.1 w.r.t. π^* . Using the analysis in Appendix H, we show that our DDPM training precisely ensures that $\hat{\pi}_\sigma = \hat{\pi} \circ W_\sigma$ in TODA minimizes (an upper bound on) the term $\sum_{h=1}^H \mathbb{E}_{\mathbf{s}_h^* \sim P_h^*} \mathbb{E}_{\mathbf{s}_h \sim W_\sigma(\mathbf{s}_h^*)} d_{\text{os},\varepsilon}(\hat{\pi}_h(\mathbf{s}_h^*) \parallel \pi_{\text{dec}}^*(\mathbf{s}_h^*))$. Finally, we combine the guarantees of Proposition 4.1, the aforementioned TVC-bound on W_σ , and Gaus-

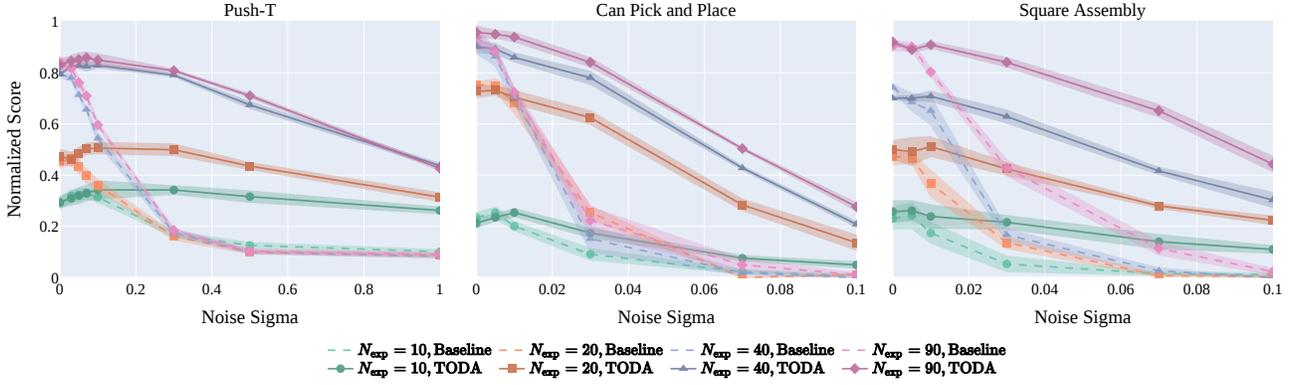


Figure 2. Performance of baseline $\hat{\pi}$ and noise-injected $\hat{\pi} \circ W_\sigma$ TODA policy for different σ . We use 4 training seeds with 50 and 22 test trajectories per seed for PushT and Can and Square Environments respectively. Mean and standard deviation of the test performance on the 3 best checkpoints across the 4 seeds are plotted. The σ values correspond to noise in the normalized $[-1, 1]$ range.

sian concentration to bound p_r with the bound in Theorem 2 to conclude. \square

Proof Sketch of Theorem 2. The proof draws inspiration from the notion of replica symmetry in statistical physics (hence, the name replica) (Mezard & Montanari, 2009). We construct a coupling between a trajectory over (s_h°, a_h°) sampled using the replica policy π_σ^* , and a trajectory (\hat{s}_h, \hat{a}_h) sampled from $\hat{\pi}_\sigma$. We introduce *teleporting* trajectories $s_{h+1}^\circ = F_h(s_h^\circ, a_h^\circ)$, and $s_{h+1}^{\text{tel}} = F_h(\hat{s}_h^{\text{tel}}, a_h^{\text{tel}})$, where \hat{s}_h^{tel} is sampled from the replica distribution of s_h^{tel} and $a_h^{\text{tel}} \sim \pi_h^*(\hat{s}_h^{\text{tel}})$; in words, s_h^{tel} teleports to an independent and identically distributed copy conditional on the noise augmentation, and draws an action from the replica policy evaluated on the new state.

The key fact of the replica distribution is that it preserves marginals, meaning that all s_h^{tel} and \hat{s}_h^{tel} both have marginals according to P_h^* . We show that s_h° tracks the teleporting trajectories, up to the IPS terms $\gamma_{\text{IPS},i}$ and concentration of the kernel, due to total variation continuity of W_σ . Because the marginals of s_h^{tel} are distributed according to P_h^* , we can argue that a (fictitious) action $\hat{a}_h^{\text{tel,inter}} \sim \hat{\pi}_\sigma(s_h^{\text{tel}})$ is close to a_h^{tel} (by the data processing inequality, it is bounded by the closeness of $\hat{\pi}_h$ and $\pi_{\text{dec},h}^*$ on $\hat{s}_h^{\text{tel}} \sim W_\sigma(s_h^{\text{tel}})$, $s_h^{\text{tel}} \sim P_h^*$). We then use total variation continuity to relate to another fictitious action $\hat{a}_h^{\circ,\text{inter}}$ to a_h° . Finally, we use input-stability and TVC again, to relate $\hat{a}_h^{\circ,\text{inter}}$ to actions $\hat{a}_h \sim \hat{\pi}_\sigma(\hat{s}_h)$. Our couplings are summarized in the following diagram:

$$\begin{array}{c}
 \underbrace{(a^\circ \leftrightarrow a^{\text{tel}})}_{\gamma_{\text{TVC}} \text{ and induction}} \rightarrow \underbrace{(a^{\text{tel}} \leftrightarrow \hat{a}^{\text{tel,inter}})}_{\text{learning and sampling}} \\
 \rightarrow \underbrace{(\hat{a}^{\text{tel,inter}} \leftrightarrow \hat{a}^{\circ,\text{inter}})}_{\gamma_{\text{TVC}} \text{ and induction}} \rightarrow \underbrace{(\hat{a}^{\circ,\text{inter}} \leftrightarrow \hat{a})}_{\gamma_{\text{TVC}} \text{ and induction}}.
 \end{array}$$

We construct conditional couplings between pairs of the aforementioned trajectories, each of which corresponds to a

certain optimal transport cost. That past trajectories can be associated to optimal couplings measurably is non-trivial, and proven in Proposition C.3. To conclude, we apply a measure theoretic result (Lemma C.2) to “glue” the pairwise couplings together and establish the main result. The full proof is given in Appendix E, relying on measure-theoretic details in Appendix C. \square

5. Simulation Study of Test-Time Noise-Injection

We empirically evaluate the effect on policy performance of our proposal to inject noise back into the dynamics at inference time. We consider three challenging robotic manipulation tasks studied in prior work: PushT block-pushing (Chi et al., 2023); Robomimic Can Pick-and-Place and Square Nut Assembly (Mandlekar et al., 2021). We explain the environments in greater detail, along with all training and computational details in Appendix K. The learned diffusion policy generates state trajectories over a $\tau_c = 8$ chunking horizon using fixed feedback gains provided by the `synth` oracle to perform position-tracking of the DDPM model output. We direct the reader to Chi et al. (2023) for an extensive empirical investigation into the performance of diffusion policies in the noiseless $\sigma = 0$ setting. We display the results of our experiments in Figure 2. Observe that the performance degradation of the replica policy from the unsmoothed $\sigma = 0$ variant is minimal across all environments and even leads to a slight but noticeable improvement in the small-noise regime for PushT (and low-data Can Pick and Place). In the presence of non-negligible noise TODA significantly outperforms the conventional policy $\hat{\pi}$ (obtained by adding augmentation at training but not test time), as predicted by our theory.

References

- 440 Abu-El-Haija, S., Kothari, N., Lee, J., Natsev, P., Toderici,
441 G., Varadarajan, B., and Vijayanarasimhan, S. Youtube-
442 8m: A large-scale video classification benchmark. *arXiv*
443 *preprint arXiv:1609.08675*, 2016.
- 444 Ajay, A., Gupta, A., Ghosh, D., Levine, S., and Agrawal,
445 P. Distributionally adaptive meta reinforcement learning.
446 *arXiv preprint arXiv:2210.03104*, 2022.
- 447 Altschuler, J. and Talwar, K. Privacy of noisy stochastic
448 gradient descent: More iterations without more privacy
449 loss. *Advances in Neural Information Processing Systems*,
450 35:3788–3800, 2022.
- 451 Altschuler, J. M. and Chewi, S. Faster high-accuracy log-
452 concave sampling via algorithmic warm starts. *arXiv*
453 *preprint arXiv:2302.10249*, 2023.
- 454 Anderson, B. D. and Moore, J. B. *Optimal control: linear*
455 *quadratic methods*. Courier Corporation, 2007.
- 456 Angel, O. and Spinka, Y. Pairwise optimal coupling of mul-
457 tiple random variables. *arXiv preprint arXiv:1903.00632*,
458 2019.
- 459 Bansal, M., Krizhevsky, A., and Ogale, A. Chauffeurnet:
460 Learning to drive by imitating the best and synthesizing
461 the worst. *arXiv preprint arXiv:1812.03079*, 2018.
- 462 Bartlett, P. L., Montanari, A., and Rakhlin, A. Deep learning:
463 a statistical viewpoint. *Acta numerica*, 30:87–201, 2021.
- 464 Block, A., Mroueh, Y., and Rakhlin, A. Generative model-
465 ing with denoising auto-encoders and langevin sampling.
466 *arXiv preprint arXiv:2002.00107*, 2020a.
- 467 Block, A., Mroueh, Y., Rakhlin, A., and Ross, J. Fast mix-
468 ing of multi-scale langevin dynamics under the manifold
469 hypothesis. *arXiv preprint arXiv:2006.11166*, 2020b.
- 470 Bojarski, M., Del Testa, D., Dworakowski, D., Firner, B.,
471 Flepp, B., Goyal, P., Jackel, L. D., Monfort, M., Muller,
472 U., Zhang, J., et al. End to end learning for self-driving
473 cars. *arXiv preprint arXiv:1604.07316*, 2016.
- 474 Chang, J. D., Uehara, M., Sreenivas, D., Kidambi, R., and
475 Sun, W. Mitigating covariate shift in imitation learning
476 via offline data without great coverage. *arXiv preprint*
477 *arXiv:2106.03207*, 2021.
- 478 Chen, S., Chewi, S., Li, J., Li, Y., Salim, A., and Zhang,
479 A. Sampling is as easy as learning the score: theory
480 for diffusion models with minimal data assumptions. In
481 *NeurIPS 2022 Workshop on Score-Based Methods*, 2022.
- 482 Chi, C., Feng, S., Du, Y., Xu, Z., Cousineau, E., Burch-
483 fiel, B., and Song, S. Diffusion policy: Visuomotor
484 policy learning via action diffusion. *arXiv preprint*
485 *arXiv:2303.04137*, 2023.
- 486 Cohen, A., Koren, T., and Mansour, Y. Learning linear-
487 quadratic regulators efficiently with only \sqrt{T} regret. In
488 *International Conference on Machine Learning*, pp. 1300–
489 1309. PMLR, 2019.
- 490 Dasari, S., Ebert, F., Tian, S., Nair, S., Bucher, B.,
491 Schmeckpeper, K., Singh, S., Levine, S., and Finn,
492 C. Robonet: Large-scale multi-robot learning. *arXiv*
493 *preprint arXiv:1910.11215*, 2019.
- 494 De Haan, P., Jayaraman, D., and Levine, S. Causal confu-
sion in imitation learning. *Advances in Neural Informa-
tion Processing Systems*, 32, 2019.
- Durrett, R. *Probability: theory and examples*, volume 49.
Cambridge university press, 2019.
- Finn, C., Yu, T., Zhang, T., Abbeel, P., and Levine, S. One-
shot visual imitation learning via meta-learning. In *Con-
ference on robot learning*, pp. 357–368. PMLR, 2017.
- Hagood, J. W. and Thomson, B. S. Recovering a func-
tion from a dini derivative. *The American Mathematical*
Monthly, 113(1):34–46, 2006.
- Havens, A. and Hu, B. On imitation learning of linear
control policies: Enforcing stability and robustness con-
straints via lmi conditions. In *2021 American Control*
Conference (ACC), pp. 882–887. IEEE, 2021.
- Hendrycks, D., Basart, S., Mu, N., Kadavath, S., Wang, F.,
Dorundo, E., Desai, R., Zhu, T., Parajuli, S., Guo, M.,
et al. Jacob steinhardt et justin gilmer. the many faces
of robustness: A critical analysis of out-of-distribution
generalization. *arXiv preprint arXiv:2006.16241*, 2020.
- Ho, J., Jain, A., and Abbeel, P. Denoising diffusion proba-
bilistic models. *Advances in Neural Information Process-
ing Systems*, 33:6840–6851, 2020.
- Hussein, A., Gaber, M. M., Elyan, E., and Jayne, C. Im-
itation learning: A survey of learning methods. *ACM*
Computing Surveys (CSUR), 50(2):1–35, 2017.
- Hussein, A., Elyan, E., Gaber, M. M., and Jayne, C. Deep
imitation learning for 3d navigation tasks. *Neural com-
puting and applications*, 29:389–404, 2018.
- Hyvärinen, A. and Dayan, P. Estimation of non-normalized
statistical models by score matching. *Journal of Machine*
Learning Research, 6(4), 2005.

- 495 Janner, M., Du, Y., Tenenbaum, J. B., and Levine, S. Plan-
496 ning with diffusion for flexible behavior synthesis. *arXiv*
497 *preprint arXiv:2205.09991*, 2022.
- 498 Jin, C., Netrapalli, P., Ge, R., Kakade, S. M., and Jordan,
499 M. I. A short note on concentration inequalities for
500 random vectors with subgaussian norm. *arXiv preprint*
501 *arXiv:1902.03736*, 2019.
- 503 Ke, L., Wang, J., Bhattacharjee, T., Boots, B., and Srinivasa,
504 S. Grasping with chopsticks: Combating covariate shift
505 in model-free imitation learning for fine manipulation.
506 In *2021 IEEE International Conference on Robotics and*
507 *Automation (ICRA)*, pp. 6185–6191. IEEE, 2021.
- 509 Kelly, M., Sidrane, C., Driggs-Campbell, K., and Kochen-
510 derfer, M. J. Hg-dagger: Interactive imitation learning
511 with human experts. In *2019 International Conference on*
512 *Robotics and Automation (ICRA)*, pp. 8077–8083. IEEE,
513 2019.
- 514 Kostrikov, I., Yarats, D., and Fergus, R. Image augmentation
515 is all you need: Regularizing deep reinforcement learning
516 from pixels. *arXiv preprint arXiv:2004.13649*, 2020.
- 518 Laskey, M., Lee, J., Fox, R., Dragan, A., and Goldberg,
519 K. Dart: Noise injection for robust imitation learning.
520 In *Conference on robot learning*, pp. 143–156. PMLR,
521 2017.
- 522 Lee, H., Lu, J., and Tan, Y. Convergence of score-based
523 generative modeling for general data distributions. In *In-*
524 *ternational Conference on Algorithmic Learning Theory*,
525 pp. 946–985. PMLR, 2023.
- 527 Mandlekar, A., Xu, D., Wong, J., Nasiriany, S., Wang,
528 C., Kulkarni, R., Fei-Fei, L., Savarese, S., Zhu, Y., and
529 Martín-Martín, R. What matters in learning from offline
530 human demonstrations for robot manipulation. *arXiv*
531 *preprint arXiv:2108.03298*, 2021.
- 533 Maurer, A. A vector-contraction inequality for rademacher
534 complexities. In *Algorithmic Learning Theory: 27th*
535 *International Conference, ALT 2016, Bari, Italy, October*
536 *19-21, 2016, Proceedings 27*, pp. 3–17. Springer, 2016.
- 537 Mezard, M. and Montanari, A. *Information, physics, and*
538 *computation*. Oxford University Press, 2009.
- 540 Misra, D. Mish: A self regularized non-monotonic activa-
541 tion function. *arXiv preprint arXiv:1908.08681*, 2019.
- 542 Nichol, A. Q. and Dhariwal, P. Improved denoising diffusion
543 probabilistic models. In *International Conference on*
544 *Machine Learning*, pp. 8162–8171. PMLR, 2021.
- 546 Pearce, T., Rashid, T., Kanervisto, A., Bignell, D., Sun,
547 M., Georgescu, R., Macua, S. V., Tan, S. Z., Momenne-
548 jad, I., Hofmann, K., et al. Imitating human behaviour
549 with diffusion models. *arXiv preprint arXiv:2301.10677*,
2023.
- Perez, E., Strub, F., De Vries, H., Dumoulin, V., and
Courville, A. Film: Visual reasoning with a general con-
ditioning layer. In *Proceedings of the AAAI Conference*
on Artificial Intelligence, volume 32, 2018.
- Pfrommer, D., Zhang, T., Tu, S., and Matni, N. Tasil: Taylor
series imitation learning. *Advances in Neural Information*
Processing Systems, 35:20162–20174, 2022.
- Pfrommer, D., Simchowitz, M., Westenbroek, T., Matni,
N., and Tu, S. The power of learned locally linear mod-
els for nonlinear policy optimization. *arXiv preprint*
arXiv:2305.09619, 2023.
- Polyanskiy, Y. and Wu, Y. *Information Theory: From Cod-*
ing to Learning. Cambridge University Press, 2022+.
- Raginsky, M., Rakhlin, A., and Telgarsky, M. Non-convex
learning via stochastic gradient langevin dynamics: a
nonasymptotic analysis. In *Conference on Learning The-*
ory, pp. 1674–1703. PMLR, 2017.
- Rakhlin, A., Sridharan, K., and Tsybakov, A. B. Empirical
entropy, minimax regret and minimax risk. *Bernoulli*
Society for Mathematical Statistics and Probability, 2017.
- Ramesh, A., Dhariwal, P., Nichol, A., Chu, C., and Chen,
M. Hierarchical text-conditional image generation with
clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- Ronneberger, O., Fischer, P., and Brox, T. U-net: Con-
volutional networks for biomedical image segmentation.
In *Medical Image Computing and Computer-Assisted*
Intervention–MICCAI 2015: 18th International Confer-
ence, Munich, Germany, October 5-9, 2015, Proceedings,
Part III 18, pp. 234–241. Springer, 2015.
- Ross, S. and Bagnell, D. Efficient reductions for imitation
learning. In *Proceedings of the thirteenth international*
conference on artificial intelligence and statistics, pp.
661–668. JMLR Workshop and Conference Proceedings,
2010.
- Ross, S., Gordon, G., and Bagnell, D. A reduction of imita-
tion learning and structured prediction to no-regret online
learning. In *Proceedings of the fourteenth international*
conference on artificial intelligence and statistics, pp.
627–635. JMLR Workshop and Conference Proceedings,
2011.
- Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., and
Ganguli, S. Deep unsupervised learning using nonequi-
librium thermodynamics. In *International Conference on*
Machine Learning, pp. 2256–2265. PMLR, 2015.

-
- 550 Song, J., Meng, C., and Ermon, S. Denoising diffusion im-
551 plicit models. *arXiv preprint arXiv:2010.02502*, 2020a.
552
- 553 Song, Y. and Ermon, S. Generative modeling by estimating
554 gradients of the data distribution. *Advances in neural*
555 *information processing systems*, 32, 2019.
- 556 Song, Y., Garg, S., Shi, J., and Ermon, S. Sliced score
557 matching: A scalable approach to density and score es-
558 timation. In *Uncertainty in Artificial Intelligence*, pp.
559 574–584. PMLR, 2020b.
- 560
- 561 Sun, X., Yang, S., and Mangharam, R. Mega-dagger: Im-
562 itation learning with multiple imperfect experts. *arXiv*
563 *preprint arXiv:2303.00638*, 2023.
- 564
- 565 Thoppilan, R., De Freitas, D., Hall, J., Shazeer, N., Kul-
566 shreshtha, A., Cheng, H.-T., Jin, A., Bos, T., Baker, L.,
567 Du, Y., et al. Lamda: Language models for dialog appli-
568 cations. *arXiv preprint arXiv:2201.08239*, 2022.
- 569
- 570 Tu, S., Robey, A., Zhang, T., and Matni, N. On the sample
571 complexity of stability constrained imitation learning. In
572 *Learning for Dynamics and Control Conference*, pp. 180–
573 191. PMLR, 2022.
- 574
- 575 Van Handel, R. Probability in high dimension. Technical
576 report, PRINCETON UNIV NJ, 2014.
- 577
- 578 Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones,
579 L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. At-
580 tention is all you need. *Advances in neural information*
581 *processing systems*, 30, 2017.
- 582
- 583 Vershynin, R. *High-dimensional probability: An introduc-*
584 *tion with applications in data science*, volume 47. Cam-
585 bridge university press, 2018.
- 586
- 587 Villani, C. *Topics in optimal transportation*, volume 58.
588 American Mathematical Soc., 2021.
- 589
- 590 Villani, C. et al. *Optimal transport: old and new*, volume
591 338. Springer, 2009.
- 592
- 593 Vincent, P. A connection between score matching and de-
594 noising autoencoders. *Neural computation*, 23(7):1661–
595 1674, 2011.
- 596
- 597 Wainwright, M. J. *High-dimensional statistics: A non-*
598 *asymptotic viewpoint*, volume 48. Cambridge university
599 press, 2019.
- 600
- 601 Westenbroek, T., Simchowitz, M., Jordan, M. I., and Sas-
602 try, S. S. On the stability of nonlinear receding horizon
603 control: a geometric perspective. In *2021 60th IEEE Con-*
604 *ference on Decision and Control (CDC)*, pp. 742–749.
IEEE, 2021.
- Wu, Y. and He, K. Group normalization. In *Proceedings of
the European conference on computer vision (ECCV)*, pp.
3–19, 2018.
- Zhang, T., McCarthy, Z., Jow, O., Lee, D., Chen, X., Gold-
berg, K., and Abbeel, P. Deep imitation learning for com-
plex manipulation tasks from virtual reality teleoperation.
In *2018 IEEE International Conference on Robotics and
Automation (ICRA)*, pp. 5628–5635. IEEE, 2018.

605	Contents	
606		
607	1 Introduction	1
608		
609	2 Setting	2
610		
611		
612	3 Algorithm and Results	4
613		
614	4 Analysis	5
615		
616	4.1 Instantiation for control	6
617		
618	4.2 A Guarantee in the Composite MDP Stability, and the derivation of Theorem 1	7
619		
620	5 Simulation Study of Test-Time Noise-Injection	8
621		
622	A Notation, Organization of Appendix, and Full Related Work	15
623		
624	A.1 Notation Summary	15
625		
626	A.2 Organization of the Appendix	15
627		
628	B Complete Related Work	16
629		
630	B.1 Comparison to prior notions of Stability.	17
631		
632		
633	I Composite MDP	18
634		
635	C Measure-Theoretic Background	18
636		
637	C.1 Kernels, Regular Conditional Probabilities and Gluing	20
638		
639	C.2 Optimal Transport and Kernel Couplings	21
640		
641	C.3 Data Processing Inequalities	22
642	C.3.1 Deferred lemmas for the data processing inequalities	23
643		
644	C.4 Proof of Proposition C.3	23
645		
646	C.5 A simple union-bound recursion.	25
647		
648	D Warmup: Analysis Without Augmentation	25
649		
650	E Imitation in the Composite MDP	27
651		
652	E.1 A generalization of Theorem 2	28
653		
654	E.2 Proof of Theorem 4	30
655	E.2.1 Proof Overview and Coupling Construction	30
656	E.2.2 Properties of smoothing, deconvolution, and replicas.	32
657	E.2.3 Formal proof of Theorem 4	34
658	E.2.4 Proof of Lemma E.9	37
659		

660	E.3	Proof of Theorem 2, and generalization to direct decompositions	39
661			
662	F	Lower Bounds	40
663			
664	F.1	Sharpness of Proposition D.1 and Theorem 2	41
665	F.2	$\pi_{\circ\sigma}^*$ and π^* induce the same marginals but different joint distributions, even with memoryless dynamics .	43
666			
667	F.3	$\pi_{\circ\sigma}^*$ and π^* can have different marginals, implying necessity of $\gamma_{\text{IPS},2}$	44
668	F.4	$\pi_{\circ\sigma}^*$ and π_{dec}^* have different marginals, even with memoryless dynamics	45
669			
670			
671	II	The Control Setting	47
672			
673			
674	G	Stability in the Control System	47
675			
676	G.1	Recalling preliminaries and assumptions.	48
677			
678	G.1.1	Properties satisfied by π^*	49
679	G.1.2	Norm notation.	49
680	G.2	Composite Problem Constants	50
681	G.3	IPS Guarantees & Proof of Proposition 4.1	50
682			
683	G.3.1	A more granular stability statement	52
684	G.3.2	Deriving Corollary G.1 from Proposition G.3	54
685	G.4	Stability guarantees for single control (sub-)trajectories.	55
686	G.5	Deriving Proposition G.3 from Proposition G.5	56
687			
688	G.5.1	Interpreting the error terms.	57
689	G.5.2	An intermediate guarantee.	58
690	G.5.3	Concluding the proof of Proposition G.3.	59
691	G.6	Proof of Lemma G.4 (state perturbation)	61
692	G.7	Proof of Proposition G.5 (input and gain perturbation)	62
693			
694	G.7.1	Deferred Claims	65
695	G.8	Ricatti synthesis of stabilizing gains.	68
696			
697	G.8.1	Proof of Proposition G.12 (Ricatti synthesis of gains)	68
698	G.9	Solutions to recursions	70
699			
700			
701			
702			
703	H	Sampling and Score Matching	73
704			
705	H.1	Denoising Diffusion Probabilistic Models	75
706	H.2	Score Estimation	77
707			
708			
709	I	End-to-end Guarantees and the Proof of Theorem 1	80
710			
711	I.1	Preliminaries	81
712			
713	I.1.1	Preliminaries for joint-distribution imitation.	82
714	I.1.2	Translating Control Imitation Losses to Composite-MDP Imitation Gaps	82

715	I.2	Proof of Theorem 1 and a general reduction	83
716	I.3	Imitation of the joint trajectory under total variation continuity of demonstrator policy	86
717			
718	I.3.1	Proof of Theorem 9	87
719	I.3.2	Proof of Proposition I.2	87
720			
721	I.4	Imitation in total variation distance	89
722	I.5	Imitation with no augmentation	89
723			
724	I.6	Consequence for expected costs	90
725	I.7	Useful Lemmata	91
726			
727	I.7.1	On the trajectories induced by π^* from \mathcal{D}_{exp}	91
728	I.7.2	Concentration and TVC of Gaussian Smoothing.	92
729	I.7.3	Total Variation Telescoping	93
730			
731			
732	J	Extensions and Further Results	93
733			
734	J.1	Noisy Dynamics	93
735	J.2	Robustness to Adversarial Perturbations	94
736	J.3	Deconvolution Policies and Total Variation Continuity	94
737			
738			
739	K	Experiment Details	97
740			
741	K.1	Compute and Codebase Details	97
742	K.2	Environment Details	97
743	K.3	DDPM Model and Training Details.	98
744			
745			
746			
747			
748			
749			
750			
751			
752			
753			
754			
755			
756			
757			
758			
759			
760			
761			
762			
763			
764			
765			
766			
767			
768			
769			

770 A. Notation, Organization of Appendix, and Full Related Work

771 In this appendix, we collect the notation we use throughout the paper, as well as providing a high level organization of the
772 appendices.
773

774 A.1. Notation Summary

775 In this section, we summarize some of the notation used throughout the work, divided by subject.
776

777 **Measure Theory** We always let \mathcal{X} denote a Polish space, $\mathcal{B}(\mathcal{X})$ the Borel-algebra on \mathcal{X} , and $\Delta(\mathcal{X})$ the set of borel
778 probability measures on \mathcal{X} . For a random variable X on \mathcal{X} , we let P_X denote the law of X . For random variables X, Y ,
779 we let $\mathcal{C}(P_X, P_Y)$ denote the set of couplings of these measures and for laws P_1, P_2 . We write $P_1 \otimes P_2$ for the product
780 measure. We will generally reserve P to denote measure, Q and W for probability kernels, and μ for a joint measure on
781 several random variables.
782

783 When $P_1, P_2 \in \Delta(\mathcal{X})$ are laws on the same space, we let $\text{TV}(P_1, P_2)$ denote the total variation distance. We write
784 $P_1 \ll P_2$ if P_1 is absolutely continuous with respect to P_2 . Given a Polish space \mathcal{X} and element $x \in \mathcal{X}$, we let $\delta_x \in \Delta(\mathcal{X})$
785 denote the dirac-delta measure supported on the set $\{x\} \in \mathcal{B}(\mathcal{X})$ (note that, in a Polish space, the singleton $\{x\}$ set is
786 closed, and therefore Borel).
787

788 **Norms and linear algebra notation.** We use bold lower case vector \mathbf{z} to denote vectors, and bold upper case \mathbf{Z} to denote
789 matrices. We let $\mathbf{z}_{1:K} = (\mathbf{z}_1, \dots, \mathbf{z}_K)$ and $\mathbf{Z}_{1:K} = (\mathbf{Z}_1, \dots, \mathbf{Z}_K)$ denote concatenations. The norms $\|\cdot\|$ denote Euclidean
790 norms on vectors and operator norms on matrices. We identify the spaces \mathcal{P}_k with Euclidean vectors in the standard sense.
791 Given a Euclidean vector $\mathbf{z} \in \mathbb{R}^d$, $\mathcal{N}(\mathbf{z}, \sigma^2 \mathbf{I})$ denote the multivariate normal distribution on \mathbb{R}^d with covariance $\sigma^2 \mathbf{I}$.
792

793 **Control notation.** We let $\mathbf{x}_t \in \mathbb{R}^{d_x}$ denote control states, $\mathbf{u}_t \in \mathbb{R}^{d_u}$ denote control inputs, and $\rho_\tau \in \mathcal{P}_\tau$ denotes
794 trajectories $(\mathbf{x}_{1:\tau+1}, \mathbf{u}_{1:\tau})$. T denotes the time horizon of imitation, so $\rho_T \sim \mathcal{P}_T$. Our dynamics are $\mathbf{x}_{t+1} = f(\mathbf{x}_t, \mathbf{u}_t)$;
795 for our main results (Section 3), we suppose $f(\mathbf{x}, \mathbf{u}) = \mathbf{x} + \eta f_\eta(\mathbf{x}, \mathbf{u})$, parametrizing dynamics in the form of an Euler
796 discretization with step $\eta > 0$.
797

798 Recall that primitive controllers κ take the form $\kappa(\mathbf{x}) = \bar{\mathbf{K}}(\mathbf{x} - \bar{\mathbf{x}}) + \bar{\mathbf{u}}$, where terms with $(\bar{\cdot})$, $\bar{\mathbf{K}}$, $\bar{\mathbf{x}}$, $\bar{\mathbf{u}}$, denote parameters
799 of the primitive controller. The space of these is \mathcal{K} .
800

801 We also recall the chunk-length τ_c and memory length τ_m satisfying $0 \leq \tau_m \leq \tau_c$. We recall the definition of the trajectory-
802 chunk $\rho_{c,h}$ and memory-chunk in $\rho_{m,h}$ in Section 2, which introduced the indexing h , such that $t_h = (h-1)\tau_c + 1$. Recall
803 also the composite actions $\mathbf{a}_h = (\kappa_{t_h:t_{h+1}-1}) \in \mathcal{A} = \mathcal{K}^{\tau_c}$ as the concatenation of τ_c primitive controllers.
804

805 **Abstractions in the composite MDP.** The composite MDP is a deterministic MDP with composite-states $\mathbf{s} \in \mathcal{S}$ and
806 composite-actions $\mathbf{a} \in \mathcal{A}$, and (possibly time-varying) deterministic transition dynamics $F_h : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$ for $1 \leq h \leq H$.
807 The goal is to imitate a policy $\pi^* = (\pi_h^*)_{1 \leq h \leq H}$, in terms of imitation gaps $\Gamma_{\text{joint}, \varepsilon}$ and $\Gamma_{\text{marg}, \varepsilon}$ defined in Definition 4.1.
808 We refer the reader to Section 4 for the relevant terminology, and to Section 4.1 for its instantiation in our original control
809 setting.
810

811 A.2. Organization of the Appendix

812 We now describe the organization of our many appendices. In Appendix B, we expand on our abbreviated discussion of
813 related work in the body as well as provide a more detailed comparison of our notion of stability Definition 4.5 with those
814 found in prior work.
815

816 After the preliminaries on organization, notation, and related work, we divide our appendices into two parts. In the first
817 part, we expand on and provide rigorous proofs of statements and results pertaining to the composite MDP as considered in
818 Section 4. We begin by providing a detailed background in Appendix C on the requisite measure theory we use to make
819 our arguments rigorous. In particular, we provide definitions of probability kernels and couplings, as well as measurability
820 properties of optimal transport couplings. In Appendix D, we provide a warmup to the proof of Theorem 2. In particular,
821 the argument substantially simplifies if we consider the case of no added augmentation (when $\sigma = 0$ in TODA) and we
822 present a coupling construction that implies the analogous bound in the presence of an additional assumption. The heart of
823 the first part of our appendices is Appendix E, where we rigorously prove a generalization of Theorem 2 by constructing a
824

sophisticated coupling between the imitator and demonstrator trajectories. We conclude the first part of our appendices by proving a number of lower bounds in the composite MDP setting in [Appendix F](#), which demonstrate the tightness of our arguments in [Appendix E](#).

We continue our appendices in the second part, which is concerned with the instantiation of the composite MDP in the control setting of interest. In [Appendix G](#), we provide a detailed proof that the control setting considered in [Section 2](#) satisfies the stability properties required by our analysis of the composite MDP and prove [Proposition 4.1](#). Of particular note are [Definition G.7](#), which provide explicit dependence of the relevant constants in [Theorem 1](#) on the parameters of interest, and [Appendix G.8](#), which explains how to synthesize stabilizing gains, as assumed in [Section 2](#). With the stability properties thus proven, we proceed in [Appendix H](#) to instantiate our conditional sampling guarantees with DDPMs. In particular, by applying earlier work, we state and prove [Theorem 6](#), which guarantees that with sufficiently many samples, in our setting we can ensure that the learned DDPM provides samples close in the relevant optimal transport distance to the expert distribution. We also explain in [Remark H.5](#) why stronger total variation guarantees on sampling are unrealistic in our setting. The heart of the second part of our appendices is [Appendix I](#), which provides the final, end-to-end guarantees and the proof of [Theorem 1](#). In that section, we prove a reduction from imitation learning to conditional sampling and derive [Theorem 1](#) as a corollary. We also provide a number of variations on this result, including stronger guarantees on imitation of the joint trajectories ([Appendix I.3](#)), guarantees on TODA under the assumption that sampling is close in total variation ([Appendix I.4](#)), and imitation with no augmentation ([Appendix I.5](#)). We also show in [Proposition I.5](#) that most natural cost functions have similar expected values on imitator and demonstrator trajectories assuming that the imitation losses are small.

We provide a number of extensions of our main results in [Appendix J](#), including to the setting of noisy dynamics ([Appendix J.1](#)). Finally, in [Appendix K](#), we expand the discussion of our experiments, including training and compute details, environment details, and a link to our code for the purpose of reproducibility.

B. Complete Related Work

Imitation Learning. Over the past few years, there has been a significant surge of interest in utilizing machine learning techniques for the execution of exceedingly intricate manipulation and control tasks. Imitation learning, whereby a policy is trained to mimic expert demonstrations, has emerged as a highly data efficient and effective method in this domain, with application to self-driving vehicles ([Hussein et al., 2017](#); [Bojarski et al., 2016](#); [Bansal et al., 2018](#)), visuomotor policies ([Finn et al., 2017](#); [Zhang et al., 2018](#)), and navigation tasks ([Hussein et al., 2018](#)). A widely acknowledged challenge of imitation learning is distribution shift: since the training and test time distributions are induced by the expert and trained policies respectively, compounding errors in imitating the expert at test-time can lead the trained policy to explore out-of-distribution states ([Ross & Bagnell, 2010](#)). This distribution shift has been shown to result in the imitator making incorrect judgements regarding observation-action causality, often with catastrophic consequences ([De Haan et al., 2019](#)). Prior work in this domain has predominantly attempted to mitigate this issue in the non-stochastic setting via online data augmentation strategies, sampling new trajectories to mitigate distribution shift ([Ross et al., 2011](#); [Ross & Bagnell, 2010](#); [Laskey et al., 2017](#)). Among this class of methods, the DAGger algorithm in particular has seen widespread adoption ([Ross & Bagnell, 2010](#); [Sun et al., 2023](#); [Kelly et al., 2019](#)). These approaches have the drawback that sampling new trajectories or performing queries on the expert is often expensive or intractable. Due to these limitations, recent developments have focused on novel algorithms and theoretical guarantees for imitation learning in an offline, non-interactive environment ([Chang et al., 2021](#); [Pfrommer et al., 2022](#)). Our work is similarly focused on the offline setting, but is capable of handling stochastic, non-Markovian demonstrators. Unlike ([Pfrommer et al., 2022](#)), we do not require our expert demonstrations to be sampled from a stabilizing expert policy, instead utilizing a synthesis oracle to stabilize around the provided demonstrations. This is a significantly weaker requirement and enables the development of high-probability guarantees for human demonstrators, where sampling new trajectories and reasoning about the stability properties is not possible.

Denoising Diffusion Probabilistic Models. Denoising Diffusion Probabilistic Models (DDPMs) ([Sohl-Dickstein et al., 2015](#); [Ho et al., 2020](#)) and their variant, Annealed Langevin Sampling ([Song & Ermon, 2019](#)), have seen enormous empirical success in recent years, especially in state-of-the-art image generation ([Ramesh et al., 2022](#); [Nichol & Dhariwal, 2021](#); [Song et al., 2020a](#)). More relevant to this paper is their application to imitation learning, where they have seen success even without the proposed data augmentation in [Janner et al. \(2022\)](#); [Chi et al. \(2023\)](#); [Pearce et al. \(2023\)](#). DDPMs rely on learning the score function of the target distribution, which is generally accomplished through some kind of denoised estimation ([Hyvärinen & Dayan, 2005](#); [Vincent, 2011](#); [Song et al., 2020b](#)). On the theoretical end, annealed Langevin sampling has been studied with score estimators under a variety of assumptions including the manifold hypothesis and some

form of dissapitivity (Raginsky et al., 2017; Block et al., 2020a;b), although these works have generally suffered from an exponential dependence on ambient dimension, which is unacceptable in our setting. Of greatest relevance to the present paper are the concurrent works of Chen et al. (2022); Lee et al. (2023) that provide polynomial guarantees on the quality of sampling using a DDPM assuming that the score functions are close in an appropriate mean squared error sense. We take advantage of these latter two works in order to provide concrete end-to-end bounds in our setting of interest. To our knowledge, ours is the first work to consider the application of DDPMs to imitation learning under a rigorous theoretical framework, although we emphasize that this does not constitute a strong technical contribution as opposed to an instantiation of earlier work for the sake of completeness and concreteness.

Smoothing Augmentations. Data augmentation with smoothing noise has become such common practice, its adoption is essentially folklore. While augmentation of actions which noise is common practice for exploration (see, e.g. (Laskey et al., 2017)), it is widely accepted that noising actions in the learned policy is not best practice, and thus it is more common to add noise to the *states* at training time, preserving target actions as fixed (Ke et al., 2021). Our work gives an interpretation of this decision as enforcing that the learned policy obey the distributional continuity property we term TVC (Definition 4.3), so that the policy selects similar actions on nearby states. Previous work has interpreted noise augmentation as providing robustness. Data augmentation has been demonstrated to provide more robustness in RL from pixels (Kostrikov et al., 2020), adaptive meta-learning (Ajay et al., 2022), in more traditional supervised learning as well (Hendrycks et al., 2020).

B.1. Comparison to prior notions of Stability.

Prior work in guarantees for imitation learning focuses either on constraining the learned policy to be stable (Havens & Hu, 2021; Tu et al., 2022) or assume the expert policy is suitably stable (Pfrommer et al., 2022).

The principal notion of stability used in these prior works is *incremental-input-to-state* stability of the closed-loop system under a deterministic controller π :

Definition B.1 (Incremental Input-to-State Stability). There exists class \mathcal{K} function γ and class \mathcal{KL} function β such that for any two initial conditions $\xi_1, \xi_2 \in \mathcal{X}$, the closed-loop dynamics under policy $\pi : \mathcal{X} \rightarrow \mathcal{U}$ given by $f_{\text{cl}}(x_t, \Delta_t) = f(x_t, \pi(x_t) + \Delta_t)$ satisfies:

$$\|x_t(\xi_1; \{\Delta_s\}_{s=0}^t) - x_t(\xi_2; \{0\}_{s=0}^t)\| \leq \beta(\|\xi_1 - \xi_2\|) + \gamma\left(\max_{0 \leq s \leq t-1} \|\Delta_s\|\right),$$

where $x_t(\xi; \{\Delta_s\}_{s=0}^{t-1})$ is the state at time t under f_{cl} with $x_0 = \xi$ and input perturbations $\{\Delta_s\}_{s=0}^{t-1}$.

This notion of stability is quite restrictive, as the β -term necessitates that the dynamics converge irrespective of initial condition. Without time-varying dynamics this can only be achieved by a policy which stabilizes to an equilibrium point, as a policy which tracks a reference trajectory is unable to “forget” the initial condition. Constraining learned policies such that they satisfy this notion of stability is also challenging. Tu et. al. (Tu et al., 2022) attempt to do so through regularization while Haven et. a. (Havens & Hu, 2021) use matrix inequalities to satisfy this stability property under linear dynamics. Pfrommer et. at. (Pfrommer et al., 2022) avoid this difficulty by relaxing the incremental stability to a local variant of stability:

Definition B.2 (η -Local Incremental Input-to-State Stability). There exists class \mathcal{K} function γ such that for any $\xi \in \mathcal{X}$, the closed-loop dynamics under policy $\pi : \mathcal{X} \rightarrow \mathcal{U}$ given by $f_{\text{cl}}(x_t, \Delta_t) = f(x_t, \pi(x_t) + \Delta_t)$ satisfies:

$$\|x_t(\xi; \{\Delta_s\}_{s=0}^t) - x_t(\xi; \{0\}_{s=0}^t)\| \leq \gamma\left(\max_{0 \leq s \leq t-1} \|\Delta_s\|\right),$$

for all $\{\Delta_s\}_{s=0}^t$ where $\max_{0 \leq s \leq t} \|\Delta_s\| \leq \eta$.

This weaker notion of incremental stability simply postulates the existence of a (local) input-perturbation to state-perturbation gain function γ . Since this stability property does not necessitate convergence across with different initial conditions and only under input perturbations of magnitude $\leq \eta$, this only necessitates that the expert policy can correct from small input perturbations.

We further weaken this assumption, which we formalize in Assumption 3.2 and abstract to the composite MDP through Definition G.4, by only requiring that a locally stabilizing controller can be synthesized per-demonstration. Through the introduction of a synthesis oracle which can generate locally stabilizing primitive controllers, we decouple the stability

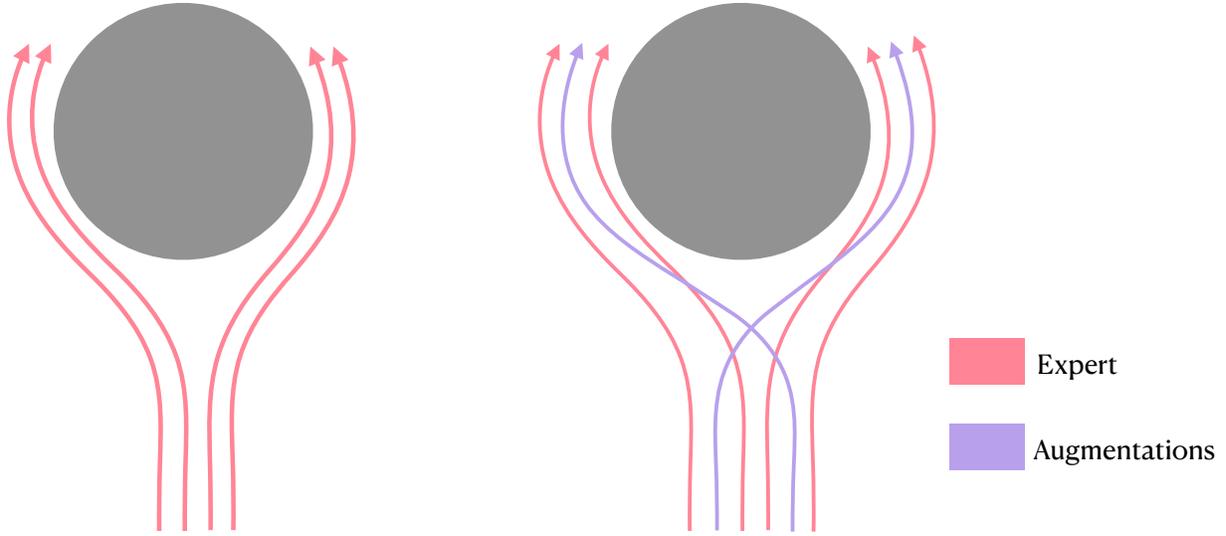


Figure 3. Instance of bifurcation, where augmentation is necessary for stability. The example on the left has an expert demonstrator bifurcating around a circular obstacle. The example on the right demonstrates the utility of augmentations, allowing for trajectories that navigate around the object in the direction farther from their starting point.

properties of the expert from the stabilizability of the underlying dynamical system. This allows for reasoning about generalization in the presence of bifurcations or conflicting demonstrations, which is precluded by Definition B.2 since an expert policy cannot simultaneously stabilize to multiple branches of a bifurcation. For a concrete example, consider Figure 3. Indeed, continuity is the *sine qua non* of stability and the example given demonstrates the necessity of augmentation to enforce the former. In detail, the figure illustrates an example where an agent is navigating around an obstacle, providing a bifurcation. Without augmentation, the demonstrator trajectories always navigate around the obstacle in the direction closer to their starting point, leading to a sharp discontinuity along a bisector of the obstacle. On the other hand, the data augmentations allow for the policy to have some probability of navigating around the obstacle in the “wrong” direction, which leads to the notion of continuity we consider: total variation continuity.

Because our notion of stability is applied in chunks, our theory is sufficiently flexible so as to allow for the learned policy to switch between expert demonstrations in a manner preserving the marginal distributions but not consistent with the joint distribution across the entire trajectory. This flexibility is illustrated in Figure 4, where we suppose that the demonstrator distribution consists both of trajectories traversing a figure “8” consistently in either a clockwise or counter-clockwise manner, with both orientations represented in the data set. Due to the multi-modality at the critical point in the trajectory, there is ambiguity about which loop to traverse next; specifically, there may exist a policy that randomly select which loop to traverse each time the critical point is visited in such a way that the marginal distributions on states and actions is the same as that induced by the demonstrator. Such a policy will, by definition, preserve the correct *marginal* distributions across states and actions; at the same time, this policy has a different *joint* distribution across all time steps from the demonstrator due to the possibility of traversing the same loop twice in a row.

Part I

Composite MDP

C. Measure-Theoretic Background

In this section, we introduce the prerequisite notions from probability theory that we use to formally construct the couplings in Appendices D and E. We begin by introducing general preliminaries, followed by kernels, regular conditional probabilities

990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044

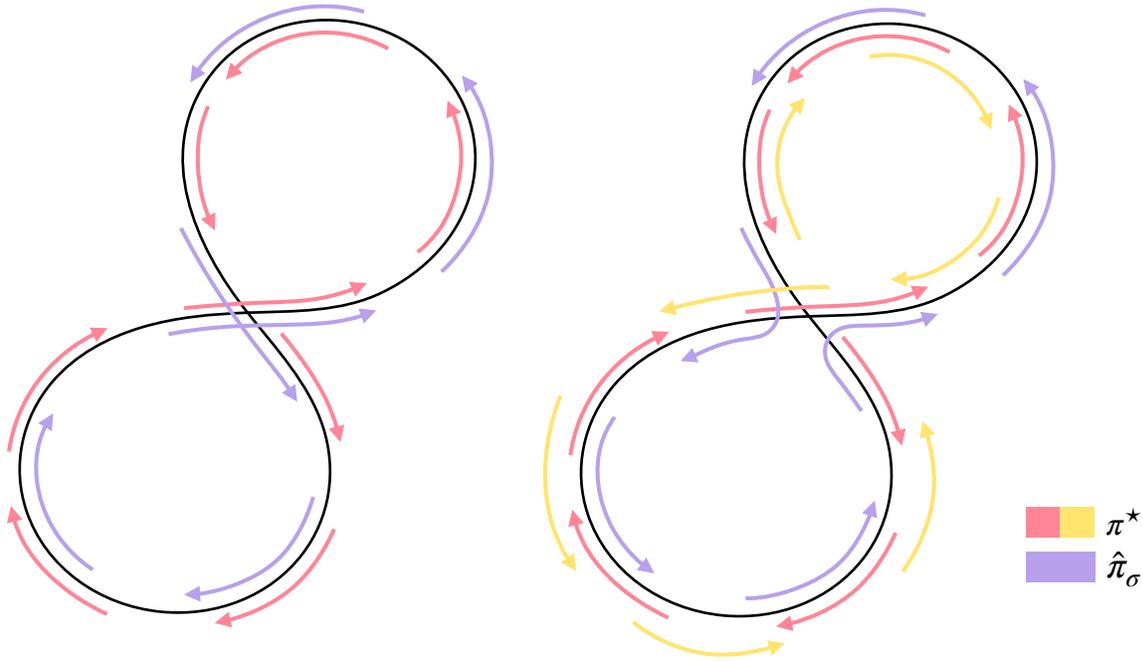


Figure 4. Instance where $\hat{\pi}_\sigma$ and π^* induce the same marginals and joint distributions (left), but in the presence of expert demonstration trajectories that traverse the figure eight both clockwise and counterclockwise directions, $\hat{\pi}_\sigma$ may switch with some probability between demonstrations where they overlap.

and a “gluing” lemma in Appendix C.1. We then show that optimal transport costs commute in an appropriate sense with conditional probabilities (Proposition C.3 in Appendix C.2). We use the preliminaries in the previous sections to derive certain optimal-transport and data processing inequalities in Appendix C.3. We prove Proposition C.3 in Appendix C.4. Finally, we state a simple union bound lemma (Lemma C.11 in Appendix C.5) of use in later appendices.

General preliminaries. We rely extensively on the exposition in Durrett (2019) and refer the reader there for a more thorough introduction. Throughout, we assume there is a Polish space Ω such that all random variables of interest are mappings $X : \Omega \rightarrow \mathcal{X}$, where \mathcal{X} is also Polish. Here, the σ -algebras are always the Borel algebras (the σ -algebra generated by open subsets), denoted $\mathcal{B}(\Omega)$ and $\mathcal{B}(\mathcal{X})$.

The space of (Borel) probability distributions on \mathcal{X} is denoted $\Delta(\mathcal{X})$, and measurability is meant in the Borel sense. Given a measure μ on a space $\mathcal{X} \times \mathcal{Y}$, we say that $X \sim P_X$ under μ if, for all $A \in \mathcal{B}(\mathcal{X})$, $\mu(X \in A) = P_X(A)$.

We adopt standard information theoretic notation to denote joint, marginal, and conditional distributions on vectors of random variables. In particular, if random variables X, Y are distributed according to P , we denote by P_X as the marginal over X , $P_{X|Y}$ as the conditional of $X|Y$ under P , and $P_{X,Y}$ as the joint distribution when this needs to be emphasized.

Definition C.1 (Couplings). Let \mathcal{X}, \mathcal{Y} be Polish spaces and let $P_X \in \Delta(\mathcal{X})$ and $P_Y \in \Delta(\mathcal{Y})$. The set of couplings $\mathcal{C}(P_X, P_Y)$ denotes the set of measure $\mu \in \Delta(\mathcal{X} \times \mathcal{Y})$ such that, $(X, Y) \sim \mu$ has marginals $X \sim P_X$ and $Y \sim P_Y$.² We let $P_X \otimes P_Y \in \mathcal{C}(P_X, P_Y)$ denote the *independent coupling* under which X and Y are independent.

It is standard that $P_X \otimes P_Y$ is always a valid coupling, and hence $\mathcal{C}(P_X, P_Y)$ is nonempty. Couplings have the advantage that they can be used to design many probability-theoretic distances. Through the paper, we use the total variation distance.

Definition C.2 (Total Variation Distance). Let $P_1, P_2 \in \Delta(\mathcal{X})$. We define the total variation distance $\text{TV}(P_1, P_2) := \sup_{A \in \mathcal{B}(\mathcal{X})} |P_1(A) - P_2(A)|$

The total variation distance can be expressed in terms of couplings as follows (Polyanskiy & Wu, 2022+).

²More pedantically, for all Borel sets $A_1 \in \mathcal{B}(\mathcal{X})$, $\mu(A_1 \times \mathcal{Y}) = P_X(A_1)$ all Borel sets $A_2 \in \mathcal{B}(\mathcal{Y})$, $\mu(\mathcal{X} \times A_2) = P_Y(A_2)$.

1045 **Lemma C.1.** Let $P_1, P_2 \in \Delta(\mathcal{X})$. Then,

1046
1047
$$\text{TV}(P_1, P_2) = \inf_{\mu \in \mathcal{C}(P_1, P_2)} \mathbb{P}_{(X_1, X_2) \sim \mu} \{X_1 \neq X_2\}.$$

1048
1049 Moreover, there exists a coupling μ_* attaining the infimum.

1051 **Support and absolute continuity.** We will also require the definition of the support of a measure.

1052 **Definition C.3.** Given a measure μ on a Borel space (Ω, \mathcal{F}) , we define the *support* $\text{supp}(\mu)$ to be the closure in the topology
1053 given by the metric of the set $\{\omega \in \Omega \mid \mu(U) > 0 \text{ for all open } U \ni \omega\}$.

1054
1055 In addition, we require the definition of absolute continuity.

1056 **Definition C.4 (Absolute Continuity).** We say that $P \in \Delta(\mathcal{X})$ is absolutely continuous with respect to law $P' \in \Delta(\mathcal{X})$,
1057 written $P \ll P'$, if for $A \in \mathcal{B}(\mathcal{X})$, $P'(A) = 0$ implies $P(A) = 0$.

1058
1059 We now go into greater detail on the kinds of couplings that we consider.

1061 C.1. Kernels, Regular Conditional Probabilities and Gluing

1063 One key technical challenge in proving results in the sequel is the fact that we need to “glue” together multiple different
1064 couplings. Specifically, while it may be the case that there exist pairwise couplings which satisfy desired properties, there
1065 exists a coupling such that the probability of the relevant event is small, it is not obvious that there exists a *single* coupling
1066 such that all of these probabilities are small *simultaneously*. There are two natural ways to do this gluing: the first, using
1067 regular conditional probabilities we provide here. The second, involving a sophisticated construction of [Angel & Spinka](#)
1068 (2019) requires stronger assumptions on the pseudo-metric, but generalizing beyond Polish spaces, we simply remark can be
1069 substituted with a loss of a constant factor.

1071 **Kernels.** We begin by introducing the notion of a kernel.

1072 **Definition C.5 (Kernels).** Let (Ω, \mathbb{P}) be a probability space and let X denote a random variable on this space. For a given
1073 σ -algebra \mathcal{G} , and map $Q : \Omega \times \mathcal{G} \rightarrow [0, 1]$, we say that Q is a probability kernel if the following two conditions are satisfied:

- 1074
1075
1076 1. For all measurable events A , the map $\omega \mapsto Q(\omega, A)$ is measurable.
1077
1078 2. For almost every $\omega \in \Omega$, the map $A \mapsto Q(\omega, A)$ is a probability measure.

1079
1080 We can combine a probability kernel with a probability measure on \mathcal{Y} to yield joint distributions over $\mathcal{X} \times \mathcal{Y}$.

1081 **Definition C.6.** Given an $P_Y \in \Delta(\mathcal{Y})$, we define the probability measure $\text{law}(Q_{X|Y}; P_Y) \in \Delta(\mathcal{X} \times \mathcal{Y})$ such that
1082 $\mu = \text{law}(Q_{X|Y}; P_Y)$ satisfies³

1083
1084
$$\mu(A \times B) = \mathbb{E}_{Y \sim P_Y} [Q_{X|Y}(A | Y) \mathbf{I}\{Y \in B\}], \quad \forall A \in \mathcal{B}(\mathcal{X}), B \in \mathcal{B}(\mathcal{Y}). \quad (\text{C.1})$$

1085
1086 We let $Q_{X|Y} \circ P_Y \in \Delta(\mathcal{X})$ denote the measure for which $\mu = Q_{X|Y} \circ P_Y$ satisfies

1087
1088
$$\mu(A) = \mathbb{E}_{Y \sim P_Y} [Q_{X|Y}(A | Y)], \quad \forall A \in \mathcal{B}(\mathcal{X})$$

1089
1090 From these, we define the space of conditional couplings as follows.

1091 **Definition C.7 (Kernel Couplings).** Let $P_Y \in \Delta(\mathcal{Y})$, and $Q_{X_i|Y} \in \Delta(\mathcal{X} | \mathcal{Y})$ for $i \in \{1, 2\}$. We let $\mathcal{C}_{P_Y}(Q_{X_1|Y}, Q_{X_2|Y})$
1092 denote the space of measures $\mu \in \Delta(\mathcal{X}_1 \times \mathcal{X}_2 \times \mathcal{Y})$ over random variables (X_1, X_2, Y) such that $(X_i, Y) \sim \text{law}(Q_{X_i|Y}; P_Y)$
1093 for $i \in \{1, 2\}$.

1094
1095 Note that a similar construction to the independent coupling ensures $\mathcal{C}_{P_Y}(Q_{X_1|Y}, Q_{X_2|Y})$ is nonempty, namely considering
1096 the measure $\mu(A_1 \times A_2 \times B_2) = \mathbb{E}_{Y \sim P_Y} [Q_{X_1|Y}(A_1 | Y) Q_{X_2|Y}(A_2) \mathbf{I}\{Y \in B\}]$.

1097
1098 ³Recall that $\mathcal{B}(\mathcal{X} \times \mathcal{Y})$ is generated by sets $A \times B \in \mathcal{B}(\mathcal{X}) \times \mathcal{B}(\mathcal{Y})$, so (C.1) defines a unique probability measure
1099

1100 **Regular Conditional Probabilities.** We now recall a standard result that conditional probabilities can be expressed
 1101 through kernels in our setting.

1102 **Theorem 3** (Theorem 5.1.9, [Durrett \(2019\)](#)). *If Ω is a Polish space and \mathbb{P} is a probability measure on the Borel sets of Ω ,
 1103 such that random variables $(X, Y) \sim \mathbb{P}$ in spaces \mathcal{X} and \mathcal{Y} , then there exists a kernel $Q(\cdot | \cdot) \in \Delta(\mathcal{X} | \mathcal{Y})$ such that, for all
 1104 $A \in \mathcal{B}(\mathcal{X})$ and \mathbb{P} -almost every y , the (standard) conditional probability $\mathbb{P}[X \in A | Y] = Q(A | y)$. We can $Q(\cdot | \cdot)$ the
 1105 regular conditional probability measure.*

1107 Regular conditional probabilities allow one to think of conditional probabilities in the most intuitive way, i.e., for two
 1108 random variables X, Y , the map $Y \mapsto \mathbb{P}(X \in A | Y)$ is a probability kernel. This will be the essential property that we use
 1109 below.

1111 **Gluing.** Finally, regular conditional probabilities allow us to “glue together” couplings which share a common random
 1112 variable.

1113 **Lemma C.2** (Gluing Lemma). *Suppose that X, Y, Z are random variables taking value in Polish spaces $\mathcal{X}, \mathcal{Y}, \mathcal{Z}$. Let $\mu_1 \in$
 1114 $\Delta(\mathcal{X} \times \mathcal{Y}), \mu_2 \in \Delta(\mathcal{Y} \times \mathcal{Z})$ be couplings of (X, Y) and (Y, Z) respectively. Then there exists a coupling $\mu \in \Delta(\mathcal{X} \times \mathcal{Y} \times \mathcal{Z})$
 1115 on (X, Y, Z) such that under μ , $(X, Y) \sim \mu_1$ and $(Y, Z) \sim \mu_2$.*

1117 *Proof.* Let $Q(\cdot | Y)$ be a regular conditional probability for Z given Y under μ_2 (who existence is ensured by [Theorem 3](#)).

1119 We construct μ by first sampling $(X, Y) \sim \mu_1$ and then sampling $Z \sim Q(\cdot | Y)$; observe that by the second property in
 1120 [Definition C.5](#), this is a valid construction. It is immediate that under μ , we have $(X, Y) \sim \mu_1$ and thus we must only show
 1121 that $(Y, Z) \sim \mu_2$ to conclude the proof. Let A, B be two measurable sets and we see that

$$\begin{aligned} 1122 \mathbb{P}_\mu((Y, Z) \in A \times B) &= \mathbb{E}_{Y \sim \mu} [\mathbb{P}_\mu((Y, Z) \in A \times B | Y)] \\ 1123 &= \mathbb{E}_{Y \sim \mu} [\mathbb{E}_{(Y, Z) \sim \mu} [\mathbf{I}[Y \in A] \cdot \mathbf{I}[Z \in B] | Y]] \\ 1124 &= \mathbb{E}_{Y \sim \mu} [\mathbf{I}[Y \in A] \cdot \mathbb{E}_\mu[\mathbf{I}[Z \in B] | Y]] \\ 1125 &= \mathbb{E}_{Y \sim \mu} [\mathbf{I}[Y \in A] \cdot \mathbb{P}_{\mu_2}(Z \in B | Y)] \\ 1126 &= \mathbb{E}_{Y \sim \mu} [\mathbf{I}[Y \in A] \cdot \mathbb{P}_{\mu_2}(Z \in B | Y)] \\ 1127 &= \mu_2((Y, Z) \in A \times B), \end{aligned}$$

1129 where the first equality follows from the tower property of expectations, the second follows by definition of conditional
 1130 probability, the third follows from the definition of conditional expectation, the fourth follows by the first property from
 1131 [Definition C.5](#), and the last follows from the fact that the marginals of Y under μ and under μ_2 are the same. The result
 1132 follows. \square

1134 C.2. Optimal Transport and Kernel Couplings

1136 As shown above for the TV distance, many measures of distributional distance can be quantified in terms of *optimal transport*
 1137 costs; these are quantities expressed as infima, over all couplings, of the expectation of a certain lower-semicontinuous
 1138 functions. We show that if the optimal transport costs between two kernels $Y \rightarrow \Delta(\mathcal{X}_i)$ are controlled pointwise, then for
 1139 any $P_Y \in \Delta(\mathcal{Y})$, is a there exists a joint distribution over (X_1, X_2, Y) which attains the minimal transport cost.

1140 **Proposition C.3.** *Let $\mathcal{X}_1, \mathcal{X}_2, \mathcal{Y}$ be Polish spaces, and let $P_Y \in \Delta(\mathcal{Y})$, and $Q_i \in \Delta(\mathcal{X}_i | \mathcal{Y})$. for $i \in \{1, 2\}$. Finally, let
 1141 $\phi : \mathcal{X}_1 \times \mathcal{X}_2 \rightarrow \mathbb{R}$ be lower semicontinuous and bounded below. Then, the following function*

$$1143 \psi(y) := \inf_{\mu \in \mathcal{C}(Q_1(y), Q_2(y))} \mathbb{E}_{(X_1, X_2) \sim \mu} [\phi(X_1, X_2)]$$

1145 *is a measurable function of y and there exists some $\mu_\star \in \mathcal{C}_{P_Y}(Q_1, Q_2)$ such that*

$$1147 \mathbb{E}_{(X_1, X_2, Y) \sim \mu_\star} [\phi(X_1, X_2)] = \mathbb{E}_{Y \sim P_Y} \psi(Y).$$

1148 *In particular it holds μ_\star -almost surely that*

$$1150 \mathbb{E}_{\mu_\star} [\phi(X_1, X_2) | Y] = \psi(Y).$$

1152 We prove the above proposition in [Appendix C.4](#). One useful consequence is the following identity for the total variation
 1153 distance.

1154

1155 **Corollary C.1.** Let \mathcal{X}, \mathcal{Y} be Polish spaces, and let $P_Y \in \Delta(\mathcal{Y})$, and $Q_i \in \Delta(\mathcal{X} | \mathcal{Y})$, for $i \in \{1, 2\}$. Then, there exists a
 1156 coupling $\mu_* \in \mathcal{C}_{P_Y}(Q_1, Q_2)$ such that

$$1157 \mathbb{P}_{\mu_*}[X_1 \neq X_2] = \mathbb{E}_{Y \sim P_Y} \text{TV}(Q_1(\cdot | Y), Q_2(\cdot | Y)),$$

1158
 1159 with the left-hand side integrand being measurable.

1160
 1161 *Proof.* Using [Lemma C.1](#), we can represent total variation as an optimal transport cost with $\phi(x_1, x_2) = \mathbf{I}\{x_1 \neq x_2\}$. Note
 1162 that $\phi(x_1, x_2)$ is lower semicontinuous, being the indicator of an open set. Thus, the result follows from [Proposition C.3](#)
 1163 with $\mathcal{X} = \mathcal{X}_1 = \mathcal{X}_2$, and $\phi(x_1, x_2) = \mathbf{I}\{x_1 \neq x_2\}$. \square

1164 C.3. Data Processing Inequalities

1165 We now derive two *inequalities*. First, we recall the classical version for the total variation distance, and check that a
 1166 well-known identity holds in our setting.

1167 **Lemma C.4** (Data Processing for Total Variation). Let $P_{Y_1}, P_{Y_2} \in \Delta(\mathcal{Y})$ and let $Q_X \in \Delta(\mathcal{X} | \mathcal{Y})$. Then,

$$1168 \text{TV}(Q_X \circ P_{Y_1}, Q_X \circ P_{Y_2}) \leq \text{TV}(\text{law}(Q_X; P_{Y_1}), \text{law}(Q_X; P_{Y_2})) = \text{TV}(P_{Y_1}, P_{Y_2}).$$

1169
 1170 *Proof.* The first inequality is just the data processing inequality ([Polyanskiy & Wu, 2022+](#), Theorem 7.7), which also shows
 1171 that $\text{TV}(\text{law}(Q_X; P_{Y_1}), \text{law}(Q_X; P_{Y_2})) \geq \text{TV}(P_{Y_1}, P_{Y_2})$. To prove the reverse inequality, we use [Lemma C.1](#) to find a
 1172 coupling μ_Y such that (P_{Y_1}, P_{Y_2}) such that $\mathbb{E}[\mathbf{I}\{Y_1 \neq Y_2\}] = \text{TV}(P_{Y_1}, P_{Y_2})$.

1173 Define a probability kernel in $\Delta(\mathcal{X} \times \mathcal{X} | \mathcal{Y}_1 \times \mathcal{Y}_2)$ via defining the set $B_ = \{(x_1, x_2) \in \mathcal{X} \times \mathcal{X} : x_1 = x_2\} \subset \mathcal{X} \times \mathcal{X}$, and
 1174 define for $A \in \mathcal{B}(\mathcal{X} \times \mathcal{X})$,

$$1175 \mathbb{Q}(A | y_1, y_2) = \begin{cases} \mathbb{Q}_X(\pi_1(A \cap B_ =) | y_1) & y_1 = y_2 \\ \mathbb{Q}_X(\cdot | y_1) \otimes \mathbb{Q}_X(\cdot | y_2)(A) & \text{otherwise} \end{cases}$$

1176 In a Polish space, [Lemmas C.6](#) and [C.7](#) imply that $A \mapsto \mathbb{Q}_X(\pi_1(A \cap B_ =) | y_1)$ for each y_1 is a valid measure, and it is
 1177 standard that the product measures $\mathbb{Q}_X(\cdot | y_1) \otimes \mathbb{Q}_X(\cdot | y_2)(A)$ are valid. Moreover, this construction ensures that for
 1178 $\mu = \text{law}(\mathbb{Q}; \mu_Y)$,

$$1179 \mathbb{P}_\mu[\{Y_1 = Y_2\} \text{ and } \{X_1 \neq X_2\}] = 0. \tag{C.2}$$

1180
 1181 Lastly, one can check that under $\mu = \text{law}(\mathbb{Q}; \mu_Y)$, that $(X_1, Y_1) \sim \text{law}(\mathbb{Q}_X; P_{Y_1})$ and $(X_2, Y_2) \sim \text{law}(\mathbb{Q}_X; P_{Y_2})$. Thus, μ
 1182 can be regarded as an element of $\mathcal{C}(\text{law}(\mathbb{Q}_X; P_{Y_1}), \text{law}(\mathbb{Q}_X; P_{Y_2}))$. Hence, [Lemma C.1](#) implies that

$$1183 \begin{aligned} \text{TV}(\text{law}(\mathbb{Q}_X; P_{Y_1}), \text{law}(\mathbb{Q}_X; P_{Y_2})) &\leq \text{TV}(\mathbb{P}_\mu[(X_1, Y_1) \neq (X_2, Y_2)]) \\ &= \mathbb{P}_\mu[Y_1 \neq Y_2] + \mathbb{P}_\mu[\{Y_1 = Y_2\} \text{ and } \{X_1 \neq X_2\}] \\ &= \mathbb{P}_{\mu_*}[Y_1 \neq Y_2] && \text{(Eq.(C.2))} \\ &= \mathbb{P}_{(Y_1, Y_2) \sim \mu_Y}[Y_1 \neq Y_2] \\ &= \text{TV}(P_{Y_1}, P_{Y_2}). && \text{(construction of } \mu_Y) \end{aligned}$$

1184
 1185 \square

1186 Next, we derive a general data processing inequality for optimal costs. This result is a corollary of [Proposition C.3](#).

1187 **Lemma C.5** (Another Data Processing Inequality for Optimal Transport). Let $\mathcal{X}_1, \mathcal{X}_2, \mathcal{Y}$ be Polish spaces, and let $P_Y \in$
 1188 $\Delta(\mathcal{Y})$, and $Q_i \in \Delta(\mathcal{Y} | \mathcal{X}_i)$. for $i \in \{1, 2\}$. Denote by $Q_i \circ P_Y$ the marginal of X_i under $(X_i, Y) \sim \text{law}(Q_i; P_Y)$. Then,

$$1189 \inf_{\mu \in \mathcal{C}(Q_1 \circ P_Y, Q_2 \circ P_Y)} \mathbb{E}_{X_1, X_2 \sim \mu} \phi(X_1, X_2) \leq \mathbb{E}_{Y \sim P_Y} \left(\inf_{\mu' \in \mathcal{C}(Q_1(Y) \circ Q_2(Y))} \mathbb{E}_{X_1, X_2 \sim \mu'} \phi(X_1, X_2) \right).$$

1190
 1191

1210 *Proof.* One can check that any coupling in $\mu \in \mathcal{C}(\mathbb{Q}_1 \circ P_Y, \mathbb{Q}_2 \circ P_Y)$ can be obtained by marginalizing Y in a certain
 1211 coupling of $\mu' \in \mathcal{C}(\text{law}(\mathbb{Q}_1; P_Y), \text{law}(\mathbb{Q}_2; P_Y))$, and any coupling in the latter can be marginalized to a coupling in the
 1212 former. Hence,

$$1214 \quad \inf_{\mu \in \mathcal{C}(\mathbb{Q}_1 \circ P_Y, \mathbb{Q}_2 \circ P_Y)} \mathbb{E}_{X_1, X_2 \sim \mu} \phi(X_1, X_2) = \inf_{\mu \in \mathcal{C}(\text{law}(\mathbb{Q}_1; P_Y), \text{law}(\mathbb{Q}_2; P_Y))} \mathbb{E}_{X_1, X_2, Y_1, Y_2 \sim \mu} \phi(X_1, X_2)$$

1216 Moreover, to every measure $\mu \in \mu_{P_Y}(\mathbb{Q}_1, \mathbb{Q}_2)$ over (X_1, X_2, Y) , [Lemma C.8](#) implies that there exists a coupling $\mu' \in$
 1217 $\mathcal{C}(\text{law}(\mathbb{Q}_1; P_Y), \text{law}(\mathbb{Q}_2; P_Y))$ over (X_1, X_2, Y_1, Y_2) such (X_1, X_2) have the same marginals under μ and μ' . Therefore,

$$1219 \quad \inf_{\mu \in \mathcal{C}(\text{law}(\mathbb{Q}_1; P_Y), \text{law}(\mathbb{Q}_2; P_Y))} \mathbb{E}_{X_1, X_2, Y_1, Y_2 \sim \mu} \phi(X_1, X_2) \leq \inf_{\mu' \in \mathcal{C}_{P_Y}(\mathbb{Q}_1, \mathbb{Q}_2)} \mathbb{E}_{X_1, X_2, Y \sim \mu'} \phi(X_1, X_2).$$

1222 Finally, the right hand side is equal to $\mathbb{E}_{Y \sim \mu_Y} (\inf_{\mu' \in \mathcal{C}(\mathbb{Q}_1(Y) \circ Q_2(Y))} \mathbb{E}_{X_1, X_2 \sim \mu'} \phi(X_1, X_2))$ by [Proposition C.3](#). \square

1224 C.3.1. DEFERRED LEMMAS FOR THE DATA PROCESSING INEQUALITIES

1226 **Lemma C.6.** *Let \mathcal{X} be a Polish space. Then, the set $\{(x_1, x_2) \in \mathcal{X} \times \mathcal{X} : x_1 \neq x_2\}$ is open in $\mathcal{X} \times \mathcal{X}$.*

1228 *Proof.* The diagonal is closed in any Polish space by definition of the topology. The result follows. \square

1230 **Lemma C.7.** *Let \mathcal{X} be a Polish space, and let $\pi_1, \pi_2 : \mathcal{X} \times \mathcal{X}$ denote the projection mappings onto each coordinate. Then,
 1231 for any $A \in \mathcal{B}(\mathcal{X} \times \mathcal{X})$, $\pi_1(A)$ and $\pi_2(A)$ are in $\mathcal{B}(\mathcal{X})$.*

1234 *Proof.* The projection map is open so the result follows immediately by definition of the Borel algebra. \square

1236 **Lemma C.8.** *Let \mathcal{X}, \mathcal{Y} be Polish spaces, and let $\mu \in \Delta(\mathcal{X} \times \mathcal{Y})$. Then, there is a measure $\mu' \in \Delta(\mathcal{X} \times \mathcal{Y} \times \mathcal{Y})$ satisfying*

$$1238 \quad \mu'(A \times \mathcal{Y}) = \mu(A), \quad \forall A \in \mathcal{B}(\mathcal{X} \times \mathcal{Y})$$

1240 and

$$1242 \quad \mu'(\mathcal{X} \times \{(y_1, y_2) : y_1 = y_2\}) = 1$$

1244 *Proof.* Define the set $B_{=} = \{(y_1, y_2) : y_1 = y_2\}$. One can check that $\mu'(A \times B) = \mu(A \times \pi_1(B \cap B_{=}))$, where
 1245 $\pi_1 : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathcal{Y}$ is the projection onto the first coordinate, is a valid measure. \square

1247 C.4. Proof of Proposition C.3

1249 In the case that $\phi(\cdot, \cdot)$ is continuous, the result follows from [Villani et al. \(2009, Corollary 5.22\)](#). For general lower-
 1250 semicontinuous ϕ , our argument adopts the strategy of ‘‘Step 3’’ of the proof of [Villani \(2021, Theorem 1.3\)](#). This shows that
 1251 there exists a sequence $\phi_n \uparrow \phi$ pointwise, such that each ϕ_n is uniformly bounded. Define

$$1253 \quad \psi_n(y) := \inf_{\mu \in \mathcal{C}(\mathbb{Q}_1(y), \mathbb{Q}_2(y))} \mathbb{E}_{(X_1, X_2) \sim \mu} [\phi_n(X_1, X_2)].$$

1255 Then, for each n , the continuous case implies that there exists a measure $\mu_{*,n} \in \mathcal{C}_{P_Y}(\mathbb{Q}_1, \mathbb{Q}_2)$ such that

$$1257 \quad \mathbb{E}_{Y \sim \nu_Y} \psi_n(Y) = \mathbb{E}_{(X_1, X_2, Y) \sim \mu_{*,n}} [\phi_n(X_1, X_2)] \tag{C.3}$$

1259 Recall now the definition

$$1261 \quad \psi(y) = \inf_{\mu \in \mathcal{C}(\mathbb{Q}_1(y), \mathbb{Q}_2(y))} \mathbb{E}_{(X_1, X_2) \sim \mu} [\phi(X_1, X_2)].$$

1263 **Claim C.9.** *$\psi(y)$ is measurable and satisfies $\psi_n(y) \uparrow \psi(y)$ pointwise.*

1264

1265 *Proof.* We can write

$$\begin{aligned}
1266 & \\
1267 & \sup_{n \geq 0} \psi_n(y) = \sup_{n \geq 0} \inf_{\mu \in \mathcal{C}(\mathcal{Q}_1(y), \mathcal{Q}_2(y))} \mathbb{E}_{(X_1, X_2) \sim \mu} [\phi_n(X_1, X_2)] \\
1268 & \\
1269 & \stackrel{(i)}{=} \inf_{\mu \in \mathcal{C}(\mathcal{Q}_1(y), \mathcal{Q}_2(y))} \mathbb{E}_{(X_1, X_2) \sim \mu} [\phi(X_1, X_2)] = \psi(y). \\
1270 & \\
1271 &
\end{aligned}$$

1272 Here, (i) follows from the ‘‘Step 3’’ in the proof of [Villani \(2021, Theorem 1.3\)](#), which shows that any optimal transport
1273 cost C of a lowersemicontinuous ϕ is equal to a limit of the costs C_n of any bounded continuous $\phi_n \uparrow \phi$. In our case, we
1274 fix each y , so $C = \psi(y)$ and $C_n = \psi_n(y)$. It is clear that $\psi_n(y)$ is increasing, so for each y , $\psi_n(y) \uparrow \psi(y)$. As ψ is the
1275 pointwise monotone limit of ψ_n , it is measurable. \square

1276 **Claim C.10.** *The set of couplings of $\mathcal{C}_{P_Y}(X_1, X_2)$ is compact in the weak topology.*

1277 *Proof.* Recall that $\Delta(\mathcal{Y} \times \mathcal{X}_1 \times \mathcal{X}_2)$ denote the set of Borel measures on $\mathcal{Y} \times \mathcal{X}_1 \times \mathcal{X}_2$. This set is also a Polish space in the
1278 weak topology. The subset $\mathcal{C}_{P_Y}(X_1, X_2) \subset \Delta(\mathcal{Y} \times \mathcal{X}_1 \times \mathcal{X}_2)$ is compact if and only if it is relatively compact and closed.

1281 To show relative compactness, Prokhorov’s theorem means that it suffices to show that $\mu_{P_Y}(\mathcal{Q}_1, \mathcal{Q}_2)$ is tight, i.e. for all
1282 $\varepsilon > 0$, there exists a compact $\mathcal{K}_\varepsilon \subset \mathcal{Y} \times \mathcal{X}_1 \times \mathcal{X}_2$ such that for any $\mu \in \mathcal{C}_{P_Y}(X_1, X_2)$, $\mathbb{P}_\mu[(Y, X_1, X_2) \in \mathcal{K}_\varepsilon] \geq 1 - \varepsilon$.
1283 This follows by setting $\mathcal{K} = \mathcal{K}_{Y, \varepsilon} \times \mathcal{K}_{X_1, \varepsilon} \times \mathcal{K}_{X_2, \varepsilon}$, where the sets are such that $\mathbb{P}_{P_Y}[Y \notin \mathcal{K}_{Y, \varepsilon}] \geq 1 - \varepsilon/3$ and
1284 $\mathbb{P}_{\mathcal{Q}_i}[X_i \notin \mathcal{K}_{X_i, \varepsilon}] \geq 1 - \varepsilon/3$, where \mathcal{Q}_i is the marginal of X_i given by $Y \sim P_Y$, $X_i \sim P_i(\cdot | Y)$ (such sets exist because
1285 $\mathcal{X}_1, \mathcal{X}_2, \mathcal{Y}$ are Polish).

1286 To check that $\mathcal{C}_{P_Y}(\mathcal{Q}_1, \mathcal{Q}_2) \subset \Delta(\mathcal{Y} \times \mathcal{X}_1 \times \mathcal{X}_2)$ is closed, it suffices to show that it is sequentially closed (as $\Delta(\mathcal{Y} \times \mathcal{X}_1 \times \mathcal{X}_2)$
1287 is Polish). To this end, consider any sequence $\mu_n \in \mathcal{C}_{P_Y}(\mathcal{Q}_1, \mathcal{Q}_2)$ such that $\mu_n \xrightarrow{\text{weak}} \mu \in \Delta(\mathcal{Y} \times \mathcal{X}_1 \times \mathcal{X}_2)$ in the weak
1288 topology. By definition, this means that for any $i \in \{1, 2\}$ and any continuous and bounded $f_i : \mathcal{Y} \times \mathcal{X}_i \rightarrow \mathbb{R}$,

$$1291 \lim_{n \rightarrow \infty} \mathbb{E}_{\mu_n} f_i(Y, X_i) = \mathbb{E}_\mu f_i(Y, X_i).$$

1293 For all $\mu_n \in \mathcal{C}_{P_Y}(\mathcal{Q}_1, \mathcal{Q}_2)$, $\mathbb{E}_{\mu_n} f_i(Y, X_i) = \mathbb{E}_{Y \sim \nu_Y} \mathbb{E}_{X_i \sim \nu_i(\cdot | Y_i)} f_i(Y, X_i)$. Thus,

$$1295 \mathbb{E}_\mu f_i(Y, X_i) = \mathbb{E}_{Y \sim \nu_Y} \mathbb{E}_{X_i \sim \nu_i(\cdot | Y_i)} f_i(Y, X_i), \quad \text{for all continuous, bounded } f_i : \mathcal{Y} \times \mathcal{X}_i \rightarrow \mathbb{R}.$$

1297 Hence, the marginal distribution of (Y, X_i) under μ must be equal to that of $(Y \sim P_Y, X_i \sim \mathcal{Q}_i(\cdot | Y))$ for $i \in \{1, 2\}$,
1298 which means $\mu \in \mathcal{C}_{P_Y}(\mathcal{Q}_1, \mathcal{Q}_2)$. \square

1300 By compactness, there exists (passing to a subsequence if necessary) a $\mu_\star \in \mathcal{C}_{P_Y}(\mathcal{Q}_1, \mathcal{Q}_2)$ such that $\mu_{\star, n} \xrightarrow{\text{weak}} \mu_\star$ in the
1301 weak topology. Then, as ϕ_m is continuous and bounded, it follows that for all m ,

$$\begin{aligned}
1303 & \\
1304 & \mathbb{E}_{(X_1, X_2, Y) \sim \mu_\star} [\phi_m(X_1, X_2)] = \limsup_{n \rightarrow \infty} \mathbb{E}_{(X_1, X_2, Y) \sim \mu_{\star, n}} [\phi_m(X_1, X_2)] && (\mu_{\star, n} \xrightarrow{\text{weak}} \mu_\star) \\
1305 & \leq \limsup_{n \rightarrow \infty} \mathbb{E}_{(X_1, X_2, Y) \sim \mu_{\star, n}} [\phi_n(X_1, X_2)] && (\phi_m \leq \phi_n \text{ for } n \geq m) \\
1306 & = \limsup_{n \rightarrow \infty} \mathbb{E}_Y \psi_n(Y) && ((\text{C.3})) \\
1307 & = \mathbb{E}_Y \lim_{n \rightarrow \infty} \psi_n(Y) && (\text{Monotone Convergence}) \\
1308 & = \mathbb{E}_Y \psi(Y). && (\text{Claim C.9})
\end{aligned}$$

1312 Thus, by the monotone convergence theorem,

$$\begin{aligned}
1314 & \\
1315 & \mathbb{E}_{(X_1, X_2, Y) \sim \mu_\star} [\phi(X_1, X_2)] = \mathbb{E}_{(X_1, X_2, Y) \sim \mu_\star} \left[\lim_{m \rightarrow \infty} \phi_m(X_1, X_2) \right] \\
1316 & = \lim_{m \rightarrow \infty} \mathbb{E}_{(X_1, X_2, Y) \sim \mu_\star} [\phi_m(X_1, X_2)] \\
1317 & \leq \lim_{m \rightarrow \infty} \mathbb{E}_Y \psi(Y) = \mathbb{E}_Y \psi(Y). \\
1318 & \\
1319 &
\end{aligned}$$

1320 Similarly, repeating some of the above steps,

$$\begin{aligned}
1321 \quad \mathbb{E}_Y \psi(Y) &= \limsup_{n \rightarrow \infty} \mathbb{E}_Y \psi_n(Y) \\
1322 &= \limsup_{n \rightarrow \infty} \mathbb{E}_{(X_1, X_2, Y) \sim \mu_{*,n}} [\phi_n(X_1, X_2)] \\
1323 &\leq \limsup_{n \rightarrow \infty} \mathbb{E}_{(X_1, X_2, Y) \sim \mu_{*,m}} [\phi_n(X_1, X_2)] && (\mu_{*,n} \text{ is the optimal coupling for } \phi_n) \\
1324 &\leq \mathbb{E}_{(X_1, X_2, Y) \sim \mu_{*,m}} \left[\lim_{n \rightarrow \infty} \phi_n(X_1, X_2) \right] && (\text{monotone convergence}) \\
1325 &\leq \mathbb{E}_{(X_1, X_2, Y) \sim \mu_{*,m}} [\phi(X_1, X_2)].
\end{aligned}$$

1330 Hence, $\mathbb{E}_Y \psi(Y) \leq \liminf_{m \geq 1} \mathbb{E}_{(X_1, X_2, Y) \sim \mu_{*,m}} [\phi(X_1, X_2)]$. By assumption, $\phi(X_1, X_2)$ is lower semicontinuous and bounded from below. Thus, the Portmanteau theorem (Durrett, 2019) implies that, as $\mu_{*,m} \xrightarrow{\text{weak}} \mu_*$, $\liminf_{m \geq 1} \mathbb{E}_{(X_1, X_2, Y) \sim \mu_{*,m}} [\phi(X_1, X_2)] = \mathbb{E}_{(X_1, X_2, Y) \sim \mu_*} [\phi(X_1, X_2)]$. Hence, $\mathbb{E}_Y \psi(Y) \leq \mathbb{E}_{(X_1, X_2, Y) \sim \mu_*} [\phi(X_1, X_2)]$, proving the reverse inequality.

1335 **Proof of the last statement.** To prove the last statement, we observe that if $\mu_* \in \mathcal{C}_{P_Y}(\mathbb{Q}_1, \mathbb{Q}_2)$ then there exists a version of $(\mu_*)_{X, X'|Y}$ that is a regular conditional probability and such that for almost every y it holds that $(\mu_*)_{X, X'|y} \in \mathcal{C}(\mathbb{Q}_1(y), \mathbb{Q}_2(y))$. Indeed, the existence of a version that is a regular conditional probability is immediate by Theorem 3. To see that this version is a valid coupling of $\mathbb{Q}_1(y)$ and $\mathbb{Q}_2(y)$, observe that under μ_* , the joint law of $(X, Y) \sim \mathbb{Q}_1$ and thus the conditional distribution under μ_* of $X|Y$ is determined up to sets of \mathbb{Q}_1 -measure 0. In particular, again by Theorem 3, there exists a regular conditional probability that is a version of $(\mu_*)_{X|y}$ and this must agree almost everywhere with $(\mathbb{Q}_1)_{X|y} = \mathbb{Q}_1(y)$. The same argument holds for X' and thus $(\mu_*)_{X, X'|y} \in \mathcal{C}(\mathbb{Q}_1(y), \mathbb{Q}_2(y))$ for almost every y . Thus, by definition of ψ as an infimum, it holds for almost every y that

$$1344 \quad \psi(y) \leq \mathbb{E}_{(X, X') \sim (\mu_*)_{X, X'|y}} [\phi(X, X')].$$

1345 By the second claim of the proposition, we also have that

$$1347 \quad \mathbb{E}_{\mu_*} [\phi(X_1, X_2)] = \mathbb{E}_{\mu_*} [\psi(Y)].$$

1349 Because the expectations are equal and one function is pointwise almost everywhere dominated by the other function, the two functions must be equal almost everywhere, concluding the proof. \square

1352 C.5. A simple union-bound recursion.

1353 Finally, we also use the following version of the union bound extensively in our recursion proofs.

1354 **Lemma C.11.** For any event \mathcal{E} and events $\mathcal{B}_1, \mathcal{B}_2, \dots, \mathcal{B}_H$, it holds that

$$1356 \quad \mathbb{P}[(\mathcal{Q} \cap \bigcap_{h=1}^H \mathcal{B}_h)^c] \leq \mathbb{P}[\mathcal{Q}^c] + \mathbb{P} \left[\exists h \in [H] \text{ s.t. } \left(\mathcal{Q} \cap \bigcap_{j=1}^{h-1} \mathcal{B}_j \cap \mathcal{B}_h^c \right) \text{ holds} \right]$$

1360 *Proof.* Note that

$$1362 \quad \left(\mathcal{Q} \cap \bigcap_{h=1}^H \mathcal{B}_h \right)^c = \mathcal{Q}^c \cup \left(\mathcal{Q} \cap \left(\bigcap_{h=1}^H \mathcal{B}_h \right)^c \right) = \mathcal{Q}^c \cup \bigcup_{h=1}^H \mathcal{Q} \cap \mathcal{B}_h \cap \bigcap_{j=1}^{h-1} \mathcal{B}_j.$$

1365 The result follows by a union bound. \square

1367 D. Warmup: Analysis Without Augmentation

1369 In this section, we give a simplified analysis that replaces the smoothing kernels W_σ with the assumption that the learner policy $\hat{\pi}$ is already total variation continuous. The removal of the coupling kernel makes the coupling construction considerably simpler while still communicating some intuition for the full proof in Appendix E.

1372 Throughout this section, we make the following assumptions on the state and action spaces, along with their associated metrics:

1374

Assumption D.1. We assume that \mathcal{S} and \mathcal{A} are Polish spaces. This means they are metrizable, but we do not annotate their metrics because, e.g. the metric on \mathcal{S} may be other than $d_{\mathcal{S}}$. We further assume that

- $d_{\mathcal{S}}, d_{\text{TVC}}$ are pseudometrics and Borel measurable function from $\mathcal{S} \times \mathcal{S} \rightarrow \mathbb{R}_{\geq 0}$
- For any $\varepsilon \geq 0$, the set $\{(a, a') \in \mathcal{A} \times \mathcal{A} : d_{\mathcal{A}}(a, a') > \varepsilon\}$ is an open subset of $\mathcal{A} \times \mathcal{A}$; i.e. $d_{\mathcal{A}}(\cdot, \cdot)$ is lower semicontinuous. In particular, this means $d_{\mathcal{A}}$ is a Borel measurable function.

Recall the definitions of total variation continuity (TVC) and input-stability in [Section 4](#). The main result of this section is as follows.

Proposition D.1. *Let π^* be input-stable w.r.t. $(d_{\mathcal{S}}, d_{\mathcal{A}})$ and let $\hat{\pi}$ be γ -TVC. Then, for all $\varepsilon > 0$, $\Gamma_{\text{joint}, \varepsilon}(\hat{\pi} \parallel \pi^*) \leq H\gamma(\varepsilon) + \sum_{h=1}^H \mathbb{E}_{s_h^* \sim P_h^*} d_{\text{os}, \varepsilon}(\hat{\pi}_h(s_h^*) \parallel \pi^*(s_h^*))$.*

Proof. The key to the proof is to construct an appropriate ‘‘interpolating sequence’’ of actions $\hat{a}_{1:H}^{\text{inter}}$ to which we couple both $(s_{1:H+1}^*, a_{1:H}^*)$ and $(\hat{s}_{1:H+1}, \hat{a}_{1:H})$. This technique will be used in a significantly more sophisticated manner in the sequel to prove the analogous result with smoothing.

Let \mathcal{F}_h denote the σ -algebra generated by $(s_{1:h}^*, a_{1:h}^*)$, $(\hat{s}_{1:h}, \hat{a}_{1:h})$, and $\hat{a}_{1:h}^{\text{inter}}$, and let \mathcal{F}_0 denote the σ -algebra generated by s_1^*, \hat{s}_1 . We construct couplings of the following form:

- The initial states are generated as $s_1^* = \hat{s}_1 \sim P_{\text{init}}$.
- The dynamics are determined by F_h :

$$s_{h+1}^* = F_h(s_h^*, a_h^*), \quad \hat{s}_{h+1} = F_h(\hat{s}_h, \hat{a}_h) \quad (\text{D.1})$$

In particular, $s_{h+1}^*, \hat{s}_{1:h+1}$ are \mathcal{F}_h measurable.

- The conditional distributions of the primitive controllers satisfy the following

$$a_h^* \mid \mathcal{F}_{h-1} \sim \pi_h^*(s_h^*), \quad \hat{a}_{h-1} \mid \mathcal{F}_{h-1} \sim \hat{\pi}_h(\hat{s}_h), \quad \hat{a}_h^{\text{inter}} \mid \mathcal{F}_h \sim \hat{\pi}_h(s_h^*). \quad (\text{D.2})$$

Note that if μ satisfies the above construction, then $(s_{1:H+1}^*, s_{1:H}^*) \sim D_{\pi^*}$ and $(\hat{s}_{1:H+1}, \hat{a}_{1:H}) \sim D_{\hat{\pi}}$.

Specifying the rest of the coupling. It remains to specify the coupling of the terms in [\(D.2\)](#). We establish our coupling sequentially. Let $\mu^{(0)}$ denote the coupling of $\hat{s}_1 = s_1^* \sim P_{\text{init}}$.

Assume we have constructed the coupling up to state $h-1$. For ease, let Y_{h-1} denote the random variable corresponding to $(s_{1:h}^*, \hat{s}_{1:h}, a_{1:h-1}^*, \hat{a}_{1:h-1}, \hat{a}_{1:h-1}^{\text{inter}})$; note that Y_{h-1} is \mathcal{F}_{h-1} -measurable (as \hat{s}_h, s_h^* are determined by the dynamics [\(D.1\)](#)). Observe that, by the assumption of $\hat{\pi}_h$ being TVC, it holds that

$$\text{TV}(\mathbb{P}_{\hat{a}_h \mid Y_{h-1}}, \mathbb{P}_{\hat{a}_h^{\text{inter}} \mid Y_{h-1}}) \leq \gamma(d_{\text{TVC}}(\hat{s}_h, s_h^*)).$$

Thus by [Lemma C.1](#), there exists a coupling $\mu_1^{(h)}$ between $Y_{h-1}, \hat{a}_h, \hat{a}_h^{\text{inter}}$, with $Y_{h-1} \sim \mu^{(h-1)}$ such that it holds that

$$\mathbb{P}[\hat{a}_h \neq \hat{a}_h^{\text{inter}}] \leq \mathbb{E}_{\mu^{(h-1)}}[\gamma(d_{\text{TVC}}(\hat{s}_h, s_h^*))].$$

Similarly by [Proposition C.3](#), there is a coupling $\mu_2^{(h)}$ of $Y_{h-1}, \hat{a}_h^{\text{inter}}, a_h^*$ such that

$$\mathbb{P}_{\mu_2^{(h)}}[d_{\mathcal{A}}(\hat{a}_h^{\text{inter}}, a_h^*) \geq \varepsilon] \leq \mathbb{E}_{s_h^* \sim \mu^{(h-1)}}[d_{\text{os}, \varepsilon}(\hat{\pi}_h(s_h^*), \pi_h^*(s_h^*))].$$

By the gluing lemma [Lemma C.2](#) and a union bound, we may construct a coupling $\mu^{(h)}$ of $Y_h, \hat{a}_h^{\text{inter}}, a_h^*, \hat{a}_h$ such that (almost surely),

$$\begin{aligned} & \mathbb{P}_{\mu^{(h)}}[\{d_{\mathcal{A}}(\hat{a}_h^{\text{inter}}, a_h^*) \geq \varepsilon\} \cup \{\hat{a}_h \neq \hat{a}_h^{\text{inter}}\} \mid \mathcal{F}_{h-1}] \\ &= \mathbb{P}_{\mu^{(h)}}[\{d_{\mathcal{A}}(\hat{a}_h^{\text{inter}}, a_h^*) \geq \varepsilon\} \cup \{\hat{a}_h \neq \hat{a}_h^{\text{inter}}\} \mid Y_{h-1}] \\ &\leq \gamma(d_{\text{TVC}}(\hat{s}_h, s_h^*)) + d_{\text{os}, \varepsilon}(\hat{\pi}_h(s_h^*), \pi_h^*(s_h^*)) \end{aligned} \quad (\text{D.3})$$

Thus inductively, we may continue this construction for $h \leq H$ and let $\mu = \mu^{(H)}$.

1430 **Concluding the proof.** Define the event $\mathcal{B}_h := \{d_{\mathcal{A}}(\mathbf{a}_h, \hat{\mathbf{a}}_h^{\text{inter}}) \leq \varepsilon\}$ and $\mathcal{C}_h = \{\hat{\mathbf{a}}_h^{\text{inter}} = \hat{\mathbf{a}}_h\}$. Then, by [Lemma C.11](#)

$$1431 \mathbb{P}_{\mu} \left[\left(\bigcap_{h=1}^H \mathcal{B}_h \cap \mathcal{C}_h \right)^c \right] \leq \sum_{h=1}^H \mathbb{P}_{\mu} \left[\left(\bigcap_{j=1}^{h-1} \mathcal{B}_j \cap \mathcal{C}_j \right) \cap (\mathcal{B}_h^c \cup \mathcal{C}_h^c) \right]. \quad (\text{D.4})$$

1432 Note first that $(\bigcap_{j=1}^{h-1} \mathcal{B}_j \cap \mathcal{C}_j)$ is \mathcal{F}_{h-1} measurable. On this event, input stability at $\hat{\mathbf{a}}_j^{\text{inter}} = \hat{\mathbf{a}}_j$, $1 \leq j \leq h-1$, implies that

$$1433 d_{\mathcal{S}}(\mathbf{s}_h^*, \hat{\mathbf{s}}_h) \leq \varepsilon.$$

1434 Thus, (D.3) implies that

$$1435 \mathbb{P}_{\mu} \left[\left(\bigcap_{j=1}^{h-1} \mathcal{B}_j \cap \mathcal{C}_j \right) \cap (\mathcal{B}_h^c \cup \mathcal{C}_h^c) \right] \leq \mathbb{E}_{\mu} [\gamma(d_{\text{TVC}}(\hat{\mathbf{s}}_h, \mathbf{s}_h^*)) \mathbf{I}\{d_{\text{TVC}}(\hat{\mathbf{s}}_h, \mathbf{s}_h^*) \leq \varepsilon\} + d_{\text{os},\varepsilon}(\hat{\pi}_h(\mathbf{s}_h^*), \pi_h^*(\mathbf{s}_h^*)) \mid \mathcal{F}_{h-1}]$$

$$1436 \leq \gamma(\varepsilon) + \mathbb{E}_{\mu} [\mathbb{E}_{\mu} [d_{\text{os},\varepsilon}(\hat{\pi}_h(\mathbf{s}_h^*), \pi_h^*(\mathbf{s}_h^*)) \mid \mathcal{F}_{h-1}]]$$

$$1437 = \gamma(\varepsilon) + \mathbb{E}_{\mu} [d_{\text{os},\varepsilon}(\hat{\pi}_h(\mathbf{s}_h^*), \pi_h^*(\mathbf{s}_h^*))]$$

$$1438 = \gamma(\varepsilon) + \mathbb{E}_{\mathbf{s}_h^* \sim \mathcal{P}_h^*} \mathbb{E}_{\mu} [d_{\text{os},\varepsilon}(\hat{\pi}_h(\mathbf{s}_h^*), \pi_h^*(\mathbf{s}_h^*))],$$

1439 where the first equality follows from the tower rule for conditional expectations and the second follows because $\mathbf{s}_h^* \sim \mathcal{P}_h^*$ under μ . Summing and applying (D.4) implies that

$$1440 \mathbb{P}_{\mu} \left[\left(\bigcap_{h=1}^H \mathcal{B}_h \cap \mathcal{C}_h \right)^c \right] \leq H\gamma(\varepsilon) + \sum_{h=1}^H \mathbb{E}_{\mathbf{s}_h^* \sim \mathcal{P}_h^*} [d_{\text{os},\varepsilon}(\hat{\pi}_h(\mathbf{s}_h^*), \pi_h^*(\mathbf{s}_h^*))].$$

1441 Again, invoking input stability and the definitions $\mathcal{B}_h := \{d_{\mathcal{A}}(\mathbf{a}_h, \hat{\mathbf{a}}_h^{\text{inter}}) \leq \varepsilon\}$ and $\mathcal{C}_h = \{\hat{\mathbf{a}}_h^{\text{inter}} = \hat{\mathbf{a}}_h\}$, $(\bigcap_{h=1}^H \mathcal{B}_h \cap \mathcal{C}_h)^c$ implies that

$$1442 \max_{1 \leq h \leq H} \max\{d_{\mathcal{S}}(\mathbf{s}_{h+1}^*, \hat{\mathbf{s}}_{h+1}), d_{\mathcal{A}}(\mathbf{a}_h^*, \hat{\mathbf{a}}_h)\} \leq \varepsilon.$$

1443 This concludes the proof. □

1444 E. Imitation in the Composite MDP

1445 In this section, we prove our imitation guarantees in the composite MDP under the full generality of data augmentation. The majority of this section is devoted to proving a more general version of [Theorem 2](#) that applies to vectorized notions of distance and helps tighten our bounds when instantiated in the control setting. In [Appendix E.1](#), we introduce some notation and state our most general result, [Theorem 4](#). We then proceed to show that [Theorem 2](#) follows from [Theorem 4](#) and in [Appendix E.2](#), we provide a detailed and rigorous proof of the main result. In [Appendix E.3](#), we show that the more general [Theorem 4](#) implies [Theorem 2](#) from the text.

1446 Throughout, we also assume \mathcal{S} admits a direct decomposition. This is useful to capture the fact that we only apply smoothing on the $\rho_{\text{m},h}$ coordinates (memory chunk), not the full trajectory chunk $\rho_{\text{c},h}$.

1447 **Definition E.1** (Direct Decomposition). Let $\mathcal{S} = \mathcal{Z} \oplus \mathcal{V}$ is a direct decomposition. We let $\phi_{\mathcal{Z}}$ and $\phi_{\mathcal{V}}$ denote projections onto the \mathcal{Z} and \mathcal{V} components, respectively. We say that the $\mathcal{S} = \mathcal{Z} \oplus \mathcal{V}$ is *compatible* with the dynamics if $F_h((\mathbf{z}, \mathbf{v}), \mathbf{a}) = F_h((\mathbf{z}, \mathbf{v}'), \mathbf{a})$ for all $\mathbf{v}, \mathbf{v}' \in \mathcal{V}$ and $\mathbf{z} \in \mathcal{Z}$, and *compatible* with policy π if $\pi_h((\mathbf{z}, \mathbf{v}), \mathbf{a}) = \pi_h((\mathbf{z}, \mathbf{v}'), \mathbf{a})$; we define compatibility of a kernel W and of a pseudometric $d(\cdot, \cdot) : \mathcal{S} \times \mathcal{S} \rightarrow \mathbb{R}_{\geq 0}$ with $\mathcal{S} = \mathcal{Z} \oplus \mathcal{V}$ similarly.

1448 We emphasize that compatibility of dynamics with a direct decomposition does not make \mathbf{v} irrelevant because $d_{\mathcal{S}}$ still depends on \mathbf{v} . For the purposes of the instantiation for control in the following appendix, we wish to control the imitation gaps on distances that do depend on \mathbf{v}_h , even though \mathbf{v}_h does not figure directly into the dynamics. Note that as defined, \mathbf{v}_h does depend on the dynamics up until time $h-1$ and thus it is necessary to deal with this component in order to provide guarantees in $d_{\mathcal{S}}$.

E.1. A generalization of Theorem 2

We now state a generalization of [Theorem 2](#), which replaces a single distance by a vector of distances of dimension K ; this will be useful for our instantiation of the composite MDP as a chunked control system in our final application (in particular, for deriving a bound on $\mathcal{L}_{\text{fin},\varepsilon}$). It also showcases the most general structure accomodated by our proof technique.

We begin by defining some notation:

- Let $K \in \mathbb{N}$ denote a dimension
- Let $\vec{\varepsilon} \in \mathbb{R}_{\geq 0}^K$ denote a vector of tolerances
- Let $\vec{d}_{\mathcal{S}}(\cdot, \cdot)$ denote a vector of pseudometrics $d_{\mathcal{S},i}$ on \mathcal{S}
- Let $\vec{d}_{\mathcal{A}}$ denote a vector of non-negative functions $d_{\mathcal{A},i} : \mathcal{A}^2 \rightarrow \mathbb{R}_{\geq 0}$, not necessarily pseuometrics.
- Let \preceq denote vector wise inequality, and let the symbols \wedge and \vee be generalized to denote entrywise minima and maxima. Similarly, addition of vectors is coordinate wise with scalars assumed to be broadcast appropriately.
- We let $d_{\mathcal{S},1} = d_{\text{TVC}}$ denote the metric we consider for evaluating total variation distance.

We generalize We assume the following measure-theoretic regularity conditions, generalizing [Assumption D.1](#) as follows.

Assumption E.1. We assume that \mathcal{S} and \mathcal{A} are Polish spaces. This means they are metrizable, but we do not annotate their metrics because, e.g. the metric on \mathcal{S} may be other than $d_{\mathcal{S}}$. We further assume that

- $d_{\mathcal{S},i}$ is a pseudometric and Borel measurable function from $\mathcal{S} \times \mathcal{S} \rightarrow \mathbb{R}_{\geq 0}$.
- For any $\varepsilon \geq 0$, the set $\{(a, a') \in \mathcal{A} \times \mathcal{A} : d_{\mathcal{A},i}(a, a') > \varepsilon\}$ is an open subset of $\mathcal{A} \times \mathcal{A}$; i.e. $d_{\mathcal{A},i}(\cdot, \cdot)$ is lower semicontinuous. In particular, this means $d_{\mathcal{A},i}$ is a Borel measurable function. Note that this implies that the

$$\{(a, a') \in \mathcal{A} \times \mathcal{A} : \vec{d}_{\mathcal{A}}(a, a') \not\preceq \vec{\varepsilon}\}.$$

is closed and thus measurable.

Note that the above assumption is the natural vectorized generalization of [Assumption D.1](#). Next, we define vector versions of our imitation errors.

Definition E.2 (Imitation Errors, vector version). Given error parameter $\vec{\varepsilon} \in \mathbb{R}_{\geq 0}^K$, define

- The **vector joint-error**

$$\vec{\Gamma}_{\text{joint},\vec{\varepsilon}}(\hat{\pi} \parallel \pi^*) := \inf_{\mu_1} \mathbb{P}_{\mu_1} \left[\exists h \in [H] : \vec{d}_{\mathcal{S}}(\hat{s}_{h+1}, s_{h+1}^*) \vee \vec{d}_{\mathcal{A}}(a_h^*, \hat{a}_h) \not\preceq \vec{\varepsilon} \right],$$

where the infimum is over trajectory couplings $((\hat{s}_{1:H+1}, \hat{a}_{1:H}), (s_{1:H+1}^*, a_{1:H}^*)) \sim \mu_1 \in \mathcal{C}(\mathcal{D}_{\hat{\pi}}, \mathcal{D}_{\pi^*})$ satisfying $\mathbb{P}_{\mu_1}[\hat{s}_1 = s_1^*] = 1$.

- The **vector marginal error**

$$\vec{\Gamma}_{\text{marg},\vec{\varepsilon}}(\hat{\pi} \parallel \pi^*) := \max_{h \in [H]} \max \left\{ \inf_{\mu_1} \mathbb{P}_{\mu_1} \left[\vec{d}_{\mathcal{S}}(\hat{s}_{h+1}, s_{h+1}^*) \not\preceq \vec{\varepsilon} \right], \inf_{\mu_1} \mathbb{P}_{\mu_1} \left[\vec{d}_{\mathcal{A}}(a_h^*, \hat{a}_h) \not\preceq \vec{\varepsilon} \right] \right\}$$

the same as the to joint-gap, with the “max” outside the probability and infimum over couplings.

- The **vector-wise one-step error**

$$\vec{d}_{\text{os},\vec{\varepsilon}}(\hat{\pi}_h(s) \parallel \pi_h^*(s)) := \inf_{\mu_2} \mathbb{P}_{\mu_2} \left[\vec{d}_{\mathcal{A}}(\hat{a}_h, a_h^*) \not\preceq \vec{\varepsilon} \right],$$

where the infimum is over $(a_h^*, \hat{a}_h) \sim \mu_2 \in \mathcal{C}(\hat{\pi}_h(s), \pi_h^*(s))$.

1540 We now describe input stability.

1541 **Definition E.3** (Input-Stability, vector version). A trajectory $(s_{1:H+1}, a_{1:H})$ is *input-stable* w.r.t. $(\vec{d}_{\mathcal{S}}, \vec{d}_{\mathcal{A}})$ if all sequences
 1542 $s'_1 = s_1$ and $s'_{h+1} = F_h(s'_h, a'_h)$ satisfy
 1543

$$1544 \quad d_{\mathcal{S},i}(s'_{h+1}, s_{h+1}) \leq \max_{1 \leq j \leq h} d_{\mathcal{A},i}(a'_j, a_j), \quad \forall h \in [H], i \in [K]$$

1545
 1546 Finally, define input process stability. A slight technicality is that, in our instantiation, π^* is taken to be a suitable regular
 1547 condition probability of the joint distribution \mathcal{D}_{exp} of expert trajectories. This means that π^* can only really satisfy desired
 1548 regularity conditions on states visited with positive probability by \mathcal{D}_{exp} . We address this subtlety by considering the
 1549 following definition generalizing [Definition 4.5](#) in the body. We also restrict the kernels under consideration to those which
 1550 produce distributions *absolutely continuous* ([Definition C.4](#)) with respect to P_h^* , and denoted with the \ll comparator. More
 1551 specifically, we only care about absolute continuity under the projections onto the \mathcal{Z} component of \mathcal{S} .
 1552

1553 **Definition E.4** (Input & Process Stability, vector version). Let $p_{\text{IPS}} \in (0, 1)$, $\vec{\gamma}_{\text{IPS}} = (\gamma_{\text{IPS},i})_{1 \leq i \leq K}$ be a collection non-
 1554 decreasing maps $\gamma_{\text{IPS},i} : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}_{\geq 0}$, let $d_{\text{IPS}} : \mathcal{S} \times \mathcal{S} \rightarrow \mathbb{R}$ be a pseudometric (possibly other than any of the $d_{\mathcal{S},i}$), and
 1555 $r_{\text{IPS}} > 0$. We say a policy π^* is $(\vec{\gamma}_{\text{IPS}}, d_{\text{IPS}}, r_{\text{IPS}}, p_{\text{IPS}})$ -*vectorwise-input-&-process stable (vIPS)* if the following holds for
 1556 any $r \in [0, r_{\text{IPS}}]$:

1557 Consider any sequence of kernels $W_h : \mathcal{S} \rightarrow \Delta(\mathcal{S})$, $1 \leq h \leq H$, satisfying

$$1558 \quad \forall h, s \in \mathcal{S} : \mathbb{P}_{\tilde{s} \sim W_h(s)}[d_{\text{IPS}}(\tilde{s}, s) \leq r] = 1, \quad \phi_{\mathcal{Z}} \circ W_h(s) \ll \phi_{\mathcal{Z}} \circ P_h^*. \quad (\text{E.1})$$

1559 Define a process $s_1 \sim P_{\text{init}}$, $\tilde{s}_h \sim W_h(s_h)$, $a_h \sim \pi_h(\tilde{s}_h)$, and $s_{h+1} := F_h(s_h, a_h)$. Then, with probability at least $1 - p_{\text{IPS}}$,

- 1560
 1561 (a) the sequence $(s_{1:H+1}, a_{1:H})$ is input-stable w.r.t. $(\vec{d}_{\mathcal{S}}, \vec{d}_{\mathcal{A}})$ (as defined by [Definition E.3](#)).
 1562
 1563 (b) $\max_{h \in [H]} d_{\mathcal{S},i}(F_h(\tilde{s}_h, a_h), s_{h+1}) \leq \gamma_{\text{IPS},i}(r)$.

1564 We can now state our desired generalization.

1565 **Theorem 4.** *Suppose that there*

- 1566
 1567 (a) π^* is $(\vec{\gamma}_{\text{IPS}}, d_{\text{IPS}}, r_{\text{IPS}}, p_{\text{IPS}})$ -vector IPS in the sense of [Definition E.4](#).
 1568
 1569 (b) There is a direct decomposition of $\mathcal{S} = \mathcal{Z} \oplus \mathcal{V}$, which associated projection maps $\phi_{\mathcal{Z}}$ and $\phi_{\mathcal{V}}$, and which is compatible
 1570 with the dynamics, and policies π^* , $\hat{\pi}$, and smoothing kernel W_{σ} , and d_{IPS} .
 1571
 1572 (c) $\phi_{\mathcal{Z}} \circ W_{\sigma}$ is γ_{σ} -TVC with respect to the pseudometric $d_{\text{TVC}} = d_{\mathcal{S},1}$.
 1573
 1574

1575 Let $\hat{\pi}_{\sigma}$ be any policy which is $\hat{\gamma}$ -TVC, also w.r.t. $d_{\text{TVC}} = d_{\mathcal{S},1}$. Finally, let $\vec{\varepsilon} \in \mathbb{R}_{\geq 0}^K$, $r \in (0, \frac{1}{2}r_{\text{IPS}}]$, and define

$$1576 \quad p_r := \sup_{\mathcal{S}} \mathbb{P}_{s' \sim W_{\sigma}(s)}[d_{\text{IPS}}(s', s) > r], \quad \vec{\varepsilon}_{\text{marg}} := \vec{\varepsilon} + \vec{\gamma}_{\text{IPS}}(2r).$$

1577 Then,

- 1578 • For any policy $\hat{\pi}$, both $\vec{\Gamma}_{\text{joint}, \vec{\varepsilon}}(\hat{\pi}_{\sigma} \parallel \pi_{\odot}^*)$ and $\vec{\Gamma}_{\text{marg}, \vec{\varepsilon}_{\text{marg}}}(\hat{\pi}_{\sigma} \parallel \pi^*)$ are upper bounded by

$$1582 \quad p_{\text{IPS}} + H(2p_r + \hat{\gamma}(\vec{\varepsilon}_1) + (\hat{\gamma} + \gamma_{\sigma}) \circ \gamma_{\text{IPS},1}(2r)) + \sum_{h=1}^H \mathbb{E}_{s_h^* \sim P_h^*} \vec{d}_{\text{os}, \vec{\varepsilon}}(\hat{\pi}_{\sigma, h}(s_h^{\text{tel}}) \parallel \pi_{\odot, h}^*(s_h^{\text{tel}})) \quad (\text{E.2})$$

- 1583
 1584 • In the special case where $\hat{\pi}_{\sigma} = \hat{\pi} \circ W_{\sigma}$, we can take $\hat{\gamma} = \gamma_{\sigma}$, and obtain that $\vec{\Gamma}_{\text{joint}, \vec{\varepsilon}}(\hat{\pi}_{\sigma} \parallel \pi_{\odot}^*)$ and $\vec{\Gamma}_{\text{marg}, \vec{\varepsilon}_{\text{marg}}}(\hat{\pi}_{\sigma} \parallel \pi^*)$ are upper bounded by

$$1585 \quad p_{\text{IPS}} + H(2p_r + 3\gamma_{\sigma}(\max\{\varepsilon, \gamma_{\text{IPS},1}(2r)\})) + \sum_{h=1}^H \mathbb{E}_{s_h^* \sim P_h^*} \mathbb{E}_{\tilde{s}_h^* \sim W_{\sigma}(s_h^*)} \vec{d}_{\text{os}, \vec{\varepsilon}}(\hat{\pi}_h(\tilde{s}_h^*) \parallel \pi_{\text{dec}}^*(\tilde{s}_h^*)). \quad (\text{E.3})$$

1586 We note that [Theorem 2](#) is a special case of [Theorem 4](#) and prove the former assuming the latter here at the end of the
 1587 section.
 1588
 1589
 1590
 1591
 1592
 1593
 1594

1595 E.2. Proof of Theorem 4

1596 In this section, we prove [Theorem 4](#). We begin with an intuitive overview of the proof and partially construct the relevant
 1597 intermediate trajectories used to define our coupling in [Appendix E.2.1](#). In [Appendix E.2.2](#), we prove several prerequisite
 1598 properties of the construction given in [Appendix E.2.1](#). Finally, in [Appendix E.2.3](#) we formally construct the coupling and
 1599 rigorously prove [Theorem 4](#).
 1600

1601 E.2.1. PROOF OVERVIEW AND COUPLING CONSTRUCTION

1602 The proof proceeds by constructing a sophisticated coupling between the law of a trajectory evolving according to $\hat{\pi}$ and
 1603 a trajectory evolving according to π_{\odot}^* by introducing several intermediate sequences of composite states and composite
 1604 actions.
 1605

1606 We partially specify this coupling below and formally construct it in [Appendix E.2.3](#). Our construction is recursive and
 1607 relies on the input and process stability as well as total variation continuity to show that if the trajectories generated by
 1608 π_{\odot}^* and $\hat{\pi}$ are close in $\bar{d}_{os, \bar{\varepsilon}}$ evaluated on states at step h , then they will remain close at step $h + 1$. There are a number of
 1609 technical subtleties involved, especially those of a measure-theoretic nature, but much of the intuition can be gleaned from
 1610 the following partial specification of the coupling μ over composite-state $(\hat{s}_{1:H}, s_{1:H}^{\odot}, s_{1:H}^{\text{tel}}, \tilde{s}_{1:H}^{\text{tel}}) \subset \mathcal{S}$, composite-actions
 1611 $(\hat{a}_{1:H}^{\odot}, \hat{a}_{1:h}, a_{1:H}^{\text{tel}}) \subset \mathcal{K}$ and interpolating composite-actions, $(\hat{a}_{1:H}^{\odot, \text{inter}}, \hat{a}_{1:H}^{\text{tel, inter}}) \subset \mathcal{A}$.
 1612

1613 To define the construction, we define the probability kernels corresponding to the replica and deconvolution policies. Note
 1614 that these are slightly different from the definitions in the body due to the use of the direct decomposition; the intuition is the
 1615 same, however.
 1616

1617 **Definition E.5** (Replica and Deconvolution Kernels). Let $P_{\text{aug}, h}^{\text{proj}}$ denote the joint distribution over $(z_h^*, s_h^*, \tilde{z}_h^*, a_h^*)$ under the
 1618 generative process

$$1619 s_h^* \sim P_h^*, \quad a_h^* \sim \pi_h^*(s_h^*), \quad z_h^* = \phi_{\mathcal{Z}}(s_h^*), \quad \tilde{z}_h^* \sim \phi_{\mathcal{Z}} \circ W_{\sigma}(s_h^*)$$

1620 For $z \in \mathcal{Z}$, let $W_{\text{dec}, \mathcal{Z}, h}^*(z)$ denote the distribution of z_h^* conditioned on $\tilde{z}_h^* = z$, under $P_{\text{aug}, h}^{\text{proj}}$. Given $s = (z, v)$, define

$$1621 W_{\text{dec}, h}^*(s) = W_{\text{dec}, \mathcal{Z}, h}^*(\phi_{\mathcal{Z}}(s)) \otimes \delta_{\phi_{\mathcal{V}}(s)},$$

$$1622 W_{\odot, h}^*(s) = W_{\text{dec}, h}^* \circ (W_{\sigma}(\phi_{\mathcal{Z}}(s)) \otimes \delta_{\phi_{\mathcal{V}}(s)}) = (W_{\text{dec}, \mathcal{Z}, h}^* \circ W_{\sigma}(\phi_{\mathcal{Z}}(s))) \otimes \delta_{\phi_{\mathcal{V}}(s)}.$$

1623 where we recall the dirac-delta δ . Equivalently, $W_{\text{dec}, h}^*(s)$ denotes the conditional sequence of (\tilde{z}, v) , where $v = \phi_{\mathcal{V}}(s)$, and
 1624 $\tilde{z} \sim W_{\text{dec}, \mathcal{Z}, h}^*(s)$; $W_{\odot, h}^*$ can be expressed similarly.
 1625

1626 We remark that $W_{\text{dec}, h}^*$ and $W_{\odot, h}^*$ are both kernels and by [Theorem 3](#), we may assume that the joint distribution over
 1627 $(s_h^*, \tilde{s}_h^{\text{tel}})$ admits a regular conditional probability and thus these constructions are well-defined.
 1628

1629 **Remark E.1.** Note that the kernels $W_{\text{dec}, h}^*$ and $W_{\odot, h}^*$ are compatible with the decomposition $\mathcal{S} = \mathcal{Z} \oplus \mathcal{V}$ by construction.
 1630 Moreover, note that if $s = (z, v)$, $\phi_{\mathcal{V}} \circ W_{\text{dec}, h}^*(s) = \phi_{\mathcal{V}} \circ W_{\odot, h}^*(s)$ is the dirac-delta distribution supported on v .
 1631

1632 **Lemma E.1.** Under our the assumption that π^* and W_{σ} are compatible with the direct decomposition,

$$1633 \pi_{\text{dec}, h}^*(s) = \pi^* \circ W_{\text{dec}, h}^*, \quad \pi_{\odot, h}^*(s) = \pi^* \circ W_{\odot, h}^*$$

1634 *Proof.* This follows immediately because π^* and W_{σ} are compatible with the direct decomposition, and by the definition of
 1635 [Definition 4.4](#). \square
 1636

1637 **A template for the coupling.** Our couplings are partially specified by the following generative process, and what remains
 1638 unspecified are couplings between random variables at each each step h . In what follows, let \mathcal{F}_0 denote the σ -algebra
 1639 generated by $\hat{s}_1 = s_1^{\odot} = s_1^{\text{tel}}$. Let \mathcal{F}_h denote the sigma-algebra generated by $(\hat{s}_{1:h}, s_{1:h}^{\odot}, s_{1:h}^{\text{tel}})$, $(a_{1:h}^{\odot}, \tilde{s}_{1:h}^{\odot}, \tilde{s}_{1:h}^{\text{tel}}, a_{1:h}^{\text{tel}}, \hat{a}_{1:h})$,
 1640 and $(\hat{a}_{1:h}^{\odot, \text{inter}}, \hat{a}_{1:h}^{\text{tel, inter}})$.
 1641

- 1642 • The initial states are drawn as

$$1643 \hat{s}_1 = s_1^{\odot} = s_1^{\text{tel}} \sim P_{\text{init}}.$$

1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649

1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704

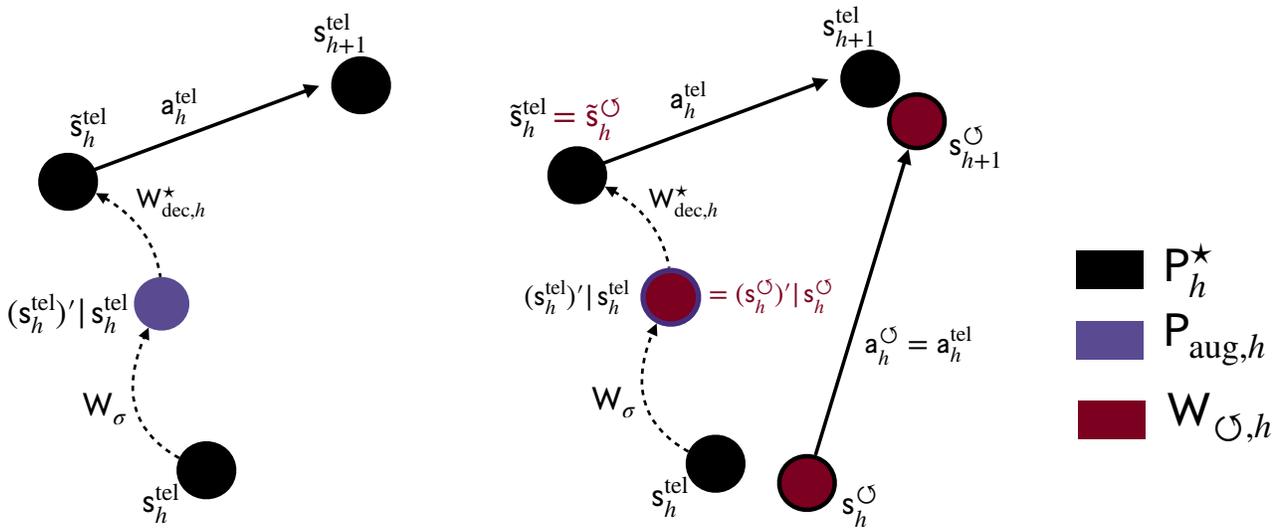


Figure 5. Graphical illustration of the coupling, in the special case where $\mathcal{Z} = \mathcal{S}$ for simplicity. **On the left** is the teleporting sequence, with $\tilde{s}_h^{\text{tel}} \sim W_{\circ,h}^*(s_h^{\text{tel}}) = W_{\text{dec},h}^* \circ W_\sigma(s_h^{\text{tel}})$. We represent the teleporting explicitly by noising s_h^{tel} to become $(s_h^{\text{tel}})'$ by applying W_σ and then applying $W_{\text{dec},h}^*$ to complete the “teleporting” to \tilde{s}_h^{tel} . We then apply $a_h^{\text{tel}} \sim \pi_h^*(\tilde{s}_h^{\text{tel}})$, and continue onto s_{h+1}^{tel} from the teleported state \tilde{s}_h^{tel} . **On the right**, we illustrate the replica sequence next to the teleporting sequence. We start with s_h° , which is close to s_h^{tel} (a consequence of our proof). We then apply the replica kernel to achieve \tilde{s}_h° . Our argument uses that $W_{\circ,h}^* = W_{\text{dec},h}^* \circ W_\sigma$ is TVC (a consequence of TVC of W_σ as shown in Lemma E.2). We depict this property pictorially: since W_σ is TVC and s_h^{tel} and s_h° are close, we can couple things in such a way that, with good probability, $(s_h^{\text{tel}})' \sim W_\sigma(s_h^{\text{tel}})$ and $(s_h^\circ)' \sim W_\sigma(s_h^\circ)$ are equal. We then extend the coupling to that $\tilde{s}_h^\circ = \tilde{s}_h^{\text{tel}}$ on the event $\{(s_h^{\text{tel}})' = (s_h^\circ)'\}$, both being drawn by applying $W_{\text{dec},h}^*$ to both of $(s_h^{\text{tel}})' = (s_h^\circ)'$. We extend the coupling once more so that $a_h^{\text{tel}} \sim \pi_h^*(\tilde{s}_h^{\text{tel}})$ and $a_h^\circ \sim \pi_h^*(\tilde{s}_h^\circ)$ are equal on this good probability event. Using our notion of stability, IPS, and the fact that s_h° and s_h^{tel} are close, the good probability event on which a_h^{tel} and a_h° are equal implies that s_{h+1}° remains close to s_{h+1}^{tel} . We remark that our actual analysis never explicitly computes the $(\cdot)'$ -terms drawn from W_σ ; rather, these terms appear implicitly in our definitions of $W_{\circ,h}^*$ and the verification of its TVC property.

- The dynamics satisfy

$$\hat{s}_{h+1} = F_h(\hat{s}_h, \hat{a}_h), \quad s_{h+1}^\circ = F_h(s_h^\circ, a_h^\circ), \quad s_{h+1}^{\text{tel}} = F_h(\tilde{s}_h^{\text{tel}}, a_h^{\text{tel}})$$

Note that determinism of the dynamics implies that s_{h+1}^{tel} , s_{h+1}° and \hat{s}_{h+1} are \mathcal{F}_h -measurable.

- We generate

$$\begin{aligned} \tilde{s}_h^\circ | \mathcal{F}_{h-1} &\sim W_{\circ,h}^*(s_h^\circ), & a_h^\circ | \mathcal{F}_{h-1}, \tilde{s}_h^\circ &\sim \pi_h^*(\tilde{s}_h^\circ), \\ \tilde{s}_h^{\text{tel}} | \mathcal{F}_{h-1} &\sim W_{\circ,h}^*(s_h^{\text{tel}}), & a_h^{\text{tel}} | \mathcal{F}_{h-1}, \tilde{s}_h^{\text{tel}} &\sim \pi_h^*(\tilde{s}_h^{\text{tel}}). \\ \hat{a}_h | \mathcal{F}_{h-1} &\sim \hat{\pi}_\sigma(\hat{s}_h) \end{aligned}$$

Importantly, we note that, marginalizing over \tilde{s}_h^{tel} and \tilde{s}_h° , respectively, $a_h^{\text{tel}} | \mathcal{F}_{h-1} \sim \pi_{\circ\sigma,h}^*(s^{\text{tel}})$ and $a_h^\circ | \mathcal{F}_{h-1} \sim \pi_{\circ\sigma,h}^*(s^\circ)$.

- Lastly, we select interpolating actions via

$$\hat{a}_h^{\circ,\text{inter}} | \mathcal{F}_{h-1} \sim \hat{\pi}_{\sigma,h}(s_h^\circ), \quad \hat{a}_h^{\text{tel},\text{inter}} | \mathcal{F}_{h-1} \sim \hat{\pi}_{\sigma,h}(s_h^{\text{tel}})$$

We will say μ is ‘‘respects the construction’’ as shorthand to mean that μ obeys the above equations. The coupling is illustrated graphically in [Figure 5](#). We now establish several key properties of the above constructions, separated into a subsection for the sake of clarity.

E.2.2. PROPERTIES OF SMOOTHING, DECONVOLUTION, AND REPLICAS.

In this section, we establish several useful properties of smoothed and replica policies. We begin by showing that smoothed policies are TVC.

Lemma E.2. *The following hold*

- For any h , $\phi_{\mathcal{Z}} \circ W_{\circ,h}^*$ and $\pi_{\circ\sigma,h}^*$ are γ_σ TVC.
- If π is any policy compatible with the direct decomposition $\mathcal{S} = \mathcal{Z} \oplus \mathcal{V}$ (in the sense of [Definition E.1](#)), then $\pi \circ W_\sigma$ is γ_σ -TVC.

Proof. We observe that $\phi_{\mathcal{Z}} \circ W_{\circ,h}^* = \phi_{\mathcal{Z}} \circ W_{\text{dec},h}^* \circ W_\sigma(s)$. Moreover, we observe $W_{\text{dec},h}^*$ satisfies $\phi_{\mathcal{Z}} \circ W_{\text{dec},h}^*(s) = W_{\text{dec},\mathcal{Z},h}^* \circ \phi_{\mathcal{Z}}$, so that $\phi_{\mathcal{Z}} \circ W_{\circ,h}^* = W_{\text{dec},\mathcal{Z},h}^* \circ \phi_{\mathcal{Z}} \circ W_\sigma(s)$. As $\phi_{\mathcal{Z}} \circ W_\sigma$ is TVC, the first claim is a consequence of the data-processing inequality [Lemma C.4](#). The second uses the fact that all listed objects involve composition of kernels with W_σ . \square

Next, we show that the replica construction preserves marginals.

Lemma E.3 (Marginal-Preservation). *There exists a coupling \mathbb{P} of $z_h \sim \phi_{\mathcal{Z}} \circ P_h^*$, $z'_h \sim \phi_{\mathcal{Z}} \circ W_\sigma(z_h, \cdot)$ (where (\cdot) denotes an irrelevant argument due to compatibility of W_σ with the direct decomposition), and $\tilde{z}_h \sim \phi_{\mathcal{Z}} \circ W_{\circ,h}^*(z_h, \cdot)$ (again, (\cdot) denotes an irrelevant argument) such that*

$$(z_h, z'_h) \stackrel{d}{=} (\tilde{z}_h, z'_h).$$

In particular, for s_h^{tel} and \tilde{s}_h^{tel} as in our construction, the marginal distributions of $\phi_{\mathcal{Z}}(s_h^{\text{tel}})$ and $\phi_{\mathcal{Z}}(\tilde{s}_h^{\text{tel}})$ are the same, where $s_h^{\text{tel}} \sim P_h^$ and $\tilde{s}_h^{\text{tel}} | s_h^{\text{tel}} \sim W_{\circ,h}^*(s_h^{\text{tel}})$.*

Proof. By [Assumption D.1](#) and [Theorem 3](#), we may assume that all joint distributions’ conditional probabilities are regular conditional probabilities and thus almost surely equal to a kernel. Moreover, since all kernels are compatible with the direct decomposition, it suffices to prove the special case of the trivial direct-decomposition where $\mathcal{Z} = \mathcal{S}$. Fix a common

1760 measure \mathbb{P} over which $s_h^{\text{tel}}, \tilde{s}_h^{\text{tel}}$, and s'_h are defined such that $s_h^{\text{tel}} \sim P_h^*$, $s'_h \sim W_\sigma(s_h^{\text{tel}})$, and $\tilde{s}_h^{\text{tel}} \sim W_{\text{dec},h}(s'_h)$. Then for any
 1761 measurable sets A, B , we have

$$\begin{aligned} 1762 \mathbb{P}(s_h^{\text{tel}} \in A, s'_h \in B) &= \mathbb{P}(s'_h \in B) \cdot \mathbb{E}_{s'_h} [\mathbf{I}[s'_h \in B] \cdot \mathbb{P}(s_h^{\text{tel}} \in A | s'_h)] \\ 1763 &= \mathbb{P}(s'_h \in B) \cdot \mathbb{E}_{s'_h} [\mathbf{I}[s'_h \in B] \cdot \mathbb{P}(\tilde{s}_h^{\text{tel}} \in A | s'_h)] \\ 1764 &= \mathbb{P}(\tilde{s}_h^{\text{tel}} \in A, s'_h \in B), \end{aligned}$$

1765 where the first equality holds by the fact that we are working with regular conditional probabilities and Bayes' rule, the
 1766 second equality holds by the definition of the deconvolution kernel above, and the last equality holds again by Bayes' rule
 1767 and the tower rule for conditional expectations.

1771 To prove the second statement, we apply induction, again assuming that $\mathcal{Z} = \mathcal{S}$ as in the proof of the first statement.
 1772 Note that $s_1^{\text{tel}} \sim P_1^* = P_{\text{init}}$, and $\tilde{s}_1^{\text{tel}} \sim W_{\odot,1}^* \circ P_1^*$. Thus, from the first part of the lemma, $\phi_{\mathcal{Z}}(s_1^{\text{tel}}) \sim \phi_{\mathcal{Z}} \circ P_1^*$. Now,
 1773 suppose the induction holds up to step h . Then, $\tilde{s}_h^{\text{tel}} \sim P_h^*$, as $a_h^{\text{tel}} \sim \pi_h^*(a_h^{\text{tel}})$, then $s_{h+1}^{\text{tel}} = F_h(\tilde{s}_h^{\text{tel}}, a_h^{\text{tel}}) \sim P_{h+1}^*$. Again
 1774 $\tilde{s}_{h+1}^{\text{tel}} \sim W_{\odot,h+1}^*(s_{h+1}^{\text{tel}})$, so that $\tilde{s}_{h+1}^{\text{tel}}$ has marginal $W_{\odot,h+1}^* \circ P_{h+1}^* = P_{h+1}^*$, as needed. \square

1777 We further show that $W_{\odot,h}$ can be defined to be absolutely continuous with respect to P_h^* .

1778 **Lemma E.4.** *The kernel $W_{\odot,h}$ satisfies that $\phi_{\mathcal{Z}} \circ W_{\odot,h} \ll \phi_{\mathcal{Z}} \circ P_h^*$ as laws, validating the second condition in (E.1). It
 1779 further holds that $\phi_{\mathcal{Z}} \circ W_{\text{dec},h} \ll \phi_{\mathcal{Z}} \circ P_h^*$.*

1782 *Proof.* The first statement follows immediately from Lemma E.3 because these distributions are the same. The second
 1783 statement follows immediately from the tower law of conditional expectation and the definition of $W_{\text{dec},h}$. \square

1785 Lastly, we establish that the replica kernel inherits all concentration properties from the smoothing kernel.

1787 **Lemma E.5 (Replica Concentration).** *Recall that*

$$1788 p_r := \sup_s \mathbb{P}_{s' \sim W_\sigma(s)}[\mathbf{d}_{\text{IPS}}(s', s) > r].$$

1791 *We then have*

$$1792 \mathbb{P}_{s_h \sim P_h^*, \tilde{s}_h \sim W_{\odot,h}^*(s_h)}[\mathbf{d}_{\text{IPS}}(\tilde{s}_h, s_h) > 2r] \leq 2p_r$$

1796 *Proof.* Again, all terms – $W_\sigma, W_{\odot,h}^*, W_{\text{dec},h}^*$ and \mathbf{d}_{IPS} – are compatible with the direct decomposition, it suffices to consider
 1797 the case of the trivial direct decomposition under which $\mathcal{Z} = \mathcal{S}$.

1798 Let \mathbb{P} denote a distribution over $s_h \sim P_h^*$, $s'_h \sim W_\sigma(s_h)$, and $\tilde{s}_h \sim W_{\text{dec},h}^*(s'_h)$. In this special case, we see that
 1799 $\tilde{s}_h | s_h \sim W_{\odot,h}^*(s_h)$ ⁴. By a union bound,

$$\begin{aligned} 1800 \mathbb{P}_{s_h \sim P_h^*, \tilde{s}_h \sim W_{\odot,h}^*(s_h)}[\mathbf{d}_{\text{IPS}}(s_h, \tilde{s}_h) > 2r] &\leq \mathbb{P}[\mathbf{d}_{\text{IPS}}(\tilde{s}_h, s'_h) > r] + \mathbb{P}[\mathbf{d}_{\text{IPS}}(s_h, s'_h) > r] \\ 1801 &= 2 \mathbb{P}[\mathbf{d}_{\text{IPS}}(s_h, s'_h) > r] \leq 2p_r, \end{aligned}$$

1805 where the equality follows from the first statment of Lemma E.3. \square

1807 **Remark E.2.** Note that, in the previous lemma, it suffices that the following weaker condition holds:
 1808 $\mathbb{P}_{s \sim P_h^*, s' \sim W_\sigma(s)}[\mathbf{d}_{\text{IPS}}(s', s) > r] \leq p_r$, i.e. for concentration to hold only in distribution over $s \sim P_h^*$, instead of *uni-*
 1809 *formly* over states.

1811 We now proceed to formally prove Theorem 4

1812 ⁴Notice that, for general $S = \mathcal{Z} \oplus \mathcal{V}$, this condition would become $\phi_{\mathcal{Z}}(\tilde{s}_h) | \phi_{\mathcal{Z}}(s_h) \sim \phi_{\mathcal{Z}} \circ W_{\odot,h}^*(\phi_{\mathcal{Z}}(s_h), \cdot)$, where the \cdot argument
 1813 is irrelevant.

1815 E.2.3. FORMAL PROOF OF THEOREM 4

1816 **Key Events.** For the random variables defined above, we define three groups of events.

- 1817
- 1818
- 1819 • The *coupling events*, denoted by \mathcal{B} , which are controlled by carefully selecting a coupling.
- 1820
- 1821 • The *inductive events*, denoted by \mathcal{C} , which we condition on when bounding the probability of the coupling events.
- 1822
- 1823 • The *stability events*, denoted by \mathcal{Q} , which take advantage of the stability properties of the imitation policy.

1824 **Definition E.6** (Coupling Events). Define the events

$$\begin{aligned}
 1825 & \mathcal{B}_{\text{tel},h} = \{ \mathbf{a}_h^\circ = \mathbf{a}_h^{\text{tel}}, \phi_{\mathcal{Z}}(\tilde{\mathbf{s}}_h^\circ) = \phi_{\mathcal{Z}}(\tilde{\mathbf{s}}_h^{\text{tel}}) \} \\
 1826 & \mathcal{B}_{\text{est},h} = \{ \vec{d}_{\mathcal{A}}(\hat{\mathbf{a}}_h^{\text{tel,inter}}, \mathbf{a}_h^{\text{tel}}) \not\leq \vec{\varepsilon} \} \\
 1827 & \mathcal{B}_{\text{inter},h} = \{ \hat{\mathbf{a}}_h^{\text{tel,inter}} = \hat{\mathbf{a}}_h^{\circ,\text{inter}} \} \\
 1828 & \mathcal{B}_{\hat{\mathbf{a}},h} = \{ \hat{\mathbf{a}}_h^{\circ,\text{inter}} = \hat{\mathbf{a}}_h \} \\
 1829 & \mathcal{B}_{\text{all},h} = \mathcal{B}_{\text{inter},h} \cap \mathcal{B}_{\text{tel},h} \cap \mathcal{B}_{\text{est},h} \cap \mathcal{B}_{\hat{\mathbf{a}},h} \\
 1830 & \bar{\mathcal{B}}_{\text{all},h} = \bigcap_{j=1}^h \mathcal{B}_{\text{all},h}
 \end{aligned}$$

1831 Notice that each of the events above are \mathcal{F}_h -measurable. Moreover, note that on $\bar{\mathcal{B}}_{\text{all},h}$, $\max_{1 \leq j \leq h} \phi_{\text{IS}}(\hat{\mathbf{a}}_j, \mathbf{a}_j^\circ) \leq \varepsilon$.

1832 **Definition E.7** (Inductive Event). Define the events

$$\begin{aligned}
 1833 & \mathcal{C}_{\tilde{\mathbf{s}},h} = \{ \vec{d}_{\mathcal{S}}(\mathbf{s}_h^\circ, \hat{\mathbf{s}}_h) \leq \vec{\varepsilon} \}, \\
 1834 & \mathcal{C}_{\text{tel},h} = \{ \vec{d}_{\mathcal{S}}(\mathbf{s}_h^\circ, \mathbf{s}_h^{\text{tel}}) \leq \vec{\gamma}_{\text{IPS}}(2r) \} \\
 1835 & \mathcal{C}_{\text{all},h} := \mathcal{C}_{\tilde{\mathbf{s}},h} \cap \mathcal{C}_{\text{tel},h} \\
 1836 & \bar{\mathcal{C}}_{\text{all},h} = \bigcap_{j=1}^h \mathcal{C}_{\text{all},j}
 \end{aligned}$$

1837 Notice that all the above events are \mathcal{F}_{h-1} -measurable, due to determinism of the dynamics. Note that also $\mathbb{P}_\mu[\bar{\mathcal{C}}_{\text{all},1}] = 1$ for any μ that respects the construction (as $\mathbf{s}_1^\circ = \mathbf{s}_1^{\text{tel}} = \hat{\mathbf{s}}_1$).

1838 **Definition E.8** (Stability Events). Define the events

$$\begin{aligned}
 1839 & \mathcal{Q}_{\text{close}} := \{ \forall h \in [H] : d_{\text{IPS}}(\mathbf{s}_h^\circ, \tilde{\mathbf{s}}_h^\circ) \leq 2r \} \\
 1840 & \mathcal{Q}_{\text{IS}} := \{ (\mathbf{s}_{1:H+1}^\circ, \mathbf{a}_{1:H}^\circ) \text{ is input-stable w.r.t. } (\vec{d}_{\mathcal{S}}, \vec{d}_{\mathcal{A}}) \} \\
 1841 & \mathcal{Q}_{\text{IPS}} := \{ \vec{d}_{\mathcal{S}}(F_h(\tilde{\mathbf{s}}_h^\circ, \mathbf{a}_h^\circ), \mathbf{s}_{h+1}^\circ) \leq \vec{\gamma}_{\text{IPS}} \circ d_{\text{IPS}}(\tilde{\mathbf{s}}_h^\circ, \mathbf{s}_h^\circ), \quad 1 \leq j \leq H \} \\
 1842 & \mathcal{Q}_{\text{all}} := \mathcal{Q}_{\text{IPS}} \cap \mathcal{Q}_{\text{close}}.
 \end{aligned}$$

1843 In words, $\mathcal{Q}_{\text{close}}$ the event on which \mathbf{s}_h° and $\tilde{\mathbf{s}}_h^\circ \sim W_{\circ,h}^*(\mathbf{s}_h^{\text{tel}})$ are close, and \mathcal{Q}_{IS} and \mathcal{Q}_{IPS} ensure consequences of (vector) input-stability and (vector) input process stability holds.

1844 **Steps of the proof.** First, we use stability to reduce the event $\bar{\mathcal{C}}_{\text{all},h+1}$ to $\bar{\mathcal{C}}_{\text{all},h} \cap \bar{\mathcal{B}}_{\text{all},h}$:

1845 **Claim E.6** (Stability Claim). *By construction,*

$$1846 \bar{\mathcal{C}}_{\text{all},h+1} \subset \mathcal{Q}_{\text{all}} \cap \bar{\mathcal{C}}_{\text{all},h} \cap \bar{\mathcal{B}}_{\text{all},h}.$$

1867

1870 *Proof.* It suffices to show that on $\mathcal{Q}_{\text{all}} \cap \bar{\mathcal{C}}_{\text{all},h} \cap \bar{\mathcal{B}}_{\text{all},h}$, $\bar{d}_{\mathcal{S}}(s_{h+1}^{\circ}, \hat{s}_{h+1}) \leq \bar{\varepsilon}$ and $\bar{d}_{\mathcal{S}}(s_{h+1}^{\circ}, s_{h+1}^{\text{tel}}) \leq \bar{\gamma}_{\text{IPS}}(2r)$. By applying
 1871 the event \mathcal{Q}_{IS} to the sequence $a'_h = \hat{a}_h$ and $s'_h = \hat{s}_h$, we have that on $\mathcal{Q}_{\text{all}} \subset \mathcal{Q}_{\text{IS}}$ that

$$1872 \quad \forall h \in [H], i \in [K], \quad d_{\mathcal{S},i}(s_{h+1}^{\circ}, \hat{s}_{h+1}) \leq \max_{1 \leq j \leq h} d_{\mathcal{A},i}(a_j^{\circ}, \hat{a}_j)$$

1873
 1874
 1875 For the next point, note that the compatibility of the dynamics with the direct decomposition $\mathcal{S} = \mathcal{Z} \oplus \mathcal{V}$ implies that there
 1876 exists a dynamics map $F_h^{\mathcal{Z}}$ for which

$$1877 \quad F_h(s, a) = F_h^{\mathcal{Z}}(\phi_{\mathcal{Z}}(s), a).$$

1878
 1879 Similarly, by applying \mathcal{Q}_{IPS} and $\mathcal{Q}_{\text{close}}$ and the event $\{\phi_{\mathcal{Z}}(\tilde{s}_h^{\circ}) = \phi_{\mathcal{Z}}(\tilde{s}_h^{\text{tel}}), a_h^{\text{tel}} = a_h^{\circ}\}$ on $\mathcal{B}_{\text{tel},h}$, it holds that on $\mathcal{Q}_{\text{all}} \cap$
 1880 $\bar{\mathcal{C}}_{\text{all},h} \cap \bar{\mathcal{B}}_{\text{all},h}$ that, for all $h \in [H]$,

$$\begin{aligned} 1881 \quad \bar{d}_{\mathcal{S}}(s_{h+1}^{\circ}, F_h(\tilde{s}_h^{\circ}, a_h^{\circ})) &= \bar{d}_{\mathcal{S}}(s_{h+1}^{\circ}, F_h^{\mathcal{Z}}(\phi_{\mathcal{Z}}(\tilde{s}_h^{\circ}), a_h^{\circ})) \\ 1882 \quad &= \bar{d}_{\mathcal{S}}(s_{h+1}^{\circ}, F_h^{\mathcal{Z}}(\phi_{\mathcal{Z}}(\tilde{s}_h^{\text{tel}}), a_h^{\text{tel}})) && (\mathcal{B}_{\text{tel},h}) \\ 1883 \quad &= \bar{d}_{\mathcal{S}}(s_{h+1}^{\circ}, F_h(\tilde{s}_h^{\text{tel}}, a_h^{\text{tel}})) \\ 1884 \quad &= \bar{d}_{\mathcal{S}}(s_{h+1}^{\circ}, s_{h+1}^{\text{tel}}) \\ 1885 \quad &\leq \bar{\gamma}_{\text{IPS}} \circ d_{\text{IPS}}(s_j^{\text{tel}}, \tilde{s}_j^{\text{tel}}) && (\mathcal{Q}_{\text{IPS}}) \\ 1886 \quad &\leq \bar{\gamma}_{\text{IPS}} \circ d_{\text{IPS}}(2r). && (\mathcal{Q}_{\text{close}}) \end{aligned}$$

1887
 1888
 1889
 1890
 1891
 1892
 1893
 1894
 1895
 1896
 1897
 1898
 1899
 1900
 1901
 1902
 1903
 1904
 1905
 1906
 1907
 1908
 1909
 1910
 1911
 1912
 1913
 1914
 1915
 1916
 1917
 1918
 1919
 1920
 1921
 1922
 1923
 1924

From [Claim E.6](#), we decompose our error probability as follows:

Lemma E.7 (Key Error Decomposition). *Let μ respect the construction (in the sense of [Appendix E.2.1](#)). Then*

$$\begin{aligned} \mathbb{P}_{\mu}[\exists h \in [H] : \max\{d_{\mathcal{S}}(s_{h+1}^{\circ}, \hat{s}_{h+1}), \phi_{\text{IS}}(a_h^{\circ}, \hat{a}_h)\} > \varepsilon] \\ \leq \mathbb{P}_{\mu}[\mathcal{Q}_{\text{all}}^c] + \sum_{h=1}^H \mathbb{P}_{\mu}[\mathcal{B}_{\text{all},h}^c \cap \bar{\mathcal{C}}_{\text{all},h} \cap \bar{\mathcal{B}}_{\text{all},h-1}] \end{aligned}$$

Hence, letting \inf_{μ} denote the infimum over couplings μ which respect the construction,

$$\begin{aligned} \bar{\Gamma}_{\text{joint},\bar{\varepsilon}}(\hat{\pi}_{\sigma} \parallel \pi_{\circ\sigma}^*) \vee \bar{\Gamma}_{\text{marg},\bar{\varepsilon}_{\text{marg}}}(\hat{\pi}_{\sigma} \parallel \pi^*) \\ \leq \inf_{\mu} \left\{ \mathbb{P}_{\mu}[\mathcal{Q}_{\text{all}}^c] + \sum_{h=1}^H \mathbb{P}_{\mu}[\mathcal{B}_{\text{all},h}^c \cap \bar{\mathcal{C}}_{\text{all},h} \cap \bar{\mathcal{B}}_{\text{all},h-1}] \right\} \end{aligned} \quad (\text{E.4})$$

Proof. Define the events $\mathcal{E}_h := \bar{\mathcal{C}}_{\text{all},h+1} \cap \bar{\mathcal{B}}_{\text{all},h}$. Observe that the events are nested: $\mathcal{E}_h \supset \mathcal{E}_{h+1}$, and that on \mathcal{E}_H , we have that for all $h \in [H]$

$$\begin{aligned} \bar{d}_{\mathcal{S}}(s_{h+1}^{\circ}, \hat{s}_{h+1}) \vee \bar{d}_{\mathcal{A}}(a_h^{\circ}, \hat{a}_h) &\leq \bar{\varepsilon} \vee \bar{d}_{\mathcal{A}}(a_h^{\circ}, \hat{a}_h) && (\mathcal{C}_{\hat{s},h+1} \supset \bar{\mathcal{C}}_{\text{all},h+1} \supset \mathcal{E}_h) \\ &\leq \bar{\varepsilon}. && (\bar{\mathcal{B}}_{\text{all},h} \supset \mathcal{E}_h) \end{aligned}$$

Thus,

$$\mathbb{P}_{\mu}[\exists h \in [H] : \bar{d}_{\mathcal{A}}(s_{h+1}^{\circ}, \hat{s}_{h+1}) \vee \bar{d}_{\mathcal{A}}(a_h^{\circ}, \hat{a}_h) \not\leq \bar{\varepsilon}] \leq \mathbb{P}_{\mu}[\mathcal{E}_H^c] \leq \mathbb{P}_{\mu}[(\mathcal{Q}_{\text{all}} \cap \mathcal{E}_H)^c] \quad (\text{E.5})$$

As $(s_{1:H+1}^{\circ}, a_{1:H}^{\circ}) \sim \mathbb{D}_{\pi_{\circ\sigma}^*}$, this shows $\bar{\Gamma}_{\text{joint},\bar{\varepsilon}}(\hat{\pi}_{\sigma} \parallel \pi_{\circ\sigma}^*) \leq \mathbb{P}_{\mu}[(\mathcal{Q}_{\text{all}} \cap \mathcal{E}_H)^c]$. Moreover, on $\mathcal{Q}_{\text{all}} \cap \mathcal{E}_H$, we have that

$$\max_h \bar{d}_{\mathcal{S}}(s_h^{\circ}, s_h^{\text{tel}}) \leq \bar{\gamma}_{\text{IPS}}(2r),$$

so that, by the inequality preceding [\(E.5\)](#), the following holds for all $h \in [H]$ on $\mathcal{Q}_{\text{all}} \cap \mathcal{E}_H$.

$$\bar{d}_{\mathcal{S}}(s_{h+1}^{\circ}, \hat{s}_{h+1}) \vee \bar{d}_{\mathcal{A}}(a_h^{\circ}, \hat{a}_h) \leq \bar{d}_{\mathcal{S}}(s_{h+1}^{\circ}, \hat{s}_{h+1}) \vee \bar{d}_{\mathcal{A}}(a_h^{\circ}, \hat{a}_h) \leq \bar{\varepsilon}. \quad (\text{E.6})$$

By construction, for each h , $\mathbf{a}_h^{\text{tel}} \mid \mathcal{F}_h \sim \pi_{\circlearrowleft \sigma, h}^*(s_h^{\text{tel}})$. Moreover, [Lemma E.3](#) implies that s_h^{tel} has the marginal distribution of $s_h^* \sim P_h^*$. Thus, for each h , s_{h+1}^{tel} and $\mathbf{a}_h^{\text{tel}}$ have the same *marginals* as the marginals under D_{π^*} . Consequently, [\(E.6\)](#) implies that,

$$\begin{aligned} \bar{\Gamma}_{\text{marg}, \bar{\varepsilon}_{\text{marg}}}(\hat{\pi}_\sigma \parallel \pi^*) &:= \max_{h \in [H]} \max_{\mu_1} \left\{ \inf_{\mu_1} \mathbb{P}_{\mu_1} \left[\bar{\mathbf{d}}_{\mathcal{S}}(\hat{s}_{h+1}, s_{h+1}^*) \not\leq \bar{\varepsilon} \right], \inf_{\mu_1} \mathbb{P}_{\mu_1} \left[\bar{\mathbf{d}}_{\mathcal{A}}(\mathbf{a}_h^*, \hat{\mathbf{a}}_h) \not\leq \bar{\varepsilon} \right] \right\} \\ &\leq \mathbb{P}_\mu[(\mathcal{Q}_{\text{all}} \cap \mathcal{E}_H)^c]. \end{aligned}$$

where above we take inf over $\mu_1 \in \mathcal{C}(D_{\hat{\pi}_\sigma}, D_{\pi^*})$. Summarizing our findings thus far,

$$\bar{\Gamma}_{\text{joint}, \bar{\varepsilon}}(\hat{\pi}_\sigma \parallel \pi_{\circlearrowleft \sigma}^*) \vee \bar{\Gamma}_{\text{marg}, \bar{\varepsilon}_{\text{marg}}}(\hat{\pi}_\sigma \parallel \pi^*) \leq \mathbb{P}_\mu[(\mathcal{Q}_{\text{all}} \cap \mathcal{E}_H)^c].$$

Let us conclude by bounding $\mathbb{P}_\mu[(\mathcal{Q}_{\text{all}} \cap \mathcal{E}_H)^c]$. Using the nesting structure $\mathcal{E}_h = \bigcap_{j=1}^h \mathcal{E}_j$, the peeling lemma, [Lemma C.11](#), and a union bound, it holds that

$$\begin{aligned} \mathbb{P}_\mu[(\mathcal{Q}_{\text{all}} \cap \mathcal{E}_H)^c] &\leq \mathbb{P}_\mu[\mathcal{Q}_{\text{all}}^c] + \mathbb{P}[\exists h \in [H] \text{ s.t. } (\mathcal{Q}_{\text{all}} \cap \mathcal{E}_{h-1} \cap \mathcal{E}_h^c) \text{ holds}] \\ &\leq \mathbb{P}_\mu[\mathcal{Q}_{\text{all}}^c] + \sum_{h=1}^H \mathbb{P}_\mu[\mathcal{Q}_{\text{all}} \cap \mathcal{E}_{h-1} \cap \mathcal{E}_h^c \text{ holds}] \\ &= \mathbb{P}_\mu[\mathcal{Q}_{\text{all}}^c] + \sum_{h=1}^H \mathbb{P}_\mu[\mathcal{Q}_{\text{all}} \cap \bar{\mathcal{B}}_{\text{all}, h-1} \cap \bar{\mathcal{C}}_{\text{all}, h} \cap (\bar{\mathcal{B}}_{\text{all}, h} \cap \bar{\mathcal{C}}_{\text{all}, h+1})^c \text{ holds}] \\ &= \mathbb{P}_\mu[\mathcal{Q}_{\text{all}}^c] + \sum_{h=1}^H \mathbb{P}_\mu[\mathcal{Q}_{\text{all}} \cap \bar{\mathcal{B}}_{\text{all}, h-1} \cap \bar{\mathcal{C}}_{\text{all}, h} \cap \bar{\mathcal{B}}_{\text{all}, h}^c] \\ &= \mathbb{P}_\mu[\mathcal{Q}_{\text{all}}^c] + \sum_{h=1}^H \mathbb{P}_\mu[\mathcal{Q}_{\text{all}} \cap \bar{\mathcal{B}}_{\text{all}, h-1} \cap \bar{\mathcal{C}}_{\text{all}, h} \cap \mathcal{B}_{\text{all}, h}^c], \end{aligned}$$

where the last step invokes [Claim E.6](#). □

Next, we bound the contribution of $\mathbb{P}_\mu[\mathcal{Q}_{\text{all}}^c]$ in [\(E.4\)](#), uniformly over all couplings.

Lemma E.8. *For all μ which respect the construction,*

$$\mathbb{P}_\mu[\mathcal{Q}_{\text{all}}^c] \leq p_{\text{IPS}} + 2Hp_r.$$

Proof. $\mathbb{P}_\mu[\mathcal{Q}_{\text{close}}^c] = \mathbb{P}_\mu[\exists h : d_{\text{IPS}}(s_h^{\text{tel}}, \tilde{s}_h^{\text{tel}}) > 2r] \leq 2Hp_r$ by [Lemma E.5](#) and a union bound.

Let us now bound $\mathbb{P}_\mu[\mathcal{Q}_{\text{close}} \cap \mathcal{Q}_{\text{IPS}}^c] \leq \mathbb{P}_\mu[\mathcal{Q}_{\text{IPS}}^c \mid \mathcal{Q}_{\text{close}}]$. Define the kernels $W_h(s)$ to be equal to the kernel $W_{\circlearrowleft, h}(s)$ conditioned on the event $s' \sim W_{\circlearrowleft, h}(s)$ satisfies $d_{\text{IPS}}(s', s) \leq 2r$. Then, conditional on $\mathcal{Q}_{\text{close}}$, we see that the sequence $(s_{1:H}^{\circlearrowleft}, \tilde{s}_{1:H}^{\circlearrowleft}, \mathbf{a}_{1:H}^{\circlearrowleft})$ obeys the generative process

$$\tilde{s}_h^{\circlearrowleft} \mid \tilde{s}_{1:h-1}^{\circlearrowleft}, s_{1:h}^{\circlearrowleft}, \mathbf{a}_{1:h-1}^{\circlearrowleft} \sim W_h(s), \quad \mathbf{a}_h^{\circlearrowleft} \mid \tilde{s}_{1:h}^{\circlearrowleft}, s_{1:h}^{\circlearrowleft}, \mathbf{a}_{1:h-1}^{\circlearrowleft} \sim \pi_h^*(\tilde{s}_h^{\circlearrowleft}), \quad s_{h+1}^{\circlearrowleft} = F_h(s_h^{\circlearrowleft}, \mathbf{a}_h^{\circlearrowleft}).$$

By construction, for each h , $\mathbb{P}_{s' \sim W_{\circlearrowleft, h}(s)}[d_{\text{IPS}}(s', s) > 2r] = 0$. Thus, the definition of (vector) input process stability ([Definition E.4](#)) and assumption $r \leq \frac{1}{2}r_{\text{IPS}}$ implies that $\mathbb{P}_\mu[\mathcal{Q}_{\text{IPS}}^c \mid \mathcal{Q}_{\text{close}}] \leq p_{\text{IPS}}$. □

The remaining step of the proof is therefore to bound the second term in [\(E.4\)](#).

Lemma E.9. *There exists a coupling μ which respects the construction and satisfies the following for any $h \in [H]$*

$$\begin{aligned} &\mathbb{P}_\mu[\mathcal{B}_{\text{all}, h}^c \mid \mathcal{F}_{h-1}] \\ &\leq \hat{\gamma} \circ d_{\text{TVC}}(s_h^{\circlearrowleft}, \hat{s}_h) + (\hat{\gamma} + \gamma_\sigma) \circ d_{\text{TVC}}(s_h^{\circlearrowleft}, s_h^{\text{tel}}) + \bar{\mathbf{d}}_{\text{os}, \bar{\varepsilon}}(\hat{\pi}_{\sigma, h}(s_h^{\text{tel}}) \parallel \pi_{\circlearrowleft \sigma, h}^*(s_h^{\text{tel}})), \quad \mu\text{-almost surely} \end{aligned}$$

Consequently, for all $h \in [H]$,

$$\begin{aligned} &\mathbb{P}_\mu[\mathcal{B}_{\text{all}, h}^c \cap \bar{\mathcal{C}}_{\text{all}, h} \cap \bar{\mathcal{B}}_{\text{all}, h-1}] \\ &\leq \hat{\gamma}(\bar{\varepsilon}_1) + (\hat{\gamma} + \gamma_\sigma) \circ \gamma_{\text{IPS}, 1}(2r) + \mathbb{E}_\mu[\bar{\mathbf{d}}_{\text{os}, \bar{\varepsilon}}(\hat{\pi}_{\sigma, h}(s_h^{\text{tel}}) \parallel \pi_{\circlearrowleft \sigma, h}^*(s_h^{\text{tel}}))] \end{aligned}$$

Moreover, $s \mapsto \bar{\mathbf{d}}_{\text{os}, \bar{\varepsilon}}(\hat{\pi}_{\sigma, h}(s) \parallel \pi_{\circlearrowleft \sigma, h}^*(s))$ is measurable.

1980 *Proof Sketch.* We begin by giving a high level overview of the construction, which is done recursively. The key technical
 1981 tool is [Lemma C.2](#) above, which allows us to transform any coupling μ between random variables (X, Y) into a probability
 1982 kernel $\mu(\cdot|X)$ mapping instances of X to probability distributions on Y such that $(X, Y) \sim \mu$ has the same law as
 1983 $(X, Y \sim \mu(\cdot|X))$. For each h , we then show that, assuming the coupling has kept the states and controls close together until
 1984 time $h - 1$, this will imply the following chain:

$$1985 \underbrace{(a^\circ \leftrightarrow a^{\text{tel}})}_{\gamma_{\text{TVC}} \text{ and induction}} \rightarrow \underbrace{(a^{\text{tel}} \leftrightarrow \hat{a}^{\text{tel,inter}})}_{\text{learning and sampling}} \rightarrow \underbrace{(\hat{a}^{\text{tel,inter}} \leftrightarrow \hat{a}^{\circ, \text{inter}})}_{\gamma_{\text{TVC}} \text{ and induction}} \rightarrow \underbrace{(\hat{a}^{\circ, \text{inter}} \leftrightarrow \hat{a})}_{\gamma_{\text{TVC}} \text{ and induction}},$$

1988 where the bidirectional arrows indicate individual couplings between the laws of the random variables that are constructed
 1989 by the method outlined in text below and the single directional arrows denote the probability kernels described above. The
 1990 full proof of the lemma is given in [Appendix E.2.4](#). \square

1992 **Concluding the proof.** Here, we finish the proof of [Theorem 4](#). Recall that we wish to bound $\vec{\Gamma}_{\text{joint}, \vec{\varepsilon}}(\hat{\pi}_\sigma \parallel \pi_{\circlearrowleft \sigma}^*) \vee$
 1993 $\vec{\Gamma}_{\text{marg}, \vec{\varepsilon}_{\text{marg}}}(\hat{\pi}_\sigma \parallel \pi^*)$. We begin by bounding $\vec{\Gamma}_{\text{joint}, \vec{\varepsilon}}(\hat{\pi}_\sigma \parallel \pi_{\circlearrowleft \sigma}^*) \vee \vec{\Gamma}_{\text{marg}, \vec{\varepsilon}_{\text{marg}}}(\hat{\pi}_\sigma \parallel \pi_{\circlearrowleft \sigma}^*)$. In light of [Lemma E.7](#), it
 1994 suffices to bound

$$1996 \mathbb{P}_\mu[\mathcal{Q}_{\text{all}}^c] + \sum_{h=1}^H \mathbb{P}_\mu[\bar{\mathcal{B}}_{\text{all}, h}^c \cap \bar{\mathcal{C}}_{\text{all}, h} \cap \bar{\mathcal{B}}_{\text{all}, h-1}],$$

1999 where μ is the coupling in [Lemma E.9](#). Applying [Lemma E.8](#) and [Lemma E.9](#),

$$2001 \mathbb{P}_\mu[\mathcal{Q}_{\text{all}}^c] + \sum_{h=1}^H \mathbb{P}_\mu[\bar{\mathcal{B}}_{\text{all}, h}^c \cap \bar{\mathcal{C}}_{\text{all}, h} \cap \bar{\mathcal{B}}_{\text{all}, h-1}]$$

$$2004 \leq p_{\text{IPS}} + 2H p_r + \sum_{h=1}^H \mathbb{P}_\mu[\bar{\mathcal{B}}_{\text{all}, h}^c \cap \bar{\mathcal{C}}_{\text{all}, h} \cap \bar{\mathcal{B}}_{\text{all}, h-1}]$$

$$2007 \leq p_{\text{IPS}} + H(2p_r + \hat{\gamma}(\vec{\varepsilon}_1) + (\hat{\gamma} + \gamma_\sigma) \circ \gamma_{\text{IPS}, 1}(2r)) + \sum_{h=1}^H \mathbb{E}_{s_h^{\text{tel}} \sim \mu} \vec{d}_{\text{os}, \vec{\varepsilon}}(\hat{\pi}_{\sigma, h}(s_h^{\text{tel}}) \parallel \pi_{\circlearrowleft \sigma, h}^*(s_h^{\text{tel}}))$$

2009 To conclude, we note that for any μ which respects the construction, [Lemma E.3](#) ensures that s_h^{tel} as the marginal distribution
 2010 of $s_h^* \sim \pi_h^*$. Thus, the above is at most

$$2012 p_{\text{IPS}} + H(2p_r + \hat{\gamma}(\vec{\varepsilon}_1) + (\hat{\gamma} + \gamma_\sigma) \circ \gamma_{\text{IPS}, 1}(2r)) + \sum_{h=1}^H \mathbb{E}_{s_h^* \sim \mathcal{P}_h^*} \vec{d}_{\text{os}, \vec{\varepsilon}}(\hat{\pi}_{\sigma, h}(s_h^*) \parallel \pi_{\circlearrowleft \sigma, h}^*(s_h^*)) \quad (\text{E.7})$$

2015 which concludes the proof of [\(E.2\)](#) for $\vec{\Gamma}_{\text{joint}, \vec{\varepsilon}}(\hat{\pi} \parallel \pi_\circlearrowleft^*)$.

2017 To prove [\(E.3\)](#) for $\vec{\Gamma}_{\text{joint}, \vec{\varepsilon}}(\hat{\pi} \parallel \pi_\circlearrowleft^*)$, we consider the special case that $\hat{\pi}_\sigma = \hat{\pi} \circ W_\sigma$. By definition, $\hat{\pi}_{\sigma, h} = \hat{\pi} \circ W_\sigma$. Thus,
 2018 the data-processing inequality for optimal transport ([Lemma C.5](#))

$$2019 \vec{d}_{\text{os}, \vec{\varepsilon}}(\hat{\pi}_{\sigma, h}(s_h^*) \parallel \pi_{\circlearrowleft \sigma, h}^*(s_h^*)) \leq \mathbb{E}_{s_h' \sim W_\sigma(s_h^*)} \vec{d}_{\text{os}, \vec{\varepsilon}}(\hat{\pi}(s_h') \parallel \pi_{\text{dec}, h}^*(s_h')),$$

2021 for all s_h^* . Substituting this into [\(E.7\)](#), and setting $\hat{\gamma} = \gamma_\sigma$ (in view of [Lemma E.2](#)), finishes the argument.

2023 E.2.4. PROOF OF LEMMA E.9

2024 Recall that [Assumption E.1](#) ensures all of the general measure-theoretic guarantees of [Appendix C](#) hold true in our
 2025 setting. Notably we need the gluing lemma ([Lemma C.2](#)) and the commuting of optimal transport metrics and conditional
 2026 probabilities ([Proposition C.3](#)).

2028 **Proof strategy.** Our proof follows along similar lines as that of [Proposition D.1](#), although with the added complication
 2029 of including the smoothing. We will inductively construct μ . A useful schematic for the construction at each step is the
 2030 following diagram:

$$2032 \underbrace{(\tilde{s}^\circ \leftrightarrow \tilde{s}^{\text{tel}}), (a^\circ \leftrightarrow a^{\text{tel}})}_{\mathcal{B}_{\text{tel}, h}} \rightarrow \underbrace{(a^{\text{tel}} \leftrightarrow \hat{a}^{\text{tel,inter}})}_{\mathcal{B}_{\text{est}, h}} \rightarrow \underbrace{(\hat{a}^{\text{tel,inter}} \leftrightarrow \hat{a}^{\circ, \text{inter}})}_{\mathcal{B}_{\text{inter}, h}} \rightarrow \underbrace{(\hat{a}^{\circ, \text{inter}} \leftrightarrow \hat{a})}_{\mathcal{B}_{\hat{a}, h}},$$

2034

where the events under each bidirectional arrow refer to the event such ensuring that there exists a coupling such that the objects are close. We then will apply [Lemma C.2](#) to glue the individual couplings together. We will then use [Lemma C.11](#) and a union bound to control the probability under our constructed coupling that any of the relevant events fail to hold, concluding the proof.

Recursive construction of μ . Let $h \geq 1$, and suppose that we have constructed the coupling $\mu^{(1:h-1)}$ for steps $1, \dots, h-1$ which respects the construction. Recall that \mathcal{F}_h denotes the sigma-algebra generated by $(\hat{s}_{1:h}, s_{1:h}^\circ, s_{1:h}^{\text{tel}}, (a_{1:h}^\circ, \tilde{s}_{1:h}^\circ, \tilde{s}_{1:h}^{\text{tel}}, a_{1:h}^{\text{tel}}, \hat{a}_{1:h}),$ and $(\hat{a}_{1:h}^{\circ, \text{inter}}, \hat{a}_{1:h}^{\text{tel, inter}})$. Notice that $s_{h+1}^{\text{tel}}, s_{h+1}^\circ, \hat{s}_{h+1}$ are determined by \mathcal{F}_h as well. Similarly, it can be seen from [Definition E.5](#) that $\phi_{\mathcal{V}}(\tilde{s}_{h+1}^{\text{tel}})$ and $\phi_{\mathcal{V}}(\tilde{s}_{h+1}^\circ)$ are also determined by \mathcal{F}_h (since the replica kernel preserves the \mathcal{V} -components). We summarize all these aforementioned variables in a random variable Y_h . Let \mathcal{F}_0 denote the filtration generated by $s_1^\circ = s_1^{\text{tel}} = \hat{s}_1$. We let $Y_0 = (s_1^\circ, s_1^{\text{tel}}, \hat{s}_1)$.

Correspondingly, let Z_h denote the random variables $(a_h^\circ, \phi_{\mathcal{Z}}(\tilde{s}_h^\circ), \phi_{\mathcal{Z}}(\tilde{s}_h^{\text{tel}}), a_h^{\text{tel}}, \hat{a}_h)$, and $(\hat{a}_h^{\circ, \text{inter}}, \hat{a}_h^{\text{tel, inter}})$ such that the joint law of these random variables respects the construction. Our goal is then to specify, for each $h \in [H]$, a joint distribution of (Y_{h-1}, Z_h) . Note that Z_h, Y_{h-1} determines Y_h , and we call this induced law $\mu^{(h)}$.

We begin by specifying joint distributions conditional on Y_{h-1} and subsets of Z_h , then glue them together by the gluing lemma. Below, we use information-theoretic notation.

- By total variation continuity of $\phi_{\mathcal{Z}} \circ W_{\circ, h}^*$ ([Lemma E.2](#)),

$$\text{TV}(\mathbb{P}_{\phi_{\mathcal{Z}}(\tilde{s}_h^\circ) | Y_{h-1}}, \mathbb{P}_{\phi_{\mathcal{Z}}(\tilde{s}_h^{\text{tel}}) | Y_{h-1}}) \leq \gamma_\sigma \circ d_{\text{TVC}}(s_h^\circ, s_h^{\text{tel}}).$$

Because $a_h^\circ \sim \pi_h^*(\tilde{s}_{h+1}^\circ)$ and $a_h^{\text{tel}} \sim \pi_h^*(\tilde{s}_{h+1}^{\text{tel}})$, and π^* is compatible with the decomposition $\mathcal{S} = \mathcal{Z} \oplus \mathcal{V}$ (i.e. $\pi_h^*(s)$ is a function of $\phi_{\mathcal{Z}}(s)$) [Lemma C.4](#) implies that (almost surely)

$$\text{TV}(\mathbb{P}_{(a_h^\circ, \phi_{\mathcal{Z}}(\tilde{s}_h^\circ) | Y_{h-1}), \mathbb{P}_{(a_h^{\text{tel}}, \phi_{\mathcal{Z}}(\tilde{s}_h^{\text{tel}}) | Y_{h-1})} \leq \gamma_\sigma \circ d_{\text{TVC}}(s_h^\circ, s_h^{\text{tel}}).$$

Hence, [Corollary C.1](#) implies that there exists a coupling $\mu_{\text{tel}}^{(h)}$ over $Y_{h-1}, (\phi_{\mathcal{Z}}(\tilde{s}_h^\circ), a_h^\circ), (\phi_{\mathcal{Z}}(\tilde{s}_h^{\text{tel}}), a_h^{\text{tel}})$ respecting the construction such that $Y_h \sim \mu^{(h-1)}$ and such that (almost surely)

$$\mathbb{E}_{\mu_{\text{tel}}^{(h)}}[\mathcal{B}_{\text{tel}, h} | Y_{h-1}] = \mathbb{P}_{\mu_{\text{tel}}^{(h)}}[(\phi_{\mathcal{Z}}(\tilde{s}_h^\circ), a_h^\circ) \neq (\phi_{\mathcal{Z}}(\tilde{s}_h^{\text{tel}}), a_h^{\text{tel}}) | Y_{h-1}] \leq d_{\text{TVC}}(s_h^\circ, s_h^{\text{tel}}).$$

- In our construction, $a_h^{\text{tel}} | Y_{h-1} \sim \pi_{\circ, h}^*(s_h^{\text{tel}})$, and $\hat{a}_h^{\text{tel, inter}} | Y_{h-1} \sim \hat{\pi}_{\sigma, h}(s_h^{\text{tel}})$. Thus, by definition of $\vec{d}_{\text{os}, \varepsilon}$, and the assumption $\mathbf{I}\{\vec{d}_{\mathcal{A}}(\cdot, \cdot) \not\leq \varepsilon\}$ is lower semicontinuous, [Proposition C.3](#) implies that we may find a coupling $\mu_{\text{est}}^{(h)}$ of $(a_h^{\text{tel}}, \hat{a}_h^{\text{tel, inter}}, Y_{h-1})$ respecting the construction such that, almost surely,

$$\begin{aligned} \mathbb{P}_{\mu_{\text{est}}^{(h)}}[\mathcal{B}_{\text{est}, h}^c | Y_{h-1}] &= \mathbb{P}_{\mu_{\text{est}}^{(h)}}[\vec{d}_{\mathcal{A}}(\hat{a}_h^{\text{tel, inter}}, a_h^{\text{tel}}) \not\leq \varepsilon | Y_{h-1}] \\ &= \vec{d}_{\text{os}, \varepsilon}(\hat{\pi}_{\sigma, h}(s_h^{\text{tel}}) \| \pi_{\circ, h}^*(s_h^{\text{tel}})). \end{aligned}$$

Moreover, that same proposition ensures measurability of $s \rightarrow \vec{d}_{\text{os}, \varepsilon}(\hat{\pi}_{\sigma, h}(s) \| \pi_{\circ, h}^*(s))$.

- Since $\hat{a}_h^{\text{tel, inter}} | \mathcal{F}_h \sim \hat{\pi}_{\sigma, h}(s_h^{\text{tel}})$ and $\hat{a}_{h+1}^{\circ, \text{inter}} | \mathcal{F}_h \sim \hat{\pi}_{\sigma, h}(s_h^\circ)$, and since $\hat{\pi}_{\sigma, h}(\cdot)$ is $\hat{\gamma}$ -TVC by assumption,

$$\text{TV}(\mathbb{P}_{\hat{a}_h^{\text{tel, inter}} | Y_{h-1}}, \mathbb{P}_{\hat{a}_{h+1}^{\circ, \text{inter}} | Y_{h-1}}) \leq \hat{\gamma} \circ d_{\text{TVC}}(s_h^\circ, s_h^{\text{tel}}).$$

[Corollary C.1](#) implies that there is a coupling $\mu_{\text{inter}}^{(h)}$ between $(\hat{a}_h^{\text{tel, inter}}, \hat{a}_h^{\circ, \text{inter}}, Y_{h-1})$ such that

$$\mathbb{P}_{\mu_{\text{inter}}^{(h)}}[\mathcal{B}_{\text{inter}, h}^c | Y_{h-1}] = \mathbb{P}_{\mu_{\text{inter}}^{(h)}}[\hat{a}_h^{\text{tel, inter}} \neq \hat{a}_h^{\circ, \text{inter}} | Y_{h-1}] \leq \hat{\gamma} \circ d_{\text{TVC}}(s_h^{\text{tel}}, s_h^\circ)$$

- Similarly, since $\hat{a}_h^{\circ, \text{inter}} | \mathcal{F}_{h-1} \sim \hat{\pi}_h(s_h^\circ)$ and $\hat{a}_{h+1} | \mathcal{F}_{h-1} \sim \hat{\pi}_h(\hat{s}_h)$, $\hat{\pi}_h(\cdot)$ is $\hat{\gamma}$ -TVC, [Corollary C.1](#) implies that there is a coupling $\mu_{\hat{a}}^{(h)}$ between $(\hat{a}_h^{\circ, \text{inter}}, \hat{a}_h, Y_{h-1})$ such that

$$\mathbb{P}_{\mu_{\hat{a}}^{(h)}}[\mathcal{B}_{\hat{a}, h}^c | Y_{h-1}] = \mathbb{P}_{\mu_{\hat{a}}^{(h)}}[\hat{a}_h \neq \hat{a}_h^{\circ, \text{inter}} | Y_{h-1}] \leq \hat{\gamma} \circ d_{\text{TVC}}(s_h^\circ, \hat{s}_h)$$

2090 We can then apply the gluing lemma (Lemma C.2) to

$$\begin{aligned}
2091 \quad X_{h,1} &= (\phi_{\mathcal{Z}}(\tilde{\mathbf{s}}_h^{\text{tel}}), \mathbf{a}_h^{\text{tel}}, Y_{h-1}) \\
2092 \quad X_{h,2} &= (\phi_{\mathcal{Z}}(\tilde{\mathbf{s}}_h^{\circ}), \mathbf{a}_h^{\circ}, Y_{h-1}) \\
2093 \quad X_{h,3} &= (\hat{\mathbf{a}}_h^{\text{tel}}, \hat{\mathbf{a}}_h^{\text{tel,inter}}, Y_{h-1}) \\
2094 \quad X_{h,4} &= (\hat{\mathbf{a}}_h^{\text{tel,inter}}, \hat{\mathbf{a}}_h^{\circ,inter}, Y_{h-1}) \\
2095 \quad X_{h,5} &= (\hat{\mathbf{a}}_h^{\circ,inter}, \hat{\mathbf{a}}_h, Y_{h-1})
\end{aligned}$$

2099 with

$$2100 \quad (X_{h,1}, X_{h,2}) \sim \mu_{\text{tel}}^{(h)}, \quad (X_{h,2}, X_{h,3}) \sim \mu_{\text{est}}^{(h)}, \quad (X_{h,3}, X_{h,4}) \sim \mu_{\text{inter}}^{(h)}, \quad (X_{h,4}, X_{h,5}) \sim \mu_{\hat{\mathbf{a}}}^{(h)}.$$

2103 Lemma C.2 guarantees the existence of a coupling $\mu^{(h)}$ consistent with all sub-couplings $\mu_{\text{tel}}^{(h)}, \mu_{\text{est}}^{(h)}, \mu_{\text{inter}}^{(h)}, \mu_{\hat{\mathbf{a}}}^{(h)}$. Then, $\mu^{(h)}$ -almost surely (and using that \mathcal{F}_{h-1} is precisely the σ -algebra generated by Y_{h-1})

$$\begin{aligned}
2104 \quad &\mathbb{P}_{\mu^{(h)}}[\mathcal{B}_{\text{all},h}^c \mid \mathcal{F}_{h-1}] \\
2105 \quad &\leq \mathbb{P}_{\mu^{(h)}}[\mathcal{B}_{\text{tel},h}^c \mid \mathcal{F}_{h-1}] + \mathbb{P}_{\mu^{(h)}}[\mathcal{B}_{\text{est},h}^c \mid \mathcal{F}_{h-1}] + \mathbb{P}_{\mu^{(h)}}[\mathcal{B}_{\text{inter},h}^c \mid \mathcal{F}_{h-1}] + \mathbb{P}_{\mu^{(h)}}[\mathcal{B}_{\hat{\mathbf{a}},h}^c \mid \mathcal{F}_{h-1}] \\
2106 \quad &\leq \hat{\gamma} \circ \mathbf{d}_{\text{TVC}}(\mathbf{s}_h^{\circ}, \hat{\mathbf{s}}_h) + (\hat{\gamma} + \gamma_{\sigma}) \circ \mathbf{d}_{\text{TVC}}(\mathbf{s}_h^{\circ}, \mathbf{s}_h^{\text{tel}}) + \vec{\mathbf{d}}_{\text{os},\varepsilon}(\hat{\pi}_{\sigma,h}(\mathbf{s}_h^{\text{tel}}) \parallel \pi_{\circ\sigma,h}^*(\mathbf{s}_h^{\text{tel}})) \\
2107 \quad &= \hat{\gamma} \circ \mathbf{d}_{\text{TVC}}(\mathbf{s}_h^{\circ}, \hat{\mathbf{s}}_h) + (\hat{\gamma} + \gamma_{\sigma}) \circ \mathbf{d}_{\text{TVC}}(\mathbf{s}_h^{\circ}, \mathbf{s}_h^{\text{tel}}) + \vec{\mathbf{d}}_{\text{os},\varepsilon}(\hat{\pi}_{\sigma,h}(\mathbf{s}_h^{\text{tel}}) \parallel \pi_{\circ\sigma,h}^*(\mathbf{s}_h^{\text{tel}}))
\end{aligned}$$

2111 This concludes the inductive construction.

2113 For the second statement, notice that the events $\bar{\mathcal{C}}_{\text{all},h} \cap \bar{\mathcal{B}}_{\text{all},h-1}$ are \mathcal{F}_h measurable (thus determined by $\mu^{(h-1)}$) and, when they hold, $\vec{\mathbf{d}}_{\mathcal{S}}(\mathbf{s}_h^{\circ}, \mathbf{s}_h^{\text{tel}}) \preceq \vec{\gamma}_{\text{IPS}}(2r)$ and $\mathbf{d}_{\mathcal{S}}(\mathbf{s}_h^{\circ}, \hat{\mathbf{s}}_h) \preceq \vec{\varepsilon}$. For our purposes, we use $\mathbf{d}_{\text{TVC}} = \mathbf{d}_{\mathcal{S},1}(\mathbf{s}_h^{\circ}, \mathbf{s}_h^{\text{tel}}) \preceq \gamma_{\text{IPS},1}(2r)$ and $\mathbf{d}_{\mathcal{S}}(\mathbf{s}_h^{\circ}, \hat{\mathbf{s}}_h) \preceq \vec{\varepsilon}_1$. Hence,

$$\begin{aligned}
2114 \quad \max_{h \in [H]} \mathbb{P}_{\mu}[\mathcal{B}_{\text{all},h}^c \cap \bar{\mathcal{C}}_{\text{all},h} \cap \bar{\mathcal{B}}_{\text{all},h-1}] &\leq \hat{\gamma}(\vec{\varepsilon}_1) + (\hat{\gamma} + \gamma_{\sigma}) \circ \gamma_{\text{IPS},1}(2r) \\
2115 \quad &+ \vec{\mathbf{d}}_{\text{os},\varepsilon}(\hat{\pi}_{\sigma,h}(\mathbf{s}_h^{\text{tel}}) \parallel \pi_{\circ\sigma,h}^*(\mathbf{s}_h^{\text{tel}})).
\end{aligned}$$

2120 The result follows.

2122 E.3. Proof of Theorem 2, and generalization to direct decompositions

2124 In this subsection, we consider the special case dealt with in Theorem 2. Note that there always exists a trivial direct decomposition that is compatible with all policies and dynamics simply by letting $\mathcal{V} = \emptyset$ and $\mathcal{S} = \mathcal{Z}$. We prove here the version of the result that involves a possibly nontrivial direct decomposition, as we will instantiate this in our control setting by letting $\mathcal{Z} = \{\boldsymbol{\rho}_{\text{m},h}\}$ and $\mathcal{S} = \{\boldsymbol{\rho}_{\text{c},h}\}$, i.e., projecting $\boldsymbol{\rho}_{\text{c},h}$ onto the last τ_{m} coordinates gives \mathbf{z}_h . We further consider a restriction of IPS to consider kernels absolutely continuous with respect to \mathbb{P}_h^* in their \mathcal{Z} component.

2129 **Definition E.9** (Restricted IPS). For a non-decreasing maps $\gamma_{\text{IPS},1}, \gamma_{\text{IPS},2} : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}_{\geq 0}$ a pseudometric $\mathbf{d}_{\text{IPS}} : \mathcal{S} \times \mathcal{S} \rightarrow \mathbb{R}$ (possibly other than $\mathbf{d}_{\mathcal{S}}$ or \mathbf{d}_{TVC}), and $r_{\text{IPS}} > 0$, we say a policy π is $(\gamma_{\text{IPS},1}, \gamma_{\text{IPS},2}, \mathbf{d}_{\text{IPS}}, r_{\text{IPS}})$ -restricted IPS if the following holds for any $r \in [0, r_{\text{IPS}}]$. Consider any sequence of kernels $\mathbf{W}_1, \dots, \mathbf{W}_H : \mathcal{S} \rightarrow \Delta(\mathcal{S})$ satisfying

$$2130 \quad \max_{h, \mathbf{s} \in \mathcal{S}} \mathbb{P}_{\tilde{\mathbf{s}} \sim \mathbf{W}_h(\mathbf{s})}[\mathbf{d}_{\text{IPS}}(\tilde{\mathbf{s}}, \mathbf{s}) \leq r] = 1, \quad \forall \mathbf{s}, \quad \phi_{\mathcal{Z}} \circ \mathbf{W}_h(\mathbf{s}_h) \ll \phi_{\mathcal{Z}} \circ \mathbb{P}_h^*.$$

2135 and define a process $\mathbf{s}_1 \sim \mathbb{P}_{\text{init}}$, $\tilde{\mathbf{s}}_h \sim \mathbf{W}_h(\mathbf{s}_h)$, $\mathbf{a}_h \sim \pi_h(\tilde{\mathbf{s}}_h)$, and $\mathbf{s}_{h+1} := F_h(\mathbf{s}_h, \mathbf{a}_h)$. Then, almost surely, (a) the sequence $(\mathbf{s}_{1:H+1}, \mathbf{a}_{1:H})$ is input-stable w.r.t $(\mathbf{d}_{\mathcal{S}}, \mathbf{d}_{\mathcal{A}})$ (b) $\max_{h \in [H]} \mathbf{d}_{\text{TVC}}(F_h(\tilde{\mathbf{s}}_h, \mathbf{a}_h), \mathbf{s}_{h+1}) \leq \gamma_{\text{IPS},1}(r)$ and (c) $\max_{h \in [H]} \mathbf{d}_{\mathcal{S}}(F_h(\tilde{\mathbf{s}}_h, \mathbf{a}_h), \mathbf{s}_{h+1}) \leq \gamma_{\text{IPS},2}(r)$.

2139 Note that the above is a slightly weaker condition than the one in Definition 4.5 in the main text and consequently, the following theorem which uses it as an assumption implies Theorem 2 in the body.

2141 **Theorem 5.** Suppose $\mathcal{S} = \mathcal{Z} \oplus \mathcal{V}$ as in Definition E.1 and projections $\phi_{\mathcal{Z}}, \phi_{\mathcal{V}}$, which is compatible with the dynamics and with given policies $\hat{\pi}, \pi^*$, smoothing kernel \mathbf{W}_{σ} , and pseudometric \mathbf{d}_{IPS} . Suppose π^* satisfies $(\gamma_{\text{IPS},1}, \gamma_{\text{IPS},2}, \mathbf{d}_{\text{IPS}}, r_{\text{IPS}})$ -restricted IPS (Definition E.9) and $\phi_{\mathcal{Z}} \circ \mathbf{W}_{\sigma}$ is γ_{σ} -TVC. Let $\varepsilon > 0$, $r \in (0, \frac{1}{2}r_{\text{IPS}}]$; define $p_r := \sup_{\mathbf{s}} \mathbb{P}_{\mathbf{s}' \sim \mathbf{W}_{\sigma}(\mathbf{s})}[\mathbf{d}_{\text{IPS}}(\mathbf{s}', \mathbf{s}) > r]$

2145 and $\varepsilon' := \varepsilon + \gamma_{\text{IPS},2}(2r)$. Then, for any policy $\hat{\pi}$, both $\Gamma_{\text{joint},\varepsilon}(\hat{\pi} \circ W_\sigma \parallel \pi_\odot^*)$ and $\Gamma_{\text{marg},\varepsilon'}(\hat{\pi} \circ W_\sigma \parallel \pi^*)$ are upper bounded
 2146 by

$$2147 \quad H(2p_r + 3\gamma_\sigma(\max\{\varepsilon, \gamma_{\text{IPS},1}(2r)\})) + \sum_{h=1}^H \mathbb{E}_{s_h^* \sim P_h^*} \mathbb{E}_{\tilde{s}_h^* \sim W_\sigma(s_h^*)} d_{\text{os},\varepsilon}(\hat{\pi}_h(\tilde{s}_h^*) \parallel \pi_{\text{dec}}^*(\tilde{s}_h^*)).$$

2148 Consider the special case $K = 2$ with $d_{S,1} = d_{\text{TVc}}$, $d_{S,2} = d_S$, $d_{\mathcal{A},1} = d_{\mathcal{A},2} = d_{\mathcal{A}}$ and $\vec{\varepsilon} = (\varepsilon, \varepsilon)$. In this case, applying
 2150 (E.3), we see that

$$2152 \quad \vec{\Gamma}_{\text{joint},\varepsilon}(\hat{\pi}_\sigma \parallel \pi_\odot^*) \vee \vec{\Gamma}_{\text{marg},\vec{\varepsilon}_{\text{marg}}}(\hat{\pi}_\sigma \parallel \pi_\odot^*) \\
 2153 \quad \leq p_{\text{IPS}} + H(2p_r + 3\gamma_\sigma(\max\{\varepsilon, \gamma_{\text{IPS},1}(2r)\})) + \sum_{h=1}^H \mathbb{E}_{s_h^* \sim P_h^*} \mathbb{E}_{\tilde{s}_h^* \sim W_\sigma(s_h^*)} \vec{d}_{\text{os},\varepsilon}(\hat{\pi}_h(\tilde{s}_h^*) \parallel \pi_{\text{dec}}^*(\tilde{s}_h^*))$$

2155 We now observe that under this convention,

$$2157 \quad \Gamma_{\text{joint},\varepsilon}(\hat{\pi}_\sigma \parallel \pi_\odot^*) = \inf_{\mu_1} \mathbb{P}_{\mu_1} \left[\max_{h \in [H]} d_S(\hat{s}_{h+1}, s_{h+1}^*) \vee d_{\mathcal{A}}(\hat{a}_h, a_h^*) > \varepsilon \right] \\
 2158 \quad \leq \inf_{\mu_1} \mathbb{P}_{\mu_1} \left[\max_{h \in [H]} (d_{\text{TVc}}(\hat{s}_{h+1}, s_{h+1}^*), d_S(\hat{s}_{h+1}, s_{h+1}^*)) \vee (d_{\mathcal{A}}(\hat{a}_h, a_h^*), d_{\mathcal{A}}(\hat{a}_h, a_h^*)) \not\leq \vec{\varepsilon} \right] \\
 2159 \quad = \vec{\Gamma}_{\text{joint},\varepsilon}(\hat{\pi}_\sigma \parallel \pi_\odot^*)$$

2163 and similarly $\Gamma_{\text{marg},\varepsilon'}(\hat{\pi}_\sigma \parallel \pi^*) \leq \vec{\Gamma}_{\text{marg},\vec{\varepsilon} + \gamma_{\text{IPS}}(2r)}(\hat{\pi}_\sigma \parallel \pi^*)$. From the construction of $\vec{d}_{\mathcal{A}}$, however, we see that
 2164 $\{\vec{d}_{\mathcal{A}}(a, a') \not\leq \vec{\varepsilon}\} = \{d_{\mathcal{A}}(a, a') > \varepsilon\}$ for all a, a' and thus for all $h \in [H]$,

$$2166 \quad \vec{d}_{\text{os},\varepsilon}(\hat{\pi}_h(\tilde{s}_h^*) \parallel \pi_h^*(\tilde{s}_h^*)) = \inf_{\mu_2} \mathbb{P}_{\mu_2} \left[\vec{d}_{\mathcal{A}}(\hat{a}_h, a_h^*) \not\leq \vec{\varepsilon} \right] \\
 2167 \quad = \inf_{\mu_2} \mathbb{P}_{\mu_2} [d_{\mathcal{A}}(\hat{a}_h, a_h^*) \geq \varepsilon] \\
 2168 \quad = d_{\text{os},\varepsilon}(\hat{\pi}_h(\tilde{s}_h^*) \parallel \pi_h^*(\tilde{s}_h^*)).$$

2171 Plugging in to (E.3) concludes the proof.

2174 F. Lower Bounds

2175 In this section, we establish lower bounds against the imitation results in the composite MDP. Specifically, we show that

- 2177 • In [Appendix F.1](#) we show that [Theorem 2](#) and [Proposition D.1](#) are sharp in the regime where $\gamma_{\text{IPS},1} = \gamma_{\text{IPS},2} = 0$.
- 2179 • In [Appendix F.2](#), we show that the marginals of an expert policy π^* and replica policy $\pi_{\odot\sigma}^*$ can coincide, but their joint
 2180 distributions can be different. By considering $\hat{\pi} = \pi_{\text{dec}}^*$ in [Theorem 2](#), this establishes the necessity of considering the
 2181 marginal imitation gap with respect to π^* .
- 2182 • In [Appendix F.3](#), we lower bound the distance between *marginal distributions* over states under π^* and $\pi_{\odot\sigma}^*$ in the
 2183 regime where $\gamma_{\text{IPS},2} \neq 0$. This example demonstrates that the dependence of $\gamma_{\text{IPS},2}$ in [Theorem 2](#) is essentially sharp.
- 2184 • In [Appendix F.4](#), we show that for an expert policy π^* and smoothing kernel W_σ , the state distributions under $\pi_{\odot\sigma}^*$ and
 2185 π_{dec}^* can have different marginals (and thus different joint distributions). By considering $\hat{\pi} = \pi_{\text{dec}}^*$ in [Theorem 2](#), this
 2186 explains why it is necessary to smooth $\hat{\pi}$ to $\hat{\pi} \circ W_\sigma$.

2189 Taken together, the above counterexamples show that our distinctions between joint and marginal distributions, decision to
 2190 add noise at inference time, and dependence on almost all problem quantities in [Section 4](#) are sharp. We do not, however,
 2191 establish necessity of $\gamma_{\text{IPS},1}$ in the interest of brevity; we believe this quantity is necessary. Still, the $\gamma_{\text{IPS},1}$ term contributes a
 2192 factor exponentially small in τ_c in [Theorem 1](#), so we deem lower bounds establishing its necessity of lesser importance.

2194 **Commonalities of construction.** In all but [Appendix F.3](#), we take the action and state spaces to be

$$2195 \quad \mathcal{S} = \mathcal{A} = \mathbb{R},$$

2197 which is the archetypal Polish space ([Durrett, 2019](#)). Throughout, we use δ_x to denote the dirac-delta distribution on $x \in \mathbb{R}$.
 2198 We let $d_S(s', s) = d_{\text{TVc}}(s', s) = |s' - s|$ and $d_{\mathcal{A}}(a', a) = |a' - a|$ all be the Euclidean distance.

2199

F.1. Sharpness of Proposition D.1 and Theorem 2

Here, we demonstrate that Proposition D.1 is tight up to constant factors, and that Theorem 2 is tight up to the terms $\gamma_{\text{IPS},1}$, $\gamma_{\text{IPS},2}$ and concentration probability p_r . Consider the simple dynamics

$$F_h(\mathbf{s}, \mathbf{a}) = \mathbf{a}.$$

Note that, as the dynamics are state-independent, we have $\gamma_{\text{IPS},1}(\cdot) = \gamma_{\text{IPS},2}(\cdot) \equiv 0$. Furthermore, let us assume policies do not depend on time index h . Let $\pi^* : \mathbf{s} \rightarrow \delta_0$ be deterministic, and let $\mathbf{P}_{\text{init}} = \delta_0$ be an initial state distribution concentrated on 0. Then, \mathbf{D}_{π^*} is the dirac distribution on the all-zero trajectory.

Fix parameters $0 < \varepsilon < \sigma$, and $p \in (0, 1)$. We consider the following smoothing-kernel

$$W_{\varepsilon, \sigma} = \begin{cases} \delta_0 & \mathbf{s} \leq 0 \\ (1 - \frac{\mathbf{s}}{\sigma})\delta_0 + \frac{\mathbf{s}}{\sigma}\delta_\sigma & \mathbf{s} \in [0, \sigma] \\ \delta_\sigma & \mathbf{s} > \sigma, \end{cases}$$

Define the candidate policy

$$\hat{\pi}_{\varepsilon, p, \sigma}(\mathbf{s}) := \begin{cases} (1-p)\delta_\varepsilon + p\delta_\sigma & \mathbf{s} \leq \frac{\varepsilon}{2} \\ \delta_\sigma & \mathbf{s} > \frac{\varepsilon}{2} \end{cases}$$

Proposition F.1. *For any $p \in (0, 1)$, $0 < \varepsilon < \sigma$, set $\bar{\pi} = \hat{\pi}_{\varepsilon, p, \sigma} \circ W_{\sigma, \varepsilon}$. Then,*

(a) π^* , π_{\circlearrowleft}^* and π_{dec}^* all map $\mathbf{s} \rightarrow \delta_0$, $\mathbf{P}_h^* = \delta_0$, and thus for any $\tilde{\pi} \in \{\pi^*, \pi_{\circlearrowleft}^*, \pi_{\text{dec}}^*\}$,

$$\mathbb{E}_{\mathbf{s}_h^* \sim \mathbf{P}_h^*} \mathbb{E}_{\mathbf{s}'_h \sim W_\sigma(\mathbf{s}_h^*)} [\mathbf{d}_{\text{os}, \varepsilon}(\hat{\pi}_{\varepsilon, p, \sigma}(\mathbf{s}'_h) \parallel \tilde{\pi}(\mathbf{s}'_h))] = \mathbb{E}_{\mathbf{s}_h^* \sim \mathbf{P}_h^*} [\mathbf{d}_{\text{os}, \varepsilon}(\bar{\pi}(\mathbf{s}_h^*) \parallel \tilde{\pi}(\mathbf{s}_h^*))] = p.$$

(b) The kernel $W_{\sigma, \varepsilon}$ is γ_σ -TVC, where $\gamma_\sigma(u) = u/\sigma$.

(c) For a universal constant $c > 0$,

$$\Gamma_{\text{joint}, \varepsilon}(\bar{\pi} \parallel \pi^*) = \Gamma_{\text{marg}, \varepsilon}(\bar{\pi} \parallel \pi^*) \geq c \min\{1, H(p + \varepsilon/\sigma)\},$$

and the same holds with π^* replaced by π_{\circlearrowleft}^* or π_{dec}^* .

In particular, the above proposition shows that

$$\Gamma_{\text{joint}, \varepsilon}(\bar{\pi} \parallel \pi^*) = \Gamma_{\text{marg}, \varepsilon}(\bar{\pi} \parallel \pi^*) \gtrsim H\gamma_\sigma(\varepsilon) + \sum_{h=1}^H \mathbb{E}_{\mathbf{s}_h^* \sim \mathbf{P}_h^*} [\mathbf{d}_{\text{os}, \varepsilon}(\bar{\pi}(\mathbf{s}_h^*) \parallel \pi^*(\mathbf{s}_h^*))],$$

verifying the sharpness of Proposition D.1 (note that $\bar{\pi} = \hat{\pi}_{\varepsilon, p, \sigma} \circ W_\sigma$ is γ_σ TVC). Similarly, our above proposition shows that,

$$\Gamma_{\text{joint}, \varepsilon}(\bar{\pi} \parallel \pi_{\circlearrowleft}^*) = \Gamma_{\text{marg}, \varepsilon}(\bar{\pi} \parallel \pi^*) \gtrsim H\gamma_\sigma(\varepsilon) + \sum_{h=1}^H \mathbb{E}_{\mathbf{s}_h^* \sim \mathbf{P}_h^*} [\mathbf{d}_{\text{os}, \varepsilon}(\hat{\pi}_{\varepsilon, p, \sigma}(\mathbf{s}_h^*) \parallel \pi_{\text{dec}, h}^*(\mathbf{s}_h^*))],$$

verifying that Theorem 2 is sharp up to the additional stability terms $\gamma_{\text{IPS},1}$, $\gamma_{\text{IPS},2}$.

Proof. We begin with a computation. Define

$$\eta(\mathbf{s}) = 1 - (1-p)(1 - \frac{\mathbf{s}}{\sigma}) = p + (1-p)\frac{\mathbf{s}}{\sigma}$$

2255 We compute

$$\begin{aligned}
2256 \bar{\pi} &= \hat{\pi}_{\varepsilon,p,\sigma} \circ W_{\sigma,\varepsilon} = \begin{cases} (1-p)\delta_\varepsilon + p\delta_\sigma & s \leq \frac{\varepsilon}{2} \\ \delta_\sigma & s > \frac{\varepsilon}{2} \end{cases} \circ \begin{cases} \delta_0 & s \leq 0 \\ (1 - \frac{s}{\sigma})\delta_0 + \frac{s}{\sigma}\delta_\sigma & s \in [0, \sigma] \\ \delta_\sigma & s \geq \sigma. \end{cases} \\
2257 & \\
2258 & \\
2259 & \\
2260 & \\
2261 & \\
2262 & \\
2263 & \\
2264 & \\
2265 & \\
2266 & \\
2267 & \\
2268 & \\
2269 & \\
2270 & \\
2271 & \\
2272 & \\
2273 & \\
2274 & \\
2275 & \\
2276 & \\
2277 & \\
2278 & \\
2279 & \\
2280 & \\
2281 & \\
2282 & \\
2283 & \\
2284 & \\
2285 & \\
2286 & \\
2287 & \\
2288 & \\
2289 & \\
2290 & \\
2291 & \\
2292 & \\
2293 & \\
2294 & \\
2295 & \\
2296 & \\
2297 & \\
2298 & \\
2299 & \\
2300 & \\
2301 & \\
2302 & \\
2303 & \\
2304 & \\
2305 & \\
2306 & \\
2307 & \\
2308 & \\
2309 &
\end{aligned}
\tag{F.1}$$

2265 In particular,

$$2266 \hat{\pi}(0) = \pi_{\varepsilon,p,\sigma}(0) = (1-p)\delta_\varepsilon + p\delta_\sigma$$

2269 **Part (a).** Notice that the support of the deconvolution and replica distributions are always in the support of P_h^* , which is always $s = 0$ under π^* . Thus, $\pi^* = \pi_{\circlearrowleft\sigma}^* = \pi_{\text{dec}}^*$. By the same token, for any policy π ,

$$2271 \mathbb{E}_{s_h^* \sim P_h^*} [d_{\text{os},\varepsilon}(\pi(s_h^*) \parallel \tilde{\pi}_*(s_h^*))] = \mathbb{P}[|\pi(0)| > \varepsilon].$$

2273 Hence, as $\bar{\pi}(0) = \hat{\pi}_{\varepsilon,p,\sigma}(0) = (1-p)\delta_\varepsilon + p\delta_\sigma$, and as $\sigma > \varepsilon$, part (a) follows.

2275 **Part (b).** Consider $s, s' \in \mathcal{S}$. We can assume, from the functional form of $W_{\varepsilon,\sigma}(\cdot)$, that $0 \leq s \leq s' \leq \sigma$. Then,

$$2277 \text{TV}(W_{\varepsilon,\sigma}(s), W_{\varepsilon,\sigma}(s')) = \text{TV}(\delta_0(1 - \frac{s}{\sigma}) + (\frac{s}{\sigma})\delta_\sigma, \delta_0(1 - \frac{s'}{\sigma}) + (\frac{s'}{\sigma})\delta_\sigma) = \frac{|s' - s|}{\sigma},$$

2279 establishing total variation continuity.

2281 **Part (c)** In view of part (a), it suffices to bound gaps relative to π^* . Let \mathbb{P} denote probabilities over $s_{1:H+1}, a_h$ under $\bar{\pi}$. Let $\mathcal{A}_{1,h}$ denote the event that at step h , $a_h = \varepsilon$, and let $\mathcal{A}_{2,h}$ denote the event that $a_h = \sigma$. As the state s_0 is absorbing and as $F_h(s, a) = a_h$, the following events are equal

$$2285 \{\exists h : |a_h| \vee |s_{h+1}| > \varepsilon\} = \mathcal{A}_{2,H}.$$

2287 Hence,

$$2288 \Gamma_{\text{joint},\varepsilon}(\bar{\pi} \parallel \pi^*) = \mathbb{P}[\mathcal{A}_{2,H}].$$

2290 Moreover, as $\mathcal{A}_{2,H}$ is measurable with respect to the marginal of a_H , we also have that

$$2292 \Gamma_{\text{marg},\varepsilon}(\bar{\pi} \parallel \pi^*) = \mathbb{P}[\mathcal{A}_{2,H}].$$

2294 It thus suffices to lower bound $\mathbb{P}[\mathcal{A}_{2,H}]$. By definition of $\bar{\pi}$, the events $\mathcal{A}_{1,h}, \mathcal{A}_{2,h}$ are exhaustive: $\mathcal{A}_{1,h}^c = \mathcal{A}_{2,h}$. Moreover, from (F.1),

$$2296 \mathbb{P}[\mathcal{A}_{2,h+1} \mid \mathcal{A}_{2,h}] = 1, \quad \mathbb{P}[\mathcal{A}_{2,h+1} \mid \mathcal{A}_{1,h}] = \eta(\varepsilon), \quad \mathbb{P}[\mathcal{A}_{1,1}] = 1 - \eta(0) \geq 1 - \eta(\varepsilon).$$

2298 Thus,

$$\begin{aligned}
2299 \mathbb{P}[\mathcal{A}_{2,H}] &= \mathbb{P}[\mathcal{A}_{2,H} \mid \mathcal{A}_{2,H-1}] \mathbb{P}[\mathcal{A}_{2,H-1}] + \mathbb{P}[\mathcal{A}_{2,H} \mid \mathcal{A}_{1,H-1}] \mathbb{P}[\mathcal{A}_{1,H-1}] \\
2300 &= \mathbb{P}[\mathcal{A}_{2,H-1}] + \eta(\varepsilon) \mathbb{P}[\mathcal{A}_{1,H-1}] \\
2301 &= \mathbb{P}[\mathcal{A}_{2,H-2}] + \eta(\varepsilon) (\mathbb{P}[\mathcal{A}_{1,H-1}] + \mathbb{P}[\mathcal{A}_{1,H-2}]) \\
2302 &= \eta(\varepsilon) \left(\sum_{h=1}^{H-1} \mathbb{P}[\mathcal{A}_{1,h}] \right) + \mathbb{P}[\mathcal{A}_{2,1}] \\
2303 &\geq \eta(\varepsilon) \left(\sum_{h=1}^{H-1} \mathbb{P}[\mathcal{A}_{1,h}] \right)
\end{aligned}$$

2310 Moreover, as s_0 is absorbing,

$$2311 \quad \mathbb{P}[\mathcal{A}_{1,h}] = \mathbb{P}[\mathcal{A}_{1,h} \mid \mathcal{A}_{1,h-1}] \mathbb{P}[\mathcal{A}_{1,h-1}] = (1 - \eta(\varepsilon)) \mathbb{P}[\mathcal{A}_{1,h-1}].$$

2312
2313
2314 Combining with $\mathbb{P}[\mathcal{A}_{1,1}] = (1 - p) \geq (1 - \eta(0)) \geq 1 - \eta(\varepsilon)$, we have $\mathbb{P}[\mathcal{A}_{1,h}] \geq (1 - \eta(\varepsilon))^h$. Hence,

$$2315 \quad \mathbb{P}[\mathcal{A}_{2,H+1}] \geq \eta(\varepsilon) \left(\sum_{h=1}^{H-1} (1 - \eta(\varepsilon))^h \right)$$

$$2316 \quad = \eta(\varepsilon) \frac{1 - \eta(\varepsilon) - (1 - \eta(\varepsilon))^H}{1 - (1 - \eta(\varepsilon))}$$

$$2317 \quad = 1 - \eta(\varepsilon) - (1 - \eta(\varepsilon))^H$$

$$2318 \quad = \Omega(\min\{1, H(\eta(\varepsilon))\})$$

2319
2320
2321
2322
2323
2324 as $\eta(\varepsilon) \downarrow 0$. Substituting in $\eta(\varepsilon) = p + (1 - p)\varepsilon/\sigma = \Omega(p + \varepsilon/\sigma)$ concludes. \square

2325 F.2. $\pi_{\circlearrowleft\sigma}^*$ and π^* induce the same marginals but different joint distributions, even with memoryless dynamics

2326 We give a simple example where $\pi_{\circlearrowleft\sigma}^*$ and π^* induce the same marginal distributions over trajectories, but different joints.
2327 As we show, this example demonstrates the necessity of measuring the marginal imitation error of a smoothed policy,
2328 $\Gamma_{\text{marg},\varepsilon}$, over the joint error, $\Gamma_{\text{joint},\varepsilon}$. A graphical (but nonrigorous) demonstration of this issue can be seen in [Figure 4](#) in
2329 [Appendix B](#).

2330 Again, let $\mathcal{S} = \mathcal{A} = \mathbb{R}$, and $F_h(\mathbf{s}, \mathbf{a}) = \mathbf{a}$. We let

$$2331 \quad W_\sigma(\cdot) = \mathcal{N}(\cdot, \sigma^2)$$

2332 denote Gaussian smoothing. Fix some $\varepsilon > 0$. Define

$$2333 \quad P_{\text{init}} = \frac{1}{2}(\delta_{-\varepsilon} + \delta_{+\varepsilon}), \quad \pi^*(\mathbf{s}) = \begin{cases} \delta_{-\varepsilon} & \mathbf{s} \leq 0 \\ \delta_{\varepsilon} & \mathbf{s} > 0 \end{cases}.$$

2334 Thus, D_{π^*} is supported on the trajectories with $(\mathbf{s}_{1:H+1}, \mathbf{a}_{1:H})$ being either all ε or all $-\varepsilon$, and

$$2335 \quad P_h^* = P_{\text{init}} = \frac{1}{2}(\delta_{-\varepsilon} + \delta_{+\varepsilon}).$$

2336 Hence, the replica and deconvolution map to distributions supported on $\{\varepsilon, -\varepsilon\}$. Let $\phi_\sigma(\cdot)$ denote the Gaussian PDF with
2337 variance σ . Then,

$$2338 \quad W_{\text{dec},h}^*(\mathbf{s}) = \frac{\delta_\varepsilon \phi_\sigma(\mathbf{s} - \varepsilon) + \delta_{-\varepsilon} \phi_\sigma(\mathbf{s} + \varepsilon)}{\phi_\sigma(\mathbf{s} - \varepsilon) + \phi_\sigma(\mathbf{s} + \varepsilon)}.$$

2339 Moreover,

$$2340 \quad W_{\circlearrowleft,h}^*(\mathbf{s}) = \mathbb{E}_{Z \sim \mathcal{N}(0, \sigma^2)} \left[\frac{\delta_\varepsilon \phi_\sigma(\mathbf{s} - \varepsilon + Z) + \delta_{-\varepsilon} \phi_\sigma(\mathbf{s} + \varepsilon + Z)}{\phi_\sigma(\mathbf{s} - \varepsilon + Z) + \phi_\sigma(\mathbf{s} + \varepsilon + Z)} \right]. \quad (\text{F.2})$$

2341 One can check that for $\varepsilon \leq \sigma$,

$$2342 \quad W_{\circlearrowleft,h}^*(u\varepsilon) = \Theta \left(\frac{(1 + \frac{c\varepsilon}{\sigma})\delta_{u\varepsilon} + (1 - \frac{c\varepsilon}{\sigma})\delta_{-u\varepsilon}}{2} \right), \quad u \in \{-1, 1\}$$

2343 for $\varepsilon \ll 1$. In particular, for $\mathbf{s} \in \{-\varepsilon, \varepsilon\}$

$$2344 \quad \mathbb{P}_{\mathbf{a} \sim \pi_{\circlearrowleft\sigma,h}^*}[\mathbf{a} = -\mathbf{s}] \geq \Omega(1). \quad (\text{F.3})$$

2365 In particular, if $(\mathbf{s}_{1:H+1}^\circ, \mathbf{a}_{1:H}^\circ) \sim \mathcal{D}_{\pi_{\circ\sigma}^*}$, then

$$2366 \mathbb{P}[\exists h : d(\mathbf{s}_h^\circ, \mathbf{s}_{h+1}^\circ) > \varepsilon] \leq \mathbb{P}[\exists h : \mathbf{s}_h^\circ = -\mathbf{s}_{h+1}^\circ]$$

$$2367 \leq \mathbb{P}[\exists h : \mathbf{s}_h^\circ = -\mathbf{a}_h^\circ] = 1 - \exp(-\Omega(H)),$$

2370 where in the last step we used (F.3) and the fact that the $\pi_{\circ\sigma}^*$ uses fresh randomness at each round. Moreover, as π^*
 2371 always commits to either an all- ε or all- $(-\varepsilon)$ -trajectory, we see that for any $\mu \in \mathcal{C}(\mathcal{D}_{\pi^*}, \mathcal{D}_{\pi_{\circ\sigma}^*})$ over $(\mathbf{s}_{1:H+1}^*, \mathbf{a}_{1:H}^*) \sim \mathcal{D}_{\pi^*}$
 2372 and $(\mathbf{s}_{1:H+1}^\circ, \mathbf{a}_{1:H}^\circ) \sim \mathcal{D}_{\pi_{\circ\sigma}^*}$,

$$2373 \Gamma_{\text{joint}, \varepsilon}(\pi_{\circ\sigma}^*, \pi^*) \geq \mathbb{P}_\mu[\exists 1 \leq h \leq H : d(\mathbf{s}_{h+1}^*, \mathbf{s}_{h+1}^\circ) > \varepsilon] \geq 1 - \exp(-\Omega(H)),$$

2374 That is, the replica and expert policies have different joint state distribution.

2375 **Remark F.1.** The above result demonstrates the necessity of measuring the marginal error between $\hat{\pi} \circ W_\sigma$ and π^* in
 2376 **Theorem 2:** if we apply that proposition with $\hat{\pi} = \pi_{\text{dec}}^*$, then for all ε , $\mathbb{E}_{\tilde{\mathbf{s}}_h^* \sim W_\sigma(\mathbf{s}_h^*)} d_{\text{os}, \varepsilon}(\hat{\pi}_h(\tilde{\mathbf{s}}_h^*) \parallel \pi_{\text{dec}}^*(\tilde{\mathbf{s}}_h^*)) = 0$. But
 2377 then $\hat{\pi} \circ W_\sigma = \pi_{\circ\sigma}^*$, and we know that $\Gamma_{\text{joint}, \varepsilon}(\pi_{\circ\sigma}^*, \pi^*) \geq \mathbb{P}_\mu[\exists 1 \leq h \leq H : d(\mathbf{s}_{h+1}^*, \mathbf{s}_{h+1}^\circ) > \varepsilon] \geq 1 - \exp(-\Omega(H))$.
 2378 Thus, we cannot hope for smoothed policies to imitate expert demonstrations in joint state distributions without additional
 2379 assumptions.

2380 **Remark F.2 (Importance of chunking).** Above we have shown that $\pi_{\circ\sigma}^*$ oscillates between ε and $-\varepsilon$ (for actions and
 2381 subsequent states). We remark that these oscillations can have very deleterious effects on performance on real control
 2382 systems. This is why it is beneficial to predict entire sequences of trajectories. Indeed, consider a modified construction
 2383 such that $\mathcal{S} = \mathcal{A} = \mathbb{R}^K$, and $F_h(\mathbf{s}, \mathbf{a}) = \mathbf{a}$. Here, we interpret \mathcal{S} as a sequence of K -control states in \mathbb{R} , and \mathbf{a} as sequence
 2384 of K -actions, denoting the i -th coordinate of \mathbf{s} via $\mathbf{s}[i]$,

$$2385 \pi^*(\mathbf{s}) = \begin{cases} \delta_{-\varepsilon \mathbf{1}} & \mathbf{s}[1] \leq 0 \\ \delta_{\varepsilon \mathbf{1}} & \mathbf{s}[1] > 0, \end{cases}$$

2386 Then, we can view the oscillations in $\pi_{\circ\sigma}^*$ as oscillations between length K trajectories, which is essentially what happens
 2387 in our analysis for $K = \tau_c$.

2388 **F.3. $\pi_{\circ\sigma}^*$ and π^* can have different marginals, implying necessity of $\gamma_{\text{IPS}, 2}$**

2389 Our construction lifts the construction in Appendix F.2 to a two-dimensional state space $\mathcal{S} = \mathbb{R}^2$, keeping one dimensional
 2390 actions $\mathcal{A} = \mathbb{R}$. Let $\mathbf{s} = (\mathbf{s}[1], \mathbf{s}[2])$ denote coordinate of $\mathbf{s} \in \mathcal{S}$. For some parameter ν , the dynamics are

$$2391 \mathbf{s}_{h+1} = F_h(\mathbf{s}_h, \mathbf{a}_h) = (\mathbf{a}_h, \nu \cdot (\mathbf{s}_h[1] - \mathbf{a}_h))$$

2392 We let $d_{\mathcal{S}} = d_{\text{TVc}} = d_{\text{IPS}}$ denote the ℓ_1 norm on $\mathcal{S} = \mathcal{R}^2$. Our initial state distribution is

$$2393 P_{\text{init}} = \frac{1}{2} (\delta_{(\varepsilon, 0)} + \delta_{(-\varepsilon, 0)})$$

2394 We let

$$2395 \pi^*(\mathbf{s}) = \begin{cases} \delta_{(-\varepsilon, 0)} & \mathbf{s} \leq 0 \\ \delta_{(\varepsilon, 0)} & \mathbf{s} > 0 \end{cases}.$$

2396 Thus, π^* induces trajectories which either stay on $\delta_{(\varepsilon, 0)}$ or $\delta_{(-\varepsilon, 0)}$.

$$2397 P_h^* = \frac{1}{2} (\delta_{(\varepsilon, 0)} + \delta_{(-\varepsilon, 0)}), \quad \forall h \geq 1.$$

2398 Let

$$2399 W_\sigma(\mathbf{s}) = \mathcal{N}(\mathbf{s}', \sigma^2)$$

2400

2420 **Proposition F.2.** *In the above construction, we can take $\gamma_{\text{IPS},2}(u) \leq \nu \cdot u$ in [Definition 4.5](#), and p_r satisfies the conditions in*
 2421 *[Theorem 2](#) for $r = 2\sigma\sqrt{\log(1/p_r)}$. Moreover, for any $\varepsilon \leq \sigma$,*

$$2422 \Gamma_{\text{marg},\varepsilon'}(\pi_{\circlearrowleft\sigma}^* \parallel \pi^*) \geq \Omega(1), \quad \varepsilon' = \nu\varepsilon$$

2423
 2424
 2425 **Remark F.3** (Sharpness of $\gamma_{\text{IPS},2}$). Before proving this proposition, we note that if we take $\varepsilon = \sigma$ and $r = 2\sigma\sqrt{\log(1/p_r)}$,
 2426 then $\nu\varepsilon = \tilde{\Omega}(\gamma_{\text{IPS},2}(2r))$, showing that our dependence on $\gamma_{\text{IPS},2}$ is sharp up to logarithmic factors. Moreover, the looseness up
 2427 to logarithmic factors in the above point is an artifact of using the Gaussian smoothing W_σ , and can be removed by replacing
 2428 W_σ with a truncated-Gaussian kernel.

2429
 2430 *Proof of [Proposition F.2](#).* To see $\gamma_{\text{IPS},2}(u) \leq \nu \cdot u$, we have $\|F_h(\mathbf{s}, \mathbf{a}) - F_h(\mathbf{s}', \mathbf{a})\| = \|(\mathbf{a}, \nu \cdot (\mathbf{s}[1] - \mathbf{a})) - (\mathbf{a}, \nu \cdot (\mathbf{s}'[1] - \mathbf{a}))\| =$
 2431 $\nu|\mathbf{s}[1] - \mathbf{s}'[1]| \leq \nu d_{\text{TVC}}(\mathbf{s}, \mathbf{s}')$. That we can take $r = 2\sigma\sqrt{\log(1/p_r)}$ follows from Gaussian concentration.

2432 To prove the final claim, one can directly generalize [\(F.2\)](#) to find that, for any $b \in \mathbb{R}$,

$$2433 W_{\circlearrowleft,h}^*(\mathbf{s}) = \mathbb{E}_{Z \sim \mathcal{N}(0, \sigma^2)} \left[\frac{\delta_{(\varepsilon,0)} \phi_\sigma(\mathbf{s}[1] - \varepsilon + Z) + \delta_{(-\varepsilon,0)} \phi_\sigma(\mathbf{s}[1] + \varepsilon + Z)}{\phi_\sigma(\mathbf{s}[1] - \varepsilon + Z) + \phi_\sigma(\mathbf{s}[1] + \varepsilon + Z)} \right].$$

2434 This follows from the observation that $W_{\circlearrowleft,h}^*$ and P_h^* have the same support, and as P_h^* always is support on vectors with
 2435 second coordinate zero, that the second coordinate of \mathbf{s} in $W_{\circlearrowleft,h}^*(\mathbf{s})$ is uninformative. For $\varepsilon \leq \sigma$, we find that

$$2436 W_{\circlearrowleft,h}^*((\varepsilon, b)) = c\delta_{(\varepsilon,0)} + (1-c)\delta_{(-\varepsilon,0)}, \quad c = \Omega(1), \quad b \in \mathbb{R}.$$

2437 and $W_{\circlearrowleft,h}^*((-\varepsilon, b))$ is defined symmetrically. Hence, under $(\mathbf{s}_{1:H+1}^\circlearrowleft, \mathbf{a}_{1:H}^\circlearrowleft) \sim \pi_{\circlearrowleft\sigma}^*$,

$$2438 \mathbb{P}[\mathbf{s}_1^\circlearrowleft \neq \mathbf{a}_1^\circlearrowleft] \geq \Omega(1)$$

2439 Moreover, when $\mathbf{s}_2^\circlearrowleft \neq \mathbf{a}_h^\circlearrowleft$, we have that $|\mathbf{s}_2^\circlearrowleft[2]| = \nu|\mathbf{s}_1^\circlearrowleft - \mathbf{a}_1^\circlearrowleft|$, which as π^* is supported on $\{\delta_{(\varepsilon,0)}, \delta_{(-\varepsilon,0)}\}$, means,
 2440 $|\mathbf{s}_2^\circlearrowleft[2]| \geq 2\nu\varepsilon$. Thus,

$$2441 \mathbb{P}[|\mathbf{s}_2^\circlearrowleft[2]| \geq 2\nu\varepsilon] \geq \Omega(1)$$

2442 On the other hand, $\mathbf{s}_2^* \sim P_h^*$ has $\mathbf{s}_2^*[2] = 0$ with probability one. Thus, for any coupling μ between $D_{\pi^*}, D_{\pi_{\circlearrowleft\sigma}^*}$,

$$2443 \mathbb{P}_\mu[d_S(\mathbf{s}_2^\circlearrowleft, \mathbf{s}_2^*) \geq 2\nu\varepsilon] \geq \Omega(1)$$

2444 Thus,

$$2445 \Gamma_{\text{marg},\nu\varepsilon}(\pi_{\circlearrowleft\sigma}^* \parallel \pi^*) \geq \Omega(1).$$

2446 □

2447 **F.4. $\pi_{\circlearrowleft\sigma}^*$ and π_{dec}^* have different marginals, even with memoryless dynamics**

2448 Here, we show how $\pi_{\circlearrowleft\sigma}^*$ and π_{dec}^* have different marginals even if the dynamics are memoryless. By considering $\hat{\pi} = \pi_{\text{dec}}^*$
 2449 in [Theorem 2](#), the discussion below demonstrates why one needs to consider $\hat{\pi} \circ W_\sigma$ in order to obtain small imitation gap.

2450 For simplicity, we use a discrete smoothing kernel W_σ , though the example extends to the Gaussian smoothing kernel in the
 2451 previous counter example. Again, let $\mathcal{S} = \mathcal{A} = \mathbb{R}$, and $F_h(\mathbf{s}, \mathbf{a}) = \mathbf{a}$. Take

$$2452 \pi^*(\mathbf{s}) = \begin{cases} \delta_{-\sigma} & \mathbf{s} \leq 0 \\ \delta_\sigma & \mathbf{s} > 0 \end{cases}$$

2453 Let us consider an asymmetric initial state distribution

$$2454 P_{\text{init}} = \frac{1}{4}\delta_{-\sigma} + \frac{3}{4}\delta_{+\sigma}.$$

2455

2475 Note then that

$$2476 \quad \forall h, \quad P_h^* = P_{\text{init}} = \frac{1}{4}\delta_{-\sigma} + \frac{3}{4}\delta_{\sigma}, \quad (F.4)$$

2477 We consider a smoothing kernel,

$$2480 \quad W_{\sigma}(s) = \begin{cases} (\frac{1}{2} + \frac{s}{4\sigma})\delta_{\sigma} + (\frac{1}{2} - \frac{s}{4\sigma})\delta_{-\sigma} & -2\sigma \leq s \leq 2\sigma \\ \delta_{\sigma} & s \geq 2\sigma \\ \delta_{-\sigma} & s \leq -2\sigma \end{cases}$$

2481 The salient part of our construction of W_{σ} is that

$$2482 \quad W_{\sigma}(\sigma) = \frac{1}{4}\delta_{-\sigma} + \frac{3}{4}\delta_{\sigma}, \quad W_{\sigma}(-\sigma) = \frac{1}{4}\delta_{\sigma} + \frac{3}{4}\delta_{-\sigma}.$$

2483 Denote the marginals of $\pi_{\circlearrowleft\sigma}^*$ and π_{dec}^* with $P_{\circlearrowleft,h}^*$ and $P_{\text{dec},h}^*$. One can show via the lack of memory in the dynamics and the structure of π^* that

$$2484 \quad P_{\circlearrowleft,h+1}^* = W_{\circlearrowleft,h}^* \circ P_{\circlearrowleft,h}^*, \quad W_{\text{dec},h+1}^* = W_{\text{dec},h}^* \circ P_{\text{dec},h}^*, \quad (F.5)$$

2485 By the replica property (Lemma E.3), $W_{\circlearrowleft,h}^* \circ P_h^* = P_h^*$ for all h . Thus, for all h , (F.4) and (F.5) imply

$$2486 \quad P_{\circlearrowleft,h}^* = P_h^* = \frac{1}{4}\delta_{-\sigma} + \frac{3}{4}\delta_{+\sigma}. \quad (F.6)$$

2487 The following claim computes $P_{\text{dec},h}^*$.

2488 **Claim F.3.** Consider any distribution of the form $P = (1-p)\delta_{\sigma} + p\delta_{-\sigma}$. Then

$$2489 \quad W_{\text{dec},h}^* \circ P = \left(\frac{9}{10} - \frac{p}{5}\right)\delta_{\sigma} + \left(\frac{1}{10} + \frac{p}{5}\right)\delta_{-\sigma}.$$

2490 Thus,

$$2491 \quad P_{\text{dec},h+1}^*[-\sigma] = \frac{1}{10} \left(\sum_{i=0}^{h-1} 5^{-i} \right) + \frac{1}{4} 5^{1-h}.$$

2492 Before proving the claim, let us remark on its implications. As $h \rightarrow \infty$,

$$2493 \quad P_{\text{dec},h}^*[-\sigma] \rightarrow \frac{1}{10} \left(\frac{1}{1-1/5} \right) = \frac{1}{10} \cdot \frac{5}{4} = \frac{1}{8}.$$

2494 Thus,

$$2495 \quad \lim_{h \rightarrow \infty} P_{\text{dec},h}^* = \frac{7}{8}\delta_{\sigma} + \frac{1}{8}\delta_{-\sigma},$$

2496 achieving a different stationary distribution that $P_h^* = P_{\circlearrowleft,h}^*$. This shows that

$$2497 \quad \lim_{H \rightarrow \infty} \Gamma_{\text{marg},\sigma}(\pi_{\circlearrowleft\sigma}^*, \pi_{\text{dec}}^*) \geq \text{TV}\left(\frac{7}{8}\delta_{\sigma} + \frac{1}{8}\delta_{-\sigma}, \frac{3}{4}\delta_{\sigma} + \frac{1}{4}\delta_{-\sigma}\right) = \frac{1}{8},$$

2498 which implies that the deconvolution policy π_{dec}^* does approximate $\pi_{\circlearrowleft\sigma}^*$. From (F.6), it also follows that $\pi_{\circlearrowleft\sigma}^*$ and π^* have identical marginals, so

$$2499 \quad \lim_{H \rightarrow \infty} \Gamma_{\text{marg},\sigma}(\pi^*, \pi_{\text{dec}}^*) \geq \text{TV}\left(\frac{7}{8}\delta_{\sigma} + \frac{1}{8}\delta_{-\sigma}, \frac{3}{4}\delta_{\sigma} + \frac{1}{4}\delta_{-\sigma}\right) = \frac{1}{8}$$

2500 as well. In particular, if we take $\hat{\pi} = \pi_{\text{dec}}^*$ in Theorem 2, we see that there is no hope to for bounding $\Gamma_{\text{marg},\varepsilon}(\pi^*, \hat{\pi})$; we must bound $\Gamma_{\text{marg},\varepsilon}(\pi^*, \hat{\pi} \circ W_{\sigma})$ (again noting that if $\hat{\pi} = \pi_{\text{dec}}^*$, $\hat{\pi} \circ W_{\sigma} = \pi_{\circlearrowleft\sigma}^*$).

2501

2530 *Proof of Claim F.3.* We have that for $s' \in \{-\sigma, \sigma\}$,

$$2531$$
$$2532 W_{\text{dec},s'|s}^* = \frac{W_\sigma(s')[s] \cdot P_h^*(s')}{W_\sigma(s')[s] \cdot P_h^*(s') + W_\sigma(-s')[s] \cdot P_h^*(-s')}$$
$$2533$$

2534 With $s = s' = \sigma$, the above is

$$2535$$
$$2536 W_{\text{dec},h}^*(s' = \sigma | s = \sigma) = \frac{\frac{3}{4} \cdot \frac{3}{4}}{\frac{3}{4} \cdot \frac{3}{4} + \frac{1}{4} \cdot \frac{1}{4}} = \frac{9}{10}.$$
$$2537$$
$$2538$$

2539 And

$$2540$$
$$2541 W_{\text{dec},h}^*(s' = \sigma | s = -\sigma) = \frac{\frac{1}{4} \cdot \frac{3}{4}}{\frac{1}{4} \cdot \frac{3}{4} + \frac{3}{4} \cdot \frac{1}{4}} = \frac{1}{2}.$$
$$2542$$

2543 Hence, for any $p \in [0, 1]$,

$$2544$$
$$2545 W_{\text{dec},h}^*(s' = \sigma | s = -\sigma)((1-p)\delta_\sigma + p\delta_{-\sigma}) = ((1-p)\frac{9}{10} + p\frac{1}{2})\delta_\sigma + (1 - ((1-p)\frac{9}{10} + p\frac{1}{2}))\delta_{-\sigma}$$
$$2546$$
$$2547 = (\frac{9}{10} - \frac{p}{5})\delta_\sigma + (\frac{1}{10} + \frac{p}{5})\delta_{-\sigma}.$$
$$2548$$
$$2549$$

2550 Consequently, by (F.5), we can unfold a recursion to compute

$$2551$$
$$2552 P_{\text{dec},h+1}^*[-\sigma] = W_{\text{dec},h}^*(s' = \sigma | s = -\sigma)P_{\text{dec},h}^*$$
$$2553 = (\frac{1}{10} + \frac{P_{\text{dec},h}^*[\sigma]}{5})$$
$$2554 = \frac{1}{10} \sum_{i=0}^{h-1} 5^{-i} + P_{\text{dec},1}^*[\sigma] \cdot 5^{1-h}$$
$$2555$$
$$2556 = \frac{1}{10} \sum_{i=0}^{h-1} 5^{-i} + P_1^*[\sigma] \cdot 5^{1-h}$$
$$2557$$
$$2558 = \frac{1}{10} \left(\sum_{i=0}^{h-1} 5^{-i} \right) + \frac{1}{4} 5^{1-h}.$$
$$2559$$
$$2560$$
$$2561$$
$$2562$$
$$2563$$
$$2564$$
$$2565$$
$$2566$$

□

2567 Part II

2568 The Control Setting

2570 G. Stability in the Control System

2571 This section proves our various stability conditions. One wrinkle in the exposition is that we are able to derive far sharper perturbation guarantees than are needed in our analysis. However, as the guarantees are rather technically burdensome to derive, we endeavor to present the sharpest possible results so that we may save others from having to rederive these bounds in future applications.

2572 Importantly, this section also contains the definition of the constants $c_1, \dots, c_5 > 0$ present in [Theorem 1](#), [Proposition 4.1](#), and other main results (see [Definition G.7](#)).

2573 The section is organized as follows:

- 2581 • [Appendix G.1](#) recalls various preliminaries.
- 2582
- 2583
- 2584

- [Appendix G.2](#) provides the definition of numerous problem-dependent constants, all of which are polynomial in $(R_{\text{dyn}}, L_{\text{dyn}}, M_{\text{dyn}})$ and $(R_{\text{stab}}, B_{\text{stab}}, L_{\text{stab}})$ defined in [Assumptions 3.1](#) and [3.2](#).
- [Appendix G.3](#) gives IPS guarantees in terms of the constants in the previous section. Specifically, it provides [Definition G.7](#), which instantiates the constants $c_1, \dots, c_5 > 0$ present in [Theorem 1](#), [Proposition 4.1](#), and other main results. We then state [Corollary G.1](#), from which we derive [Proposition 4.1](#) used in the body. This corollary is derived from a sharper guarantee, [Proposition G.3](#) (whose improvements over the corollary are detailed in [Remark G.2](#)).
- The results in [Appendix G.3](#) are derived from two building blocks in [Appendix G.4](#): [Lemma G.4](#) which bounds sensitive of regular trajectories to initial state, and [Proposition G.5](#) which addresses perturbations of control inputs and gain.
- [Proposition G.3](#) is derived from [Proposition G.5](#) in [Appendix G.5](#). [Lemma G.4](#) and [Proposition G.5](#) are proven in [Appendix G.7](#) in [Appendix G.7](#), respectively.
- [Appendix G.8](#) explains how to implement a synthesis oracle which produces Jacobian Stabilizing primitive controllers from trajectories which satisfy a natural stabilizability condition.
- Finally, [Appendix G.9](#) gives the solutions to various scalar recursions used in the proofs of [Lemma G.4](#) and [Proposition G.5](#).

G.1. Recalling preliminaries and assumptions.

Recall the following definitions.

- A length- K control trajectory is denoted $\boldsymbol{\rho} = (x_{1:K+1}, u_{1:K}) \in \mathcal{P}_K = (\mathbb{R}^{d_x})^{K+1} \times (\mathbb{R}^{d_u})^K$.
- Its Jacobian linearizations are denoted $\mathbf{A}_k(\boldsymbol{\rho}) := \frac{\partial}{\partial \mathbf{x}} f_\eta(\mathbf{x}_k, \mathbf{u}_k)$ and $\mathbf{B}_k(\boldsymbol{\rho}) := \frac{\partial}{\partial \mathbf{u}} f_\eta(\mathbf{x}_k, \mathbf{u}_k)$ for $k \in [K]$.
- Recalling our dynamics map $f(\cdot, \cdot)$, and step size $\eta > 0$, we say $\boldsymbol{\rho}$ is *feasible* if, for all $k \in [K]$,

$$\mathbf{x}_{k+1} = f(\mathbf{x}_k, \mathbf{u}_k), \quad \text{where } f(\mathbf{x}, \mathbf{u}) = \mathbf{x} + \eta f_\eta(\mathbf{x}, \mathbf{u}).$$

We regular the definition of regular trajectories from [Section 3](#).

Definition G.1. A control path $\boldsymbol{\rho} = (\mathbf{x}_{1:K+1}, \mathbf{u}_{1:K})$ is $(R_{\text{dyn}}, L_{\text{dyn}}, M_{\text{dyn}})$ -regular if for all $k \in [K]$ and all $(\mathbf{x}'_k, \mathbf{u}'_k) \in \mathbb{R}^{d_x} \times \mathbb{R}^{d_u}$ such that $\|\mathbf{x}'_k - \mathbf{x}_k\| \vee \|\mathbf{u}_k - \mathbf{u}'_k\| \leq R_{\text{dyn}}$,⁵

$$\|\nabla f_\eta(\mathbf{x}'_k, \mathbf{u}'_k)\|_{\text{op}} \leq L_{\text{dyn}}, \quad \|\nabla^2 f_\eta(\mathbf{x}'_k, \mathbf{u}'_k)\|_{\text{op}} \leq M_{\text{dyn}}.$$

We also recall the definitions around Jacobian stabilization. We start with a definition of Jacobian stabilization for feedback gains, from which we then recover the definition of Jacobian stabilization for primitive controllers given in the body.

Definition G.2. Consider $R_{\text{stab}}, L_{\text{stab}}, B_{\text{stab}} \geq 1$. Consider sequence of gains $\mathbf{K}_{1:K} \in (\mathbb{R}^{d_u \times d_u})^K$ and trajectory $\boldsymbol{\rho} = (\mathbf{x}_{1:K+1}, \mathbf{u}_{1:K}) \in \mathcal{P}_K$. We say that $(\boldsymbol{\rho}, \mathbf{K}_{1:K})$ -is $(R_{\text{stab}}, B_{\text{stab}}, L_{\text{stab}})$ -Jacobian Stable if $\max_k \|\mathbf{K}_k\|_{\text{op}} \leq B_{\text{stab}}$, and if the closed-loop transition operators defined by

$$\Phi_{\text{cl},k,j} := (\mathbf{I} + \eta \mathbf{A}_{\text{cl},k-1}) \cdot (\mathbf{I} + \eta \mathbf{A}_{\text{cl},k-2}) \cdot (\dots) \cdot (\mathbf{I} + \eta \mathbf{A}_{\text{cl},j})$$

with $\mathbf{A}_{\text{cl},k} = \mathbf{A}_k(\boldsymbol{\rho}) + \mathbf{B}_{k-1}(\boldsymbol{\rho})\mathbf{K}_{k-1}$ satisfies the following inequality

$$\|\Phi_{\text{cl},k,j}\|_{\text{op}} \leq B_{\text{stab}} \left(1 - \frac{\eta}{L_{\text{stab}}}\right)^{k-j}.$$

The definition of Jacobian stability of primitive controllers in [Section 3](#) may be recovered as follows.

Definition G.3. Consider $R_{\text{stab}}, L_{\text{stab}}, B_{\text{stab}} \geq 1$. Consider a sequence of primitive controllers $\kappa_{1:K} \in \mathcal{K}^K$, each expressed as $\kappa_k(\mathbf{x}) = \bar{\mathbf{u}}_k = \bar{\mathbf{K}}_k(\mathbf{x}_k - \bar{\mathbf{x}}_k)$ and $\boldsymbol{\rho} = (\mathbf{x}_{1:K+1}, \mathbf{u}_{1:K}) \in \mathcal{P}_K$. We say $(\boldsymbol{\rho}, \kappa_{1:K})$ is Jacobian Stable if $\kappa_{1:K}$ is consistent with $\boldsymbol{\rho}$, and if $(\boldsymbol{\rho}, \bar{\mathbf{K}}_{1:K})$ is $R_{\text{stab}}, L_{\text{stab}}, B_{\text{stab}} > 0$ -Jacobian stable.

Note that in Jacobian stability (both with primitive controllers and with gain-matrices), we take all parameters to be no less than one.

⁵Here, $\|\nabla^2 f_\eta(\mathbf{x}'_t, \mathbf{u}'_t)\|_{\text{op}}$ denotes the operator-norm of a three-tensor.

2640 G.1.1. PROPERTIES SATISFIED BY π^*

2641 Finally, we show that actions produced by π^* in our control instantiation of the composite MDP satisfy the assumptions in
 2642 [Assumptions 3.1](#) and [3.2](#).

2644 **Lemma G.1.** *Suppose [Assumptions 3.1](#) and [3.2](#) hold. Let $\pi^* = (\pi_h^*)_{1 \leq h \leq H}$ denote the policy constructed as a regular
 2645 conditional probability from the conditionals of \mathcal{D}_{exp} . Furthermore, let $\mathcal{D}_{\text{exp}, \rho_{m,h}}$ denote the distribution over $\rho_{m,h}$
 2646 corresponding to $\rho_T \sim \mathcal{D}_{\text{exp}}$. Then, with probability one over $\rho_{m,h} \sim \mathcal{D}_{\text{exp}, \rho_{m,h}}$ and $\mathbf{a}_h \sim \pi_h^*(\rho_{m,h})$, expressed as
 2647 $\rho_{m,h} = (\mathbf{x}_{t_h:t_h-\tau_m+1}, \mathbf{u}_{t_h-1:t_h-\tau_m+1})$, and $\mathbf{a}_h = \kappa_{t_h:t_{h+1}-1}$. Consider the unique feasible trajectory for which*

$$2648 \quad \rho_{c,h+1} = (\mathbf{x}'_{t_h:t_{h+1}}, \mathbf{u}'_{t_h:t_{h+1}-1}), \quad \mathbf{x}'_{t_h} = \mathbf{x}_{t_h}, \quad \mathbf{u}_t = \kappa_t(\mathbf{x}_t), \quad t_h \leq t < t_{h+1}.$$

2650 Then,

- 2651 • $\rho'_{c,h+1}$ is $(R_{\text{dyn}}, L_{\text{dyn}}, M_{\text{dyn}})$ -regular
- 2652 • $(\rho'_{c,h+1}, \kappa_{t_h:t_{h+1}-1})$ is $(R_{\text{stab}}, B_{\text{stab}}, L_{\text{stab}})$ -Jacobian stable.

2653 *Proof.* Since π^* in [Definition I.3](#) is constructed as the regular conditional probability of $\mathbf{a}_h \mid \rho_{m,h}$ under \mathcal{D}_{exp} , $(\mathbf{a}_h, \rho_{m,h})$
 2654 is the above lemma have the same joint distribution as under \mathcal{D}_{exp} . Thus, the lemma follows from the assumptions
 2655 [Assumptions 3.1](#) and [3.2](#) placed on \mathcal{D}_{exp} . \square

2656 The following is a direct consequence of the above lemma.

2661 **Lemma G.2.** *Consider the instantiation of the composite MDP for the control setting as in [Section 4.1](#) and in [Appendix I](#),
 2662 with π^* as in [Definition I.3](#), and $\phi_{\mathcal{Z}}$ as in [Definition E.1](#). Suppose that $W_1, \dots, W_h : \mathcal{S} \rightarrow \Delta(\mathcal{S})$ satisfy⁶*

$$2663 \quad \phi_{\mathcal{Z}} \circ W_h(\mathbf{s}) \ll \phi_{\mathcal{Z}} \circ P_h^*, \tag{G.1}$$

2664 Consider a sequence of actions $\mathbf{s}_{1:H+1}, \mathbf{a}_{1:H}$ generated via

$$2665 \quad \mathbf{a}_h \sim \pi_h^*(\tilde{\mathbf{s}}_h), \quad \tilde{\mathbf{s}}_h \sim W_h(\mathbf{s}_h), \quad \mathbf{s}_{h+1} = F_h(\mathbf{s}_h, \mathbf{a}_h), \quad \mathbf{s}_1 \sim P_{\text{init}}.$$

2666 Let $\tilde{\mathbf{s}}_h = (\tilde{\mathbf{x}}_{t_{h-1}:t_h}, \tilde{\mathbf{u}}_{t_{h-1}:t_h-1})$ and $\mathbf{a}_h = \kappa_{t_h:t_{h+1}-1}$. Then, with probability one, for each h , the unique feasible trajectory
 2667 for which

$$2668 \quad \rho_{c,h+1} = (\mathbf{x}'_{t_h:t_{h+1}}, \mathbf{u}'_{t_h:t_{h+1}-1}), \quad \mathbf{x}'_{t_h} = \mathbf{x}_{t_h}, \quad \mathbf{u}_t = \kappa_t(\mathbf{x}_t), \quad t_h \leq t < t_{h+1}.$$

2669 satisfies

- 2670 • $\mathbf{x}'_{t_h} = \tilde{\mathbf{x}}_{t_h}$, and $\rho'_{c,h+1}$ is feasible and $(R_{\text{dyn}}, L_{\text{dyn}}, M_{\text{dyn}})$ -regular
- 2671 • $(\rho'_{c,h+1}, \kappa_{t_h:t_{h+1}-1})$ is $(R_{\text{stab}}, B_{\text{stab}}, L_{\text{stab}})$ -Jacobian Stable.

2672 **Remark G.1** (On the absolute continuity constraint in [\(G.1\)](#)). Recall that $\phi_{\mathcal{Z}}$ as defined in [Definition I.1](#) simply extracts the
 2673 memory chunk $\rho_{m,h}$ from the trajectory chunk $\rho_{c,h}$. The condition $\text{supp}(\phi_{\mathcal{Z}} \circ W_h(\mathbf{s})) \subset \text{supp}(\phi_{\mathcal{Z}} \circ P_h^*)$ just means that
 2674 the distribution of the memory chunk-components from $W_h(\mathbf{s})$ is absolutely continuous with respect to the memory-chunks
 2675 $\rho_{m,h}$ under \mathcal{D}_{exp} .

2683 G.1.2. NORM NOTATION.

2684 Lastly, given our parameter $\eta > 0$, we define two types of norms. First, for sequences of vectors $\mathbf{z}_{1:K} \in (\mathbb{R}^d)^K$ and matrices
 2685 $(\mathbf{X}_{1:K}) \in \mathbb{R}^{d_1 \times d_2)^K$, define

$$2686 \quad \|\mathbf{z}_{1:K}\|_{\ell_2} = \left(\eta \sum_{k=1}^K \|\mathbf{z}_k\|^2 \right), \quad \|\mathbf{X}_{1:K}\|_{\ell_2, \text{op}} = \left(\eta \sum_{k=1}^K \|\mathbf{X}_k\|^2 \right),$$

2687 where again the standard $\|\cdot\|$ notation denotes Euclidean norm for vectors and operator norm for matrices. We also define

$$2688 \quad \|\mathbf{z}_{1:K}\|_{\max, 2} = \max_{1 \leq k \leq K} \|\mathbf{z}_k\|, \quad \|\mathbf{X}_{1:K}\|_{\max, \text{op}} = \max_{1 \leq k \leq K} \|\mathbf{X}_k\|.$$

2693 ⁶Recall the absolute-continuity comparator \ll defined in [Definition C.4](#).

2695 G.2. Composite Problem Constants

2696 We begin by writing down numerous problem constants, all of which are polynomial in the quantities $(R_{\text{dyn}}, L_{\text{dyn}}, M_{\text{dyn}})$
 2697 and $(R_{\text{stab}}, B_{\text{stab}}, L_{\text{stab}})$. First, we define the *stability exponent*,

$$2699 \beta_{\text{stab}} := \left(1 - \frac{\eta}{L_{\text{stab}}}\right) \in (0, 1).$$

2700 **Definition G.4.** Given the regularity parameters $R_{\text{dyn}}, L_{\text{dyn}}, M_{\text{dyn}}$, stability parameters $R_{\text{stab}}, B_{\text{stab}}, L_{\text{stab}}$, and the step
 2701 size $\eta > 0$, we define the ‘‘little-c’’ constants

$$2704 c_{\mathbf{u}} = 12B_{\text{stab}}\sqrt{L_{\text{stab}}L_{\text{dyn}}}, \quad c_{\mathbf{K}} = 2B_{\text{stab}} + 12B_{\text{stab}}L_{\text{stab}}^{1/2}L_{\text{dyn}}, \quad c_{\Delta} = 6B_{\text{stab}}$$

2706 as well as ‘‘big-C’’ constants

$$2707 C_{\mathbf{u}} := \min \left\{ \frac{\sqrt{L_{\text{stab}}L_{\text{dyn}}}}{M_{\text{dyn}}}, \frac{1}{256B_{\text{stab}}^2 M_{\text{dyn}} L_{\text{dyn}} L_{\text{stab}}^{3/2}} \right\}$$

$$2710 C_{\Delta} := \frac{1}{4 \cdot 324B_{\text{stab}}^2 M_{\text{dyn}} L_{\text{stab}}}$$

$$2711 C_{\hat{\mathbf{x}}} := \min \left\{ \frac{R_{\text{dyn}}}{2R_{\text{stab}}B_{\text{stab}}}, \frac{1}{16L_{\text{stab}}M_{\text{dyn}}R_{\text{stab}}^2 B_{\text{stab}}^2} \right\}$$

$$2712 C_{\mathbf{K}} := \min \left\{ \frac{1}{24\sqrt{L_{\text{stab}}B_{\text{stab}}L_{\text{dyn}}}}, \frac{C_{\hat{\mathbf{x}}}^{-1}L_{\text{dyn}}}{8 \cdot 324B_{\text{stab}}^2 M_{\text{dyn}} L_{\text{stab}}^{3/2}} \right\}$$

$$2713 C_{\mathbf{K}, \hat{\mathbf{x}}} := \frac{L_{\text{dyn}}}{8 \cdot 324B_{\text{stab}}^2 M_{\text{dyn}} L_{\text{stab}}^{3/2}}.$$

2721 The ‘‘little-c’’ constants enter directly into our error bounds, where as the ‘‘big-C’’ constants function as constraints on errors,
 2722 above which we lose guarantees. We define some additional ‘‘big-C’’ constants which take in a radius argument R_0 .

2723 **Definition G.5** (Final Stability Constants). In term of the constants in [Definition G.5](#), we define the following final stability
 2724 constants, as functions of a parameter R_0 :

$$2726 C_{\text{stab},1}(R_0) := \min \left\{ C_{\mathbf{u}}, \frac{C_{\Delta}}{4c_{\mathbf{u}}}, \frac{R_{\text{dyn}}}{R_0} \cdot \frac{1}{48c_{\mathbf{u}}c_{\Delta}} \right\}$$

$$2727 C_{\text{stab},2}(R_0) := \min \left\{ C_{\mathbf{K}}, \frac{\beta_{\text{stab}}^{-\tau_c/3} C_{\Delta}}{4c_{\mathbf{u}}}, \frac{R_{\text{dyn}}}{R_0} \cdot \frac{\beta_{\text{stab}}^{-\tau_c/3}}{48c_{\mathbf{K}}c_{\Delta}} \right\}$$

$$2728 C_{\text{stab},3}(R_0) := \frac{R_{\text{dyn}}}{12R_0 c_{\mathbf{u}} \sqrt{L_{\text{stab}}} + 3}$$

$$2729 C_{\text{stab},4}(R_0) := \min \left\{ C_{\hat{\mathbf{x}}}, \frac{C_{\mathbf{K}, \hat{\mathbf{x}}}}{C_{\mathbf{K}}}, \frac{R_{\text{dyn}}}{R_0} \cdot \frac{1}{12c_{\mathbf{K}}} \right\}$$

2736 G.3. IPS Guarantees & Proof of [Proposition 4.1](#)

2738 Here we provide our main stability guarantees for the learned policy π^* under [Assumptions 3.1](#) and [3.2](#), from which we
 2739 derive [Proposition 4.1](#). This section adopts the notation from [Section 4.1](#).

2740 We begin by introducing a functional form for our distances.

2741 **Definition G.6** (Distances). Let τ_c be given, and let $0 \leq \tau \leq \tau_c$. For $h > 1$ and chunk-states $\mathbf{s}_h = (\mathbf{x}_{t_{h-1}:t_h}, \mathbf{u}_{t_{h-1}:t_h-1}) \in$
 2742 \mathcal{P}_{τ_c} and $\mathbf{s}'_h = (\mathbf{x}'_{t_{h-1}:t_h}, \mathbf{u}'_{t_{h-1}:t_h-1})$, define

$$2744 d_{\mathcal{S}, \mathbf{x}, \tau}(\mathbf{s}_h, \mathbf{s}'_h) := \max_{t \in [t_h - \tau, t_h]} \|\mathbf{x}_t - \mathbf{x}'_t\|$$

$$2745 d_{\mathcal{S}, \mathbf{u}, \tau}(\mathbf{s}_h, \mathbf{s}'_h) := \max_{t \in [t_h - \tau, t_h - 1]} \|\mathbf{u}_t - \mathbf{u}'_t\|$$

$$2746 d_{\mathcal{S}, \tau}(\mathbf{s}_h, \mathbf{s}'_h) := \max \{d_{\mathcal{S}, \mathbf{x}, \tau}(\mathbf{s}_h, \mathbf{s}'_h), d_{\mathcal{S}, \mathbf{u}, \tau}(\mathbf{s}_h, \mathbf{s}'_h)\},$$

For $h = 1$, $\mathbf{s}_1 = \bar{\mathbf{x}}_1$ and $\mathbf{s}'_1 = \bar{\mathbf{x}}'_1$, we define $d_{\mathcal{S}}(\mathbf{s}_1, \mathbf{s}'_1) = \|\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}'_1\|$. Note therefore that

$$d_{\mathcal{S}, \tau_c}(\cdot, \cdot) = d_{\mathcal{S}}(\cdot, \cdot), \quad d_{\mathcal{S}, \tau_m - 1}(\cdot, \cdot) = d_{\text{TV}}(\cdot, \cdot), \quad d_{\mathcal{S}, 0}(\cdot, \cdot) = d_{\text{IPS}}(\cdot, \cdot)$$

Next, we introduce five problem-dependent constants c_1, \dots, c_5 , all of which are stated in terms of the constants in [Appendix G.2](#); one can readily check that these are all polynomial in the constants $(R_{\text{dyn}}, L_{\text{dyn}}, M_{\text{dyn}})$ and $(R_{\text{stab}}, B_{\text{stab}}, L_{\text{stab}})$ in [Assumptions 3.1](#) and [3.2](#).

Definition G.7 (IPS Constants). In terms of constants in [Appendix G.2](#), we define the IPS constants as follows:

$$\begin{aligned} c_1 &:= (6 \max\{R_{\text{stab}}(1 + 2c_{\mathbf{u}}\sqrt{L_{\text{stab}}}), B_{\text{stab}} + \sqrt{L_{\text{stab}}}c_{\mathbf{K}}\}). \\ c_2 &:= \min\left\{\frac{C_{\text{stab},1}(2R_{\text{stab}})}{4R_{\text{stab}}\sqrt{L_{\text{stab}}}}, \frac{C_{\text{stab},2}(2R_{\text{stab}})}{\sqrt{L_{\text{stab}}}}, \frac{C_{\text{stab},3}(2R_{\text{stab}})}{4R_{\text{stab}}}, \frac{1}{2R_{\text{stab}}}\right\}. \end{aligned} \tag{G.2}$$

We further define

$$c_3 = 3L_{\text{stab}} \log(2c_{\Delta}), \quad c_4 = \min\{1, C_{\text{stab},4}(2R_{\text{stab}})\}, \quad c_5 = 2(1 + R_{\text{stab}})B_{\text{stab}}.$$

In terms of the constants $c_1, c_2 > 0$ above, we introduce a family of distance-like functions on $\bar{d}_{\mathcal{A}, \tau}(\mathbf{a}, \mathbf{a}' | r) : \mathcal{A} \times \mathcal{A} \rightarrow \mathbb{R}_{\geq 0} \cup \{\infty\}$, defined as follows.

Definition G.8. Consider $\mathbf{a} = (\bar{\mathbf{u}}_{1:\tau_c}, \bar{\mathbf{x}}_{1:\tau_c}, \bar{\mathbf{K}}_{1:\tau_c})$ and $\mathbf{a}' = (\bar{\mathbf{u}}'_{1:\tau_c}, \bar{\mathbf{x}}'_{1:\tau_c}, \bar{\mathbf{K}}'_{1:\tau_c})$.

$$\begin{aligned} \bar{d}_{\mathcal{A}, \tau}(\mathbf{a}, \mathbf{a}' | r) &:= c_1 \max_{1 \leq k \leq \tau_c} \left(\|\bar{\mathbf{u}}_k - \bar{\mathbf{u}}'_k\| + \|\bar{\mathbf{x}}_k - \bar{\mathbf{x}}'_k\| + r e^{-\frac{\eta(\tau_c - \tau)}{3L_{\text{stab}}}} \|\bar{\mathbf{K}}_k - \bar{\mathbf{K}}'_k\| \right) \\ &\quad + \mathbf{I}_{0, \infty} \left\{ \max_{1 \leq k \leq \tau_c} (\max\{\|\bar{\mathbf{u}}_k - \bar{\mathbf{u}}'_k\|, \|\bar{\mathbf{x}}_k - \bar{\mathbf{x}}'_k\|, \|\bar{\mathbf{K}}_k - \bar{\mathbf{K}}'_k\|\}) \leq c_2 \right\}, \end{aligned}$$

where $\mathbf{I}_{0, \infty}(\mathcal{E})$ is 0 if class \mathcal{E} is true and ∞ otherwise.

In words, $\bar{d}_{\mathcal{A}, \tau}(\mathbf{a}, \mathbf{a}' | r)$ measures the maximal differences between $\bar{\mathbf{u}}_k - \bar{\mathbf{u}}'_k$, $\bar{\mathbf{x}}_k - \bar{\mathbf{x}}'_k$, and $\bar{\mathbf{K}}_k - \bar{\mathbf{K}}'_k$, subject to a constraint that each of these quantities is within some bound c_2 . One the latter threshold is met, the dependence on $\|\bar{\mathbf{K}}_k - \bar{\mathbf{K}}'_k\|$ is scaled down by r , and is also exponentially small in $\tau_c - \tau$; this latter bit is not necessary for our results, but illustrates an interesting feature of our stability guarantees: *they are far less sensitive to errors in $\bar{\mathbf{K}}$ than to errors in $\bar{\mathbf{u}}$.*

In terms of $\bar{d}_{\mathcal{A}, \tau}(\mathbf{a}, \mathbf{a}' | r)$ defined above, we can now ensure the following stability guarantee.

Corollary G.1. Suppose that $\tau_c \geq c_3/\eta$ and $r \leq c_4$, and consider any sequence of kernels $\{W_h\}_{h=1}^H$, where $W_h : \mathcal{S} \rightarrow \Delta(\mathcal{S})$, and⁷

$$\max_{h, \mathbf{s} \in \mathcal{S}} \mathbb{P}_{\tilde{\mathbf{s}} \sim W_h(\mathbf{s})} [d_{\text{IPS}}(\tilde{\mathbf{s}}, \mathbf{s}) \leq r] = 1, \quad \phi_{\mathcal{Z}} \circ W_h(\mathbf{s}) \ll \phi_{\mathcal{Z}} \circ P_h^*,$$

for $\phi_{\mathcal{Z}}$ is from the direct decomposition instantiated in [Definition I.1](#), and where P_h^* denotes the law of $\rho_{c, h}$ under \mathcal{D}_{exp} as in [Definition I.3](#).

Define a process $\mathbf{s}_1 \sim P_{\text{init}}$, $\tilde{\mathbf{s}}_h \sim W_h(\mathbf{s}_h)$, $\mathbf{a}_h \sim \pi_h^*(\tilde{\mathbf{s}}_h)$, and $\mathbf{s}_{h+1} := F_h(\mathbf{s}_h, \mathbf{a}_h)$. Then, almost surely, the following hold for all $0 \leq \tau \leq \tau_c$:

- For each $1 \leq h \leq H$, $d_{\mathcal{S}, \tau}(F_h(\tilde{\mathbf{s}}_h, \mathbf{a}_h), \mathbf{s}_h) \leq c_5 r e^{-\frac{\eta(\tau_c - \tau)}{L_{\text{stab}}}}$.
- For any sequence $(\mathbf{a}'_{1:H})$, the dynamics $\mathbf{s}'_1 = \mathbf{s}_1$, $\mathbf{s}_{h+1} = F_h(\mathbf{s}'_h, \mathbf{a}'_h)$ satisfy

$$\max_{1 \leq h \leq H+1} d_{\mathcal{S}, \tau}(\mathbf{s}_h, \mathbf{s}'_h) \leq \max_{1 \leq h \leq H} \bar{d}_{\mathcal{A}, \tau}(\mathbf{a}_h, \mathbf{a}'_h | r).$$

[Corollary G.1](#) is derived in [Appendix G.3.2](#) from an even more granular result stated just below. Before continuing, we explain how [Proposition 4.1](#) follows.

Proof of [Proposition 4.1](#). This follows directly from the above corollary notice that c_4 is define to be at most 1, so we always invoke the corollary with $r \leq 1$, and thus $\bar{d}_{\mathcal{A}, \tau}(\mathbf{a}_h, \mathbf{a}'_h | r) \leq \bar{d}_{\mathcal{A}, \tau}(\mathbf{a}_h, \mathbf{a}'_h | 1) \leq d_{\mathcal{A}}$. We remark that the guarantee only applies to kernels for which \square

⁷See [Remark G.1](#) above for interpretation of this condition below.

2805 G.3.1. A MORE GRANULAR STABILITY STATEMENT

2806 Here, we state an even more granular stability guarantee. The notation is rather onerous, but captures another nice feature of
 2807 our bound: that our stability depends not on the maximal errors over $\bar{\mathbf{u}}_k - \bar{\mathbf{u}}'_k$, $\bar{\mathbf{x}}_k - \bar{\mathbf{x}}'_k$, $\bar{\mathbf{K}}_k - \bar{\mathbf{K}}'_k$, but rather on ℓ_2 -errors.
 2808 Again, not necessary for our guarantees, but it speaks to the sharpness of our perturbation bounds. See [Remark G.2](#) at the
 2809 end of the section for more discussion.

2811 **Definition G.9** (Action Differences (inputs and gains)). Consider $\mathbf{a} = (\bar{\mathbf{u}}_{1:\tau_c}, \bar{\mathbf{x}}_{1:\tau_c}, \bar{\mathbf{K}}_{1:\tau_c})$ and $\mathbf{a}' = (\bar{\mathbf{u}}'_{1:\tau_c}, \bar{\mathbf{x}}'_{1:\tau_c}, \bar{\mathbf{K}}'_{1:\tau_c})$.
 2812 We define

$$\begin{aligned} 2813 \quad d_{\mathcal{A},\mathbf{u},\ell_2}(\mathbf{a}, \mathbf{a}') &= \max_{1 \leq k \leq \tau_c} \left(\eta \sum_{j=1}^k \beta_{\text{stab}}^{k-j} \|\bar{\mathbf{u}}_j - \bar{\mathbf{u}}'_j\|^2 \right)^{1/2} \\ 2814 \quad d_{\mathcal{A},\mathbf{x},\ell_2}(\mathbf{a}, \mathbf{a}') &= \max_{1 \leq k \leq \tau_c} \left(\eta \sum_{j=1}^k \beta_{\text{stab}}^{k-j} \|\bar{\mathbf{x}}_j - \bar{\mathbf{x}}'_j\|^2 \right)^{1/2} \\ 2815 \quad d_{\mathcal{A},\mathbf{K},\ell_2}(\mathbf{a}, \mathbf{a}') &= \max_{1 \leq k \leq \tau_c} \left(\eta \sum_{j=1}^k \beta_{\text{stab}}^{k-j} \|\bar{\mathbf{K}}_j - \bar{\mathbf{K}}'_j\|^2 \right)^{1/2}. \end{aligned}$$

2824 We further define

$$\begin{aligned} 2826 \quad d_{\mathcal{A},\mathbf{u},\infty}(\mathbf{a}, \mathbf{a}') &:= \max_{1 \leq k \leq \tau_c} \|\bar{\mathbf{u}}_k - \bar{\mathbf{u}}'_k\| = \|\bar{\mathbf{u}}_{1:\tau_c} - \bar{\mathbf{u}}'_{1:\tau_c}\|_{\max,2} \\ 2827 \quad d_{\mathcal{A},\mathbf{x},\infty}(\mathbf{a}, \mathbf{a}') &:= \max_{1 \leq k \leq \tau_c} \|\bar{\mathbf{x}}_k - \bar{\mathbf{x}}'_k\| = \|\bar{\mathbf{x}}_{1:\tau_c} - \bar{\mathbf{x}}'_{1:\tau_c}\|_{\max,2} \\ 2828 \quad d_{\mathcal{A},\mathbf{K},\infty}(\mathbf{a}, \mathbf{a}') &:= \max_{1 \leq k \leq \tau_c} \|\bar{\mathbf{K}}_k - \bar{\mathbf{K}}'_k\| = \|\bar{\mathbf{K}}_{1:\tau_c} - \bar{\mathbf{K}}'_{1:\tau_c}\|_{\max,\text{op}}. \end{aligned}$$

2832 and

$$2833 \quad \text{rad}_{\mathbf{K}}(\mathbf{a}) := \max_{1 \leq k \leq \tau_c} \|\bar{\mathbf{K}}_k\| = \|\bar{\mathbf{K}}_{1:\tau_c}\|_{\max,\text{op}}.$$

2836 We note further that as $\beta_{\text{stab}} \in (0, 1)$ and $\eta \sum_{i \geq 0} \beta_{\text{stab}}^i = L_{\text{stab}}$, we have

$$\begin{aligned} 2837 \quad d_{\mathcal{A},\mathbf{u},\ell_2}(\mathbf{a}, \mathbf{a}') &\leq \|\bar{\mathbf{u}}_{1:\tau_c} - \bar{\mathbf{u}}'_{1:\tau_c}\|_{\ell_2} \wedge \sqrt{L_{\text{stab}}} \|\bar{\mathbf{u}}_{1:\tau_c} - \bar{\mathbf{u}}'_{1:\tau_c}\|_{\max,2} \\ 2838 \quad &\leq \|\bar{\mathbf{u}}_{1:\tau_c} - \bar{\mathbf{u}}'_{1:\tau_c}\|_{\ell_2} \wedge \sqrt{L_{\text{stab}}} d_{\mathcal{A},\mathbf{u},\infty}(\mathbf{a}, \mathbf{a}') \end{aligned} \quad (\text{G.3})$$

2841 and analogously for $d_{\mathcal{A},\mathbf{x},\ell_2}(\mathbf{a}, \mathbf{a}')$ and $d_{\mathcal{A},\mathbf{K},\ell_2}(\mathbf{a}, \mathbf{a}')$.

2843 Next, recall the constants $\{C_{\text{stab},i}(R_0)\}_{i=1}^4$ in [Definition G.5](#), and $c_{\mathbf{u}}, c_{\mathbf{K}}$ in [Definition G.4](#), all of which are polynomial in
 2844 relevant problem parameters $(R_{\text{dyn}}, L_{\text{dyn}}, M_{\text{dyn}})$, $(R_{\text{stab}}, B_{\text{stab}}, L_{\text{stab}})$, and argument R_0 . We now define a very general
 2845 distance-like function between actions.

2846 **Definition G.10** (Action Divergences). Define, for $R_0 \geq 1$, the following

$$2847 \quad d_{\mathcal{A},R_0,\tau}(\mathbf{a}_h, \mathbf{a}'_h \mid r) := 2((1 + R_0)d_{\mathcal{A},R_0,\tau,\mathbf{x}}(\mathbf{a}_h, \mathbf{a}'_h \mid r) + d_{\mathcal{A},R_0,\tau,\mathbf{u}}(\mathbf{a}_h, \mathbf{a}'_h \mid r)),$$

2850 where

$$2851 \quad d_{\mathcal{A},R_0,\tau,\mathbf{u}}(\mathbf{a}_h, \mathbf{a}'_h \mid r) := d_{\mathcal{A},\mathbf{u},\infty}(\mathbf{a}_h, \mathbf{a}'_h) + R_0 d_{\mathcal{A},\mathbf{x},\infty}(\mathbf{a}_h, \mathbf{a}'_h) + 2r B_{\text{stab}} \beta_{\text{stab}}^{\tau_c - \tau} d_{\mathcal{A},\mathbf{K},\infty}(\mathbf{a}_h, \mathbf{a}'_h),$$

2853 and where

$$\begin{aligned} 2854 \quad d_{\mathcal{A},R_0,\tau,\mathbf{x}}(\mathbf{a}, \mathbf{a}' \mid r) &= 2c_{\mathbf{u}}(d_{\mathcal{A},\mathbf{u},\ell_2}(\mathbf{a}, \mathbf{a}') + R_0 d_{\mathcal{A},\mathbf{x},\ell_2}(\mathbf{a}, \mathbf{a}')) + 2c_{\mathbf{K}} r (\beta_{\text{stab}})^{\frac{\tau_c - \tau}{3}} \cdot d_{\mathcal{A},\mathbf{K},\ell_2}(\mathbf{a}, \mathbf{a}') \\ 2855 \quad &+ \mathbf{I}_{0,\infty} \left\{ \bigcap_{i=1}^3 \mathcal{E}_{\text{close},R_0,i} \right\} + \mathbf{I}_{0,\infty} \{ \text{rad}_{\mathbf{K}}(\mathbf{a}) \vee \text{rad}_{\mathbf{K}}(\mathbf{a}') \leq R_0 \}, \end{aligned}$$

2859

2860 where $\mathbf{I}_{0,\infty}\{\mathcal{E}\}$ denotes 0 if clause \mathcal{E} is true, and ∞ otherwise, and where we define the clauses

$$2861 \mathcal{E}_{\text{close},R_0,1}(\mathbf{a}, \mathbf{a}') = \{(\mathbf{d}_{\mathcal{A},\mathbf{u},\ell_2}(\mathbf{a}, \mathbf{a}') + R_0 \mathbf{d}_{\mathcal{A},\mathbf{x},\ell_2}(\mathbf{a}, \mathbf{a}') \leq C_{\text{stab},1}(R_0)\}$$

$$2862 \mathcal{E}_{\text{close},R_0,2}(\mathbf{a}, \mathbf{a}') = \{\mathbf{d}_{\mathcal{A},\mathbf{K},\ell_2}(\mathbf{a}, \mathbf{a}') \leq C_{\text{stab},2}(R_0)\}$$

$$2863 \mathcal{E}_{\text{close},R_0,3}(\mathbf{a}, \mathbf{a}') = \{\mathbf{d}_{\mathcal{A},\mathbf{u},\infty}(\mathbf{a}, \mathbf{a}') + R_0 \mathbf{d}_{\mathcal{A},\mathbf{x},\infty}(\mathbf{a}, \mathbf{a}') \leq C_{\text{stab},3}(R_0)\}.$$

2864 Again, we see that aside from the $\mathbf{I}_{0,\infty}\{\cdot\}$ terms, our distances $\mathbf{d}_{\mathcal{A},R_0,\tau,\mathbf{x}}(\mathbf{a}, \mathbf{a}' | r)$ depends only on ℓ_2 -guarantees.

2865 We may now state our most general stability guarantee.

2866 **Proposition G.3** (Main Stability Guarantees). *Suppose that*

$$2867 \tau_c \geq 3L_{\text{stab}} \log(2c_{\Delta})/\eta.$$

2868 *In addition, fix an $R_0 > 0$, and r_{\max} such that $r_{\max} \leq C_{\text{stab},4}(R_0)$. Consider any sequence of kernels $\{W_h\}_{h=1}^H$, where $W_h : \mathcal{S} \rightarrow \Delta(\mathcal{S})$ and⁸*

$$2869 \max_{h,\mathbf{s} \in \mathcal{S}} \mathbb{P}_{\tilde{\mathbf{s}} \sim W_h(\mathbf{s})}[\mathbf{d}_{\text{IPS}}(\tilde{\mathbf{s}}, \mathbf{s}) \leq r] = 1, \quad \phi_{\mathcal{Z}} \circ W_h(\mathbf{s}) \ll \phi_{\mathcal{Z}} \circ \mathbf{P}_h^*, \quad (\text{G.4})$$

2870 *and define a process $\mathbf{s}_1 \sim \mathbf{P}_{\text{init}}$, $\tilde{\mathbf{s}}_h \sim W_h(\mathbf{s}_h)$, $\mathbf{a}_h \sim \pi_h^*(\tilde{\mathbf{s}}_h)$, and $\mathbf{s}_{h+1} := F_h(\mathbf{s}_h, \mathbf{a}_h)$. Then, almost surely, the following hold for all $0 \leq \tau \leq \tau_c$:*

- 2871 • For each $1 \leq h \leq H$, $\mathbf{d}_{\mathcal{S},\tau}(F_h(\tilde{\mathbf{s}}_h, \mathbf{a}_h), \mathbf{s}_h) \leq 2(1 + R_{\text{stab}})B_{\text{stab}}r\beta_{\text{stab}}^{(\tau_c - \tau)}$.
- 2872 • For any sequence $(\mathbf{a}'_{1:H})$, the dynamics $\mathbf{s}'_1 = \mathbf{s}_1$, $\mathbf{s}_{h+1} = F_h(\mathbf{s}'_h, \mathbf{a}'_h)$ satisfy

$$2873 \max_{1 \leq h \leq H+1} \mathbf{d}_{\mathcal{S},\tau}(\mathbf{s}_h, \mathbf{s}'_h) \leq \max_{1 \leq h \leq H} \mathbf{d}_{\mathcal{A},R_0,\tau,\mathbf{x}}(\mathbf{a}_h, \mathbf{a}'_h | r).$$

2874 and

$$2875 \max_{1 \leq h \leq H+1} \mathbf{d}_{\mathcal{S},\tau}(\mathbf{s}_h, \mathbf{s}'_h) \leq \max_{1 \leq h \leq H} \mathbf{d}_{\mathcal{A},R_0,\tau}(\mathbf{a}_h, \mathbf{a}'_h | r).$$

2876 The above proposition is proven in [Appendix G.5](#), where it is derived from two key guarantees given in [Appendix G.4](#) below.

2877 **Remark G.2** (Remark on the Scaling). We now justify the extreme granularity of the above result. We demonstrate that our guarantees satisfy the following favorable properties:

- 2878 • As in [Corollary G.1](#), the dependence of $\bar{\mathbf{K}}_k - \bar{\mathbf{K}}'_k$ in the non $\mathbf{I}_{0,\infty}\{\cdot\}$ scales down with r and with $r \cdot \beta_{\text{stab}}^{(\tau_c - \tau)/3}$, so that errors in $\bar{\mathbf{K}}_k$ become less relevant as $\tau \rightarrow \tau_c$ and as $r \rightarrow 0$.
- 2879 • If we restrict our attention only to errors in states, captured by $\mathbf{d}_{\mathcal{S},\tau}$, the non- $\mathbf{I}_{0,\infty}\{\cdot\}$ terms depend only on ℓ_2 -errors rather than maximal ∞ -norm ones.
- 2880 • In the special case where $R_{\text{dyn}} = \infty$, i.e., the regularity properties in [Assumption 3.2](#) hold globally, then all terms $C_{\text{stab},i}(R_0)$ defined in [Definition G.5](#) no longer need depend on R_0 , as the terms in which R_0 appears have an $R_{\text{dyn}} = \infty$ in the numerator, and each $C_{\text{stab},i}(R_0)$ serves as an upper bound on a certain quantity of interest. Hence, we can drop the dependence on R_0 in all of these terms. $C_{\text{stab},i}(R_0)$ (
- 2881 • In particular, the term $C_{\text{stab},3}(R_0)$ if equal to ∞ when $R_{\text{dyn}} = \infty$. Thus, for $R_{\text{dyn}} = \infty$, we can drop the indicator of $\mathcal{E}_{\text{close},R_0,3}(\mathbf{a}, \mathbf{a}') := \{\mathbf{d}_{\mathcal{A},\mathbf{u},\infty}(\mathbf{a}, \mathbf{a}') + R_0 \mathbf{d}_{\mathcal{A},\mathbf{x},\infty}(\mathbf{a}, \mathbf{a}') \leq C_{\text{stab},3}(R_0)\}$, and hence each $\mathbf{d}_{\mathcal{S},\tau}$ has depends only on ℓ_2 -type errors.

2882 The proof of [Proposition G.3](#) is given in [Appendix G.5](#), derived from the results in the subsection directly below. Before we do this, we first demonstrate how [Corollary G.1](#) follows from [Proposition G.3](#).

2883 ⁸Again, we refer to [Remark G.1](#) for explanation of the second condition in the display [\(G.4\)](#)

G.3.2. DERIVING COROLLARY G.1 FROM PROPOSITION G.3

The proof is mostly notational bookkeeping.

By assumption $\phi_{\mathcal{Z}} \circ W_h(s) \ll \phi_{\mathcal{Z}} \circ P_h^*$ and Lemma G.2, and the R_{stab} -term in $(R_{\text{stab}}, B_{\text{stab}}, L_{\text{stab}})$ -Jacobian stability, the action \mathbf{a}_h with $\text{rad}_{\mathbf{K}}(\mathbf{a}) \leq R_{\text{stab}}$. Further, notice that the parameter β_{stab} used throughout this section can be bounded by at most

$$\beta_{\text{stab}} := \left(1 - \frac{\eta}{L_{\text{stab}}}\right) \leq \exp(-\eta/L_{\text{stab}}).$$

Hence, Corollary G.1 from Proposition G.3 as soon as we show that

$$\forall \mathbf{a} \text{ s.t. } \text{rad}_{\mathbf{K}}(\mathbf{a}) \leq R_{\text{stab}}, \quad d_{\mathcal{A}, R_0, \tau}(\mathbf{a}_h, \mathbf{a}'_h \mid r) \leq \bar{d}_{\mathcal{A}, \tau}(\mathbf{a}, \mathbf{a}' \mid r).$$

Consider the action divergences in Definition G.10. Take $R_0 = 2R_{\text{stab}}$, where $R_{\text{stab}} \geq 1$ by assumption. and upper bound $d_{\mathcal{A}, \mathbf{u}, \ell_2}(\cdot, \cdot) \leq \sqrt{L_{\text{stab}}} d_{\mathcal{A}, \mathbf{u}, \infty}(\cdot)$ (as in (G.3)), and similarly for $d_{\mathcal{A}, \mathbf{x}, \ell_2}(\cdot, \cdot)$ and $d_{\mathcal{A}, \mathbf{K}, \ell_2}(\cdot, \cdot)$. Then,

$$\begin{aligned} d_{\mathcal{A}, R_0, \tau}(\mathbf{a}_h, \mathbf{a}'_h \mid r) &:= 2((1 + R_0)d_{\mathcal{A}, R_0, \tau, \mathbf{x}}(\mathbf{a}_h, \mathbf{a}'_h \mid r) + d_{\mathcal{A}, R_0, \tau, \mathbf{u}}(\mathbf{a}_h, \mathbf{a}'_h \mid r)) \\ &= d_{\mathcal{A}, \mathbf{u}, \infty}(\mathbf{a}_h, \mathbf{a}'_h) + R_0 d_{\mathcal{A}, \mathbf{x}, \infty}(\mathbf{a}_h, \mathbf{a}'_h) + 2r B_{\text{stab}} \beta_{\text{stab}}^{\tau_c - \tau} d_{\mathcal{A}, \mathbf{K}, \infty}(\mathbf{a}_h, \mathbf{a}'_h) \\ &\quad + 2c_{\mathbf{u}}(d_{\mathcal{A}, \mathbf{u}, \ell_2}(\mathbf{a}, \mathbf{a}') + R_0 d_{\mathcal{A}, \mathbf{x}, \ell_2}(\mathbf{a}, \mathbf{a}')) + 2c_{\mathbf{K}} r (\beta_{\text{stab}})^{\frac{\tau_c - \tau}{3}} \cdot d_{\mathcal{A}, \mathbf{K}, \ell_2}(\mathbf{a}, \mathbf{a}') \\ &\quad + \mathbf{I}_{0, \infty} \left\{ \bigcap_{i=1}^3 \mathcal{E}_{\text{close}, R_0, i} \right\} + \mathbf{I}_{0, \infty} \{ \text{rad}_{\mathbf{K}}(\mathbf{a}) \vee \text{rad}_{\mathbf{K}}(\mathbf{a}') \leq R_0 \} \\ &\leq 2R_{\text{stab}}(1 + 2c_{\mathbf{u}}\sqrt{L_{\text{stab}}})(d_{\mathcal{A}, \mathbf{u}, \infty}(\mathbf{a}_h, \mathbf{a}'_h) + d_{\mathcal{A}, \mathbf{x}, \infty}(\mathbf{a}_h, \mathbf{a}'_h)) \\ &\quad + (2B_{\text{stab}} + 2\sqrt{L_{\text{stab}}}c_{\mathbf{K}})r \exp^{-\frac{\eta(\tau_c - \tau)}{3L_{\text{stab}}}} d_{\mathcal{A}, \mathbf{K}, \infty}(\mathbf{a}_h, \mathbf{a}'_h) \\ &\quad + \mathbf{I}_{0, \infty} \left\{ \bigcap_{i=1}^3 \mathcal{E}_{\text{close}, R_0, i} \right\} + \mathbf{I}_{0, \infty} \{ \text{rad}_{\mathbf{K}}(\mathbf{a}) \vee \text{rad}_{\mathbf{K}}(\mathbf{a}') \leq 2R_{\text{stab}} \} \\ &\leq \frac{c_1}{3}(d_{\mathcal{A}, \mathbf{u}, \infty}(\mathbf{a}_h, \mathbf{a}'_h) + d_{\mathcal{A}, \mathbf{x}, \infty}(\mathbf{a}_h, \mathbf{a}'_h)) + r \exp^{-\frac{\eta(\tau_c - \tau)}{3L_{\text{stab}}}} d_{\mathcal{A}, \mathbf{K}, \infty}(\mathbf{a}_h, \mathbf{a}'_h) \\ &\quad + \mathbf{I}_{0, \infty} \left\{ \bigcap_{i=1}^3 \mathcal{E}_{\text{close}, R_0, i} \right\} + \mathbf{I}_{0, \infty} \{ \text{rad}_{\mathbf{K}}(\mathbf{a}) \vee \text{rad}_{\mathbf{K}}(\mathbf{a}') \leq 2R_{\text{stab}} \}, \end{aligned}$$

where we recall from (G.2)

$$c_1 := 6 \max\{R_{\text{stab}}(1 + 2c_{\mathbf{u}}\sqrt{L_{\text{stab}}}), B_{\text{stab}} + \sqrt{L_{\text{stab}}}c_{\mathbf{K}}\}.$$

Let's now simplify the indicators. Restricting our attention to \mathbf{a} with $\text{rad}_{\mathbf{K}}(\mathbf{a}) \leq R_{\text{stab}}$, $\text{rad}_{\mathbf{K}}(\mathbf{a}') \leq R_{\text{stab}} + d_{\mathcal{A}, \mathbf{K}, \infty}(\mathbf{a}_h, \mathbf{a}'_h)$ by the triangle inequality. Thus, we can replace $\mathbf{I}_{0, \infty} \{ \text{rad}_{\mathbf{K}}(\mathbf{a}) \vee \text{rad}_{\mathbf{K}}(\mathbf{a}') \leq 2R_{\text{stab}} \}$ with $\mathbf{I}_{0, \infty} \{ d_{\mathcal{A}, \mathbf{K}, \infty}(\mathbf{a}_h, \mathbf{a}'_h) \leq R_{\text{stab}} \}$.

We now recall the definitions

$$\begin{aligned} \mathcal{E}_{\text{close}, R_0, 1}(\mathbf{a}, \mathbf{a}') &= \{(d_{\mathcal{A}, \mathbf{u}, \ell_2}(\mathbf{a}, \mathbf{a}') + R_0 d_{\mathcal{A}, \mathbf{x}, \ell_2}(\mathbf{a}, \mathbf{a}') \leq C_{\text{stab}, 1}(R_0)\} \\ \mathcal{E}_{\text{close}, R_0, 2}(\mathbf{a}, \mathbf{a}') &= \{d_{\mathcal{A}, \mathbf{K}, \ell_2}(\mathbf{a}, \mathbf{a}') \leq C_{\text{stab}, 2}(R_0)\} \\ \mathcal{E}_{\text{close}, R_0, 3}(\mathbf{a}, \mathbf{a}') &= \{d_{\mathcal{A}, \mathbf{u}, \infty}(\mathbf{a}, \mathbf{a}') + R_0 d_{\mathcal{A}, \mathbf{x}, \infty}(\mathbf{a}, \mathbf{a}') \leq C_{\text{stab}, 3}(R_0)\}. \end{aligned}$$

Again, recall that we take $R_0 = 2R_{\text{stab}}$. Again, invoke the upper bounds of the form $d_{\mathcal{A}, \mathbf{u}, \ell_2}(\cdot, \cdot) \leq \sqrt{L_{\text{stab}}} d_{\mathcal{A}, \mathbf{u}, \infty}(\cdot)$ (as in (G.3)). Thus, $\bigcap_{i=1}^3 \mathcal{E}_{\text{close}, R_0, i} \cap \{ \text{rad}_{\mathbf{K}}(\mathbf{a}) \vee \text{rad}_{\mathbf{K}}(\mathbf{a}') \leq 2R_{\text{stab}} \}$ holds as soon as

$$\max\{d_{\mathcal{A}, \mathbf{u}, \infty}(\mathbf{a}, \mathbf{a}'), d_{\mathcal{A}, \mathbf{x}, \infty}(\mathbf{a}, \mathbf{a}')d_{\mathcal{A}, \mathbf{K}, \infty}(\mathbf{a}, \mathbf{a}')\} \leq c_2,$$

where we recall from

$$c_2 := \min \left\{ \frac{C_{\text{stab}, 1}(2R_{\text{stab}})}{4R_{\text{stab}}\sqrt{L_{\text{stab}}}}, \frac{C_{\text{stab}, 2}(2R_{\text{stab}})}{\sqrt{L_{\text{stab}}}}, \frac{C_{\text{stab}, 3}(2R_{\text{stab}})}{4R_{\text{stab}}}, \frac{1}{2R_{\text{stab}}} \right\}.$$

2970 In sum, for any \mathbf{a} for which $\text{rad}_{\mathbf{K}}(\mathbf{a}) \leq R_{\text{stab}}$, we have

$$2971 \quad d_{\mathcal{A}, R_0, \tau}(\mathbf{a}_h, \mathbf{a}'_h \mid r) \leq \frac{c_1}{3} (d_{\mathcal{A}, \mathbf{u}, \infty}(\mathbf{a}_h, \mathbf{a}'_h) + d_{\mathcal{A}, \mathbf{x}, \infty}(\mathbf{a}_h, \mathbf{a}'_h) + r \exp^{-\frac{\eta(\tau_c - \tau)}{3L_{\text{stab}}}} d_{\mathcal{A}, \mathbf{K}, \infty}(\mathbf{a}_h, \mathbf{a}'_h))$$

$$2972 \quad + \mathbf{I}_{0, \infty} \{ \max \{ d_{\mathcal{A}, \mathbf{u}, \infty}(\mathbf{a}, \mathbf{a}'), d_{\mathcal{A}, \mathbf{x}, \infty}(\mathbf{a}, \mathbf{a}'), d_{\mathcal{A}, \mathbf{K}, \infty}(\mathbf{a}, \mathbf{a}') \} \leq c_2 \}.$$

2975 To conclude, we observe that, for any nonnegative coefficients $a_1, a_2, a_3 > 0$ and sequences $v_{1,1:n}, v_{2,1:n}, v_{3,n} \geq 0$ in \mathbb{R}^n ,

$$2976 \quad \sum_{i=1}^3 a_i (\max_{j \in [n]} v_{i,j}) \leq 3 \max_{j \in [n]} \sum_{i=1}^3 a_i v_{i,j}.$$

2977 Thus, if we express $\mathbf{a} = (\bar{\mathbf{u}}_{1:\tau_c}, \bar{\mathbf{x}}_{1:\tau_c}, \bar{\mathbf{K}}_{1:\tau_c})$ and $\mathbf{a}' = (\bar{\mathbf{u}}'_{1:\tau_c}, \bar{\mathbf{x}}'_{1:\tau_c}, \bar{\mathbf{K}}'_{1:\tau_c})$, we can bound

$$2978 \quad d_{\mathcal{A}, R_0, \tau}(\mathbf{a}_h, \mathbf{a}'_h \mid r) \leq \bar{d}_{\mathcal{A}, \tau}(\mathbf{a}, \mathbf{a}' \mid r)$$

$$2979 \quad := c_1 \max_{1 \leq k \leq \tau_c} \left(\|\bar{\mathbf{u}}_k - \bar{\mathbf{u}}'_k\| + \|\bar{\mathbf{x}}_k - \bar{\mathbf{x}}'_k\| + r e^{-\frac{\eta(\tau_c - \tau)}{3L_{\text{stab}}}} \|\bar{\mathbf{K}}_k - \bar{\mathbf{K}}'_k\| \right)$$

$$2980 \quad + \mathbf{I}_{0, \infty} \left\{ \max_{1 \leq k \leq \tau_c} (\max \{ \|\bar{\mathbf{u}}_k - \bar{\mathbf{u}}'_k\|, \|\bar{\mathbf{x}}_k - \bar{\mathbf{x}}'_k\|, \|\bar{\mathbf{K}}_k - \bar{\mathbf{K}}'_k\| \}) \leq c_2 \right\}.$$

2981 \square

2982 G.4. Stability guarantees for single control (sub-)trajectories.

2983 At the heart of the IPS guarantees in [Appendix G.3](#) above are two building blocks: one controller the perturbation of initial state around a regular (in the sense of [Assumption 3.1](#)) trajectory, and the second extending this guarantee to perturbations of control inputs and gains.

2984 **Lemma G.4** (Stability to State Perturbation). *Let $\bar{\rho} = (\bar{\mathbf{x}}_{1:K+1}, \bar{\mathbf{u}}_{1:K}) \in \mathcal{P}_K$ be an $(R_{\text{dyn}}, L_{\text{dyn}}, M_{\text{dyn}})$ -regular and feasible path, and let $\mathbf{K}_{1:K}$ be gains such that $(\bar{\rho}, \mathbf{K}_{1:K})$ is $(R_{\text{stab}}, B_{\text{stab}}, L_{\text{stab}})$ -stable. Assume, that $R_{\text{stab}} \geq 1$, $L_{\text{stab}} \geq 2\eta$. Fix another \mathbf{x}_1 and define another trajectory ρ via*

$$2985 \quad \mathbf{u}_k = \bar{\mathbf{u}}_k + \mathbf{K}_k(\mathbf{x}_k - \bar{\mathbf{x}}_k), \quad \mathbf{x}_{k+1} = \bar{\mathbf{x}}_k + \eta f_{\eta}(\mathbf{x}_k, \mathbf{u}_k)$$

2986 Then, if $\|\mathbf{x}_1 - \bar{\mathbf{x}}_1\| \leq \min\{(16L_{\text{stab}}M_{\text{dyn}}R_{\text{stab}}^2B_{\text{stab}}^2)^{-1}, \frac{R_{\text{dyn}}}{2R_{\text{stab}}B_{\text{stab}}}\}$, then

- 2987 • $\|\mathbf{x}_{k+1} - \bar{\mathbf{x}}_{k+1}\| \leq 2B_{\text{stab}}\|\mathbf{x}_1 - \bar{\mathbf{x}}_1\|\beta_{\text{stab}}^k$.
- 2988 • $(\rho, \mathbf{K}_{1:K})$ is $(R_{\text{stab}}, 2B_{\text{stab}}, L_{\text{stab}})$ -stable.
- 2989 • $\|\mathbf{B}_k(\rho)\| \leq L_{\text{dyn}}$.

2990 This lemma is proven in [Appendix G.6](#), and the following proposition in [Appendix G.7](#).

2991 **Proposition G.5** (Single Trajectory Stability Guarantee). *Let $\bar{\rho} = (\bar{\mathbf{x}}_{1:K+1}, \bar{\mathbf{u}}_{1:K}) \in \mathcal{P}_K$ be $(R_{\text{dyn}}, L_{\text{dyn}}, M_{\text{dyn}})$ -regular and feasible, and let $\mathbf{K}_{1:K}$ be such that $(\bar{\rho}, \mathbf{K}_{1:K})$ is $(R_{\text{stab}}, B_{\text{stab}}, L_{\text{stab}})$ -stable. Assume $R_{\text{stab}} \geq 1$, $L_{\text{stab}} \geq 2\eta$, and given another $\mathbf{x}_1, \mathbf{x}'_1 \in \mathcal{X}$, $\bar{\mathbf{u}}'_{1:K}$ and $\mathbf{K}'_{1:K}$, define trajectories $\rho = (\mathbf{x}_{1:K+1}, \mathbf{u}_{1:K})$ and $\rho' = (\mathbf{x}'_{1:K+1}, \mathbf{u}'_{1:K})$*

$$2992 \quad \mathbf{x}_{k+1} = \mathbf{x}_k + \eta f_{\eta}(\mathbf{x}_k, \mathbf{u}_k), \quad \mathbf{u}_k = \bar{\mathbf{u}}_k + \mathbf{K}_k(\mathbf{x}_k - \bar{\mathbf{x}}_k)$$

$$2993 \quad \mathbf{x}'_{k+1} = \mathbf{x}'_k + \eta f_{\eta}(\mathbf{x}'_k, \mathbf{u}'_k), \quad \mathbf{u}'_k = \bar{\mathbf{u}}'_k + \mathbf{K}'_k(\mathbf{x}'_k - \bar{\mathbf{x}}_k)$$

2994 Let all constants be as defined in [Definition G.4](#), and define (recalling the stability exponent $\beta_{\text{stab}} := (1 - \frac{\eta}{L_{\text{stab}}})$) the terms

$$2995 \quad \text{Err}_{\mathbf{u}} := \max_{k \in [K]} \left(\eta \sum_{j=1}^k \beta_{\text{stab}}^{k-j} \|\bar{\mathbf{u}}_j - \bar{\mathbf{u}}'_j\|^2 \right)^{1/2}, \quad \text{Err}_{\mathbf{K}} := \max_{k \in [K]} \left(\eta \sum_{j=1}^k \beta_{\text{stab}}^{k-j} \|\mathbf{K}_j - \mathbf{K}'_j\|^2 \right)^{1/2}$$

2996 Then, the conclusions of [Lemma G.4](#) applies to the trajectory ρ , and moreover, for all $1 \leq k \leq K$,

$$2997 \quad \|\mathbf{x}_{k+1} - \mathbf{x}'_{k+1}\| \leq c_{\mathbf{u}} \text{Err}_{\mathbf{u}} + (c_{\mathbf{K}} \text{Err}_{\mathbf{K}} \|\mathbf{x}_1 - \bar{\mathbf{x}}_1\| + c_{\Delta} \|\mathbf{x}_1 - \mathbf{x}'_1\|) \beta_{\text{stab}}^{k/3},$$

2998 provided that the following two conditions hold:

2999

- The above error terms satisfy

$$\text{Err}_{\mathbf{u}} \leq C_{\mathbf{u}}, \quad \text{Err}_{\mathbf{K}} \leq C_{\mathbf{K}}, \quad \|\mathbf{x}_1 - \mathbf{x}'_1\| \leq C_{\Delta}, \quad \|\mathbf{x}_1 - \bar{\mathbf{x}}_1\| \leq C_{\bar{\mathbf{x}}}, \quad \text{Err}_{\mathbf{K}} \|\mathbf{x}_1 - \bar{\mathbf{x}}_1\| \leq C_{\mathbf{K}, \bar{\mathbf{x}}}$$

- In addition, if $R_{\text{dyn}} < \infty$, $\tilde{R}_{\text{stab}} := \max\{1, R_{\text{stab}}, \max_{1 \leq j \leq K} \|\mathbf{K}'_j\|\}$ and $\Delta_{\mathbf{u}, \infty} := \max_j \|\bar{\mathbf{u}}_j - \bar{\mathbf{u}}'_j\|$ satisfy

$$R_{\text{dyn}} \geq (4\tilde{R}_{\text{stab}}c_{\mathbf{u}}\sqrt{L_{\text{stab}}} + 1)\Delta_{\mathbf{u}, \infty} + 4\tilde{R}_{\text{stab}}c_{\mathbf{K}}\|\mathbf{x}_1 - \bar{\mathbf{x}}_1\| + 4\tilde{R}_{\text{stab}}c_{\Delta}\|\mathbf{x}_1 - \mathbf{x}'_1\|.$$

The proofs of both this proposition and the lemma before it consist of translating the differences in trajectories into recursions satisfying certain functional forms. Taking norms, we obtain scalar recursions whose solutions are upper bounded in a series of technical lemmas detailed in [Appendix G.9](#). We believe these [Proposition G.5](#) and [Lemma G.4](#) are useful more broadly in the study of perturbation of non-linear control systems.

Notice that, for convenience, both the \mathbf{x} and \mathbf{x}' trajectories are stabilizing around the same $\bar{\mathbf{x}}$. This is for convenience, and simplifies the analysis. Indeed, difference generalizing to accomodate \mathbf{x}' stabilizing around $\bar{\mathbf{x}}'$ can be accomplished by a change of variables in the $\bar{\mathbf{u}}'$, which is precisely what is done in deriving [Proposition G.3](#) in the section that follows.

G.5. Deriving [Proposition G.3](#) from [Proposition G.5](#)

The majority of this proof is (also) notational bookkeeping, whereby we convert two trajectories (in the abstract states/actions notation) into separate trajectories for each a sequence of $h = 1, 2, \dots, H = T/\tau_c$, to each of which we apply [Proposition G.5](#).

Constructing the (perturbed) expert trajectory We begin by unfolding the generative process for abstract-states s_1, \dots, s_H in our proposition. Recall further that $s_h = \rho_{c,h}$ corresponds to the trajectory-chunk.

We let the (control) states and inputs for the corresponding sequence be denote as $(\mathbf{x}_{1:T+1}, \mathbf{u}_{1:T})$ be generated as follows. Start with

$$\mathbf{x}_1 \leftarrow s_1$$

drawn from the initial state distribution. Assume that we have constructed the states s_1, \dots, s_{h-1} ; this means in particular that we have constructed $\mathbf{x}_{1:t_h}, \mathbf{u}_{1:t_h-1}$, as well as the memory-chunks $\rho_{m,1}, \dots, \rho_{m,j-1}$. We extend the construction to step $h+1$ as follows:

- Define $\mathbf{x}_{h,1} = \mathbf{x}_{t_h}$
- Select a perturbation of the state $\tilde{s}_h = \tilde{\rho}_{c,h} = (\tilde{\mathbf{x}}_{t_{h-1}:t_h}, \tilde{\mathbf{u}}_{t_{h-1}:t_h-1})$, with corresponding memory-chunk $\tilde{\rho}_{m,h} = (\tilde{\mathbf{x}}_{t_h-\tau_m+1:t_h}, \tilde{\mathbf{u}}_{t_h-\tau_m+1:t_h-1})$. As per the proposition, $d_{\text{ips}}(s_h, \tilde{s}_h) \leq r$. This means that $\|\mathbf{x}_{t_h} - \tilde{\mathbf{x}}_{t_h}\| \leq r$.
- Draw $\mathbf{a}_h = \kappa_{t_h:t_h:t_h+\tau_m-1} \sim \pi_h^*(\tilde{\rho}_{m,h})$. We express

$$\kappa_t(\mathbf{x}) = \bar{\mathbf{u}}_t + \bar{\mathbf{K}}_t(\mathbf{x} - \bar{\mathbf{x}}_t), \quad t_h \leq t \leq t_{h+1} - 1,$$

and reindexed trajectory

$$\kappa_{h,k} = \kappa_{t_h+k-1}.$$

Denote

$$\bar{\mathbf{x}}_{h,k} = \bar{\mathbf{x}}_{t_h+k-1}, \quad \bar{\mathbf{u}}_{h,k} = \bar{\mathbf{u}}_{t_h+k-1}, \quad \bar{\mathbf{K}}_{h,k} = \bar{\mathbf{K}}_{t_h+k-1}$$

and

$$\bar{\rho}_{[h+1]} = (\bar{\mathbf{x}}_{h,1:\tau_c+1}, \bar{\mathbf{u}}_{h,1:\tau_c}).$$

- Moreover, because we assume $\phi_{\mathcal{Z}} \circ W_h \ll \phi_{\mathcal{Z}} \circ P_h^*$, we inherit the conclusions of [Lemma G.2](#). Hence, $\bar{\rho}_{h+1}$ such be feasible, $(R_{\text{dyn}}, L_{\text{dyn}}, M_{\text{dyn}})$ -regular, and $(\bar{\rho}_h, \kappa_{h,1:\tau_c})$ is $(R_{\text{stab}}, B_{\text{stab}}, L_{\text{stab}})$ -stable. In addition, [Lemma G.2](#) ensures $\bar{\mathbf{x}}_{h,1} = \tilde{\mathbf{x}}_{t_h}$. Consequently, we have that the composite action map F_h satisfies

$$F_h(\tilde{s}_h, \mathbf{a}_h) = \bar{\rho}_{[h+1]} = (\bar{\mathbf{x}}_{h,1:\tau_c+1}, \bar{\mathbf{u}}_{h,1:\tau_c}). \quad (\text{G.5})$$

- We execute \mathbf{a}_h for τ_c steps from our *actual* state \mathbf{x}_t (not $\tilde{\mathbf{x}}_t$), giving states and actions

$$\mathbf{x}_{t+1} = f(\mathbf{x}_t, \mathbf{u}_t), \quad \mathbf{u}_t = \kappa_t(\mathbf{x}_t), \quad 1 \leq t \leq \tau_c.$$

And define

$$\mathbf{x}_{h,k+1} = \mathbf{x}_{t_h+k}, \quad \mathbf{u}_{h,k} = \mathbf{u}_{t_h+k-1}, \quad 1 \leq k \leq \tau_c.$$

- Finally, define the chunks the trajectories $\boldsymbol{\rho}_{[h+1]} = (\mathbf{x}_{h,1:\tau_c+1}, \mathbf{x}_{h,1:\tau_c})$, which is equal to the next abstract-state

$$\mathbf{s}_{h+1} = (\mathbf{x}_{h,1:\tau_c+1}, \mathbf{u}_{h,1:\tau_c}) = (\mathbf{x}_{t_h:t_{h+1}}, \mathbf{u}_{t_h:t_{h+1}-1}) \quad (\text{G.6})$$

Construction of the imitation trajectory. We now construct the imitation trajectory by setting $\mathbf{x}'_1 = \mathbf{x}_1$, and

- For each h , select $\mathbf{a}'_h = (\kappa'_{t_h:t_h+\tau_c-1}) \in \mathcal{K}^{\tau_c}$. Define the re-indexed primitive controllers

$$\kappa'_{h,k} = \kappa'_{t_h+k-1},$$

and express

$$\kappa'_{h,k}(\mathbf{x}) = \bar{\mathbf{K}}'_{h,k}(\bar{\mathbf{x}} - \bar{\mathbf{x}}_{k,h}) + \bar{\mathbf{u}}'_{h,k}.$$

- Execute \mathbf{a}'_h for τ_c steps, giving states and actions

$$\mathbf{x}'_{t+1} = f(\mathbf{x}'_t, \mathbf{u}'_t), \quad \mathbf{x}'_t = \kappa'_t(\mathbf{x}'_t), \quad 1 \leq t \leq \tau_c.$$

And define

$$\mathbf{x}'_{h,k+1} = \mathbf{x}'_{t_h+k}, \quad \mathbf{u}'_{h,k} = \mathbf{u}'_{t_h+k-1}, \quad 1 \leq k \leq \tau_c.$$

- Finally, define the chunks

$$\mathbf{s}'_h = (\mathbf{x}'_{t_h:t_{h+1}}, \mathbf{u}'_{t_h:t_{h+1}-1}) = (\mathbf{x}'_{h,1:\tau_c+1}, \mathbf{u}'_{h,1:\tau_c}).$$

Further Notation. Let's define the following errors analogous to [Proposition G.5](#).

$$\text{Err}_{\bar{\mathbf{u}},h}^2 = \max_{1 \leq k \leq \tau_c} \eta \sum_{j=1}^k \beta_{\text{stab}}^{k-j} \|\bar{\mathbf{u}}_{h,j} - \kappa'_{h,j}(\bar{\mathbf{x}}_{h,j})\|^2$$

$$\text{Err}_{\bar{\mathbf{K}},h}^2 = \max_{1 \leq k \leq \tau_c} \eta \sum_{j=1}^k \beta_{\text{stab}}^{k-j} \|\bar{\mathbf{K}}_{h,j} - \bar{\mathbf{K}}'_{h,j}\|^2$$

$$\Delta_{\bar{\mathbf{u}},\infty,h} := \max_k \|\bar{\mathbf{u}}_{h,k} - \kappa'_{h,j}(\bar{\mathbf{x}}_{h,j})\|.$$

Importantly, in [Proposition G.5](#), it is assumed that other the primed and unprimed sequence stabilize to the same $\mathbf{x}_{h,k}$, whereas here, the primed sequence stabilized to $\mathbf{x}_{h,k'}$. This is addressed by replacing the role of $\mathbf{u}'_{h,k}$ with $\kappa'_{h,k}(\bar{\mathbf{x}}_{h,k})$.

G.5.1. INTERPRETING THE ERROR TERMS.

First, we unpack the above error terms.

Lemma G.6. *Suppose $\max_h d_{\mathcal{A},R_0,\tau,\mathbf{x}}(\mathbf{a}_h, \mathbf{a}'_h \mid r)$ is finite. Then,*

$$\text{Err}_{\bar{\mathbf{K}},h} = d_{\mathcal{A},\mathbf{K},\ell_2}(\mathbf{a}_h, \mathbf{a}'_h)$$

$$\text{Err}_{\bar{\mathbf{u}},h} \leq d_{\mathcal{A},\mathbf{u},\ell_2}(\mathbf{a}_h, \mathbf{a}'_h) + R_0 d_{\mathcal{A},\mathbf{x},\ell_2}(\mathbf{a}_h, \mathbf{a}'_h)$$

$$\Delta_{\bar{\mathbf{u}},\infty,h} = d_{\mathcal{A},\mathbf{u},\infty}(\mathbf{a}_h, \mathbf{a}'_h) + R_0 d_{\mathcal{A},\mathbf{x},\infty}(\mathbf{a}_h, \mathbf{a}'_h)$$

3135 *Proof of Lemma G.6.* The equality of $\text{Err}_{\bar{\mathbf{K}},h}$ follows from the reindexing $\bar{\mathbf{K}}_{h,k} = \bar{\mathbf{K}}_{t_h+k-1}$ and the definition of [Definition G.9](#). Next, unpacking our notation of $\bar{\mathbf{u}}, \bar{\mathbf{u}}'$, we compute

$$\begin{aligned} 3137 \quad \bar{\mathbf{u}}_{h,j} - \kappa_{h,j}(\bar{\mathbf{x}}_{h,j}) &= \bar{\mathbf{u}}_{h,j} - \bar{\mathbf{K}}'_{h,j}(\bar{\mathbf{x}}'_{h,j} - \bar{\mathbf{x}}_{h,j}) \\ 3138 \quad &= \mathbf{u}'_{t_h+k-1} - \mathbf{u}'_{t_h+k-1} - \bar{\mathbf{K}}'_{t_h+k-1}(\mathbf{x}_{t_h+k-1} - \mathbf{x}'_{t_h+k-1}) \end{aligned}$$

3139 So that as long as $d_{\mathcal{A},R_0,\tau,\mathbf{x}}(\mathbf{a}_h, \mathbf{a}'_h \mid r)$ for all h , then $\|\bar{\mathbf{K}}_{h,j}\| \leq R_0$. Thus

$$3140 \quad \|\bar{\mathbf{u}}_{h,j} - \bar{\mathbf{u}}'_{h,j}\| \leq \|\mathbf{u}'_{t_h+k-1} - \mathbf{u}'_{t_h+k-1}\| + R_0 \|\mathbf{x}_{t_h+k-1} - \mathbf{x}'_{t_h+k-1}\|$$

3141 and thus by the triangle and moving the max outside the sum,

$$\begin{aligned} 3142 \quad \text{Err}_{\bar{\mathbf{u}},h} &= \max_{1 \leq k \leq \tau_c} \left(\eta \sum_{j=1}^k \beta_{\text{stab}}^{k-j} \|\bar{\mathbf{u}}_{h,j} - \bar{\mathbf{u}}'_{h,j}\|^2 \right)^{1/2} \\ 3143 \quad &\leq \max_{1 \leq k \leq \tau_c} \left(\eta \sum_{j=1}^k \beta_{\text{stab}}^{k-j} \|\mathbf{u}'_{t_h+k-1} - \mathbf{u}'_{t_h+k-1}\|^2 \right)^{1/2} \\ 3144 \quad &\quad + R_0 \max_{1 \leq k \leq \tau_c} \left(\eta \sum_{j=1}^k \beta_{\text{stab}}^{k-j} \|\mathbf{x}'_{t_h+k-1} - \mathbf{x}'_{t_h+k-1}\|^2 \right)^{1/2} \\ 3145 \quad &\leq d_{\mathcal{A},\mathbf{u},\ell_2}(\mathbf{a}_h, \mathbf{a}'_h) + R_0 d_{\mathcal{A},\mathbf{x},\ell_2}(\mathbf{a}_h, \mathbf{a}'_h). \end{aligned}$$

3146 The inequality $\Delta_{\bar{\mathbf{u}},\infty,h} \leq d_{\mathcal{A},\mathbf{u},\infty}(\mathbf{a}_h, \mathbf{a}'_h) + R_0 d_{\mathcal{A},\mathbf{x},\infty}(\mathbf{a}_h, \mathbf{a}'_h)$ follows similarly. □

3147 G.5.2. AN INTERMEDIATE GUARANTEE.

3148 Next, we establish an intermediate guarantee, from which [Proposition G.3](#) is readily derived.

3149 **Lemma G.7.** *Suppose $\max_h d_{\mathcal{A},R_0,\tau,\mathbf{x}}(\mathbf{a}_h, \mathbf{a}'_h \mid r)$ is finite, and further that $\tau_c \geq 3L_{\text{stab}} \log(2c_{\Delta})/\eta$. Then,*

- 3150 • For all $k \in \{0, \dots, \tau_c\}$ and $h \in [H]$,

$$3151 \quad \|\mathbf{x}_{h,k+1} - \mathbf{x}'_{h+1,k+1}\| \leq \max_{h'} \left(2c_{\mathbf{u}} \text{Err}_{\bar{\mathbf{u}},h'} + 2rc_{\mathbf{K}} \text{Err}_{\bar{\mathbf{K}},h'} \beta_{\text{stab}}^{k/3} \right)$$

- 3152 • For all $h \in [H]$ and $1 \leq k \leq \tau_c$,

$$3153 \quad \|\mathbf{x}_{h,k} - \bar{\mathbf{x}}_{h,k}\| \leq 2B_{\text{stab}} r \beta_{\text{stab}}^{k-1}.$$

3154 *Proof of Lemma G.7.* First, an algebraic computation. Observe that $\log(1/\beta_{\text{stab}}) = \log(1/(1 - \frac{\eta}{L_{\text{stab}}})) = -\log(1 - \frac{\eta}{L_{\text{stab}}}) \geq \frac{\eta}{L_{\text{stab}}}$. Hence, if $\tau_c \geq 3L_{\text{stab}} \log(2c_{\Delta})/\eta$, we have $\tau_c \geq 3 \log(2c_{\Delta})/\log(1/\beta_{\text{stab}})$, so that

$$3155 \quad c_{\Delta} \beta_{\text{stab}}^{\tau_c/3} \leq 1/2. \tag{G.7}$$

3156 We continue. Suppose $\max_h d_{\mathcal{A},R_0,\tau,\mathbf{x}}(\mathbf{a}_h, \mathbf{a}'_h \mid r)$ is finite. Then, from the definition of $d_{\mathcal{A},R_0,\tau,\mathbf{x}}$ in [Definition G.10](#), the constants [Definition G.5](#), and the inequalities in [Lemma G.6](#) above, we can check that

$$\begin{aligned} 3157 \quad \max_h \text{Err}_{\bar{\mathbf{u}},h} &\leq \max_h (d_{\mathcal{A},\mathbf{u},\ell_2}(\mathbf{a}_h, \mathbf{a}'_h) + R_0 d_{\mathcal{A},\mathbf{x},\ell_2}(\mathbf{a}_h, \mathbf{a}'_h)) \leq C_{\text{stab},1}(R_0) \\ 3158 \quad \max_h \text{Err}_{\bar{\mathbf{K}},h} &= \max_h d_{\mathcal{A},\mathbf{K},\ell_2}(\mathbf{a}_h, \mathbf{a}'_h) \leq C_{\text{stab},2}(R_0) \\ 3159 \quad \max_h \Delta_{\bar{\mathbf{u}},\infty,h} &= \max_h d_{\mathcal{A},\mathbf{u},\infty}(\mathbf{a}_h, \mathbf{a}'_h) \leq C_{\text{stab},3}(R_0) \\ 3160 \quad r &\leq C_{\text{stab},4}(R_0). \end{aligned}$$

3190 We begin with an induction on states $\|\mathbf{x}_{h,1} - \mathbf{x}'_{h,1}\|$ for $h \geq 1$. Recall the assumption that $c_{\Delta} \beta_{\text{stab}}^{\tau_c/3} \leq 1/2$. We prove
3191 inductively that

$$3192 \quad \forall h \geq 1, \|\mathbf{x}'_{h,1} - \mathbf{x}_{h,1}\| \leq \max_{h'} \left(2c_{\mathbf{u}} \text{Err}_{\bar{\mathbf{u}},h'} + 2rc_{\mathbf{K}} \text{Err}_{\bar{\mathbf{K}},h'} \beta_{\text{stab}}^{\tau_c/3} \right) \quad (\text{G.8})$$

3196 For the base case, we have $\mathbf{x}_{1,1} = \mathbf{x}'_{1,1}$. Now, suppose the result holds up to some $h \geq 1$. Using the definitions of various
3197 constants in [Definitions G.4](#) and [G.5](#), and $r \geq \|\bar{\mathbf{x}}_{h,1} - \mathbf{x}_{h,1}\|$, as well as our inductive hypothesis, one can check that

$$3199 \quad \begin{aligned} \text{Err}_{\bar{\mathbf{u}},h} &\leq C_{\mathbf{u}}, \quad \text{Err}_{\bar{\mathbf{K}},h} \leq C_{\mathbf{K}} \\ 3201 \quad \|\mathbf{x}_{h,1} - \mathbf{x}'_{h,1}\| &\leq \max_{h'} \left(2c_{\mathbf{u}} \text{Err}_{\bar{\mathbf{u}},h'} + 2rc_{\mathbf{K}} \text{Err}_{\bar{\mathbf{K}},h'} \beta_{\text{stab}}^{\tau_c/3} \right) \leq C_{\Delta} \\ 3203 \quad \|\bar{\mathbf{x}}_{h,1} - \mathbf{x}_{h,1}\| &\leq C_{\bar{\mathbf{x}}}, \quad \|\bar{\mathbf{x}}_{h,1} - \mathbf{x}_{h,1}\| \text{Err}_{\mathbf{K}} \leq C_{\mathbf{K},\bar{\mathbf{x}}} \\ 3204 \quad (4R_0 c_{\mathbf{u}} \sqrt{L_{\text{stab}}} + 1) \Delta_{\mathbf{u},\infty} &+ 4R_0 c_{\mathbf{K}} \|\mathbf{x}_{h,1} - \bar{\mathbf{x}}_{h,1}\| + 4R_0 c_{\Delta} \|\mathbf{x}_{h,1} - \mathbf{x}'_{h,1}\| \leq R_{\text{dyn}}. \end{aligned}$$

3206 Then, by [Proposition G.5](#),

$$3209 \quad \begin{aligned} \|\mathbf{x}_{h+1,1} - \mathbf{x}'_{h+1,1}\| &= \|\mathbf{x}_{h,\tau_c+1} - \mathbf{x}'_{h,\tau_c+1}\| \\ 3210 \quad &\leq c_{\mathbf{u}} \text{Err}_{\bar{\mathbf{u}},h+1} + (c_{\mathbf{K}} \text{Err}_{\bar{\mathbf{K}},h+1} \|\mathbf{x}_1 - \bar{\mathbf{x}}_1\| + c_{\Delta} \|\mathbf{x}_{h,1} - \mathbf{x}'_{h,1}\|) \beta_{\text{stab}}^{\tau_c/3} \\ 3211 \quad &\leq c_{\mathbf{u}} \text{Err}_{\bar{\mathbf{u}},h+1} + (c_{\mathbf{K}} \text{Err}_{\bar{\mathbf{K}},h+1} r + c_{\Delta} \|\mathbf{x}_{h,1} - \mathbf{x}'_{h,1}\|) \beta_{\text{stab}}^{\tau_c/3} \quad (c_{\Delta} \beta_{\text{stab}}^{\tau_c/3} \leq \frac{1}{2}, \text{ as established in } (\text{G.7})) \\ 3212 \quad &\leq c_{\mathbf{u}} \text{Err}_{\bar{\mathbf{u}},h} + rc_{\mathbf{K}} \text{Err}_{\bar{\mathbf{K}},h} \beta_{\text{stab}}^{\tau_c/3} + \frac{1}{2} \max_{h'} \left(2c_{\mathbf{u}} \text{Err}_{\bar{\mathbf{u}},h'} + 2rc_{\mathbf{K}} \text{Err}_{\bar{\mathbf{K}},h'} \beta_{\text{stab}}^{\tau_c/3} \right) \\ 3213 \quad &\leq c_{\mathbf{u}} \text{Err}_{\bar{\mathbf{u}},h} + rc_{\mathbf{K}} \text{Err}_{\bar{\mathbf{K}},h} \beta_{\text{stab}}^{\tau_c/3} + \frac{1}{2} \max_{h'} \left(2c_{\mathbf{u}} \text{Err}_{\bar{\mathbf{u}},h'} + 2rc_{\mathbf{K}} \text{Err}_{\bar{\mathbf{K}},h'} \beta_{\text{stab}}^{\tau_c/3} \right) \\ 3214 \quad &\leq c_{\mathbf{u}} \text{Err}_{\bar{\mathbf{u}},h} + rc_{\mathbf{K}} \text{Err}_{\bar{\mathbf{K}},h} \beta_{\text{stab}}^{\tau_c/3} + \frac{1}{2} \max_{h'} \left(2c_{\mathbf{u}} \text{Err}_{\bar{\mathbf{u}},h'} + 2rc_{\mathbf{K}} \text{Err}_{\bar{\mathbf{K}},h'} \beta_{\text{stab}}^{\tau_c/3} \right) \\ 3215 \quad &\leq \max_{h'} \left(2c_{\mathbf{u}} \text{Err}_{\bar{\mathbf{u}},h'} + 2rc_{\mathbf{K}} \text{Err}_{\bar{\mathbf{K}},h'} \beta_{\text{stab}}^{\tau_c/3} \right) \quad (\text{inductive hypothesis}) \\ 3216 \quad &\leq \max_{h'} \left(2c_{\mathbf{u}} \text{Err}_{\bar{\mathbf{u}},h'} + 2rc_{\mathbf{K}} \text{Err}_{\bar{\mathbf{K}},h'} \beta_{\text{stab}}^{\tau_c/3} \right) \end{aligned}$$

3219 This establishes [\(G.8\)](#). A second invocation of [Proposition G.5](#) gives

$$3221 \quad \begin{aligned} \|\mathbf{x}_{h,k+1} - \mathbf{x}'_{h+1,k+1}\| &\leq c_{\mathbf{u}} \text{Err}_{\bar{\mathbf{u}},h+1} + (c_{\mathbf{K}} \text{Err}_{\bar{\mathbf{K}},h+1} \|\mathbf{x}_1 - \bar{\mathbf{x}}_1\| + c_{\Delta} \|\mathbf{x}_{h,1} - \mathbf{x}'_{h,1}\|) \beta_{\text{stab}}^{k/3} \\ 3222 \quad &\leq c_{\mathbf{u}} \text{Err}_{\bar{\mathbf{u}},h+1} + (c_{\mathbf{K}} \text{Err}_{\bar{\mathbf{K}},h+1} r + c_{\Delta} \|\mathbf{x}_{h,1} - \mathbf{x}'_{h,1}\|) \beta_{\text{stab}}^{k/3} \\ 3223 \quad &\leq c_{\mathbf{u}} \text{Err}_{\bar{\mathbf{u}},h+1} + rc_{\mathbf{K}} \text{Err}_{\bar{\mathbf{K}},h+1} \beta_{\text{stab}}^{k/3} + \frac{1}{2} \|\mathbf{x}_{h,1} - \mathbf{x}'_{h,1}\| \\ 3224 \quad &\leq c_{\mathbf{u}} \text{Err}_{\bar{\mathbf{u}},h+1} + rc_{\mathbf{K}} \text{Err}_{\bar{\mathbf{K}},h} \beta_{\text{stab}}^{k/3} + \max_{h'} \left(c_{\mathbf{u}} \text{Err}_{\bar{\mathbf{u}},h'} + rc_{\mathbf{K}} \text{Err}_{\bar{\mathbf{K}},h'} \beta_{\text{stab}}^{\tau_c/3} \right) \\ 3225 \quad &\leq c_{\mathbf{u}} \text{Err}_{\bar{\mathbf{u}},h+1} + rc_{\mathbf{K}} \text{Err}_{\bar{\mathbf{K}},h} \beta_{\text{stab}}^{k/3} + \max_{h'} \left(c_{\mathbf{u}} \text{Err}_{\bar{\mathbf{u}},h'} + rc_{\mathbf{K}} \text{Err}_{\bar{\mathbf{K}},h'} \beta_{\text{stab}}^{\tau_c/3} \right) \\ 3226 \quad &\leq \max_{h'} \left(2c_{\mathbf{u}} \text{Err}_{\bar{\mathbf{u}},h'} + 2rc_{\mathbf{K}} \text{Err}_{\bar{\mathbf{K}},h'} \beta_{\text{stab}}^{k/3} \right). \end{aligned}$$

3232 Moreover, as [Proposition G.5](#) implies that the conclusions of [Lemma G.4](#) also hold, we further find that

$$3235 \quad \|\mathbf{x}_{h,k} - \bar{\mathbf{x}}_{h,k}\| \leq 2B_{\text{stab}} \|\mathbf{x}_{h,k} - \bar{\mathbf{x}}_{h,k}\| \beta_{\text{stab}}^{k-1} \leq 2B_{\text{stab}} r \beta_{\text{stab}}^{k-1},$$

3237 as needed. □

3241 G.5.3. CONCLUDING THE PROOF OF [PROPOSITION G.3](#).

3242 *Completing the proof of [Proposition G.3](#).* Let us start with the first item, bound $d_{\mathcal{S},\mathbf{x},\tau}$. We may assume that
3243 $d_{\mathcal{A},R_0,\tau,\mathbf{x}}(\mathbf{a}_h, \mathbf{a}'_h \mid r)$ is finite for all h .
3244

3245 **Controlling** $d_{\mathcal{S},\mathbf{x},\tau}(s_h, s'_h)$. Not that $s_1 = s'_1$. For any $2 \leq h \leq H + 1$,

$$\begin{aligned}
3246 & \\
3247 & d_{\mathcal{S},\mathbf{x},\tau}(s_h, s'_h) := \max_{t \in [t_h - \tau : t_h]} \|\mathbf{x}_t - \mathbf{x}'_t\| \\
3248 & \\
3249 & = \max_{\tau_c - \tau \leq k \leq \tau_c} \|\mathbf{x}_{h-1,1+k} - \mathbf{x}'_{h-1,k+1}\| \quad (\text{our indexing scheme}) \\
3250 & \\
3251 & \leq \max_{\tau_c - \tau \leq k \leq \tau_c} \max_{h'} \left(2c_{\mathbf{u}} \text{Err}_{\bar{\mathbf{u}},h'} + 2rc_{\mathbf{K}} \text{Err}_{\bar{\mathbf{K}},h'} \beta_{\text{stab}}^{k/3} \right) \quad (\text{Lemma G.7}) \\
3252 & \\
3253 & = \max_{h'} \left(2c_{\mathbf{u}} \text{Err}_{\bar{\mathbf{u}},h'} + 2rc_{\mathbf{K}} \text{Err}_{\bar{\mathbf{K}},h'} \beta_{\text{stab}}^{(\tau_c - \tau)/3} \right) \\
3254 & \\
3255 & \leq \max_{h'} \left(2c_{\mathbf{u}} (d_{\mathcal{A},\mathbf{u},\ell_2}(\mathbf{a}_{h'}, \mathbf{a}'_{h'})) + R_0 d_{\mathcal{A},\mathbf{x},\ell_2}(\mathbf{a}_{h'}, \mathbf{a}'_{h'}) \right) + 2rc_{\mathbf{K}} d_{\mathcal{A},\mathbf{K},\ell_2}(\mathbf{a}_h, \mathbf{a}'_h) \beta_{\text{stab}}^{(\tau_c - \tau)/3} \quad (\text{Lemma G.6}) \\
3256 & \\
3257 & \leq \max_{h'} d_{\mathcal{A},R_0,\tau,\mathbf{x}}(\mathbf{a}_{h'}, \mathbf{a}'_{h'} \mid r) \quad (\text{Definition G.10}) \\
3258 &
\end{aligned}$$

3259 That is,

$$3260 \\
3261 d_{\mathcal{S},\mathbf{x},\tau}(s_h, s'_h) \leq \max_h d_{\mathcal{A},R_0,\tau,\mathbf{x}}(\mathbf{a}_h, \mathbf{a}'_h \mid r) \quad (\text{G.9}) \\
3262 \\
3263$$

3264 **Bounding** $d_{\mathcal{S},\tau}$. To bound $d_{\mathcal{S},\tau}$, we also need to account for differents in inputs. We have

$$\begin{aligned}
3265 & \\
3266 & \mathbf{u}_{h,k} - \mathbf{u}'_{h,k} = \kappa_{h,k}(\mathbf{x}_{h,k}) - \kappa'_{h,k}(\mathbf{x}'_{h,k}) \\
3267 & = \bar{\mathbf{u}}_{h,k} - \bar{\mathbf{u}}'_{h,k} + \bar{\mathbf{K}}_t(\mathbf{x}_{h,k} - \bar{\mathbf{x}}_{h,k}) - \bar{\mathbf{K}}'_{h,k}(\mathbf{x}'_{h,k} - \bar{\mathbf{x}}'_{h,k}) \\
3268 & = \bar{\mathbf{u}}_{h,k} - \bar{\mathbf{u}}'_{h,k} + (\bar{\mathbf{K}}_{h,k} - \bar{\mathbf{K}}'_{h,k})(\mathbf{x}_{h,k} - \bar{\mathbf{x}}_{h,k}) + \bar{\mathbf{K}}'_{h,k}(\mathbf{x}_t - \mathbf{x}'_{h,k} - (\bar{\mathbf{x}}_{h,k} - \bar{\mathbf{x}}'_{h,k})) \\
3269 & \\
3270 & \\
3271 & \\
3272 &
\end{aligned}$$

3273 Where $d_{\mathcal{A},R_0,\tau,\mathbf{x}}(\mathbf{a}_h, \mathbf{a}'_h \mid r)$ is finite for all h , then $\|\bar{\mathbf{K}}'_{h,k}\| = \|\bar{\mathbf{K}}_{t_h+k-1}\| \leq R_0$. Thus,

$$\begin{aligned}
3274 & \\
3275 & \|\mathbf{u}_{h,k} - \mathbf{u}'_{h,k}\| \\
3276 & \leq \|\bar{\mathbf{u}}_{h,k} - \bar{\mathbf{u}}'_{h,k}\| + R_0 \|\bar{\mathbf{x}}_{h,k} - \bar{\mathbf{x}}'_{h,k}\| + R_0 \|\bar{\mathbf{x}}_{h,k} - \bar{\mathbf{x}}'_{h,k}\| + \|\bar{\mathbf{K}}_{h,k} - \bar{\mathbf{K}}'_{h,k}\| \|\mathbf{x}_{h,k} - \bar{\mathbf{x}}_{h,k}\| \\
3277 & \leq d_{\mathcal{A},\mathbf{u},\infty}(\mathbf{a}_h, \mathbf{a}'_h) + R_0 d_{\mathcal{A},\mathbf{x},\infty}(\mathbf{a}_h, \mathbf{a}'_h) + R_0 \|\bar{\mathbf{x}}_{h,k} - \bar{\mathbf{x}}'_{h,k}\| + d_{\mathcal{A},\mathbf{K},\infty}(\mathbf{a}_h, \mathbf{a}'_h) \|\mathbf{x}_{h,k} - \bar{\mathbf{x}}_{h,k}\| \\
3278 & \leq d_{\mathcal{A},\mathbf{u},\infty}(\mathbf{a}_h, \mathbf{a}'_h) + R_0 d_{\mathcal{A},\mathbf{x},\infty}(\mathbf{a}_h, \mathbf{a}'_h) + 2r B_{\text{stab}} \beta_{\text{stab}}^{k-1} d_{\mathcal{A},\mathbf{K},\infty}(\mathbf{a}_h, \mathbf{a}'_h) + R_0 \|\bar{\mathbf{x}}_{h,k} - \bar{\mathbf{x}}'_{h,k}\| \\
3279 & \\
3280 &
\end{aligned}$$

3281 Hence,

$$\begin{aligned}
3282 & \\
3283 & \max_h d_{\mathcal{S},\tau}(s_h, s'_h) \\
3284 & = \max_h d_{\mathcal{S},\mathbf{x},\tau}(\mathbf{a}_h, \mathbf{a}'_h) \vee \max_h \max_{\tau_c - \tau \leq k \leq \tau_c - 1} \|\mathbf{u}_{h,k+1} - \mathbf{u}'_{h,k+1}\| \\
3285 & = \max_h d_{\mathcal{S},\mathbf{x},\tau}(\mathbf{a}_h, \mathbf{a}'_h) + R_0 \max_h \max_{\tau_c - \tau \leq k \leq \tau_c - 1} \|\bar{\mathbf{x}}_{h,k+1} - \bar{\mathbf{x}}'_{h,k+1}\| \\
3286 & \\
3287 & + \max_h \max_{\tau_c - \tau \leq k \leq \tau_c - 1} (d_{\mathcal{A},\mathbf{u},\infty}(\mathbf{a}_h, \mathbf{a}'_h) + R_0 d_{\mathcal{A},\mathbf{x},\infty}(\mathbf{a}_h, \mathbf{a}'_h) + 2r B_{\text{stab}} \beta_{\text{stab}}^k d_{\mathcal{A},\mathbf{K},\infty}(\mathbf{a}_h, \mathbf{a}'_h)) \\
3288 & = (1 + R_0) \max_h d_{\mathcal{S},\mathbf{x},\tau}(\mathbf{a}_h, \mathbf{a}'_h) \\
3289 & \\
3290 & + \max_h \left(\underbrace{d_{\mathcal{A},\mathbf{u},\infty}(\mathbf{a}_h, \mathbf{a}'_h) + R_0 d_{\mathcal{A},\mathbf{x},\infty}(\mathbf{a}_h, \mathbf{a}'_h) + 2r B_{\text{stab}} \beta_{\text{stab}}^{\tau_c - \tau} d_{\mathcal{A},\mathbf{K},\infty}(\mathbf{a}_h, \mathbf{a}'_h)}_{:= \max_h d_{\mathcal{A},R_0,\tau,\mathbf{u}}(\mathbf{a}_h, \mathbf{a}'_h \mid r)} \right) \\
3291 & \\
3292 & \leq (1 + R_0) \max_h d_{\mathcal{A},R_0,\tau,\mathbf{x}}(\mathbf{a}_h, \mathbf{a}'_h \mid r) + \max_h d_{\mathcal{A},R_0,\tau,\mathbf{u}}(\mathbf{a}_h, \mathbf{a}'_h \mid r) \quad ((\text{G.9})) \\
3293 & \\
3294 & \leq \max_h 2 \left((1 + R_0) d_{\mathcal{A},R_0,\tau,\mathbf{x}}(\mathbf{a}_h, \mathbf{a}'_h \mid r) + d_{\mathcal{A},R_0,\tau,\mathbf{u}}(\mathbf{a}_h, \mathbf{a}'_h \mid r) \right). \\
3295 & \\
3296 & \\
3297 & \\
3298 & \\
3299 &
\end{aligned}$$

3300 **Bounding** $d_{\mathcal{S},\tau}(s_{h+1}, F_h(\tilde{s}_h, \mathbf{a}_h))$. Next, by (G.5) and (G.6) that

$$\begin{aligned}
3301 \quad d_{\mathcal{S},\tau}(s_{h+1}, F_h(\tilde{s}_h, \mathbf{a}_h)) &= \max_{\tau_c - \tau \leq k \leq \tau_c} \|\mathbf{x}_{h,k+1} - \bar{\mathbf{x}}_{h,k+1}\| && ((\text{G.5}) \text{ and } (\text{G.6})) \\
3302 \quad &= \max_{\tau_c - \tau \leq k \leq \tau_c} 2B_{\text{stab}} r \beta_{\text{stab}}^k && (\text{Lemma G.7}) \\
3303 \quad &= 2B_{\text{stab}} r \beta_{\text{stab}}^{(\tau_c - \tau)}. \\
3304 \quad & && \\
3305 \quad & && \\
3306 \quad & &&
\end{aligned}$$

3307 Thus, We have

$$\begin{aligned}
3308 \quad d_{\mathcal{S},\tau}(s_{h+1}, F_h(\tilde{s}_h, \mathbf{a}_h)) &= d_{\mathcal{S},\mathbf{x},\tau}(s_{h+1}, F_h(\tilde{s}_h, \mathbf{a}_h)) \vee \max_{\tau_c - \tau \leq k \leq \tau_c - 1} \|\mathbf{u}_{h,k+1} - \bar{\mathbf{u}}_{h,k+1}\| \\
3309 \quad &= d_{\mathcal{S},\mathbf{x},\tau}(s_{h+1}, F_h(\tilde{s}_h, \mathbf{a}_h)) \vee \max_{\tau_c - \tau \leq k \leq \tau_c - 1} \|\kappa_{h,k+1}(\mathbf{x}_{h,k+1}) - \bar{\mathbf{u}}_{h,k+1}\| \\
3310 \quad &= d_{\mathcal{S},\mathbf{x},\tau}(s_{h+1}, F_h(\tilde{s}_h, \mathbf{a}_h)) \vee \max_{\tau_c - \tau \leq k \leq \tau_c - 1} \|(\bar{\mathbf{K}}_{h,k+1}(\mathbf{x}_{h,k+1} - \bar{\mathbf{x}}_{h,k+1}) + \bar{\mathbf{u}}_{h,k+1}) - \bar{\mathbf{u}}_{h,k+1}\| \\
3311 \quad &\leq d_{\mathcal{S},\mathbf{x},\tau}(s_{h+1}, F_h(\tilde{s}_h, \mathbf{a}_h)) \vee \max_{\tau_c - \tau \leq k \leq \tau_c - 1} \|\bar{\mathbf{K}}_{h,k+1}\| \|\mathbf{x}_{h,k+1} - \bar{\mathbf{x}}_{h,k+1}\| \\
3312 \quad &\stackrel{(i)}{\leq} d_{\mathcal{S},\mathbf{x},\tau}(s_{h+1}, F_h(\tilde{s}_h, \mathbf{a}_h)) \vee \max_{\tau_c - \tau \leq k \leq \tau_c - 1} R_{\text{stab}} \|\mathbf{x}_{h,k+1} - \bar{\mathbf{x}}_{h,k+1}\| \\
3313 \quad &\leq (1 + R_{\text{stab}}) d_{\mathcal{S},\mathbf{x},\tau}(s_{h+1}, F_h(\tilde{s}_h, \mathbf{a}_h)) \\
3314 \quad &\leq 2(1 + R_{\text{stab}}) B_{\text{stab}} r \beta_{\text{stab}}^{(\tau_c - \tau)} \\
3315 \quad & &&
\end{aligned}$$

3316 where in (i), we used $\|\bar{\mathbf{K}}_{h,k+1}\| \leq R_{\text{stab}}$ because $(\bar{\rho}_{[h+1]}, \kappa_{h,1;\tau_c})$ is $(R_{\text{stab}}, B_{\text{stab}}, L_{\text{stab}})$ -stable, so that the gains are bounded in operator norm by R_{stab} . \square

3324 G.6. Proof of Lemma G.4 (state perturbation)

3325 Define $\bar{\Delta}_{\mathbf{x},k} = \mathbf{x}_k - \bar{\mathbf{x}}_k$. Then

$$\begin{aligned}
3326 \quad \bar{\Delta}_{\mathbf{x},k+1} &= \bar{\Delta}_{\mathbf{x},k} + \eta (f_\eta(\mathbf{x}_k, \bar{\mathbf{u}}_k + \bar{\mathbf{K}}_k(\mathbf{x}_k - \bar{\mathbf{x}}_k)) - f_\eta(\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k)) \\
3327 \quad &= \bar{\Delta}_{\mathbf{x},k} + \eta(\mathbf{A}_k + \mathbf{B}_k \mathbf{K}_k) \bar{\Delta}_{\mathbf{x},k} + \text{rem}_k, \\
3328 \quad & && (\text{G.10}) \\
3329 \quad & &&
\end{aligned}$$

3330 where

$$\text{rem}_k = f_\eta(\mathbf{x}_k, \bar{\mathbf{u}}_k + \mathbf{K}_k(\mathbf{x}_k - \bar{\mathbf{x}}_k)) - f_\eta(\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k) - (\mathbf{A}_k + \mathbf{B}_k \mathbf{K}_k) \bar{\Delta}_{\mathbf{x},k}.$$

3331 **Claim G.8.** Take $R_{\text{stab}} \geq 1$, and suppose $\|\bar{\Delta}_{\mathbf{x},k}\| \leq R_{\text{dyn}}/R_{\text{stab}}$. Then,

$$\|\bar{\mathbf{x}}_k - \mathbf{x}_k\| \vee \|\bar{\mathbf{u}}_k - \mathbf{u}_k\| \leq R_{\text{dyn}}, \tag{G.11}$$

3332 and $\|\text{rem}_k\| \leq M_{\text{dyn}} R_{\text{stab}}^2 \|\bar{\Delta}_{\mathbf{x},k}\|^2$.

3333 *Proof.* Let $\mathbf{u}_k = \bar{\mathbf{u}}_k + \mathbf{K}_k(\mathbf{x}_k - \bar{\mathbf{x}}_k)$. The conditions of the claim imply $\|\mathbf{u}_k - \bar{\mathbf{u}}_k\| \vee \|\mathbf{x}_k - \bar{\mathbf{x}}_k\| \leq R_{\text{dyn}}$. From Taylor's theorem and the fact that $\bar{\rho}$ is $(R_{\text{dyn}}, L_{\text{dyn}}, M_{\text{dyn}})$ -regular imply that

$$\begin{aligned}
3334 \quad \|f_\eta(\mathbf{x}_k, \mathbf{u}_k) - f_\eta(\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k)\| &\leq \frac{1}{2} M_{\text{dyn}} (\|\mathbf{x}_k - \bar{\mathbf{x}}_k\|^2 + \|\mathbf{u}_k - \bar{\mathbf{u}}_k\|^2) \\
3335 \quad &\leq \frac{1}{2} (1 + R_{\text{stab}}^2) M_{\text{dyn}} \|\mathbf{x}_k - \bar{\mathbf{x}}_k\|^2 \leq R_{\text{stab}}^2 M_{\text{dyn}} \|\bar{\Delta}_{\mathbf{x},k}\|^2, \\
3336 \quad & &&
\end{aligned}$$

3337 where again use $R_{\text{stab}} \geq 1$ above. \square

3338 Solving the recursion from (G.10), we have

$$\bar{\Delta}_{\mathbf{x},k+1} = \eta \sum_{j=1}^k \Phi_{\text{cl},k+1,j+1} \text{rem}_j + \Phi_{\text{cl},k+1,1} \bar{\Delta}_{\mathbf{x},1}.$$

3355 Set $\beta_{\text{stab}} := (1 - \frac{\eta}{L_{\text{stab}}})$, so that $M := \frac{\eta}{\beta_{\text{stab}}^{-1} - 1} = L_{\text{stab}}$. Further, recall $R_0 \leq R_{\text{dyn}}/R_{\text{stab}}$. By assumption, $\Phi_{\text{cl},k,j} \leq$
 3356 $B_{\text{stab}}\beta_{\text{stab}}^{k-j}$, so using [Claim G.8](#) implies that, if $\max_{j \in [k]} \|\bar{\Delta}_{\mathbf{x},j}\| \leq R_0 \leq R_{\text{dyn}}/R_{\text{stab}}$ for all $j \in [k]$,

$$3357 \quad \|\bar{\Delta}_{\mathbf{x},k+1}\| \leq \eta \sum_{j=1}^k B_{\text{stab}} M_{\text{dyn}} R_{\text{stab}}^2 \beta_{\text{stab}}^{k-j} \|\bar{\Delta}_{\mathbf{x},j}\|^2 + B_{\text{stab}} \beta_{\text{stab}}^k \|\bar{\Delta}_{\mathbf{x},1}\|.$$

3361 Applying [Lemma G.17](#) with $\alpha = 0$, $C_1 = B_{\text{stab}} M_{\text{dyn}} R_{\text{stab}}^2$, and $C_2 = B_{\text{stab}} \geq 1$ and $M = L_{\text{stab}}$ (noting $\beta_{\text{stab}} \geq 1/2$), it
 3362 holds that for $\|\bar{\Delta}_{\mathbf{x},1}\| = \varepsilon_1 \leq 1/4 M C_1 C_3 = 1/4 L_{\text{stab}} M_{\text{dyn}} R_{\text{stab}}^2 B_{\text{stab}}^2$,

$$3363 \quad \|\bar{\Delta}_{\mathbf{x},k+1}\| \leq 2B_{\text{stab}} \|\bar{\Delta}_{\mathbf{x},1}\| (1 - \frac{\eta}{L_{\text{stab}}})^k.$$

3364 To ensure the inductive hypothesis that $\max_{j \in [k]} \|\bar{\Delta}_{\mathbf{x},j}\| \leq R_{\text{dyn}} R_{\text{stab}}$, it suffices to ensure that $2B_{\text{stab}} \|\bar{\Delta}_{\mathbf{x},1}\| \leq R_0$,
 3365 which is assumed by the lemma. Thus, we have shown that, if

$$3366 \quad \|\bar{\Delta}_{\mathbf{x},1}\| \leq \min\{1/2 B_{\text{stab}} R_0, 1/8 L_{\text{stab}} M_{\text{dyn}} R_{\text{stab}}^2 B_{\text{stab}}^2\},$$

3367 it holds that $\|\bar{\Delta}_{\mathbf{x},k+1}\| \leq 2B_{\text{stab}} \|\bar{\Delta}_{\mathbf{x},1}\| (1 - \frac{\eta}{L_{\text{stab}}})^k \leq R_0$ for all k .

3368 Next, we address the stability of the gains for the perturbed trajectory ρ . Using $(R_{\text{dyn}}, L_{\text{dyn}}, M_{\text{dyn}})$ -regularity of $\bar{\rho}$ and
 3369 [\(G.11\)](#),

$$\begin{aligned} 3370 & \|\mathbf{A}_k(\rho) + \mathbf{B}_k(\rho)\mathbf{K}_k - \mathbf{A}_k(\bar{\rho}) + \mathbf{B}_k(\bar{\rho})\mathbf{K}_k\| \\ 3371 &= \left\| \begin{bmatrix} \mathbf{A}_k(\rho) - \mathbf{A}_k(\bar{\rho}) & \hat{\mathbf{B}}_k(\rho) - \mathbf{B}_k(\bar{\rho}) \end{bmatrix} \begin{bmatrix} \mathbf{I} \\ \mathbf{K}_k \end{bmatrix} \right\| \\ 3372 &= \left\| (\nabla f_{\eta}(\hat{\mathbf{x}}_k, \mathbf{u}_k) - \nabla f_{\eta}(\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k)) \begin{bmatrix} \mathbf{I} \\ \mathbf{K}_k \end{bmatrix} \right\| \\ 3373 &\leq M_{\text{dyn}} \|(\mathbf{x}_k - \bar{\mathbf{x}}_k, \mathbf{K}_k(\mathbf{x}_k - \bar{\mathbf{x}}_k))\| \left\| \begin{bmatrix} \mathbf{I} \\ \mathbf{K}_k \end{bmatrix} \right\| \\ 3374 &= M_{\text{dyn}} \|\mathbf{x}_k - \bar{\mathbf{x}}_k\| \left\| \begin{bmatrix} \mathbf{I} \\ \mathbf{K}_k \end{bmatrix} \right\|^2 \leq M_{\text{dyn}} \|\mathbf{x}_k - \bar{\mathbf{x}}_k\| (1 + \|\mathbf{K}_k\|_{\text{op}}^2) \\ 3375 &= M_{\text{dyn}} \|\mathbf{x}_k - \bar{\mathbf{x}}_k\| \left\| \begin{bmatrix} \mathbf{I} \\ \mathbf{K}_k \end{bmatrix} \right\|^2 \leq M_{\text{dyn}} \|\mathbf{x}_k - \bar{\mathbf{x}}_k\| (1 + \|\mathbf{K}_k\|_{\text{op}}^2) \\ 3376 &\leq 2R_{\text{stab}}^2 M_{\text{dyn}} \|\mathbf{x}_k - \bar{\mathbf{x}}_k\| \\ 3377 &\leq 4B_{\text{stab}} R_{\text{stab}}^2 M_{\text{dyn}} \|\mathbf{x}_1 - \bar{\mathbf{x}}_1\| \beta_{\text{stab}}^{k-1}, \quad \beta_{\text{stab}} = (1 - \frac{\eta}{L_{\text{stab}}}). \end{aligned}$$

3378 Invoking [Lemma G.20](#) with $\beta_{\text{stab}} \geq 1/2$, $\|\hat{\Phi}_{\text{cl},k,j}\| \leq 2B_{\text{stab}}\beta_{\text{stab}}^{k-j}$ for all j, k provided that $4B_{\text{stab}} R_{\text{stab}}^2 M_{\text{dyn}} \|\mathbf{x}_1 - \bar{\mathbf{x}}_1\| \leq$
 3379 $1/4 B_{\text{stab}} L_{\text{stab}}$, which requires $\|\mathbf{x}_1 - \bar{\mathbf{x}}_1\| \leq 1/16 B_{\text{stab}}^2 R_{\text{stab}}^2 L_{\text{stab}} M_{\text{dyn}}$.

3380 The last part of the lemma uses $(R_{\text{dyn}}, L_{\text{dyn}}, M_{\text{dyn}})$ -regularity of $\bar{\rho}$ and [\(G.11\)](#).

3381 G.7. Proof of [Proposition G.5](#) (input and gain perturbation)

3382 Recall the trajectories $\bar{\mathbf{x}}_{k+1} = \bar{\mathbf{x}}_k + \eta f_{\eta}(\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k)$, and

$$\begin{aligned} 3383 & \mathbf{x}_{k+1} = \mathbf{x}_k + \eta f_{\eta}(\mathbf{x}_k, \mathbf{u}_k), \quad \mathbf{u}_k = \bar{\mathbf{u}}_k + \mathbf{K}_k(\mathbf{x}_k - \bar{\mathbf{x}}_k) \\ 3384 & \mathbf{x}'_{k+1} = \mathbf{x}'_k + \eta f_{\eta}(\mathbf{x}'_k, \mathbf{u}'_k), \quad \mathbf{u}'_k = \bar{\mathbf{u}}'_k + \mathbf{K}'_k(\mathbf{x}'_k - \bar{\mathbf{x}}_k). \end{aligned}$$

3385 Further introduce the shorthand $\hat{\mathbf{A}}_k = \mathbf{A}_k(\hat{\rho})$, $\hat{\mathbf{B}}_k = \mathbf{B}_k(\hat{\rho})$, $\hat{\mathbf{A}}_{\text{cl},k} = \hat{\mathbf{A}}_k + \hat{\mathbf{B}}_k + \mathbf{K}_k$, as well as

$$\begin{aligned} 3386 & \Delta_{\mathbf{x},k} = \mathbf{x}'_k - \mathbf{x}_k, \quad \Delta_{\mathbf{u},k} = \bar{\mathbf{u}}'_k - \bar{\mathbf{u}}_k, \quad \Delta_{\mathbf{K},k} = \mathbf{K}'_k - \mathbf{K}_k \\ 3387 & \bar{\Delta}_{\mathbf{x},k} = \mathbf{x}'_k - \bar{\mathbf{x}}_k, \quad \bar{\Delta}_{\mathbf{x},k} = \mathbf{x}_k - \bar{\mathbf{x}}_k, \end{aligned}$$

3388

3410 Then,

$$\begin{aligned}
3411 & \Delta_{\mathbf{x},k+1} = \Delta_{\mathbf{x},k} + \eta \left(f(\mathbf{x}_k + \Delta_{\mathbf{x},k}, \bar{\mathbf{u}}_k + \Delta_{\mathbf{u},k} + \tilde{\mathbf{K}}_k \tilde{\Delta}_{\mathbf{x},k}) - f(\mathbf{x}_k, \bar{\mathbf{u}}_k + \mathbf{K}_k \bar{\Delta}_{\mathbf{x},k}) \right) \\
3412 & = \Delta_{\mathbf{x},k} + \eta \left(f(\mathbf{x}_k + \Delta_{\mathbf{x},k}, \bar{\mathbf{u}}_k + \Delta_{\mathbf{u},k} + \mathbf{K}_k \tilde{\Delta}_{\mathbf{x},k}) - f(\mathbf{x}_k, \bar{\mathbf{u}}_k + \mathbf{K}_k \bar{\Delta}_{\mathbf{x},k}) \right) \\
3413 & \quad + \eta(\text{rem}_{k,1}) \\
3414 & = \Delta_{\mathbf{x},k} + \eta \left(\underbrace{\frac{\partial}{\partial \mathbf{x}} f(\mathbf{x}_k, \mathbf{u}_k)}_{=\hat{\mathbf{A}}_k} \Delta_{\mathbf{x},k} + \underbrace{\frac{\partial}{\partial \mathbf{u}} f(\mathbf{x}_k, \mathbf{u}_k)}_{=\hat{\mathbf{B}}_k} (\Delta_{\mathbf{u},k} + \mathbf{K}_k \underbrace{\tilde{\Delta}_{\mathbf{x},k} - \bar{\Delta}_{\mathbf{x},k}}_{\Delta_{\mathbf{x},k}}) \right) \\
3415 & \quad + \eta(\text{rem}_{k,1} + \text{rem}_{k,2}) \\
3416 & = \Delta_{\mathbf{x},k} + \eta \left(\hat{\mathbf{A}}_{\text{cl},k} \Delta_{\mathbf{x},k} + \hat{\mathbf{B}}_k \Delta_{\mathbf{u},k} \right) + \eta(\text{rem}_{k,1} + \text{rem}_{k,2}).
\end{aligned}$$

3425 where, above

$$\begin{aligned}
3426 & \text{rem}_{k,1} = f_\eta(\mathbf{x}_k + \Delta_{\mathbf{x},k}, \bar{\mathbf{u}}_k + \Delta_{\mathbf{u},k} + \mathbf{K}'_k \tilde{\Delta}_{\mathbf{x},k}) - f_\eta(\mathbf{x}_k + \Delta_{\mathbf{x},k}, \bar{\mathbf{u}}_k + \Delta_{\mathbf{u},k} + \mathbf{K}_k \tilde{\Delta}_{\mathbf{x},k}) \\
3427 & \text{rem}_{k,2} = f_\eta(\mathbf{x}_k + \Delta_{\mathbf{x},k}, \bar{\mathbf{u}}_k + \Delta_{\mathbf{u},k} + \mathbf{K}_k \tilde{\Delta}_{\mathbf{x},k}) - f(\mathbf{x}_k, \bar{\mathbf{u}}_k + \mathbf{K}_k \bar{\Delta}_{\mathbf{x},k}) \\
3428 & \quad - \frac{\partial}{\partial \mathbf{x}} f_\eta(\mathbf{x}_k, \mathbf{u}_k) \Delta_{\mathbf{x},k} + \frac{\partial}{\partial \mathbf{u}} f_\eta(\mathbf{x}_k, \mathbf{u}_k) (\Delta_{\mathbf{u},k} + \mathbf{K}_k (\tilde{\Delta}_{\mathbf{x},k} - \bar{\Delta}_{\mathbf{x},k})).
\end{aligned}$$

3433 Solving the recursion,

$$\Delta_{\mathbf{x},k+1} = \sum_{j=1}^k \hat{\Phi}_{\text{cl},k+1,j+1} (\hat{\mathbf{B}}_j \Delta_{\mathbf{u},j} + \eta(\text{rem}_{j,1} + \text{rem}_{j,2})) + \hat{\Phi}_{\text{cl},k+1,1} \Delta_{\mathbf{x},1}$$

3439 Recall that [Lemma G.4](#) implies $(\mathbf{K}_{1:K}, \boldsymbol{\rho})$ is $(R_{\text{stab}}, 2B_{\text{stab}}, L_{\text{stab}})$ -stable. Thus, recalling $\beta_{\text{stab}} = (1 - \frac{\eta}{L_{\text{stab}}}) \in [1/2, 1)$,
3440 we have

$$\begin{aligned}
3441 & \|\Delta_{\mathbf{x},k+1}\| \leq \eta \sum_{j=1}^k \|\hat{\Phi}_{\text{cl},k+1,j+1}\| (\|\hat{\mathbf{B}}_j\| \|\Delta_{\mathbf{u},j}\| + \|\text{rem}_{j,1}\| + \|\text{rem}_{j,2}\|) + \|\hat{\Phi}_{\text{cl},k+1,1}\| \|\Delta_{\mathbf{x},1}\| \\
3442 & \leq \eta \sum_{j=1}^k 2B_{\text{stab}} \beta_{\text{stab}}^{k-j} (L_{\text{dyn}} \|\Delta_{\mathbf{u},j}\| + \|\text{rem}_{j,1}\| + \|\text{rem}_{j,2}\|) + 2B_{\text{stab}} \beta_{\text{stab}}^k \|\Delta_{\mathbf{x},1}\|.
\end{aligned}$$

3449 Let us now bound each of these remainder terms. The following claim, as well as all subsequent claims, is proven at the end
3450 of the section.

3451 **Claim G.9.** *Suppose that it holds that for a given k , it holds that*

$$\|\Delta_{\mathbf{x},k}\| \leq c_{\mathbf{u}} \text{Err}_{\mathbf{u}} + c_{\mathbf{K}} \text{Err}_{\mathbf{K}} \|\mathbf{x}_1 - \bar{\mathbf{x}}_1\| + c_{\Delta} \|\mathbf{x}_1 - \mathbf{x}'_1\| \tag{G.12}$$

3456 Then,

$$\begin{aligned}
3457 & \|\text{rem}_{k,1}\| \leq L_{\text{dyn}} \|\Delta_{\mathbf{K},k}\| (\|\bar{\Delta}_{\mathbf{x},k}\| + \|\Delta_{\mathbf{x},k}\|) \\
3458 & \|\text{rem}_{k,2}\| \leq \frac{3}{2} M_{\text{dyn}} R_{\text{stab}}^2 \|\Delta_{\mathbf{x},k}\|^2 + M_{\text{dyn}} \|\Delta_{\mathbf{u},k}\|^2
\end{aligned}$$

3463 We now proceed by strong induction on the condition in [\(G.12\)](#). Observe that if this condition holds for all $1 \leq j \leq k$, we

3465 have

$$\begin{aligned}
3466 \quad \|\Delta_{\mathbf{x},k+1}\| &\leq \eta \sum_{j=1}^k 2B_{\text{stab}}\beta_{\text{stab}}^{k-j} (L_{\text{dyn}}\|\Delta_{\mathbf{u},j}\| + \|\text{rem}_{j,1}\| + \|\text{rem}_{j,2}\|) + 2B_{\text{stab}}\beta_{\text{stab}}^k \|\Delta_{\mathbf{x},1}\| \\
3467 & \\
3468 & \\
3469 & \\
3470 &\leq \underbrace{\eta \sum_{j=1}^k 2B_{\text{stab}}\beta_{\text{stab}}^{k-j} (L_{\text{dyn}}\|\Delta_{\mathbf{u},j}\| + M_{\text{dyn}}\|\Delta_{\mathbf{u},j}\|^2)}_{=\text{Term}_{1,k}} \\
3471 & \\
3472 & \\
3473 & \\
3474 &+ \eta \sum_{j=1}^k \beta_{\text{stab}}^{k-j} \left(\underbrace{3B_{\text{stab}}M_{\text{dyn}}\|\Delta_{\mathbf{x},j}\|^2}_{C_1} + \underbrace{2B_{\text{stab}}L_{\text{dyn}}\|\Delta_{\mathbf{K},j}\|\|\Delta_{\mathbf{x},j}\|}_{C_2} \right) \\
3475 & \\
3476 & \\
3477 & \\
3478 &+ \underbrace{2B_{\text{stab}}\beta_{\text{stab}}^k \|\Delta_{\mathbf{x},1}\| + \eta \sum_{j=1}^k 2L_{\text{dyn}}B_{\text{stab}}\beta_{\text{stab}}^{k-j} \|\Delta_{\mathbf{K},j}\|\|\bar{\Delta}_{\mathbf{x},j}\|}_{\text{Term}_{2,k}} \\
3479 & \tag{G.13} \\
3480 & \\
3481 & \\
3482 &
\end{aligned}$$

3483 Define the terms

$$\begin{aligned}
3484 & \\
3485 \quad C_1 &:= 3B_{\text{stab}}M_{\text{dyn}}, \quad C_2 := 2B_{\text{stab}}L_{\text{dyn}}, \\
3486 \quad \alpha &:= 2B_{\text{stab}}\text{Err}_{\mathbf{u}} \left(M_{\text{dyn}}\text{Err}_{\mathbf{u}} + \sqrt{L_{\text{stab}}L_{\text{dyn}}} \right) \\
3487 & \\
3488 \quad \bar{\varepsilon}_1 &:= 2B_{\text{stab}} \left(\|\Delta_{\mathbf{x},1}\| + 2L_{\text{stab}}^{1/2}\text{Err}_{\mathbf{K}}\|\bar{\Delta}_{\mathbf{x},1}\| \right) \\
3489 &
\end{aligned}$$

3490 where above

$$\text{Err}_{\mathbf{u}} := \max_{k \in [K]} \left(\eta \sum_{j=1}^k \beta_{\text{stab}}^{k-j} \|\Delta_{\mathbf{u},j}\|^2 \right)^{1/2}, \quad \text{Err}_{\mathbf{K}} := \max_{k \in [K]} \left(\eta \sum_{j=1}^k \beta_{\text{stab}}^{k-j} \|\Delta_{\mathbf{K},j}\|^2 \right)^{1/2}.$$

3496 We bound the two underlined terms in the above display.

3497 **Claim G.10.** Recall $\text{Err}_{\mathbf{u}} = \max_{k \in [K]} \left(\eta \sum_{j=1}^k \beta_{\text{stab}}^{k-j} \|\Delta_{\mathbf{u},j}\|^2 \right)^{1/2}$. Then, for any k ,

$$\text{Term}_{1,k} \leq \alpha := 2B_{\text{stab}}\text{Err}_{\mathbf{u}} \left(M_{\text{dyn}}\text{Err}_{\mathbf{u}} + \sqrt{L_{\text{stab}}L_{\text{dyn}}} \right)$$

3503 **Claim G.11.** Assume $\beta_{\text{stab}} \in [1/2, 1)$ and recall $\text{Err}_{\mathbf{K}} := \max_{k \in [K]} \left(\eta \sum_{j=1}^k \beta_{\text{stab}}^j \|\Delta_{\mathbf{K},j}\|^2 \right)^{1/2}$. Then,

$$\text{Term}_2 \leq \bar{\varepsilon}_1 \beta_{\text{stab}}^{k/2}, \quad \bar{\varepsilon}_1 := 2B_{\text{stab}} \left(\|\Delta_{\mathbf{x},1}\| + 2L_{\text{stab}}^{1/2}L_{\text{dyn}}\text{Err}_{\mathbf{K}}\|\bar{\Delta}_{\mathbf{x},1}\| \right)$$

3508 The previous two claims and (G.13) show that as soon as (G.12) holds for all indices $1 \leq j \leq k$,

$$\|\Delta_{\mathbf{x},k+1}\| \leq \alpha + \bar{\varepsilon}_1 \beta_{\text{stab}}^{k/2} + \eta \sum_{j=1}^k \beta_{\text{stab}}^{k-j} (C_1\|\Delta_{\mathbf{x},j}\|^2 + C_2\|\Delta_{\mathbf{K},j}\|\|\Delta_{\mathbf{x},j}\|)$$

3514 Set $\varepsilon_j = \|\Delta_{\mathbf{x},j}\|$. Note that $\bar{\varepsilon}_1 \geq \varepsilon_1$, $\beta_{\text{stab}} \in [1/2, 1)$, we can apply Lemma G.19 with $\delta_j \leftarrow \|\Delta_{\mathbf{K},j}\|$ and $M \leftarrow \frac{\eta}{1-\beta} = L_{\text{stab}}$ to find that

$$\|\Delta_{\mathbf{x},k+1}\| = \varepsilon_{k+1} \leq 3(\alpha + \bar{\varepsilon}_1)\beta_{\text{stab}}^{k/3}$$

3519

3520 provided it holds that (we take $L_{\text{stab}} \geq 1, B_{\text{stab}} \geq 1$)

$$\begin{aligned}
3521 \quad & 2B_{\text{stab}}\text{Err}_{\mathbf{u}} \left(M_{\text{dyn}}\text{Err}_{\mathbf{u}} + \sqrt{L_{\text{stab}}}L_{\text{dyn}} \right) = \alpha \leq \frac{1}{18C_1L_{\text{stab}}} = \frac{1}{64B_{\text{stab}}M_{\text{dyn}}L_{\text{stab}}} \\
3522 \quad & \\
3523 \quad & \\
3524 \quad & 2B_{\text{stab}} \left(\|\Delta_{\mathbf{x},1}\| + 2L_{\text{dyn}}L_{\text{stab}}^{1/2}\text{Err}_{\mathbf{K}}\|\bar{\Delta}_{\mathbf{x},1}\| \right) = \bar{\varepsilon}_1 \leq \frac{1}{108C_1L_{\text{stab}}} = \frac{1}{324B_{\text{stab}}M_{\text{dyn}}L_{\text{stab}}} \\
3525 \quad & \\
3526 \quad & \text{Err}_{\mathbf{K}} \leq \frac{1}{12\sqrt{L_{\text{stab}}}\max\{C_2, 1\}} \leq \frac{1}{24\sqrt{L_{\text{stab}}}B_{\text{stab}}L_{\text{dyn}}}.
\end{aligned}$$

3528 For these first two equation, it is enough that

$$\begin{aligned}
3531 \quad & \text{Err}_{\mathbf{u}} \leq \min \left\{ \frac{\sqrt{L_{\text{stab}}}L_{\text{dyn}}}{M_{\text{dyn}}}, \frac{1}{256B_{\text{stab}}^2M_{\text{dyn}}L_{\text{dyn}}L_{\text{stab}}^{3/2}} \right\} \\
3532 \quad & \\
3533 \quad & \\
3534 \quad & \text{Err}_{\mathbf{K}} \leq \frac{1}{24\sqrt{L_{\text{stab}}}B_{\text{stab}}L_{\text{dyn}}} \\
3535 \quad & \\
3536 \quad & \|\Delta_{\mathbf{x},1}\| \leq \frac{1}{4 \cdot 324B_{\text{stab}}^2M_{\text{dyn}}L_{\text{stab}}} \\
3537 \quad & \\
3538 \quad & \text{Err}_{\mathbf{K}}\|\bar{\Delta}_{\mathbf{x},1}\| \leq \frac{L_{\text{dyn}}}{8 \cdot 324B_{\text{stab}}^2M_{\text{dyn}}L_{\text{stab}}^{3/2}}
\end{aligned}$$

3541 for which $\text{Err}_{\mathbf{u}} \leq C_{\mathbf{u}}, \|\Delta_{\mathbf{x},1}\| \leq C_{\Delta}, \|\bar{\Delta}_{\mathbf{x},1}\| \leq C_{\bar{\mathbf{x}}}, \text{Err}_{\mathbf{K}} \leq C_{\mathbf{K}}, \text{Err}_{\mathbf{K}}\|\bar{\Delta}_{\mathbf{x},1}\| \leq C_{\mathbf{K},\bar{\mathbf{x}}}$. Moreover, under the above

3542 condition on $\text{Err}_{\mathbf{u}}$, we have

$$\begin{aligned}
3543 \quad & \|\Delta_{\mathbf{x},k+1}\| \leq 3(\alpha + \bar{\varepsilon}_1)\beta_{\text{stab}}^{k/3} \\
3544 \quad & \\
3545 \quad & \leq 12B_{\text{stab}}\sqrt{L_{\text{stab}}}L_{\text{dyn}}\text{Err}_{\mathbf{u}} + 2B_{\text{stab}} \left(\|\Delta_{\mathbf{x},1}\| + 2L_{\text{stab}}^{1/2}L_{\text{dyn}}\text{Err}_{\mathbf{K}}\|\bar{\Delta}_{\mathbf{x},1}\| \right) \beta_{\text{stab}}^{k/3} \\
3546 \quad & \\
3547 \quad & \leq 12B_{\text{stab}}\sqrt{L_{\text{stab}}}L_{\text{dyn}}\text{Err}_{\mathbf{u}} + 2B_{\text{stab}} \left(\|\Delta_{\mathbf{x},1}\| + 2L_{\text{stab}}^{1/2}L_{\text{dyn}}\text{Err}_{\mathbf{K}}\|\bar{\Delta}_{\mathbf{x},1}\| \right) \beta_{\text{stab}}^{k/3} \\
3548 \quad & \\
3549 \quad & \leq c_{\mathbf{u}}\text{Err}_{\mathbf{u}} + (c_{\mathbf{K}}\text{Err}_{\mathbf{K}}\|\bar{\Delta}_{\mathbf{x},1}\| + c_{\Delta}\|\Delta_{\mathbf{x},1}\|) \beta_{\text{stab}}^{k/3}.
\end{aligned}$$

3550 This in turn shows that the inductive hypothesis (G.12) holds, completing the induction.

3552 G.7.1. DEFERRED CLAIMS

3553 *Proof of Claim G.9.* We argue in steps. Recall also \tilde{R}_{stab} be such that $\tilde{R}_{\text{stab}} \geq \max_k \{\|\mathbf{K}_k\|, \|\mathbf{K}'_k\|, 1\}$.

3555 **Ensuring within radius of regularity.** Our first step is to establish that the maximum of the following three terms is at

3556 most R_{dyn} :

$$\begin{aligned}
3558 \quad & \|(\mathbf{x}_k + \Delta_{\mathbf{x},k}, \bar{\mathbf{u}}_k + \Delta_{\mathbf{u},k} + \mathbf{K}'_k\tilde{\Delta}_{\mathbf{x},k}) - (\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k)\| \\
3559 \quad & \vee \|(\mathbf{x}_k + \Delta_{\mathbf{x},k}, \bar{\mathbf{u}}_k + \Delta_{\mathbf{u},k} + \mathbf{K}'_k\tilde{\Delta}_{\mathbf{x},k}) - (\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k)\| \\
3560 \quad & \vee \|(\mathbf{x}_k, \mathbf{u}_k) - (\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k)\| \leq R_{\text{dyn}}
\end{aligned}$$

3562 First, we observe

$$\begin{aligned}
3563 \quad & \|(\mathbf{x}_k + \Delta_{\mathbf{x},k}, \bar{\mathbf{u}}_k + \Delta_{\mathbf{u},k} + \mathbf{K}'_k\tilde{\Delta}_{\mathbf{x},k}) - (\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k)\| \\
3564 \quad & \leq \|(\mathbf{x}_k + \Delta_{\mathbf{x},k}, \bar{\mathbf{u}}_k + \Delta_{\mathbf{u},k} + \mathbf{K}'_k\tilde{\Delta}_{\mathbf{x},k}) - (\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k)\| + \|\Delta_{\mathbf{K},k}\|\|\tilde{\Delta}_{\mathbf{x},k}\| \\
3565 \quad & \leq \|(\mathbf{x}_k + \Delta_{\mathbf{x},k}, \bar{\mathbf{u}}_k + \Delta_{\mathbf{u},k} + \mathbf{K}_k\tilde{\Delta}_{\mathbf{x},k}) - (\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k)\| + \|\Delta_{\mathbf{K},k}\|\|\Delta_{\mathbf{x},k}\| + \|\Delta_{\mathbf{K},k}\|\|\tilde{\Delta}_{\mathbf{x},k}\| \\
3566 \quad & \leq \|(\mathbf{x}_k + \Delta_{\mathbf{x},k}, \bar{\mathbf{u}}_k + \Delta_{\mathbf{u},k} + \mathbf{K}_k\tilde{\Delta}_{\mathbf{x},k}) - (\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k)\| + \|\Delta_{\mathbf{K},k}\|\|\Delta_{\mathbf{x},k}\| + \|\Delta_{\mathbf{K},k}\|\|\tilde{\Delta}_{\mathbf{x},k}\| + \|\mathbf{K}_k\|\|\Delta_{\mathbf{x},k}\| \\
3567 \quad & \leq \|(\mathbf{x}_k, \bar{\mathbf{u}}_k + \mathbf{K}_k\tilde{\Delta}_{\mathbf{x},k}) - (\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k)\| + (1 + \|\Delta_{\mathbf{K},k}\|)\|\Delta_{\mathbf{x},k}\| + \|\Delta_{\mathbf{K},k}\|\|\tilde{\Delta}_{\mathbf{x},k}\| + \|\mathbf{K}_k\|\|\Delta_{\mathbf{x},k}\| + \|\Delta_{\mathbf{u},k}\| \\
3568 \quad & \leq \|\bar{\Delta}_{\mathbf{x},k}\|(1 + \|\mathbf{K}\|) + (1 + \|\Delta_{\mathbf{K},k}\| + \|\mathbf{K}_k\|)\|\Delta_{\mathbf{x},k}\| + \|\bar{\Delta}_{\mathbf{x},k}\| + \|\Delta_{\mathbf{u},k}\| \\
3569 \quad & \leq (1 + \|\mathbf{K}_k\| + \|\Delta_{\mathbf{K},k}\|)(\|\Delta_{\mathbf{x},k}\| + \|\bar{\Delta}_{\mathbf{x},k}\|) + \|\Delta_{\mathbf{u},k}\| \\
3570 \quad & \leq (2R_{\text{stab}} + \max_j \|\mathbf{K}_k - \mathbf{K}'_j\|)(\|\Delta_{\mathbf{x},k}\| + \|\bar{\Delta}_{\mathbf{x},k}\|) + \|\Delta_{\mathbf{u},k}\|
\end{aligned}$$

3574

3575 Recall the notation

3576

$$3577 \quad \Delta_{\mathbf{K},\infty} := \max_j \|\mathbf{K}_j - \mathbf{K}'_j\|, \quad \Delta_{\mathbf{u},\infty} := \max_j \|\bar{\mathbf{u}}_j - \bar{\mathbf{u}}'_j\|.$$

3578

3579 Hence, it is enough that

3580

$$3581 \quad (2R_{\text{stab}} + \Delta_{\mathbf{K},\infty})(\|\Delta_{\mathbf{x},k}\| + \|\bar{\Delta}_{\mathbf{x},k}\|) + \Delta_{\mathbf{u},\infty} \leq R_{\text{dyn}}.$$

3582

3583 Thus, since $\|\Delta_{\mathbf{x},k}\| \leq c_{\mathbf{u}}\text{Err}_{\mathbf{u}} + (c_{\mathbf{K}}\text{Err}_{\mathbf{K}} - 2B_{\text{stab}})\|\bar{\Delta}_{\mathbf{x},1}\| + c_{\Delta}\|\Delta_{\mathbf{x},1}\|$ due to (G.12) and $\|\bar{\Delta}_{\mathbf{x},k}\| \leq 2B_{\text{stab}}\|\bar{\Delta}_{\mathbf{x},1}\|$ by
 3584 Lemma G.4

3585

$$3586 \quad \|\Delta_{\mathbf{x},k}\| \leq c_{\mathbf{u}}\text{Err}_{\mathbf{u}} + c_{\mathbf{K}}\text{Err}_{\mathbf{K}}\|\bar{\Delta}_{\mathbf{x},1}\| + c_{\Delta}\|\Delta_{\mathbf{x},1}\|$$

3587

3588

3589 Hence, it is enough that

3590

$$3591 \quad R_{\text{dyn}} \geq (2R_{\text{stab}} + \Delta_{\mathbf{K},\infty})(c_{\mathbf{u}}\text{Err}_{\mathbf{u}} + c_{\mathbf{K}}\text{Err}_{\mathbf{K}}\|\bar{\Delta}_{\mathbf{x},1}\| + c_{\Delta}\|\Delta_{\mathbf{x},1}\|) + \Delta_{\mathbf{u},\infty},$$

3592

3593 We can bound $2R_{\text{stab}} + \Delta_{\mathbf{K},\infty} \leq 4\tilde{R}_{\text{stab}}$, and solving the geometric series, bound $\text{Err}_{\mathbf{u}} \leq \sqrt{L_{\text{stab}}}\Delta_{\mathbf{u},\infty}$ and $\text{Err}_{\mathbf{K}} \leq$
 3594 $\sqrt{L_{\text{stab}}}\Delta_{\mathbf{K},\infty} \leq 2\sqrt{L_{\text{stab}}}\tilde{R}_{\text{stab}}$. Thus, it is enough that

3595

$$3596 \quad R_{\text{dyn}} \geq (4\tilde{R}_{\text{stab}}c_{\mathbf{u}}\sqrt{L_{\text{stab}}} + 1)\Delta_{\mathbf{u},\infty} + 4\tilde{R}_{\text{stab}}c_{\mathbf{K}}\|\bar{\Delta}_{\mathbf{x},1}\| + 4\tilde{R}_{\text{stab}}c_{\Delta}\|\Delta_{\mathbf{x},1}\|.$$

3597

3598

3599 which is ensured by Proposition G.5.

3600

3601 **Controlling the first remainder.** Using that the relevant terms are within the radius of regularity,

3602

$$3603 \quad \|\text{rem}_{k,1}\| = \|f_{\eta}(\mathbf{x}_k + \Delta_{\mathbf{x},k}, \bar{\mathbf{u}}_k + \Delta_{\mathbf{u},k} + \mathbf{K}'_k \tilde{\Delta}_{\mathbf{x},k}) - f_{\eta}(\mathbf{x}_k + \Delta_{\mathbf{x},k}, \bar{\mathbf{u}}_k + \Delta_{\mathbf{u},k} + \mathbf{K}_k \tilde{\Delta}_{\mathbf{x},k})\|$$

$$3604 \quad \leq L_{\text{dyn}}\|\mathbf{K}'_k - \mathbf{K}_k\|\tilde{\Delta}_{\mathbf{x},k}\|$$

$$3605 \quad \leq L_{\text{dyn}}\Delta_{\mathbf{K},k}(\|\bar{\Delta}_{\mathbf{x},k}\| + \|\Delta_{\mathbf{x},k}\|).$$

3606

3607

3608 **Controlling the first remainder.** Using the definitions of $\mathbf{x}_k = \bar{\mathbf{x}}_k + \bar{\Delta}_{\mathbf{x},k}$ $\mathbf{u}_k = \bar{\mathbf{u}}_k + \mathbf{K}_k \bar{\Delta}_{\mathbf{x},k}$, and the fact that $(\mathbf{x}_k, \mathbf{u}_k)$
 3609 is in the radius of regularity around $(\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k)$, a Taylor expansion implies

3610

$$3611 \quad \|\text{rem}_{k,2}\| = \left\| \begin{aligned} & f_{\eta}(\mathbf{x}_k + \Delta_{\mathbf{x},k}, \bar{\mathbf{u}}_k + \Delta_{\mathbf{u},k} + \mathbf{K}_k \tilde{\Delta}_{\mathbf{x},k}) - f(\mathbf{x}_k, \bar{\mathbf{u}}_k + \mathbf{K}_k \bar{\Delta}_{\mathbf{x},k}) \\ & - \frac{\partial}{\partial \mathbf{x}} f_{\eta}(\mathbf{x}_k, \mathbf{u}_k) \Delta_{\mathbf{x},k} + \frac{\partial}{\partial \mathbf{u}} f_{\eta}(\mathbf{x}_k, \mathbf{u}_k) (\Delta_{\mathbf{u},k} + \mathbf{K}_k (\tilde{\Delta}_{\mathbf{x},k} - \bar{\Delta}_{\mathbf{x},k})) \end{aligned} \right\|$$

$$3612 \quad = \left\| \begin{aligned} & f_{\eta}(\mathbf{x}_k + \Delta_{\mathbf{x},k}, \bar{\mathbf{u}}_k + \Delta_{\mathbf{u},k} + \mathbf{K}_k \tilde{\Delta}_{\mathbf{x},k}) - f(\mathbf{x}_k, \mathbf{u}_k) \\ & - \frac{\partial}{\partial \mathbf{x}} f_{\eta}(\mathbf{x}_k, \mathbf{u}_k) \Delta_{\mathbf{x},k} + \frac{\partial}{\partial \mathbf{u}} f_{\eta}(\mathbf{x}_k, \mathbf{u}_k) (\bar{\mathbf{u}}_k + \Delta_{\mathbf{u},k} + \mathbf{K}_k \tilde{\Delta}_{\mathbf{x},k} - \mathbf{u}_k) \end{aligned} \right\|$$

$$3613 \quad \leq \frac{M_{\text{dyn}}}{2} \left(\|\Delta_{\mathbf{x},k}\|^2 + \|\bar{\mathbf{u}}_k + \Delta_{\mathbf{u},k} + \mathbf{K}_k \tilde{\Delta}_{\mathbf{x},k} - \mathbf{u}_k\|^2 \right)$$

$$3614 \quad = \frac{M_{\text{dyn}}}{2} \left(\|\Delta_{\mathbf{x},k}\|^2 + \|\Delta_{\mathbf{u},k} + \mathbf{K}_k (\tilde{\Delta}_{\mathbf{x},k} - \bar{\Delta}_{\mathbf{x},k})\|^2 \right)$$

$$3615 \quad = \frac{M_{\text{dyn}}}{2} \left(\|\Delta_{\mathbf{x},k}\|^2 + \|\Delta_{\mathbf{u},k} + \mathbf{K}_k \Delta_{\mathbf{x},k}\|^2 \right)$$

$$3616 \quad = \frac{M_{\text{dyn}}}{2} \left((1 + 2\|\mathbf{K}_k\|^2) \|\Delta_{\mathbf{x},k}\|^2 + 2\|\Delta_{\mathbf{u},k}\|^2 \right)$$

$$3617 \quad = \frac{3}{2} M_{\text{dyn}} R_{\text{stab}}^2 \|\Delta_{\mathbf{x},k}\|^2 + M_{\text{dyn}} \|\Delta_{\mathbf{u},k}\|^2$$

3618

3619

3620

3621

3622

3623

3624

3625

□

3630 *Proof of Claim G.10.* Recall $\text{Err}_{\mathbf{u}} = \max_{k \in [K]} \left(\eta \sum_{j=1}^k \beta_{\text{stab}}^{k-j} \|\Delta_{\mathbf{u},j}\|^2 \right)^{1/2}$ and $\beta_{\text{stab}} = 1 - \frac{\eta}{L_{\text{stab}}}$. Then,

$$\begin{aligned}
3631 & \\
3632 & \\
3633 & \\
3634 & \text{Term}_{1,k} = \eta \sum_{j=1}^k 2B_{\text{stab}} \beta_{\text{stab}}^{k-j} (L_{\text{dyn}} \|\Delta_{\mathbf{u},j}\| + M_{\text{dyn}} \|\Delta_{\mathbf{u},j}\|^2) \\
3635 & \\
3636 & \leq 2B_{\text{stab}} \left(M_{\text{dyn}} \text{Err}_{\mathbf{u}}^2 + L_{\text{dyn}} \cdot \eta \sum_{j=1}^k \beta_{\text{stab}}^{k-j} \|\Delta_{\mathbf{u},j}\| \right) \\
3637 & \\
3638 & \leq 2B_{\text{stab}} \left(M_{\text{dyn}} \text{Err}_{\mathbf{u}}^2 + L_{\text{dyn}} \cdot \left(\eta \sum_{j=1}^k \beta_{\text{stab}}^{k-j} \right)^{1/2} \left(\eta \sum_{j=1}^k \beta_{\text{stab}}^{k-j} \|\Delta_{\mathbf{u},j}\|^2 \right)^{1/2} \right) \\
3639 & \\
3640 & \leq 2B_{\text{stab}} \left(M_{\text{dyn}} \text{Err}_{\mathbf{u}}^2 + L_{\text{dyn}} \cdot \left(\eta \sum_{j=1}^k \beta_{\text{stab}}^{k-j} \right)^{1/2} \text{Err}_{\mathbf{u}} \right) \\
3641 & \\
3642 & \leq 2B_{\text{stab}} \left(M_{\text{dyn}} \text{Err}_{\mathbf{u}}^2 + L_{\text{dyn}} \cdot \underbrace{\left(\eta \frac{1}{\beta_{\text{stab}}^{-1} - 1} \right)^{1/2}}_{=L_{\text{stab}}} \text{Err}_{\mathbf{u}} \right) \\
3643 & \\
3644 & = 2B_{\text{stab}} \text{Err}_{\mathbf{u}} \left(M_{\text{dyn}} \text{Err}_{\mathbf{u}} + \sqrt{L_{\text{stab}}} L_{\text{dyn}} \right) \\
3645 & \\
3646 & \\
3647 & \\
3648 & \\
3649 & \\
3650 & \\
3651 & \\
3652 & \\
3653 & \square
\end{aligned}$$

3654
3655
3656 *Proof of Claim G.11.*

$$\begin{aligned}
3657 & \\
3658 & \\
3659 & L_{\text{dyn}}^{-1} (\text{Term}_{2,k} - 2B_{\text{stab}} \beta_{\text{stab}}^k \|\Delta_{\mathbf{x},1}\|) = \eta \sum_{j=1}^k 2B_{\text{stab}} \beta_{\text{stab}}^{k-j} \|\Delta_{\mathbf{K},j}\| \|\bar{\Delta}_{\mathbf{x},j}\| \\
3660 & \\
3661 & = \eta \sum_{j=1}^k 2B_{\text{stab}} \beta_{\text{stab}}^{k-j} \beta_{\text{stab}}^{j-1} \|\Delta_{\mathbf{K},j}\| \|\bar{\Delta}_{\mathbf{x},1}\| \\
3662 & \\
3663 & = 2B_{\text{stab}} \|\bar{\Delta}_{\mathbf{x},1}\| \cdot \eta \sum_{j=1}^k \beta_{\text{stab}}^{k-1} \|\Delta_{\mathbf{K},j}\| \\
3664 & \\
3665 & = 2B_{\text{stab}} \beta_{\text{stab}}^{\frac{k}{2}-1} \|\bar{\Delta}_{\mathbf{x},1}\| \cdot \eta \sum_{j=1}^k \beta_{\text{stab}}^{(k-j)/2} \beta_{\text{stab}}^{j/2} \|\Delta_{\mathbf{K},j}\| \\
3666 & \\
3667 & \leq 2B_{\text{stab}} \beta_{\text{stab}}^{\frac{k}{2}-1} \|\bar{\Delta}_{\mathbf{x},1}\| \leq \left(\eta \sum_{j=1}^k \beta_{\text{stab}}^j \right) \left(\eta \sum_{j=1}^k \beta_{\text{stab}}^j \|\Delta_{\mathbf{K},j}\|^2 \right)^{1/2} \\
3668 & \\
3669 & \leq 2B_{\text{stab}} \beta_{\text{stab}}^{\frac{k}{2}-1} \|\bar{\Delta}_{\mathbf{x},1}\| \leq \left(\eta \sum_{j=1}^k \beta_{\text{stab}}^j \right)^{1/2} \left(\eta \sum_{j=1}^k \beta_{\text{stab}}^j \|\Delta_{\mathbf{K},j}\|^2 \right)^{1/2} \\
3670 & \\
3671 & \leq 2B_{\text{stab}} L_{\text{stab}}^{1/2} \beta_{\text{stab}}^{\frac{k}{2}-1} \|\bar{\Delta}_{\mathbf{x},1}\| \cdot \underbrace{\left(\eta \sum_{j=1}^k \beta_{\text{stab}}^j \|\Delta_{\mathbf{K},j}\|^2 \right)^{1/2}}_{=\text{Err}_{\mathbf{K}}} \\
3672 & \\
3673 & \leq 4B_{\text{stab}} L_{\text{stab}}^{1/2} \beta_{\text{stab}}^{\frac{k}{2}} \text{Err}_{\mathbf{K}} \|\bar{\Delta}_{\mathbf{x},1}\|. \quad (\beta_{\text{stab}} \geq 1/2) \\
3674 & \\
3675 & \\
3676 & \\
3677 & \\
3678 & \\
3679 & \\
3680 & \\
3681 & \\
3682 & \\
3683 & \\
3684 &
\end{aligned}$$

3685 Thus,

$$\begin{aligned}
3686 \text{Term}_{2,k} &\leq 4B_{\text{stab}}L_{\text{dyn}}L_{\text{stab}}^{1/2}\beta_{\text{stab}}^{\frac{k}{2}}\text{Err}_{\mathbf{K}}\|\bar{\Delta}_{\mathbf{x},1}\| + 2B_{\text{stab}}\beta_{\text{stab}}^k\|\Delta_{\mathbf{x},1}\| \\
3687 &\leq \beta_{\text{stab}}^{\frac{k}{2}}\left(4L_{\text{dyn}}B_{\text{stab}}L_{\text{stab}}^{1/2}\text{Err}_{\mathbf{K}}\|\bar{\Delta}_{\mathbf{x},1}\| + 2B_{\text{stab}}\|\Delta_{\mathbf{x},1}\|\right)
\end{aligned}$$

3690 □

3692 G.8. Ricatti synthesis of stabilizing gains.

3694 In this section, we show that under a certain *stabilizability* condition, it is always possible to synthesize primitive controllers
3695 satisfying [Assumption 3.2](#) with reasonable constants. We begin by defining our notion of stabilizability; we adopt the
3696 formulation based on Jacobian linearizations of non-linear systems the discrete analogue of the senses proposed in which is
3697 consistent with ([Pfrommer et al., 2023](#); [Westenbroek et al., 2021](#)).

3698 **Definition G.11** (Stabilizability). A control trajectory $\rho = (\mathbf{x}_{1:K+1}, \mathbf{u}_{1:K}) \in \mathcal{P}_K$ is $L_{\mathcal{V}}$ -Jacobian-Stabilizable if
3699 $\max_k \mathcal{V}_k(\rho) \leq L_{\mathcal{V}}$, where for $k \in [K+1]$, $\mathcal{V}_k(\rho)$ is defined by

$$\begin{aligned}
3700 \mathcal{V}_k(\rho) &:= \sup_{\xi: \|\xi\| \leq 1} \left(\inf_{\tilde{\mathbf{u}}_{1:s}} \|\tilde{\mathbf{x}}_{K+1}\|^2 + \eta \sum_{j=k}^K \|\tilde{\mathbf{x}}_j\|^2 + \|\tilde{\mathbf{u}}_j\|^2 \right) \\
3701 &\text{s.t. } \tilde{\mathbf{x}}_k = \xi, \quad \tilde{\mathbf{x}}_{j+1} = \tilde{\mathbf{x}}_j + \eta(\mathbf{A}_j(\rho)\tilde{\mathbf{x}}_j + \mathbf{B}_j(\rho)\tilde{\mathbf{u}}_j),
\end{aligned}$$

3706 Here, for simplicity, we use Euclidean-norm costs, though any Mahalanobis-norm cost induced by a positive definite matrix
3707 would suffice. We propose to synthesize gain matrices by performing a standard Ricatti update, normalized appropriately to
3708 take account of the step size $\eta > 0$ (see, e.g. Appendix F in ([Pfrommer et al., 2023](#))).

3709 **Definition G.12** (Ricatti update). Given a path $\rho \in \mathcal{P}_k$ with $\mathbf{A}_k = \mathbf{A}_k(\rho)$, $\mathbf{B}_k = \mathbf{B}_k(\rho)$ we define

$$\begin{aligned}
3710 \mathbf{P}_{K+1}^{\text{ric}}(\rho) &= \mathbf{I}, \quad \mathbf{P}_k^{\text{ric}}(\rho) = (\mathbf{I} + \eta\mathbf{A}_{\text{cl},k}(\rho))^\top \mathbf{P}_{k+1}^{\text{ric}}(\rho)(\mathbf{I} + \eta\mathbf{A}_{\text{cl},k}(\rho)) + \eta(\mathbf{I} + \mathbf{K}_k(\rho)\mathbf{K}_k(\rho)^\top) \\
3711 \mathbf{K}_k^{\text{ric}}(\rho) &= (\mathbf{I} + \eta\mathbf{B}_k^\top \mathbf{P}_{k+1}^{\text{ric}}(\rho)\mathbf{B}_k)^{-1}(\mathbf{B}_k^\top \mathbf{P}_{k+1}^{\text{ric}}(\rho))(\mathbf{I} + \eta\mathbf{A}_k) \\
3712 \mathbf{A}_{\text{cl},k}^{\text{ric}}(\rho) &= \mathbf{A}_k + \mathbf{B}_k\mathbf{K}_k(\rho).
\end{aligned}$$

3716 The main result of this section is that the parameters $(R_{\text{stab}}, B_{\text{stab}}, L_{\text{stab}})$ in [Assumption 3.2](#) can be bounded in terms of
3717 L_{dyn} in [Assumption 3.1](#), and the bound $L_{\mathcal{V}}$ defined above.

3718 **Proposition G.12** (Instantiating the Lyapunov Lemma). Let $L_{\text{dyn}}, L_{\mathcal{V}} \geq 1$, and let $\rho = (\mathbf{x}_{1:K+1}, \mathbf{u}_{1:K})$ be
3719 $(R_{\text{dyn}}, L_{\text{dyn}}, M_{\text{dyn}})$ -regular and $L_{\mathcal{V}}$ -Jacobian Stabilizable. Suppose further that $\eta \leq 1/5L_{\text{dyn}}^2L_{\mathcal{V}}$. Then, $(\rho, \mathbf{K}_{1:K}^{\text{ric}})$ -
3720 is $(R_{\text{stab}}, B_{\text{stab}}, L_{\text{stab}})$ -Jacobian Stable, where

$$3721 R_{\text{stab}} = \frac{4}{3}L_{\mathcal{V}}L_{\text{dyn}}, \quad B_{\text{stab}} = \sqrt{5}L_{\text{dyn}}L_{\mathcal{V}}, \quad L_{\text{stab}} = 2L_{\mathcal{V}}$$

3725 [Proposition G.12](#) is proven in [Appendix G.8.1](#) below. A consequence of the above proposition is that, given access to a
3726 smooth local model of dynamics, one can implement the synthesis oracle by computing linearizations around demonstrated
3727 trajectories, and solving the corresponding Ricatti equations as per the above discussions to synthesize the correct gains.

3729 G.8.1. PROOF OF [PROPOSITION G.12](#) (RICATTI SYNTHESIS OF GAINS)

3730 Throughout, we use the shorthand $\mathbf{A}_k = \mathbf{A}_k(\rho)$ and $\mathbf{B}_k = \mathbf{B}_k(\rho)$, recall that $\|\cdot\|$ denotes the operator norm when applied
3731 to matrices. We also recall our assumptions that $L_{\text{dyn}}, L_{\mathcal{V}} \geq 1$. We begin by translating our stabilizability assumption
3732 ([Definition G.11](#)) into the the \mathbf{P} -matrices in [Definition G.12](#). The following statement recalls Lemma F.1 in ([Pfrommer](#)
3733 [et al., 2023](#)), an instantiation of well-known solutions to linear-quadratic dynamic programming (e.g. ([Anderson & Moore,](#)
3734 [2007](#))).

3735 **Lemma G.13** (Equivalence of stabilizability and Ricatti matrices). Consider a trajectory $(\mathbf{x}_{1:K}, \mathbf{u}_{1:K})$, and define the
3736 parameter $\Theta := (\mathbf{A}_{\text{jac}}(\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k), \mathbf{B}_{\text{jac}}(\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k))_{k \in [K]}$. Then, for all $k \in [K]$,

$$3737 \forall k \in [K], \quad \mathcal{V}_k(\rho) = \|\mathbf{P}_k(\Theta)\|_{\text{op}}$$

3739

3740 Hence, if ρ is $L_{\mathcal{V}}$ -stabilizable,

$$3741 \max_{k \in [K+1]} \|\mathbf{P}_k(\Theta)\|_{\text{op}} \leq L_{\mathcal{V}}.$$

3742 **Lemma G.14** (Lyapunov Lemma, Lemma F.10 in (Pfrommer et al., 2023)). Let $\mathbf{X}_{1:K}, \mathbf{Y}_{1:K}$ be matrices of appropriate
3743 dimension, and let $\mathbf{Q} \succeq \mathbf{I}$ and $\mathbf{Y}_k \succeq 0$. Define $\Lambda_{1:K+1}$ as the solution of the recursion

$$3744 \Lambda_{K+1} = \mathbf{Q}, \quad \Lambda_k = \mathbf{X}_k^{\top} \Lambda_{k+1} \mathbf{X}_k + \eta \mathbf{Q} + \mathbf{Y}_k$$

3745 Define the operator $\Phi_{j+1,k} = \mathbf{X}_j \cdot \mathbf{X}_{j-1} \cdot \dots \cdot \mathbf{X}_k$, with the convention $\Phi_{k,k} = \mathbf{I}$. Then, if $\max_k \|\mathbf{I} - \mathbf{X}_k\|_{\text{op}} \leq \kappa \eta$ for
3746 some $\kappa \leq 1/2\eta$,

$$3747 \|\Phi_{j,k}\|^2 \leq \max\{1, 2\kappa\} \max_{k \in [K+1]} \|\Lambda_k\| (1 - \eta \alpha)^{j-k}, \quad \alpha := \frac{1}{\max_{k \in [K+1]} \|\Lambda_{1:K+1}\|}.$$

3748 **Claim G.15.** If ρ is $(0, L_{\text{dyn}}, \infty)$ -regular, then for all k , $\mathbf{A}_k = \mathbf{A}_k(\rho)$ and $\mathbf{B}_k = \mathbf{B}_k(\rho)$ satisfy
3749 $\max_{k \in [K]} \max\{\|\mathbf{A}_k\|, \|\mathbf{B}_k\|\} \leq L_{\text{dyn}}$.

3750 *Proof.* For any $k \in [K]$,

$$3751 \max\{\|\mathbf{A}_k\|, \|\mathbf{B}_k\|\} = \max \left\{ \left\| \frac{\partial}{\partial x} f(\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k) \right\|, \left\| \frac{\partial}{\partial u} f(\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k) \right\| \right\} \leq \|\nabla f(\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k)\| \leq L_{\text{dyn}},$$

3752 where the last inequality follows by regularity. □

3753 **Claim G.16.** Recall $\mathbf{K}_k^{\text{ric}}(\rho) = (\mathbf{I} + \eta \mathbf{B}_k^{\top} \mathbf{P}_{k+1}^{\text{ric}}(\rho) \mathbf{B}_k)^{-1} (\mathbf{B}_k^{\top} \mathbf{P}_{k+1}^{\text{ric}}(\rho)) (\mathbf{I} + \eta \mathbf{A}_k)$. Then, if ρ is $L_{\mathcal{V}}$ -stabilizable and
3754 $(0, L_{\text{dyn}}, \infty)$ -regular, and if $\eta \leq 1/3L_{\text{dyn}}$,

$$3755 \|\mathbf{K}_k^{\text{ric}}(\rho)\| \leq \frac{4}{3} L_{\mathcal{V}} L_{\text{dyn}}$$

3756 *Proof.* We bound

$$\begin{aligned} 3757 \|\mathbf{K}_k^{\text{ric}}(\rho)\| &\leq \|\mathbf{B}_k\| \|\mathbf{P}_{k+1}^{\text{ric}}(\rho)\| (1 + \eta \|\mathbf{A}_k\|) \\ 3758 &\leq L_{\text{dyn}} (1 + \eta L_{\text{dyn}}) \|\mathbf{P}_{k+1}^{\text{ric}}(\rho)\| && \text{(Claim G.15)} \\ 3759 &\leq L_{\mathcal{V}} L_{\text{dyn}} (1 + \eta L_{\text{dyn}}) && \text{(Lemma G.13, } L_{\mathcal{V}} \geq 1) \\ 3760 &\leq \frac{4}{3} L_{\mathcal{V}} L_{\text{dyn}} && (\eta \leq 1/3L_{\text{dyn}}). \end{aligned}$$

3761 □

3762 *Proof of Proposition G.12.* We want to show that $\mathbf{K}_{1:K}^{\text{ric}}(\rho)$ is $(R_{\text{stab}}, B_{\text{stab}}, L_{\text{stab}})$ -stabilizing. Claim G.16 has already
3763 established that $\max_{k \in [K]} \|\mathbf{K}_k^{\text{ric}}(\rho)\| \leq R_{\text{stab}} = \frac{4}{3} L_{\mathcal{V}} L_{\text{dyn}}$.

3764 To prove the other conditions, we apply Lemma G.14 with $\mathbf{Y}_k = \mathbf{K}_k(\Theta) \mathbf{K}_k(\Theta)$, $\mathbf{Q} = \mathbf{I}$, and $\mathbf{X}_k = \mathbf{I} + \eta \mathbf{A}_{\text{cl},k}(\Theta)$. From
3765 Definition G.12, let have that the term Λ_k in Lemma G.14 is precise equal to $\mathbf{P}_k(\Theta)$. From Lemma G.13,

$$3766 \max_{k \in [K+1]} \|\mathbf{P}_k(\Theta)\|_{\text{op}} = \max_{k \in [K+1]} \mathcal{V}_k(\rho) \leq L_{\mathcal{V}}.$$

3767 This implies that if $\max_k \|\mathbf{X}_k - \mathbf{I}\| \leq \kappa \eta \leq 1/2$, we have

$$3768 \|\Phi_{\text{cl},j,k}(\Theta)\|^2 = \|(\mathbf{X}_j \cdot \mathbf{X}_{j-1} \cdot \dots \cdot \mathbf{X}_k)\| \leq \max\{1, 2\kappa\} L_{\mathcal{V}} \left(1 - \frac{\eta}{L_{\mathcal{V}}}\right)^{j-k}.$$

3769

3795 It suffices to find an appropriate upper bound κ . We have

$$\begin{aligned}
3796 \quad \|\mathbf{X}_k - \mathbf{I}\| &= \|\eta \mathbf{A}_{\text{cl},k}(\Theta)\| \leq \eta(\|\mathbf{A}_k\| + \|\mathbf{B}_k\| \|\mathbf{K}_k(\Theta)\|) \\
3797 &\leq \eta L_{\text{dyn}}(1 + \|\mathbf{K}_k(\Theta)\|) \\
3798 &\leq \eta L_{\text{dyn}}\left(1 + \frac{4}{3}L_{\text{dyn}}L_{\mathcal{V}}\right) \quad (\text{Claim G.16}) \\
3800 &\leq \frac{7}{3}\eta L_{\text{dyn}}^2 L_{\mathcal{V}} \quad (L_{\mathcal{V}}, L_{\text{dyn}} \geq 1)
\end{aligned}$$

3804 Setting $\kappa = \frac{7}{3}L_{\text{dyn}}^2 L_{\mathcal{V}}$, we have that as $\eta \leq \frac{1}{5L_{\text{dyn}}^2 L_{\mathcal{V}}} \leq \min\{\frac{3}{14L_{\text{dyn}}^2 L_{\mathcal{V}}}, \frac{1}{3L_{\text{dyn}}}\}$ (recall $L_{\text{dyn}}, L_{\mathcal{V}} \geq 1$), we can bound

$$3807 \quad \max\{1, 2\kappa\} \leq \max\left\{1, \frac{14}{3}L_{\text{dyn}}^2 L_{\mathcal{V}}\right\} \leq \max\{1, 5L_{\text{dyn}}^2 L_{\mathcal{V}}\} = 5L_{\text{dyn}}^2 L_{\mathcal{V}}^2,$$

3809 where again recall $L_{\mathcal{V}}, L_{\text{dyn}} \geq 1$. In sum, for $\eta \leq \frac{1}{5L_{\text{dyn}}^2 L_{\mathcal{V}}}$, we have

$$3812 \quad \|\Phi_{\text{cl},j,k}\|^2 \leq 5L_{\text{dyn}}^2 L_{\mathcal{V}}^2 \left(1 - \frac{\eta}{L_{\mathcal{V}}}\right)^{j-k}.$$

3815 Hence, using the elementary inequality $\sqrt{1-a} \leq (1-a/2)$,

$$3818 \quad \|\Phi_{\text{cl},j,k}\| \leq \sqrt{5}L_{\text{dyn}}L_{\mathcal{V}} \left(1 - \frac{\eta}{L_{\mathcal{V}}}\right)^{(j-k)/2} \leq \sqrt{5}L_{\text{dyn}}L_{\mathcal{V}} \left(1 - \frac{\eta}{2L_{\mathcal{V}}}\right)^{j-k},$$

3820 which shows that we can select $B_{\text{stab}} = \sqrt{5}L_{\text{dyn}}L_{\mathcal{V}}$ and $L_{\text{stab}} = 2L_{\mathcal{V}}$. □

3822 G.9. Solutions to recursions

3824 This section contains the solutions to various recursions used in the proof of the two results in [Appendix G.4: Proposition G.5](#) (whose proof is given in [Appendix G.7](#)) and [Lemma G.4](#) (whose proof is given in [Appendix G.6](#)).

3827 **Lemma G.17** (First Key Recursion). *Let $C_1 > 0, C_2 \geq 1/2, \beta_{\text{stab}} \in (0, 1)$, and suppose $\varepsilon_1, \varepsilon_2, \dots$ is a sequence satisfying $\varepsilon_1 \leq \bar{\varepsilon}_1$, and*

$$3830 \quad \varepsilon_{k+1} \leq C_2 \beta_{\text{stab}}^k \bar{\varepsilon}_1 + C_1 \eta \sum_{j=1}^k \beta_{\text{stab}}^{k-j} \varepsilon_j^2$$

3833 *Then, as long as $C_1 \leq \beta(1-\beta)/2\eta$, it holds that $\varepsilon_k \leq 2C_2 \beta_{\text{stab}}^{k-1} \bar{\varepsilon}_1$ for all k .*

3836 *Proof.* Consider the sequence $\nu_k = 2C_2 \beta_{\text{stab}}^{k-1} \bar{\varepsilon}_1$. Since $C_2 \geq 1/2$, we have $\nu_1 \geq \bar{\varepsilon}_1 \geq \varepsilon_1$. Moreover,

$$\begin{aligned}
3838 \quad C_2 \beta_{\text{stab}}^k \bar{\varepsilon}_1 + C_1 \sum_{j=1}^k \beta_{\text{stab}}^{k-j} \nu_j &= C_2 \beta_{\text{stab}}^k \bar{\varepsilon}_1 + 2C_1 C_2 \sum_{j=1}^k \beta_{\text{stab}}^{k+j-2} \bar{\varepsilon}_1 \\
3840 &= C_2 \beta_{\text{stab}}^k \bar{\varepsilon}_1 \left(1 + \frac{2C_1}{\beta} \sum_{j=0}^{k-1} \beta_{\text{stab}}^j\right) \\
3842 &\leq C_2 \beta_{\text{stab}}^k \bar{\varepsilon}_1 \left(1 + \frac{2C_1 \eta}{\beta(1-\beta)}\right)
\end{aligned}$$

3847 Hence, for $C_1 \leq \beta(1-\beta)/2\eta$, we have $C_2 \beta_{\text{stab}}^k \bar{\varepsilon}_1 + C_1 \sum_{j=1}^k \beta_{\text{stab}}^{k-j} \nu_j \leq 2C_2 \beta_{\text{stab}}^{k-1} \bar{\varepsilon}_1 \leq \nu_{k+1}$. This shows that the (ν_k) sequence dominates the (ε_k) sequence, as needed. □

3850 **Lemma G.18** (Second Key Recursion). *Let $c, \Delta, \eta > 0$, $\beta_{\text{stab}} \in (0, 1)$ and let $\varepsilon_1, \varepsilon_2, \dots$ satisfy $\varepsilon_1 \leq c$ and*

$$3851 \varepsilon_{k+1} \leq c\beta_{\text{stab}}^k + c\eta\Delta\beta_{\text{stab}}^{k-1} \sum_{j=1}^k \varepsilon_j.$$

3852
3853
3854
3855 *Then, if $\Delta \leq \frac{\beta(1-\beta)}{2c\eta}$, $\varepsilon_{k+1} \leq 2c\beta_{\text{stab}}^k$ for all k .*

3856
3857 *Proof.* Consider the sequence $\nu_k = 2c\beta_{\text{stab}}^{k-1}$. Since $\varepsilon_1 \leq c$, $\nu_1 \geq \varepsilon_1$. Moreover,

$$3858 c\beta_{\text{stab}}^k + c\eta\Delta\beta_{\text{stab}}^{k-1} \sum_{j=1}^k \nu_j \leq c\beta_{\text{stab}}^k + 2c^2\eta\Delta\beta_{\text{stab}}^{k-1} \sum_{j=1}^k \beta_{\text{stab}}^{j-1}$$

$$3859 \leq c\beta_{\text{stab}}^k + 2c^2\eta\Delta\beta_{\text{stab}}^{k-1} \frac{1}{1-\beta}$$

$$3860 \leq c\beta_{\text{stab}}^k \left(1 + 2c\Delta \frac{\eta}{\beta(1-\beta)} \right).$$

3861
3862
3863
3864
3865
3866
3867 Hence, for $\Delta \leq \frac{\beta(1-\beta)}{2c\eta}$, the above is at most $2c\beta_{\text{stab}}^k \leq \nu_{k+1}$. This shows that the (ν_k) sequence dominates the (ε_k) sequence, as needed. \square

3868
3869 **Lemma G.19** (Third Key Recursion). *Let $C_1, C_2 > 0$, $\alpha \geq 0$, $\beta_{\text{stab}} \in (1/2, 1)$, and let $\varepsilon_1, \varepsilon_2, \dots$, and $\delta_1, \delta_2, \dots$, and $\bar{\varepsilon}_1 \geq \varepsilon_1$ and be a sequence of real numbers satisfying*

$$3870 \varepsilon_{k+1} \leq \alpha + \eta \sum_{j=1}^k \beta_{\text{stab}}^{k-j} (C_1 \varepsilon_j^2 + C_2 \varepsilon_j \delta_j) + \beta_{\text{stab}}^{k/3} \bar{\varepsilon}_1$$

3871
3872 *Defin, $\text{Err}_\delta := \max_k \eta \sum_{j=1}^k \beta_{\text{stab}}^{(k-j)} \delta_j^2$ and $M = \eta/(1-\beta)$. Then, as long as*

$$3873 \alpha \leq \frac{1}{18C_1M}, \quad \bar{\varepsilon}_1 \leq \frac{1}{108C_1M}, \quad \text{Err}_\delta \leq \frac{1}{12\sqrt{M} \max\{C_2, 1\}},$$

3874
3875
3876 *the following holds for all $k \geq 0$:*

$$3877 \varepsilon_{k+1} \leq 3\alpha + 3\bar{\varepsilon}_1 \beta_{\text{stab}}^{k/3}.$$

3878
3879 *Proof of Lemma G.19.* Consider a sequence

$$3880 \nu_{k+1} = \alpha_* + c_* \beta_*^k \bar{\varepsilon}_1, \quad \alpha_* = 3\alpha, c_* = 3, \beta_* = \beta_{\text{stab}}^{1/3}$$

3881 defined for $k \geq 0$, for some $\alpha_* \geq \alpha$, $\beta_* \in (\beta, 1)$, and $c_* \geq 1$. Then, $\nu_1 \geq \bar{\varepsilon}_1, \geq \varepsilon_1$. Let us define the term B_k via

$$3882 B_k = \alpha + \eta \sum_{j=1}^k \beta_{\text{stab}}^{k-j} (C_1 \nu_j^2 + C_2 \nu_j \delta_j) + \beta_{\text{stab}}^{k/3} \bar{\varepsilon}_1.$$

3883
3884
3885
3886 It suffices to show $B_k \leq \nu_{k+1}$ for all k . Introduce $\text{Term}_{\nu,k} = \left(\eta \sum_{j=1}^k \beta_{\text{stab}}^{k-j} \nu_j^2 \right)^{1/2}$ and $\text{Err}_\delta = \max_k \left(\eta \sum_{j=1}^k \beta_{\text{stab}}^{k-j} \delta_j^2 \right)^{1/2}$. Then, by Cauch-Schwartz,

$$3887 B_k = \alpha + \eta \sum_{j=1}^k \beta_{\text{stab}}^{k-j} (C_1 \nu_j^2 + C_2 \nu_j \delta_j) + \beta_{\text{stab}}^{k/3} \bar{\varepsilon}_1$$

$$3888 \leq \alpha + C_1 \text{Term}_{\nu,k}^2 + C_2 \text{Term}_{\nu,k} \text{Err}_\delta + \beta_{\text{stab}}^{k/3} \bar{\varepsilon}_1.$$

3889
3890
3891
3892
3893
3894
3895
3896
3897
3898
3899
3900
3901
3902
3903
3904

3905 We now bound

3906

3907

3908

3909

3910

3911

3912

3913

3914

3915

3916

3917

3918

$$\begin{aligned}
\text{Term}_{\nu,k}^2 &= \eta \sum_{j=1}^k \beta_{\text{stab}}^{k-j} \nu_j^2 \\
&= \eta \sum_{j=1}^k \beta_{\text{stab}}^{k-j} (\alpha_* + c_* \bar{\varepsilon}_1 \beta_*^{j-1})^2 \\
&\leq 2\eta \sum_{j=1}^k \beta_{\text{stab}}^{k-j} \alpha_*^2 + 2\eta c_*^2 \bar{\varepsilon}_1^2 \sum_{j=1}^k \beta_{\text{stab}}^{k-j} \beta_*^{2(j-1)} \\
&\leq \frac{2\eta \alpha_*^2}{1-\beta} + 2\eta c_*^2 \bar{\varepsilon}_1^2 \sum_{j=1}^k \beta_{\text{stab}}^{k-j} \beta_*^{2(j-1)}.
\end{aligned}$$

3919 Now, recalling $\beta_* = \beta_{\text{stab}}^{1/3}$, we have

3920

3921

3922

3923

3924

3925

3926

3927

3928

3929

3930

3931

$$\begin{aligned}
\sum_{j=1}^k \beta_{\text{stab}}^{k-j} \beta_*^{j-1} &= \sum_{j=1}^k \beta_*^{3k-3j} \beta_*^{2(j-1)} = \sum_{j=1}^k \beta_*^{3k-j-2} \\
&= \beta_*^{2k-2} \sum_{j=1}^k \beta_*^{k-j} = \beta_*^{2k-2} \sum_{j \geq 0} \beta_*^j \\
&\leq 3\beta_*^{2k-2} \sum_{j \geq 0} \beta_*^{3j} = 3\beta_*^{2k} \beta_{\text{stab}}^{-2/3} \sum_{j \geq 0} \beta_{\text{stab}}^j \\
&= \frac{3}{1-\beta} \beta_*^{2k} \beta_{\text{stab}}^{-2/3} \leq \frac{3}{\beta(1-\beta)} \beta_*^{2k}.
\end{aligned}$$

3932 Thus, adopting the shorthand $M = \eta/(1-\beta)$, and using the assumption $\beta_{\text{stab}} \geq 1/2$,

3933

3934

$$\text{Term}_{\nu,k}^2 \leq 2\alpha_*^2 M + 12M c_*^2 \bar{\varepsilon}_1^2 \beta_*^{2k}.$$

3935 Thus,

3936

3937

3938

3939

3940

3941

3942

3943

3944

$$\begin{aligned}
B_k &\leq \alpha + C_1 \text{Term}_{\nu,k}^2 + C_2 \text{Term}_{\nu,k} \text{Err}_\delta + \beta_{\text{stab}}^k \bar{\varepsilon}_1 \\
&\leq \alpha + 2C_1 \alpha_*^2 M + 12C_1 M c_*^2 \bar{\varepsilon}_1^2 \beta_*^{2k} + \text{Err}_\delta C_2 \sqrt{2M} \alpha_* + \text{Err}_\delta C_2 \sqrt{12M} c_* \bar{\varepsilon}_1 \beta_*^k + \beta_{\text{stab}}^{k/3} \bar{\varepsilon}_1 \\
&= \alpha \left(1 + 2C_1 \frac{\alpha_*^2}{\alpha} M + \frac{\alpha_*}{\alpha} \text{Err}_\delta C_2 \sqrt{2M} \right) + \beta_*^k \bar{\varepsilon}_1 \left(12C_1 M c_*^2 \bar{\varepsilon}_1 + E_\delta \sqrt{12M} c_* + \beta_{\text{stab}}^{k/3} \beta_*^{-k} \right) \\
&\leq \alpha \left(1 + 2C_1 \frac{\alpha_*^2}{\alpha} M + \frac{\alpha_*}{\alpha} \text{Err}_\delta C_2 \sqrt{2M} \right) + \beta_*^k \bar{\varepsilon}_1 \left(12C_1 M c_*^2 \bar{\varepsilon}_1 + E_\delta \sqrt{12M} c_* \right)
\end{aligned}$$

3945 where in the last line, we use $\beta_* = \beta_{\text{stab}}^{1/3} \leq 1$. Recalling $\alpha_* = 3\alpha$ and $c = 3$, we have $B_k \leq \alpha_* + c_* \bar{\varepsilon}_1 \beta_*^k = \nu_{k+1}$ as soon

3946 as

3947

3948

3949

3950

3951

3952

$$\begin{aligned}
1 &\geq 2C_1 \frac{\alpha_*^2}{\alpha} M \vee \frac{\alpha_*}{\alpha} \text{Err}_\delta C_2 \sqrt{2M} \vee 12C_1 M c_*^2 \bar{\varepsilon}_1 + E_\delta \sqrt{12M} c_* \\
&= 18\alpha C_1 M \vee 3\text{Err}_\delta C_2 \sqrt{2M} \vee 108C_1 M \bar{\varepsilon}_1 + 3E_\delta \sqrt{12M} \\
&= 18\alpha C_1 M \vee 108C_1 M \bar{\varepsilon}_1 \vee \text{Err}_\delta (3C_2 \sqrt{2M} \vee 3\sqrt{12M}).
\end{aligned}$$

3953 Thus, it suffices that

3954

3955

3956

$$\alpha \leq \frac{1}{18C_1 M}, \quad \bar{\varepsilon}_1 \leq \frac{1}{108C_1 M}, \quad \text{Err}_\delta \leq \frac{1}{12\sqrt{M} \max\{C_2, 1\}},$$

3957 as needed.

3958

3959

□

3960 **Lemma G.20** (Matrix Product Perturbation). *Define matrix products*

$$3961 \Phi_{k,j} = \mathbf{X}_{k-1} \cdot \mathbf{X}_{k-2} \cdots \mathbf{X}_j, \quad \Phi'_{k,j} = \mathbf{X}'_{k-1} \cdot \mathbf{X}'_{k-2} \cdots \mathbf{X}'_j.$$

3964 Let $\eta, \Delta, c > 0$ and $\beta_{\text{stab}} \in (0, 1)$. If (a) $\Phi_{k,j} \leq \beta_{\text{stab}}^{k-j}$ for all $j \leq k$, (b) $\|\mathbf{X}_j - \mathbf{X}'_j\| \leq \eta \Delta \beta_{\text{stab}}^{j-1}$ for all $j \geq 1$ and (c)
3965 $\Delta \leq \frac{\beta(1-\beta)}{2c\eta}$, then, for all $j \leq k$, $\|\Phi'_{k,j}\| \leq 2c\beta_{\text{stab}}^{k-j}$.

3967 *Proof.* Without loss of generality, take $j = 1$. Then, letting $\Delta_k = (\mathbf{X}'_k - \mathbf{X}_k)$,

$$\begin{aligned} 3969 \Phi'_{k+1,1} &= \mathbf{X}'_k \cdot \mathbf{X}'_{k-2} \cdots \mathbf{X}'_1 \\ 3970 &= \mathbf{X}'_k \cdot \Phi'_{k,1} \\ 3971 &= \Delta_k \Phi'_{k,1} + \mathbf{X}_k \Phi'_{k,1} \\ 3972 &= \Delta_k \Phi'_{k,1} + \mathbf{X}_k \Delta_{k-1} \Phi'_{k-2,1} + \mathbf{X}_k \mathbf{X}_{k-1} \Phi'_{k-2,1} \\ 3973 &= \Phi_{k+1,k+1} \Delta_k \Phi'_{k,1} + \Phi_{k+1,k} \Delta_{k-1} \Phi'_{k-2,1} + \Phi_{k+1,k} \Phi'_{k-2,1} \\ 3974 &= \sum_{j=1}^k \Phi_{k+1,j+1} \Delta_j \Phi'_{j,1} + \Phi_{k+1,1}. \end{aligned}$$

3979 Thus,

$$\begin{aligned} 3982 \|\Phi'_{k+1,1}\|_{\text{op}} &\leq c\eta \sum_{j=1}^k \beta_{\text{stab}}^{k-j} \|\mathbf{X}_j - \mathbf{X}'_j\| \|\Phi'_{j,1}\| + c\beta_{\text{stab}}^k \\ 3983 &\leq c\eta \beta_{\text{stab}}^{k-1} \Delta \sum_{j=1}^k \|\Phi'_{j,1}\| + c\beta_{\text{stab}}^k. \end{aligned} \quad (\|\mathbf{X}_j - \mathbf{X}'_j\| \leq \eta \Delta \beta_{\text{stab}}^{j-1})$$

3988 Define $\varepsilon_j = \|\Phi'_{j,1}\|$. Then, $\varepsilon_1 = 1 \leq c$, so [Lemma G.18](#) implies that for $\Delta \leq \frac{(1-\beta)\beta}{2\eta}$, $\|\Phi'_{k,1}\| := \varepsilon_k \leq 2c\beta_{\text{stab}}^k$ for all
3989 k . \square

3991 H. Sampling and Score Matching

3993 In this section, we provide a rigorous guarantee on the quality of sampling from the learned DDPM under [Assumption 3.3](#).
3994 We organize the section as follows:

- 3996 • In [Definition H.1](#) we provide the main notion of function class complexity, a vectorized Rademacher complexity that
3997 also appears in some form in [Block et al. \(2020a\)](#); [Maurer \(2016\)](#).
- 3999 • We then state the main result of the section, [Theorem 6](#), which provides a high probability upper bound on the number
4000 of samples n required in order to sample from DDPM trained on a given score estimate such that the sample is close in
4001 our optimal transport metric to the target distribution.
- 4002 • In particular, in [\(H.1\)](#), we give the exact polynomial dependence of the sampling parameters α and J on the parameters
4003 of the problem.
- 4005 • We break the proof of [Theorem 6](#) into two sections. First, in [Appendix H.1](#), we recall a result of [Chen et al. \(2022\)](#);
4006 [Lee et al. \(2023\)](#) that shows that it suffices to accurately learn the score in the sense that if the score estimate is accurate
4007 in the appropriate sense, then the DDPM will adequately sample from a distribution close to the target.
- 4009 • In [Remark H.5](#), we emphasize the conditions that would be required to sample in total variation and explain why they
4010 do not hold in our setting.
- 4011 • Then, in [Appendix H.2](#), we apply statistical learning techniques, similar to those in [Block et al. \(2020a\)](#), to show that
4012 with sufficiently many samples, we can effectively learn the score. We detail in [Remark H.7](#) how the realizability part
4013 of [Assumption 3.3](#) can be relaxed.

- Finally, we conclude the proof of [Theorem 6](#) by combining the two intermediate results detailed above.

To begin, we define our notion of statistical complexity:

Definition H.1 (Complexity of Θ Complexity). Define the vector- Rademacher complexity of a function class $\{\mathbf{s}_\theta | \theta \in \Theta_j\}$ by:

$$\mathcal{R}_n(\Theta_j) = \mathbb{E} \left[\sup_{\theta \in \Theta_j} \frac{1}{n} \sum_{k=1}^n \sum_{i=1}^d \varepsilon_{k,i} \mathbf{s}_\theta^{(i)}(\mathbf{a}_k, \boldsymbol{\rho}_{m,k}, j) \right],$$

where $\mathbf{s}_\theta^{(i)}$ denotes the i^{th} coordinate of \mathbf{s}_θ and the expectation is over $(\kappa, \boldsymbol{\rho}_{m,h}) \sim q_{[t]}$ and independent Rademacher random variables $\varepsilon_{k,i}$, with $q_{[t]}$ as in [Section 2](#).

We now state the main result of this section.

Theorem 6. Fix $1 \leq h \leq H$ and suppose that $(\mathbf{a}_i, \boldsymbol{\rho}_{m,h,i}) \sim q$ are independent for $1 \leq i \leq n$. Suppose that the projection of q onto the first coordinate has support (as defined in [Definition C.3](#)) contained in the euclidean ball of radius $R \geq 1$ in \mathbb{R}^d . For $\varepsilon > 0$, set

$$J = c \frac{d^3 R^4 (R + \sqrt{d})^4 \log\left(\frac{dR}{\varepsilon}\right)}{\varepsilon^{20}}, \quad \alpha = c \frac{\varepsilon^8}{d^2 R^2 (R + \sqrt{d})^2}. \quad (\text{H.1})$$

for some universal constant $c > 0$. Suppose that for all $1 \leq j \leq J$, the following hold:

- There exists a function class Θ_j containing some θ_j^* such that $\mathbf{s}_*(\cdot, \cdot, j\alpha) = \mathbf{s}_{\theta_j^*}(\cdot, \cdot, j\alpha) = \nabla \log q_{[j\alpha]}(\cdot | \cdot)$, where $q_{[\cdot]}$ is defined in [Section 2](#).
- The following holds for some $\delta > 0$:

$$\sup_{\substack{\theta, \theta' \in \Theta_j \\ \|\mathbf{a}\| \vee \|\mathbf{a}'\| \leq R + \sqrt{d \log\left(\frac{2nd}{\delta}\right)} \\ \boldsymbol{\rho}_{m,h}}} \left\| \mathbf{s}_\theta(\mathbf{a}, \boldsymbol{\rho}_{m,h}, t) - \mathbf{s}_{\theta'}(\mathbf{a}', \boldsymbol{\rho}_{m,h}, t) \right\| \leq c \frac{d^2 (R + \sqrt{d \log\left(\frac{2nd}{\delta}\right)})^2}{\varepsilon^8}.$$

- [Assumption 3.3](#) holds and thus, for all $j \in [J]$, it holds that $\mathcal{R}_n(\Theta_j) \leq C_\Theta n^{-1/\nu}$ for some $\nu \geq 2$ and all $n \in \mathbb{N}$.
- The parameter $\hat{\theta} = \hat{\theta}_{1:J}$ is defined to be the empirical minimizer of $\mathcal{L}_{\text{DDPM}}$ from [Section 3](#).

If

$$n \geq c \left(\frac{C_\Theta d R (R \vee \sqrt{d}) \log(dn)}{\varepsilon^4} \right)^{4\nu} \vee \left(\frac{d^6 (R^4 \vee d^2 \log^3\left(\frac{ndR}{\delta\varepsilon}\right))}{\varepsilon^{24}} d^2 \right)^{4\nu},$$

then with probability at least $1 - \delta$, it holds that

$$\mathbb{E}_{\boldsymbol{\rho}_{m,h} \sim q_{\boldsymbol{\rho}_{m,h}}} \left[\inf_{\mu \in \mathcal{C}(\text{DDPM}(\mathbf{s}_{\hat{\theta}}, \boldsymbol{\rho}_{m,h}), q(\cdot | \boldsymbol{\rho}_{m,h}))} \mathbb{P}_{(\hat{\mathbf{a}}, \mathbf{a}^*) \sim \mu} (\|\hat{\mathbf{a}} - \mathbf{a}^*\| \geq \varepsilon) \right] \leq 3\varepsilon.$$

Remark H.1. We emphasize that the exact value of the polynomial dependence (and in particular its pessimism) stem from the guarantees of [Chen et al. \(2022\)](#); [Lee et al. \(2023\)](#) regarding the quality of sampling with DDPMs. We remark below that the learning process itself does not incur such poor polynomial dependence except via these guarantees. Furthermore, we do not expect the sampling guarantees of those two works to be tight in any sense and such a poor polynomial dependence is not observed in practice. Rather, we include the bounds of [Chen et al. \(2022\)](#); [Lee et al. \(2023\)](#) so as to provide a fully rigorous end-to-end guarantee showing that polynomially many samples suffice to do imitation learning under our assumptions.

Remark H.2. A subtle difference between the presentation in the body and that here is the dependence of the complexity of Θ on the parameter α . We phrase the complexity guarantee as we did in the body in order to emphasize the dependence on the algorithmic parameter. If we let C'_Θ denote a constant such that $\mathcal{R}_n(\Theta) \leq C'_\Theta (\alpha/n)^{-1/\nu}$, then the sample complexity above becomes

$$n \geq c \left(\frac{C'_\Theta \log(dn)}{\alpha} \right)^{4\nu} \vee \left(\frac{d^2 (R^2 \vee d^2 \log^3(\frac{ndR}{\varepsilon\delta}))}{\alpha^2 \varepsilon^{16}} \right)^{4\nu}.$$

Remark H.3. We observe that while at first it may seem that the upper bound on the osculation of \mathbf{s}_θ is limiting, and, indeed, it is not obvious that this assumption does not contradict the realizability assumption immediately preceding it, it follows immediately from [Lemma H.2](#) that if the preceding assumptions are satisfied, then the true score function \mathbf{s}_* automatically satisfies the bound on osculation. Moreover, the boundedness of the function class is only assumed for the sake of convenience and could be substantially relaxed to an assumption requiring finiteness of moments of the envelope of the class ([Wainwright, 2019](#); [Rakhlin et al., 2017](#)). For the sake of simplicity, we do not further remark on this.

Critically, the guarantee of the quality of our DDPM is not in TV, but rather an optimal transport distance tailored to the problem at hand. As remarked in [Section 3](#), it is precisely this weaker guarantee that makes the problem challenging.

We begin by recalling the basic motivation for Denoising Diffusion Probabilistic Models (DDPMs) and explain how we train them. We then apply results from [Chen et al. \(2022\)](#) to show that if we have learned the conditional score function, then sampling can be done efficiently. While [Block et al. \(2020a\)](#) demonstrated that unconditional score learning can be learned through standard statistical learning techniques, we generalize these results to the case of conditional score learning and conclude the section by proving that with sufficiently many samples, we can efficiently sample from a distribution close to our target. In this section, we drop the subscript h for clarity, as our theoretical analysis treats each $\mathbf{s}_{\theta,h}$ separately; while empirically one sees better success in training the score estimates jointly, the focus of this paper is not on sampling and score estimation and so we make the simplifying assumption for the sake of convenience.

H.1. Denoising Diffusion Probabilistic Models

We begin by motivating the sampling procedure described in [\(2.2\)](#), which is derived by fixing a horizon T and considering the continuum limit as $\alpha \downarrow 0$ and $J = \frac{T}{\alpha}$. More precisely, consider a forward process satisfying the stochastic differential equation (SDE) for $0 \leq t \leq T$:

$$d\mathbf{a}^t = -\mathbf{a}^t dt + \sqrt{2} dB_t, \quad \mathbf{a}^0 \sim q,$$

where B_t is a Brownian motion on \mathbb{R}^d and \mathbf{a}^0 is sampled from the target density. Applying the standard time reversal to this process results in the following SDE:

$$d\mathbf{a}_{\leftarrow}^{T-t} = (\mathbf{a}_{\leftarrow}^t + 2\nabla \log q_{T-t}(\mathbf{a}_{\leftarrow}^t)) dt + \sqrt{2} dB_t, \quad \mathbf{a}_{\leftarrow}^0 \sim q_T,$$

where q_t is the law of \mathbf{a}^t . Because the forward process mixes exponentially quickly to a standard Gaussian, in order to approximately sample from q , the learner may sample $\tilde{\mathbf{a}}_{\leftarrow}^0 \sim \mathcal{N}(0, \mathbf{I})$ and evolving $\tilde{\mathbf{a}}_{\leftarrow}^t$ according to the SDE above. Note that the classical Euler-Maruyama discretization of the above procedure is exactly [\(2.2\)](#), but with the true score $\nabla \log q_{T-t}$ replaced by score estimates $\mathbf{s}_\theta(\cdot, T-t) : \mathbb{R}^d \rightarrow \mathbb{R}^d$; we may hope that if $\mathbf{s}_\theta(\cdot, T-t) \approx \nabla \log q_{T-t}$ as functions, then the procedure in [\(2.2\)](#) produces a sample close in law to q . Indeed, the following result provides a quantitative bound:

Theorem 7 (Corollary 4, [Chen et al. \(2022\)](#)). *Suppose that a distribution q on \mathbb{R}^d is supported on some ball of radius $R \geq 1$. Let C be a universal constant, fix $\varepsilon > 0$, and let α, J be set as in [\(H.1\)](#). If we have a score estimator $\mathbf{s}_\theta : \mathbb{R}^d \times [\tau] \rightarrow \mathbb{R}^d$ such that*

$$\max_{j \in [J]} \mathbb{E}_{\mathbf{a} \sim q_{[\alpha j]}} \left[\left\| \mathbf{s}_\theta(\mathbf{a}, j) - \nabla \log q_{[\alpha j]}(\mathbf{a}) \right\|^2 \right] \leq \varepsilon^4,$$

then

$$\sup_{f: \|f\|_\infty \vee \|\nabla f\|_\infty \leq 1} \mathbb{E}_{\hat{\mathbf{a}} \sim \text{Law}(\mathbf{a}^J)} [f(\hat{\mathbf{a}})] - \mathbb{E}_{\mathbf{a}^* \sim q} [f(\mathbf{a}^*)] \leq \varepsilon^2,$$

where \mathbf{a}^J is sampled as in [\(2.2\)](#).

4125 **Remark H.4.** As a technical aside, we note that [Chen et al. \(2022, Corollary 4\)](#) applies to an “early stopped” DDPM, in the
4126 sense that the denoising is stopped in slightly fewer than J steps. On the other hand, for the choice of α given above, [Chen](#)
4127 [et al. \(2022, Lemma 20 \(a\)\)](#) demonstrates that this distribution is ε^2 -close in Wasserstein distance to the sample produced by
4128 using all J steps and so by multiplying C above by a factor of 2 the guarantee is preserved. Because in practice we do not
4129 stop the DDPM early, we phrase [Theorem 7](#) in the way above as opposed to the more complicated version with the early
4130 stopping.

4131 **Remark H.5.** While ([Chen et al., 2022; Lee et al., 2023](#)) show that if \mathbf{s}_θ is close to the $\mathbf{s}_{*,h}$ in $L^2(q_{[\alpha j]})$ and q satisfies mild
4132 regularity properties, then the law of $\hat{\mathbf{a}}_h^j$ will be close in total variation to q . Unfortunately, the required regularity of q , that
4133 the score is Lipschitz, is too strong to hold in many of our applications, such as when the data lie close to a low-dimensional
4134 manifold. In such cases, [Chen et al. \(2022\)](#) provided guarantees in a weaker metric on distributions. We emphasize that even
4135 with full dimensional support, the Lipschitz constant of $\nabla \log q$ is likely large and thus the dependence on this constant
4136 appearing in [Chen et al. \(2022, Theorem 2\)](#) is unacceptable. In particular, this subtle point is what necessitates the intricate
4137 construction of our paper; as remarked in [Section 3](#), if we could expect the score to be sufficiently regular and producing a
4138 sample close in total variation to the target distribution were feasible, the problem would be trivial.
4139

4140 While [Theorem 7](#) applies to unconditional sampling, it is easy to derive conditional sampling guarantees as a corollary.

4141 **Corollary H.1.** *Suppose that q is a joint distribution on actions \mathbf{a} and observations $\boldsymbol{\rho}_{m,h} \in \mathbb{R}^d$. Further assume that the*
4142 *marginals over \mathbb{R}^d are fully supported in a ball of radius $R \geq 1$. Then there exists a universal constant C such that for all*
4143 *small $\varepsilon > 0$, if J and α are set as in [\(H.1\)](#) and*
4144

$$4145 \mathbb{E}_{\boldsymbol{\rho}_{m,h} \sim q_{\boldsymbol{\rho}_{m,h}}} \left[\max_{j \in [J]} \mathbb{E}_{\mathbf{a} \sim q_{[\alpha j]}(\cdot | \boldsymbol{\rho}_{m,h})} \left[\left\| \mathbf{s}_\theta(\mathbf{a}, j, \boldsymbol{\rho}_{m,h}) - \nabla \log q_{[\alpha j]}(\mathbf{a} | \boldsymbol{\rho}_{m,h}) \right\|^2 \right] \right] \leq \varepsilon^4, \quad (\text{H.2})$$

4148 then

$$4150 \mathbb{E}_{\boldsymbol{\rho}_{m,h} \sim q_{\boldsymbol{\rho}_{m,h}}} \left[\inf_{\mu \in \mathcal{C}(\text{DDPM}(\mathbf{s}_\theta, \boldsymbol{\rho}_{m,h}), q(\cdot | \boldsymbol{\rho}_{m,h}))} \mathbb{P}_{(\hat{\mathbf{a}}, \mathbf{a}^*) \sim \mu} (\|\hat{\mathbf{a}} - \mathbf{a}^*\| \geq \varepsilon) \right] \leq 3\varepsilon$$

4153 *Proof.* We begin by showing an intermediate result,

$$4155 \mathbb{E}_{\boldsymbol{\rho}_{m,h} \sim q_{\boldsymbol{\rho}_{m,h}}} \left[\sup_{f: \|f\|_\infty \vee \|\nabla f\|_\infty \leq 1} \mathbb{E}_{\hat{\mathbf{a}} \sim \text{DDPM}(\mathbf{s}_\theta, \boldsymbol{\rho}_{m,h})} [f(\hat{\mathbf{a}})] - \mathbb{E}_{\mathbf{a}^* \sim q(\cdot | \boldsymbol{\rho}_{m,h})} [f(\mathbf{a}^*)] \right] \leq 3\varepsilon^2. \quad (\text{H.3})$$

4158 using [Theorem 7](#). Let

$$4160 \mathcal{A} = \left\{ \max_{j \in [J]} \mathbb{E}_{\mathbf{a} \sim q_{[\alpha j]}(\cdot | \boldsymbol{\rho}_{m,h})} \left[\left\| \mathbf{s}_\theta(\mathbf{a}, j, \boldsymbol{\rho}_{m,h}) - \nabla \log q_{[\alpha j]}(\mathbf{a} | \boldsymbol{\rho}_{m,h}) \right\|^2 \right] \leq \varepsilon^2 \right\}.$$

4162 By Markov’s inequality and [\(H.2\)](#), it holds that

$$4165 \mathbb{P}_{\boldsymbol{\rho}_{m,h} \sim q_{\boldsymbol{\rho}_{m,h}}} (\mathcal{A}^c) \leq \frac{\varepsilon^4}{\varepsilon^2} = \varepsilon^2$$

4167 and thus

$$4169 \mathbb{E}_{\boldsymbol{\rho}_{m,h} \sim q_{\boldsymbol{\rho}_{m,h}}} \left[\sup_{f: \|f\|_\infty \vee \|\nabla f\|_\infty \leq 1} \mathbb{E}_{\hat{\mathbf{a}} \sim \text{DDPM}(\mathbf{s}_\theta, \boldsymbol{\rho}_{m,h})} [f(\hat{\mathbf{a}})] - \mathbb{E}_{\mathbf{a}^* \sim q(\cdot | \boldsymbol{\rho}_{m,h})} [f(\mathbf{a}^*)] \right]$$

$$4172 = \mathbb{E}_{\boldsymbol{\rho}_{m,h} \sim q_{\boldsymbol{\rho}_{m,h}}} \left[\mathbf{I}[\mathcal{A}] \sup_{f: \|f\|_\infty \vee \|\nabla f\|_\infty \leq 1} \mathbb{E}_{\hat{\mathbf{a}} \sim \text{DDPM}(\mathbf{s}_\theta, \boldsymbol{\rho}_{m,h})} [f(\hat{\mathbf{a}})] - \mathbb{E}_{\mathbf{a}^* \sim q(\cdot | \boldsymbol{\rho}_{m,h})} [f(\mathbf{a}^*)] \right]$$

$$4173 + \mathbb{E}_{\boldsymbol{\rho}_{m,h} \sim q_{\boldsymbol{\rho}_{m,h}}} \left[\mathbf{I}[\mathcal{A}^c] \sup_{f: \|f\|_\infty \vee \|\nabla f\|_\infty \leq 1} \mathbb{E}_{\hat{\mathbf{a}} \sim \text{DDPM}(\mathbf{s}_\theta, \boldsymbol{\rho}_{m,h})} [f(\hat{\mathbf{a}})] - \mathbb{E}_{\mathbf{a}^* \sim q(\cdot | \boldsymbol{\rho}_{m,h})} [f(\mathbf{a}^*)] \right]$$

$$4176 \leq \mathbb{E}_{\boldsymbol{\rho}_{m,h} \sim q_{\boldsymbol{\rho}_{m,h}}} \left[\mathbf{I}[\mathcal{A}] \inf_{q' \in \Delta(\mathbb{R}^d)} W_2(q(\cdot | \boldsymbol{\rho}_{m,h}), q') + \text{TV}(q', \text{Law}(\pi^\tau)) \right] + 2\varepsilon^2.$$

4179

4180 For each $\rho_{m,h}$, we may apply [Theorem 7](#) and observe that for $\rho_{m,h} \in \mathcal{A}$,

$$4181 \sup_{f: \|f\|_\infty \vee \|\nabla f\|_\infty \leq 1} \mathbb{E}_{\hat{\mathbf{a}} \sim \text{DDPM}(\mathbf{s}_\theta, \rho_{m,h})} [f(\hat{\mathbf{a}})] - \mathbb{E}_{\mathbf{a}^* \sim q(\cdot | \rho_{m,h})} [f(\mathbf{a}^*)] \leq \varepsilon^2,$$

4184 which proves [\(H.3\)](#). Now, for any fixed $\rho_{m,h}$, by Markov's inequality and the definition of Wasserstein distance,

$$4185 \inf_{\mu \in \mathcal{C}(\text{DDPM}(\mathbf{s}_\theta, \rho_{m,h}), q(\cdot | \rho_{m,h}))} \mathbb{P}_{(\hat{\mathbf{a}}, \mathbf{a}^*) \sim \mu} (\|\hat{\mathbf{a}} - \mathbf{a}^*\| \geq \varepsilon) \leq \frac{W_1(\text{DDPM}(\mathbf{s}_\theta, \rho_{m,h}), q(\cdot | \rho_{m,h}))}{\varepsilon}.$$

4189 The result follows. \square

4191 Note that the guarantee in [Corollary H.1](#) is precisely what we need to control the one step imitation error in [Theorem 2](#); thus, the problem of conditional sampling has been reduced to estimating the score. In the subsequent section, we will apply standard statistical learning techniques to provide a nonasymptotic bound on the quality of a score estimator.

4195 H.2. Score Estimation

4196 In the previous section we have shown that conditional sampling can be reduced to the problem of learning the conditional score. While there exist non-asymptotic bounds for learning the unconditional score ([Block et al., 2020a](#)), they apply to a slightly different score estimator than is typically used in practice. Here we upper bound the estimation error in terms of the complexity of the space of parameters Θ .

4201 Observe that in order to apply [Corollary H.1](#), we need a guarantee on the error of our score estimate in $L^2(q_{[\alpha j]})$ for all $j \in [J]$.

4202 Ideally, then, for fixed $\rho_{m,h}$ and $t = \alpha j$, we would like to minimize $\mathbb{E}_{\mathbf{a} \sim q_{[t]}} \left[\left\| \mathbf{s}_\theta(\mathbf{a}, \rho_{m,h}, t) - \nabla \log q_{[t]}(\mathbf{a} | \rho_{m,h}) \right\|^2 \right]$, where the inner norm is the Euclidean norm on \mathbb{R}^d . Unfortunately, because $q_{[t]}$ itself is unknown, we cannot even take an empirical version of this loss. Instead, through a now classical integration by parts ([Hyvärinen & Dayan, 2005](#); [Vincent, 2011](#); [Song & Ermon, 2019](#)), this objective can be shown to be equivalent to minimizing

$$4207 \mathcal{L}_{\text{DDPM}}(\theta, \mathbf{a}, \rho_{m,h}, t) = \mathbb{E}_{\mathbf{a} \sim q_{[t]}} \left[\left\| \mathbf{s}_\theta \left(e^{-t} \mathbf{a} + \sqrt{1 - e^{-2t}} \boldsymbol{\gamma}, \rho_{m,h}, t \right) + \frac{1}{\sqrt{1 - e^{-2t}}} \boldsymbol{\gamma} \right\|^2 \right].$$

4211 Because we are really interested in the expectation over the joint distribution $(\mathbf{a}, \rho_{m,h})$, we may take the expectation over $\rho_{m,h}$ and recover [\(3.1\)](#) as the empirical approximation. We now prove the following result for a single time step t :

4213 **Proposition H.1.** *Suppose that q is a distribution such that $q(\cdot | \rho_{m,i})$ is supported on a ball of radius R for q -almost every $\rho_{m,h}$. For fixed $j \in [J]$ and α from [\(H.1\)](#), let $t = j\alpha$ and suppose that there is some $\theta^* \in \Theta_j$ such that $\mathbf{s}_*(\cdot, \cdot, t) = \mathbf{s}_{\theta^*}(\cdot, \cdot, t) = \nabla \log q_{[t]}(\cdot | \cdot)$, i.e., \mathbf{s}_θ is rich enough to represent the true score at time t . Suppose further that the class of functions $\{\mathbf{s}_\theta | \theta \in \Theta_j\}$ satisfies for all $\theta \in \Theta_j$,*

$$4218 \sup_{\substack{\theta, \theta' \in \Theta_j \\ \|\mathbf{a}\| \vee \|\mathbf{a}'\| \leq R \\ \rho_{m,h}}} \left\| \mathbf{s}_\theta(\mathbf{a}, \rho_{m,h}, t) - \mathbf{s}_{\theta'}(\mathbf{a}', \rho_{m,h}, t) \right\| \leq c \frac{d^2 (R + \sqrt{d \log(\frac{2nd}{\delta})})^2}{\varepsilon^8}$$

4223 for some universal constant $c > 0$. Recall the Rademacher term $\mathcal{R}_n(\Theta_j)$ defined in [Definition H.1](#), and let

$$4225 \hat{\theta} \in \arg \min_{\theta \in \Theta} \sum_{i=1}^n \mathcal{L}_{\text{DDPM}}(\theta, \mathbf{a}_i, \rho_{m,i}, t)$$

4228 for independent and identically distributed $(\mathbf{a}_i, \rho_{m,i}) \sim q$. Then it holds with probability at least $1 - \delta$ over the data that

$$4230 \mathbb{E}_{(\mathbf{a}_t, \rho_{m,h}) \sim q_{[t]}} \left[\left\| \mathbf{s}_{\hat{\theta}}(\mathbf{a}_t, \rho_{m,h}, t) - \nabla \log q_{[t]}(\mathbf{a}_t | \rho_{m,h}) \right\|^2 \right]$$

$$4231 \leq c \cdot \sqrt{\frac{\log(dn)}{1 - e^{-2t}}} \left(\mathcal{R}_n(\Theta) + \frac{d^2 (R + \sqrt{d \log(\frac{2nd}{\delta})})^2}{\varepsilon^8} \cdot \sqrt{\frac{d \log(\frac{4dn}{\delta})}{n}} \right).$$

4235 Before we provide a proof, we recall the following result:

4236 **Lemma H.2.** Suppose that q is supported in a ball of radius R and let $t \geq \alpha$ for α as in (H.1). Then $\nabla \log q_{[t]}(\cdot|\cdot)$ is
4237 L -Lipschitz with respect to the first parameter for
4238

$$4239 L = \frac{dR^2(R \vee \sqrt{d})^2}{\varepsilon^8}.$$

4240 In particular,

$$4241 \sup_{\substack{\|\mathbf{a}\| \vee \|\mathbf{a}'\| \leq R \\ \boldsymbol{\rho}_{m,h}}} \left| \nabla \log q_{[t]}(\mathbf{a}|\boldsymbol{\rho}_{m,h}) - \nabla \log q_{[t]}(\mathbf{a}'|\boldsymbol{\rho}_{m,h}) \right| \leq 2LR$$

4242 and there exists some assignment of Θ and \mathbf{s}_θ that satisfies the boundedness condition in Proposition H.1.

4243 *Proof.* The first statement follows from replacing the ε in Chen et al. (2022, Lemma 20 (c)) by ε^2 . The second statement
4244 follows immediately from the first. \square

4245 We also require the following standard result:

4246 **Lemma H.3.** If $\mathcal{R}_n(\Theta_j)$ is defined as in Definition H.1, then

$$4247 \mathbb{E}_{\boldsymbol{\gamma}_1, \dots, \boldsymbol{\gamma}_n} \left[\sup_{\substack{\theta \in \Theta_j \\ 1 \leq j \leq J}} \frac{1}{n} \cdot \sum_{i=1}^n \langle \mathbf{s}_\theta(\mathbf{a}, \boldsymbol{\rho}_{m,i}, j), \boldsymbol{\gamma}_i \rangle \right] \leq \sqrt{\pi \log(dn)} \cdot \mathcal{R}_n(\Theta_j)$$

4248 *Proof.* This statement is classical and follows immediately from the fact that the norm of a Gaussian is independent from its
4249 sign as well as the fact that $\mathbb{E}[\max_{i,j}(\boldsymbol{\gamma}_i)_j] \leq \sqrt{\pi \log(dn)}$ by classical Gaussian concentration. See Van Handel (2014) for
4250 more details. \square

4251 *Proof of Proposition H.1.* Let P_n denote the empirical measure on n independent samples $\{(\mathbf{a}_i, \boldsymbol{\rho}_{m,i}, \boldsymbol{\gamma}_i)\}$ and let $\mathbf{a}_i^t =$
4252 $e^{-t}\mathbf{a}_i + \sqrt{1 - e^{-2t}}\boldsymbol{\gamma}_i$. Let $C_t = \sqrt{1 - e^{-2t}}$ and observe that by definition and realizability,

$$4253 P_n \left(\|C_t \mathbf{s}_{\hat{\theta}}(\mathbf{a}^t, \boldsymbol{\rho}_{m,h}, t) - \boldsymbol{\gamma}\|^2 \right) \leq P_n \left(\|C_t \nabla \log q_{[t]}(\mathbf{a}^t|\boldsymbol{\rho}_{m,h}) - \boldsymbol{\gamma}\|^2 \right). \quad (\text{H.4})$$

4254 We emphasize that by Lemma H.2, realizability does not make the result vacuous. Adding and subtracting
4255 $C_t \nabla \log q_{[t]}(\mathbf{a}^t|\boldsymbol{\rho}_{m,h})$ from the left hand inequality, expanding and rearranging, we see that

$$4256 C_t^2 P_n \left(\|\mathbf{s}_{\hat{\theta}}(\mathbf{a}^t, \boldsymbol{\rho}_{m,h}, t) - \nabla \log q_{[t]}(\mathbf{a}^t|\boldsymbol{\rho}_{m,h})\|^2 \right) \leq 2C_t \cdot P_n \left(\langle \mathbf{s}_{\hat{\theta}}(\mathbf{a}^t, \boldsymbol{\rho}_{m,h}, t) - \nabla \log q_{[t]}(\mathbf{a}^t|\boldsymbol{\rho}_{m,h}), \boldsymbol{\gamma} \rangle \right) \\
4257 \leq 2C_t \cdot P_n \left(\sup_{\theta \in \Theta_j} \langle \mathbf{s}_\theta(\mathbf{a}^t, \boldsymbol{\rho}_{m,h}, t) - \nabla \log q_{[t]}(\mathbf{a}^t|\boldsymbol{\rho}_{m,h}), \boldsymbol{\gamma} \rangle \right).$$

4258 We now claim that with probability at least $1 - \delta$, it holds that

$$4259 P_n \left(\sup_{\theta \in \Theta} \langle \mathbf{s}_\theta(\mathbf{a}^t, \boldsymbol{\rho}_{m,h}, t) - \nabla \log q_{[t]}(\mathbf{a}^t|\boldsymbol{\rho}_{m,h}), \boldsymbol{\gamma} \rangle \right) \leq \mathbb{E} \left[P_n \left(\sup_{\theta \in \Theta_j} \langle \mathbf{s}_\theta(\mathbf{a}^t, \boldsymbol{\rho}_{m,h}, t) - \nabla \log q_{[t]}(\mathbf{a}^t|\boldsymbol{\rho}_{m,h}), \boldsymbol{\gamma} \rangle \right) \right] \\
4260 + B \cdot \sqrt{\frac{d \log \left(\frac{2d}{\delta} \right)}{n}},$$

4261 where

$$4262 B = c \frac{d^2(R + \sqrt{d \log \left(\frac{2nd}{\delta} \right)})^2}{\varepsilon^8} \quad (\text{H.5})$$

4263

4290 for some universal constant $c > 0$. To see this, we claim that with probability at least $1 - \frac{\delta}{2}$, it holds that $\|\mathbf{a}_i^t\| \leq$
4291 $c \left(R + \sqrt{d \log \left(\frac{2nd}{\delta} \right)} \right)$ for all $1 \leq i \leq n$. Indeed, this follows by Gaussian concentration in [Jin et al. \(2019, Lemmata](#)
4292 [1 & 2\)](#). We may now apply [Lemma H.2](#) to a bound on the osculation of $\mathbf{s}_\theta - \nabla \log q_{[t]}$ in the ball of the above radius.
4293 Conditioning on the event that $\|\mathbf{a}_i^t\|$ is bounded by the above, we may argue as in [Wainwright \(2019, Theorem 4.10\)](#) that if
4294 we let the function
4295

$$4296 \quad G = G(\mathbf{a}_1, \boldsymbol{\rho}_{m,1}, \dots, \mathbf{a}_n, \boldsymbol{\rho}_{m,n}) = P_n \left(\sup_{\theta \in \Theta_j} \langle \mathbf{s}_\theta(\mathbf{a}^t, \boldsymbol{\rho}_{m,h}, t) - \nabla \log q_{[t]}(\mathbf{a}^t | \boldsymbol{\rho}_{m,h}), \boldsymbol{\gamma} \rangle \right),$$

4300 then for any i , on the event of bounded norm, replacing $(\mathbf{a}_i, \boldsymbol{\rho}_{m,i})$ with $(\mathbf{a}'_i, \boldsymbol{\rho}'_{m,i})$ and leaving other terms unchanged changes
4301 ensures that $|G - G'| \leq \frac{2B}{n} \gamma_i$. Thus by [Jin et al. \(2019, Corollary 7\)](#) and a union bound, the claim holds. Because $\boldsymbol{\gamma}$ is
4302 mean zero, we have
4303

$$4304 \quad \mathbb{E} \left[P_n \left(\sup_{\theta \in \Theta} \langle \mathbf{s}_\theta(\mathbf{a}^t, \boldsymbol{\rho}_{m,h}, t) - \nabla \log q_{[t]}(\mathbf{a}^t | \boldsymbol{\rho}_{m,h}), \boldsymbol{\gamma} \rangle \right) \right] \leq \mathbb{E} \left[P_n \left(\sup_{\theta \in \Theta} \langle \mathbf{s}_\theta(\mathbf{a}^t, \boldsymbol{\rho}_{m,h}, t), \boldsymbol{\gamma} \rangle \right) \right]$$

$$4305 \quad \leq \sqrt{\pi \log(dn)} \cdot \mathcal{R}_n(\Theta_j),$$

4308 where the last inequality follows by [Lemma H.3](#) and the fact that $t = jJ$. Summing up the argument until this point and
4309 rearranging tells us that with probability at least $1 - \delta$, it holds that
4310

$$4311 \quad P_n \left(\left\| \mathbf{s}_{\hat{\theta}}(\mathbf{a}^t, \boldsymbol{\rho}_{m,h}, t) - \nabla \log q_{[t]}(\mathbf{a}^t | \boldsymbol{\rho}_{m,h}) \right\|^2 \right) \leq \frac{2}{C_t} \sqrt{\pi \log(nd)} \cdot \mathcal{R}_n(\Theta) + \frac{B}{C_t} \cdot \sqrt{\frac{d \log \left(\frac{2nd}{\delta} \right)}{n}},$$

4315 with B given in [\(H.5\)](#). We now use a uniform norm comparison between population and empirical norms to conclude the
4316 proof. Indeed, it holds by [Rakhlin et al. \(2017, Lemma 8.i & 9\)](#) that there exists a critical radius
4317

$$4318 \quad r_n \leq cB \log^3(n) \mathcal{R}_n(\Theta_j)^2$$

4319 such that with probability at least $1 - \delta$,

$$4320 \quad \mathbb{E}_{(\mathbf{a}^t, \boldsymbol{\rho}_{m,h}) \sim q_{[t]}} \left[\left\| \mathbf{s}_{\hat{\theta}}(\mathbf{a}^t, \boldsymbol{\rho}_{m,h}, t) - \nabla \log q_{[t]}(\mathbf{a}^t | \boldsymbol{\rho}_{m,h}) \right\|^2 \right]$$

$$4321 \quad \leq 2 \cdot P_n \left(\left\| \mathbf{s}_{\hat{\theta}}(\mathbf{a}^t, \boldsymbol{\rho}_{m,h}, t) - \nabla \log q_{[t]}(\mathbf{a}^t | \boldsymbol{\rho}_{m,h}) \right\|^2 \right) + cr_n + c \frac{\log \left(\frac{1}{\delta} \right) + \log \log n}{n},$$

4326 where again c is some universal constant. Combining this with our earlier bound on the empirical distance and a union
4327 bound, after rescaling δ , we have that
4328

$$4329 \quad \mathbb{E}_{(\mathbf{a}^t, \boldsymbol{\rho}_{m,h}) \sim q_{[t]}} \left[\left\| \mathbf{s}_{\hat{\theta}}(\mathbf{a}^t, \boldsymbol{\rho}_{m,h}, t) - \nabla \log q_{[t]}(\mathbf{a}^t | \boldsymbol{\rho}_{m,h}) \right\|^2 \right] \leq \frac{4}{C_t} \sqrt{\pi \log(nd)} \cdot \mathcal{R}_n(\Theta_j) + \frac{2B}{C_t} \cdot \sqrt{\frac{d \log \left(\frac{4nd}{\delta} \right)}{n}}$$

$$4330 \quad + cB \log^3(n) \cdot \mathcal{R}_n^2(\Theta_j) + c \frac{\log \left(\frac{2}{\delta} \right) + \log \log(n)}{n}$$

4334 with probability at least $1 - \delta$. This concludes the proof. \square

4336 **Remark H.6.** For the sake of simplicity, in the proof of [Proposition H.1](#) we applied uniform deviations and recovered
4337 the “slow rate” of $\mathcal{R}_n(\Theta)$, up to logarithmic factors. Indeed, if we were to further assume that the score function class is
4338 star-shaped around the true score, we could recover a faster rate, as was done in the case of unconditional sampling in [Block](#)
4339 [et al. \(2020a\)](#) with a slightly different loss. While in our proof the appeal to [Rakhlin et al. \(2017\)](#) to control the population
4340 norm by the empirical norm could be replaced with a simpler uniform deviations argument because we have already given
4341 up on the fast rate, such an argument is necessary in the more refined analysis. As the focus of this paper is not on the
4342 sampling portion of the end-to-end analysis, we do not include a rigorous proof of the case of fast rates for the sake of
4343 simplicity and space.
4344

Remark H.7. While we assumed for simplicity that the score was realizable with respect to our function class for every time $t = \alpha j$, this condition can be relaxed to approximate realizability in a standard way. In particular, if the score is ε -far away from some function representable by our class in a pointwise sense, then we can add an ε to the right hand side of (H.4) with minimal modification to the proof.

With Proposition H.1, and a union bound, we recover the following result:

Proposition H.4. Suppose that the conditions on \mathbf{s}_θ in Proposition H.1 continue to hold. Let J and α be as in (H.1) and suppose that $\alpha \leq \frac{1}{2}$. Then, with probability at least $1 - \delta$ over \mathcal{D}' , it holds that

$$\begin{aligned} & \mathbb{E}_{\rho_{m,h} \sim q_{\rho_{m,h}}} \left[\max_{j \in [J]} \mathbb{E}_{\mathbf{a} \sim q_{[\alpha j]}(\cdot | \rho_{m,h})} \left[\left\| \mathbf{s}_\theta(\mathbf{a}, j, \rho_{m,h}) - \nabla \log q_{[\alpha j]}(\mathbf{a} | \rho_{m,h}) \right\|^2 \right] \right] \\ & \leq c \frac{dR(R \vee \sqrt{d}) \log(dn)}{\varepsilon^4} \mathcal{R}_n(\Theta) + c \frac{d^3 (R^2 + d \log(\frac{ndR}{\delta\varepsilon}))}{\varepsilon^{12}} \sqrt{\frac{d \log(\frac{4dnR}{\delta\varepsilon})}{n}} \end{aligned}$$

In particular if

$$\mathcal{R}_n(\Theta_j) \leq C_\Theta n^{-1/\nu}$$

for some $\nu \geq 2$ and all $j \in [J]$, then for

$$n \geq c \left(\frac{C_\Theta dR(R \vee \sqrt{d}) \log(dn)}{\varepsilon^4} \right)^{4\nu} \vee \left(\frac{d^6 (R^4 \vee d^2 \log^3(\frac{ndR}{\delta\varepsilon}))}{\varepsilon^{24}} d^2 \right)^{4\nu}$$

it holds that with probability at least $1 - \delta$,

$$\mathbb{E}_{\rho_{m,h} \sim q_{\rho_{m,h}}} \left[\max_{j \in [J]} \mathbb{E}_{\mathbf{a} \sim q_{[\alpha j]}(\cdot | \rho_{m,h})} \left[\left\| \mathbf{s}_\theta(\mathbf{a}, j, \rho_{m,h}) - \nabla \log q_{[\alpha j]}(\mathbf{a} | \rho_{m,h}) \right\|^2 \right] \right] \leq \varepsilon^4.$$

Proof. We begin by noting that

$$1 - e^{-2t} \geq 1 - e^{-2\alpha} \geq \alpha$$

because $2\alpha \leq 1$. We now apply Proposition H.1 and take a union bound over $j \in [J]$. The result follows. \square

We note that in our simplified analysis, we have assumed that $N_{\text{aug}} = 1$, i.e., for each sample, we take a single noise level from the path. In practice, we use many augmentations per sample. Again, as the focus of our paper is not on score estimation and sampling, we treat this as a simple convenience and leave open to future work the problem of rigorously demonstrating that multiple augmentations indeed help with learning. Finally, for a discussion on bounding $\mathcal{R}_n(\Theta)$, see Wainwright (2019).

Proof of Theorem 6. We note that the proof follows immediately from combining Corollary H.1 with Proposition H.4. \square

I. End-to-end Guarantees and the Proof of Theorem 1

In this section, we provide a number of end-to-end guarantees for the learned imitation policy under various assumptions. The core of the section is Theorem 8, which provides the basis for the final proof of Theorem 1 in the body by uniting the analysis in the composite MDP from Appendix E, the control theory from Appendix G, and the sampling guarantees from Appendix H. We now summarize the organisation of the appendix:

- In Appendix I.1, we recall the association between the control setting and the composite MDP presented in Section 4, as well as rigorously instantiating the direct decomposition and the expert policy.
- In Appendix I.2, we state a reduction from imitation learning to conditional sampling, which we then use to derive a proof of Theorem 1.

- In [Appendix I.3](#), we demonstrate that if the demonstrator policy is assumed to be TVC, then we can recover stronger guarantees than those provided in [Theorem 1](#) without this assumption; in particular, we show that we can bound the *joint* imitation loss as well as the marginal and final versions.
- In [Appendix I.4](#), we show that if we were able to produce samples from a distribution close in *total variation* to the expert policy distribution, as opposed to the weaker optimal transport metric that we consider in the rest of the paper, then without any further assumptions, imitation learning is easily achievable.
- In [Appendix I.5](#), we show that if we remove the data augmentation from TODA, i.e., we set $\sigma = 0$, then we can recover similar guarantees under the assumption that the imitator policy $\hat{\pi}$ is TVC. In this way, we show that in some sense, total variation continuity is the important property imparted by smoothing.
- In [Appendix I.6](#), we demonstrate the utility of our imitation losses, showing that for Lipschitz cost functions decomposing in natural ways, our imitation losses as defined in [Definition 2.2](#) provide control over the difference in expected cost under expert and imitated distributions.
- Finally, in [Appendix I.7](#), we collect a number of useful lemmata that we use throughout the appendix.

I.1. Preliminaries

Here, we state various preliminaries to the end-to-end theorems. For simplicity, to avoid complications with the boundary effects at $h = 1$, we re-define $h = 1$ -memory chunks $\rho_{m,1}$ as elements \mathcal{P}_{τ_m-1} by prepending the necessary zeros – i.e. $\rho_{m,1} = (0, 0, \dots, 0, \mathbf{x}_1)$ – and similarly modifying $\rho_{c,1} \in \mathcal{P}_{\tau_c}$ by prepending zeros. We first recall the definitions of the composite-states and -actions from [Section 4](#). The prepending of zeros in the $h = 1$ case is mentioned above. For $h > 1$, recall that $s_h = (\mathbf{x}_{t_{h-1}:t_h}, \mathbf{u}_{t_{h-1}:t_h-1})$ and that $\mathbf{a}_h = \kappa_{t_h:t_{h+1}-1}$, where we again emphasize that \mathbf{a}_h begins at the same t that s_{h+1} does. We further recall that $d_S(s_h, s'_h) = \max_{t \in [t_{h-1}:t_h]} \|\mathbf{x}_t - \mathbf{x}'_t\| \vee \max_{t \in [t_{h-1}:t_h-1]} \|\mathbf{u}_t - \mathbf{u}'_t\|$, $d_{\text{TVC}}(s_h, s'_h) = \max_{t \in [t_h-\tau_m:t_h]} \|\mathbf{x}_t - \mathbf{x}'_t\| \vee \max_{t \in [t_h-\tau_m:t_h-1]} \|\mathbf{u}_t - \mathbf{u}'_t\|$, and $d_{\text{IPS}}(s_h, s'_h) = \|\mathbf{x}_{t_h} - \mathbf{x}'_{t_h}\|$. Finally, for $\mathbf{a} = (\bar{\mathbf{u}}_{1:\tau_c}, \bar{\mathbf{x}}_{1:\tau_c}, \bar{\mathbf{K}}_{1:\tau_c})$ and $\mathbf{a}' = (\bar{\mathbf{u}}'_{1:\tau_c}, \bar{\mathbf{x}}'_{1:\tau_c}, \bar{\mathbf{K}}'_{1:\tau_c})$, recall from [\(4.2\)](#) that

$$d_{\mathcal{A}}(\mathbf{a}, \mathbf{a}') := c_1 \max_{1 \leq k \leq \tau_c} (\|\bar{\mathbf{u}}_k - \bar{\mathbf{u}}'_k\| + \|\bar{\mathbf{x}}_k - \bar{\mathbf{x}}'_k\| + \|\bar{\mathbf{K}}_k - \bar{\mathbf{K}}'_k\|) + \mathbf{I}_{0,\infty}\{\mathcal{E}\},$$

where we $\mathcal{E} = \{\max_{1 \leq k \leq \tau_c} \max\{\|\bar{\mathbf{u}}_k - \bar{\mathbf{u}}'_k\|, \|\bar{\mathbf{x}}_k - \bar{\mathbf{x}}'_k\|, \|\bar{\mathbf{K}}_k - \bar{\mathbf{K}}'_k\|\} \leq c_2\}$, $\mathbf{I}_{0,\infty}$ is the indicator taking infinite value when the event fails to hold, and c_1 and c_2 are given in [Definition G.5](#).

Direct Decomposition and Smoothing Kernel. This section will invoke the generalizations [Theorem 2](#) which requires TVC only subspace of the state space. This invokes the direct decomposition explained in [Appendix E](#).

Definition I.1 (Direct Decomposition and Smoothing Kernel). We consider the decomposition of $\mathcal{S} = \mathcal{Z} \oplus \mathcal{V}$, where $\mathcal{Z} = \mathcal{P}_{\tau_m-1}$ are the coordinates of $\rho_{c,h}$ corresponding to the memory chunk $\rho_{m,h}$, and \mathcal{V} are all remaining coordinates. We let $\phi_{\mathcal{Z}} : \mathcal{S} \rightarrow \mathcal{Z}$ denote the projection onto the coordinates in \mathcal{Z} . We instantiate the smoothing kernel W_{σ} as follows: For $\mathbf{s} = \rho_{c,h} \in \mathcal{P}_{\tau_c}$, we let

$$W_{\sigma}(\mathbf{s}) = \mathcal{N}\left(\rho_{c,h}, \begin{bmatrix} \sigma^2 \mathbf{I}_{\mathcal{Z}} & 0 \\ 0 & 0 \end{bmatrix}\right),$$

where $\mathbf{I}_{\mathcal{Z}}$ denotes the identity supported on the coordinates in \mathcal{Z} as described above.

We note that the above direct decomposition satisfies the requisite compatibility assumptions explained in [Appendix E](#). Note also that d_{IPS} and W_{σ} are compatible with the above direct decomposition.

Chunking Policies. We continue by centralizing a definition of chunking policies.

Definition I.2 (Policy and Initial-State Distributions). Given an *chunking policy* $\pi = (\pi_h)_{h=1}^H$ with $\pi_h : \mathcal{P}_{\tau_m-1} \rightarrow \Delta(\mathcal{A})$, we let \mathcal{D}_{π} denote the distribution over ρ_T and $\mathbf{a}_{1:H}$ induced by selecting $\mathbf{a}_h \sim \pi_h(\rho_{m,h})$, and rolling out the dynamics as described in [Section 2](#). We extend chunking policies to maps $\pi_h : \mathcal{S} = \mathcal{P}_{\tau_c} \rightarrow \Delta(\mathcal{A})$ by expressing $\pi_h = \pi_h \circ \phi_{\mathcal{Z}}$ (i.e., projection $\rho_{c,h}$ onto its $\rho_{m,h}$ -components). Further, we let \mathbf{P}_{init} denote the distribution of \mathbf{x}_1 under $\rho_T \sim \mathcal{D}_{\text{exp}}$.

Remark I.1. The notation \mathcal{D}_{π} denotes the special case of chunking policies in the control setting of [Section 2](#), whereas we reserve the serif font \mathbf{D}_{π} for the distribution induced by policies in the abstract MDP. For composite MDPs instantiated as in [Section 4.1](#), the two exactly coincide.

4455 **Construction of π^* for composite MDP.** We now explain how to extract π^* from \mathcal{D}_{exp} in the composite MDP.

4456 **Definition I.3** (Policies corresponding to \mathcal{D}_{exp}). Define the following sequence kernels $\pi^* = (\pi_h^*)_{h=1}^H$ and $\pi_{\text{dec}}^* =$
 4457 $(\pi_{\text{dec},h}^*)_{h=1}^H$ via the following process. Let $\rho_T \sim \mathcal{D}_{\text{exp}}$, and let $\mathbf{a}_{1:H} = \text{synth}(\rho_T)$; further, let $\rho_{\text{m},1:H}$ be the corresponding
 4458 memory-chunks from ρ_T . Let

- 4459 • $\pi_h^*(\cdot) : \mathcal{P}_{\tau_{\text{m}}-1} \rightarrow \mathcal{A}$ denote a regular conditional probability corresponding to the distribution over \mathbf{a}_h given $\rho_{\text{m},h}$ in
 4460 the above construction.
- 4461 • $\pi_{\text{dec},h}^*(\cdot) : \mathcal{P}_{\tau_{\text{m}}-1} \rightarrow \mathcal{A}$ denote a regular conditional probability corresponding to the distribution over \mathbf{a}_h given an
 4462 augmented $\tilde{\rho}_{\text{m},h} \sim \mathcal{N}(\rho_{\text{m},h}, \sigma^2 \mathbf{I})$.

4463 Finally, for π^* as constructed above, \mathbb{P}_h^* denotes the distribution over $\rho_{c,h}$ under \mathcal{D}_{π^*} . By [Lemma I.6](#), this is in fact equal to
 4464 the distribution over $\rho_{c,h}$ under \mathcal{D}_{exp} . Notice further, therefore, that $\phi_{\mathcal{Z}} \circ \mathbb{P}_h^*$ is precisely the distribution of $\rho_{\text{m},h}$ under
 4465 \mathcal{D}_{exp} .

4466 **Remark I.2.** We remark that by [Theorem 3](#), π_h^* is unique up to a measure zero set of $\rho_{\text{m},h}$ as distributed as above, and
 4467 $\pi_{\text{dec},h}^*$ is unique almost surely for $\tilde{\rho}_{\text{m},h}$ distributed as above. In particular, since the latter has density with respect to the
 4468 Lebesgue measure and infinite support, $\pi_{\text{dec},h}^*$ is unique in a Lebesgue almost everywhere sense.

4469 **Instantiation of the distance $d_{\mathcal{A}}$ for pairs of actions.** We recall the instantiation of the distance $d_{\mathcal{A}}$:

4470 **Definition I.4** (Instantiation of $d_{\mathcal{A}}$). We recall $d_{\mathcal{A}} : \mathcal{A} \times \mathcal{A} \rightarrow \mathbb{R}_{\geq 0}$ as defined in [\(4.2\)](#):

$$4471 \quad d_{\mathcal{A}}(\mathbf{a}, \mathbf{a}') := c_1 \max_{1 \leq k \leq \tau_c} (\|\bar{\mathbf{u}}_k - \bar{\mathbf{u}}'_k\| + \|\bar{\mathbf{x}}_k - \bar{\mathbf{x}}'_k\| + \|\bar{\mathbf{K}}_k - \bar{\mathbf{K}}'_k\|) + \mathbf{I}_{0,\infty}\{\mathcal{E}\},$$

4472 where we define $\mathcal{E} := \{\max_{1 \leq k \leq \tau_c} \max\{\|\bar{\mathbf{u}}_k - \bar{\mathbf{u}}'_k\|, \|\bar{\mathbf{x}}_k - \bar{\mathbf{x}}'_k\|, \|\bar{\mathbf{K}}_k - \bar{\mathbf{K}}'_k\|\} \leq c_2\}$, $\mathbf{I}_{0,\infty}$ is the indicator taking infinite
 4473 value when the event fails to hold, and c_1 and c_2 are constants depending polynomially on all problem parameters, given in
 4474 [Definition G.7](#).

4475 I.1.1. PRELIMINARIES FOR JOINT-DISTRIBUTION IMITATION.

4476 This section introduces a further *joint imitation gap*, which we can make small under a stronger bounded-memory assumption
 4477 on \mathcal{D}_{exp} stated below.

4478 **Definition I.5** (Joint Imitation Gap). Given a chunking policy π' , we let

$$4479 \quad \mathcal{L}_{\text{joint},\varepsilon}(\pi) := \inf_{\mu} \mathbb{P}_{\mu} \left[\max_{t \in [T]} \max \{ \|\mathbf{x}_{t+1}^{\text{exp}} - \mathbf{x}_{t+1}^{\pi}\|, \|\mathbf{u}_t^{\text{exp}} - \mathbf{u}_t^{\pi}\| \} > \varepsilon \right],$$

4480 where the infimum is over all couplings between the distribution of ρ_T under \mathcal{D}_{exp} and that induced by the policy π .

4481 Controlling $\mathcal{L}_{\text{joint},\varepsilon}(\pi)$ requires various additional stronger assumptions (*which we do not require in [Theorem 1](#)*), one of
 4482 which is that the demonstrator has bounded memory:

4483 **Definition I.6.** We say that the demonstration distribution, synthesis oracle pair $(\mathcal{D}_{\text{exp}}, \text{synth})$ have τ -bounded memory if
 4484 under $\rho_T = (\mathbf{x}_{1:T+1}, \mathbf{u}_{1:T}) \sim \mathcal{D}_{\text{exp}}$ and $\mathbf{a}_{1:H} = \text{synth}(\rho_T)$, the conditional distribution of \mathbf{a}_h and $\mathbf{x}_{1:t_h-\tau}, \mathbf{u}_{1:t_h-\tau}$ are
 4485 conditionally independence given $(\mathbf{x}_{t_h-\tau+1:t_h}, \mathbf{u}_{t_h-\tau+1:t_h-1})$.

4486 We note that enforcing [Definition I.6](#) can be relaxed to a mixing time assumption (see [Remark I.4](#)). Moreover, we stress that
 4487 we *do not* need the condition in [Definition I.6](#) if we only seek imitation of marginal distributions (as captured by $\mathcal{L}_{\text{marg},\varepsilon}$
 4488 and $\mathcal{L}_{\text{fin},\varepsilon}$), as in [Theorem 1](#).

4489 I.1.2. TRANSLATING CONTROL IMITATION LOSSES TO COMPOSITE-MDP IMITATION GAPS

4490 **Lemma I.1.** Recall the imitation losses [Definitions 2.2 and I.5](#), and the composite-MDP imitation gaps [Definition 4.1](#).
 4491 Further consider, the substitutions defined in [Section 4.1](#), with π^* instantiated as in [Definition I.3](#). Given policies $\pi = (\pi_h)$
 4492 with $\pi_h : \mathcal{P}_{\tau_{\text{m}}-1} \rightarrow \mathcal{A}$, we can extend $\pi_h : \mathcal{S} = \mathcal{P}_{\tau_c} \rightarrow \mathcal{A}$ by the natural embedding of $\mathcal{P}_{\tau_{\text{m}}-1}$ into \mathcal{P}_{τ_c} . Then, for any
 4493 $\varepsilon > 0$,

$$4494 \quad \mathcal{L}_{\text{marg},\varepsilon}(\pi) \leq \Gamma_{\text{marg},\varepsilon}(\pi \parallel \pi^*).$$

4495

4510 If we instead consider the the substitutions defined in [Section 4.1](#), but set d_S to equal d_{IPS} , which only measures distance in
 4511 the final coordinate of each trajectory chunk $\rho_{c,h}$,

$$4512 \mathcal{L}_{\text{fin},\varepsilon}(\pi) \leq \Gamma_{\text{marg},\varepsilon}(\pi \parallel \pi^*), \quad d_S(\cdot, \cdot) \leftarrow d_{\text{IPS}}(\cdot, \cdot) \quad (\text{I.1})$$

4514 Finally, if \mathcal{D}_{exp} has $\tau \leq \tau_m$ -bounded memory,

$$4516 \mathcal{L}_{\text{joint},\varepsilon}(\pi) \leq \Gamma_{\text{joint},\varepsilon}(\pi \parallel \pi^*).$$

4518 *Proof.* Let's start with the first bound, let superscript exp denote objects from \mathcal{D}_{exp} and superscript π from \mathcal{D}_π , the
 4519 distribution induced by chunking policy π . Letting \inf_μ denote infima over couplings between the two, we have

$$\begin{aligned} 4520 \mathcal{L}_{\text{marg},\varepsilon}(\pi) &:= \max_{t \in [T]} \inf_\mu \left\{ \mathbb{P}_\mu \left[\|\mathbf{x}_{t+1}^{\text{exp}} - \mathbf{x}_{t+1}^\pi\| > \varepsilon \right], \mathbb{P}_\mu \left[\|\mathbf{u}_t^{\text{exp}} - \mathbf{u}_t^\pi\| > \varepsilon \right] \right\} \\ 4521 &:= \max_{t \in [T]} \inf_\mu \left\{ \mathbb{P}_\mu \left[\|\mathbf{x}_{t+1}^{\text{exp}} - \mathbf{x}_{t+1}^\pi\| \vee \|\mathbf{u}_t^{\text{exp}} - \mathbf{u}_t^\pi\| > \varepsilon \right] \right\} \\ 4522 &\leq \max_{h \in [H]} \inf_\mu \left\{ \mathbb{P}_\mu \left[\max_{0 \leq i \leq \tau_c} \|\mathbf{x}_{t_h-i}^{\text{exp}} - \mathbf{x}_{t_h-i}^\pi\| \vee \max_{1 \leq i \leq \tau_c} \|\mathbf{u}_{t_h-i}^{\text{exp}} - \mathbf{u}_{t_h-i}^\pi\| \right] \right\} \\ 4523 &\leq \max_{h \in [H]} \inf_\mu \left\{ \mathbb{P}_\mu \left[d_S(\rho_{c,h}^{\text{exp}}, \rho_{c,h}^\pi) \right] \right\}, \end{aligned}$$

4529 From [Lemma I.6](#), $\rho_{c,h}^{\text{exp}}$ has the same marginal distribution as $\rho_{c,h}^{\pi^*}$, the distribution induced by π^* in [Definition I.3](#). Note the
 4530 subtlety that the joint distribution of these may defer because π^* has limited trajectories. Still, letting $\inf_{\mu'}$ denote infimum
 4531 over couplings between \mathcal{D}_π and \mathcal{D}_{π^*} , equality of marginals suffices to ensure

$$4532 \mathcal{L}_{\text{marg},\varepsilon}(\pi) = \max_{h \in [H]} \inf_{\mu'} \left\{ \mathbb{P}_{\mu'} \left[d_S(\rho_{c,h}^{\pi^*}, \rho_{c,h}^\pi) \right] \right\},$$

4535 which is at most $\Gamma_{\text{marg},\varepsilon}(\pi \parallel \pi^*)$ by definition [Definition 4.1](#).

4536 For the final-state imitation loss,

$$\begin{aligned} 4537 \mathcal{L}_{\text{fin},\varepsilon}(\pi) &:= \inf_\mu \mathbb{P}_\mu \left[\|\mathbf{x}_{T+1}^{\text{exp}} - \mathbf{x}_{T+1}^\pi\| > \varepsilon \right] \\ 4538 &\leq \max_{h \in [H]} \inf_\mu \left\{ \mathbb{P}_\mu \left[d_{\text{IPS}}(\rho_{c,h}^{\text{exp}}, \rho_{c,h}^\pi) \right] \right\}, \end{aligned}$$

4542 where again d_{IPS} only measures error in the final state of $\rho_{c,h}$. The corresponding bound in [\(I.1\)](#) follows similarly.

4543 Finally, we have

$$4544 \mathcal{L}_{\text{joint},\varepsilon}(\pi) := \inf_\mu \mathbb{P}_\mu \left[\max_{t \in [T]} \max \left\{ \|\mathbf{x}_{t+1}^{\text{exp}} - \mathbf{x}_{t+1}^\pi\|, \|\mathbf{u}_t^{\text{exp}} - \mathbf{u}_t^\pi\| \right\} > \varepsilon \right],$$

4548 When \mathcal{D}_{exp} has $\tau \leq \tau_m$ -bounded memory, then, the expert and π^* -induced trajectories are identically distributed. Therefore,
 4549 directly from this observation and [Definition 4.1](#),

$$4550 \mathcal{L}_{\text{joint},\varepsilon}(\pi) = \inf_\mu \mathbb{P}_\mu \left[\max_{t \in [T]} \max \left\{ \|\mathbf{x}_{t+1}^{\pi^*} - \mathbf{x}_{t+1}^\pi\|, \|\mathbf{u}_t^{\pi^*} - \mathbf{u}_t^\pi\| \right\} > \varepsilon \right] \leq \Gamma_{\text{joint},\varepsilon}(\pi \parallel \pi^*).$$

4553 □

4555 I.2. Proof of [Theorem 1](#) and a general reduction

4556 We now state a reduction from which [Theorem 1](#) is readily derived from our statistical learning analysis of score estimation.

4557 **Theorem 8** (Reduction from trajectory imitation to conditional sampling). *Consider applying TODA with $\sigma > 0$, and let*
 4558 *$\delta \in (0, 1)$, and define*

$$4560 \Delta_h(\varepsilon) := \mathbb{E}_{\rho_{m,h} \sim \mathcal{D}_{\text{exp}}} \mathbb{E}_{\tilde{\rho}_{m,h} \sim \mathcal{N}(\rho_{m,h}, \sigma^2 \mathbf{I})} \inf_{\mu \in \mathcal{C}(\pi_{\text{dec},h}^*, \tilde{\pi}_h(\tilde{\rho}_{m,h}))} \mathbb{P}_{(\mathbf{a}, \mathbf{a}') \sim \mu} [d_A(\mathbf{a}, \mathbf{a}')]. \quad (\text{I.2})$$

4562 where d_A is as in [Definition I.4](#), $\rho_{m,h} \sim \mathcal{D}_{\text{exp}}$ is shorthand for $\rho_T \sim \mathcal{D}_{\text{exp}}$, and $\rho_{m,h}$ denotes the corresponding h -th
 4563 memory chunk of \mathcal{D}_{exp} . Consider the following setup:

4564

- Suppose that [Assumptions 3.1 and 3.2](#) hold.
- Let c_1, \dots, c_5 be the constants defined [Definition G.7](#), which we recall are polynomial in the terms in [Assumptions 3.1 and 3.2](#)
- Define $d = \tau_c(d_u + d_x + d_u d_x)$.
- Suppose that $\tau_c \geq c_3/\eta$
- The parameters $\varepsilon, \sigma > 0$ satisfy $5d_x + \log\left(\frac{4\sigma}{\varepsilon}\right) \leq c_4^2/(16\sigma^2)$,

For all we have

$$\begin{aligned} & \mathcal{L}_{\text{marg}, \varepsilon_1}(\hat{\pi}_\sigma) \vee \mathcal{L}_{\text{fin}, \varepsilon_2}(\hat{\pi}_\sigma) \\ & \leq H \left(\frac{2\varepsilon}{\sigma} + 6c_5 \sqrt{5d_x + 2 \log\left(\frac{4\sigma}{\varepsilon}\right)} e^{-\frac{\eta(\tau_c - \tau_m)}{L_{\text{stab}}}} \right) + \sum_{h=1}^H \Delta_h(\varepsilon) \end{aligned}$$

where

$$\begin{aligned} \varepsilon_1 &= \varepsilon + 4c_5\sigma \cdot \sqrt{5d_x + 2 \log\left(\frac{4\sigma}{\varepsilon}\right)} \\ \varepsilon_2 &= \varepsilon + 4c_5\sigma e^{-\frac{\eta\tau_c}{L_{\text{stab}}}} \cdot \sqrt{5d_x + 2 \log\left(\frac{4\sigma}{\varepsilon}\right)} \end{aligned} \tag{I.3}$$

We first demonstrate how [Theorem 1](#) follows from [Theorem 8](#) and [Theorem 6](#):

Proof of [Theorem 1](#). From [Theorem 8](#), it suffices to show that with probability at least $1 - \delta$, it holds that $\Delta_h \leq \frac{3\varepsilon}{\sigma}$ for all $h \in [H]$. Note that by [Assumption 3.2](#) it holds \mathcal{D}_{exp} -almost surely that $\|a_h\| \leq R_{\text{stab}}$ and thus the condition on q in [Theorem 6](#) holds for $R = R_{\text{stab}}$. Moreover, for $d = \tau_c(d_x + d_u + d_x d_u)$, we have that $a \in \mathbb{R}^d$. By [Assumption 3.3](#), the conditions on the score class \mathfrak{s}_θ hold for us to apply [Theorem 6](#). Note that by assumption,

$$N_{\text{exp}} \geq c \left(\frac{C_\Theta d R (R \vee \sqrt{d}) \log(dn)}{(\varepsilon/\sigma)^4} \right)^{4\nu} \vee \left(\frac{d^6 (R^4 \vee d^2 \log^3(\frac{HndR\sigma}{\delta\varepsilon}))}{(\varepsilon/\sigma)^{24}} d^2 \right)^{4\nu},$$

where we note that the right hand side is $\text{poly}(C_\Theta, \varepsilon/\sigma, R_{\text{stab}}, d, \log(H/\delta))^\nu$, and J and α are set as in [\(H.1\)](#). Taking a union bound over $h \in [H]$ and applying [Theorem 6](#) tells us that with probability at least $1 - \delta$, for all $h \in [H]$, it holds that

$$\mathbb{E}_{\rho_{m,h} \sim q_{\rho_{m,h}}} \left[\inf_{\mu \in \mathcal{C}(\text{DDPM}(\mathfrak{s}_{\hat{\theta}}, \rho_{m,h}), q(\cdot | \rho_{m,h}))} \mathbb{P}_{(\hat{a}, a^*) \sim \mu} (\|\hat{a} - a^*\| \geq \varepsilon/\sigma) \right] \leq \frac{3\varepsilon}{\sigma}.$$

Thus it holds that with probability at least $1 - \delta$,

$$\sum_{h=1}^H \Delta_h(\varepsilon/\sigma) \leq \frac{3H\varepsilon}{\sigma}.$$

Plugging this in to [Theorem 8](#) concludes the proof. \square

Proof of [Theorem 8](#). Lets begin by bounding $\mathcal{L}_{\text{marg}, \varepsilon}(\pi)$. Recall the definitions of $d_S, d_{\text{TVC}}, d_{\text{IPS}}$ in [Section 4](#), and let $\mathbf{s}_{1:H+1}^*$ and $\mathbf{s}_{1:H+1}$ denote the composite states corresponding to a trajectory $(\mathbf{x}_{1:T+1}^{\pi^*}, \mathbf{u}_{1:T}^{\pi^*})$ under π^* and $(\mathbf{x}_{1:T+1}^\pi, \mathbf{u}_{1:T}^\pi)$, respectively, under the instantiation of the composite MDP in [Section 4.1](#). We can view π^* and π (which depend only on memory chunks $\rho_{m,h}$) as policies in the composite MDP which are compatible with the decomposition [Definition E.1](#). We make the following points:

- In light of [Lemma I.1](#),

$$\mathcal{L}_{\text{marg},\varepsilon_1}(\pi \parallel \pi^*) \leq \Gamma_{\text{marg},\varepsilon_1}(\pi \parallel \pi^*).$$

- By [Lemma I.8](#), a consequence of Pinsker's inequality, it holds that the Gaussian kernel W_σ used in TODA is γ_σ -TVC (w.r.t. d_{TVC}) with $\gamma_\sigma(u) = \frac{u\sqrt{\tau_m+1}}{2\sigma}$.

- Note that $d_{\text{IPS}}(s_h, s'_h) = \|\mathbf{x}_{t_h} - \mathbf{x}'_{t_h}\|$ measures Euclidean distance between the last \mathbf{x} -coordinates of s_h, s'_h . Moreover, if $s'_h \sim W_\sigma(s_h)$ the last coordinate \mathbf{x}'_{t_h} of s' is distributed as $\mathcal{N}(\mathbf{x}_{t_h}, \sigma^2 I)$. By [Lemma I.7](#) with $d = d_x$, that for

$$r = 2\sigma \cdot \sqrt{5d_x + 2 \log\left(\frac{1}{p}\right)}$$

$$\mathbb{P}_{s' \sim W_\sigma(s)}[d_{\text{IPS}}(s, s') > r] \leq p.$$

- As (a) s_h^* corresponds to $\rho_{c,h}$ from $\rho_T \sim \mathcal{D}_{\text{exp}}$, (b) as $\hat{\pi}, \pi_{\text{dec}}^*$ are functions of $\rho_{m,h}$, and (c) by recalling the definition of $d_{\text{os},\varepsilon}$ in [Definition 4.1](#),

$$\begin{aligned} & \mathbb{E}_{s_h^* \sim P_h^*} \mathbb{E}_{\tilde{s}_h^* \sim W_\sigma(s_h^*)} d_{\text{os},\varepsilon}(\pi_h(\tilde{s}_h^*) \parallel \pi_{\text{dec}}^*(\tilde{s}_h^*)) \\ &= \mathbb{E}_{\rho_{m,h} \sim \mathcal{D}_{\text{exp}}} \mathbb{E}_{\tilde{\rho}_{m,h} \sim \mathcal{N}(\rho_{m,h}, \sigma^2 I)} \inf_{\mu \in \mathcal{C}(\pi_{\text{dec}}^*(\tilde{\rho}_{m,h}), \hat{\pi}(\tilde{\rho}_{m,h}))} \mathbb{P}_{(a,a') \sim \mu}[d_{\mathcal{A}}(a, a') \geq \varepsilon], \end{aligned}$$

which is at most $\Delta_h(\varepsilon_0)$ by assumption.

- Finally, [Proposition 4.1](#) ensures that under our assumption $\tau_c \geq c_3/\eta$, and let $r_{\text{IPS}} = c_4$, $\gamma_{\text{IPS},1}(u) = c_5 u \exp(-\eta(\tau_c - \tau_m)/L_{\text{stab}})$, $\gamma_{\text{IPS},2}(u) = c_5 u$ for c_3, c_4, c_5 given in [Definition G.7](#). Then, for $d_{\mathcal{S}}, d_{\text{TVC}}, d_{\text{IPS}}$ as above, we have that π^* is $(\gamma_{\text{IPS},1}, \gamma_{\text{IPS},2}, d_{\text{IPS}}, r_{\text{IPS}})$ -IPS.

Consequently, for $r = 2\sigma \cdot \sqrt{5d_x + 2 \log\left(\frac{4\sigma}{\varepsilon}\right)} \in (0, \frac{1}{2}r_{\text{IPS}})$, [Theorem 5](#) (which, we recall, generalizes [Theorem 2](#) to account for the direct decomposition structure) implies

$$\begin{aligned} \mathcal{L}_{\text{marg},\varepsilon+2rc_5}(\hat{\pi}_\sigma) &= \mathcal{L}_{\text{marg},\varepsilon+2rc_5}(\hat{\pi}_\sigma \parallel \pi^*) \leq \Gamma_{\text{marg},\varepsilon+2rc_5}(\hat{\pi}_\sigma \parallel \pi^*) \\ &\leq H\sqrt{2\tau_m-1} \left(\frac{\varepsilon}{2\sigma} + \frac{3}{2\sigma} \left(\max \left\{ \varepsilon, 2rc_5 e^{-\frac{\eta(\tau_c-\tau_m)}{L_{\text{stab}}}} \right\} \right) \right) \\ &\quad + \sum_{h=1}^H \mathbb{E}_{s_h^* \sim P_h^*} \mathbb{E}_{\tilde{s}_h^* \sim W_\sigma(s_h^*)} d_{\text{os},\varepsilon}(\pi_h(\tilde{s}_h^*) \parallel \pi_{\text{dec}}^*(\tilde{s}_h^*)) \\ &\leq H\sqrt{2\tau_m-1} \left(\frac{2\varepsilon}{\sigma} + 6c_5 \sqrt{5d_x + 2 \log\left(\frac{4\sigma}{\varepsilon}\right)} e^{-\frac{\eta(\tau_c-\tau_m)}{L_{\text{stab}}}} \right) + \sum_{h=1}^H \Delta_h(\varepsilon) \end{aligned}$$

Substituting in $\varepsilon_1 = \varepsilon + 2rc_5 = \varepsilon + 4c_5\sigma \cdot \sqrt{5d_x + 2 \log\left(\frac{4\sigma}{\varepsilon}\right)}$ the bound on $\mathcal{L}_{\text{marg},\varepsilon_1}$ is proved.

To show $\mathcal{L}_{\hat{\pi}_m,\varepsilon_2}(\hat{\pi}_\sigma)$ satisfies the same bound, we replace $d_{\mathcal{S}}$ in the above argument (as defined in [Section 4.1](#)) with $d_{\mathcal{S}}(\cdot, \cdot) \leftarrow d_{\text{IPS}}(\cdot, \cdot)$, where again we recall that $d_{\text{IPS}}(s_s, s'_s) = \|\mathbf{x}_{t_h} - \mathbf{x}'_{t_h}\|$ measures differences in the final associated control state. From [Corollary G.1](#), which is a generalization of [Proposition 4.1](#), it follows that we can replace $\gamma_{\text{IPS},2}(u) = c_5 u$ as used above with the considerable smaller quantity $\gamma_{\text{IPS},2}(u) = c_5 u e^{-\frac{\eta\tau_c}{L_{\text{stab}}}}$. Thus, we can replace ε_1 above with $\varepsilon_2 := \varepsilon + 4c_5 e^{-\eta\tau_c/L_{\text{stab}}} \sigma \cdot (5d_x + 2 \log\left(\frac{1}{\varepsilon}\right))^{1/2}$. This concludes the proof that

$$\mathcal{L}_{\text{marg},\varepsilon_2}(\hat{\pi}_\sigma) \leq H\sqrt{2\tau_m-1} \left(6c_5 \sqrt{5d_x + 2 \log\left(\frac{4\sigma}{\varepsilon}\right)} e^{-\frac{\eta(\tau_c-\tau_m)}{L_{\text{stab}}}} + \frac{2\varepsilon}{\sigma} \right),$$

as needed. \square

4675 I.3. Imitation of the joint trajectory under total variation continuity of demonstrator policy

4676 Here, we show that if the demonstrator policy satisfies a certain continuity property in total variation distance, then we
 4677 can imitate the *joint distribution* over trajectories, not just marginals. Recall the joint imitation loss from $\mathcal{L}_{\text{joint},\varepsilon}$ from
 4678 [Definition I.5](#).

4680 **Theorem 9.** Consider the setting [Theorem 8](#), with $\Delta_h(\varepsilon)$ as in [\(I.2\)](#), and suppose all the assumptions of that theorem are
 4681 met. Suppose that, in addition, there is a strictly increasing function $\gamma(\cdot)$ such that for all $\rho_{m,h}, \rho'_{m,h} \in \mathcal{P}_{\tau_m-1}$,

$$4682 \text{TV}(\pi^*(\rho_{m,h}), \pi^*(\rho'_{m,h})) \leq \gamma(\|\rho_{m,h} - \rho'_{m,h}\|),$$

4684 where π^* is defined is the conditional in [Definition I.3](#). Further, suppose that \mathcal{D}_{exp} has $\tau \leq \tau_m$ bounded memory
 4685 ([Definition I.6](#)). Then, with $\varepsilon_1 := \varepsilon + 4c_5\sigma \cdot \sqrt{5d_x + 2 \log \left(\frac{4\sigma}{\varepsilon}\right)}$ as in [\(I.3\)](#),

$$4687 \mathcal{L}_{\text{joint},\varepsilon_1}(\hat{\pi}_\sigma) \leq H \cdot \text{ERRTVC}(\sigma, \gamma) \\
 4688 + H\sqrt{2\tau_m - 1} \left(\frac{2\varepsilon}{\sigma} + 6c_5\sqrt{5d_x + 2 \log \left(\frac{4\sigma}{\varepsilon}\right)} e^{-\frac{\eta(\tau_c - \tau_m)}{L_{\text{stab}}}} \right) + \sum_{h=1}^H \Delta_h(\varepsilon),$$

4692 where we define $d_0 = \tau_m d_x + (\tau_m - 1)d_u$ and $u_0 = \gamma(8\sigma\sqrt{d_0 \log(9)})$, and

$$4694 \text{ERRTVC}(\sigma, \gamma) = \begin{cases} 2c\sigma\sqrt{d_0} & \text{linear } \gamma(u) = c \cdot u, c > 0 \\ u_0 + \int_{u_0}^{\infty} e^{-\frac{\gamma^{-1}(u)^2}{64\sigma^2}} du & \text{general } \gamma(\cdot) \end{cases}. \quad (\text{I.4})$$

4698 In particular, under [Assumption 3.3](#), if

$$4700 N_{\text{exp}} \geq c \left(\frac{C_\Theta dR(R \vee \sqrt{d}) \log(dn)}{(\varepsilon/\sigma)^4} \right)^{4\nu} \vee \left(\frac{d^6 (R^4 \vee d^2 \log^3 \left(\frac{HndR\sigma}{\delta\varepsilon}\right))}{(\varepsilon/\sigma)^{24}} d^2 \right)^{4\nu},$$

4703 then with probability at least $1 - \delta$, it holds that

$$4704 \mathcal{L}_{\text{joint},\varepsilon_1}(\hat{\pi}_\sigma) \leq H \cdot \text{ERRTVC}(\sigma, \gamma) + H\sqrt{2\tau_m - 1} \left(\frac{3\varepsilon}{\sigma} + 6c_5\sqrt{5d_x + 2 \log \left(\frac{4\sigma}{\varepsilon}\right)} e^{-\frac{\eta(\tau_c - \tau_m)}{L_{\text{stab}}}} \right).$$

4708 **Remark I.3.** The second term in our bound on $\mathcal{L}_{\text{joint},\varepsilon}(\pi)$ is identical to the bound in [Theorem 8](#). The term ERRTVC
 4709 captures the additional penalty we pay to strengthen for imitation of marginals to imitation of joint distributions. Notice that
 4710 if $\lim_{u \rightarrow 0} \gamma(u) \rightarrow 0$ and $\gamma(u)$ is sufficiently integrable, then, $\lim_{\sigma \rightarrow 0} \text{Err}(\sigma, \gamma) = 0$. This is most clear in the linear $\gamma(\cdot)$
 4711 case, where $\text{Err}(\sigma, \gamma) = \mathcal{O}(\sigma)$.

4713 The proof is given in [Appendix I.3.1](#); it mirrors that of [Theorem 8](#), but replaces [Theorem 2](#) with the following imitation
 4714 guarantee in the composite MDP abstraction of [Section 4](#), which bounds the joint imitation gap relative to π^* if π^* is TVC.

4715 **Proposition I.2.** Consider the set-up of [Section 4](#), and suppose that the assumptions of [Theorem 5](#), but that, in addition, the
 4716 expert policy π^* is $\tilde{\gamma}(\cdot)$ -TVC with respect to the pseudometric d_{TVC} , where $\tilde{\gamma} : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}_{\geq 0}$ is strictly increasing. Then, for
 4717 all parameters as in [Theorem 2](#), and any $\tilde{r} > 0$,

$$4718 \Gamma_{\text{joint},\varepsilon}(\hat{\pi} \circ W_\sigma \parallel \pi^*) \leq H \int_0^\infty \max_s \mathbb{P}_{s' \sim W_\sigma(s)}[d_{\text{TVC}}(s, s') > \tilde{\gamma}^{-1}(u)/2] du \\
 4719 + H(2p_r + 3\gamma_\sigma(\max\{\varepsilon, \gamma_{\text{IPS},1}(2r)\})) + \sum_{h=1}^H \mathbb{E}_{s_h^* \sim P_h^*} \mathbb{E}_{\tilde{s}_h^* \sim W_\sigma(s_h^*)} d_{\text{os},\varepsilon}(\hat{\pi}_h(\tilde{s}_h^*) \parallel \pi_{\text{dec}}^*(\tilde{s}_h^*)),$$

4722 where the *term in color* on the first line is the only term that differs from the bound in [Theorem 2](#).

4724 Moreover, in the special case where all of the distributions of $d_{\text{TVC}}(s, s') \mid s' \sim W_\sigma(s)$ are stochastically dominated by a
 4725 common random variable Z , and further more $\tilde{\gamma}(u) = \tilde{c} \cdot u$ for some constant \tilde{c} , then our bound may be simplified to

$$4727 \Gamma_{\text{joint},\varepsilon}(\hat{\pi} \circ W_\sigma \parallel \pi^*) \leq 2\tilde{c}H\mathbb{E}[Z] \\
 4728 + H(2p_r + 3\gamma_\sigma(\max\{\varepsilon, \gamma_{\text{IPS},1}(2r)\})) + \sum_{h=1}^H \mathbb{E}_{s_h^* \sim P_h^*} \mathbb{E}_{\tilde{s}_h^* \sim W_\sigma(s_h^*)} d_{\text{os},\varepsilon}(\hat{\pi}_h(\tilde{s}_h^*) \parallel \pi_{\text{dec}}^*(\tilde{s}_h^*)).$$

4730 *Proof Sketch.* [Proposition I.2](#) is derived below in [Appendix I.3.2](#). It is corollary of [Theorem 2](#), combined with adjoining the
4731 coupling constructed therein to a TV distance coupling between $\pi_{\odot\sigma}^*$ (whose joints we *can always* imitate) and π^* . Coupling
4732 trajectories induced by $\pi_{\odot\sigma}^*$ and π^* relies on the TVC of π^* , as well as concentration of W_σ . \square
4733

4734 Using the above proposition, we can derive the following consequences for imitation of the joint distribution.
4735

4736 I.3.1. PROOF OF THEOREM 9

4737 The proof is nearly identical to that of [Theorem 8](#), with the modifications that we replace our use of [Theorem 2](#) with
4738 [Proposition I.2](#) taking $\tilde{\gamma} \leftarrow \gamma$. By [Lemma I.1](#) and the assumption that \mathcal{D}_{exp} has $\tau \leq \tau_m$ -bounded memory, it suffices to
4739 bound the joint-gap in the composite MDP:
4740

$$4741 \mathcal{L}_{\text{joint},\varepsilon}(\pi) \leq \Gamma_{\text{joint},\varepsilon}(\pi \parallel \pi^*).$$

4743 We bound this directly from [Proposition I.2](#). The final statement follows from [Theorem 6](#) in the same way that it does in the
4744 proof of [Theorem 1](#).
4745

4746 The only remaining modification, then, is to evaluate the additional additive terms colored in purple in [Proposition I.2](#); we will
4747 show that ERRTVC as defined in [\(I.4\)](#) suffices as an upper bound. We have two cases. In both, let $d_0 = \tau_m d_x + (\tau_m - 1)d_u$.
4748 As d_{TVC} measures the distance between the chunks $\mathbf{p}_{m,h} = \phi_{\mathcal{Z}}(s_h)$, $\hat{\mathbf{p}}_{m,h} = \phi_{\mathcal{Z}}(s'_h)$, which have dimension d_0 , and since
4749 we $\phi_{\mathcal{Z}} \circ W_\sigma(\cdot) = \mathcal{N}(\cdot, \sigma^2 \mathbf{I}_{d_0})$, we have
4750

$$4751 d_{\text{TVC}}(\phi_{\mathcal{Z}} \circ s, \phi_{\mathcal{Z}} \circ s') \mid s' \sim W_\sigma(s) \stackrel{\text{dist}}{=} \|\gamma\|, \quad \gamma \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_{d_0}) \quad (\text{I.5})$$

4753 **General $\gamma(\cdot)$.** Eq. [\(I.5\)](#) and [Lemma I.7](#) imply that
4754

$$4755 \mathbb{P}_{s' \sim W_\sigma(s)}[d_{\text{TVC}}(s, s')] \leq \exp(-r^2/16\sigma^2), \quad r \geq 4\sigma d_0 \log(9).$$

4756 Hence, if $u_0 = \gamma(8\sigma d_0 \log(9))$, then
4757

$$4758 \mathbb{P}[d_{\text{TVC}}(s, s') > \gamma^{-1}(u)/2] \leq \exp(-\gamma^{-1}(u)^2/64\sigma^2), \quad u \geq u_0.$$

4760 Thus, as probabilities are at most one,
4761

$$4762 \int_0^\infty \max_s \mathbb{P}_{s' \sim W_\sigma(s)}[d_{\text{TVC}}(s, s') > \gamma^{-1}(u)/2] du \leq u_0 + \int_{u_0}^\infty e^{-\frac{\gamma^{-1}(u)^2}{64\sigma^2}} du,$$

4765 as needed.
4766

4767 **Linear $\gamma(\cdot)$.** In the special case where $\gamma(u) = c(u)$, [Eq. \(I.5\)](#) implies that we can take $Z = \|\gamma\|$ where $\gamma \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_{d_0})$
4768 in the second part of [Proposition I.2](#). The corresponding additive term is then $2Hc\mathbb{E}[\|\gamma\|]$. By Jensen's inequality,
4769 $\mathbb{E}[\|\gamma\|] \leq \sqrt{\mathbb{E}[\|\gamma\|^2]} = \sqrt{\sigma^2 d_0} = \sigma\sqrt{d_0}$, as needed. \square
4770

4771 I.3.2. PROOF OF PROPOSITION I.2

4772 Define the shorthand
4773

$$4774 B := H(2p_r + 3\gamma_\sigma(\max\{\varepsilon, \gamma_{\text{IPS},1}(2r)\})) + \sum_{h=1}^H \mathbb{E}_{s_h^* \sim \mathbf{P}_h^*} \mathbb{E}_{\tilde{s}_h^* \sim W_\sigma(s_h^*)} d_{\text{os},\varepsilon}(\hat{\pi}_h(\tilde{s}_h^*) \parallel \pi_{\text{dec}}^*(\tilde{s}_h^*)),$$

4776 and recall that [Theorem 2](#) ensures $\Gamma_{\text{joint},\varepsilon}(\hat{\pi} \circ W_\sigma \parallel \pi_{\odot\sigma}^*) \leq B$. Further, recall from [Definition 4.1](#) that
4777

$$4778 \Gamma_{\text{joint},\varepsilon}(\hat{\pi} \circ W_\sigma \parallel \pi_{\odot\sigma}^*) = \inf_{\mu_1} \mathbb{P}_{\mu_1} \left[\max_{h \in [H]} \max\{d_{\mathcal{S}}(s_{h+1}^\odot, \hat{s}_{h+1}), d_{\mathcal{A}}(\mathbf{a}_h^\odot, \hat{\mathbf{a}}_h)\} > \varepsilon \right],$$

4781 where the infimum is over all couplings μ_1 of $(\hat{s}_{1:H+1}, \hat{\mathbf{a}}_{1:H}) \sim D_{\hat{\pi} \circ W_\sigma}$ and $(s_{1:H+1}^\odot, \mathbf{a}_{1:H}^\odot) \sim D_{\pi_{\odot\sigma}^*}$ with $\mathbb{P}_{\mu_1}[\hat{s}_1 = s_1^\odot] = 1$.
4782 For any coupling μ_1 , we can consider another coupling μ_2 of $(s_{1:H+1}^*, \mathbf{a}_{1:H}^*) \sim D_{\pi^*}$ and $(s_{1:H+1}^\odot, \mathbf{a}_{1:H}^\odot) \sim D_{\pi_{\odot\sigma}^*}$ with
4783 $\mathbb{P}_{\mu_2}[s_1^* = s_1^\odot] = 1$. By the ‘‘gluing lemma’’ ([Lemma C.2](#)), we can construct a combined coupling μ which respects the
4784

4785 marginals of μ_1 and μ_2 . This combined coupling induces a joint coupling $\tilde{\mu}_1$ of $D_{\hat{\pi} \circ W_\sigma}$ and D_{π^*} which, by a union bound,
 4786 satisfies $\mathbb{P}_{\tilde{\mu}_1}[\hat{s}_1 = s_1^*] = 1$. Thus, by a union bound, we can bound

$$\begin{aligned}
 4787 \Gamma_{\text{joint},\varepsilon}(\hat{\pi} \circ W_\sigma \parallel \pi^*) &\leq \mathbb{P}_{\tilde{\mu}_1} \left[\max_{h \in [H]} \max\{d_S(s_{h+1}^*, \hat{s}_{h+1}), d_A(a_h^*, \hat{a}_h)\} > \varepsilon \right] \\
 4788 &\leq \mathbb{P}_{\mu_1} \left[\max_{h \in [H]} \max\{d_S(s_{h+1}^\circ, \hat{s}_{h+1}), d_A(a_h^\circ, \hat{a}_h)\} > \varepsilon \right] \\
 4789 &\quad + \mathbb{P}_{\mu_2} \left[(s_{1:H+1}^*, a_{1:H}^*) \neq (s_{1:H+1}^\circ, a_{1:H}^\circ) \right].
 \end{aligned}$$

4794 Passing to the infimum over μ_1, μ_2 ,

$$4795 \Gamma_{\text{joint},\varepsilon}(\hat{\pi} \circ W_\sigma \parallel \pi^*) \leq \underbrace{\Gamma_{\text{joint},\varepsilon}(\hat{\pi} \circ W_\sigma \parallel \pi_{\circlearrowleft}^*)}_{\leq B} + \inf_{\mu_2} \mathbb{P}_{\mu_2} \left[(s_{1:H+1}^*, a_{1:H}^*) \neq (s_{1:H+1}^\circ, a_{1:H}^\circ) \right],$$

4796 where again μ_2 quantify couplings of $(s_{1:H+1}^*, a_{1:H}^*) \sim D_{\pi^*}$ and $(s_{1:H+1}^\circ, a_{1:H}^\circ) \sim D_{\pi_{\circlearrowleft}^*}$ with $\mathbb{P}_{\mu_2}[s_1^* = s_1^\circ] = 1$. Bounding
 4797 the infimum over μ_2 with [Proposition I.4](#), we have

$$4800 \Gamma_{\text{joint},\varepsilon}(\hat{\pi} \circ W_\sigma \parallel \pi^*) \leq B + \sum_{h=1}^H \mathbb{E}_{s_h^*} \text{TV}(\pi_h^*(s_h^*), \pi_{\circlearrowleft,\sigma,h}^*(s_h^*))$$

4804 To conclude, it suffices to show the following bound:

4805 **Claim I.3.** For any $s \in \mathcal{S}$, $h \in [H]$, and $\tilde{r} \geq 0$, $\text{TV}(\pi_h^*(s), \pi_{\circlearrowleft,\sigma,h}^*(s)) \leq \int_0^\infty \max_s \max_{s'} \mathbb{P}_{s' \sim W_\sigma(s)}[d_{\text{TV}}(s, s') >$
 4806 $\tilde{\gamma}^{-1}(u)/2]$.

4807 *Proof.* To show this claim, we note that we can represent (via the notation in [Appendix E.3](#)) $\pi_{\circlearrowleft,\sigma,h}^*(s) = \pi_h^* \circ W_{\circlearrowleft,h}^*(s)$, where
 4808 $W_{\circlearrowleft,h}^*$ is the replica-kernel defined in [Definition E.5](#). Thus, we can construct a coupling of $a^* \sim \pi_h^*(s)$ and $a^\circ \sim \pi_{\circlearrowleft,\sigma,h}^*(s)$
 4809 by introducing an intermediate state $s' \sim W_{\circlearrowleft,h}^*(s)$ and $a^\circ \sim \pi^*(s')$. By [Lemma C.4](#), the fact that TV distance is bounded
 4810 by one, and the assumption that π^* is $\tilde{\gamma}$ -TVC, we then have

$$4811 \text{TV}(\pi_h^*(s), \pi_{\circlearrowleft,\sigma,h}^*(s)) \leq \mathbb{E}_{s' \sim W_{\circlearrowleft,h}^*(s)} \text{TV}(\pi_h^*(s), \pi_h^*(s')).$$

4812 Recall the well-known formula that, for a non-negative random variable X , $\mathbb{E}[X] = \int_0^\infty \mathbb{P}[X > u] du$ ([Durrett, 2019](#)). From
 4813 this formula, we find

$$\begin{aligned}
 4814 \text{TV}(\pi_h^*(s), \pi_{\circlearrowleft,\sigma,h}^*(s)) &\leq \int_0^\infty \mathbb{P}[\text{TV}(\pi_h^*(s), \pi_h^*(s')) > u] du \\
 4815 &\stackrel{(i)}{\leq} \int_0^\infty \mathbb{P}[d_{\text{TV}}(s, s') > \tilde{\gamma}^{-1}(u)] du
 \end{aligned}$$

4816 where in (i) we used that $\text{TV}(\pi_h^*(s), \pi_h^*(s')) \leq \tilde{\gamma}(d_{\text{TV}}(s, s'))$ and that, as $\tilde{\gamma}(\cdot)$ is strictly increasing, we have the equality
 4817 of events $\{\text{TV}(\pi_h^*(s), \pi_h^*(s')) > u\} = \{d_{\text{TV}}(s, s') > \tilde{\gamma}^{-1}(u)\}$. Arguing as in the proof of [Lemma E.5](#), we have that
 4818 $\mathbb{P}_{s' \sim W_\sigma(s)}[d_{\text{TV}}(s, s') > \tilde{\gamma}^{-1}(u)] \leq \max_s \mathbb{P}_{s' \sim W_\sigma(s)}[d_{\text{TV}}(s, s') > \tilde{\gamma}^{-1}(u)/2]$. Hence, we conclude

$$4819 \text{TV}(\pi_h^*(s), \pi_{\circlearrowleft,\sigma,h}^*(s)) \leq \int_0^\infty \max_s \mathbb{P}_{s' \sim W_\sigma(s)}[d_{\text{TV}}(s, s') > \tilde{\gamma}^{-1}(u)/2] du$$

4820 which proves the first guarantee. \square

4821 With the above claim proven, we conclude the proof of the first statement of [Proposition I.2](#). For the second statement, we
 4822 observe that under the stated stochastic domination assumption by Z , and if $\tilde{\gamma}(u) = \tilde{c} \cdot u$, then $\max_s \mathbb{P}_{s' \sim W_\sigma(s)}[d_{\text{TV}}(s, s') >$
 4823 $\tilde{\gamma}^{-1}(u)/2] \leq \mathbb{P}[Z > \frac{u}{2\tilde{c}}]$. Hence, by a change of variables $u = \frac{t}{2\tilde{c}}$,

$$4824 \int_0^\infty \max_s \mathbb{P}_{s' \sim W_\sigma(s)}[d_{\text{TV}}(s, s') > \tilde{\gamma}^{-1}(u)/2] du \leq \int_0^\infty \mathbb{P}[Z > \frac{u}{2\tilde{c}}] = 2\tilde{c} \int_0^\infty \mathbb{P}[Z > u] = 2\tilde{c}\mathbb{E}[Z],$$

4825 where again we invoke that Z must be nonnegative (to stochastically dominate non-negative random variables), and thus
 4826 used the expectation formula referenced above. \square

4839

4840 I.4. Imitation in total variation distance

4841 Here, we notice that estimating the score in TV distance facilitates estimation in the composite MDP, with no smoothing:

4842 **Theorem 10.** For a chunking policy $\hat{\pi}$, suppose that there are terms $(\bar{\Delta}_h)_{1 \leq h \leq H}$ such that

$$4843 \mathbb{E}_{\boldsymbol{\rho}_{m,h} \sim \mathcal{D}_{\text{exp}}} \text{TV}(\pi^*(\boldsymbol{\rho}_{m,h}), \hat{\pi}(\boldsymbol{\rho}_{m,h})) \leq \bar{\Delta}_h,$$

4844 Then, under no additional assumption (not even those in Section 3), we have

$$4845 \mathcal{L}_{\text{fin}, \varepsilon=0}(\hat{\pi}) \leq \mathcal{L}_{\text{marg}, \varepsilon=0}(\hat{\pi}) \leq \sum_{h=1}^H \bar{\Delta}_h$$

4846 In addition π^* has τ -bounded memory (Definition 1.6) for $\tau \leq \tau_m$, then for $\mathcal{L}_{\text{joint}, \varepsilon}$ as in Definition 1.5,

$$4847 \mathcal{L}_{\text{joint}, \varepsilon=0}(\hat{\pi}) \leq \sum_{h=1}^H \bar{\Delta}_h$$

4848 The above theorem is a direct consequence of the result below in the composite MDP, together with the correct instantiations
4849 for control, and Lemma 1.1 to convert $\mathcal{L}_{\text{marg}, \varepsilon}$ and $\mathcal{L}_{\text{fin}, \varepsilon}$ into $\Gamma_{\text{marg}, \varepsilon} \leq \Gamma_{\text{joint}, \varepsilon}$, and $\Gamma_{\text{joint}, \varepsilon}$, respectively.

4850 **Proposition 1.4.** Consider the composite MDP setting of Section 4. Then, there exists a coupling

$$4851 \text{TV}(\mathbb{D}_{\hat{\pi}}, \mathbb{D}_{\pi^*}) \leq \sum_{h=1}^H \mathbb{E}_{\mathbf{s}_h^* \sim \mathbb{P}_h^*} \text{TV}(\pi_h^*(\mathbf{s}_h^*), \hat{\pi}_h(\mathbf{s}_h^*))$$

4852 Thus, there exists a couple $\mu \in \mathcal{C}(\mathbb{D}_{\pi^*}, \mathbb{D}_{\hat{\pi}})$ of $(\mathbf{s}_{1:H+1}^*, \mathbf{a}_{1:H}^*) \sim \mathbb{D}_{\pi^*}$ and $(\hat{\mathbf{s}}_{1:H+1}, \hat{\mathbf{a}}_{1:H}) \sim \mathbb{D}_{\hat{\pi}}$ such that
4853 $\mathbb{P}_{\mu}[(\mathbf{s}_{1:H+1}^*, \mathbf{a}_{1:H}^*) \neq (\hat{\mathbf{s}}_{1:H+1}, \hat{\mathbf{a}}_{1:H})]$ is bounded by the right-hand side of the above display. Moreover, this coupling
4854 can be constructed such that $\mathbb{P}_{\mu}[\mathbf{s}_1^* = \hat{\mathbf{s}}_1]$.

4855 *Proof of Proposition 1.4.* This is a direct consequence of Lemma 1.9, with $\mathbb{P}_1 \leftarrow \mathbb{P}_{\text{init}}$, and \mathbb{Q}_{h+1} corresponding to the
4856 kernel for sampling $\mathbf{a}_h^* \sim \pi^*(\mathbf{s}_h^*)$ and incrementing the dynamics $\mathbf{s}_{h+1}^* = F_h(\mathbf{s}_h^*, \mathbf{a}_h^*)$, and \mathbb{Q}'_h the same for $\hat{\mathbf{a}}_h \sim \hat{\pi}(\hat{\mathbf{s}}_h)$, and
4857 similar incrementing of the dynamics. \square

4874 I.5. Imitation with no augmentation

4875 **Theorem 11.** Let $\hat{\pi}$ be a learner policy, and define

$$4876 \Delta_h^*(\varepsilon) := \mathbb{E}_{\boldsymbol{\rho}_{m,h} \sim \mathcal{D}_{\text{exp}}} \mathbb{E}_{\tilde{\boldsymbol{\rho}}_{m,h} \sim \mathcal{N}(\boldsymbol{\rho}_{m,h}, \sigma^2 \mathbf{I})} \inf_{\mu \in \mathcal{C}(\pi_h^*(\tilde{\boldsymbol{\rho}}_{m,h}), \hat{\pi}_h(\tilde{\boldsymbol{\rho}}_{m,h}))} \mathbb{P}_{(\mathbf{a}, \mathbf{a}') \sim \mu} [\mathbf{d}_{\mathcal{A}}(\mathbf{a}, \mathbf{a}')],$$

4877 which we note defers from $\Delta_h(\varepsilon)$ in Eq. (1.2) in that it measures error with respect to π_h^* , rather than $\pi_{\text{dec}, h}^*$. Suppose that
4878 there is a non-decreasing function $\gamma(\cdot)$ such that for all $\boldsymbol{\rho}_{m,h}, \boldsymbol{\rho}'_{m,h} \in \mathcal{P}_{\tau_m-1}$

$$4879 \text{TV}(\hat{\pi}(\boldsymbol{\rho}_{m,h}), \hat{\pi}(\boldsymbol{\rho}'_{m,h})) \leq \gamma(\|\boldsymbol{\rho}_{m,h} - \boldsymbol{\rho}'_{m,h}\|),$$

4880 where π^* is defined as the conditional in Definition 1.3. Then, the loss of $\hat{\pi}$, without smoothing, is bounded by

$$4881 \mathcal{L}_{\text{marg}, \varepsilon}(\hat{\pi}) \leq H\gamma(\varepsilon\sqrt{2\tau_m-1}) + \sum_{h=1}^H \Delta_h^*(\varepsilon),$$

4882 Further if \mathcal{D}_{exp} has $\tau \leq \tau_m$ bounded memory (Definition 1.6), then it also holds that

$$4883 \mathcal{L}_{\text{joint}, \varepsilon}(\hat{\pi}) \leq H\gamma(\varepsilon\sqrt{2\tau_m-1}) + \sum_{h=1}^H \Delta_h^*(\varepsilon)$$

4894

4895 *Proof.* The above is a direct consequence of the following points. First, with our instantiation of the composite MDP, we
4896 can bound $\mathcal{L}_{\text{marg},\varepsilon}(\hat{\pi}) \leq \Gamma_{\text{marg},\varepsilon}(\hat{\pi} \parallel \pi^*) \leq \Gamma_{\text{joint},\varepsilon}(\hat{\pi} \parallel \pi^*)$ due to [Lemma I.1](#); and moreover, we have $\mathcal{L}_{\text{joint},\varepsilon}(\hat{\pi}) \leq$
4897 $\Gamma_{\text{joint},\varepsilon}(\hat{\pi} \parallel \pi^*)$ when \mathcal{D}_{exp} has $\tau \leq \tau_m$ -bounded memory.

4898 Next, bounding $\|\rho_{m,h} - \rho'_{m,h}\| \leq \sqrt{2\tau_m - 1} d_{\text{TVC}}(\rho_{m,h}, \rho'_{m,h})$, we see $\hat{\pi}$ is $\tilde{\gamma}(\cdot)$ -TVC w.r.t. d_{TVC} , where $\tilde{\gamma}(u) =$
4899 $\gamma(u\sqrt{2\tau_m - 1})$. The bound now follows from [Proposition D.1](#), and the fact that [Proposition 4.1](#) verifies the input-stability
4900 property. \square

4903 I.6. Consequence for expected costs

4904 Finally, we prove [Proposition I.5](#), which shows that it is sufficient to control the imitation losses in [Definition 2.2](#) if we wish
4905 to control the difference of a Lipschitz cost function between the learned policy and the expert distribution:

4906 **Proposition I.5.** *Recall the marginal and final imitation losses in [Definition 2.2](#), and also the joint imitation loss in*
4907 *[Definition I.5](#). Consider a cost function $\mathfrak{J} : \mathcal{P}_T \rightarrow \mathbb{R}$ on trajectories $\rho_T \in \mathcal{P}_T$. Finally, let $\rho_T \sim \mathcal{D}_{\text{exp}}$, and let $\rho'_T \sim \mathcal{D}_\pi$*
4908 *be under the distribution induced by π . Then,*

4910 (a) *If $\max_{\rho_T} |\mathfrak{J}(\rho_T)| \leq B$, and ρ_T is L Lipschitz in the Euclidean norm⁹ (treating ρ_T as Euclidean vector in*
4911 *$\mathbb{R}^{(T+1)d_x + Td_u}$), then*

$$4913 \quad |\mathbb{E}_{\mathcal{D}_{\text{exp}}}[\mathfrak{J}(\rho_T)] - \mathbb{E}_{\mathcal{D}_\pi}[\mathfrak{J}(\rho'_T)]| \leq \sqrt{2T}L\varepsilon + 2B\mathcal{L}_{\text{joint},\varepsilon}(\pi).$$

4915 (b) *If \mathfrak{J} decomposes into a sum of costs, $\mathfrak{J}(\rho) = \ell_{T+1,1}(\mathbf{x}_{1:T}) + \sum_{t=1}^T \ell_{t,1}(\mathbf{x}_t) + \ell_{t,2}(\mathbf{u}_t)$, where $\ell_{t,1}(\cdot), \ell_{t,2}(\cdot)$ are*
4916 *L -Lipschitz and bounded in magnitude in B . Then,*

$$4918 \quad |\mathbb{E}_{\mathcal{D}_{\text{exp}}}[\mathfrak{J}(\rho_T)] - \mathbb{E}_{\mathcal{D}_\pi}[\mathfrak{J}(\rho'_T)]| \leq 4TB\mathcal{L}_{\text{marg},\varepsilon}(\pi) + 2TL\varepsilon.$$

4920 (c) *$\mathfrak{J}(\rho) = \ell_{T+1,1}(\mathbf{x}_{T+1})$ depends only on \mathbf{x}_{T+1} , then*

$$4922 \quad |\mathbb{E}_{\mathcal{D}_{\text{exp}}}[\mathfrak{J}(\rho_T)] - \mathbb{E}_{\mathcal{D}_\pi}[\mathfrak{J}(\rho'_T)]| \leq +2B\mathcal{L}_{\text{fin},\varepsilon}(\pi) + L\varepsilon$$

4923 Thus, for our imitation guarantees to apply to most natural cost functions used in practice, it suffices to control the imitation
4924 losses defined above.

4926 *Proof of [Proposition I.5](#).* Let $\rho_T = (\mathbf{x}_{1:T+1}, \mathbf{u}_{1:T}) \sim \mathcal{D}_{\text{exp}}$, and let $\rho'_T = (\mathbf{x}'_{1:T+1}, \mathbf{u}'_{1:T})$ be under the distribution induced
4927 by π .

4929 **Part (a).** For any coupling μ between the two under which $\mathbf{x}_1 = \mathbf{x}'_1$, and let $\mathcal{E}_\varepsilon := \{\max_t \|\mathbf{x}_{t+1} - \mathbf{x}'_{t+1}\| \vee \|\mathbf{u}_t - \mathbf{u}'_t\| \leq \varepsilon\}$.

$$4931 \quad |\mathbb{E}[\mathfrak{J}(\rho_T)] - \mathbb{E}[\mathfrak{J}(\rho'_T)]| = |\mathbb{E}_\mu[\mathfrak{J}(\rho_T) - \mathfrak{J}(\rho'_T)]|$$

$$4932 \quad \leq \mathbb{E}_\mu[|\mathfrak{J}(\rho_T) - \mathfrak{J}(\rho'_T)|]$$

$$4933 \quad \leq 2B\mathbb{P}_\mu[\mathcal{E}_\varepsilon^c] + \mathbb{E}_\mu[|\mathfrak{J}(\rho_T) - \mathfrak{J}(\rho'_T)|\mathbf{I}\{\mathcal{E}_\varepsilon\}]$$

4935 By passing to an infimum over couplings, $\inf_\mu \mathbb{P}_\mu[\mathcal{E}_\varepsilon^c] \leq \mathcal{L}_{\text{joint},\varepsilon}(\pi)$. Moreover, we observe that under μ , $\mathbf{x}_1 = \mathbf{x}'_1$, and the
4936 remaining coordinates, $(\mathbf{x}_{2:T+1}, \mathbf{u}_{1:T})$ and $(\mathbf{x}'_{2:T+1}, \mathbf{u}'_{1:T})$ are the concatenation of $2T$ vectors, so the Euclidean norm of
4937 the concatenations $\|\rho_T - \rho'_T\|$ is at most $\sqrt{2T} \max_t \|\mathbf{x}_{t+1} - \mathbf{x}'_{t+1}\| \vee \|\mathbf{u}_t - \mathbf{u}'_t\|$, which on \mathcal{E}_ε is at most $\sqrt{2T}\varepsilon$. Using
4938 Lipschitz-ness of \mathfrak{J} concludes.

4940 **Part (b)** . Using the adaptive decomposition of the cost and the fact that \mathbf{x}_1 and \mathbf{x}'_1 have the same distributions,

$$4942 \quad |\mathbb{E}[\mathfrak{J}(\rho_T)] - \mathbb{E}[\mathfrak{J}(\rho'_T)]| = \left| \sum_{t=1}^T (\mathbb{E}[\ell_{t,1}(\mathbf{x}_{t+1})] - \mathbb{E}[\ell_{t,1}(\mathbf{x}'_{t+1})]) + (\mathbb{E}[\ell_{t,2}(\mathbf{u}_t)] - \mathbb{E}[\ell_{t,2}(\mathbf{u}'_t)]) \right|$$

$$4943 \quad \leq \sum_{t=1}^T |\mathbb{E}[\ell_{t,1}(\mathbf{x}_{t+1})] - \mathbb{E}[\ell_{t,1}(\mathbf{x}'_{t+1})]| + |\mathbb{E}[\ell_{t,2}(\mathbf{u}_t)] - \mathbb{E}[\ell_{t,2}(\mathbf{u}'_t)]|$$

4948 ⁹Of course, Lipschitzness in other norms can be derived, albeit with different T dependence
4949

4950 Applying similar arguments as in part (a) to each term, we can bound

$$4951 \max \{ |\mathbb{E}[\ell_{t,1}(\mathbf{x}_{t+1})] - \mathbb{E}[\ell_{t,1}(\mathbf{x}'_{t+1})]|, |\mathbb{E}[\ell_{t,2}(\mathbf{u}_t)] - \mathbb{E}[\ell_{t,2}(\mathbf{u}'_t)]| \} \leq 2B\mathcal{L}_{\text{marg},\varepsilon}(\pi) + L\varepsilon.$$

4952
4953 Summing over the $2T$ terms concludes.

4954
4955 **Part (c).** Follows similar to part (b). □

4956 I.7. Useful Lemmata

4957 I.7.1. ON THE TRAJECTORIES INDUCED BY π^* FROM \mathcal{D}_{exp}

4958 The key step in all of our proofs is to relate the expert distribution over trajectories $\rho_T \sim \mathcal{D}_{\text{exp}}$ to the distribution induced by
4959 the chunking policy π^* in [Definition I.3](#)

4960 **Lemma I.6.** *There exists a sequence of probability kernels π_h^* mapping $\rho_{m,h} \rightarrow \Delta(\mathcal{A})$ such that the chunking policy
4961 $\pi^* = (\pi_h^*)_{1 \leq h \leq H}$ satisfies the following:*

4962 (a) $\pi_h^*(\rho_{m,h})$ is equal to the almost-sure conditional probability of \mathbf{a}_h conditioned on $\rho_{m,h}$ under $\rho_T \sim \mathcal{D}_{\text{exp}}$ and
4963 $\mathbf{a}_{1:H} = \text{synth}(\rho_T)$.

4964 (b) The marginal distribution over each $\rho_{c,h}$ is the same as the marginals of each \mathbf{x}_t and \mathbf{u}_t under $\rho_T \sim \mathcal{D}_{\text{exp}}$.

4965 (c) If \mathcal{D}_{exp} has τ -bounded memory ([Definition I.6](#)) and if $\tau \leq \tau_m$, then the joint distribution of ρ_T induced by π^* is equal
4966 to the joint distribution over ρ_T under \mathcal{D}_{exp} .

4967 **Remark I.4** (Replacing τ -bounded memory with mixing). We can replace that τ -bounded memory condition to the
4968 following mixing assumption. Define the chunk $\rho_{i \leq j} = (\mathbf{x}_{i:j}, \mathbf{u}_{i:j-1})$. Define the measures

$$4969 \begin{aligned} 4970 Q_h(\rho_{m,h}) &= \mathbb{P}_{\mathbf{a}_{1:h-1}, \rho_{1:t_h - \tau_m - 1}, \mathbf{a}_{h:H}, \rho_{t_h:T+1} | \rho_{m,h}} \\ 4971 Q_h^\otimes(\rho_{m,h}) &= \mathbb{P}_{\mathbf{a}_{1:h-1}, \rho_{1:t_h - \tau_m - 1} | \rho_{m,h}} \otimes \mathbb{P}_{\mathbf{a}_{h:H}, \rho_{t_h:T+1} | \rho_{m,h}} \end{aligned}$$

4972 which describes the conditional distribution of the whole trajectory without $\rho_{m,h}$ and the product-distribution of the
4973 conditional distributions of the before- $\rho_{m,h}$ part of the trajectory, and after $\rho_{m,h}$ -part. Under the condition

$$4974 \mathbb{E}_{\rho_{m,h} \text{ from } \rho_T \sim \mathcal{D}_{\text{exp}}} \text{TV}(Q_h(\rho_{m,h}), Q_h^\otimes(\rho_{m,h})) \leq \varepsilon_{\text{mix}}(\tau_m),$$

4975 which measures how close the before- and after- $\rho_{m,h}$ parts of the trajectory are to being conditionally independent, one can
4976 leverage [Lemma I.9](#) to show that

$$4977 \text{TV}(\mathcal{D}_{\pi^*}, \mathcal{D}_{\text{exp}}) \leq H\varepsilon_{\text{mix}}(\tau_m)$$

4978 [Lemma I.6](#) corresponds to the special when when $\varepsilon_{\text{mix}} = 0$.

4979 *Proof of [Lemma I.6](#).* We prove each part in sequence

4980 **Part (a).** follows from the fact that all random variables are in real vector spaces, and thus Polish spaces. Hence, we can
4981 invoke the existence of regular conditional probabilities by [Theorem 3](#).

4982 **Part (b).** This follows by marginalization and Markovianity of the dynamics. Specifically, let $(\rho_T^*, \mathbf{a}_{1:H}^*)$ be a trajectory
4983 and composite actions induced by the chunking policy π^* , and let $(\rho_T, \mathbf{a}_{1:H})$ be the same induced by \mathcal{D}_{exp} . Let $\rho_{m,h}^*$
4984 denote memory chunks of ρ_T^* , and let $\rho_{m,h}$ memory chunks of ρ_T (length $\tau_m - 1$); similarly, denote by $\rho_{c,h}^*$ and $\rho_{c,h}$ the
4985 respective trajectory chunks (length $\tau_c \geq \tau_m$).

4986 We argue inductively that the trajectory chunks $\rho_{c,h}^*$ and $\rho_{c,h}$ are identically distributed for each h . For $h = 1$, $\rho_{c,1}^*$ and
4987 $\rho_{c,1}$ are identically distributed according to $\mathcal{D}_{\mathbf{x}_1}$. Now assume we have show that $\rho_{c,h}^*$ and $\rho_{c,h}$ are identically distributed.
4988 As memory chunks are sub-chunks of trajectory chunks, this means that $\rho_{m,h}^*$ and $\rho_{m,h}$ are identically distributed. By
4989 part (a), it follows that $(\rho_{m,h}^*, \mathbf{a}_h^*)$ and $(\rho_{m,h}, \mathbf{a}_h)$ are identically distributed. In particular, $(\mathbf{x}_{t_h}^*, \mathbf{a}_h^*)$ and $(\mathbf{x}_{t_h}, \mathbf{a}_h)$
4990 are identically distributed, where $\mathbf{x}_{t_h}^*$ (resp \mathbf{x}_{t_h}) these denote the t_h -th control state under π^* (resp. \mathcal{D}_{exp}). By Markovianity of
4991 the dynamics, $\rho_{c,h+1}^*$ and $\rho_{c,h+1}$ are functions of $(\mathbf{x}_{t_h}^*, \mathbf{a}_h^*)$ and $(\mathbf{x}_{t_h}, \mathbf{a}_h)$, respectively, $\rho_{c,h+1}^*$ and $\rho_{c,h+1}$ are identically
4992 distributed, as needed.

5005 **Part (c).** When \mathcal{D}_{exp} has τ -bounded memory and $\tau \leq \tau_m$, then we have the almost-sure equality

$$5006 \mathbb{P}_{\mathcal{D}_{\text{exp}}}[\mathbf{a}_h \in \cdot \mid \mathbf{x}_{1:t_h}, \mathbf{u}_{1:t_h}] = \mathbb{P}_{\mathcal{D}_{\text{exp}}}[\mathbf{a}_h \in \cdot \mid \boldsymbol{\rho}_{m,h}] = \pi_h^*(\boldsymbol{\rho}_{m,h})[\mathbf{a}_h \in \cdot].$$

5007 Finally, $\mathbf{x}_{t_h+1:t_{h+1}}$, $\mathbf{u}_{t_h:t_{h+1}-1}$ are determined by \mathbf{x}_{t_h} and \mathbf{a}_h , this inductively establishes equality of the joint-trajectory

5008 distributions.

5009 □

5013 I.7.2. CONCENTRATION AND TVC OF GAUSSIAN SMOOTHING.

5014 We now include two easy lemmata necessary for the proof. The first shows that p_r is small when r is $\Theta(\sigma)$ by elementary

5015 Gaussian concentration:
5016 **Lemma I.7.** Suppose that $\boldsymbol{\gamma} \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$ is a centred Gaussian vector with covariance $\sigma^2 \mathbf{I}$ in \mathbb{R}^d for some $\sigma > 0$. Then

5017 for all $p > 0$, it holds with probability at least $1 - p$ that

$$5018 \|\boldsymbol{\gamma}\| \leq 2\sigma \cdot \sqrt{2d \log(9) + 2 \log\left(\frac{1}{p}\right)} \leq 2\sigma \cdot \sqrt{5d + 2 \log\left(\frac{1}{p}\right)}$$

5019 Moreover, for $r \geq 4\sigma \sqrt{d \log(9)}$, $\mathbb{P}[\|\boldsymbol{\gamma}\| \geq r] \leq \exp(-r^2/16\sigma^2)$.

5020 *Proof.* We apply the standard covering based argument as in, e.g., [Vershynin \(2018, Section 4.2\)](#). Note that

$$5021 \|\boldsymbol{\gamma}\| = \sup_{\mathbf{w} \in \mathcal{S}^{d-1}} \langle \boldsymbol{\gamma}, \mathbf{w} \rangle,$$

5022 where \mathcal{S}^{d-1} is the unit sphere in \mathbb{R}^d . Let \mathcal{U} denote a minimal $(1/4)$ -net on \mathcal{S}^{d-1} and observe that a simple computation

5023 tells us that

$$5024 \sup_{\mathbf{w} \in \mathcal{S}^{d-1}} \langle \boldsymbol{\gamma}, \mathbf{w} \rangle \leq 2 \cdot \max_{\mathbf{w} \in \mathcal{U}} \langle \mathbf{w}, \boldsymbol{\gamma} \rangle.$$

5025 A classical volume argument (see for example, [Vershynin \(2018, Section 4.2\)](#)) tells us that $|\mathcal{U}| \leq 9^d$. A classical Gaussian

5026 tail bound tells us that for any $\mathbf{w} \in \mathcal{S}^{d-1}$, it holds that for any $r > 0$,

$$5027 \mathbb{P}(\langle \mathbf{w}, \boldsymbol{\gamma} \rangle > r) \leq e^{-\frac{r^2}{2\sigma^2}}.$$

5028 Thus by a union bound, we have

$$5029 \mathbb{P}(\|\boldsymbol{\gamma}\| > r) \leq |\mathcal{U}| \cdot \max_{\mathbf{w} \in \mathcal{U}} \mathbb{P}\left(\|\boldsymbol{\gamma}\| > \frac{r}{2}\right) \leq 9^d \cdot e^{-\frac{r^2}{8\sigma^2}}.$$

5030 Inverting concludes the proof.

5031 □

5032 The second lemma shows that the relevant smoothing kernel is TVC:

5033 **Lemma I.8.** For any $\sigma > 0$, let $\phi_{\mathcal{Z}}$ and W_σ be as in [Definition I.1](#) kernel, then W_σ is γ_{TVC} -TVC for with respect to d_{TVC} (as

5034 defined in [Section 4.1](#))

$$5035 \gamma_{\text{TVC}}(u) = \frac{u\sqrt{2\tau_m - 1}}{2\sigma}.$$

5036 *Proof.* Recall that $\phi_{\mathcal{Z}}$ denotes projection onto the \mathcal{Z} -component of the direct decomposition in [Definition E.1](#), i.e. projects

5037 onto the memory chunk $\boldsymbol{\rho}_{m,h}$. We apply Pinsker's inequality ([Polyanskiy & Wu, 2022+](#)): Then, for for $\mathbf{s}, \mathbf{s}' \in \mathbb{R}^p$, we have

$$5038 \text{TV}(\phi_{\mathcal{Z}} \circ W_\sigma(\mathbf{s}), \phi_{\mathcal{Z}} \circ W_\sigma(\mathbf{s}')) \leq \sqrt{\frac{1}{2} \cdot \text{D}_{\text{KL}}(\phi_{\mathcal{Z}} \circ W_\sigma(\mathbf{s}) \parallel \phi_{\mathcal{Z}} \circ W_\sigma(\mathbf{s}'))}.$$

5039

5060 Note that for $\mathbf{s} = \boldsymbol{\rho}_{c,h}$ with corresponding memory chunk $\boldsymbol{\rho}_{m,h}$ $\phi_{\mathcal{Z}} \circ W_{\sigma}(\mathbf{s}) \sim \mathcal{N}(\boldsymbol{\rho}_{m,h}, \sigma^2 \mathbf{I})$. Similarly, for $\boldsymbol{\rho}'_{m,h}$
 5061 corresponding to \mathbf{s}' , $\phi_{\mathcal{Z}} \circ W_{\sigma}(\mathbf{s}') \sim \mathcal{N}(\boldsymbol{\rho}'_{m,h}, \sigma^2 \mathbf{I})$. Hence,
 5062

$$5063 \quad \text{D}_{\text{KL}}(\phi_{\mathcal{Z}} \circ W_{\sigma}(\mathbf{s}) \parallel \phi_{\mathcal{Z}} \circ W_{\sigma}(\mathbf{s}')) \leq \frac{\|\boldsymbol{\rho}_{m,h} - \boldsymbol{\rho}'_{m,h}\|^2}{2\sigma^2}.$$

5064 Thus, we conclude $\text{TV}(\phi_{\mathcal{Z}} \circ W_{\sigma}(\mathbf{s}), \phi_{\mathcal{Z}} \circ W_{\sigma}(\mathbf{s}')) \leq \frac{\|\boldsymbol{\rho}_{m,h} - \boldsymbol{\rho}'_{m,h}\|}{2\sigma}$. Finally, we upper bound the Euclidean norm
 5065 $\|\boldsymbol{\rho}_{m,h} - \boldsymbol{\rho}'_{m,h}\|$ of vectors consisting of $2\tau_m - 1$ sub-vectors via d_{TVc} (which is the maximum Euclidean norm of
 5066 these subvectors) via $\|\boldsymbol{\rho}_{m,h} - \boldsymbol{\rho}'_{m,h}\| \leq \sqrt{2\tau_m - 1} d_{\text{TVc}}(\mathbf{s}, \mathbf{s}')$. \square
 5067

5072 I.7.3. TOTAL VARIATION TELESCOPING

5073 **Lemma I.9** (Total Variation Telescoping). *Let $\mathcal{Y}_1, \dots, \mathcal{Y}_H, \mathcal{Y}_{H+1}$ be Polish spaces. Let $P_1 \in \Delta(\mathcal{Y}_1)$, and let $Q_h, Q'_h \in$
 5074 $\Delta(\mathcal{Y}_h \mid \mathcal{X}, \mathcal{Y}_{1:h-1})$, $h > 1$. Define $P'_1 = P_1$, and recursively define*

$$5075 \quad P_h = \text{law}(Q_h; P_{h-1}), \quad P'_h = \text{law}(Q'_h; P'_{h-1}), \quad h > 1.$$

5076 Then,

$$5077 \quad \text{TV}(P_{H+1}, P'_{H+1}) \leq \sum_{h=1}^H \mathbb{E}_{Y_{1:h} \sim P_h} \text{TV}(Q_{h+1}(\cdot \mid Y_{1:h}), Q'_{h+1}(\cdot \mid Y_{1:h}))$$

5078 Moreover, there exists a coupling of $\mu \in \mathcal{C}(P_{H+1}, P'_{H+1})$ over $Y_{1:H+1} \sim P_{H+1}$ and $Y_{1:H+1} \sim P'_{H+1}$ such that

$$5079 \quad \mathbb{P}_{\mu}[Y_1 = Y'_1] = 1, \quad \mathbb{P}_{\mu}[Y_{1:H+1} \neq Y'_{1:H+1}] \leq \sum_{h=1}^H \mathbb{E}_{Y_{1:h} \sim P_h} \text{TV}(Q_{h+1}(\cdot \mid Y_{1:h}), Q'_{h+1}(\cdot \mid Y_{1:h})).$$

5080 *Proof.* To prove the first part of the lemma, define $Q'_{i,j}$ for $2 \leq i \leq j \leq H+1$ by $Q'_{i,i} = Q_i$ define $Q'_{i,j}$ by appending $Q'_{i,j}$
 5081 to $Q'_{i,j-1}$ and $\text{law}(Q'_{i,j}; (\cdot)) = \text{law}(Q'_j; \text{law}(Q_{i,j-1}; (\cdot)))$. We now define

$$5082 \quad P^{(i)} = \text{law}(Q'_{i+1,H+1}; P_i),$$

5083 with the convention $\text{law}(Q'_{H+2,H+1}; P_{H+1}) = P_{H+1}$. Note that $P^{(H+1)} = P_{H+1}$, and $P^{(1)} = P'_{H+1}$. Then, because TV
 5084 distance is a metric,

$$5085 \quad \text{TV}(P_{H+1}, P'_{H+1}) \leq \sum_{h=1}^H \text{TV}(P^{(h)}, P^{(h+1)})$$

5086 Moreover, we can write $P^{(i)} = \text{law}(Q'_{i+2,H+1}; \text{law}(Q'_{i+1}; P_i))$ and $P_{i+1} = \text{law}(Q_{i+1}; P_i)$. Thus,

$$5087 \quad \text{TV}(P^{(i)}, P^{(i+1)}) = \text{TV}(\text{law}(Q'_{i+2,H+1}; \text{law}(Q'_{i+1}; P_i)), \text{law}(Q'_{i+2,H+1}; \text{law}(Q_{i+1}; P_i)) \quad (\text{Lemma C.4})$$

$$5088 \quad = \text{TV}(\text{law}(Q'_{i+1}; P_i), \text{law}(Q_{i+1}; P_i))$$

$$5089 \quad = \mathbb{E}_{Y_{1:i} \sim P_i} \text{TV}(Q'_i(Y_{1:i}), Q_i(Y_{1:i})). \quad (\text{Corollary C.1})$$

5090 This completes the first part of the demonstration (noting symmetry of TV). The second part follows from [Corollary C.1](#), by
 5091 letting $Y \leftarrow Y_1$, and $X \leftarrow Y_{2:H+1}$ in that lemma. \square
 5092

5108 J. Extensions and Further Results

5109 J.1. Noisy Dynamics

5110 We can directly extend our imitation guarantees in the composite MDP to settings with noise:

$$5111 \quad \mathbf{s}_{h+1} \sim F_h^{\text{noise}}(\mathbf{s}_h, \mathbf{a}_h, \mathbf{w}_h), \quad \mathbf{w}_h \sim P_{\text{noise},h}, \quad (\text{J.1})$$

5112

5115 where the noises are independent of states and of each other. Indeed, (J.1) can be directly reduced to the no-noise setting by
 5116 lifting “actions” to pairs (a_h, w_h) , and policies π to encompass their distribution of actions, and over noise.

5117 Another approach is instead to condition on the noises $w_{1:H}$ first, and treat the noise-conditioned dynamics as deterministic.
 5118 Then one can take expectation over the noises and conclude. The advantage of this approach is that the couplings constructed
 5119 thereby is that the trajectories experience identical sequences of noise with probability one.
 5120

5121 Extending the control setting to incorporate noise is doable but requires more effort:

- 5122 • If the *demonstrations are noiseless*, then one can still appeal to the synthesis oracle to synthesis stabilizing gains. How-
 5123 ever, one needs to (ever so slightly) generalize the proofs of the various stability properties (e.g. IPS in [Proposition 4.1](#))
 5124 to accomodate system noise.
- 5125 • If the demonstrations themselves have noise, one may need to modify the synthesis oracle setup somewhat. This
 5126 is because the synthesis oracle, if it synthesizes stabilizing gains, will attempt to get the learner to stabilize to a
 5127 noise-perturbed trajectory. This can perhaps be modified by synthesizing controllers which stabilize to smoothed
 5128 trajectories, or by collecting demonstrations of desired trajectories (e.g. position control), and stabilizing to the these
 5129 states than than to actual states visited in demonstrations.
 5130
 5131
 5132

5133 J.2. Robustness to Adversarial Perturbations

5134 Our results can accomodate an even more general framework where there are both noises as well adversarial perturbations.
 5135 We explain this generalization in the composite MDP.

5136 Specificial, consider a space \mathcal{E} of adversarial perturbations, as well as \mathcal{W} of noises as above. We may posit a dynamics
 5137 function $F^{\text{adv}} : \mathcal{S} \times \mathcal{A} \times \mathcal{W} \times \mathcal{A} \rightarrow \mathcal{S}$, and consider the evolution of an imitator policy $\hat{\pi}$ under the adversary
 5138

$$\begin{aligned}
 5139 \hat{s}_{h+1} &= F_h^{\text{adv}}(\hat{s}_h, \hat{a}_h, w_h, e_h), \quad w_h \sim P_{\text{noise},h} \\
 5140 \hat{a}_h &\sim \hat{\pi}_h(s_h) \\
 5141 e_h &\sim \pi_h^{\text{adv}}(\hat{s}_{1:h}, a_{1:h}, w_{1:h}, e_{1:h-1}), \\
 5142 \hat{s}_1 &\sim \pi_0^{\text{adv}}(s_1), \quad s_1 \sim P_{\text{init}}.
 \end{aligned}$$

5143 By constrast, we can model the demonstrator trajectory as arising from noisy, but otherwise unperturbed trajectories:

$$5144 s_{h+1}^* \sim F_h^{\text{adv}}(s_h^*, a_h^*, w_h, 0), \quad w_h \sim P_{\text{noise},h}, \quad a_h^* \sim \pi_h^*(s_h^*), \quad s_1^* \sim P_{\text{init}}.$$

5145 To reduce the composite-MDP in [Section 4](#), we can view the combination of adveryary π^{adv} and imitator $\hat{\pi}$ as a combined
 5146 policy, and the π^* with zero augmentation as another policy; here, we would them treat actions as $\tilde{a} = (a, e)$. Then, one
 5147 can consider modified senses of stability which preserve trajectory tracking, as well as a modification of $d_{\mathcal{A}}$ to a function
 5148 measuring distances between $\tilde{a} = (a, e)$ and $\tilde{a}' = (a', e')$. The extension is rather mechanical, and we fit details. Note
 5149 further that, by including a $\pi_0^{\text{adv}}(s_1)$, we can modify the analysis to allow for subtle differences in initial state distribution.
 5150 This would in turn require strengthening our stability assumptions to allow stability to initial state (e.g., the definition of
 5151 incremental stability as exposted by ([Pfrommer et al., 2022](#))).
 5152

5153 J.3. Deconvolution Policies and Total Variation Continuity

5154 While our strongest guarantees hold for the replica policies, where we add noise both as a data augmentation at training
 5155 time *and* at test time, many practitioners have seen some success with the deconvolution policies where noise is only added
 5156 at training time. We note that [Proposition D.1](#) holds when the learned policy is TVC; without noise at training time this
 5157 certainly will not hold when the expert policy is not TVC. We show here that the deconvolution expert policy is TVC under
 5158 mild assumptions, which lends some credence to the empirical success of deconvolution policies.
 5159

5160 Precisely, we show that, under reasonable conditions, deconvolution is total variation continuous. In particular, suppose
 5161 that $\mu \in \Delta(\mathbb{R}^d)$ is a Borel probabily measure and p is a density with respect to μ . Further suppose that Q is a density
 5162 with respect to the Lebesgue measure on \mathbb{R}^d . Suppose that $\mathbf{x} \sim p$, $\mathbf{w} \sim Q$, and let $\tilde{\mathbf{x}} = \mathbf{x} + \mathbf{w}$. We will show that the
 5163 deconvolution measure $p(\mathbf{x}|\tilde{\mathbf{x}})$ is continuous in TV.
 5164
 5165
 5166
 5167
 5168
 5169

5170 **Proposition J.1.** Let $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$ be fixed, let $p : \mathbb{R}^d \rightarrow \mathbb{R}$ denote a probability density, and let $Q : \mathbb{R}^d \rightarrow \mathbb{R}$ denote a
5171 function such that $\nabla^2 Q$ and $\nabla \log Q$ exist and are continuous on the set

$$5172 \mathcal{X} = \{(1-t)\tilde{\mathbf{x}} + t\tilde{\mathbf{x}}' - \mathbf{x} \mid \mathbf{x} \in \text{supp } p \text{ and } t \in [0, 1]\}$$

5175 Then it holds that

$$5176 \text{TV}(p(\cdot|\tilde{\mathbf{x}}), p(\cdot|\tilde{\mathbf{x}}')) \leq \|\tilde{\mathbf{x}} - \tilde{\mathbf{x}}'\| \cdot \sup_{\mathbf{x} \in \mathcal{X}} \|\nabla \log Q(\mathbf{x})\|.$$

5180 By [Lemma C.4](#), any policy composed with the total variation kernel is thus total variation continuous with a linear γ_{TVC} ;
5181 moreover, the Lipschitz constant is given by the maximal norm of the score of the noise distribution. For example, if Q is
5182 the density of a Gaussian with variance σ^2 , then $\gamma_{\text{TVC}}(u) \leq \frac{\sup_{\mathbf{x}} \|\mathbf{x}\|}{\sigma^2}$ is dimension independent.

5183 **Remark J.1.** Note that our notation is intentionally different from that in the body to emphasize that this is a general fact
5184 about abstract probability measures. We may instantiate the guarantee in the control setting of interest by letting $\mathbf{x} = \boldsymbol{\rho}_{\mathbf{m},h}$
5185 and consider Q to be a Gaussian (for example) kernel. In this case, we see that the deconvolution policy of [Definition 3.1](#) is
5186 automatically TVC.

5188 To prove [Proposition J.1](#), we begin with the following lemma:

5189 **Lemma J.2.** Let $\tilde{\mathbf{x}} \in \mathbb{R}^d$ be fixed and suppose that $\nabla \log Q(\tilde{\mathbf{x}} - \mathbf{x})$ exists for all $\mathbf{x} \in \text{supp } p$. Then, for all $\mathbf{x} \in \text{supp } p$, it
5190 holds that $\nabla_{\tilde{\mathbf{x}}} p(\mathbf{x}|\tilde{\mathbf{x}})$ exists. Furthermore,

$$5192 \int \|\nabla p(\mathbf{x}|\tilde{\mathbf{x}})\| d\mu(\mathbf{x}) \leq 2 \sup_{\mathbf{x} \in \text{supp } p} \|\nabla \log Q(\tilde{\mathbf{x}} - \mathbf{x})\|,$$

5195 where the gradient above is with respect to $\tilde{\mathbf{x}}$.

5198 *Proof.* We begin by noting that if $\nabla \log Q(\tilde{\mathbf{x}} - \mathbf{x})$ exists, then so does $\nabla Q(\tilde{\mathbf{x}} - \mathbf{x})$. By Bayes' rule,

$$5200 p(\mathbf{x}|\tilde{\mathbf{x}}) = \frac{p(\mathbf{x})Q(\tilde{\mathbf{x}} - \mathbf{x})}{\int Q(\tilde{\mathbf{x}} - \mathbf{x}')p(\mathbf{x}')d\mu(\mathbf{x}')}$$

5203 We can then compute directly that

$$5205 \nabla p(\mathbf{x}|\tilde{\mathbf{x}}) = \frac{p(\mathbf{x})\nabla Q(\tilde{\mathbf{x}} - \mathbf{x})}{\int Q(\tilde{\mathbf{x}} - \mathbf{x}')p(\mathbf{x}')d\mu(\mathbf{x}')} - \frac{p(\mathbf{x})Q(\tilde{\mathbf{x}} - \mathbf{x}) \cdot \int \nabla Q(\tilde{\mathbf{x}} - \mathbf{x}')p(\mathbf{x}')d\mu(\mathbf{x}')}{\left(\int Q(\tilde{\mathbf{x}} - \mathbf{x}')p(\mathbf{x}')d\mu(\mathbf{x}')\right)^2},$$

5208 where the exchange of the gradient and the integral is justified by Lebesgue dominated convergence and the assumption of
5209 differentiability of Q and thus existence is ensured. We have now that

$$\begin{aligned} 5211 \|\nabla p(\mathbf{x}|\tilde{\mathbf{x}})\| &= \frac{p(\mathbf{x})Q(\tilde{\mathbf{x}} - \mathbf{x})}{\int Q(\tilde{\mathbf{x}} - \mathbf{x}')p(\mathbf{x}')d\mu(\mathbf{x}')} \cdot \left\| \nabla \log Q(\tilde{\mathbf{x}} - \mathbf{x}) - \frac{\int \nabla Q(\tilde{\mathbf{x}} - \mathbf{x}')p(\mathbf{x}')d\mu(\mathbf{x}')}{\int Q(\tilde{\mathbf{x}} - \mathbf{x}')p(\mathbf{x}')d\mu(\mathbf{x}')} \right\| \\ 5212 &= \frac{p(\mathbf{x})Q(\tilde{\mathbf{x}} - \mathbf{x})}{\int Q(\tilde{\mathbf{x}} - \mathbf{x}')p(\mathbf{x}')d\mu(\mathbf{x}')} \cdot \left\| \nabla \log Q(\tilde{\mathbf{x}} - \mathbf{x}) - \frac{\int (\nabla \log Q(\tilde{\mathbf{x}} - \mathbf{x}')) \cdot Q(\tilde{\mathbf{x}} - \mathbf{x})p(\mathbf{x}')d\mu(\mathbf{x}')}{\int Q(\tilde{\mathbf{x}} - \mathbf{x}')p(\mathbf{x}')d\mu(\mathbf{x}')} \right\| \\ 5213 &\leq \left(\sup_{\mathbf{x} \in \text{supp } p} \|\nabla \log Q(\tilde{\mathbf{x}} - \mathbf{x})\| \right) \cdot \frac{p(\mathbf{x})Q(\tilde{\mathbf{x}} - \mathbf{x})}{\int Q(\tilde{\mathbf{x}} - \mathbf{x}')p(\mathbf{x}')d\mu(\mathbf{x}')} \cdot \left(1 + \frac{\int Q(\tilde{\mathbf{x}} - \mathbf{x})p(\mathbf{x}')d\mu(\mathbf{x}')}{\int Q(\tilde{\mathbf{x}} - \mathbf{x}')p(\mathbf{x}')d\mu(\mathbf{x}')} \right) \\ 5214 &= \left(2 \sup_{\mathbf{x} \in \text{supp } p} \|\nabla \log Q(\tilde{\mathbf{x}} - \mathbf{x})\| \right) \cdot \frac{p(\mathbf{x})Q(\tilde{\mathbf{x}} - \mathbf{x})}{\int Q(\tilde{\mathbf{x}} - \mathbf{x}')p(\mathbf{x}')d\mu(\mathbf{x}')} \end{aligned}$$

5221 Now, integrating over \mathbf{x} makes the second factor 1, concluding the proof. □

5223 We will now make use of the theory of Dini derivatives ([Hagood & Thomson, 2006](#)) to prove a bound on total variation.

5225 **Lemma J.3.** For fixed $\tilde{\mathbf{x}}, \tilde{\mathbf{x}}'$ and $0 \leq t \leq 1$, let the upper Dini derivative

$$5226 \quad D^+ \text{TV}(p(\cdot|\tilde{\mathbf{x}}), p(\cdot|\tilde{\mathbf{x}}_t)) = \limsup_{h \downarrow 0} \frac{\text{TV}(p(\cdot|\tilde{\mathbf{x}}), p(\cdot|\tilde{\mathbf{x}}_{t+h})) - \text{TV}(p(\cdot|\tilde{\mathbf{x}}), p(\cdot|\tilde{\mathbf{x}}_t))}{h},$$

5228 where

$$5229 \quad \tilde{\mathbf{x}}_t = (1-t)\tilde{\mathbf{x}} + t\tilde{\mathbf{x}}'.$$

5230 If $\nabla \log Q(\tilde{\mathbf{x}}_t - \mathbf{x})$ exists and is finite for all $\mathbf{x} \in \text{supp } p$ and $t \in [0, 1]$, then

$$5231 \quad \text{TV}(p(\cdot|\tilde{\mathbf{x}}), p(\cdot|\tilde{\mathbf{x}}')) \leq \int_0^1 D^+ \text{TV}(p(\cdot|\tilde{\mathbf{x}}), p(\cdot|\tilde{\mathbf{x}}_t)) dt. \quad (\text{J.2})$$

5232 *Proof.* We compute:

$$5233 \quad \begin{aligned} 5234 \quad 2 |\text{TV}(p(\cdot|\tilde{\mathbf{x}}), p(\cdot|\tilde{\mathbf{x}}_{t+h})) - \text{TV}(p(\cdot|\tilde{\mathbf{x}}), p(\cdot|\tilde{\mathbf{x}}_t))| &= \left| \int |p(\mathbf{x}|\tilde{\mathbf{x}}) - p(\mathbf{x}|\tilde{\mathbf{x}}_{t+h})| - |p(\mathbf{x}|\tilde{\mathbf{x}}) - p(\mathbf{x}|\tilde{\mathbf{x}}_t)| d\mu(\mathbf{x}) \right| \\ 5235 \quad &\leq \int |p(\mathbf{x}|\tilde{\mathbf{x}}_{t+h}) - p(\mathbf{x}|\tilde{\mathbf{x}}_t)| d\mu(\mathbf{x}). \end{aligned} \quad (\text{J.3})$$

5236 Observe that by the assumption on Q and [Lemma J.2](#), $p(\mathbf{x}|\tilde{\mathbf{x}}_t)$ is differentiable and thus continuous in $\tilde{\mathbf{x}}_t$. We therefor see that the function

$$5237 \quad t \mapsto \text{TV}(p(\cdot|\tilde{\mathbf{x}}), p(\cdot|\tilde{\mathbf{x}}_t))$$

5238 is continuous as $\tilde{\mathbf{x}}_t$ is linear in t . By [Hagood & Thomson \(2006, Theorem 10\)](#), [\(J.2\)](#) holds. □

5239 We now bound the Dini derivatives:

5240 **Lemma J.4.** Let $\tilde{\mathbf{x}}, \tilde{\mathbf{x}}' \in \mathbb{R}^d$ such that for all $t \in [0, 1]$ it holds that

$$5241 \quad \sup_{\mathbf{x} \in \text{supp } p} \left| \frac{d^2}{dt^2} (p(\mathbf{x}|\tilde{\mathbf{x}}_t)) \right| = C < \infty,$$

5242 where the derivative is applied on $\tilde{\mathbf{x}}_t$. If the assumptions of [Lemmas J.2](#) and [J.4](#) hold, then

$$5243 \quad D^+ \text{TV}(p(\cdot|\tilde{\mathbf{x}}), p(\cdot|\tilde{\mathbf{x}}_t)) \leq \|\tilde{\mathbf{x}} - \tilde{\mathbf{x}}'\| \cdot \sup_{\substack{\mathbf{x} \in \text{supp } p \\ t \in [0, 1]}} \|\nabla \log Q(\tilde{\mathbf{x}}_t - \mathbf{x})\|.$$

5244 *Proof.* By definition,

$$5245 \quad D^+ \text{TV}(p(\cdot|\tilde{\mathbf{x}}), p(\cdot|\tilde{\mathbf{x}}_t)) = \limsup_{h \downarrow 0} \frac{\text{TV}(p(\cdot|\tilde{\mathbf{x}}), p(\cdot|\tilde{\mathbf{x}}_{t+h})) - \text{TV}(p(\cdot|\tilde{\mathbf{x}}), p(\cdot|\tilde{\mathbf{x}}_t))}{h}.$$

5246 Fix some t and some small h . By [\(J.3\)](#), it holds that

$$5247 \quad |\text{TV}(p(\cdot|\tilde{\mathbf{x}}), p(\cdot|\tilde{\mathbf{x}}_{t+h})) - \text{TV}(p(\cdot|\tilde{\mathbf{x}}), p(\cdot|\tilde{\mathbf{x}}_t))| \leq \frac{1}{2} \cdot \int |p(\mathbf{x}|\tilde{\mathbf{x}}_{t+h}) - p(\mathbf{x}|\tilde{\mathbf{x}}_t)| d\mu(\mathbf{x}).$$

5248 By Taylor's theorem, it holds that

$$5249 \quad p(\mathbf{x}|\tilde{\mathbf{x}}_{t+h}) - p(\mathbf{x}|\tilde{\mathbf{x}}_t) = h \cdot \frac{d}{dt} (p(\mathbf{x}|\tilde{\mathbf{x}}_t)) + h^2 \cdot \frac{d^2}{dt^2} (p(\mathbf{x}|\tilde{\mathbf{x}}_{t'}))$$

5250 for some $t' \in [0, 1]$. By the chain rule, we have

$$5251 \quad \frac{d}{dt} (p(\mathbf{x}|\tilde{\mathbf{x}}_t)) = \langle \tilde{\mathbf{x}}' - \tilde{\mathbf{x}}, \nabla p(\mathbf{x}|\tilde{\mathbf{x}}_t) \rangle,$$

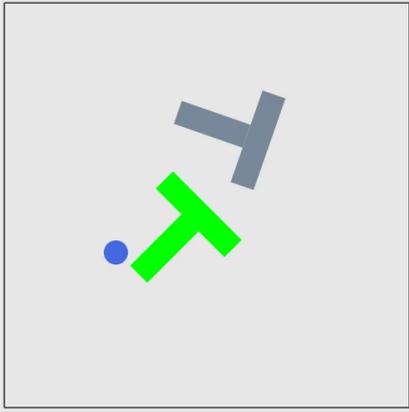
5252 and thus,

$$5253 \quad |p(\mathbf{x}|\tilde{\mathbf{x}}_{t+h}) - p(\mathbf{x}|\tilde{\mathbf{x}}_t)| \leq h \cdot \|\tilde{\mathbf{x}} - \tilde{\mathbf{x}}'\| \cdot \|\nabla p(\mathbf{x}|\tilde{\mathbf{x}}_t)\| + h^2 C$$

5254 Now, applying [Lemma J.2](#) and plugging into the previous computation concludes the proof. □

5255

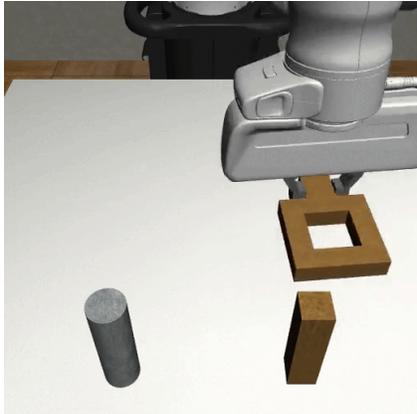
5280
5281
5282
5283
5284
5285
5286
5287
5288
5289
5290
5291
5292
5293
5294
5295
5296
5297
5298
5299
5300
5301
5302
5303
5304
5305
5306
5307
5308
5309
5310
5311
5312
5313
5314
5315
5316
5317
5318
5319
5320
5321
5322
5323
5324
5325
5326
5327
5328
5329
5330
5331
5332
5333
5334



(a) PushT Environment (Chi et al., 2023). The blue circle is the manipulation agent, while the green area is the target position which the agent must push the blue T block into.



(b) Can Pick-and-Place Environment (Mandlekar et al., 2021). The grasper must pick up a can from the left bin and place it into the correct bin on the right side.



(c) Square Nut Assembly Environment (Mandlekar et al., 2021). The grasper must pick up the square nut (the position of which is randomized) and place it over the square peg.

Figure 6. Environment Visualizations.

We are finally ready to state and prove our main result:

Proof of Proposition J.1. Note that

$$\frac{d^2}{dt^2} (p(\mathbf{x}|\tilde{\mathbf{x}}_t)) = (\tilde{\mathbf{x}} - \tilde{\mathbf{x}}')^T \nabla^2 p(\mathbf{x}|\tilde{\mathbf{x}}_t) (\tilde{\mathbf{x}} - \tilde{\mathbf{x}}')$$

and thus is bounded if and only if $\nabla^2 p(\mathbf{x}|\tilde{\mathbf{x}}_t)$ is bounded. An elementary computation shows that if $\nabla^2 Q$ exists and is continuous on \mathcal{X} , then $\nabla^2 p(\mathbf{x}|\tilde{\mathbf{x}}_t)$ is bounded in operator norm on \mathcal{X} . Thus the assumption in Lemma J.4 holds. Applying Lemma J.3 then concludes the proof. \square

K. Experiment Details

K.1. Compute and Codebase Details

Code. For our experiments we build on the existing PyTorch-based codebase and standard environment set provided by Chi et al. (2023) as well as the robomimic demonstration dataset Mandlekar et al. (2021).¹⁰

Compute. We ran all experiments using 4 Nvidia V100 GPUs on an internal cluster node. For each environment running all experiments depicted in Figure 2 took 12 hours to complete with 20 workers running simultaneously for a total of approximately 10 days worth of compute-hours. Between all 20 workers, peak system RAM consumption totaled about 500 GB.

K.2. Environment Details

For simplicity the stabilizatin oracle `synth` is built into the environment so that the diffusion policy effectively only performs positional control. See Appendix K for visualizations of the environments.

PushT. The PushT environment introduced in (Chi et al., 2023) is a 2D manipulation problem simulated using the PyMunk physics engine. It consists of pushing a T-shaped block from a randomized start position into a target position using a controllable circular agent. The synthesis oracle runs a low-level feedback controller at a 10 times higher to stabilize the

¹⁰The modified codebase with instructions for running the experiments is available at the following anonymous link: https://www.dropbox.com/s/vzw0gvk1fd3yadw/diffusion_policy.zip?dl=0. We will provide a public github repository for the final release.

5335 agent’s position towards a desired target position at each point in time via acceleration control. Similar to [Chi et al. \(2023\)](#),
5336 we use a position-error gain of $k_p = 100$ and velocity-error gain of $k_v = 20$. The observation provided to the DDPM model
5337 consists of the x,y coordinates of 9 keypoints on the T block in addition to the x,y coordinates of the manipulation agent, for
5338 a total observation dimensionality of 20.

5339 For rollouts on this environment we used trajectories of length $T = 300$. Policies were scored based on the maximum
5340 coverage between the goal area and the current block position, with > 95 percent coverage considered an “successful”
5341 (score = 1) demonstration and the score linearly interpolating between 0 and 1 for less coverage. A total of 206 human
5342 demonstrations were collected, out of which we use a subset of 90 for training.
5343

5344 **Can Pick-and-Place.** This environment is based on the Robomimic ([Mandlekar et al., 2021](#)) project, which in turn uses the
5345 MuJoCo physics simulator. For the low-level control synthesis we use the feedback controller provided by the Robomimic
5346 package. The position-control action space is 7 dimensional, including the desired end manipulator position, rotation, and
5347 gripper position, while the observation space includes the object pose, rotation in addition to position and rotation of all
5348 linkages for a total of 23 dimensions. Demonstrations are given a score of 1 if they successfully complete the pick-and-place
5349 task and a score of 0 otherwise. We roll out 400 timesteps during evaluation and for training use a subset of up to 90 of the
5350 200 “proficient human” demonstrations provided.
5351

5352 **Square Nut Assembly.** For Square Nut Assembly, which is also Robomimic-based ([Mandlekar et al., 2021](#)), we use
5353 the same setup as the Can Pick and Place task in terms of training data, demonstration scoring, and low-level positional
5354 controller. The observation, action spaces are also equivalent to the Can Pick-and-Place task with 23 and 7 dimensions
5355 respectively.
5356

5357 K.3. DDPM Model and Training Details.

5358
5359 For our DDPM we use the same 1-D convolutional UNet-style ([Ronneberger et al., 2015](#)) architecture employed by ([Chi
5360 et al., 2023](#)), which is in turn adapted from [Janner et al. \(2022\)](#). This principally consists of 3 sets of downsampling
5361 1-dimensional convolution operations using Mish activation functions ([Misra, 2019](#)), Group Normalization (with 8 groups)
5362 ([Wu & He, 2018](#)), and skip connections with 64, 128, and 256 channels followed by transposed convolutions and activations
5363 in the reversed order. The observation and timestep were provided to the model with Feature-wise Linear Modulation
5364 (FiLM) ([Perez et al., 2018](#)), with the timestep encoded using sin-positional encoding into a 64 dimensional vector.
5365

5366 During training and evaluation we utilize a squared cosine noise schedule ([Nichol & Dhariwal, 2021](#)) with 100 timesteps.
5367 For training we use the AdamW optimizer with linear warmup of 500 steps, followed by an initial learning rate of 1×10^{-4}
5368 combined with cosine learning rate decay over the rest of the training horizon. For PushT models we train for 800 epochs and
5369 evaluate test trajectories every 200 epochs while for Can Pick-and-Place and Square Nut Assembly we evaluate performance
5370 every 250 epochs and train for a total of 1500 epochs.

5371 In both environments the diffusion models are conditioned on the previous two observations trained to predict a sequence of
5372 16 target manipulator positions, starting at the first timestep in the conditional observation sequence. The 2rd (corresponding
5373 to the target position for the current timestep) through 9th generated actions are emitted as the $\tau_c = 8$ length action sequence
5374 and the rest is discarded. Extracting a subsequence of a longer prediction horizon in this manner has been shown to improve
5375 performance over just predicting the $H = 8$ action sequence directly ([Chi et al., 2023](#)).
5376

5377 For $\sigma > 0$ we generate new perturbed observations per training iteration, effectively using $N_{\text{aug}} = N_{\text{epoch}}$ augmentations.
5378 We find this to be easier than generating and storing N_{aug} augmentations with little impact on the training and validation
5379 error. Noise is injected after the observations have been normalized such that all components lie within $[-1, 1]$ range.
5380 Performing noise injection post normalization ensures that the magnitude of noise injected is not affected by different units
5381 or magnitudes. .
5382

5383
5384
5385
5386
5387
5388
5389