

---

# The Power of Learned Locally Linear Models for Nonlinear Policy Optimization

---

Daniel Pfrommer<sup>\*1</sup> Max Simchowitz<sup>\*1</sup> Tyler Westenbroek<sup>2</sup> Nikolai Matni<sup>3</sup> Stephen Tu<sup>4</sup>

## Abstract

A common pipeline in learning-based control is to iteratively estimate a model of system dynamics, and apply a trajectory optimization algorithm - e.g. iLQR - on the learned model to minimize a target cost. This paper conducts a rigorous analysis of a simplified variant of this strategy for general nonlinear systems. We analyze an algorithm which iterates between estimating local linear models of nonlinear system dynamics and performing iLQR-like policy updates. We demonstrate that this algorithm attains sample complexity polynomial in relevant problem parameters, and, by synthesizing locally stabilizing gains, overcomes exponential dependence in problem horizon. Experimental results validate the performance of our algorithm, and compare to natural deep-learning baselines.

## 1. Introduction

Machine learning methods such as model-based reinforcement learning have lead to a number of breakthroughs in key applications across robotics and control (Kocijan et al., 2004; Tassa et al., 2012; Nagabandi et al., 2019). A popular technique in these domains is learning-based model-predictive control (MPC) (Morari & Lee, 1999; Williams et al., 2017), wherein a model learned from data is used to repeatedly solve online planning problems to control the real system. It has long been understood that solving MPC *exactly*—both with perfectly accurate dynamics and minimization to globally optimality for each planning problem—enjoys numerous beneficial control-theoretic properties (Jadbabaie & Hauser, 2001).

Unfortunately, the above situation is not reflective of prac-

---

<sup>\*</sup>Equal contribution <sup>1</sup>Massachusetts Institute of Technology <sup>2</sup>University of Texas, Austin <sup>3</sup>University of Pennsylvania <sup>4</sup>Google Brain. Correspondence to: Daniel Pfrommer <dpfrom@mit.edu>.

tice. For one, most systems of practical interest are *nonlinear*, and therefore exact global recovery of system dynamics suffers from a curse of dimensionality. And second, the nonlinear dynamics render any natural trajectory planning problem nonconvex, making global optimality elusive. In this work, we focus on learning-based trajectory optimization, the “inner-loop” in MPC. We ask *when can we obtain rigorous guarantees about the solutions to nonlinear trajectory optimization under unknown dynamics?*

We take as our point of departure the iLQR algorithm (Li & Todorov, 2004). Initially proposed under known dynamics, iLQR solves a planning objective by solving an iterative linear control problem around a first-order Taylor expansion (the *Jacobian linearization*) of the dynamics, and second-order Taylor expansion of the control costs. In solving this objective, iLQR synthesizes a sequence of locally-stabilizing feedback gains, and each iLQR-update can be interpreted as a gradient-step through the closed-loop linearized dynamics in feedback with these gains. This has the dual benefit of proposing a locally stabilizing policy (not just an open-loop trajectory), and of stabilizing the gradients to circumvent exponential blow-up in planning horizon. iLQR, and its variants (Todorov & Li, 2005; Williams et al., 2017), are now ubiquitous in robotics and control applications; and, when dynamics are unknown or uncertain, one can simply substitute the exact dynamics model with an estimate (e.g. Levine & Koltun (2013)). In this case, dynamics are typically estimated with neural networks. Thus, Jacobian linearizations can be computed by automated differentiation (AutoDiff) through the learned model.

**Contributions.** We propose and analyze an alternative to the aforementioned approach of first learning a deep neural model of dynamics, and then performing AutoDiff to conduct the iLQR update. We consider a simplified setting with fixed initial starting condition. Our algorithm maintains a *policy*, specified by an open-loop input sequence and a sequence of stabilizing gains, and loops two steps: **(a)** it learns local linear model of the closed-loop linearized dynamics (in feedback with these gains), which we use to perform a gradient update; **(b)** it re-estimates a linear model after the gradient step, and synthesizes a new set of set gains from this new model. In contrast to past approaches,

our algorithm *only ever estimates linear models of system dynamics*.

For our analysis, we treat the underlying system dynamics as continuous and policy as discrete; this reflects real physical systems, is representative of discrete-time simulated environments which update on smaller timescales than learned policies, and renders explicit the effect of discretization size on sample complexity. We consider an interaction model where we query an oracle for trajectories corrupted with measurement (but not process) noise. Our approach enjoys the following theoretical properties. **1.** Using a number of iterations and oracle queries *polynomial* in relevant problem parameters and tolerance  $\epsilon$ , it computes a policy  $\pi$  whose input sequence is an  $\epsilon$ -first order stationary point for the iLQR approximation of the planning objective (i.e., the gradient through the closed-loop linearized dynamics has norm  $\leq \epsilon$ ). Importantly, learning the linearized model at each iteration obviates the need for global dynamics models, allowing for sample complexity polynomial in dimension.

**2.** We show that contribution 1 implies convergence to a local-optimality criterion we call an  $\epsilon$ -approximate *Jacobian Stationary Point* ( $\epsilon$ -JSP); this roughly equates to the open-loop trajectory under  $\pi$  having cost within  $\epsilon$ -globally optimal for the linearized dynamics about its trajectory.

JSPs are purely a property of the open-loop inputs, allowing comparison of the quality of the open-loop plan with differing gains. Moreover, the results of Westebroek et al. (2021) show that approximate JSPs for certain planning objective enjoy favorable *global properties*, despite (as we show) being computable from (local) gradient-based search (see Appendix B.2 for elaboration).

**Experimental Findings.** We validate our algorithms on standard models of the quadrotor and inverted pendulum, finding an improved performance as iteration number increases, and that the synthesized gains prescribed by iLQR yield improved performance over vanilla gradient updates.

**Proof Techniques.** Central to our analysis are novel perturbation bounds for controlled nonlinear differential equations. Prior results primarily focus on the open-loop setting (Polak, 2012, Theorem 5.6.9), and implicitly hide an exponential dependence on the time horizon for open-loop unstable dynamics. We provide what is to the best of our knowledge the first analysis which demonstrates that local feedback can overcome this pathology. Specifically, we show that if the feedback gains stabilize the Jacobian-linearized dynamics, then (a) the Taylor-remainder of the first-order approximation of the dynamics does *not* scale exponentially on problem horizon (Proposition 4.3), and (b) small perturbations to the nominal input sequence preserve closed-loop stability of the linearized dynamics.

These findings are detailed in Appendix A.6, and enable us to bootstrap the many recent advances in statistical learning for linear systems to our nonlinear setting.

### 1.1. Related Work

iLQR (Li & Todorov, 2004) is a more computationally expedient variant of differential dynamic programming (DPP) (Jacobson & Mayne, 1970); numerous variants exist, notably iLQG (Todorov & Li, 2005) and iLQR (Li & Todorov, 2004), which better address problem stochasticity. iLQR is a predominant approach for the “inner loop” trajectory optimization step in MPC, with applications in robotics (Tassa et al., 2012), quadrotors (Torrente et al., 2021), and autonomous racing (Kabzan et al., 2019).

A considerable literature has combined iLQR with learned dynamics models; here, the Jacobian linearization matrices are typically derived through automated differentiation (Levine & Koltun, 2013; Levine & Abbeel, 2014; Koller et al., 2018), though local kernel least squares regression has also been studied (Rosolia & Borrelli, 2019; Papadimitriou et al., 2020). In these works, the dynamics models are refined/re-estimated as the policy is optimized; thus, these approaches are one instantiation of the broader iterative learning control (ILC) paradigm (Arimoto et al., 1984); other instantiations of ILC include (Kocijan et al., 2004; Dai et al., 2021; Aswani et al., 2013; Bechtel et al., 2020).

Recent years have seen multiple rigorous guarantees for learning system identification and control (Dean et al., 2017; Simchowitz et al., 2018; Oymak & Ozay, 2019; Agarwal et al., 2019; Simchowitz & Foster, 2020), though a general theory of learning for nonlinear control remains elusive. Recent progress includes nonlinear imitation learning (Pfrommer et al., 2022), learning systems with known nonlinearities in the dynamics (Sattar & Oymak, 2022; Foster et al., 2020; Mania et al., 2020) or perception model (Mhammedi et al., 2020; Dean & Recht, 2021).

Lastly, there has been recent theoretical attention given to the study of first-order trajectory optimization methods. Roulet et al. (2019) perform an extension theoretical study of the convergence properties of iLQR, iLQG, and DPP with *exact* dynamics models, and corroborate their findings experimentally. Westebroek et al. (2021) show further that for certain classes of nonlinear systems, all  $\epsilon$ -first order stationary points of a suitable trajectory optimization objective induce trajectories which converge exponentially to desired system equilibria. In some cases, there may be multiple spurious local minima, each of which is nevertheless exponentially stabilizing. Examining the proof (Westebroek et al., 2021) shows the result holds more generally for all  $\epsilon$ -JSPs, and therefore we use their work justify the JSP criterion proposed in this paper.

## 2. Setting

We consider a continuous-time nonlinear control system with state  $\mathbf{x}(t) \in \mathbb{R}^{d_x}$ , input  $\mathbf{u}(t) \in \mathbb{R}^{d_u}$  with finite horizon  $T > 0$ , and fixed initial condition  $\xi_{\text{init}} \in \mathbb{R}^{d_x}$ . We denote the space of bounded input signals  $\mathcal{U} := \{\mathbf{u}(\cdot) : [0, T] \rightarrow \mathbb{R}^{d_u} : \sup_{t \in [0, T]} \|\mathbf{u}(t)\| < \infty\}$ . We endow  $\mathcal{U}$  with an inner product  $\langle \mathbf{u}(\cdot), \mathbf{u}'(\cdot) \rangle_{\mathcal{L}_2(\mathcal{U})} := \int_0^T \langle \mathbf{u}(s), \mathbf{u}'(s) \rangle ds$ , where  $\langle \cdot, \cdot \rangle$  is the standard Euclidean inner product, which induces a norm  $\|\mathbf{u}(\cdot)\|_{\mathcal{L}_2(\mathcal{U})}^2 := \langle \mathbf{u}(\cdot), \mathbf{u}(\cdot) \rangle_{\mathcal{L}_2(\mathcal{U})}$ . For  $\mathbf{u} \in \mathcal{U}$ , the open-loop dynamics are governed by the ordinary differential equation (ODE)

$$\frac{d}{dt} \mathbf{x}(t | \mathbf{u}) = f_{\text{dyn}}(\mathbf{x}(t | \mathbf{u}), \mathbf{u}(t)), \quad \mathbf{x}(0 | \mathbf{u}) = \xi_{\text{init}},$$

where  $f_{\text{dyn}} : \mathbb{R}^{d_x} \times \mathbb{R}^{d_u} \rightarrow \mathbb{R}^{d_x}$  a  $\mathcal{C}^2$  map. Given a terminal cost  $V(\cdot) : \mathbb{R}^{d_x} \rightarrow \mathbb{R}$  and running  $Q(\cdot, \cdot, \cdot) : \mathbb{R}^{d_x} \times \mathbb{R}^{d_u} \times [0, T] \rightarrow \mathbb{R}$ , we optimize the control objective

$$\mathcal{J}_T(\mathbf{u}) := V(\mathbf{x}(T | \mathbf{u})) + \int_{t=0}^T Q(\mathbf{x}(t | \mathbf{u}), \mathbf{u}(t), t) dt.$$

We make the common assumption that the costs are strongly  $\mathcal{C}^2$ , and that  $Q$  is strongly convex:

**Assumption 2.1.** For all  $t \in [0, T]$ ,  $V(\cdot)$  and  $Q(\cdot, \cdot, t)$  are twice-continuously differentiable ( $\mathcal{C}^2$ ), and  $x \mapsto V(x)$  and  $(x, u) \mapsto Q(x, u, t) - \frac{\alpha}{2}(\|x\|^2 + \|u\|^2)$  are convex.

Given a continuously differentiable function  $\mathcal{F} : \mathcal{U} \rightarrow \mathbb{R}^n$  and perturbation  $\delta \mathbf{u} \in \mathcal{U}$ , we define its *directional derivative*  $D\mathcal{F}(\mathbf{u})[\delta \mathbf{u}] := \lim_{\eta \rightarrow 0} \eta^{-1}(\mathcal{F}(\mathbf{u} + \eta \delta \mathbf{u}) - \mathcal{F}(\mathbf{u}))$ . The *gradient*  $\nabla \mathcal{F}(\mathbf{u}) \in \mathcal{U}$  is the (almost-everywhere) unique element of  $\mathcal{U}$  such that  $\forall \delta \mathbf{u} \in \mathcal{U}$ ,  $\int_0^T \nabla \mathcal{F}(\mathbf{u})(t) \delta \mathbf{u}(t) dt = D\mathcal{F}(\mathbf{u})[\delta \mathbf{u}]$ . We denote the gradients of  $\mathbf{u} \mapsto \mathbf{x}(t | \mathbf{u})$  as  $\nabla_{\mathbf{u}} \mathbf{x}(t | \mathbf{u})$ , and of  $\mathbf{u} \mapsto \mathcal{J}_T(\mathbf{u})$  as  $\nabla_{\mathbf{u}} \mathcal{J}_T(\mathbf{u})$ .

**Discretization and Feedback Policies.** Because digital controllers cannot represent continuous open-loop inputs, we compute  $\epsilon$ -JSPs  $\mathbf{u} \in \mathcal{U}$  which are the zero-order holds of discrete-time control sequences. We let  $\tau \in (0, T]$  be a discretization size, and set  $K = \lfloor T/\tau \rfloor$ . Going forward, we denote discrete-time quantities in **colored**, **bold-serif font**.

For  $k \geq 1$ , define  $t_k = (k-1)\tau$ , and define the intervals  $\mathcal{I}_k = [t_k, t_{k+1})$ . For  $t \in [0, T]$ , let  $k(t) := \sup\{k : t_k \leq t\}$ . We let  $\mathbf{U} := (\mathbb{R}^{d_u})^K$ , whose elements are denoted  $\bar{\mathbf{u}} = \mathbf{u}_{1:K}$ , and let  $\text{ct} : \mathbf{U} \rightarrow \mathcal{U}$  denote the natural inclusion  $\text{ct}(\bar{\mathbf{u}})(t) := \mathbf{u}_{k(t)}$ .

Next, to mitigate the curse of horizon, we study *policies* which (a) have discrete-time open-loop inputs and (b) have discrete-time feedback gains to stabilize around the trajectories induced by the nominal inputs. In this work,  $\Pi_\tau$  denotes the set of all policies  $\pi = (\mathbf{u}_{1:K}^\pi, \mathbf{K}_{1:K}^\pi)$  defined by a discrete-time open-loop policy  $\mathbf{u}_{1:K}^\pi \in \mathbf{U}$ , and a sequence of feedback gains  $(\mathbf{K}_k^\pi)_{k \in [K]} \in (\mathbb{R}^{d_x d_u})^K$ . A policy  $\pi$  induces nominal dynamics  $\mathbf{u}^\pi(\cdot) = \text{ct}(\mathbf{u}_{1:K}^\pi)$  and

$\mathbf{x}^\pi(t) = \mathbf{x}(t | \mathbf{u}^\pi)$ ; we set  $\mathbf{x}_k^\pi = \mathbf{x}^\pi(t_k)$ . It also induces the following dynamics by stabilizing around the policy.

**Definition 2.1.** Given a continuous-time input  $\bar{\mathbf{u}} \in \mathcal{U}$ , we define the *stabilized trajectory*  $\tilde{\mathbf{x}}^{\pi, \text{ct}}(t | \bar{\mathbf{u}}) := \mathbf{x}(t | \tilde{\mathbf{u}}^{\pi, \text{ct}})$ , where  $\tilde{\mathbf{u}}^{\pi, \text{ct}}(t | \bar{\mathbf{u}}) := \bar{\mathbf{u}}(t) + \mathbf{K}_{k(t)}^\pi(\tilde{\mathbf{x}}^{\pi, \text{ct}}(t_{k(t)} | \bar{\mathbf{u}}) - \mathbf{x}_{k(t)}^\pi)$ . This induces a stabilized objective:  $\mathcal{J}_T^\pi(\bar{\mathbf{u}}) := V(\tilde{\mathbf{x}}^{\pi, \text{ct}}(T | \bar{\mathbf{u}})) + \int_0^T Q(\tilde{\mathbf{x}}^{\pi, \text{ct}}(t | \bar{\mathbf{u}}), \tilde{\mathbf{u}}^{\pi, \text{ct}}(t | \bar{\mathbf{u}}), t) dt$ . We define the shorthand  $\nabla \mathcal{J}_T^\pi(\pi) := \nabla_{\bar{\mathbf{u}}} \mathcal{J}_T^\pi(\bar{\mathbf{u}})|_{\bar{\mathbf{u}}=\mathbf{u}^\pi}$

Notice that, while  $\pi$  is specified by *discrete-time* inputs,  $\tilde{\mathbf{x}}^{\pi, \text{ct}}(\cdot)$ ,  $\tilde{\mathbf{u}}^{\pi, \text{ct}}(\cdot)$  are *continuous-time* inputs and trajectories stabilized by  $\pi$  and the gradient  $\nabla \mathcal{J}_T^\pi(\cdot)$  is defined over *continuous-time perturbations*.

**Optimization Criteria.** Due to nonlinear dynamics, the objectives  $\mathcal{J}_T, \mathcal{J}_T^\pi$  are nonconvex, so we can only aim for local optimality. Approximate first-order stationary points (FOS) are a natural candidate (Roulet et al., 2019).

**Definition 2.2.** We say  $\mathbf{u}$  is an  $\epsilon$ -FOS of a function  $\mathcal{F} : \mathcal{U} \rightarrow \mathbb{R}$  if  $\|\nabla_{\mathbf{u}} \mathcal{F}(\mathbf{u})\|_{\mathcal{L}_2(\mathcal{U})} \leq \epsilon$ . We say  $\pi$  is  $\epsilon$ -stationary if  $\|\nabla \mathcal{J}_T(\pi)\|_{\mathcal{L}_2(\mathcal{U})} := \|\nabla_{\bar{\mathbf{u}}} \mathcal{J}_T^\pi(\bar{\mathbf{u}})|_{\bar{\mathbf{u}}=\mathbf{u}^\pi}\|_{\mathcal{L}_2(\mathcal{U})} \leq \epsilon$ .

Our primary criterion is to compute  $\epsilon$ -stationary policies  $\pi$ . However, this depends both on the policy inputs  $\mathbf{u}^\pi$  (and induced trajectory  $\mathbf{x}^\pi$ ), as well as the gains. We therefore propose a secondary optimization criterion which depends only on the policies inputs/trajectory. It might be tempting to hope that  $\mathbf{u}^\pi$  is an  $\epsilon$ -FOS of the original objective  $\mathcal{J}_T(\mathbf{u})$ . However, when the Jacobian linearized trajectory (Definition 2.3 below) of the dynamics around  $(\mathbf{x}^\pi, \mathbf{u}^\pi)$  are unstable, the open-loop gradient  $\|\nabla \mathcal{J}_T(\mathbf{u}^\pi)\|_{\mathcal{L}_2(\mathcal{U})}$  can be a factor of  $e^T$  larger than the stabilized gradient  $\|\nabla \mathcal{J}_T(\pi)\|_{\mathcal{L}_2(\mathcal{U})}$  despite the fact that, definitionally,  $\mathcal{J}_T(\mathbf{u}^\pi) = \mathcal{J}_T^\pi(\mathbf{u}^\pi)$  (see Appendix B.1). We therefore propose an alternative definition in terms of Jacobian-linearized trajectory.

**Definition 2.3.** Given  $\mathbf{u}, \bar{\mathbf{u}} \in \mathcal{U}$ , define the Jacobian-linearized (JL) trajectory  $\mathbf{x}^{\text{jac}}(t | \bar{\mathbf{u}}; \mathbf{u}) = \mathbf{x}(t | \mathbf{u}) + \langle \nabla_{\mathbf{u}} \mathbf{x}(t | \mathbf{u}), \bar{\mathbf{u}} - \mathbf{u} \rangle$ , and cost  $\mathcal{J}_T^{\text{jac}}(\bar{\mathbf{u}}; \mathbf{u}) := V(\mathbf{x}^{\text{jac}}(T | \bar{\mathbf{u}}; \mathbf{u})) + \int_{t=0}^T Q(\mathbf{x}^{\text{jac}}(t | \bar{\mathbf{u}}; \mathbf{u}), \bar{\mathbf{u}}(t), t) dt$ .

In words, the JL trajectory is just the first-order Taylor expansion of the dynamics around an input  $\mathbf{u} \in \mathcal{U}$ , and the cost is the cost functional applied to those JL dynamics. We propose an optimization criterion which requires that  $\mathbf{u}$  is near-globally optimal for the JL dynamics around  $\mathbf{u}$ :

**Definition 2.4.** We say  $\mathbf{u} \in \mathcal{U}$  is an  $\epsilon$ -Jacobian Stationary Point (JSP) if  $\mathcal{J}_T(\mathbf{u}) \leq \inf_{\bar{\mathbf{u}} \in \mathcal{U}} \mathcal{J}_T^{\text{jac}}(\bar{\mathbf{u}}; \mathbf{u}) + \epsilon$ .

The consideration of JSPs has three advantages: (1) as noted above, JSPs depend only on a trajectory and not on feedback gains; (2) a JSP is sufficient to ensure that the exponential-stability guarantees derived in Westenbroek

et al. (2021) (and mentioned in the introduction above) hold for certain systems; this provides a link between the local optimality derived in this work and *global* trajectory behavior (see Appendix B.2 for further discussion); (3) despite the potentially exponential-in-horizon gap between gradients of  $\mathcal{J}_T$  and  $\mathcal{J}_T^\pi$ , the following result enables us to compare stationary points of the two objectives in a manner that is *independent* of the horizon  $T$ .

**Proposition 4.1 (informal).** *Suppose  $\pi$  is  $\epsilon$ -stationary, and  $\tau$  is sufficiently small. Then,  $\mathbf{u}^\pi$  is an  $\epsilon'$ -JSP of  $\mathcal{J}_T$ , where  $\epsilon' = \mathcal{O}(\epsilon^2 / (2\alpha(1 + \max_k \|\mathbf{K}_k^\pi\|^2)))$ .*

**Oracle Model and Problem Desideratum.** In light of the above discussion, we aim to compute a approximately stationary policy, whose open-loop is therefore an approximate JSP for the original objective. To do so, we assume access to an oracle which can perform feedback with respect to gains  $\mathbf{K}_k^\pi$ .

**Definition 2.5 (Oracle Dynamics).** Given  $\vec{\mathbf{u}} = \mathbf{u}_{1:K} \in \mathbf{U}$ , we define the *oracle dynamics*  $\mathbf{x}_{\text{orac}}^\pi(t \mid \vec{\mathbf{u}}) := \mathbf{x}(t \mid \text{ct}(\mathbf{u}_{\text{orac},1:K}^\pi(\vec{\mathbf{u}})))$ , where we define  $\mathbf{u}_{\text{orac},k}^\pi(\vec{\mathbf{u}}) := \mathbf{u}_k + \mathbf{K}_k^\pi \mathbf{x}_{\text{orac}}^\pi(t_k \mid \vec{\mathbf{u}})$ , and define  $\mathbf{x}_{\text{orac},k}^\pi(\vec{\mathbf{u}}) := \mathbf{x}_{\text{orac}}^\pi(t_k \mid \vec{\mathbf{u}})$ .

**Oracle 2.1.** *We assume access to an oracle `orac` with variance  $\sigma_{\text{orac}}^2 > 0$ , which given any  $\pi \in \Pi_\tau$  and  $\vec{\mathbf{u}} = \mathbf{u}_{1:K}$ , returns,  $\text{orac}_{\pi,x}(\vec{\mathbf{u}}) \sim \mathcal{N}(\mathbf{x}_{\text{orac},1:K+1}^\pi(\vec{\mathbf{u}}), \mathbf{I}_{(K+1)d_x} \sigma_{\text{orac}}^2)$  and  $\text{orac}_{\pi,u}(\vec{\mathbf{u}}) = \mathbf{u}_{\text{orac},1:K}^\pi(\vec{\mathbf{u}})$*

In words, Oracle 2.1 returns entire trajectories by applying feedback along the gains  $\mathbf{K}_k^\pi$ . The addition of measurement noise is to introduce statistical tradeoffs that prevent near-exact zero-order differentiation; we discuss extensions to process noise in Appendix B.4. Because of this, the oracle trajectory in Definition 2.5 differs from the trajectory dynamics in Definition 2.1 in that the feedback does not subtract off the normal  $\mathbf{x}_k^\pi$ ; thus, the oracle can be implemented without noiseless access to the nominal trajectory. Still, we assume that the feedback applied by the oracle is exact. Having defined our oracle, we specify the following problem desideratum (note below that  $M$  is scaled by  $1/\tau$  to capture the computational burden of finer discretization).

**Desideratum 1.** *Given  $\epsilon, \epsilon'$  and unknown dynamical system  $f_{\text{dyn}}(\cdot, \cdot)$ , compute a policy  $\pi$  for which (a)  $\pi$  is  $\epsilon$ -stationary, and (b)  $\mathbf{u}^\pi$  is an  $\epsilon'$ -JSP of  $\mathcal{J}_T$ , using  $M$  calls to Oracle 2.1, where  $M/\tau$  is polynomial in  $1/\epsilon, 1/\epsilon'$ , and relevant problem parameters.*

**Notation.** We let  $[j : k] := \{j, j+1, \dots, k\}$ , and  $[k] = [1 : k]$ . We use standard-bold for continuous-time quantities ( $\mathbf{x}, \mathbf{u}$ ), and bold-serif for discrete (e.g.  $\mathbf{u}_k^\pi$ ). We let  $\mathbf{u}_{j:k}^\pi = (\mathbf{u}_j, \mathbf{u}_{j+1}, \dots, \mathbf{u}_k)$ . Given vector  $\mathbf{v}$  and matrices  $\mathbf{X}$ , let  $\|\mathbf{v}\|$  and  $\|\mathbf{X}\|$  Euclidean and operator norm, respectively; for clarity, we write  $\|\mathbf{u}_{j:k}^\pi\|_{\ell_2}^2 = \sum_{i=j}^k \|\mathbf{u}_i^\pi\|^2$ . As denoted above,  $\langle \cdot, \cdot \rangle_{\mathcal{L}_2(\mathcal{U})}$  and  $\|\cdot\|_{\mathcal{L}_2(\mathcal{U})}$  denote inner prod-

ucts and norms in  $\mathcal{L}_2(\mathcal{U})$ . We let  $x \vee y := \max\{x, y\}$ , and  $x \wedge y := \min\{x, y\}$ .

### 3. Algorithm

Our iterative approach is summarized in Algorithm 1 and takes in a time step  $\tau > 0$ , horizon  $T$ , a per iteration sample size  $N$ , iteration number  $n_{\text{iter}}$ , a noise variance  $\sigma_w$ , a gradient step size  $\eta > 0$  and a controllability parameter  $k_0$ . The algorithm produces a sequence of policies  $\pi^{(n)} = (\mathbf{u}_{1:K}^{(n)}, \mathbf{K}_{1:K}^{(n)})$ , where  $K = \lfloor T/\tau \rfloor$  is the number of time steps per roll-out. Our algorithm uses the primitive ESTMARKOV( $\pi; N, \sigma_w$ ) (Algorithm 2), which makes  $N$  calls to the oracle to produce estimates  $\hat{\mathbf{x}}_{1:K+1}$  of the nominal state trajectory, and another  $N$  calls with randomly-perturbed inputs of perturbation-variance  $\sigma_w$  to produce estimates  $(\hat{\Psi}_{j,k})_{k < j}$  of the closed-loop Markov parameters associated to the current policy  $\Psi_{\text{cl},j,k}^\pi$ , defined in Definition 4.6. We use a method-of-moments estimator for simplicity. At each iteration  $n$ , Algorithm 1 calls ESTMARKOV( $\pi; N, \sigma_w$ ) first to produce an estimate of the gradient of the closed-loop objective with respect to the current discrete-time nominal inputs. The gradient with respect to the  $k$ -th input  $\mathbf{u}_k^{(n)}$  is given by:

$$\hat{\nabla}_k^{(n)} = \hat{\Psi}_{K+1,k}^\top V_x(\hat{\mathbf{x}}_{K+1}) + Q_u(\hat{\mathbf{x}}_k, \mathbf{u}_k^\pi, t_k) + \tau \sum_{j=k+1}^K \hat{\Psi}_{j,k}^\top (Q_x(\hat{\mathbf{x}}_j, \mathbf{u}_j^\pi, t_j) + (\mathbf{K}_j^\pi)^\top Q_u(\hat{\mathbf{x}}_j, \mathbf{u}_j^\pi, t_j)) \quad (3.1)$$

The form of this estimate corresponds to a natural plug-in estimate of the gradient of the discrete-time objective defined in Definition 4.4. We use this gradient in Eq. (3.1) to update the current input; this update is rolled-out in feedback with the current feedback controller to produce the nominal input  $\mathbf{u}_{1:K}^{(n+1)}$  for the next iteration (Algorithm 1, Line 5). Finally, we call ESTGAINS (Algorithm 3), which synthesizes gains for the new policy using a Riccati-type recursion along a second estimate of the linearized dynamics, produced by unrolling the system with the new nominal input and old gains described above. The algorithm then terminates at  $n_{\text{iter}}$  iterations and chooses the policy with the smallest estimated gradient that was observed.

### 4. Algorithm Analysis

For simplicity, we assume  $K = \lfloor T/\tau \rfloor \in \mathbb{N}$  is integral. In order to state uniform regularity conditions on the dynamics and costs, we fix an *feasible radius*  $R_{\text{feas}} > 0$  and restrict to states and inputs bounded thereby.

**Definition 4.1.** We say  $(x, u) \in \mathbb{R}^{d_x \times d_u}$  are *feasible* if  $\|x\| \vee \|u\| \leq R_{\text{feas}}$ . We say a policy  $\pi$  is *feasible* if  $(2\mathbf{x}^\pi(t), 2\mathbf{u}^\pi(t))$  are feasible for all  $t \in [0, T]$ .

We adopt the following boundedness condition.

**Algorithm 1** Trajectory Optimization

- 1: **Initialize** time step  $\tau > 0$ , horizon  $T \geq \tau$ ,  $K \leftarrow \lceil T/\tau \rceil$ , initial policy  $\pi^{(1)}$ , sample size  $N$ , noise variance  $\sigma_w$ , gradient step size  $\eta$ , controllability parameter  $k_0$ , iteration number  $n_{\text{iter}}$ .
- 2: **for** iterations  $n = 1, 2, \dots, n_{\text{iter}}$  **do**
- 3:  $(\hat{\Psi}_{j,k})_{k < j}, \hat{\mathbf{x}}_{1:K+1} = \text{ESTMARKOV}(\pi; N, \sigma_w)$ .
- 4: Compute  $\hat{\nabla}_k^{(n)}$  in Eq. (3.1)
- 5: Gradient update  $\mathbf{u}_{1:K}^{(n+1)} \leftarrow \text{orac}_{\pi^{(n)}, \mathbf{u}}(\tilde{\mathbf{u}}_{1:K}^{(n)})$ , where  $\tilde{\mathbf{u}}_k^{(n)} := \mathbf{u}_k^{(n)} - \frac{\eta}{\tau} \hat{\nabla}_k^{(n)} - \mathbf{K}_k^{\pi^{(n)}} \hat{\mathbf{x}}_k$ .
- 6: Estimate  $\mathbf{K}_{1:K}^{(n+1)} = \text{ESTGAINS}(\tilde{\pi}^{(n)}; \sigma_w, N, k_0)$ , where  $\tilde{\pi}^{(n)} = (\mathbf{u}^{(n+1)}(\cdot), \mathbf{K}_{1:K}^{(n)})$
- 7: Update policy  $\pi^{(n+1)} = (\mathbf{u}_{1:K}^{(n+1)}, \mathbf{K}_{1:K}^{(n+1)})$
- return**  $\pi^{(n_{\text{out}})}, n_{\text{out}} \in \arg \min_{n \in [n_{\text{iter}}]} \|\hat{\nabla}_k^{(n)}\|$ .

**Algorithm 2** ESTMARKOV( $\pi; N, \sigma_w$ )

- ```

% estimate nominal trajectory
1: for samples  $i = 1, 2, \dots, N$  do
2:   Collect trajectory  $\mathbf{x}_{1:K+1}^{(i)} \sim \text{TrajOrac}_{\pi}(\mathbf{u}_{1:K}^{\pi})$ .
3: Average  $\hat{\mathbf{x}}_{1:K+1} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_{1:K+1}^{(i)}$ 
% estimate perturbed trajectory
4: for samples  $i = 1, 2, \dots, N$  do
5:   Draw  $\mathbf{w}_{1:K}^{(i)}$  uniformly from  $\sigma_w \cdot (\{-1, 1\}^{d_u})^K$ .
6:   Let  $\mathbf{u}_k^{(i)} = \mathbf{u}_k^{\pi} + \mathbf{w}_k^{(i)} - \mathbf{K}_k^{\pi} \hat{\mathbf{x}}_k$ , for  $k \in [K]$ 
7:   Collect trajectory  $\mathbf{y}_{1:K+1}^{(i)} \sim \text{orac}_{\pi, \mathbf{x}}(\mathbf{u}_{1:K}^{(i)})$ .
8: Estimate transition operators  $\hat{\Psi}_{j,k} := \frac{1}{N\sigma_w^2} \sum_{i=1}^N (\mathbf{y}_j^{(i)} - \hat{\mathbf{x}}_j)(\mathbf{w}_k^{(i)})^{\top}$ ,  $k < j$ 
9: return  $(\hat{\Psi}_{j,k})_{k < j}, \hat{\mathbf{x}}_{1:K+1}$ 
    
```

**Condition 4.1.** For all  $n$ , the policies  $\pi^{(n)}$  and  $\tilde{\pi}^{(n)}$  produced by Algorithm 1 are feasible.

If  $\pi$  and  $\tilde{\pi}^{(n)}$  produce bounded inputs, and the resulting state trajectories also remain bounded, then Condition 4.1 will hold for  $R_{\text{feas}} > 0$  sufficiently large. This is a common assumption in the control literature (see e.g. Jadbabaie & Hauser (2001)), as physical systems, such as those with Lagrangian dynamics, will remain bounded under bounded inputs (see Appendix B.3 for discussion).

**Assumption 4.1** (Dynamics regularity).  $f_{\text{dyn}}$  is  $\mathcal{C}^2$ , and for all feasible  $(x, u)$ , the following hold  $\|f_{\text{dyn}}(x, u)\| \leq \kappa_f$ ,  $\|\partial_x f_{\text{dyn}}(x, u)\| \vee \|\partial_u f_{\text{dyn}}(x, u)\| \leq L_f$ ,  $\|\nabla^2 f_{\text{dyn}}(x, u)\| \leq M_f$ .

**Assumption 4.2** (Cost regularity). For all feasible  $(x, u)$ , the following hold  $0 \leq V(x) \vee Q(x, u, t) \leq \kappa_{\text{cost}}$ ,  $\|\partial_x V(x)\| \vee \|\partial_x Q(x, u, t)\| \vee \|\partial_u Q(x, u, t)\| \leq L_{\text{cost}}$ ,  $\|\nabla^2 V(x)\| \vee \|\nabla^2 Q(x, u, t)\| \leq M_{\text{cost}}$ .

To take advantage of stabilizing gains, we require two ad-

**Algorithm 3** ESTGAINS( $\pi; N, \sigma_w, k_0$ )

- 1: **Initialize** number of samples  $N$ , noise variance  $\sigma_w$ , (discrete) controllability window  $k_0 \in \mathbb{N}$
- 2: **Estimate** Markov Parameters  $(\hat{\Psi}_{j,k})_{k < j} = \text{ESTMARKOV}(\pi; N, \sigma_w)$
- 3: **for**  $k = k_0, k_0 + 1, \dots, K$  **do**
- 4: Define  $\hat{\mathbf{B}}_k = \hat{\Psi}_{k+1,k}$
- 5: **Define**  $\hat{\mathbf{C}}_{k,\text{in}} := \hat{\mathbf{C}}_{k-1|k-1, k-k_0+1}$ ,  $\hat{\mathbf{C}}_{k,\text{out}} := \hat{\mathbf{C}}_{k|k-1, k-k_0+1}$
- 6: Define  $\hat{\mathbf{A}}_k := \hat{\mathbf{C}}_{k,\text{out}} \hat{\mathbf{C}}_{k,\text{in}}^{\dagger} - \hat{\mathbf{B}}_k \mathbf{K}_k^{\pi}$
- 7: Set  $\hat{\mathbf{P}}_{K+1} = \mathbf{I}_{d_x}$ .
- 8: **for**  $k = K, K-1, \dots, k_0$  **do**
- 9:  $\hat{\mathbf{K}}_k := (\mathbf{I}_{d_u} + \hat{\mathbf{B}}_k^{\top} \hat{\mathbf{P}}_{k+1} \hat{\mathbf{B}}_k)^{-1} (\hat{\mathbf{B}}_k^{\top} \hat{\mathbf{P}}_{k+1} \hat{\mathbf{A}}_k)$ .
- 10:  $\hat{\mathbf{P}}_k = (\hat{\mathbf{A}}_k + \hat{\mathbf{B}}_k \hat{\mathbf{K}}_k)^{\top} \hat{\mathbf{P}}_{k+1} (\hat{\mathbf{A}}_k + \hat{\mathbf{B}}_k \hat{\mathbf{K}}_k) + \tau (\mathbf{I}_{d_x} + \hat{\mathbf{K}}_k^{\top} \hat{\mathbf{K}}_k)$ .
- 11: Set  $\hat{\mathbf{K}}_k = 0$  for  $k \leq k_0$ .
- 12: **Return**  $\hat{\mathbf{K}}_{1:K}$ .

ditional assumptions, which are defined in terms of the JL dynamics.

**Definition 4.2** (Open-Loop Linearized Dynamics). We define the (open-loop) JL dynamic matrices about  $\pi$  as  $\mathbf{A}_{\text{ol}}^{\pi}(t) = \partial_x f_{\text{dyn}}(\mathbf{x}^{\pi}(t), \mathbf{u}^{\pi}(t))$  and  $\mathbf{B}_{\text{ol}}^{\pi}(t) = \partial_u f_{\text{dyn}}(\mathbf{x}^{\pi}(t), \mathbf{u}^{\pi}(t))$ . We define the *open-loop* JL transition function  $\Phi_{\text{ol}}^{\pi}(s, t)$ , defined for  $t \geq s$  as the solution to  $\frac{d}{ds} \Phi_{\text{ol}}^{\pi}(s, t) = \mathbf{A}_{\text{ol}}^{\pi}(s) \Phi_{\text{ol}}^{\pi}(s, t)$ , with initial condition  $\Phi_{\text{ol}}^{\pi}(t, t) = \mathbf{I}$ .

We first require that stabilizing gains can be synthesized; this is formulated in terms of an upper bound on the cost-to-go for the LQR control problem (Anderson & Moore (2007, Section 2)) induced by the JL dynamics.

**Assumption 4.3** (Stabilizability). Given a *policy*  $\pi$ , and a sequence of controls  $\tilde{\mathbf{u}}(\cdot) \in \mathcal{U}$ , let  $V^{\pi}(t | \tilde{\mathbf{u}}, \xi) = \int_{s=t}^T (\|\tilde{\mathbf{x}}(s)\|^2 + \|\tilde{\mathbf{u}}(s)\|^2) ds + \|\tilde{\mathbf{x}}(T)\|^2$ , under the linearized dynamics  $\frac{d}{ds} \tilde{\mathbf{x}}(s) = \mathbf{A}_{\text{ol}}^{\pi}(s) \tilde{\mathbf{x}}(s) + \mathbf{B}_{\text{ol}}^{\pi}(s) \tilde{\mathbf{u}}(s)$ ,  $\tilde{\mathbf{x}}(t) = \xi$ . We assume that, for all feasible policies,  $\sup_{t \in [0, T]} V^{\pi}(t | \tilde{\mathbf{u}}, \xi) \leq \mu_{\text{ric}} \|\xi\|^2$ . Moreover, we assume (for simplicity) that the initial policy has (a) no gains:  $\mathbf{K}_k^{\pi^{(1)}} = 0$  for all  $k \in [K]$ , and (b) satisfies  $V^{\pi^{(1)}}(t | 0, \xi) \leq \mu_{\text{ric}} \|\xi\|^2$ .

The assumption on  $\pi^{(1)}$  can easily be generalized to accommodate initial policies with stabilizing gains. Our final assumption is controllability (see e.g. Anderson & Moore (2007, Appendix B)), which is necessary for identification of system parameters to synthesize stabilizing gains.

**Assumption 4.4** (Controllability). There exists constants  $t_{\text{ctrl}}, \nu_{\text{ctrl}} > 0$  such that, for all feasible  $\pi$  and  $t \in [t_{\text{ctrl}}, T]$ ,  $\int_{s=t-t_{\text{ctrl}}}^t \Phi_{\text{ol}}^{\pi}(t, s) \mathbf{B}_{\text{ol}}^{\pi}(s) \mathbf{B}_{\text{ol}}^{\pi}(s)^{\top} \Phi_{\text{ol}}^{\pi}(t, s)^{\top} ds \succeq \nu_{\text{ctrl}} \mathbf{I}_{d_x}$

For simplicity, we assume  $k_{\text{ctrl}} := t_{\text{ctrl}}/\tau$  is integral. Finally, to state our theorem, we adopt an asymptotic notation which suppresses all parameters except  $\{T, \tau, \alpha\}$ .

**Definition 4.3** (Asymptotic Notation). We let  $\mathcal{O}_*(\cdot)$  term a term which hides polynomial dependences on  $d_x, d_u, R_{\text{feas}}, \kappa_f, M_f, L_f, \kappa_{\text{cost}}, L_{\text{cost}}, M_{\text{cost}}, \mu_{\text{ric}}, \nu_{\text{ctrl}}, t_{\text{ctrl}}$ , and on  $\exp(t_0 L_f)$ , where  $t_0 = \tau k_0 \geq t_{\text{ctrl}}$ .

Notice that we suppress an *exponential* dependence on our proxy  $t_0$  for the controllability horizon  $t_{\text{ctrl}}$ ; this is because the system cannot be stabilized until the dynamics can be accurately estimated, which requires waiting as long as the controllability window (Chen & Hazan, 2021; Tsiamis et al., 2022). We discuss this dependence further in Appendix B.5. Finally, we state a logarithmic term which addresses high-probability confidence:

$$\iota(\delta) := \log \frac{24T^2 n_{\text{iter}} \max\{d_x, d_u\}}{\tau^2 \delta}. \quad (4.1)$$

We can now state our main theorem, which establishes that, with high probability, for a small enough step size  $\tau$ , and large enough sample size  $N$  and iteration number  $n_{\text{iter}}$ , we obtain an  $\epsilon$ -stationary policy and  $\epsilon'$ -JSP, where  $\epsilon^2, \epsilon'$  scale as  $\text{poly}(T)(\tau^2 + \frac{1}{\tau^2 \sqrt{N}})$ :

**Theorem 1.** Fix  $\delta \in (0, 1)$ , and suppose for the sake of simplicity that  $\tau \leq 1 \leq T$ . Then, there are constants  $c_1, \dots, c_5 = \mathcal{O}_*(1)$  such that if we tune  $\eta = 1/c_1 \sqrt{T}$ ,  $\sigma_w = (\sigma_{\text{orac}}^2 \iota(\delta)/N)^{\frac{1}{4}}$  and  $k_0 \geq k_{\text{ctrl}} + 2$ , then as long as

$$\tau \leq \frac{1}{c_2}, \quad N \geq c_3 \iota(\delta) \max \left\{ \frac{T^2}{\tau^2}, \frac{1}{\tau^4}, \sigma_{\text{orac}}^2 \frac{T^4}{\tau^2}, \frac{\sigma_{\text{orac}}^2}{\tau^8} \right\}.$$

Then, with probability  $1 - \delta$ , if Condition 4.1 and all aforementioned Assumptions hold,

(a) For all  $n \in [n_{\text{iter}}]$ , and  $\pi' \in \{\pi^{(n)}, \tilde{\pi}^{(n)}\}$ ,  $\mu_{\pi', \star} \leq 8\mu_{\text{ric}}$  and  $L_{\pi'} \leq 6 \max\{1, L_f\} \mu_{\text{ric}}$ .

(b)  $\pi = \pi^{(n_{\text{out}})}$  is  $\epsilon$ -stationary, where  $\epsilon^2 = c_4(T\tau^2 + \frac{T^{\frac{3}{2}}}{n_{\text{iter}}}) + c_4(\frac{T^{\frac{7}{2}}}{\tau^2} (\frac{\iota(\delta)^2}{N} + \sigma_{\text{orac}} \sqrt{\frac{\iota(\delta)}{N}}) + \sigma_{\text{orac}}^2 \frac{T^{\frac{3}{2}} \iota(\delta)^2}{N})$ .

(c) For  $\pi = \pi^{(n_{\text{out}})}$ ,  $\mathbf{u}^\pi$  is an  $\epsilon'$ -JSP, where  $\epsilon' = c_5 \frac{\epsilon^2}{\alpha}$ .

As a corollary, we achieve Desideratum 1.

**Corollary 4.1.** For any  $\epsilon, \epsilon' > 0$  and  $\delta \in (0, 1)$ , there exists an appropriate choices of  $\{\tau, N, \eta, \sigma_w\}$  such that Algorithm 1 finds, with probability  $\geq 1 - \delta$ , an  $\epsilon$ -stationary policy  $\pi$  with  $\mathbf{u}^\pi$  being an  $\epsilon'$ -JSP using at most  $M$  oracle calls, where  $M/\tau = \mathcal{O}_*(\text{poly}(T, 1/\epsilon, 1/\epsilon', \log(1/\delta)))$ .

#### 4.1. Analysis Overview

In this section, we provide a high-level sketch of the analysis. Appendix A provides the formal proof, and carefully outlines the organization of the subsequent appendices which establish the subordinate results.

As our policies consists of zero-order hold discrete-time inputs, our analysis is mostly performed in discrete-time.

**Definition 4.4** (Stabilized trajectories, discrete-time inputs). Let  $\bar{\mathbf{u}} \in \mathcal{U}$ , and recall the continuous-input trajectories  $\tilde{\mathbf{x}}^{\pi, \text{ct}}, \tilde{\mathbf{u}}^{\pi, \text{ct}}$  in Definition 2.1. We define  $\tilde{\mathbf{x}}^\pi(t | \bar{\mathbf{u}}) := \tilde{\mathbf{x}}^{\pi, \text{ct}}(t | \text{ct}(\bar{\mathbf{u}}))$  and  $\tilde{\mathbf{u}}^\pi(t | \bar{\mathbf{u}}) := \tilde{\mathbf{u}}^{\pi, \text{ct}}(t | \text{ct}(\bar{\mathbf{u}}))$ , and their discrete samplings  $\tilde{\mathbf{x}}_k^\pi(\bar{\mathbf{u}}) := \tilde{\mathbf{x}}^\pi(t_k | \bar{\mathbf{u}})$  and  $\tilde{\mathbf{u}}_k^\pi(\bar{\mathbf{u}}) := \tilde{\mathbf{u}}^\pi(t_k | \bar{\mathbf{u}})$ . We define a discretized objective  $\mathcal{J}_T^{\pi, \text{disc}}(\bar{\mathbf{u}}) := V(\tilde{\mathbf{x}}_{K+1}^\pi(\bar{\mathbf{u}})) + \tau \sum_{k=1}^K Q(\tilde{\mathbf{x}}_k^\pi(\bar{\mathbf{u}}), \tilde{\mathbf{u}}_k^\pi(\bar{\mathbf{u}}), t_k)$ , and the shorthand  $\mathcal{J}_T^{\text{disc}}(\pi) = \mathcal{J}_T^{\pi, \text{disc}}(\mathbf{u}_{1:K}^\pi)$  and  $\nabla \mathcal{J}_T^{\pi, \text{disc}}(\pi) := \nabla_{\bar{\mathbf{u}}} \mathcal{J}_T^{\pi, \text{disc}}(\bar{\mathbf{u}})|_{\bar{\mathbf{u}}=\mathbf{u}_{1:K}^\pi}$ .

What we shall show is that our algorithm (a) finds a policy  $\pi$  such that  $\|\nabla \mathcal{J}_T^{\text{disc}}(\pi)\|_{\ell_2} \leq \epsilon$  is small, (b) by discretization,  $\|\nabla \mathcal{J}_T^{\text{disc}}(\mathbf{u}^\pi)\|_{\mathcal{L}_2(\mathcal{U})} \leq \epsilon + \mathcal{O}(\tau)$  is small (i.e.  $\pi$  is approximately stationary), and that (c) this implies that  $\mathbf{u}^\pi$  is an approximate-JSP of  $\mathcal{J}_T(\mathbf{u}^\pi)$ . Part (a) requires the most effort, part (b) is a tedious discretization, and part (c) is by Proposition 4.1 stated below. Key in these steps are certain regularity conditions on the policy  $\pi$ . The first is the magnitude of the gains:

**Definition 4.5.** We define an upper bound on the gains of policy  $\pi$  as  $L_\pi := \max\{1, \max_{k \in [K]} \|\mathbf{K}_k^\pi\|\}$ .

This term suffices to translate stationary policies to JSPs:

**Proposition 4.1.** Suppose Assumptions 2.1, 4.1 and 4.2,  $\pi$  is feasible,  $\tau \leq \frac{1}{16L_\pi L_f}$ . Then, if  $\|\nabla \mathcal{J}_T(\pi)\|_{\mathcal{L}_2(\mathcal{U})} \leq \epsilon$ ,  $\mathbf{u}^\pi(t)$  is an  $\epsilon'$ -JSP of  $\mathcal{J}_T$  for  $\epsilon' = 64\epsilon^2 L_\pi^2 / \alpha$ .

*Proof Sketch.* We construct a Jacobian linearization  $\mathcal{J}_T^{\pi, \text{jac}}$  of  $\mathcal{J}_T^\pi$  by analogy to  $\mathcal{J}^{\text{jac}}$ , and define  $\epsilon$ -JSPs of  $\mathcal{J}_T^\pi$  analogously. We show by inverting the gains that an  $\epsilon$ -JSP of  $\mathcal{J}_T^\pi$  is precisely an  $\epsilon$ -JSP of  $\mathcal{J}_T$ . We then establish strong convexity of  $\mathcal{J}_T^{\pi, \text{jac}}$  (non-trivial due to the gains), and use the PL inequality for strongly convex functions to conclude. The formal proof is given in Appendix H.2.  $\square$

To establish parts (a) and (b), we need to measure the stability of the policies. To this end, we first introduce *closed-loop* (discrete-time) linearizations of the dynamics, in terms of which we define a Lyapunov stability modulus.

**Definition 4.6** (Closed-Loop Linearizations). We discretize the open-loop linearizations in Definition 4.2 defining  $\mathbf{A}_{\text{ol}, k}^\pi = \Phi_{\text{ol}}^\pi(t_{k+1}, t_k)$  and  $\mathbf{B}_{\text{ol}, k}^\pi := \int_{s=t_k}^{t_{k+1}} \Phi_{\text{ol}}^\pi(t_{k+1}, s) \mathbf{B}_{\text{ol}}^\pi(s) ds$ . We define an *discrete-time closed-loop* linearization  $\mathbf{A}_{\text{cl}, k}^\pi := \mathbf{A}_{\text{ol}, k}^\pi + \mathbf{B}_{\text{ol}, k}^\pi \mathbf{K}_k^\pi$ , and a discrete closed-loop *transition operator* is defined, for  $1 \leq k_1 \leq k_2 \leq K+1$ ,  $\Phi_{\text{cl}, k_2, k_1}^\pi = \mathbf{A}_{\text{cl}, k_2-1}^\pi \cdot \mathbf{A}_{\text{cl}, k_2-2}^\pi \cdots \mathbf{A}_{\text{cl}, k_1}^\pi$ , with the convention  $\Phi_{\text{cl}, k_1, k_1}^\pi = \mathbf{I}$ . For  $1 \leq k_1 < k_2 \leq K+1$ , we define the closed-loop *Markov operator*  $\Psi_{\text{cl}, k_2, k_1}^\pi := \Phi_{\text{cl}, k_2, k_1+1}^\pi \mathbf{B}_{\text{ol}, k_1}^\pi$ .

**Definition 4.7** (Lyapunov Stability Modulus). Given a policy  $\pi$ , define  $\Lambda_{K+1}^\pi = \mathbf{I}$ , and  $\Lambda_k^\pi = (\mathbf{A}_{\text{cl},k}^\pi)^\top \Lambda_{k+1}^\pi \mathbf{A}_{\text{cl},k}^\pi + \tau \mathbf{I}$ . We define  $\mu_{\pi,\star} := \max_{k \in \{k_0, \dots, K+1\}} \|\Lambda_k^\pi\|$ .

Notice that the stability modulus is taken after step  $k_0$ , which is where we terminate the Riccati recursion in [Algorithm 3](#). We shall show that, with high probability, [Algorithm 1](#) synthesizes policies  $\pi$  which satisfy

$$L_\pi \leq 6 \max\{1, L_f\} \mu_{\text{ric}}, \quad \mu_{\pi,\star} \leq 8 \mu_{\text{ric}}, \quad (4.2)$$

so that  $L_\pi, \mu_{\pi,\star} = \mathcal{O}_\star(1)$ . Going forward, we let  $\mathcal{O}_\pi(\cdot)$  denote a term suppressing polynomials in  $L_\pi, \mu_{\pi,\star}$  and terms  $\mathcal{O}_\star(1)$ ; when  $\pi$  satisfies [Eq. \(4.2\)](#), then  $\mathcal{O}_\pi(\cdot) = \mathcal{O}_\star(\cdot)$ . We say  $x \leq 1/\mathcal{O}_\pi(y)$ , if  $x \leq 1/y'$ , where  $y' = \mathcal{O}_\pi(y)$ . In [Appendix I.3](#), we translate discrete-time stationary points to continuous-time ones, establishing part [\(b\)](#) of the argument.

**Proposition 4.2.** *For  $\pi$  feasible,  $\|\nabla \mathcal{J}_T(\pi)\|_{\mathcal{L}_2(\mathcal{U})} \leq \frac{1}{\sqrt{\tau}} \|\nabla \mathcal{J}_T^{\text{disc}}(\pi)\|_{\ell_2} + \mathcal{O}_\pi(\tau\sqrt{T})$ .*

A more precise statement and explanation of the proof are given in [Appendix A.5](#). The rest of the analysis boils down to [\(a\)](#): finding an approximate stationary point of the time-discretized objective.

## 4.2. Finding a stationary point of $\mathcal{J}_T^{\pi, \text{disc}}$

**Taylor expansion of the dynamics.** To begin, we derive perturbation bounds for solutions to the stabilized ordinary differential equations. Specifically, we provide bounds for when  $\mathbf{u}_{1:K}^\pi$  is perturbed by a sufficiently small input  $\delta \mathbf{u}_{1:K}$ . Our formal proposition, [Appendix A.6](#) states perturbations in both the  $\ell_\infty$  and normalized  $\ell_2$ -norms; for simplicity, state the special case for  $\ell_\infty$ -perturbation.

**Proposition 4.3.** *Let  $\mathbf{u}_{1:K} = \mathbf{u}_k^\pi + \delta \mathbf{u}_{1:K}$ , and suppose  $\max_k \|\delta \mathbf{u}_k\| \leq B_\infty \leq 1/\mathcal{O}_\pi(1)$ . Then, for all  $k \in [K+1]$ ,  $\|\tilde{\mathbf{x}}_k^\pi[\mathbf{u}_{1:K}] - \mathbf{x}_k^\pi - \sum_{j=1}^{k-1} \Psi_{\text{cl},k,j}^\pi \delta \mathbf{u}_j\| \leq \mathcal{O}_\pi(B_\infty^2)$ .*

We also show, that if  $B_\infty = 1/\mathcal{O}_\pi(T)$ , then the policy with  $\pi'$  with the same gains  $\mathbf{K}_k^\pi = \mathbf{K}_k^{\pi'}$  as  $\pi$ , but the perturbed inputs  $\mathbf{u}_k^{\pi'} = \mathbf{u}_k$  at most double its Lyapunov stability modulus  $\mu_{\pi',\star} \leq 2\mu_{\pi,\star}$ . This allows small gradient steps to preserve stability.

**Estimation of linearizations and gradients.** We then argue that by making  $\sigma_w$  small, then to first order, the estimation procedure in [Algorithm 2](#) recovers the *linearization* of the dynamics. The proof combines standard method-of-moments analysis based on matrix Chernoff concentration ([Tropp, 2012](#)) and [Proposition 4.3](#) to argue the dynamics can be approximated by their linearization. Specifically, [Appendix A.7](#) argues that, for all rounds  $n \in [n_{\text{iter}}]$  and  $1 \leq j < k \leq K+1$ , it holds that  $\|\Psi_{\text{cl},k,j}^\pi - \hat{\Psi}_{k,j}\| \leq \text{Err}_\Psi(\delta)$  where  $\text{Err}_\Psi(\delta) = \mathcal{O}_\pi(\sqrt{\frac{\ell(\delta)}{N}}(1 + \frac{\sigma_{\text{orac}}}{\sigma_w} + \sigma_w))$ ,

which can be made to scales as  $N^{-\frac{1}{4}}$  by tuning  $\sigma_w = (\sigma_{\text{orac}}^2 \ell(\delta)/N)^{\frac{1}{4}}$ . From the Markov-recovery error, as well as a simpler bound for recovering  $\mathbf{x}_{1:K}^\pi$  in [Algorithm 2](#) (Lines 1-3), we show accurate recovery of the gradients:  $\max_k \|\hat{\mathbf{V}}_k^{(n)} - (\nabla \mathcal{J}_T^{\text{disc}}(\pi^{(n)}))_k\| \leq T \mathcal{O}_\pi(\text{Err}_\Psi(\delta))$ .

The last step here is to argue that we also approximately recover  $\mathbf{A}_{\text{ol},k}^\pi, \mathbf{B}_{\text{ol},k}^\pi$  in [Algorithm 3](#) for synthesizing the gains: for all  $k \geq k_0$ ,

$$\|\hat{\mathbf{B}}_k - \mathbf{B}_{\text{ol},k}^\pi\| \vee \|\hat{\mathbf{A}}_k - \mathbf{A}_{\text{ol},k}^\pi\| \leq \mathcal{O}_\pi\left(\frac{\text{Err}_{\Psi,\pi}(\delta)}{\tau}\right).$$

This consists of two steps: (1) using controllability to show the matrices  $\hat{\mathbf{C}}_{k,\text{in}}$  in [Algorithm 3](#) are well-conditioned and (2) using closeness of the Markov operators to show that  $\hat{\mathbf{C}}_{k,\text{in}}$  and  $\hat{\mathbf{C}}_{k,\text{out}}$  concentrate around their idealized values. Crucially, we only estimate system matrices for  $k \geq k_0$  to ensure  $\hat{\mathbf{C}}_{k,\text{in}}$  is well-defined, and we use window  $k_0 \geq k_{\text{ctrl}} + 2$  to ensure  $\hat{\mathbf{C}}_{k,\text{out}}$  is sufficiently well-conditioned.

**Concluding the proof.** [Appendices A.8](#) and [A.9](#) conclude the proof with two steps: first, we show that cost-function decreases during the gradient step [Algorithm 1](#) ([Line 5](#)) at round  $n \in [n_{\text{iter}}]$  in proportion to  $-\|\hat{\mathbf{V}}_k^{(n)}\|^2$  (a consequence of the standard smooth descent argument). Here, we also apply the aforementioned result that small gradient steps preserve stability:  $\mu_{\tilde{\pi}^{(n)},\star} \leq 2\mu_{\pi^{(n)},\star}$ . Second, we argue that the gains synthesized by [Algorithm 3](#) ensure that the Lyapunov stability modulus of  $\pi^{(n+1)}$  and the magnitude of its gains stay bounded by an algorithm-independent constant:  $\mu_{\pi^{(n+1)},\star} \leq 4\mu_{\text{ric}} = \mathcal{O}_\star(1)$  and  $L_{\pi^{(n+1)}} \leq \mathcal{O}_\star(1)$ ; we use a novel certainty-equivalence analysis for discretized, time-varying linear systems which may be of independent interest ([Appendix F](#)). By combining these two results, we inductively show that all policies constructed satisfy [\(4.2\)](#), namely they have  $\mu_{\pi,\star}$  and  $L_\pi$  at most  $\mathcal{O}_\star(1)$ . We then combine this with the typical analysis of nonconvex smooth gradient descent to argue that the policy  $\pi^{(n_{\text{out}})}$  has small discretized gradient, as needed.

## 5. Experiments

Our experiments evaluate the performance of our proposed trajectory optimization algorithm ([Algorithm 1](#)) and compare it with the well-established model-based baseline of trajectory optimization (iLQR) on top of learned dynamics (e.g. [Levine & Koltun \(2013\)](#)). Though our analysis considers a fixed horizon, we perform experiments in a receding horizon control (RHC) fashion. We consider two control tasks: [\(a\)](#) a pendulum swing up task, and [\(b\)](#) a 2D quadrotor stabilization task. We implement our experiments using the `jax` ([Bradbury et al., 2018](#)) ecosystem. More details regarding the environments, tasks, and experimental setup details are found in [Appendix J](#). Though our analysis considers the noisy oracle model, all experiments assume *noiseless* observations.

**Least-squares vs. Method-of-Moments.** Algorithm 1 prescribes the method-of-moments estimator to simplify the analysis; in our implementation, we find that estimating the transition operators using regularized least-squares instead yields to more sample efficient gradient estimation. This choice can also be analyzed with minor modifications (see e.g. Oymak & Ozay (2019); Simchowitz et al. (2019)).

**iLQR baseline.** We first collect a training dataset according to a prescribed exploration strategy, then train a neural network dynamics model on these dynamics, and finally optimize our policy by applying the iLQR algorithm directly on the learned model. We consider several variants of our iLQR baseline which use different exploration strategies and different supervision signals for model learning.

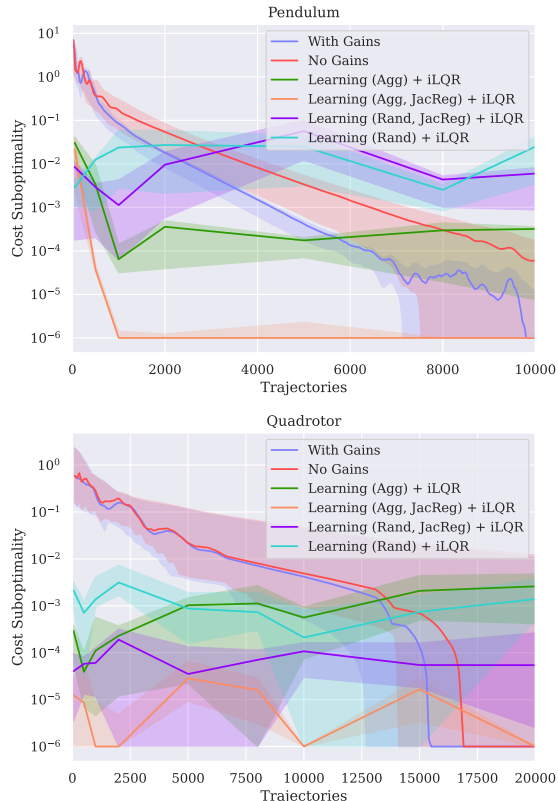
**(1) Sampling strategies:** We consider two sampling strategies; (a) *Agg* alternates between collecting data, fitting a dynamics model, and running iLQR to collect more data, and (b) *Rand* executes rollouts with random inputs starting from random initial conditions. The rationale is that the *Agg* strategy provides better data coverage for the desired task than *Rand*.

**(2) Loss supervision:** The standard loss supervision for learning dynamics is to regress against the next state transition. Inspired by our analysis, we also consider an idealized oracle that augments the supervision to also include noiseless the Jacobians of the ground truth model with respect to both the state and control input; we refer to this augmentation as *JacReg*.

**(2) Model architecture:** We use a fully connected three layer MLP network to for fitting the dynamics of the environment. Specifically, our model takes in input  $(\mathbf{x}_k, \mathbf{u}_k)$  and predicts the state difference  $\mathbf{x}_{k+1} - \mathbf{x}_k$ .

Figure 1 shows the results of Algorithm 1 compared with several iLQR baselines on the pendulum and quadrotor tasks, respectively. In these figures, the x-axis plots the number of trajectories available to each algorithm, and the y-axis plots the cost suboptimality  $(\mathcal{J}_T^{\text{alg}} - \mathcal{J}_T^*)/\mathcal{J}_T^*$  incurred by each algorithm; where  $\mathcal{J}_T^{\text{alg}}$  is algorithmic cost and  $\mathcal{J}_T^*$  is the cost obtained via iLQR with the ground truth dynamics. The error bars in the plot are median, first and third quartile intervals computed over 20 different evaluation seeds.

**Discussion.** We observe that Algorithm 1 with feedback-gains consistently outperforms Algorithm 1 without gains, validating the important of locally-stabilized dynamics. Second, we see that the performance of the iLQR baselines does not significantly improve as more trajectory data is collected. We find that our learned models achieve very low train and test error, over the sampling distribution (i.e., *Agg* or *Rand*) used for learning. For *Rand*, we postulate that the distribution shift incurred by perform-



**Figure 1:** Cost suboptimality  $(\mathcal{J}_T^{\text{alg}} - \mathcal{J}_T^*)/\mathcal{J}_T^*$  versus number of trajectories available to both Algorithm 1 and iLQR baselines. For visualization, the suboptimality is clipped to  $(10^{-6}, \infty)$ .

ing RHC via trajectory optimization on the learned model limits the closed-loop performance of our baseline. However, we note that *Agg+JacReg* achieves stellar performance early on, suggesting that (a) the *Agg* data collection method suffices for strong closed-loop performance (notice that *Rand+JacReg* fares far worse), and (b) that a second limiting factor is that estimating *dynamics* and performing automated differentiation is less favorable than directly estimating *Jacobians*, which are the fundamental quantities used by the iLQR algorithm. This gap between estimation of dynamics and derivatives has been observed in prior work (Pfrommer et al., 2022).

Though we find that our method outperforms deep-learning baselines (excluding *Agg+JacReg*) on the simpler inverted pendulum environment, the learning+iLQR approaches fare better on the quadrotor in the  $\leq 10000$  trajectories regime. We suspect that this is attributable to data-reuse, as Algorithm 1 estimates an entirely new model of system dynamics at each iteration. We believe that finding a way to combine the advantages of directly estimating linearized dynamics (observed in Algorithm 1, as well as *Agg+JacReg*) with the advantages of data-reuse would yield significant sample efficiency improvements.



**References**

- Agarwal, N., Bullins, B., Hazan, E., Kakade, S., and Singh, K. Online control with adversarial disturbances. In *International Conference on Machine Learning*, pp. 111–119. PMLR, 2019.
- Anderson, B. D. and Moore, J. B. *Optimal control: linear quadratic methods*. Courier Corporation, 2007.
- Arimoto, S., Kawamura, S., and Miyazaki, F. Bettering operation of robots by learning. *Journal of Robotic systems*, 1(2):123–140, 1984.
- Aswani, A., Gonzalez, H., Sastry, S. S., and Tomlin, C. Provably safe and robust learning-based model predictive control. *Automatica*, 49(5):1216–1226, 2013.
- Babuschkin, I., Baumli, K., Bell, A., Bhupatiraju, S., Bruce, J., Buchlovsky, P., Budden, D., Cai, T., Clark, A., Danihelka, I., Fantacci, C., Godwin, J., Jones, C., Hemsley, R., Hennigan, T., Hessel, M., Hou, S., Kapturowski, S., Keck, T., Kemaev, I., King, M., Kunesch, M., Martens, L., Merzic, H., Mikulik, V., Norman, T., Quan, J., Papamakarios, G., Ring, R., Ruiz, F., Sanchez, A., Schneider, R., Sezener, E., Spencer, S., Srinivasan, S., Wang, L., Stokowiec, W., and Viola, F. The DeepMind JAX Ecosystem, 2020. URL <http://github.com/deepmind>.
- Bechtle, S., Lin, Y., Rai, A., Righetti, L., and Meier, F. Curious ilqr: Resolving uncertainty in model-based rl. In *Conference on Robot Learning*, pp. 162–171. PMLR, 2020.
- Boucheron, S., Lugosi, G., and Massart, P. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford university press, 2013.
- Bradbury, J., Frostig, R., Hawkins, P., Johnson, M. J., Leary, C., Maclaurin, D., Necula, G., Paszke, A., VanderPlas, J., Wanderman-Milne, S., and Zhang, Q. JAX: composable transformations of Python+NumPy programs, 2018. URL <http://github.com/google/jax>.
- Chen, X. and Hazan, E. Black-box control for linear dynamical systems. In *Conference on Learning Theory*, pp. 1114–1143. PMLR, 2021.
- Dai, H., Landry, B., Yang, L., Pavone, M., and Tedrake, R. Lyapunov-stable neural-network control. *arXiv preprint arXiv:2109.14152*, 2021.
- Dean, S. and Recht, B. Certainty equivalent perception-based control. In *Learning for Dynamics and Control*, pp. 399–411. PMLR, 2021.
- Dean, S., Mania, H., Matni, N., Recht, B., and Tu, S. On the sample complexity of the linear quadratic regulator working draft. 2017.
- Foster, D., Sarkar, T., and Rakhlin, A. Learning nonlinear dynamical systems from a single trajectory. In *Learning for Dynamics and Control*, pp. 851–861. PMLR, 2020.
- Frostig, R., Sindhvani, V., Singh, S., and Tu, S. trajax: differentiable optimal control on accelerators, 2021. URL <http://github.com/google/trajax>.
- Hennigan, T., Cai, T., Norman, T., and Babuschkin, I. Haiku: Sonnet for JAX, 2020. URL <http://github.com/deepmind/dm-haiku>.
- Jacobson, D. H. and Mayne, D. Q. *Differential dynamic programming*. Number 24. Elsevier Publishing Company, 1970.
- Jadbabaie, A. and Hauser, J. On the stability of unconstrained receding horizon control with a general terminal cost. In *Proceedings of the 40th IEEE Conference on Decision and Control (Cat. No. 01CH37228)*, volume 5, pp. 4826–4831. IEEE, 2001.
- Kabzan, J., Hewing, L., Liniger, A., and Zeilinger, M. N. Learning-based model predictive control for autonomous racing. *IEEE Robotics and Automation Letters*, 4(4):3363–3370, 2019.
- Karimi, H., Nutini, J., and Schmidt, M. Linear convergence of gradient and proximal-gradient methods under the polyak-łojasiewicz condition. In *Joint European conference on machine learning and knowledge discovery in databases*, pp. 795–811. Springer, 2016.
- Kocijan, J., Murray-Smith, R., Rasmussen, C. E., and Girard, A. Gaussian process model based predictive control. In *Proceedings of the 2004 American control conference*, volume 3, pp. 2214–2219. IEEE, 2004.
- Koller, T., Berkenkamp, F., Turchetta, M., and Krause, A. Learning-based model predictive control for safe exploration. In *2018 IEEE conference on decision and control (CDC)*, pp. 6059–6066. IEEE, 2018.
- Levine, S. and Abbeel, P. Learning neural network policies with guided policy search under unknown dynamics. *Advances in neural information processing systems*, 27, 2014.
- Levine, S. and Koltun, V. Guided policy search. In *International conference on machine learning*, pp. 1–9. PMLR, 2013.
- Li, W. and Todorov, E. Iterative linear quadratic regulator design for nonlinear biological movement systems. In *ICINCO (1)*, pp. 222–229. Citeseer, 2004.

- Mania, H., Jordan, M. I., and Recht, B. Active learning for nonlinear system identification with guarantees. *arXiv preprint arXiv:2006.10277*, 2020.
- Mhammedi, Z., Foster, D. J., Simchowitz, M., Misra, D., Sun, W., Krishnamurthy, A., Rakhlin, A., and Langford, J. Learning the linear quadratic regulator from nonlinear observations. *Advances in Neural Information Processing Systems*, 33:14532–14543, 2020.
- Morari, M. and Lee, J. H. Model predictive control: past, present and future. *Computers & Chemical Engineering*, 23(4-5):667–682, 1999.
- Nagabandi, A., Konogle, K., Levine, S., and Kumar, V. Deep dynamics models for learning dexterous manipulation. arxiv. *arXiv preprint arXiv:1909.11652*, 10, 2019.
- Oymak, S. and Ozay, N. Non-asymptotic identification of lti systems from a single trajectory. In *2019 American control conference (ACC)*, pp. 5655–5661. IEEE, 2019.
- Papadimitriou, D., Rosolia, U., and Borrelli, F. Control of unknown nonlinear systems with linear time-varying mpc. In *2020 59th IEEE Conference on Decision and Control (CDC)*, pp. 2258–2263. IEEE, 2020.
- Pfrommer, D., Zhang, T. T., Tu, S., and Matni, N. Tasil: Taylor series imitation learning. *arXiv preprint arXiv:2205.14812*, 2022.
- Polak, E. *Optimization: algorithms and consistent approximations*, volume 124. Springer Science & Business Media, 2012.
- Rosolia, U. and Borrelli, F. Learning how to autonomously race a car: a predictive control approach. *IEEE Transactions on Control Systems Technology*, 28(6):2713–2719, 2019.
- Roulet, V., Srinivasa, S., Drusvyatskiy, D., and Harchaoui, Z. Iterative linearized control: stable algorithms and complexity guarantees. In *International Conference on Machine Learning*, pp. 5518–5527. PMLR, 2019.
- Sattar, Y. and Oymak, S. Non-asymptotic and accurate learning of nonlinear dynamical systems. *Journal of Machine Learning Research*, 23(140):1–49, 2022.
- Simchowitz, M. and Foster, D. Naive exploration is optimal for online lqr. In *International Conference on Machine Learning*, pp. 8937–8948. PMLR, 2020.
- Simchowitz, M., Mania, H., Tu, S., Jordan, M. I., and Recht, B. Learning without mixing: Towards a sharp analysis of linear system identification. In *Conference On Learning Theory*, pp. 439–473. PMLR, 2018.
- Simchowitz, M., Boczar, R., and Recht, B. Learning linear dynamical systems with semi-parametric least squares. In *Conference on Learning Theory*, pp. 2714–2802. PMLR, 2019.
- Stewart, G. W. On the perturbation of pseudo-inverses, projections and linear least squares problems. *SIAM review*, 19(4):634–662, 1977.
- Tassa, Y., Erez, T., and Todorov, E. Synthesis and stabilization of complex behaviors through online trajectory optimization. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pp. 4906–4913. IEEE, 2012.
- Todorov, E. and Li, W. A generalized iterative lqg method for locally-optimal feedback control of constrained nonlinear stochastic systems. In *Proceedings of the 2005, American Control Conference, 2005.*, pp. 300–306. IEEE, 2005.
- Torrente, G., Kaufmann, E., Föhn, P., and Scaramuzza, D. Data-driven mpc for quadrotors. *IEEE Robotics and Automation Letters*, 6(2):3769–3776, 2021.
- Tropp, J. A. User-friendly tail bounds for sums of random matrices. *Foundations of computational mathematics*, 12(4):389–434, 2012.
- Tsiamis, A., Ziemann, I. M., Morari, M., Matni, N., and Pappas, G. J. Learning to control linear systems can be hard. In *Conference on Learning Theory*, pp. 3820–3857. PMLR, 2022.
- Vershynin, R. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.
- Westenbroek, T., Simchowitz, M., Jordan, M. I., and Sastry, S. S. On the stability of nonlinear receding horizon control: a geometric perspective. In *2021 60th IEEE Conference on Decision and Control (CDC)*, pp. 742–749. IEEE, 2021.
- Williams, G., Wagener, N., Goldfain, B., Drews, P., Rehg, J. M., Boots, B., and Theodorou, E. A. Information theoretic mpc for model-based reinforcement learning. In *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 1714–1721. IEEE, 2017.
- Xu, X. On the perturbation of the moore–penrose inverse of a matrix. *Applied Mathematics and Computation*, 374:124920, 2020.

## Contents

|          |                                                                                              |           |
|----------|----------------------------------------------------------------------------------------------|-----------|
| <b>1</b> | <b>Introduction</b>                                                                          | <b>1</b>  |
| 1.1      | Related Work . . . . .                                                                       | 2         |
| <b>2</b> | <b>Setting</b>                                                                               | <b>3</b>  |
| <b>3</b> | <b>Algorithm</b>                                                                             | <b>4</b>  |
| <b>4</b> | <b>Algorithm Analysis</b>                                                                    | <b>4</b>  |
| 4.1      | Analysis Overview . . . . .                                                                  | 6         |
| 4.2      | Finding a stationary point of $\mathcal{J}_T^{\pi, \text{disc}}$ . . . . .                   | 7         |
| <b>5</b> | <b>Experiments</b>                                                                           | <b>7</b>  |
| <b>I</b> | <b>Analysis</b>                                                                              | <b>13</b> |
| <b>A</b> | <b>Formal Analysis</b>                                                                       | <b>13</b> |
| A.1      | Organization of the Appendix . . . . .                                                       | 13        |
| A.2      | Notation Review . . . . .                                                                    | 14        |
| A.3      | Full Statement of Main Result . . . . .                                                      | 15        |
| A.4      | Problem Parameters . . . . .                                                                 | 16        |
| A.4.1    | Stability Constants . . . . .                                                                | 16        |
| A.4.2    | Discretization Step Magnitudes . . . . .                                                     | 16        |
| A.4.3    | Taylor Expansion Constants. . . . .                                                          | 17        |
| A.4.4    | Estimation Error Terms. . . . .                                                              | 17        |
| A.5      | Gradient Discretization . . . . .                                                            | 18        |
| A.6      | Main Taylor Expansion Results . . . . .                                                      | 19        |
| A.7      | Estimation Errors . . . . .                                                                  | 19        |
| A.8      | Descent and Stabilization . . . . .                                                          | 20        |
| A.9      | Concluding the proof. . . . .                                                                | 21        |
| A.9.1    | Translating Theorem 3 into Theorem 2 . . . . .                                               | 23        |
| A.9.2    | Proof of Theorem 3 . . . . .                                                                 | 24        |
| <b>B</b> | <b>Discussion and Extensions</b>                                                             | <b>25</b> |
| B.1      | Separation between and Open-Loop and Closed-Loop Gradients . . . . .                         | 25        |
| B.2      | Global Stability Guarantees of JSPs and Consequences of (Westenbroek et al., 2021) . . . . . | 26        |
| B.3      | Projections to ensure boundedness. . . . .                                                   | 26        |
| B.4      | Extensions to include Process Noise . . . . .                                                | 27        |

|          |                                                                                                                                                                |           |
|----------|----------------------------------------------------------------------------------------------------------------------------------------------------------------|-----------|
| B.5      | Discussion of the $\exp(L_f t_0)$ dependence. . . . .                                                                                                          | 27        |
| <b>C</b> | <b>Jacobian Linearizations</b>                                                                                                                                 | <b>27</b> |
| C.1      | Preliminaries . . . . .                                                                                                                                        | 27        |
| C.1.1    | Exact Trajectories . . . . .                                                                                                                                   | 27        |
| C.1.2    | Trajectory Linearizations . . . . .                                                                                                                            | 27        |
| C.1.3    | Jacobian Linearized Dynamics . . . . .                                                                                                                         | 28        |
| C.2      | Characterizations of the Jacobian Linearizations . . . . .                                                                                                     | 29        |
| C.3      | Gradient Computations . . . . .                                                                                                                                | 30        |
| C.4      | Technical Tools . . . . .                                                                                                                                      | 31        |
| <b>D</b> | <b>Taylor Expansions of the Dynamics</b>                                                                                                                       | <b>32</b> |
| D.1      | Proof of Proposition A.5 . . . . .                                                                                                                             | 32        |
| D.2      | Taylor Expansion of the Cost (Lemma A.6) . . . . .                                                                                                             | 37        |
| D.3      | Proof of Lemma A.7 . . . . .                                                                                                                                   | 39        |
| D.4      | Proof of Lemma A.8 . . . . .                                                                                                                                   | 41        |
| <b>E</b> | <b>Estimation Proofs</b>                                                                                                                                       | <b>41</b> |
| E.1      | Estimation of Markov Parameters: Proof of Proposition A.9 . . . . .                                                                                            | 41        |
| E.2      | Error in the Gradient (Proof of Lemma A.10) . . . . .                                                                                                          | 44        |
| E.3      | Discrete-Time Closed-Loop Controllability (Proposition A.11) . . . . .                                                                                         | 45        |
| E.4      | Recovery of State-Transition Matrix (Proposition A.12) . . . . .                                                                                               | 48        |
| <b>F</b> | <b>Certainty Equivalence</b>                                                                                                                                   | <b>51</b> |
| F.1      | Proof Overview of Theorem 4 . . . . .                                                                                                                          | 52        |
| F.1.1    | Proof of Theorem 4 . . . . .                                                                                                                                   | 53        |
| F.2      | Proof of Lemma F.3 . . . . .                                                                                                                                   | 55        |
| F.3      | Proof of Lemma F.4 . . . . .                                                                                                                                   | 57        |
| F.4      | Perturbation on the gains (Lemma F.5) . . . . .                                                                                                                | 58        |
| F.5      | Proof of Lemma F.2 . . . . .                                                                                                                                   | 58        |
| F.6      | Perturbation bounds for Lyapunov Solutions . . . . .                                                                                                           | 59        |
| <b>G</b> | <b>Instantiations of Certainty Equivalence Bound</b>                                                                                                           | <b>62</b> |
| G.1      | Proof of Proposition G.4 . . . . .                                                                                                                             | 63        |
| G.2      | Proof of Proposition G.2 . . . . .                                                                                                                             | 65        |
| G.2.1    | Preliminaries. . . . .                                                                                                                                         | 65        |
| G.2.2    | Controlling the rate of change of $\mathbf{K}(t)$ . . . . .                                                                                                    | 66        |
| G.2.3    | Controlling differences in $\ \mathbf{x}_k^{\text{ct}} - \mathbf{x}_k^{\text{sub}}\ $ and $\ \mathbf{y}_k^{\text{ct}} - \mathbf{y}_k^{\text{sub}}\ $ . . . . . | 67        |

|           |                                                                     |           |
|-----------|---------------------------------------------------------------------|-----------|
| G.2.4     | Concluding the proof of Proposition G.2 . . . . .                   | 69        |
| G.3       | Proof of Lemma G.3 . . . . .                                        | 70        |
| G.4       | Proof of Lemma A.1 . . . . .                                        | 70        |
| <b>H</b>  | <b>Optimization Proofs</b>                                          | <b>73</b> |
| H.1       | Proof of Descent Lemma (Lemma A.13) . . . . .                       | 73        |
| H.2       | Proof of Proposition 4.1 . . . . .                                  | 74        |
| H.2.1     | Proof of Lemma H.3 . . . . .                                        | 75        |
| <b>I</b>  | <b>Discretization Arguments</b>                                     | <b>77</b> |
| I.1       | Discretization of Open-Loop Linearizations . . . . .                | 78        |
| I.2       | Discretization of Transition and Markov Operators . . . . .         | 78        |
| I.3       | Discretization of the Gradient (Proof of Proposition A.4) . . . . . | 81        |
| <b>II</b> | <b>Experiments</b>                                                  | <b>83</b> |
| <b>J</b>  | <b>Experiments Details</b>                                          | <b>83</b> |
| J.1       | Implementation Details . . . . .                                    | 83        |
| J.2       | Environments . . . . .                                              | 84        |
| J.2.1     | Pendulum . . . . .                                                  | 84        |
| J.2.2     | Quadrotor . . . . .                                                 | 84        |
| J.3       | Neural network training . . . . .                                   | 84        |
| J.4       | Least Squares . . . . .                                             | 85        |
| J.5       | Scaling the gain matrix . . . . .                                   | 85        |

## Part I

# Analysis

### A. Formal Analysis

#### A.1. Organization of the Appendix

First, we begin with an outline of [Appendix A](#):

- [Appendix A.2](#) reviews essential notation.
- [Appendix A.3](#) gives a restatement of our main result, [Theorem 1](#), as [Theorem 2](#).

The rest of [Appendix A](#) carries out the proof of [Theorem 2](#). Specially,

- [Appendix A.4](#) defines numerous problem parameters on which our arguments depend.

- [Appendix A.5](#) proves [Corollary A.1](#), a precise statement of [Proposition 4.2](#) in the main text. It does so via an intermediate result, [Proposition A.4](#), which bounds the  $\mathcal{L}_\infty$  difference between the continuous-time gradient, and the imagine of the discrete-time gradient under the continuous-time inclusion map  $\text{ct}(\cdot)$ .
- [Appendix A.6](#) states key results based on Taylor expansions of dynamics around their linearizations, and norms of various derivative-like quantities.
- [Appendix A.7](#) contains the main statements of the various estimation guarantees, notably, the recovery of nominal trajectories, Markov operators, discretized gradients, and linearized transition matrices  $(\mathbf{A}_{\text{ol},k}^\pi, \mathbf{B}_{\text{ol},k}^\pi)$ .
- [Appendix A.8](#) leverages the previous section to demonstrate (a) a certain descent condition holds for each gradient step and (b) that sufficiently accurate estimates of transition matrices lead to the synthesis of gains for which the corresponding policies have bounded stability moduli.
- Finally, [Appendix A.9](#) concludes the proof, as well as states a more granular guarantee in terms of specific problem parameters and not general  $\mathcal{O}_*(\cdot)$  notation.

The rest of [Part I](#) of the Appendix provides the proofs of constituent results. Specifically,

- [Appendix B](#) presents various discussion of main results, as well as gesturing to extensions. Specifically, [Appendix B.1](#) describes the exponential gap between FOSs of  $\mathcal{J}_T$  and JSPs, and [Appendix B.2](#) explains the consequences of combining our result with ([Westenbroek et al., 2021](#)). We discuss how to implement a projection step to ensure [Definition 4.7](#) in [Appendix B.3](#). Finally, we discuss extensions to an oracle with process noise in [Appendix B.4](#).
- [Appendix C](#) presents various computations of Jacobian linearizations, establishing that they do accurately capture first-order expansions.
- [Appendix D](#) proves all the Taylor-expansion like results listed in [Appendix A.6](#).
- [Appendix E](#) proves all the estimation-error bounds in [Appendix A.7](#).
- [Appendix F](#) provides a general certainty-equivalence and Lyapunov stability perturbation results for time-varying, discrete-time linear systems, in the regime that naturally arises when the state matrices are derived from discretizations of continuous-time dynamics.
- [Appendix G](#) instantiates the bounds in [Appendix F](#) to show that the gains synthesized by [Algorithm 2](#) do indeed lead to policies with bounded stability modulus.
- [Appendix H](#) contains the proofs of optimization-related results: the proof of the descent lemma ([Lemma A.13](#) (in [Appendix H.1](#))) and the proof of the conversion between stationary points and JSPs, [Proposition 4.1](#) (in [Appendix H.2](#)).
- Finally, [Appendix I](#) contains various time-discretization arguments, and in particular establishes the aforementioned [Proposition A.4](#) from [Appendix A.5](#).

## A.2. Notation Review

In this section, we review our basic notation.

**Dynamics.** Recall the nominal system dynamics are given by

$$\frac{d}{dt}\mathbf{x}(t \mid \mathbf{u}) = f_{\text{dyn}}(\mathbf{x}(t \mid \mathbf{u}), \mathbf{u}(t)), \quad \mathbf{x}(0 \mid \mathbf{u}) = \xi_{\text{init}}.$$

We recall the definition of various stabilized dynamics.

**Definition 2.1.** Given a continuous-time input  $\bar{\mathbf{u}} \in \mathcal{U}$ , we define the *stabilized trajectory*  $\tilde{\mathbf{x}}^{\pi, \text{ct}}(t \mid \bar{\mathbf{u}}) := \mathbf{x}(t \mid \tilde{\mathbf{u}}^{\pi, \text{ct}})$ , where  $\tilde{\mathbf{u}}^{\pi, \text{ct}}(t \mid \bar{\mathbf{u}}) := \bar{\mathbf{u}}(t) + \mathbf{K}_{k(t)}^\pi (\tilde{\mathbf{x}}^{\pi, \text{ct}}(t_{k(t)} \mid \bar{\mathbf{u}}) - \mathbf{x}_{k(t)}^\pi)$ . This induces a stabilized objective:  $\mathcal{J}_T^\pi(\bar{\mathbf{u}}) := V(\tilde{\mathbf{x}}^{\pi, \text{ct}}(t \mid \bar{\mathbf{u}}) + \int_0^T Q(\tilde{\mathbf{x}}^{\pi, \text{ct}}(t \mid \bar{\mathbf{u}}), \tilde{\mathbf{u}}^{\pi, \text{ct}}(t \mid \bar{\mathbf{u}}), t) dt$ . We define the shorthand  $\nabla \mathcal{J}_T(\pi) := \nabla_{\bar{\mathbf{u}}} \mathcal{J}_T^\pi(\bar{\mathbf{u}})|_{\bar{\mathbf{u}}=\mathbf{u}^\pi}$

**Definition 4.4** (Stabilized trajectories, discrete-time inputs). Let  $\bar{\mathbf{u}} \in \mathcal{U}$ , and recall the continuous-input trajectories  $\tilde{\mathbf{x}}^{\pi, \text{ct}}, \tilde{\mathbf{u}}^{\pi, \text{ct}}$  in Definition 2.1. We define  $\tilde{\mathbf{x}}^\pi(t | \bar{\mathbf{u}}) := \tilde{\mathbf{x}}^{\pi, \text{ct}}(t | \text{ct}(\bar{\mathbf{u}}))$  and  $\tilde{\mathbf{u}}^\pi(t | \bar{\mathbf{u}}) := \tilde{\mathbf{u}}^{\pi, \text{ct}}(t | \text{ct}(\bar{\mathbf{u}}))$ , and their discrete samplings  $\tilde{\mathbf{x}}_k^\pi(\bar{\mathbf{u}}) := \tilde{\mathbf{x}}^\pi(t_k | \bar{\mathbf{u}})$  and  $\tilde{\mathbf{u}}_k^\pi(\bar{\mathbf{u}}) := \tilde{\mathbf{u}}^\pi(t_k | \bar{\mathbf{u}})$ . We define a discretized objective  $\mathcal{J}_T^{\pi, \text{disc}}(\bar{\mathbf{u}}) := V(\tilde{\mathbf{x}}_{K+1}^\pi(\bar{\mathbf{u}})) + \tau \sum_{k=1}^K Q(\tilde{\mathbf{x}}_k^\pi(\bar{\mathbf{u}}), \tilde{\mathbf{u}}_k^\pi(\bar{\mathbf{u}}), t_k)$ , and the shorthand  $\mathcal{J}_T^{\text{disc}}(\pi) = \mathcal{J}_T^{\pi, \text{disc}}(\mathbf{u}_{1:K}^\pi)$  and  $\nabla \mathcal{J}_T^{\pi, \text{disc}}(\pi) := \nabla_{\bar{\mathbf{u}}} \mathcal{J}_T^{\pi, \text{disc}}(\bar{\mathbf{u}})|_{\bar{\mathbf{u}}=\mathbf{u}_{1:K}^\pi}$ .

**Linearizations.** Next, we recall the definition of the various linearizations.

**Definition 4.2** (Open-Loop Linearized Dynamics). We define the (open-loop) JL dynamic matrices about  $\pi$  as  $\mathbf{A}_{\text{ol}}^\pi(t) = \partial_x f_{\text{dyn}}(\mathbf{x}^\pi(t), \mathbf{u}^\pi(t))$  and  $\mathbf{B}_{\text{ol}}^\pi(t) = \partial_u f_{\text{dyn}}(\mathbf{x}^\pi(t), \mathbf{u}^\pi(t))$ . We define the *open-loop* JL transition function  $\Phi_{\text{ol}}^\pi(s, t)$ , defined for  $t \geq s$  as the solution to  $\frac{d}{ds} \Phi_{\text{ol}}^\pi(s, t) = \mathbf{A}_{\text{ol}}^\pi(s) \Phi_{\text{ol}}^\pi(s, t)$ , with initial condition  $\Phi_{\text{ol}}^\pi(t, t) = \mathbf{I}$ .

**Definition 4.6** (Closed-Loop Linearizations). We discretize the open-loop linearizations in Definition 4.2 defining  $\mathbf{A}_{\text{ol}, k}^\pi = \Phi_{\text{ol}}^\pi(t_{k+1}, t_k)$  and  $\mathbf{B}_{\text{ol}, k}^\pi := \int_{s=t_k}^{t_{k+1}} \Phi_{\text{ol}}^\pi(t_{k+1}, s) \mathbf{B}_{\text{ol}}^\pi(s) ds$ . We define an *discrete-time closed-loop* linearization  $\mathbf{A}_{\text{cl}, k}^\pi := \mathbf{A}_{\text{ol}, k}^\pi + \mathbf{B}_{\text{ol}, k}^\pi \mathbf{K}_k^\pi$ , and a discrete closed-loop *transition operator* is defined, for  $1 \leq k_1 \leq k_2 \leq K+1$ ,  $\Phi_{\text{cl}, k_2, k_1}^\pi = \mathbf{A}_{\text{cl}, k_2-1}^\pi \cdot \mathbf{A}_{\text{cl}, k_2-2}^\pi \cdots \mathbf{A}_{\text{cl}, k_1}^\pi$ , with the convention  $\Phi_{\text{cl}, k_1, k_1}^\pi = \mathbf{I}$ . For  $1 \leq k_1 < k_2 \leq K+1$ , we define the closed-loop *Markov operator*  $\Psi_{\text{cl}, k_2, k_1}^\pi := \Phi_{\text{cl}, k_2, k_1+1}^\pi \mathbf{B}_{\text{ol}, k_1}^\pi$ .

We also recall the definition of stationary policies and JSPs.

**Definition 2.2.** We say  $\mathbf{u}$  is an  $\epsilon$ -FOS of a function  $\mathcal{F} : \mathcal{U} \rightarrow \mathbb{R}$  if  $\|\nabla_{\mathbf{u}} \mathcal{F}(\mathbf{u})\|_{\mathcal{L}_2(\mathcal{U})} \leq \epsilon$ . We say  $\pi$  is  $\epsilon$ -stationary if  $\|\nabla \mathcal{J}_T(\pi)\|_{\mathcal{L}_2(\mathcal{U})} := \|\nabla_{\bar{\mathbf{u}}} \mathcal{J}_T^\pi(\bar{\mathbf{u}})|_{\bar{\mathbf{u}}=\mathbf{u}^\pi}\|_{\mathcal{L}_2(\mathcal{U})} \leq \epsilon$ .

**Definition 2.4.** We say  $\mathbf{u} \in \mathcal{U}$  is an  $\epsilon$ -Jacobian Stationary Point (JSP) if  $\mathcal{J}_T(\mathbf{u}) \leq \inf_{\bar{\mathbf{u}} \in \mathcal{U}} \mathcal{J}_T^{\text{jac}}(\bar{\mathbf{u}}; \mathbf{u}) + \epsilon$ .

**Problem Constants.** We recall the dynamics-constants  $\kappa_f, L_f, M_f$  defined in Assumption 4.1,  $\kappa_{\text{cost}}, L_{\text{cost}}, M_{\text{cost}}$  in Assumption 4.2, the strong-convexity parameter  $\alpha$  in Assumption 2.1, the controllability parameters  $t_{\text{ctrl}}, \nu_{\text{ctrl}}$  from Assumption 4.4, with  $k_{\text{ctrl}} := t_{\text{ctrl}}/\tau$ , and the Riccati parameter  $\mu_{\text{ric}}$  from Assumption 4.3. Finally, we recall the feasibility radius from Condition 4.1. We also recall

**Definition 4.1.** We say  $(x, u) \in \mathbb{R}^{d_x \times d_u}$  are *feasible* if  $\|x\| \vee \|u\| \leq R_{\text{feas}}$ . We say a policy  $\pi$  is feasible if  $(2\mathbf{x}^\pi(t), 2\mathbf{u}^\pi(t))$  are feasible for all  $t \in [0, T]$ .

**Gradient and Cost Shorthands.** Notably, we bound out the following shorthand for gradients and costs:

$$\nabla \mathcal{J}_T(\pi) := \nabla_{\mathbf{u}} \mathcal{J}_T^\pi(\mathbf{u})|_{\mathbf{u}=\mathbf{u}^\pi}, \quad \nabla \mathcal{J}_T^{\text{disc}}(\pi) := \nabla_{\bar{\mathbf{u}}} \mathcal{J}_T^{\pi, \text{disc}}(\bar{\mathbf{u}})|_{\bar{\mathbf{u}}=\mathbf{u}_{1:K}^\pi}, \quad \mathcal{J}_T^{\text{disc}}(\pi) := \mathcal{J}_T^{\pi, \text{disc}}(\mathbf{u}_{1:K}^\pi). \quad (\text{A.1})$$

### A.3. Full Statement of Main Result

The following is a slightly more general statement of Theorem 1, which implies Theorem 1 for appropriate choice of  $\eta \leftarrow \frac{1}{c_1} \min \left\{ \frac{1}{\sqrt{T}}, 1 \right\}$ ,  $\sigma_w \leftarrow (\sigma_{\text{orac}}^2 \iota(\delta)/N)^{\frac{1}{4}}$ , and with the simplifications  $\tau \leq 1 \leq T$ .

**Theorem 2.** Fix  $\delta \in (0, 1)$ , define  $\iota(\delta) := \log \frac{24T^2 n_{\text{iter}} \max\{d_x, d_u\}}{\tau^2 \delta}$  and  $\text{Err}_0(\delta) := \sqrt{\iota(\delta)/N}$ , where  $N$  is the sample size, and suppose we select  $\sigma_w = c(\sigma_{\text{orac}}^2 \iota(\delta)/N)^{\frac{1}{4}}$  for any  $c \in [\frac{1}{\mathcal{O}_*(1)}, \mathcal{O}_*(1)]$ . Then, there exists constants  $c_1, c_2, \dots, c_5 = \mathcal{O}_*(1)$  depending on  $c$  such that the following holds. Suppose that

$$\eta \leq \frac{1}{c_1} \min \left\{ \frac{1}{\sqrt{T}}, 1 \right\}, \quad \tau \leq \frac{1}{c_2} \quad N \geq c_3 \iota(\delta) \max \left\{ 1, \frac{T^2}{\tau^2}, \frac{1}{\tau^4}, T, \sigma_{\text{orac}}^2 \frac{T^4}{\tau^2}, \frac{\sigma_{\text{orac}}^2}{\tau^8} \right\}. \quad (\text{A.2})$$

Then, with probability  $1 - \delta$ , if Condition 4.1 and all listed Assumptions hold,

(a) For all  $n \in [n_{\text{iter}}]$ , and  $\pi' \in \{\pi^{(n)}, \tilde{\pi}^{(n)}\}$ ,  $\mu_{\pi', \star} \leq 8\mu_{\text{ric}}$  and  $L_{\pi'} \leq 6 \max\{1, L_f\} \mu_{\text{ric}}$ .

(a) For  $\pi = \pi^{(n_{\text{out}})}$ ,  $\pi$  is  $\epsilon$ -stationary where

$$\epsilon^2 = c_4 T \left( \tau^2 + \frac{1}{\eta} \left( \frac{1}{n_{\text{iter}}} + \left( \eta T^2 + \frac{T^2}{\tau^2} \right) \left( \frac{\iota(\delta)^2}{N} + \sigma_{\text{orac}} \sqrt{\frac{\iota(\delta)}{N}} \right) + \sigma_{\text{orac}}^2 \frac{\iota(\delta)^2}{N} \right) \right).$$

(c) For  $\pi = \pi^{(n_{\text{out}})}$ ,  $\mathbf{u}^\pi$  is an  $\epsilon'$ -JSP, where  $\epsilon' = c_5 \frac{\epsilon^2}{\alpha}$ .

#### A.4. Problem Parameters

In this section, we provide all definitions of various problem parameters. The notation is extensive, but we maintain the following conventions:

1.  $\mu_{(\cdot)}$  refers to upper bounds on Lyapunov operators,  $\kappa_{(\cdot)}$  to upper bounds on zero-order terms (e.g.  $\|f_{\text{dyn}}(x, u)\|$ ) or magnitudes of transition operators,  $M_{(\cdot)}$  to bounds on second-order derivatives,  $L_{(\cdot)}$  to bounds on first-order derivatives,  $B_{(\cdot)}$  to upper bounds on radii,  $\tau_{(\cdot)}$  to step sizes,  $\text{Err}_{(\cdot)}$  to error terms.
2.  $q \in \{1, 2, \infty\}$  corresponds to  $\ell_q$  norms
3. Subscripts  $\text{tay}$  denote relevance to Taylor expansions of the dynamics.
4. Terms with have a subscript  $\pi$  hide dependence on  $L_\pi$ ,  $\mu_{\pi, \star}$  and  $\kappa_q$  for  $q \in \{1, 2, \infty\}$

**Remark A.1** (Reminder on Asymptotic Notation). We let  $\mathcal{O}_\star(x)$  denote a term which suppresses polynomial dependence on all the constants in [Assumptions 4.1](#) and [4.2](#), as well as  $\mu_{\text{ric}}$  in [Assumption 4.3](#), and  $\nu_{\text{ctrl}}, t_0 \geq t_{\text{ctrl}}$  and  $e^{L_f t_0} \geq e^{L_f t_{\text{ctrl}}}$ , where  $t_0 = \tau k_0$ , and  $\nu_{\text{ctrl}}, t_{\text{ctrl}}$  are given in [Assumption 4.4](#). We let  $\mathcal{O}_\pi(x)$  suppress all of these constants, as well as polynomials in  $L_\pi$  and  $\mu_{\pi, \star}$ .

##### A.4.1. STABILITY CONSTANTS

We begin by recalling the primary constants controlling the stability of a policy  $\pi$ .

**Definition 4.7** (Lyapunov Stability Modulus). Given a policy  $\pi$ , define  $\Lambda_{K+1}^\pi = \mathbf{I}$ , and  $\Lambda_k^\pi = (\mathbf{A}_{\text{cl}, k}^\pi)^\top \Lambda_{k+1}^\pi \mathbf{A}_{\text{cl}, k}^\pi + \tau \mathbf{I}$ . We define  $\mu_{\pi, \star} := \max_{k \in \{k_0, \dots, K+1\}} \|\Lambda_k^\pi\|$ .

It is more convenient to prove bounds in terms of the following three quantity, which are defined in terms of the magnitudes of the closed-loop transition operators.

**Definition A.1** (Norms of  $\pi$ ). We define the constants  $\kappa_{\pi, \infty} := \max_{1 \leq j \leq k \leq K+1} \|\Phi_{\text{cl}, k, j}^\pi\|$ , and

$$\begin{aligned} \kappa_{\pi, 1} &:= \max_{k \in [K+1]} \tau \left( \sum_{j=1}^k \|\Phi_{\text{cl}, k, j}^\pi\| \vee \sum_{j=k}^{K+1} \|\Phi_{\text{cl}, j, k}^\pi\| \right) \\ \kappa_{\pi, 2}^2 &:= \max_{k \in [K+1]} \tau \left( \sum_{j=1}^k \|\Phi_{\text{cl}, k, j}^\pi\|^2 \vee \sum_{j=k}^{K+1} \|\Phi_{\text{cl}, j, k}^\pi\|^2 \right) \end{aligned}$$

We also define the following upper bounds on these quantities:

$$\begin{aligned} \kappa_\infty(\mu, L) &:= \sqrt{\max\{1, 6L_f L\} \mu \exp(t_0 L_f)} \\ \kappa_2(\mu, L) &:= \max\{1, 6L_f L\} \mu (t_0 \exp(2t_0 L_f) + \mu) \\ \kappa_1(\mu, L) &:= \sqrt{\max\{1, 6L_f L\} \mu (t_0 \exp(t_0 L_f) + 2\mu)} \end{aligned}$$

The following lemma is proven in [Appendix G.4](#), and shows that each of the above terms is  $\mathcal{O}_\pi(1)$ .

**Lemma A.1.** *Let  $\pi$  be any policy. Recall  $t_0 = \tau k_0$ . Then, as long as  $\tau \leq 1/6L_f L_\pi$ ,*

$$\mu_{\pi, q} \leq \mu_q(\mu_{\pi, \star}, L_\pi) = \mathcal{O}_\pi(1).$$

##### A.4.2. DISCRETIZATION STEP MAGNITUDES

Next, we introduce various maximal discretization step sizes for which our discrete-time dynamics are sufficiently faithful to the continuous ones. The first is a general condition for the dynamics to be “close”, the second is useful for closeness of solutions of Riccati equations, the third for the discrete-time dynamics to admit useful Taylor expansions, and the fourth for discrete-time controllability. We note that the first two do not depend on  $\pi$ , while the second two do.



**Definition A.2** (Discretization Sizes). We define

$$\begin{aligned}\tau_{\text{dyn}} &:= \frac{1}{4L_f} \\ \tau_{\text{ric}} &:= \frac{1}{4\mu_{\text{ric}}^2 \left( 3M_f \kappa_f \mu_{\text{ric}} L_f + 13L_f^2 (1 + L_f \mu_{\text{ric}})^2 \right)} \\ \tau_{\text{tay}, \pi} &:= \min \left\{ \frac{1}{16L_f L_\pi}, \frac{1}{8\kappa_f} \right\} \leq \tau_{\text{dyn}}. \\ \tau_{\text{ctrl}, \pi} &:= \frac{\nu_{\text{ctrl}}}{8L_\pi^2 K_\pi^2 \gamma_{\text{ctrl}}^3 \exp(2\gamma_{\text{ctrl}}) \left( \kappa_f M_f + 2L_f^2 \right)}, \quad \gamma_{\text{ctrl}} := \max\{1, L_f t_{\text{ctrl}}\}\end{aligned}$$

We note that  $\tau_{\text{dyn}}, \tau_{\text{ric}} = 1/\mathcal{O}_*(1)$  and  $\tau_{\text{tay}, \pi}, \tau_{\text{ctrl}, \pi} = 1/\mathcal{O}_\pi(1)$ .

#### A.4.3. TAYLOR EXPANSION CONSTANTS.

We now define the relevant constants in terms of which we bound our Taylor expansions.

**Definition A.3** (Taylor Expansion Constants, Policy Dependent). We define  $L_{\text{tay}, \infty, \pi} = 2L_f \kappa_{\pi, 1}$ ,  $L_{\text{tay}, 2, \pi} := 2L_f \kappa_{\pi, 2}$ , and

$$\begin{aligned}M_{\text{tay}, 2, \pi} &:= 8M_f (\kappa_{\pi, \infty} + 10L_\pi^2 L_f^2 \kappa_{\pi, 2}^2 \kappa_{\pi, 1}) \\ M_{\text{tay}, \text{inf}, \pi} &:= 8M_f (\kappa_{\pi, 1} + 10L_\pi^2 L_f^2 \kappa_{\pi, 1}^3) \\ B_{\text{tay}, 2, \pi} &= \min \left\{ \frac{1}{\sqrt{40M_f L_\pi^2 \kappa_{\pi, 1} M_{\text{tay}, 2, \pi}}}, \frac{L_f \kappa_{\pi, 2}}{2M_{\text{tay}, 2, \pi}}, \frac{R_{\text{feas}}}{16L_\pi L_f \kappa_{\pi, 2}} \right\} \\ B_{\text{tay}, \text{inf}, \pi} &= \min \left\{ \frac{1}{40L_\pi^2 \kappa_{\pi, 1} M_{\text{tay}, \text{inf}, \pi}}, \frac{L_f \kappa_{\pi, 1}}{2M_{\text{tay}, \text{inf}, \pi}}, \frac{R_{\text{feas}}}{16L_\pi L_f \kappa_{\pi, 1}} \right\}\end{aligned}$$

We also define

$$\begin{aligned}M_{\mathcal{J}, \text{tay}, \pi} &:= 2M_{\text{cost}} L_f^2 \kappa_{\pi, 2}^2 (1 + 3L_\pi^2 T) M_{\text{tay}, 2, \pi} + L_{\text{cost}} (1 + 2L_\pi T) M_{\text{tay}, 2, \pi} + 2L_\pi L_{\text{cost}}, \\ B_{\text{stab}, \pi} &:= (\max\{6, 36L_f L_\pi\} \mu_{\pi, \star} \cdot 12T M_f L_\pi (1 + L_f K_\pi) B_\infty)^{-1} \\ L_{\nabla, \pi, \infty} &:= L_{\text{cost}} \left( 1 + \frac{3L_f}{2} \kappa_{\pi, \infty} + 3L_\pi \kappa_{\pi, 1} \right)\end{aligned}$$

The following is a consequence of [Lemma A.1](#).

**Lemma A.2.** By [Lemma A.1](#),  $M_{\text{tay}, 2, \pi}, M_{\text{tay}, \text{inf}, \pi}, L_{\nabla, \pi, \infty}, L_{\text{tay}, q, \pi} = \mathcal{O}_\pi(1)$ ,  $B_{\text{tay}, 2, \pi}, B_{\text{tay}, \text{inf}, \pi} = 1/\mathcal{O}_\pi(1)$ ,  $M_{\mathcal{J}, \text{tay}, \pi} = T \cdot \mathcal{O}_\pi(1)$ , and  $B_{\text{stab}, \pi} = \frac{1}{T} \cdot \mathcal{O}_\pi(1)$ .

The first group of four constants arises in Taylor expansions of the dynamics, the fifth in a Taylor expansion of the cost functional, and the sixth in controlling the stability of policies under changes to the input, and the last upper bounds the norm of the gradient.

#### A.4.4. ESTIMATION ERROR TERMS.

Finally, we define the following error terms which arise in the errors of the estimated nominal trajectories, Markov operators, and gradients. Note that the first term has no dependence on  $\pi$ , while the latter two do.

**Definition A.4** (Error Terms). Define  $\iota(\delta) := \log \frac{24T^2 n_{\text{iter}} \max\{d_x, d_u\}}{\tau^2 \delta} = \log \frac{24K^2 n_{\text{iter}} d_\star}{\delta}$ , where  $d_\star := \max\{d_x, d_u\}$ .

Further, define

$$\begin{aligned}\text{Err}_{\hat{x}}(\delta) &:= \sigma_{\text{orac}} \sqrt{2 \frac{d_{\star} \iota(\delta)}{N}} \\ \text{Err}_{\Psi, \pi}(\delta) &:= \sqrt{\frac{\iota(\delta)}{N}} \left( \frac{2\sigma_{\text{orac}}}{\sigma_w} d_{\star}^{3/2} + 8L_{\text{tay}, \infty, \pi} d_{\star} \right) + 4\sigma_w M_{\text{tay}, 2, \pi} d_{\star}^{3/2} = \mathcal{O}_{\pi} \left( \sqrt{\frac{\iota(\delta)}{N}} \left( 1 + \frac{\sigma_{\text{orac}}}{\sigma_w} + \sigma_w \right) \right) \\ \text{Err}_{\nabla, \pi}(\delta) &:= (L_{\text{cost}} \text{Err}_{\Psi, \pi}(\delta) + (1 + \kappa_{\pi, \infty}) M_{\text{cost}} \text{Err}_{\hat{x}}(\delta)) (1 + 2TL_{\pi}).\end{aligned}$$

We note that, in view of [Lemma A.1](#),

By [Lemmas A.1](#) and [A.2](#), we have

**Lemma A.3.** Define  $\text{Err}_0(\delta) = \sqrt{\frac{\iota(\delta)}{N}}$ . Then,

$$\begin{aligned}\text{Err}_{\hat{x}}(\delta) &= \sigma_{\text{orac}} \sqrt{2d_{\star}} \text{Err}_0(\delta) \leq \mathcal{O}_{\star}(\text{Err}_0(\delta)) \\ \text{Err}_{\Psi}(\delta) &\leq \mathcal{O}_{\pi} \left( \text{Err}_0(\delta) \left( 1 + \frac{\sigma_{\text{orac}}}{\sigma_w} \right) + \sigma_w \right) \\ \text{Err}_{\nabla, \pi}(\delta) &\leq \mathcal{O}_{\pi} \left( T \left( \text{Err}_0(\delta) \left( 1 + \sigma_{\text{orac}} + \frac{\sigma_{\text{orac}}}{\sigma_w} \right) + \sigma_w \right) \right).\end{aligned}\tag{A.3}$$

If we further tune  $\sigma_w = c\sqrt{\sigma_{\text{orac}} \text{Err}_0(\delta)}$  for any  $c \in [1/\mathcal{O}_{\star}(1), \mathcal{O}_{\star}(1)]$ , then

$$\begin{aligned}\text{Err}_{\hat{x}}(\delta) &\leq \mathcal{O}_{\star}(\sigma_{\text{orac}} \text{Err}_0(\delta)) \\ \text{Err}_{\Psi}(\delta) &\leq \mathcal{O}_{\pi} \left( \text{Err}_0(\delta) + \sqrt{\sigma_{\text{orac}} \text{Err}_0(\delta)} \right) \\ \text{Err}_{\nabla, \pi}(\delta) &\leq \mathcal{O}_{\pi} \left( T \left( \text{Err}_0(\delta) + \sqrt{\sigma_{\text{orac}} \text{Err}_0(\delta)} \right) \right).\end{aligned}\tag{A.4}$$

### A.5. Gradient Discretization

We begin by stating with the precise statement of [Proposition 4.2](#), which relates norms of gradients of the discretized objective to that of the continuous-time one. We begin with the following proposition which bounds the difference between the continuous-time gradient, and a (normalized) embedding of the discrete-time gradient into continuous-time. We define the constant

$$\kappa_{\nabla} := \left( (1 + L_f) M_{\text{cost}} (1 + \kappa_f) + L_{\text{cost}} (3\kappa_f M_f + 8L_f^2 + L_f) \right) = \mathcal{O}_{\star}(\cdot) 1\tag{A.5}$$

**Proposition A.4** (Discretization of the Gradient). *Let  $\pi$  be feasible, and let  $\tilde{\nabla} \mathcal{J}_T(\pi) = \frac{1}{\tau} \tau (\nabla \mathcal{J}_T^{\text{disc}}(\pi))$  is the continuous-time inclusion of the discrete-time gradient, normalized by  $\tau^{-1}$ . Then,*

$$\sup_{t \in [0, T]} \|\nabla \mathcal{J}_T(\pi)(t) - \tilde{\nabla} \mathcal{J}_T(\pi)(t)\| \leq \tau e^{\tau L_f} \max\{\kappa_{\pi, \infty}, \kappa_{\pi, 1}, 1\} L_{\pi} \kappa_{\nabla},$$

The above result is proven in [Appendix I.3](#). By integrating, we see that  $\|\nabla \mathcal{J}_T(\pi) - \tilde{\nabla} \mathcal{J}_T(\pi)\|_{\mathcal{L}_2(\mathcal{U})} \leq \sqrt{T} \tau e^{\tau L_f} \max\{\kappa_{\pi, \infty}, \kappa_{\pi, 1}, 1\} L_{\pi} \kappa_{\nabla}$ , and thus the triangle inequality gives  $\|\nabla \mathcal{J}_T(\pi)\|_{\mathcal{L}_2(\mathcal{U})} \leq \|\tilde{\nabla} \mathcal{J}_T(\pi)\|_{\mathcal{L}_2(\mathcal{U})} + \sqrt{T} \tau e^{\tau L_f} \max\{\kappa_{\pi, \infty}, \kappa_{\pi, 1}, 1\} L_{\pi} \kappa_{\nabla}$ . We can see that for any  $\bar{\mathbf{u}} = \mathbf{u}_{1:K} \in \mathcal{U}$ ,  $\|\bar{\mathbf{u}}\|_{\ell_2}^2 = \sum_{k=1}^K \|\mathbf{u}_k\|^2 = \frac{1}{\tau} \int_0^T \|\mathbf{u}_k(t)\|^2 = \frac{1}{\tau} \|\text{ct}(\bar{\mathbf{u}})\|_{\mathcal{L}_2(\mathcal{U})}^2$ . Hence, in particular,  $\|\nabla \mathcal{J}_T(\pi)\|_{\mathcal{L}_2(\mathcal{U})} \leq \frac{1}{\sqrt{\tau}} \|\nabla \mathcal{J}_T^{\text{disc}}(\pi)\|_{\ell_2} + \sqrt{T} \tau e^{\tau L_f} \max\{\kappa_{\pi, \infty}, \kappa_{\pi, 1}, 1\} L_{\pi} \kappa_{\nabla}$ . From this, and from using [Lemma A.1](#) to bound  $\kappa_{\pi, \infty}, \kappa_{\pi, 1} = \mathcal{O}_{\pi}(1)$ , we obtain the following corollary, which is a precise statement of [Proposition 4.2](#).

**Corollary A.1.** *Suppose  $\pi$  is feasible. Then, recalling  $\kappa_{\nabla}$  from [Eq. \(A.5\)](#),*

$$\|\nabla \mathcal{J}_T(\pi)\|_{\mathcal{L}_2(\mathcal{U})} \leq \frac{1}{\sqrt{\tau}} \|\nabla \mathcal{J}_T^{\text{disc}}(\pi)\|_{\ell_2} + \sqrt{T} \tau e^{\tau L_f} \max\{\kappa_{\pi, \infty}, \kappa_{\pi, 1}, 1\} L_{\pi} \kappa_{\nabla}.$$

In particular, for  $\tau \leq 1/4L_f$ ,

$$\|\nabla \mathcal{J}_T(\pi)\|_{\mathcal{L}_2(\mathcal{U})} \leq \frac{1}{\sqrt{\tau}} \|\nabla \mathcal{J}_T^{\text{disc}}(\pi)\|_{\ell_2} + \sqrt{T}\tau \cdot \underbrace{2 \max\{\kappa_{\pi,\infty}, \kappa_{\pi,1}, 1\} L_{\pi} \kappa_{\nabla}}_{=\mathcal{O}_{\pi}(1)}.$$

### A.6. Main Taylor Expansion Results

We now state various bounds on Taylor-expansion like terms. All the following results are proven in [Appendix D](#). The first is a Taylor expansion of the dynamics (proof in [Appendix D.1](#)).

**Proposition A.5.** *Let  $\pi$  be feasible,  $\tau \leq \tau_{\text{tay},\pi}$ . Fix a  $\mathbf{u}_{1:K} \in \mathbf{U}$ , and define the perturbation  $\delta \mathbf{u}_{1:K} := \mathbf{u}_{1:K} - \mathbf{u}_{1:K}^{\pi}$ , and define*

$$B_2 := \sqrt{\tau} \|\delta \mathbf{u}_{1:K}\|_{\ell_2}, \quad B_{\infty} := \max_k \|\delta \mathbf{u}_k\|.$$

Then, if  $B_{\infty} \leq R_{\text{feas}}/8$ , and if for either  $q \in \{2, \infty\}$ , it holds that  $B_q \leq B_{\text{tay},q,\pi}$ , then

(a) The following bounds hold for all  $k \in [K+1]$

$$\|\tilde{\mathbf{x}}_k^{\pi}(\mathbf{u}_{1:K}) - \mathbf{x}_k^{\pi} - \sum_{j=1}^{k-1} \Psi_{\text{cl},k,j}^{\pi} \delta \mathbf{u}_j\| \leq M_{\text{tay},q,\pi} B_q^2, \quad \|\tilde{\mathbf{x}}_k^{\pi}(\mathbf{u}_{1:K}) - \mathbf{x}_k^{\pi}\| \leq L_{\text{tay},q,\pi} B_q,$$

(b) Moreover, for all  $k \in [K+1]$  and  $t \in [0, T]$ ,

$$\max\{\|\tilde{\mathbf{x}}_k^{\pi}(\mathbf{u}_{1:K})\|, \|\tilde{\mathbf{u}}_k^{\pi}(\mathbf{u}_{1:K})\|\} \leq \frac{3R_{\text{feas}}}{4}, \text{ and } \|\tilde{\mathbf{x}}^{\pi}(t \mid \mathbf{u}_{1:K})\| \leq R_{\text{feas}}.$$

Next, we provide a Taylor expansion of the discrete-time cost functional (proof in [Appendix D.2](#)).

**Lemma A.6.** *Consider the setting of [Proposition A.5](#), and suppose  $B_{\infty} \leq R_{\text{feas}}/8$  and  $B_2 \leq B_{\text{tay},2,\pi}$ . Then,*

$$\|\mathcal{J}_T^{\pi,\text{disc}}(\delta \mathbf{u}_{1:K} + \mathbf{u}_{1:K}^{\pi}) - \mathcal{J}_T^{\pi,\text{disc}}(\mathbf{u}_{1:K}^{\pi}) - \langle \delta \mathbf{u}_k, \nabla \mathcal{J}_T^{\pi,\text{disc}}(\mathbf{u}_{1:K}^{\pi}) \rangle\| \leq M_{\mathcal{J},\text{tay},\pi} B_2^2.$$

Next, we show sufficiently small perturbations of the nominal input preserve stability of the dynamics (proof in [Appendix D.3](#)).

**Lemma A.7.** *Again consider the setting of [Proposition A.5](#), and suppose  $B_{\infty} \leq \min\{R_{\text{feas}}/8, B_{\text{tay},\text{inf},\pi}, B_{\text{stab},\pi}\}$ . Then,*

$$\mu_{\pi',\star} \leq (1 + B_{\infty}/B_{\text{stab},\pi}) \mu_{\pi,\star} \leq 2\mu_{\pi,\star}, \quad L_{\pi'} = L_{\pi}.$$

Lastly, we bound the norm of the discretized gradient ([Appendix D.4](#)).

**Lemma A.8.** *Let  $\pi$  be feasible, and let  $\tau \leq \tau_{\text{dyn}}$ . Then*

$$\max_{k \in [K]} \|(\nabla \mathcal{J}_T^{\text{disc}}(\pi))_k\| \leq \tau L_{\nabla,\pi,\infty}$$

### A.7. Estimation Errors

In this section, we bound the various estimation errors. All the proofs are given in [Appendix E](#). We begin with a simple condition we need for estimation of Markov parameters to go through.

**Definition A.5.** We say  $\pi$  is *estimation-friendly* if  $\pi$  is feasible, and if

$$\sigma_{\text{orac}} \sqrt{\frac{\iota(\delta)}{2NL_{\pi}}} \leq \sigma_w \leq \frac{B_{\text{tay},\text{inf},\pi}}{2\sqrt{d_{\star}}}, \quad \tau \leq \tau_{\text{tay},\pi}$$

Our first result is recovery of the nominal trajectory and Markov operators. Recovery of the nominal trajectory follows from Gaussian concentration, and recovery of the Markov operator for the Matrix Hoeffding inequality (Tropp (2012, Theorem 1.4)) combined with the Taylor expansion of the dynamics due to Proposition A.5. The following is proven in Appendix E.1. To state the bound, we recall the estimation error terms in Definition A.4.

**Proposition A.9.** Fix  $\delta \in (0, 1)$  and suppose that  $N$  is large enough that  $\pi$  is estimation friendly. Then, for any estimation-friendly ESTMARKOV( $\pi; N, \sigma_w$ ) (Algorithm 2) returns estimates with such that, with probability  $1 - \delta/2n_{\text{iter}}$ .

$$\max_{1 \leq j < k \leq K+1} \|\Psi_{\text{cl},k,j}^\pi - \hat{\Psi}_{k,j}\|_{\text{op}} \leq \text{Err}_{\Psi,\pi}(\delta) \quad \max_{k \in [K+1]} \|\hat{\mathbf{x}}_k - \mathbf{x}_k^\pi\| \leq \text{Err}_{\hat{\mathbf{x}}}(\delta) \quad (\text{A.6})$$

Let  $\Pi_{\text{alg}} := \{\pi^{(n)}, \tilde{\pi}^{(n)} : n \in [n_{\text{iter}}]\}$  denote the set of policies constructed by the algorithm, and note that ESTMARKOV is called once for each policy in  $\Pi_{\text{alg}}$ . We define the good estimation event as

$$\mathcal{E}_{\text{est}}(\delta) := \bigcap_{n=1}^{\infty} (\mathcal{E}_n(\delta) \cap \tilde{\mathcal{E}}_n(\delta)), \quad (\text{A.7})$$

$$\mathcal{E}_n(\delta) := \{\text{Eq. (A.6) holds for } \pi = \pi^{(n)} \text{ if } \pi^{(n)} \text{ is estimation friendly}\} \quad (\text{A.8})$$

$$\tilde{\mathcal{E}}_n(\delta) := \{\text{Eq. (A.6) holds for } \tilde{\pi} = \pi^{(n)} \text{ if } \tilde{\pi}^{(n)} \text{ is estimation-friendly}\} \quad (\text{A.9})$$

By Proposition A.9 and a union bound implies

$$\mathbb{P}[\mathcal{E}_{\text{est}}(\delta)] \geq 1 - \delta.$$

We now show that on the good estimation event, the error of the gradient is bounded. The proof is Appendix E.2.

**Lemma A.10** (Gradient Error). *On the event  $\mathcal{E}_{\text{est}}(\delta)$ , it holds that that if  $\pi^{(n)}$  is estimation-friendly, then Algorithm 1 (Line 4) produces*

$$\max_k \|\hat{\nabla}_k^{(n)} - (\nabla \mathcal{J}_T^{\text{disc}}(\pi^{(n)}))_k\| \leq \text{Err}_{\nabla,\pi^{(n)}}(\delta).$$

We also bound the error in the recovery of the system parameters used for synthesizing the stabilizing gains. Recovery of said parameters requires first establishing controllability of the discrete-time Markov operator. We prove the following in Appendix E.3:

**Proposition A.11.** Define  $\gamma_{\text{ctr}} := \max\{1, L_f t_{\text{ctrl}}\}$ , and suppose that  $\tau \leq \min\{\tau_{\text{ctrl},\pi}, \tau_{\text{dyn}}\}$ . Then, for  $k \geq k_{\text{ctrl}} + 1$ , it holds that

$$\lambda_{\min} \left( \sum_{j=k-k_{\text{ctrl}}}^{k-1} \Psi_{\text{cl},k,j}^\pi (\Psi_{\text{cl},k,j}^\pi)^\top \right) \geq \tau \cdot \frac{\nu_{\text{ctrl}}}{8L_\pi^2 \gamma_{\text{ctr}}^2 \exp(2\gamma_{\text{ctr}})}$$

With this result, Appendix E.4 upper bounds the estimation error for the discrete-time system matrices.

**Proposition A.12.** Suppose  $\mathcal{E}_{\text{est}}(\delta)$  holds, fix  $n \in n_{\text{iter}}$ , and let  $\pi = \tilde{\pi}^{(n)}$ . Then, suppose that  $\tau \leq \min\{\tau_{\text{ctrl},\pi}, \tau_{\text{dyn}}\}$ ,  $k_0 \geq k_{\text{ctrl}} + 2$ , and

$$\text{Err}_\Psi(\delta) \leq \tau \frac{\sqrt{\nu_{\text{ctrl}}/t_{\text{ctrl}}}}{2\sqrt{2}L_\pi \gamma_{\text{ctr}} \exp(\gamma_{\text{ctr}})}, \quad (\text{A.10})$$

Then, on  $\mathcal{E}_{\text{est}}(\delta)$ , if  $\pi$  is estimation-friendly, the estimates from the call of ESTGAINS( $\pi; N, \sigma$ ) satisfy

$$\|\hat{\mathbf{B}}_k - \mathbf{B}_{\text{ol},k}^\pi\| \vee \|\hat{\mathbf{A}}_k - \mathbf{A}_{\text{ol},k}^\pi\| \leq \frac{\text{Err}_{\Psi,\pi}(\delta)}{\tau} \cdot t_0 \kappa_{\pi,\infty} L_\pi^2 \frac{192\gamma_{\text{ctr}}^3 \exp(2\gamma_{\text{ctr}})}{\nu_{\text{ctrl}}}.$$

## A.8. Descent and Stabilization

In this section, we leverage the estimation results in the previous section to demonstrate the two key features of the algorithm: descent on the discrete-time objective, and stability after the synthesized gains. We begin with a standard first-order descent lemma, whose proof is given in Appendix H.1. This lemma also ensures, by invoking Lemma A.1, that the step size is sufficiently small to control the stability of  $\tilde{\pi}^{(n)}$ , which uses the same gains as  $\pi^{(n)}$  but has a slightly perturbed control input.

**Lemma A.13** (Descent Lemma). *Suppose  $\pi = \pi^{(n)}$  is estimation friendly, let  $M \geq M_{\mathcal{J}, \text{tay}, \pi}$ , and suppose*

$$\eta \leq \frac{1}{4M}, \quad (\eta(L_{\nabla, \pi, \infty} + \frac{1}{\tau} \text{Err}_{\nabla, \pi^{(n)}}(\delta)) + \text{Err}_{\hat{x}}(\delta)) \leq \min \left\{ \frac{R_{\text{feas}}}{8}, B_{\text{stab}, \pi}, B_{\text{tay}, \text{inf}, \pi}, \frac{B_{\text{tay}, 2, \pi}}{\sqrt{T}} \right\}.$$

*Then, on event  $\mathcal{E}_{\text{est}}(\delta)$ , it holds (again setting  $\pi \leftarrow \pi^{(n)}$  on the right-hand side)*

$$\mathcal{J}_T^{\text{disc}}(\tilde{\pi}^{(n)}) - \mathcal{J}_T^{\text{disc}}(\pi^{(n)}) \leq -\frac{\eta}{2\tau} \|\hat{\nabla}_{1:K}^{(n)}\|_{\ell_2}^2 + T \left( \frac{\text{Err}_{\nabla, \pi}(\delta)^2}{4\tau^2 M} + \text{Err}_{\hat{x}}(\delta) L_{\nabla, \pi, \infty} + M \text{Err}_{\hat{x}}(\delta)^2 \right).$$

and that

$$L_{\tilde{\pi}^{(n)}} = L_{\pi^{(n)}}, \quad \mu_{\tilde{\pi}^{(n)}, \star} \leq 2\mu_{\pi^{(n)}, \star}.$$

The next step is to establish a stability guarantee for the certainty-equivalent gains synthesized. We begin with a generic guarantee, whose proof is given in [Appendix G](#).

**Proposition A.14** (Certainty Equivalence Bound). *Let  $\hat{\mathbf{A}}_k^\pi$  and  $\hat{\mathbf{B}}_k^\pi$  be estimates of  $\mathbf{A}_{\text{ol}, k}^\pi$  and  $\mathbf{B}_{\text{ol}, k}^\pi$ , and let  $\hat{\mathbf{K}}_k$  denote the corresponding certainty equivalence controller synthesized by [Algorithm 3](#) (Lines 7 and 10). Suppose that  $\tau \leq \min\{\tau_{\text{ric}}, \tau_{\text{dyn}}\}$  and*

$$\max_{k \in [k_0:K]} \|\hat{\mathbf{A}}_k^\pi - \mathbf{A}_{\text{ol}, k}^\pi\|_{\text{op}} \vee \|\hat{\mathbf{B}}_k^\pi - \mathbf{B}_{\text{ol}, k}^\pi\|_{\text{op}} \leq \tau(2^{17} \mu_{\text{ric}}^4 \max\{1, L_f^3\})^{-1}$$

*Then, if  $\pi' = (\mathbf{u}_{1:K}^\pi, \hat{\mathbf{K}}_{1:K})$ , we have*

$$\mu_{\pi', \star} \leq 4\mu_{\text{ric}}, \quad L_{\pi'} \leq 6 \max\{1, L_f\} \mu_{\text{ric}}.$$

As a direct corollary of the above proposition and [Proposition A.12](#), we obtain the following:

**Lemma A.15.** *Suppose  $\mathcal{E}_{\text{est}}(\delta)$  holds, fix  $n \in n_{\text{iter}}$ , and let  $\pi = \tilde{\pi}^{(n)}$ . Then, suppose that  $\tau \leq \min\{\tau_{\text{ctrl}, \pi}, \tau_{\text{dyn}}\}$ ,  $\pi$  is estimation-friendly,  $k_0 \geq k_{\text{ctrl}} + 2$ , and*

$$\text{Err}_{\Psi, \pi}(\delta) \leq \tau^2 \left( 2^{25} \mu_{\text{ric}}^4 \max\{1, L_f^3\} \cdot t_0 \kappa_{\pi, \infty} L_\pi^2 \frac{\gamma_{\text{ctr}}^3 \exp(2\gamma_{\text{ctr}})}{\nu_{\text{ctrl}}} \right)^{-1} \quad (\text{A.11})$$

Then,

$$\mu_{\pi^{(n+1)}, \star} \leq 4\mu_{\text{ric}}, \quad L_{\pi^{(n+1)}} \leq 6 \max\{1, L_f\} \mu_{\text{ric}}.$$

*Proof.* One can check that, as  $\mu_{\text{ric}}, L_\pi \geq 1$ , [Eq. \(A.11\)](#) implies [Eq. \(A.10\)](#). Thus, the lemma follows directly from [Propositions A.12](#) and [A.14](#), as well as noting  $192 \cdot 2^{17} \leq 2^{25}$   $\square$

## A.9. Concluding the proof.

In this section, we conclude the proof. First, we define uniform upper bounds on all  $\pi$ -dependent parameters.

**Uniform upper bounds on parameters.** To begin, define

$$\bar{\mu} = 8\mu_{\text{ric}}, \quad \bar{L} = 6 \max\{L_f, 1\} \mu_{\text{ric}}. \quad (\text{A.12})$$

Next, for  $q \in \{1, 2, \infty\}$ , define  $\bar{\kappa}_q := \kappa_q(\bar{\mu}, \bar{L})$  defined in [Definition A.1](#). We define  $\bar{\tau}_{\text{tay}}, \bar{\tau}_{\text{ctrl}}$  analogously to  $\tau_{\text{tay}, \pi}, \tau_{\text{ctrl}, \pi}$  in [Definition A.2](#) with  $\kappa_{\pi, \infty}$  replaced by  $\bar{\kappa}_\infty$  and  $L_\pi$  with  $\bar{L}$ . For  $q \in \{2, \infty\}$ , we define  $\bar{M}_{\text{tay}, q}, \bar{M}_{\mathcal{J}, \text{tay}}, \bar{L}_{\text{tay}, q}, \bar{L}_{\nabla, \infty}, \bar{B}_{\text{tay}, q}, \bar{B}_{\text{stab}}$  analogously to  $M_{\text{tay}, q, \pi}, M_{\mathcal{J}, \text{tay}, \pi}, L_{\text{tay}, q, \pi}, L_{\nabla, \pi, \infty}, B_{\text{tay}, q, \pi}, B_{\text{stab}, \pi}$  in [Definition A.3](#), with all occurrences of  $\kappa_{\pi, \infty}, \kappa_{\pi, 1}, \kappa_{\pi, 2}$  replaced by  $\bar{\kappa}_\infty, \bar{\kappa}_1, \bar{\kappa}_2$  and all occurrences of  $L_\pi$  replaced by  $\bar{L}$ . Finally,

we define  $\overline{\text{Err}}_\Psi, \overline{\text{Err}}_\nabla$  to be analogous to  $\text{Err}_{\Psi,\pi}, \text{Err}_{\nabla,\pi}$  but with the same above substitutions. From [Lemmas A.1](#) and [A.2](#), we have

$$\begin{aligned}\bar{\kappa}_q, \bar{M}_{\text{tay},q}, \bar{L}_{\text{tay},q}, \bar{L}_{\nabla,\infty}, \bar{B}_{\text{tay},q} &= \mathcal{O}_*(1) \\ \tau_{\text{dyn}}, \tau_{\text{ric}}, \bar{\tau}_{\text{tay}}, \bar{\tau}_{\text{ctrl}} &= 1/\mathcal{O}_*(1) \\ \bar{M}_{\mathcal{J},\text{tay}} &= T \cdot \mathcal{O}_*(1). \\ \bar{B}_{\text{stab}} &= \frac{1}{T} \mathcal{O}_*(1)\end{aligned}$$

Moreover, recalling  $\text{Err}_0(\delta) := \sqrt{\iota(\delta)/N}$ , and setting  $\sigma_w = c\sqrt{\text{Err}_0\sigma_{\text{orac}}}$  for any  $c \in [1/\mathcal{O}_*(1), \mathcal{O}_*(1)]$ , [Lemma A.3](#) gives

$$\begin{aligned}\text{Err}_{\hat{x}}(\delta) &= \mathcal{O}_*(\sigma_{\text{orac}}\text{Err}_0(\delta)) \\ \text{Err}_{\Psi,\pi}(\delta) &= \mathcal{O}_*(\text{Err}_0(\delta) + \sqrt{\sigma_{\text{orac}}\text{Err}_0(\delta)}) \\ \text{Err}_{\nabla,\pi}(\delta) &= \mathcal{O}_*(T(\text{Err}_0(\delta) + \sqrt{\sigma_{\text{orac}}\text{Err}_0(\delta)})).\end{aligned}\tag{A.13}$$

**Statement of Main Guarantee, Explicit Constants.** We begin by stating our main guarantee, first with explicit constants. We then translate into a  $\mathcal{O}_*(1)$  notation. To begin, define the following descent error term:

$$\overline{\text{Err}}_{\text{dec}}(\delta) := T \left( \frac{\overline{\text{Err}}_\nabla(\delta)^2}{4\tau^2 \bar{M}_{\mathcal{J},\text{tay}}} + \text{Err}_{\hat{x}}(\delta) \bar{L}_{\nabla,\infty} + \bar{M}_{\mathcal{J},\text{tay}} \text{Err}_{\hat{x}}(\delta)^2 \right)\tag{A.14}$$

And note that for  $\sigma_w = c\sqrt{\text{Err}_0\sigma_{\text{orac}}}$  for  $c \in [1/\mathcal{O}_*(1), \mathcal{O}_*(1)]$  (using numerous simplifications, such as  $T/\tau \geq 1$ )

$$\overline{\text{Err}}_{\text{dec}}(\delta) := \mathcal{O}_*(1) \cdot \left( \frac{T^3}{\tau^2} (\text{Err}_0(\delta)^2 + \sigma_{\text{orac}}\text{Err}_0(\delta)) + T\sigma_{\text{orac}}^2 \text{Err}_0(\delta)^2 \right).$$

**Theorem 3.** Fix  $\delta \in (0, 1)$ , and suppose that  $\eta \leq \frac{1}{4M_{\mathcal{J},\text{tay}}}$ ,  $k_0 \geq k_{\text{ctrl}} + 2$ , and suppose

$$\sigma_{\text{orac}} \sqrt{\frac{\iota(\delta)}{2NL}} \leq \sigma_w \leq \frac{\bar{B}_{\text{tay},\infty}}{2\sqrt{d_*}},\tag{A.15a}$$

$$(\eta(\bar{L}_{\nabla,\infty} + \frac{1}{\tau}\overline{\text{Err}}_\nabla(\delta)) + \text{Err}_{\hat{x}}(\delta)) \leq \min \left\{ \frac{R_{\text{feas}}}{8}, \bar{B}_{\text{stab}}, \bar{B}_{\text{tay},\infty}, \frac{\bar{B}_{\text{tay},2}}{\sqrt{T}} \right\}\tag{A.15b}$$

$$\overline{\text{Err}}_\Psi(\delta) \leq \tau^2 \left( 2^{25} \mu_{\text{ric}}^4 \max\{1, L_f^3\} \cdot t_{\text{ctrl}} \kappa_{\pi,\infty} \bar{L}^2 \frac{\gamma_{\text{ctr}}^3 \exp(2\gamma_{\text{ctr}})}{\nu_{\text{ctrl}}} \right)^{-1} = \frac{\tau^2}{\mathcal{O}_*(1)}\tag{A.15c}$$

$$\tau \leq \min\{\bar{\tau}_{\text{tay}}, \bar{\tau}_{\text{ctrl}}, \tau_{\text{ric}}\} = \frac{1}{\mathcal{O}_*(1)}\tag{A.15d}$$

Then, for  $\pi = \pi^{(n_{\text{out}})}$  returned by [Algorithm 1](#) satisfies all four properties with probability  $1 - \delta$ :

(a)  $\mu_{\pi,*} \leq 8\mu_{\text{ric}}$  and  $L_\pi \leq 6 \max\{1, \kappa_f\} \mu_{\text{ric}} = \bar{L}$ . In fact, for all  $n \in [n_{\text{iter}}]$ , and  $\pi' \in \{\pi^{(n)}, \tilde{\pi}^{(n)}\}$ ,  $\mu_{\pi',*} \leq \bar{\mu} = 8\mu_{\text{ric}}$  and  $L_{\pi'} \leq \bar{L} = 6 \max\{1, L_f\} \mu_{\text{ric}}$ .

(b) The discrete-time stabilized gradient is bounded by

$$\begin{aligned}\tau \|\mathcal{J}_T^{\text{disc}}(\pi)\|_{\ell_2}^2 &\leq 2T\overline{\text{Err}}_\nabla(\delta)^2 + \frac{2}{\eta} \left( \frac{2(1+T)\kappa_{\text{cost}}}{n_{\text{iter}}} + \overline{\text{Err}}_{\text{dec}}(\delta) \right) \\ &= \frac{1}{\eta} \mathcal{O}_*(1) \cdot \left( \frac{T}{n_{\text{iter}}} + \left( \eta T^3 + \frac{T^3}{\tau^2} \right) (\text{Err}_0(\delta)^2 + \sigma_{\text{orac}}\text{Err}_0(\delta)) + T\sigma_{\text{orac}}^2 \text{Err}_0(\delta)^2 \right),\end{aligned}$$

where the last line holds when  $\sigma_w = c\sqrt{\sigma_{\text{orac}}\text{Err}_0(\delta)}$  for some  $c \in [\frac{1}{\mathcal{O}_*(1)}, \mathcal{O}_*(1)]$ .

(c) Recall  $\kappa_{\nabla} := \left( (1 + L_f)M_{\text{cost}}(1 + \kappa_f) + L_{\text{cost}}(3\kappa_f M_f + 8L_f^2 + L_f) \right) = \mathcal{O}_*(1)$  from Eq. (A.5). Then  $\pi$  is  $\epsilon$ -stationary for

$$\begin{aligned} \epsilon^2 &= 4T\overline{\text{Err}}_{\nabla}(\delta)^2 + \frac{4}{\eta} \left( \frac{2(1+T)\kappa_{\text{cost}}}{n_{\text{iter}}} + \overline{\text{Err}}_{\text{dec}}(\delta) \right) + 4T\tau^2 (\max\{\bar{\kappa}_{\infty}, \bar{1}_{\infty}, 1\} \bar{L}_{\kappa_{\nabla}})^2 \\ &= \mathcal{O}_*(1) \cdot T \left( \tau^2 + \frac{1}{\eta} \left( \frac{1}{n_{\text{iter}}} + (\eta T^2 + \frac{T^2}{\tau^2}) (\text{Err}_0(\delta)^2 + \sigma_{\text{orac}} \text{Err}_0(\delta)) + \sigma_{\text{orac}}^2 \text{Err}_0(\delta)^2 \right) \right). \end{aligned}$$

where the last line holds when  $\sigma_w = c\sqrt{\sigma_{\text{orac}} \text{Err}_0(\delta)}$  for some  $c \in [\frac{1}{\mathcal{O}_*(1)}, \mathcal{O}_*(1)]$ .

(d)  $\mathbf{u}^{\pi}$  is an  $\epsilon'$ -JSP, where  $\epsilon' = 64\epsilon^2 \bar{L}^2 / \alpha = \mathcal{O}_*(1) \cdot \frac{\epsilon^2}{\alpha}$ .

We prove Theorem 3 from the above results in Appendix A.9.2 just below. Appendix A.9.1 below translates the above theorem into Theorem 2 which uses  $\mathcal{O}_*(\cdot)$  notation.

#### A.9.1. TRANSLATING THEOREM 3 INTO THEOREM 2

*Proof.* It suffices to translate the conditions Eqs. (A.15a) to (A.15d) into  $\mathcal{O}_*(\cdot)$  notation. Again, recall  $\text{Err}_0(\delta) = \sqrt{\iota(\delta)/N}$ , and take  $\sigma_w = c\sqrt{\text{Err}_0(\delta)\sigma_{\text{orac}}}$  for  $c \in [1/\mathcal{O}_*(1), \mathcal{O}_*(1)]$ . Then, Eq. (A.15a) holds for  $\text{Err}_0(\delta) \leq 1/c_1$ , where  $c_1 = \mathcal{O}_*(1)$ . Next, to make Eq. (A.15b) hold, it suffices that

$$\max \left\{ (\eta \bar{L}_{\nabla, \infty}), \frac{\eta \overline{\text{Err}}_{\nabla}(\delta)}{\tau}, \text{Err}_{\hat{x}}(\delta) \right\} \leq \frac{1}{3} \min \left\{ \frac{R_{\text{feas}}}{8}, \bar{B}_{\text{stab}}, \bar{B}_{\text{tay}, \infty}, \frac{\bar{B}_{\text{tay}, 2}}{\sqrt{T}} \right\},$$

The term  $\eta \bar{L}_{\nabla, \infty}$  is sufficiently bounded where  $\eta \leq \frac{1}{c_2 \sqrt{T}}$  for  $c_2 = \mathcal{O}_*(1)$ . Recalling  $\text{Err}_{\nabla, \pi}(\delta) \leq \mathcal{O}_*(T(\text{Err}_0(\delta) + \sqrt{\sigma_{\text{orac}} \text{Err}_0(\delta)}))$  from Eq. (A.13), and that  $\eta \leq \frac{1}{c_2 \sqrt{T}}$ , it is enough that  $(\text{Err}_0(\delta) + \sqrt{\sigma_{\text{orac}} \text{Err}_0(\delta)}) \leq \frac{c_2 \tau}{c_3 T}$  for  $c_3 = \mathcal{O}_*(1)$ . Finally  $\text{Err}_{\hat{x}}(\delta)$  is bounded for  $\text{Err}_{\hat{x}}(\delta) = \text{Err}_0(\delta) \leq 1/c_4 \sqrt{T}$ , where  $c_4 = \mathcal{O}_*(1)$ . Collecting these conditions, we have that for  $c_1, c_2, c_3, c_4 = \mathcal{O}_*(1)$ , Eqs. (A.15a) and (A.15b) hold for

$$\eta \leq \frac{1}{c_2 \sqrt{T}}, \quad \text{Err}_0(\delta) + \sqrt{\sigma_{\text{orac}} \text{Err}_0(\delta)} \leq \frac{c_2 \tau}{c_3 T},$$

Next, as  $\text{Err}_{\Psi, \pi}(\delta) = \mathcal{O}_*(\text{Err}_0(\delta) + \sqrt{\sigma_{\text{orac}} \text{Err}_0(\delta)})$  from Eq. (A.13), Eq. (A.15c) holds as long as  $\text{Err}_0(\delta) + \sqrt{\sigma_{\text{orac}} \text{Err}_0(\delta)} \leq \tau^2/c_5$  for a  $c_5 = \mathcal{O}_*(1)$ . Combining,

$$\eta \leq \frac{1}{c_2 \sqrt{T}}, \quad \text{Err}_0(\delta) + \sqrt{\sigma_{\text{orac}} \text{Err}_0(\delta)} \leq \min \left\{ \frac{c_2 \tau}{c_3 T}, \frac{\tau^2}{c_5} \right\} \quad \text{Err}_0(\delta) \leq \min \left\{ \frac{1}{c_1}, \frac{1}{c_4 \sqrt{T}} \right\}$$

Finally, Eq. (A.15d) requires  $\tau \leq 1/c_6$ , for  $c_6 = \mathcal{O}_*(1)$ , and that  $\eta \leq 1/c_7$  where  $c_7 = 4\bar{M}_{\mathcal{J}, \text{tay}} = \mathcal{O}_{\pi}(1)$ . By shrinking constants if necessary, this can be simplified into

$$\eta \leq \min \left\{ \frac{1}{c_7}, \frac{1}{c_2 \sqrt{T}} \right\}, \quad \text{Err}_0(\delta) \leq \min \left\{ \min \left\{ \frac{c_2 \tau}{c_3 T}, \frac{\tau^2}{c_5}, \frac{1}{c_1}, \frac{1}{c_4 \sqrt{T}} \right\}, \frac{1}{\sigma_{\text{orac}}} \min \left\{ \frac{c_2 \tau}{c_3 T}, \frac{\tau^2}{c_5} \right\}^2 \right\}.$$

And recall  $\text{Err}_0(\delta) = \sqrt{\iota(\delta)/N}$ , this becomes

$$\tau \leq \frac{1}{c_6}, \quad \eta \leq \frac{1}{c_2} \min \left\{ 1, \frac{1}{\sqrt{T}} \right\}, \quad N \geq \iota(\delta) \min \left\{ \min \left\{ \frac{1}{c_1}, \frac{c_2 \tau}{c_3 T}, \frac{\tau^2}{c_5}, \frac{1}{c_4 \sqrt{T}} \right\}, \frac{1}{\sigma_{\text{orac}}} \min \left\{ \frac{c_2 \tau}{c_3 T}, \frac{\tau^2}{c_5} \right\}^2 \right\}^{-2}.$$

By consolidating constants and relabeling  $c_1, c_2 = \mathcal{O}_*(1)$  as needed, it suffices that

$$\begin{aligned} \eta &\leq c_1 \min \left\{ \frac{1}{\sqrt{T}}, 1 \right\}, \quad \tau \leq \frac{1}{c_2}, \quad N \geq c_3 \iota(\delta) \min \left\{ \min \left\{ 1, \frac{\tau}{T}, \tau^2, \frac{1}{\sqrt{T}} \right\}, \frac{1}{\sigma_{\text{orac}}} \min \left\{ \frac{\tau}{T}, \tau^2 \right\}^2 \right\}^{-2} \\ &= c_3 \iota(\delta) \max \left\{ 1, \frac{T^2}{\tau^2}, \frac{1}{\tau^4}, T, \sigma_{\text{orac}}^2 \frac{T^4}{\tau^2}, \frac{\sigma_{\text{orac}}^2}{\tau^8} \right\}. \end{aligned}$$

Having shown that the above conditions suffice to ensure Theorem 3 holds, the bound follows (again replacing  $\text{Err}_0(\delta)$  with  $\sqrt{\iota(\delta)/N}$ ).

□

## A.9.2. PROOF OF THEOREM 3

We shall show the following invariant. At each step  $n$ ,

$$\mu_{\pi^{(n)}, \star} \leq \bar{\mu}/2, \quad L_{\pi^{(n)}} \leq \bar{L}. \quad (\text{A.16})$$

**Lemma G.3** shows that Eq. (A.16) holds for  $n = 1$ . Next, for  $n \geq 1$ , directly combining **Lemmas A.13** and **A.15** imply the following per-round guarantee.

**Lemma A.16** (Per-Round Lemma). *Suppose that  $\eta \leq \frac{1}{4M_{\mathcal{J}, \text{tay}}}$ ,  $k_0 \geq k_{\text{ctrl}} + 2$ , Then if  $\pi^{(n)}$  satisfies Eq. (A.16) and Eqs. (A.15a) to (A.15d). Then, on  $\mathcal{E}_{\text{est}}(\delta)$ ,*

$$(a) \max_k \|\hat{\nabla}_k^{(n)} - (\nabla \mathcal{J}_T^{\text{disc}}(\pi^{(n)}))_k\| \leq \overline{\text{Err}}_{\nabla}(\delta); \text{ thus } \tau \|\hat{\nabla}_{1:K}^{(n)} - (\nabla \mathcal{J}_T^{\text{disc}}(\pi^{(n)}))\|_{\ell_2}^2 \leq T \overline{\text{Err}}_{\nabla}(\delta)^2.$$

(b) The following descent guarantee holds

$$\mathcal{J}_T^{\text{disc}}(\tilde{\pi}^{(n)}) - \mathcal{J}_T^{\text{disc}}(\pi^{(n)}) \leq -\frac{\eta}{2\tau} \|\hat{\nabla}_{1:K}^{(n)}\|_{\ell_2}^2 + \overline{\text{Err}}_{\text{dec}}(\delta)$$

$$(c) L_{\tilde{\pi}^{(n)}} = L_{\pi^{(n)}} \leq \bar{L} \text{ and } \mu_{\tilde{\pi}^{(n)}, \star} \leq 2\mu_{\pi^{(n)}, \star} \leq \bar{\mu}.$$

$$(d) \mu_{\pi^{(n+1)}, \star} \leq 4\mu_{\text{ric}} = \bar{L}/2, \quad L_{\pi^{(n+1)}} \leq 6 \max\{1, L_f\} \mu_{\text{ric}} = \bar{L}; \text{ that is } \pi^{(n+1)} \text{ satisfies Eq. (D.4)}$$

*Proof.* Part (a) follows from **Lemma A.10**, and parts (b) and (c) follow from **Lemma A.13**, with the necessary replacement of  $\pi$ -dependent terms with  $(\cdot)$  terms. Part (b) allows us to make the same substitutions in **Lemma A.15**, which gives part (c).  $\square$

*Proof of Theorem 3.* Under the conditions of this lemma, **Lemma A.16** holds. As **Lemma G.3** shows that Eq. (A.16) holds for  $n = 1$ , induction implies **Lemma G.3** holds for all  $n \in [n_{\text{iter}}]$  on  $\mathcal{E}_{\text{est}}(\delta)$ , an event which occurs with probability  $1 - \delta$ . We now prove each part of the present theorem in sequence.

**Part (a).** Directly from **Lemma A.16**(d)

**Part (b).** Notice that, since  $\tilde{\pi}^{(n)}$  and  $\pi^{(n+1)}$  differ only in their gains,  $\mathcal{J}_T^{\text{disc}}(\pi^{(n+2)}) = \mathcal{J}_T^{\text{disc}}(\tilde{\pi}^{(n)})$ . Therefore, summing up the descent guarantee in **Lemma A.16**(b), we have

$$\begin{aligned} \mathcal{J}_T^{\text{disc}}(\pi^{(n_{\text{iter}}+1)}) - \mathcal{J}_T^{\text{disc}}(\pi^{(1)}) &\leq -\frac{\eta}{2\tau} \sum_{n=1}^{n_{\text{iter}}} \|\hat{\nabla}_{1:K}^{(n)}\|_{\ell_2}^2 + n_{\text{iter}} \overline{\text{Err}}_{\text{dec}}(\delta) \\ &\leq -\frac{\eta}{2\tau} n_{\text{iter}} \min_{n \in [n_{\text{iter}}]} \|\hat{\nabla}_{1:K}^{(n)}\|_{\ell_2}^2 + n_{\text{iter}} \overline{\text{Err}}_{\text{dec}}(\delta) \\ &= -\frac{\eta}{2\tau} \|\hat{\nabla}_{1:K}^{(n_{\text{out}})}\|_{\ell_2}^2 + n_{\text{iter}} \overline{\text{Err}}_{\text{dec}}(\delta) \end{aligned}$$

where we recall that our algorithm selects to output  $\pi^{(n_{\text{out}})}$ , where  $n_{\text{out}}$  minimizes  $\|\hat{\nabla}_{1:K}^{(n)}\|_{\ell_2}^2$ . Recall that  $\mathcal{J}_T^{\pi, \text{disc}}(\pi) = \mathcal{J}_T^{\pi, \text{disc}}(\mathbf{u}_{1:K}^{\pi})$  where  $\mathcal{J}_T^{\pi, \text{disc}}(\bar{\mathbf{u}}) := V(\tilde{\mathbf{x}}_{K+1}^{\pi}(\bar{\mathbf{u}})) + \tau \sum_{k=1}^K Q(\tilde{\mathbf{x}}_k^{\pi}(\bar{\mathbf{u}}), \tilde{\mathbf{u}}_k^{\pi}(\bar{\mathbf{u}}), t_k)$ . Hence, for all feasible  $\pi$ , **Assumption 4.2** implies  $0 \leq \mathcal{J}_T^{\pi, \text{disc}}(\pi) \leq \kappa_{\text{cost}}(1 + \tau K) = (1 + T)\kappa_{\text{cost}}$ . By **Condition 4.1**,  $\pi^{(n+1)}$  and  $\pi^{(1)}$  are by feasible, and thus  $\mathcal{J}_T^{\text{disc}}(\pi^{(n+1)}) - \mathcal{J}_T^{\text{disc}}(\pi^{(1)}) \geq -(1 + T)\kappa_{\text{cost}}$ . Therefore, by rearranging the previous display,

$$\tau \|\hat{\nabla}_{1:K}^{(n_{\text{out}})}\|_{\ell_2}^2 \leq \frac{1}{\eta} \left( \frac{2(1 + T)\kappa_{\text{cost}}}{n_{\text{iter}}} + \overline{\text{Err}}_{\text{dec}}(\delta) \right).$$

By **Lemma A.16**(a), and AM-GM imply then

$$\tau \|\mathcal{J}_T^{\text{disc}}(\pi^{(n_{\text{out}})})\|_{\ell_2}^2 \leq 2T \overline{\text{Err}}_{\nabla}(\delta)^2 + \frac{2}{\eta} \left( \frac{2(1 + T)\kappa_{\text{cost}}}{n_{\text{iter}}} + \overline{\text{Err}}_{\text{dec}}(\delta) \right).$$



**Part (c).** Note that  $\tau \leq \bar{\tau}_{\text{tay}}$  implies  $\tau \leq 1/4L_f$ . From [Corollary A.1](#), and for  $\kappa_{\nabla} = \mathcal{O}_*(1)$  as in [Eq. \(A.5\)](#), the following holds for any feasible  $\pi$ :

$$\begin{aligned} \|\nabla \mathcal{J}_T(\pi)\|_{\mathcal{L}_2(\mathcal{U})}^2 &\leq \left( \frac{1}{\sqrt{\tau}} \|\nabla \mathcal{J}_T^{\text{disc}}(\pi)\|_{\ell_2} + \sqrt{T}\tau \cdot 2 \max\{\kappa_{\pi, \infty}, \kappa_{\pi, 1}, 1\} L_{\pi} \kappa_{\nabla} \right)^2 \\ &\leq \frac{2}{\tau} \|\nabla \mathcal{J}_T^{\text{disc}}(\pi)\|_{\ell_2}^2 + 4T\tau^2 (\max\{\kappa_{\pi, \infty}, \kappa_{\pi, 1}, 1\} L_{\pi} \kappa_{\nabla})^2. \end{aligned}$$

Apply the above with  $\pi = \pi^{(n_{\text{out}})}$  gives part *c*, and upper bound  $\kappa_{\pi, \infty}, \kappa_{\pi, 1}, L_{\pi}$  by  $\bar{\kappa}_{\infty}, \bar{\kappa}_1, \bar{L}$  concludes.

**Part (d).** This follows directly from [Proposition 4.1](#), noting that  $L_{\pi} \leq \bar{L}$  for  $\pi = \pi^{(n_{\text{out}})}$ , and that  $\bar{\tau}_{\text{tay}} = \frac{1}{16\bar{L}L_f}$ , so that the step-size condition of [Proposition 4.1](#) is met.  $\square$

## B. Discussion and Extensions

### B.1. Separation between and Open-Loop and Closed-Loop Gradients

In this section, we provided an illustrative example as to why a approximation JSP is more natural than canonical stationary points. Fix an  $\epsilon \in (0, 1]$ , and consider the system with dynamic map

$$f_{\epsilon}(x, u) = 2x + u - \epsilon.$$

Let  $\mathbf{x}_{\epsilon}(t | \mathbf{u})$  denote the scalar trajectory with

$$\frac{d}{dt} \mathbf{x}_{\epsilon}(t | \mathbf{u}) = f_{\epsilon}(\mathbf{x}_{\epsilon}(t | \mathbf{u}), \mathbf{u}(t)), \quad \mathbf{x}_{\epsilon}(0 | \mathbf{u}) = \epsilon.$$

Then,  $\mathbf{x}_{\epsilon}(t | \mathbf{0}) = \epsilon$  for all  $t$ . We can now consider the following planning objective

$$\mathcal{J}_{T, \epsilon}(\mathbf{u}) = \frac{1}{2} \int_0^T (\mathbf{x}_{\epsilon}(t | \mathbf{u})^2 + \mathbf{u}(t)^2) dt. \quad (\text{B.1})$$

Since the dynamics  $f_{\epsilon}$  are affine, we find that

$$\mathbf{u} \text{ is an } \epsilon' \text{-JSP of } \mathcal{J}_{T, \epsilon} \iff \mathbf{u} \leq \inf_{\mathbf{u}'} \mathcal{J}_{T, \epsilon}(\mathbf{u}') + \epsilon'.$$

In particular, as  $\mathcal{J}_{T, \epsilon}(\mathbf{0}) = \frac{1}{2}T\epsilon^2$ , and as  $\mathcal{J}_{T, \epsilon} \geq 0$ ,

$$\mathbf{u} = \mathbf{0} \text{ is an } \frac{T\epsilon^2}{2} \text{-JSP of } \mathcal{J}_{T, \epsilon}. \quad (\text{B.2})$$

However, we show that the magnitude of the gradient at  $\mathbf{u} = \mathbf{0}$  is much larger. We compute the following shortly below.

**Lemma B.1.** *For  $T \geq 1$ , we have  $\|\nabla \mathcal{J}_{T, \epsilon}(\mathbf{0})(t)\|_{\mathcal{L}_2(\mathcal{U})} \geq \sqrt{T}\epsilon e^T / 4\sqrt{2}$ .*

Thus, the magnitude of the gradient (through open-loop dynamics) is exponentially larger than the suboptimality of the cost. This suggests that gradients through open-loop dynamics are poor proxy for global optimality, motivating instead the JSP. Moreover, one can easily compute that if  $\pi$  has inputs  $\mathbf{u}_k^{\pi} = 0$  and stabilizing gains  $\mathbf{K}_k = -3$ , then for sufficiently small step sizes, the gradients of  $\mathbf{J}_T^{\pi}(\mathbf{u})|_{\mathbf{u}=\mathbf{0}}$  scale only as  $c\epsilon\sqrt{T}$  for a universal  $c > 0$ , and do not depend exponentially on the horizon.

*Proof of Lemma B.1.* We have from [Lemma C.5](#) that

$$\nabla \mathcal{J}_{T, \epsilon}(\mathbf{0})(t) = \int_{s=t}^T \underbrace{\mathbf{x}_{\epsilon}(s | \mathbf{0})}_{=\epsilon} \cdot \Phi(s, t) \mathbf{B}(t),$$

where  $\Phi(s, t)$  solves the ODE  $\Phi(t, t) = 1$  and  $\frac{d}{ds}\Phi(s, t) = 2\Phi(s, t)$ . Thus,  $\Phi(s, t) = \exp(2(t - s))$ . Moreover,  $\mathbf{B}(t) = 1$ . Hence,

$$\begin{aligned}\nabla\mathcal{J}_{T,\epsilon}(\mathbf{0})(t) &= \epsilon \int_{s=t}^T \exp(2(t-s)) \\ &= \epsilon \frac{1}{2} (e^{2(T-t)} - 1)\end{aligned}$$

Hence, for  $t \leq T/2$  and  $T \geq 1$ ,

$$|\nabla\mathcal{J}_{T,\epsilon}(\mathbf{0})(t)| \geq \epsilon \frac{1}{2} (e^T - 1) \geq \epsilon \frac{1}{2} (e^T - 1) \geq \frac{\epsilon e^T}{4}.$$

Hence,

$$\|\nabla\mathcal{J}_{T,\epsilon}(\mathbf{u})(t)\|_{\mathcal{L}_2(\mathcal{U})}^2 \geq \frac{T}{2} \left( \frac{\epsilon e^T}{4} \right)^2,$$

so  $\|\nabla\mathcal{J}_{T,\epsilon}(\mathbf{u})(t)\|_{\mathcal{L}_2(\mathcal{U})} \geq \sqrt{T}\epsilon e^T / 4\sqrt{2}$ .  $\square$

## B.2. Global Stability Guarantees of JSPs and Consequences of (Westenbroek et al., 2021)

(Westenbroek et al., 2021) demonstrate that, for a certain class of nonlinear systems whose Jacobian Linearizations satisfy various favorable properties, an  $\epsilon$ -FOS point  $\mathbf{u}$  of the objective  $\mathcal{J}_T$  corresponds to a trajectory which converges exponentially to a desired equilibrium. Examining their proof, the first step follows from Westenbroek et al. (2021, Lemma 2), which establishes that  $\mathbf{u}$  is an  $\epsilon' = \epsilon^2/2\alpha$ -JSP, and it is this property (rather than the  $\epsilon$ -FOS) that is used throughout the rest of the proof. Hence, their result extends from FOSs to JSPs. Hence, the *local* optimization guarantees established in this work imply, via Westenbroek et al. (2021, Theorem 1), exponentially stabilizing *global* behavior.

## B.3. Projections to ensure boundedness.

Let us describe one way to ensure the feasibility condition, Definition 4.7. Suppose that  $f_{\text{dyn}}$  has the following stability property, which can be thought of as the state-output analogue of BIBO stability, and is common in the control literature (Jadbabaie & Hauser, 2001). For example, we may consider the following assumption.

**Assumption B.1.** There exists some function  $\phi_{\text{sys}} : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}_{\geq 0}$  such that that, if  $\|\mathbf{u}(t)\| \leq R$  for all  $t \in [0, T]$ , then  $\|\mathbf{x}(t | \mathbf{u})\| \leq \phi_{\text{sys}}(R, T)$  for all  $t \in [0, T]$ .

Next, fix a bound  $R_u > 0$ , and set

$$R_{\text{feas}} := 2 \max\{R_u, \phi_{\text{sys}}(R_u, T)\}$$

Then, it follows that for any policy for which

$$\mathbf{u}_k^\pi \leq R_u \tag{B.3}$$

for all  $k \in [K]$  is feasible in the sense of Definition 4.1. We therefore modify Algorithm 1, Line 5 to the *projected gradient step*

$$\mathbf{u}_{1:K}^{(n+1)} \leftarrow \text{Proj}_{(\mathcal{B}_{d_u}(R_u))^K} \left[ \text{orac}_{\pi^{(n)}, u}(\tilde{\mathbf{u}}_{1:K}^{(n)}) \right], \quad \text{where again } \tilde{\mathbf{u}}_k^{(n)} := \mathbf{u}_k^{(n)} - \frac{\eta}{\tau} \hat{\nabla}_k^{(n)} - K_k^{\pi^{(n)}} \hat{\mathbf{x}}_k,$$

where we let  $\text{Proj}_{(\mathcal{B}_{d_u}(R_u))^K}$  denote the orthogonal-projection on the  $K$ -fold project of  $d_u$ -dimensional balls of Euclidean radius  $R_u$ ,  $\mathcal{B}_{d_u}(R_u)$ . This projection is explicitly given by

$$\left( \text{Proj}_{(\mathcal{B}_{d_u}(R_u))^K} [\mathbf{u}_{1:K}] \right)_k = \mathbf{u}_k \cdot \min \left\{ 1, \frac{R_u}{\|\mathbf{u}_k\|} \right\},$$

here using the convention that when  $\mathbf{u}_k = 0$ , the above evaluates to 0. In this case, our algorithm converges (up to gradient estimation error) to a stationary-point of the projected gradient descent algorithm (see, e.g. the note <https://damek.github.io/teaching/orie6300/lec22.pdf> for details). We leave the control-theoretic interpretation of such stationary points to future work.

## B.4. Extensions to include Process Noise

As explained in [Section 2](#), [Oracle 2.1](#) only adds observation noise but not process noise. Process noise somewhat complicates the analysis, because then our method will only learn the Jacobians dynamics up to a noise floor determined by the process noise. However, by generalization our Taylor expansion of the dynamics (e.g. [Proposition A.5](#)), we can show that as the process noise magnitude decreases, we would achieve better and better accuracy, recovering the noiseless case in the limit. In addition, process noise may warrant greater algorithmic modifications: for example we may want to incorporate higher-order Taylor expansions of the dynamics (not just the Jacobian linearization), or more sophisticated gradient updates (i.e. iLQG ([\(Todorov & Li, 2005\)](#))) better tuned to handle process noise.

## B.5. Discussion of the $\exp(L_f t_0)$ dependence.

There are two sources of the exponential dependence on  $t_0 = \tau k_0$  that arises in our analysis. First, we translate open-loop controllability ([Assumption 4.4](#)) to closed-loop controllability needed for recovery of system matrices, in an argument based on ([Chen & Hazan, 2021](#)), and which incurs dependent on  $\exp(L_f t_{\text{ctrl}}) \leq \exp(L_f t_0)$ . Second, we only consider a stability modulus ([Definition 4.7](#)) for a Lyapunov equation terminating at  $k = k_0$ , because we do not estimate  $A_{\text{ol},k}^\pi, B_{\text{ol},k}^\pi$ , and therefore cannot synthesize the system gains, for  $k \leq k_0$ . This means that (see [Lemma A.1](#)) that many natural bounds on the discretized transition operators  $\|\Phi_{\text{cl},k,j}^\pi\|$  scale as  $\text{poly}(\mu_{\pi,*}, \exp(t_0 L_f))$ , yielding exponential dependence on  $t_0 L_f$ .

## C. Jacobian Linearizations

### C.1. Preliminaries

Recall  $\mathcal{U}$  denotes the space of continuous-time inputs  $\mathbf{u} : [0, T] \rightarrow \mathbb{R}^{d_u}$ , and  $\mathbf{U}$  continuous-time inputs  $\bar{\mathbf{u}} \in (\mathbb{R}^{d_u})^K$

#### C.1.1. EXACT TRAJECTORIES

We recall definitions of various trajectories.

**Definition C.1** (Open-Loop Trajectories and Nominal Trajectories). For a  $\mathbf{u} \in \mathcal{U}$ , we define  $\mathbf{x}(t \mid \mathbf{u})$  as the curve given by

$$\frac{d}{dt}\mathbf{x}(t \mid \mathbf{u}) = f_{\text{dyn}}(\mathbf{x}(t \mid \mathbf{u}), \mathbf{u}(t)), \quad \mathbf{x}(0 \mid \mathbf{u}) = \xi_{\text{init}}.$$

For a policy  $\pi = (\mathbf{u}_{1:K}^\pi, K_{1:K}^\pi)$ , we define  $\mathbf{u}^\pi = \text{ct}(\mathbf{u}_{1:K}^\pi)$ ,  $\mathbf{x}^\pi(t) = \mathbf{x}(t \mid \mathbf{u}^\pi)$ , and  $\mathbf{x}_k^\pi = \mathbf{x}^\pi(t_k)$ .

Similarly, we present a summary of the definition of various stabilized trajectories, consistent with [Definitions 2.1](#) and [4.4](#).

**Definition C.2** (Stabilized Trajectories). For  $\bar{\mathbf{u}} \in \mathcal{U}$  and a policy  $\pi$ , we define continuous-time perturbations of the dynamics with feedback

$$\tilde{\mathbf{x}}^{\pi, \text{ct}}(t \mid \bar{\mathbf{u}}) := \mathbf{x}(t \mid \tilde{\mathbf{u}}^{\pi, \text{ct}}), \quad \tilde{\mathbf{u}}^{\pi, \text{ct}}(t \mid \bar{\mathbf{u}}) := \bar{\mathbf{u}}(t) + K_{k(t)}^\pi (\tilde{\mathbf{x}}^{\pi, \text{ct}}(t_{k(t)} \mid \bar{\mathbf{u}}) - \mathbf{x}_{k(t)}^\pi),$$

there specialization to discrete-time inputs  $\bar{\mathbf{u}} \in \mathbf{U}$

$$\tilde{\mathbf{x}}^\pi(t \mid \bar{\mathbf{u}}) := \tilde{\mathbf{x}}^{\pi, \text{ct}}(t \mid \text{ct}(\bar{\mathbf{u}})), \quad \tilde{\mathbf{u}}^\pi(t \mid \bar{\mathbf{u}}) := \tilde{\mathbf{u}}^{\pi, \text{ct}}(t \mid \text{ct}(\bar{\mathbf{u}})),$$

and their discrete samplings

$$\tilde{\mathbf{x}}_k^\pi(\bar{\mathbf{u}}) := \tilde{\mathbf{x}}^\pi(t_k \mid \bar{\mathbf{u}}), \quad \tilde{\mathbf{u}}_k^\pi(\bar{\mathbf{u}}) := \tilde{\mathbf{u}}^\pi(t_k \mid \bar{\mathbf{u}})$$

#### C.1.2. TRAJECTORY LINEARIZATIONS

**Definition C.3** (Open-Loop Jacobian Linearizations of Trajectories). We define the continuous-time Jacobian linearizations

$$\begin{aligned} \mathbf{x}^{\text{jac}}(t \mid \bar{\mathbf{u}}) &:= \mathbf{x}(t \mid \mathbf{u}) + \langle \nabla_{\mathbf{u}} \mathbf{x}(t \mid \mathbf{u}) \Big|_{\mathbf{u}=\mathbf{u}^\pi}, \bar{\mathbf{u}} - \mathbf{u}^\pi \rangle_{\mathcal{L}_2(\mathcal{U})} \\ \delta \mathbf{x}^{\text{jac}}(t \mid \mathbf{u}; \bar{\mathbf{u}}) &:= \mathbf{x}^{\text{jac}}(t \mid \mathbf{u}; \bar{\mathbf{u}}) - \mathbf{x}(t \mid \mathbf{u}) \end{aligned} \tag{C.1}$$

**Definition C.4** (Closed-Loop Jacobian Linearizations of Trajectories, Discrete-Time).

$$\begin{aligned}\tilde{\mathbf{x}}^{\pi,\text{jac}}(t \mid \bar{\mathbf{u}}) &:= \tilde{\mathbf{x}}^{\pi,\text{ct}}(t \mid \mathbf{u}) + \langle \nabla_{\mathbf{u}} \tilde{\mathbf{x}}^{\pi,\text{ct}}(t \mid \mathbf{u}) \Big|_{\mathbf{u}=\mathbf{u}^\pi}, \bar{\mathbf{u}} - \mathbf{u}^\pi \rangle_{\mathcal{L}_2(\mathcal{U})} \\ \tilde{\mathbf{u}}^{\pi,\text{jac}}(t \mid \bar{\mathbf{u}}) &:= \tilde{\mathbf{u}}^{\pi,\text{ct}}(t \mid \mathbf{u}) + \langle \nabla_{\mathbf{u}} \tilde{\mathbf{u}}^{\pi,\text{ct}}(t \mid \mathbf{u}) \Big|_{\mathbf{u}=\mathbf{u}^\pi}, \bar{\mathbf{u}} - \mathbf{u}^\pi \rangle_{\mathcal{L}_2(\mathcal{U})}\end{aligned}\tag{C.2}$$

We further define the linearized differences

$$\delta \tilde{\mathbf{x}}^{\pi,\text{jac}}(t \mid \bar{\mathbf{u}}) := \tilde{\mathbf{x}}^{\pi,\text{jac}}(t \mid \bar{\mathbf{u}}) - \mathbf{x}^\pi(t), \quad \delta \tilde{\mathbf{u}}^{\pi,\text{jac}}(t \mid \bar{\mathbf{u}}) := \tilde{\mathbf{u}}^{\pi,\text{jac}}(t \mid \bar{\mathbf{u}}) - \mathbf{u}^\pi(t)\tag{C.3}$$

**Definition C.5** (Jacobian Linearization, with gains, discrete Time). Given  $\bar{\mathbf{u}} \in \mathcal{U}$ , we define

$$\begin{aligned}\tilde{\mathbf{x}}_k^{\pi,\text{jac}}(\bar{\mathbf{u}}) &:= \tilde{\mathbf{x}}_k^\pi(\bar{\mathbf{u}}) + \langle \nabla_{\mathbf{u}} \tilde{\mathbf{x}}_k^\pi(t \mid \bar{\mathbf{u}}) \Big|_{\bar{\mathbf{u}}=\mathbf{u}_{1:K}^\pi}, \bar{\mathbf{u}} - \mathbf{u}_{1:K}^\pi \rangle \\ \tilde{\mathbf{u}}_k^{\pi,\text{jac}}(\bar{\mathbf{u}}) &:= \tilde{\mathbf{u}}_k^\pi(\bar{\mathbf{u}}) + \langle \nabla_{\bar{\mathbf{u}}} \tilde{\mathbf{u}}_k^\pi(t \mid \mathbf{u}) \Big|_{\bar{\mathbf{u}}=\mathbf{u}_{1:K}^\pi}, \bar{\mathbf{u}} - \mathbf{u}_{1:K}^\pi \rangle_{\mathcal{L}_2(\mathcal{U})}\end{aligned}\tag{C.4}$$

We further define the linearized differences

$$\delta \tilde{\mathbf{x}}_k^{\pi,\text{jac}}(\bar{\mathbf{u}}) := \tilde{\mathbf{x}}_k^{\pi,\text{jac}}(\bar{\mathbf{u}}) - \mathbf{x}_k^\pi, \quad \delta \tilde{\mathbf{u}}_k^{\pi,\text{jac}}(\bar{\mathbf{u}}) := \tilde{\mathbf{u}}_k^{\pi,\text{jac}}(\bar{\mathbf{u}}) - \mathbf{u}_k^\pi\tag{C.5}$$

### C.1.3. JACOBIAN LINEARIZED DYNAMICS

We now recall the definitions of various linearizations, consistent with [Definition 4.6](#).

**Definition C.6** (Open-Loop, On-Policy Linearized Dynamics). We define the open-loop, on-policy linearization around a policy  $\pi$  via

$$\mathbf{A}_{\text{ol}}^\pi(t) = \partial_x f_{\text{dyn}}(\mathbf{x}^\pi(t), \mathbf{u}^\pi(t)), \quad \mathbf{B}_{\text{ol}}^\pi(t) = \partial_u f_{\text{dyn}}(\mathbf{x}^\pi(t), \mathbf{u}^\pi(t)).$$

**Definition C.7** (Open-Loop, On-Policy Linearized Transition, Markov Operators, and Discrete-Dynamics). We define the linearized transition function  $\Phi_{\text{ol}}^\pi(s, t)$  defined for  $s > t$  as the solution to  $\frac{d}{ds} \Phi_{\text{ol}}^\pi(s, t) = \mathbf{A}_{\text{ol}}^\pi(s) \Phi_{\text{ol}}^\pi(s, t)$ , with initial condition  $\Phi_{\text{ol}}^\pi(t, t) = \mathbf{I}$ . We *discretize* the open-loop transition function by define

$$\mathbf{A}_{\text{ol},k}^\pi = \Phi_{\text{ol}}^\pi(t_{k+1}, t_k), \quad \mathbf{B}_{\text{ol},k}^\pi := \int_{s=t_k}^{t_{k+1}} \Phi_{\text{ol}}^\pi(t_{k+1}, s) \mathbf{B}_{\text{ol}}^\pi(s) ds.$$

**Definition C.8** (Closed-Loop Jacobian Linearization, Discrete-Time). We define a *discrete-time closed-loop* linearization

$$\mathbf{A}_{\text{cl},k}^\pi := \mathbf{A}_{\text{ol},k}^\pi + \mathbf{B}_{\text{ol},k}^\pi \mathbf{K}_k^\pi = \Phi_{\text{ol}}^\pi(t_{k+1}, t) + \int_{s=t}^{t_{k+1}} \Phi_{\text{ol}}^\pi(t_{k+1}, s) \mathbf{B}_{\text{ol}}^\pi(s) \mathbf{K}_k^\pi,$$

and a discrete closed-loop *transition operator* is defined, for  $1 \leq k_1 \leq k_2 \leq K+1$ ,  $\Phi_{\text{cl},k_2,k_1}^\pi = \mathbf{A}_{\text{cl},k_2-1}^\pi \cdot \mathbf{A}_{\text{cl},k_2-2}^\pi \cdots \mathbf{A}_{\text{cl},k_1}^\pi$ , with the convention  $\Phi_{\text{cl},k_1,k_1}^\pi = \mathbf{I}$ . Finally, we define the closed-loop *markov operator* via  $\Psi_{\text{cl},k_2,k_1}^\pi := \Phi_{\text{cl},k_2,k_1+1}^\pi \mathbf{B}_{\text{ol},k_1}^\pi$  for  $1 \leq k_1 < k_2 \leq K+1$ .

**Definition C.9** (Closed-Loop Jacobian Linearizations, Continuous-Time). We define

$$\Phi_{\text{cl}}^\pi(s, t) := \begin{cases} \Phi_{\text{ol}}^\pi(s, t) & s, t \in \mathcal{I}_k \\ \tilde{\Phi}_{\text{cl}}^\pi(s, t_{k_2}) \cdot \Phi_{\text{cl},k_2,k_1}^\pi \cdot \Phi_{\text{ol}}^\pi(t_{k_1+1}, t) & t \in \mathcal{I}_{k_1}, s \in \mathcal{I}_{k_2}, k_2 > k_1, \end{cases}$$

where above, we define

$$\tilde{\Phi}_{\text{cl}}^\pi(s, t_k) = \Phi_{\text{ol}}^\pi(s, t_k) + \left( \int_{s'=t_k}^s \Phi_{\text{ol}}^\pi(s, s') \mathbf{B}_{\text{ol}}^\pi(s') ds \right) \mathbf{K}_k.$$

Lastly, we define

$$\Psi_{\text{cl}}^\pi(s, t) = \Phi_{\text{cl}}^\pi(s, t) \mathbf{B}(t).$$

## C.2. Characterizations of the Jacobian Linearizations

In this section we provide characterizations of the Jacobian Linearizations of the open-loop and closed-loop trajectories.

**Lemma C.1** (Implicit Characterization of the linearizations in open-loop). *Given  $\mathbf{u}, \bar{\mathbf{u}} \in \mathcal{U}$ , define  $\delta\bar{\mathbf{u}}(t) = \bar{\mathbf{u}}(t) - \mathbf{u}(t)$ . Then,*

$$\frac{d}{dt}\delta\mathbf{x}^{\text{jac}}(t | \bar{\mathbf{u}}; \mathbf{u}) = \mathbf{A}(t | \mathbf{u})\delta\mathbf{x}^{\text{jac}}(t | \bar{\mathbf{u}}; \mathbf{u}) + \mathbf{B}(t | \mathbf{u})\delta\bar{\mathbf{u}}(t)$$

with initial condition  $\delta\mathbf{x}^{\text{jac}}(0 | \bar{\mathbf{u}}) = 0$ , where

$$\mathbf{A}(t | \mathbf{u}) = \partial_x f_{\text{dyn}}(x, u)|_{x=\mathbf{x}(t|\mathbf{u}), u=\mathbf{u}(t)} \quad \text{and} \quad \mathbf{B}(t | \mathbf{u}) = \partial_u f_{\text{dyn}}(x, u)|_{x=\mathbf{x}(t|\mathbf{u}), u=\mathbf{u}(t)}.$$

*Proof.* The result follows directly from [Lemma C.8](#) and the definition of  $\delta\mathbf{x}^{\text{jac}}(t | \bar{\mathbf{u}}; \mathbf{u})$ .  $\square$

**Lemma C.2** (Implicit Characterization of the linearizations in closed-loop). *Given a policy  $\pi$  and  $\bar{\mathbf{u}} \in \mathcal{U}$ , set  $\delta\bar{\mathbf{u}}^\pi(t) = \mathbf{u}^\pi(t) - \mathbf{u}(t)$ . Then, recalling  $\delta\tilde{\mathbf{x}}^{\pi, \text{jac}}(t | \bar{\mathbf{u}}) = \tilde{\mathbf{x}}^{\pi, \text{jac}}(t | \bar{\mathbf{u}}) - \mathbf{x}^\pi(t)$ ,*

$$\begin{aligned} \frac{d}{dt}\delta\tilde{\mathbf{x}}^{\pi, \text{jac}}(t | \bar{\mathbf{u}}) &= \mathbf{A}_{\text{ol}}^\pi(t)\delta\tilde{\mathbf{x}}^{\pi, \text{ct}}(t | \bar{\mathbf{u}}) + \mathbf{B}_{\text{ol}}^\pi(t)\delta\tilde{\mathbf{u}}^{\pi, \text{jac}}(t) \\ \delta\tilde{\mathbf{u}}^{\pi, \text{jac}} &:= \delta\bar{\mathbf{u}}^\pi(t) + \mathbf{K}_{k(t)}^\pi\delta\tilde{\mathbf{x}}^{\pi, \text{jac}}(t_{k(t)} | \bar{\mathbf{u}}), \end{aligned}$$

with initial condition  $\delta\tilde{\mathbf{x}}^{\pi, \text{jac}}(0 | \bar{\mathbf{u}}) = 0$ .

*Proof.* The result follows directly from [Lemma C.8](#) and the definitions of  $\delta\tilde{\mathbf{x}}^{\pi, \text{jac}}(t | \bar{\mathbf{u}})$  and the construction of the perturbed input  $\delta\tilde{\mathbf{u}}^{\pi, \text{jac}}$ .  $\square$

**Lemma C.3** (Explicit Characterizations of Linearizations, Continuous-Time). *For a policy  $\pi$ , we have:*

$$\begin{aligned} \delta\mathbf{x}^{\text{jac}}(t | \mathbf{u}^\pi + \delta\mathbf{u}; \mathbf{u}^\pi) &= \int_{s=0}^t \Phi_{\text{ol}}^\pi(t, s)\mathbf{B}_{\text{ol}}^\pi(s)\delta\mathbf{u}(s)ds. \\ \delta\tilde{\mathbf{x}}^{\pi, \text{jac}}(t | \mathbf{u}^\pi + \delta\mathbf{u}) &= \int_{s=0}^t \Phi_{\text{cl}}^\pi(t, s)\mathbf{B}_{\text{ol}}^\pi(s)\delta\mathbf{u}(s)ds. \end{aligned}$$

*Proof.* The first condition follows directly from the characterization of the evolution of  $\delta\mathbf{x}^{\text{jac}}(t | \bar{\mathbf{u}}; \mathbf{u})$  and [Lemma C.8](#). For the second condition, we will directly argue that the proposed formula satisfied the differential equation in [Lemma C.2](#). By the Leibniz integral rule we have:

$$\begin{aligned} \frac{d}{dt}\left(\int_{s=0}^t \Phi_{\text{cl}}^\pi(t, s)\mathbf{B}_{\text{ol}}^\pi(s)\delta\mathbf{u}(s)ds\right) &= \int_{s=0}^t \frac{d}{dt}\left(\Phi_{\text{cl}}^\pi(t, s)\mathbf{B}_{\text{ol}}^\pi(s)\delta\mathbf{u}(s)\right)ds + \Phi_{\text{cl}}^\pi(t, t)\mathbf{B}_{\text{ol}}^\pi(t)\delta\mathbf{u}(t) \\ &= \int_{s=0}^t \frac{d}{dt}\left(\Phi_{\text{cl}}^\pi(t, s)\mathbf{B}_{\text{ol}}^\pi(s)\delta\mathbf{u}(s)\right)ds + \mathbf{B}_{\text{ol}}^\pi(t)\delta\mathbf{u}(t). \end{aligned}$$

Using the expression for  $\Phi_{\text{cl}}^\pi(t, s)$  in [Definition C.9](#), we can similarly calculate:

$$\frac{d}{dt}\Phi_{\text{cl}}^\pi(t, s) = \mathbf{A}_{\text{ol}}^\pi(t)\Phi_{\text{cl}}^\pi(t, s) + \mathbf{B}_{\text{ol}}^\pi(t)\mathbf{K}_k\Phi_{\text{cl}}^\pi(t_{k(t)}, s).$$

Together the preceding quantities demonstrate that:

$$\begin{aligned} \frac{d}{dt}\left(\int_{s=0}^t \Phi_{\text{cl}}^\pi(t, s)\mathbf{B}_{\text{ol}}^\pi(s)\delta\mathbf{u}(s)ds\right) &= \mathbf{A}_{\text{ol}}^\pi(t) \cdot \left(\int_{s=0}^t \Phi_{\text{cl}}^\pi(t, s)\mathbf{B}_{\text{ol}}^\pi(s)\delta\mathbf{u}(s)ds\right) + \mathbf{B}_{\text{ol}}^\pi(t)\delta\mathbf{u}(t) \\ &\quad + \mathbf{K}_{k(t)}^\pi \cdot \left(\int_{s=0}^{t_{k(t)}} \Phi_{\text{cl}}^\pi(t, s)\mathbf{B}_{\text{ol}}^\pi(s)\delta\mathbf{u}(s)ds\right), \end{aligned}$$

which demonstrates the proposed solutions satisfies the desired differential equation.  $\square$

**Lemma C.4** (Explicit Characterizations of Linearizations, Discrete-Time). *For a policy  $\pi$ , and perturbation  $\delta \mathbf{u}_{1:K} \in \mathbf{U}$ ,*

$$\delta \tilde{\mathbf{x}}_k^{\pi, \text{jac}}(\mathbf{u}_{1:K}^\pi + \delta \mathbf{u}_{1:K}) = \sum_{j=1}^{k-1} \Psi_{\text{cl}, k, j}^\pi \delta \mathbf{u}_j = \sum_{j=1}^{k-1} \Phi_{\text{cl}, k, j+1}^\pi \mathbf{B}_{\text{ol}, j}^\pi \delta \mathbf{u}_j.$$

*Proof.* The proof follows directly from [Definition C.9](#), [Definition C.8](#) and [Lemma C.3](#).  $\square$

### C.3. Gradient Computations

**Lemma C.5** (Computation of Continuous-Time Gradient, Open-Loop). *Fix  $\xi = \xi^\pi$ . Define  $Q_{(\cdot)}^\pi(t) := \partial_{(\cdot)} Q(\mathbf{x}^\pi(t), \mathbf{u}^\pi(t), t)$ . Then,*

$$\nabla \mathcal{J}_T(\mathbf{u})(t) \Big|_{\mathbf{u}=\mathbf{u}^\pi} = Q_u^\pi(t) + \int_{s=t}^T \Phi_{\text{ol}}^\pi(s, t) \mathbf{B}_{\text{ol}}^\pi(t) ds.$$

and

$$\langle \nabla \mathcal{J}_T(\mathbf{u})(t) \Big|_{\mathbf{u}=\mathbf{u}^\pi}, \delta \mathbf{u} \rangle = \int_0^T (\langle Q_u^\pi(t), \delta \mathbf{u}(t) \rangle + \langle Q_x^\pi(t), \delta \mathbf{x}^{\text{jac}}(t | \mathbf{u} + \delta \mathbf{u}) \rangle) dt.$$

*Proof.* For a given perturbation  $\delta \mathbf{u}$ , by the chain rule we have:

$$D\mathcal{J}_T(\mathbf{u})[\delta \mathbf{u}] = \int_0^T (\langle Q_u^\pi(t), \delta \mathbf{u}(t) \rangle + \langle Q_x^\pi(t), \delta \mathbf{x}^{\text{jac}}(t | \mathbf{u} + \delta \mathbf{u}) \rangle) dt + \langle \partial_x V(\mathbf{x}^\pi(T)), \delta \mathbf{x}^{\text{jac}}(T | \mathbf{u} + \delta \mathbf{u}) \rangle$$

Because  $\delta \mathbf{u}$  is arbitrary, an application of [Lemma C.3](#) demonstrates the desired results.  $\square$

**Lemma C.6** (Computation of Continuous-Time Gradient, Closed-Loop). *Fix  $\xi = \xi^\pi$ . Define  $Q_{(\cdot)}^\pi(t) := \partial_{(\cdot)} Q(\mathbf{x}^\pi(t), \mathbf{u}^\pi(t), t)$ . Then,*

$$\begin{aligned} \nabla \mathcal{J}_T(\pi)(t) := \nabla_{\bar{\mathbf{u}}} \mathcal{J}_T^\pi(\bar{\mathbf{u}}) \Big|_{\bar{\mathbf{u}}=\mathbf{u}^\pi}(t) &= Q_u^\pi(t) + \Psi_{\text{cl}}^\pi(T, t)^\top (\partial_x V(\mathbf{x}^\pi(T))) \\ &+ \int_{s=t}^T \Psi_{\text{cl}}^\pi(s, t)^\top Q_x^\pi(s) ds + \int_{s=t_{k(t)+1}}^T \Psi_{\text{cl}}^\pi(t_{k(s)}, t)^\top \mathbf{K}_{k(s)}^\top Q_u^\pi(s) ds. \end{aligned}$$

*Proof.* The proof follows the steps of [Lemma C.5](#), but replaces the open-loop state and input perturbations with the appropriate closed-loop perturbations, as defined in [Lemma C.3](#) and calculated in [Lemma C.2](#).  $\square$

Similarly, we can compute the gradient of the discrete-time objective. Its proof is analogous to the previous two.

**Lemma C.7** (Computation of Discrete-Time Gradient).

$$\begin{aligned} (\nabla \mathcal{J}_T^{\text{disc}}(\pi))_k &= \tau Q_u(\mathbf{x}_k^\pi, \mathbf{u}_k^\pi, t_k) + (\Psi_{\text{cl}, K+1, k}^\pi)^\top V_x(\mathbf{x}_{K+1}^\pi) \\ &+ \tau \sum_{j=k+1}^K (\Psi_{\text{cl}, j, k}^\pi)^\top (Q_x(\mathbf{x}_j^\pi, \mathbf{u}_j^\pi, t_j) + (\mathbf{K}_j^\pi)^\top Q_u(\mathbf{x}_j^\pi, \mathbf{u}_j^\pi, t_j)) \end{aligned}$$

Moreover, defining the shorthand  $\delta \tilde{\mathbf{x}}_k^{\text{jac}} = \delta \tilde{\mathbf{x}}_k^{\pi, \text{jac}}(\mathbf{u}_{1:K}^\pi + \delta \mathbf{u}_{1:K})$ ,

$$\begin{aligned} \langle \delta \mathbf{u}_{1:K}, \nabla \mathcal{J}_T^{\text{disc}}(\pi) \rangle &= \langle \partial_x V(\mathbf{x}_{K+1}^\pi), \delta \tilde{\mathbf{x}}_k^{\text{jac}} \rangle + \tau \sum_{k=1}^K \langle \partial_x Q(\mathbf{x}_k^\pi, \mathbf{u}_k^\pi, t_k), \delta \tilde{\mathbf{x}}_k^{\text{jac}} \rangle + \langle \partial_u Q(\mathbf{x}_k^\pi, \mathbf{u}_k^\pi, t_k), \delta \mathbf{u}_k + \mathbf{K}_k \delta \tilde{\mathbf{x}}_k^{\text{jac}} \rangle. \end{aligned}$$

#### C.4. Technical Tools

The first supportive lemma is a standard result from variational calculus, and characterizes how the solution to the controlled differential equation changes under perturbations to the input. Note that this result does not depend on how the input is generated, namely, whether the perturbation is generated in open-loop or closed-loop. Concretely, the statement of the following result is equivalent to Theorem 5.6.9 from (Polak, 2012).

**Lemma C.8** (State Variation of Controlled CT Systems). *For each nominal input  $\mathbf{u} \in \mathcal{U}$  and perturbation  $\delta\mathbf{u} \in \mathcal{U}$  we have:*

$$\langle \nabla_{\mathbf{u}} \mathbf{x}(t \mid \mathbf{u}), \delta\mathbf{u} \rangle = \delta x(t), \quad (\text{C.6})$$

where the curve  $\delta x(\cdot)$  satisfies:

$$\frac{d}{dt} \delta x(t) = \mathbf{A}(t \mid \mathbf{u}) \delta x(t) + \mathbf{B}(t \mid \mathbf{u}) \delta\mathbf{u}(t), \quad (\text{C.7})$$

where we recall that:

$$\mathbf{A}(t \mid \mathbf{u}) = \partial_x f_{\text{dyn}}(x, u) \Big|_{x=\mathbf{x}(t|\mathbf{u}), u=\mathbf{u}(t)} \quad \text{and} \quad \mathbf{B}(t \mid \mathbf{u}) = \partial_u f_{\text{dyn}}(x, u) \Big|_{x=\mathbf{x}(t|\mathbf{u}), u=\mathbf{u}(t)}.$$

Moreover, we have

$$\delta x(t) = \int_{s=0}^t \Phi(t, s) \mathbf{B}(s \mid \mathbf{u}) \delta\mathbf{u}(s),$$

where the transition operator satisfies  $\frac{d}{dt} \Phi(t, s) = \mathbf{A}(t \mid \mathbf{u}) \Phi$  and  $\Phi(s, s) = \mathbf{I}$ .

The following result is equivalent to Lemma 5.6.2 from (Polak, 2012).

**Lemma C.9** (Picard Lemma). *Consider two dynamical  $\frac{d}{dt} y(t) = \phi(y(t), t)$ ,  $i \in \{1, 2\}$ , and suppose that  $y \mapsto \phi(y, s)$  is  $L$ -Lipschitz for each  $s$  fixed. Let  $z(t)$  be any other absolutely continuous curve. Then,*

$$\|y(t) - z(t)\| \leq \exp(tL) \cdot \left( \|y(0) - z(0)\| + \int_{s=0}^t \left\| \frac{d}{ds} z(s) - f(z(s), s) \right\| \right).$$

**Lemma C.10** (Solution to Affine ODEs). *Consider an affine ODE given by  $\mathbf{y}(0) = y_0$ ,  $\frac{d}{dt} \mathbf{y}(t) = \mathbf{A}(t) \mathbf{y}(t) + \mathbf{B}(t) u$ . Then,*

$$y(t) = \Phi(t, 0) y_0 + \left( \int \Phi(t, s) \mathbf{B}(s) ds \right) u,$$

where  $\Phi(t, s)$  solves the ODE  $\Phi(s, s) = \mathbf{I}$  and  $\frac{d}{dt} \Phi(t, s) = \mathbf{A}(t) \Phi(t, s)$ .

## D. Taylor Expansions of the Dynamics

### D.1. Proof of Proposition A.5

Recall  $\delta \mathbf{u}_k = \mathbf{u}_k - \mathbf{u}_k^\pi$ , and define  $\mathbf{u} = \text{ct}(\mathbf{u}_{1:K})$  and  $\delta \mathbf{u} := \text{ct}(\mathbf{u}_{1:K}^\pi)$ . We define shorthand for relevant continuous curves and their discretizations:

$$\begin{aligned} \mathbf{y}(t) &= \tilde{\mathbf{x}}^\pi(t \mid \mathbf{u}^\pi + \delta \mathbf{u}) = \tilde{\mathbf{x}}^\pi(t \mid \mathbf{u}^\pi) \\ \mathbf{y}_k &= \mathbf{y}(t_k) = \tilde{\mathbf{x}}_k^\pi(\mathbf{u}_{1:K}^\pi + \delta \mathbf{u}_{1:K}) = \tilde{\mathbf{x}}_k^\pi(\mathbf{u}_{1:K}) \\ \mathbf{y}^{\text{jac}}(t) &= \tilde{\mathbf{x}}^{\pi, \text{jac}}(t \mid \mathbf{u}^\pi + \delta \mathbf{u}), \quad \mathbf{y}_k^{\text{jac}} = \mathbf{y}^{\text{jac}}(t_k) \end{aligned} \quad (\text{D.1})$$

We also define their differences from the nominal as

$$\delta \mathbf{y}(t) = \mathbf{y}(t) - \mathbf{x}^\pi(t), \quad \delta \mathbf{y}^{\text{jac}}(t) = \mathbf{y}^{\text{jac}}(t) - \mathbf{x}^\pi(t), \quad \delta \mathbf{y}_k = \mathbf{y}_k - \mathbf{x}_k^\pi, \quad \delta \mathbf{y}_k^{\text{jac}} = \mathbf{y}_k^{\text{jac}} - \mathbf{x}_k^\pi.$$

And the Jacobian error

$$\mathbf{e}^{\text{jac}}(t) := \mathbf{y}(t) - \mathbf{y}^{\text{jac}}(t), \quad \mathbf{e}_k^{\text{jac}} := \mathbf{e}^{\text{jac}}(t_k) = \mathbf{y}_k - \mathbf{y}_k^{\text{jac}}.$$

The main challenge is recursively controlling  $\|\mathbf{e}_k^{\text{jac}}\|$ . We begin with a computation which is immediate from Definition C.1 (the first equality) and (the second equality):

**Lemma D.1** (Curve Computations). *For  $t \in \mathcal{I}_k$ ,*

$$\begin{aligned} \frac{d}{dt} \mathbf{y}(t) &= f_{\text{dyn}}(\mathbf{y}(t), \mathbf{u}_k^\pi + \delta \mathbf{u}_k + \mathbf{K}_k^\pi(\delta \mathbf{y}_k)) \\ \frac{d}{dt} \mathbf{y}^{\text{jac}}(t) &= \frac{d}{dt} \mathbf{x}^\pi(t) + \mathbf{A}_{\text{ol}}^\pi(t) \delta \mathbf{y}^{\text{jac}}(t) + \mathbf{B}_{\text{ol}}^\pi(t) (\delta \mathbf{u}_k + \mathbf{K}_k(\delta \mathbf{y}_k^{\text{jac}})) \end{aligned}$$

**Computing the Jacobian Linearization.** The first step of the proof is a computation of the Jacobian linearization and a bound on its magnitude.

**Lemma D.2** (Computation of Jacobian Linearization).

$$\mathbf{y}_k^{\text{jac}} = \sum_{j=1}^{k-1} \Phi_{\text{cl}, k, j+1}^\pi \mathbf{B}_{\text{ol}, j}^\pi \delta \mathbf{u}_j = \sum_{j=1}^{k-1} \Psi_{\text{cl}, k, j}^\pi \delta \mathbf{u}_j \quad (\text{D.2})$$

Therefore,

$$\max_{k \in [K+1]} \|\mathbf{y}_k^{\text{jac}}\| \leq L_{\text{ol}} L_f \min\{B_2 \kappa_{\pi, 2}, B_\infty \kappa_{\pi, 1}\} \quad (\text{D.3})$$

*Proof.* Eq. (D.2) follows from Lemma C.6. Eq. (D.3) follows from Cauchy Schwartz/Holder's inequality, and the bound  $\|\mathbf{B}_{\text{ol}, j}^\pi\| \leq \tau L_{\text{ol}} L_f$  due to Lemma I.3.  $\square$

**Recursion on proximity to Jacobian linearization.** Next, we argue that the true dynamics  $\mathbf{y}(t)$  remain close to  $\mathbf{y}^{\text{jac}}(t)$ . We establish a recursion under the following invariant:

$$\begin{aligned} \|\mathbf{y}_k^{\text{jac}}\| \vee \|\mathbf{y}_k\| \vee \|\mathbf{u}_k^\pi + \delta \mathbf{u}_k + \mathbf{K}_k^\pi(\delta \mathbf{y}_k)\| \vee \|\mathbf{u}_k^\pi + \delta \mathbf{u}_k + \mathbf{K}_k^\pi(\delta \mathbf{y}_k^{\text{jac}})\| &\leq \frac{3}{4} R_{\text{feas}} \\ \tau &\leq \min \left\{ \frac{1}{16 L_f L_\pi}, \frac{1}{8 \kappa_f} \right\} \end{aligned} \quad (\text{D.4})$$

We prove the following recursion:

**Lemma D.3** (Recursion on Error of Linearization). *Suppose Eq. (D.4) holds. Let  $\tilde{\Phi}_{\text{cl}, \pi}(t, t_k) = \Phi_{\text{ol}}^\pi(t, t_k) + (\int_{s=t_k}^t \Phi_{\text{ol}}^\pi(t, s) \mathbf{B}_{\text{ol}}^\pi(s) ds) \mathbf{K}_k^\pi$ . Then, the following bound holds:*

$$\sup_{t \in \mathcal{K}} \|\mathbf{e}^{\text{jac}}(t) - \tilde{\Phi}_{\text{cl}, \pi}(t, t_k) \mathbf{e}_k^{\text{jac}}\| \leq M_f \tau \left( 4 \|\delta \mathbf{u}_k\|^2 + 20 L_\pi^2 \|\mathbf{e}_k^{\text{jac}}\|^2 + 20 L_\pi^2 \|\mathbf{y}_k^{\text{jac}} - \mathbf{y}_k\|^2 \right).$$

In particular, we have

$$\|\mathbf{e}_{k+1}^{\text{jac}} - \mathbf{A}_{\text{cl}, k}^\pi \mathbf{e}_k^{\text{jac}}\| \leq M_f \tau \left( 4 \|\delta \mathbf{u}_k\|^2 + 20 L_\pi^2 \|\mathbf{e}_k^{\text{jac}}\|^2 + 20 L_\pi^2 \|\mathbf{y}_k^{\text{jac}} - \mathbf{y}_k\|^2 \right).$$



To do prove [Lemma D.3](#), we introduce another family of curves  $\check{\mathbf{y}}_k^{\text{jac}}(t)$ , defined for  $t \geq t_k$ , which begin at  $\mathbf{y}(t_k)$  but evolve according to the Jacobian linearization:

$$\begin{aligned} \frac{d}{dt} \check{\mathbf{y}}_k^{\text{jac}}(t) &= \frac{d}{dt} \mathbf{x}^\pi(t) + \mathbf{A}_{\text{ol}}^\pi(t) \delta \check{\mathbf{y}}_k^{\text{jac}}(t) + \mathbf{B}_{\text{ol}}^\pi(t) (\delta \mathbf{u}_k + \mathbf{K}_k(\delta \mathbf{y}_k)) \\ \check{\mathbf{y}}_k^{\text{jac}}(t_k) &= \mathbf{y}(t_k) = \mathbf{y}_k, \quad \delta \mathbf{y}_k^{\text{jac}}(t) = \mathbf{y}_k^{\text{jac}}(t) - \mathbf{x}^\pi(t). \end{aligned}$$

We begin by establishing feasibility of all relevant continuous-time curves on the interval  $\mathcal{I}_k$ .

**Lemma D.4.** *Suppose that [Eq. \(D.4\)](#) holds. Then, for all  $t \in \mathcal{I}_k$ ,*

$$\|\check{\mathbf{y}}_k^{\text{jac}}(t)\| \vee \|\mathbf{y}(t)\| \vee \|\mathbf{y}^{\text{jac}}(t)\| \leq R_{\text{feas}}.$$

*Proof.* Let us start with  $\mathbf{y}(t)$ . Define the shorthand  $\bar{\mathbf{u}}_k := \mathbf{u}_k^\pi + \delta \mathbf{u}_k + \mathbf{K}_k^\pi(\delta \mathbf{y}_k)$ , so that  $\frac{d}{dt} \mathbf{y}(t) = f_{\text{dyn}}(\mathbf{y}(t), \mathbf{u}_k)$ . Under [Eq. \(D.4\)](#),  $\|\bar{\mathbf{u}}_k\| \leq \frac{3}{4} R_{\text{feas}}$ . At  $t = t_k$ ,  $\|\mathbf{y}(t)\| \leq \frac{3}{4} R_{\text{feas}}$ . Moreover, if at a given  $t$ ,  $\|\mathbf{y}(t)\| \leq R_{\text{feas}}$ , then  $(\mathbf{y}(t), \mathbf{u}_k)$  is feasible, so  $\|\frac{d}{dt} \mathbf{y}(t)\| = \|f_{\text{dyn}}(\mathbf{y}(t), \bar{\mathbf{u}}_k)\| \leq \kappa_f$ . Thus, letting  $t_\star := \sup\{t \in \mathcal{I}_K : \|\mathbf{y}(t)\| \leq R_{\text{feas}}\}$ , we see that if  $\tau \leq \frac{R_{\text{feas}}}{4\kappa_f}$ ,  $t_\star = t_{k+1}$ .

The arguments for  $\check{\mathbf{y}}_k^{\text{jac}}(t)$  and  $\mathbf{y}^{\text{jac}}(t)$  are similar: if, say,  $\|\check{\mathbf{y}}_k^{\text{jac}}(t)\| \leq R_{\text{feas}}$  for a given  $t \in \mathcal{I}_k$ , then by

$$\begin{aligned} \left\| \frac{d}{dt} \check{\mathbf{y}}_k^{\text{jac}}(t) \right\| &= \left\| \frac{d}{dt} \mathbf{x}^\pi(t) + \mathbf{A}_{\text{ol}}^\pi(t) \delta \check{\mathbf{y}}_k^{\text{jac}}(t) + \mathbf{B}_{\text{ol}}^\pi(t) \mathbf{u}_k \right\| \\ &\leq \kappa_f + L_f (\|\delta \check{\mathbf{y}}_k^{\text{jac}}(t)\| + \|\mathbf{u}_k\|) \\ &\leq \kappa_f + 2R_{\text{feas}} L_f. \end{aligned}$$

where above we use [Eq. \(D.4\)](#) to bound  $\|\mathbf{u}_k\| \leq R_{\text{feas}}$ , feasibility of  $\pi$ . As  $\|\delta \check{\mathbf{y}}_k^{\text{jac}}(t)\| \leq \frac{3}{4} R_{\text{feas}}$ , integrating (specifically, again considering  $t_\star := \sup\{t \in \mathcal{I}_K : \|\check{\mathbf{y}}_k^{\text{jac}}(t)\| \leq R_{\text{feas}}\}$ ) shows that as long as  $\tau(\kappa_f + 2L_f R_{\text{feas}}) \leq R_{\text{feas}}/4$ ,  $\|\check{\mathbf{y}}_k^{\text{jac}}(t)\| \leq R_{\text{feas}}$  for all  $t \in \mathcal{I}_k$ . For this, it suffices that  $\tau \leq \min\{\frac{1}{16L_f}, 1/8\kappa_f\}$ .  $\square$

We continue with a crude bound on the difference  $\delta \check{\mathbf{y}}_k^{\text{jac}}(t) = \check{\mathbf{y}}_k^{\text{jac}}(t) - \mathbf{x}^\pi(t)$ .

**Lemma D.5.** *Suppose [Eq. \(D.4\)](#) holds all  $t \in \mathcal{I}_k$*

$$\|\delta \check{\mathbf{y}}_k^{\text{jac}}(t)\| \leq (\tau L_{\text{ol}} L_f \|\delta \mathbf{u}_k\| + L_{\text{ol}}(1 + \tau L_\pi L_f) \|\delta \mathbf{y}_k\|)$$

*Similarly,*

$$\|\delta \mathbf{y}(t)\| \leq (\tau L_{\text{ol}} L_f \|\delta \mathbf{u}_k\| + L_{\text{ol}}(1 + \tau L_\pi L_f) \|\delta \mathbf{y}_k\|).$$

*Proof.* Then, Picard's Lemma ([Lemma C.9](#)), feasibility of  $\pi$  and [Assumption 4.1](#) implies that, for any  $t \in \mathcal{I}_k$

$$\|\delta \check{\mathbf{y}}_k^{\text{jac}}(t)\| \leq \exp((t - t_k)L_f) \epsilon_1$$

where

$$\begin{aligned} \epsilon_1 &:= \|\delta \mathbf{y}_k\| + \int_{s=t_k}^t \|\mathbf{B}_{\text{ol}}^\pi(s) (\delta \mathbf{u}_k + \mathbf{K}_k(\delta \mathbf{y}_k))\| ds \\ &\leq \|\delta \mathbf{y}_k\| + \int_{s=t_k}^{t_{k+1}} L_f (\|\delta \mathbf{u}_k\| + \|\mathbf{K}_k^\pi \delta \mathbf{y}_k\|) ds \quad (\text{Assumption 4.1}) \\ &\leq \tau(L_f \|\delta \mathbf{u}_k\| + L_\pi L_f \|\delta \mathbf{y}_k\|), \quad (\text{Definition 4.7}) \end{aligned}$$

Bounding  $\exp((t - t_k)L_f) \leq \exp(\tau L_f) = L_{\text{ol}}$  concludes the first part. The second part follows from a similar argument, using Lipschitzness of  $f_{\text{dyn}}$  in accordance with [Assumption 4.1](#), and the feasibility of  $(\mathbf{y}(t), \mathbf{u}_k^\pi + \delta \mathbf{u}_k + \mathbf{K}_k^\pi(\delta \mathbf{y}_k^{\text{jac}}))$  for  $t \in \mathcal{I}_k$ , as ensured by [Eq. \(D.4\)](#) and [Lemma D.4](#).  $\square$

We are now ready to prove [Lemma D.3](#).

*Proof of Lemma D.3.* Observe that, for  $t \in \mathcal{I}_k$

$$\frac{d}{dt}(\tilde{\mathbf{y}}_k^{\text{jac}}(t) - \mathbf{y}^{\text{jac}}(t)) = \mathbf{A}_{\text{ol}}^\pi(t)(\tilde{\mathbf{y}}_k^{\text{jac}}(t) - \mathbf{y}^{\text{jac}}(t)) + \mathbf{B}_{\text{ol}}^\pi(t)\mathbf{K}_k^\pi(\tilde{\mathbf{y}}_k^{\text{jac}}(t_k) - \mathbf{y}^{\text{jac}}(t_k))$$

By solving the affine ODE define  $\tilde{\mathbf{y}}_k^{\text{jac}}(t) - \mathbf{y}^{\text{jac}}(t)$  (applying Lemma C.10), and recalling all various definitions,

$$\tilde{\mathbf{y}}_k^{\text{jac}}(t) - \mathbf{y}^{\text{jac}}(t) = \tilde{\Phi}_{\text{cl},\pi}(t, t_k)(\mathbf{y}(t_k) - \mathbf{y}^{\text{jac}}(t_k)). \quad (\text{D.5})$$

We now bound  $\tilde{\mathbf{y}}_k^{\text{jac}}(t) - \mathbf{y}(t)$ . By applying Picard's Lemma (Lemma C.9) and Assumption 4.1 with  $L_{\text{ol}} = \exp(\tau L_f)$  to control the Lipschitz constant contribution, and using the agreement of initial conditions  $\tilde{\mathbf{y}}_k^{\text{jac}}(t_k) = \mathbf{y}(t_k)$ ,

$$\|\tilde{\mathbf{y}}_k^{\text{jac}}(t) - \mathbf{y}(t)\| \leq L_{\text{ol}} \int_{s=t_k}^t \|\Delta(s)\| ds. \quad (\text{D.6})$$

where

$$\Delta(s) = f_{\text{dyn}}(\tilde{\mathbf{y}}_k^{\text{jac}}(s), \mathbf{u}_k^\pi + \delta \mathbf{u}_k + \mathbf{K}_k^\pi \delta \mathbf{y}_k) - \frac{d}{ds} \tilde{\mathbf{y}}_k^{\text{jac}}(s).$$

By a Taylor expansion, we have

$$\begin{aligned} \Delta(s) &= f_{\text{dyn}}(\mathbf{x}^\pi(s), \mathbf{u}_k^\pi) - \frac{d}{ds} \tilde{\mathbf{y}}_k^{\text{jac}}(s) + \partial_x f_{\text{dyn}}(\mathbf{x}^\pi(s), \mathbf{u}_k^\pi) \delta \tilde{\mathbf{y}}_k^{\text{jac}}(s) + \partial_u f_{\text{dyn}}(\mathbf{x}^\pi(s), \mathbf{u}_k^\pi) (\delta \mathbf{u}_k + \mathbf{K}_k^\pi \delta \mathbf{y}_k) \\ &\quad + \frac{1}{2} \text{remainder}, \end{aligned}$$

where we bound

$$\begin{aligned} \|\text{remainder}(s)\| &\leq \sup_{\alpha \in [0,1]} \|\nabla^2 f_{\text{dyn}}(\alpha \mathbf{x}^\pi(s) + (1-\alpha) \mathbf{y}_k^{\text{jac}}(s), \alpha \mathbf{u}_k^\pi + (1-\alpha)(\mathbf{u}_k^\pi + \delta \mathbf{u}_k + \mathbf{K}_k^\pi \delta \mathbf{y}_k))\| \\ &\quad \cdot (\|\delta \tilde{\mathbf{y}}_k^{\text{jac}}(s)\|^2 + \|\delta \mathbf{u}_k + \mathbf{K}_k^\pi \delta \mathbf{y}_k\|^2). \end{aligned}$$

From by feasibility of  $\pi$ ,  $\|\mathbf{u}_k^\pi\| \vee \|\mathbf{x}^\pi(s)\| \leq R_{\text{feas}}$ . Moreover,  $\|\mathbf{u}_k^\pi + \delta \mathbf{u}_k + \mathbf{K}_k^\pi \delta \mathbf{y}_k\| \leq R_{\text{feas}}$  by Eq. (D.4) and  $\|\mathbf{y}_k^{\text{jac}}(s)\| \leq R_{\text{feas}}$  by Lemma D.4. Thus, for  $\alpha \in [0, 1]$ ,

$$\|\alpha \mathbf{x}^\pi(s) + (1-\alpha) \mathbf{y}_k^{\text{jac}}(s)\| \vee \|\alpha \mathbf{u}_k^\pi + (1-\alpha)(\mathbf{u}_k^\pi + \delta \mathbf{u}_k + \mathbf{K}_k^\pi \delta \mathbf{y}_k)\| \leq R_{\text{feas}}.$$

Hence, as  $\|\nabla^2 f_{\text{dyn}}(x, u)\| \leq M_f$  for feasible  $(x, u)$ , Assumption 4.1 implies

$$\begin{aligned} \|\text{remainder}(s)\| &\leq M_f (\|\delta \tilde{\mathbf{y}}_k^{\text{jac}}(s)\|^2 + 2\|\delta \mathbf{u}_k\|^2 + 2L_\pi^2 \|\delta \mathbf{y}_k\|^2) && (\text{AM-GM and Definition 4.7}) \\ &\leq M_f (\|(L_{\text{ol}} L_f \tau \|\delta \mathbf{u}_k\| + L_{\text{ol}}(1 + \tau L_\pi L_f)) \|\delta \mathbf{y}_k\|\|^2 + 2\|\delta \mathbf{u}_k\|^2 + 2L_\pi^2 \|\delta \mathbf{y}_k\|^2) && (\text{ALemma D.5}) \\ &\leq 2M_f ((\tau^2 L_{\text{ol}}^2 L_f^2 \|\delta \mathbf{u}_k\|^2) + L_{\text{ol}}^2 (1 + \tau L_f L_\pi)^2 \|\delta \mathbf{y}_k\|^2 + \|\delta \mathbf{u}_k\|^2 + L_\pi^2 \|\delta \mathbf{y}_k\|^2) && (\text{AM-GM}) \\ &= 2M_f ((1 + \tau^2 L_{\text{ol}}^2 L_f^2) \|\delta \mathbf{u}_k\|^2) + (L_\pi^2 + L_{\text{ol}}^2 (1 + \tau L_f L_\pi)^2) \|\delta \mathbf{y}_k\|^2. \end{aligned}$$

Finally, we conclude by noting that

$$\begin{aligned} &f_{\text{dyn}}(\mathbf{x}^\pi(s), \mathbf{u}_k^\pi) \\ &\quad + \partial_x f_{\text{dyn}}(\mathbf{x}^\pi(s), \mathbf{u}_k^\pi) \delta \tilde{\mathbf{y}}_k^{\text{jac}}(s) + \partial_u f_{\text{dyn}}(\mathbf{x}^\pi(s), \mathbf{u}_k^\pi) (\delta \mathbf{u}_k + \mathbf{K}_k^\pi \delta \mathbf{y}_k) \\ &= \frac{d}{dt} \mathbf{x}^\pi(s) + \mathbf{A}_{\text{ol}}^\pi(s) \delta \tilde{\mathbf{y}}_k^{\text{jac}}(s) + \mathbf{B}_{\text{ol}}^\pi(s) f_{\text{dyn}}(\mathbf{x}^\pi(s), \mathbf{u}_k^\pi) (\delta \mathbf{u}_k + \mathbf{K}_k^\pi \delta \mathbf{y}_k) = \frac{d}{ds} \tilde{\mathbf{y}}_k^{\text{jac}}(s), \end{aligned}$$

so that

$$\begin{aligned} \|\Delta(s)\| = \frac{1}{2} \|\text{remainder}(s)\| &\leq M_f ((1 + \tau^2 L_{\text{ol}}^2 L_f^2) \|\delta \mathbf{u}_k\|^2) + (L_\pi^2 + L_{\text{ol}}^2 (1 + \tau L_f L_\pi)^2) \|\delta \mathbf{y}_k\|^2 \\ &\leq M_f ((1 + 2\tau^2 L_f^2) \|\delta \mathbf{u}_k\|^2) + (L_\pi^2 + 2(1 + \tau L_f L_\pi)^2) \|\delta \mathbf{y}_k\|^2 && (\tau \leq 1/4 L_f, \text{ so } L_{\text{ol}}^2 \leq 2) \\ &\leq M_f (2\|\delta \mathbf{u}_k\|^2) + (L_\pi^2 + 4) \|\delta \mathbf{y}_k\|^2, && (\text{again } \tau \leq 1/4 L_f, \text{ and when as } \tau \leq 4/L_f L_\pi) \\ &\leq M_f (2\|\delta \mathbf{u}_k\|^2) + 5L_\pi^2 \|\delta \mathbf{y}_k\|^2 && (L_\pi \geq 1) \end{aligned}$$

where in the last line, we use  $\tau \leq 1/4L_f$ , so  $L_{o1}^2 \leq 2$ . Hence, from Eq. (D.6), for all  $t \in \mathcal{I}_k$ ,

$$\begin{aligned} \|\tilde{\mathbf{y}}_k^{\text{jac}}(t) - \mathbf{y}(t)\| &\leq \tau L_{o1} M_f M_f (2\|\delta \mathbf{u}_k\|^2) + (L_\pi^2 + 4)\|\delta \mathbf{y}_k\|^2) \\ &\leq 2M_f \tau (2\|\delta \mathbf{u}_k\|^2) + (L_\pi^2 + 4)\|\delta \mathbf{y}_k\|^2 \\ &= M_f \tau (4\|\delta \mathbf{u}_k\|^2) + 10L_\pi^2 \|\delta \mathbf{y}_k\|^2). \end{aligned}$$

where above we bound  $L_{o1} \leq 2$  again. And thus, from Eq. (D.5),

$$\begin{aligned} \|\mathbf{y}(t) - \mathbf{y}^{\text{jac}}(t) - \tilde{\Phi}_{\text{cl},\pi}(t, t_k)(\mathbf{y}(t_k) - \mathbf{y}^{\text{jac}}(t_k))\| &\leq M_f \tau (4\|\delta \mathbf{u}_k\|^2) + 10L_\pi^2 \|\delta \mathbf{y}_k\|^2) ., \\ &\leq M_f \tau (4\|\delta \mathbf{u}_k\|^2) + 10L_\pi^2 \|\|\mathbf{y}_k^{\text{jac}} - \mathbf{y}_k\|^2\|) \\ &\leq 2M_f \tau (4\|\delta \mathbf{u}_k\|^2 + 20L_\pi^2 \|\delta \mathbf{y}_k^{\text{jac}}\|^2 + 20L_\pi^2 \|\mathbf{y}_k^{\text{jac}} - \mathbf{y}_k\|^2). \end{aligned}$$

Substituting in  $\mathbf{e}_k^{\text{jac}} := \mathbf{y}_k - \mathbf{y}_k^{\text{jac}}$  concludes.  $\square$

**Solving the recursion.** To upper bound the recursion, assume an inductive hypothesis that, for some  $R$  to be chosen

$$\max_{j \leq k} \|\mathbf{e}_k^{\text{jac}}\| \leq R, \quad \text{and} \quad \forall j \leq k, \text{ Eq. (D.4) holds.} \quad (\text{D.7})$$

Note this hypothesis is true for  $k = 1$ , where  $\mathbf{e}_1^{\text{jac}} = 0$ , and all terms in Eq. (D.4) coincide with  $(\mathbf{x}_1^\pi, \mathbf{u}_1^\pi)$ , which is feasible. Now assume Eq. (D.7) holds. Define

$$\mathbf{v}_k := \mathbf{e}_{k+1}^{\text{jac}} - \mathbf{A}_{\text{cl},k}^\pi \mathbf{e}_k^{\text{jac}},$$

and note that Lemma D.3, followed by our induction hypothesis, implies that for  $c_1 = 4$  and  $c_2 = 20L_\pi^2$ ,

$$\begin{aligned} \|\mathbf{v}_k\| &\leq \tau M_f (4\|\delta \mathbf{u}_k\|^2 + 20L_\pi^2 \|\mathbf{e}_k^{\text{jac}}\|^2 + c_2 \|\mathbf{y}_k^{\text{jac}}\|^2) \\ &\leq \tau M_f (4\|\delta \mathbf{u}_k\|^2 + 20L_\pi^2 R^2 + 20L_\pi^2 \|\mathbf{y}_k^{\text{jac}}\|^2). \end{aligned}$$

By unfolding the recursion for  $\mathbf{v}_k := \mathbf{e}_{k+1}^{\text{jac}} - \mathbf{A}_{\text{cl},k}^\pi \mathbf{e}_k^{\text{jac}}$ , we have

$$\begin{aligned} \mathbf{e}_{k+1}^{\text{jac}} &= \mathbf{v}_k + \mathbf{A}_{\text{cl},k}^\pi \mathbf{e}_k^{\text{jac}} \\ &= \underbrace{\mathbf{I}}_{=\Phi_{\text{cl},k+1,k+1}^\pi} \mathbf{v}_k + \underbrace{\mathbf{A}_{\text{cl},k}^\pi}_{=\Phi_{\text{cl},k+1,k}^\pi} \mathbf{v}_{k-1} + \underbrace{\mathbf{A}_{\text{cl},k}^\pi \mathbf{A}_{\text{cl},k-1}^\pi}_{=\Phi_{\text{cl},k+1,k-1}^\pi} \mathbf{e}_{k-1}^{\text{jac}} \\ &= \sum_{j=1}^{k+1} \underbrace{\Phi_{\text{cl},k+1,j}^\pi}_{=0} \mathbf{v}_k + \underbrace{\Phi_{\text{cl},k+1,1}^\pi}_{=0} \mathbf{e}_1^{\text{jac}}. \end{aligned}$$

Thus, under our inductive hypothesis, recalling  $B_2 := \tau \|\delta \mathbf{u}_{1:K}\|_{\ell_2}^2$ , and  $B_\infty := \tau \max_k \|\delta \mathbf{u}_k\|^2$ ,

$$\begin{aligned} \|\mathbf{e}_{k+1}^{\text{jac}}\| &\leq M_f \sum_{j=1}^{k+1} \|\Phi_{\text{cl},k+1,j}^\pi\| \tau (c_1 \|\delta \mathbf{u}_k\|^2 + c_2 R^2 + c_2 \|\mathbf{y}_k^{\text{jac}}\|^2) \\ &\leq 4M_f \min\{\kappa_{\pi,\infty} B_2^2, \kappa_{\pi,1} B_\infty^2\} + 20M_f L_\pi^2 \kappa_{\pi,1} (R^2 + \max_k \|\mathbf{y}_k^{\text{jac}}\|^2) \\ &\leq 4M_f \min\{\kappa_{\pi,\infty} B_2^2, \kappa_{\pi,1} B_\infty^2\} + 20M_f L_\pi^2 \kappa_{\pi,1} (R^2 + L_{o1}^2 L_f^2 \min\{B_2 \kappa_{\pi,2}, B_\infty \kappa_{\pi,1}\}^2) \quad (\text{Lemma D.2}) \\ &\leq 4M_f \min\{B_2^2 (\kappa_{\pi,\infty} + 5L_\pi^2 L_{o1}^2 L_f^2 \kappa_{\pi,2}^2 \kappa_{\pi,1}), B_\infty^2 (\kappa_{\pi,1} + 5L_\pi^2 L_{o1}^2 L_f^2 \kappa_{\pi,1}^3)\} + 20M_f L_\pi^2 \kappa_{\pi,1} R^2. \\ &\leq 4M_f \min\{B_2^2 (\kappa_{\pi,\infty} + 10L_\pi^2 L_f^2 \kappa_{\pi,2}^2 \kappa_{\pi,1}), B_\infty^2 (\kappa_{\pi,1} + 10L_\pi^2 L_f^2 \kappa_{\pi,1}^3)\} + 20M_f L_\pi^2 \kappa_{\pi,1} R^2, \end{aligned}$$

where in the last step we use  $\tau \leq 1/4L_f$  to bound  $L_{o1}^2 \leq 2$ . Hence, if we select

$$\begin{aligned} R &= 8M_f \min\{B_2^2 (\kappa_{\pi,\infty} + 10L_\pi^2 L_f^2 \kappa_{\pi,2}^2 \kappa_{\pi,1}), B_\infty^2 (\kappa_{\pi,1} + 10L_\pi^2 L_f^2 \kappa_{\pi,1}^3)\} \\ &:= \min\{B_2^2 M_{\text{tay},2,\pi}, B_\infty^2 M_{\text{tay},\text{inf},\pi}\}, \end{aligned}$$

where we recall

$$\begin{aligned} M_{\text{tay},2,\pi} &:= 8M_f(\kappa_{\pi,\infty} + 10L_\pi^2 L_f^2 \kappa_{\pi,2}^2 \kappa_{\pi,1}) \\ M_{\text{tay},\text{inf},\pi} &:= 8M_f(\kappa_{\pi,1} + 10L_\pi^2 L_f^2 \kappa_{\pi,1}^3), \end{aligned}$$

we get

$$\|\mathbf{e}_{k+1}^{\text{jac}}\| \leq \frac{R}{2} + 20M_f L_\pi^2 \kappa_{\pi,1} R^2.$$

Thus, if  $R \leq \frac{1}{40M_f L_\pi^2 \kappa_{\pi,1}}$ , it holds that

$$\|\mathbf{e}_{k+1}^{\text{jac}}\| \leq \min\{B_2^2 M_{\text{tay},2,\pi}, B_\infty^2 M_{\text{tay},\text{inf},\pi}\}.$$

Lastly, for the condition  $R \leq \frac{1}{40M_f L_\pi^2 \kappa_{\pi,1}}$  to hold, it suffices

$$B_2^2 \leq \frac{1}{40M_f L_\pi^2 \kappa_{\pi,1} M_{\text{tay},2,\pi}}, \quad \text{or } B_\infty^2 \leq \frac{1}{40L_\pi^2 \kappa_{\pi,1} M_{\text{tay},\text{inf},\pi}} \quad (\text{D.8})$$

Notice that these conditions are met for  $B_q^2 \leq B_{\text{tay},q,\pi}^2$ . Moreover,

$$\begin{aligned} \|\mathbf{y}_{k+1} - \mathbf{x}_{k+1}^\pi\| \vee \|\mathbf{y}_{k+1}^{\text{jac}} - \mathbf{x}_{k+1}^\pi\| &\leq \|\mathbf{y}_{k+1} - \mathbf{y}_{k+1}^{\text{jac}}\| + \|\mathbf{y}_{k+1}^{\text{jac}} - \mathbf{x}_{k+1}^\pi\| \\ &= \|\mathbf{e}_{k+1}^{\text{jac}}\| + \left\| \sum_{j=1}^k \Phi_{\text{cl},k+1,j+1}^\pi \mathbf{B}_{\text{ol},j}^\pi \delta \mathbf{u}_j \right\| \\ &\leq \min_{q \in \{2, \infty\}} M_{\text{tay},q,\pi} B_q^2 + \left\| \sum_{j=1}^{k-1} \Phi_{\text{cl},k+1,j+1}^\pi \mathbf{B}_{\text{ol},j}^\pi \delta \mathbf{u}_j \right\| \\ &\leq \min_{q \in \{2, \infty\}} M_{\text{tay},q,\pi} B_q^2 + L_{\text{ol}} L_f \min\{B_2 \kappa_{\pi,2}, B_\infty \kappa_{\pi,1}\} \quad (\text{Lemma D.3}) \\ &\leq \min\{B_2(L_{\text{ol}} L_f \kappa_{\pi,2} + M_{\text{tay},2,\pi} B_2), B_\infty(L_{\text{ol}} L_f \kappa_{\pi,1} + M_{\text{tay},\text{inf},\pi} B_\infty)\} \\ &\leq \min\{B_2(1.5L_f \kappa_{\pi,2} + M_{\text{tay},2,\pi} B_2), B_\infty(1.5L_f \kappa_{\pi,1} + M_{\text{tay},\text{inf},\pi} B_\infty)\} \\ &\quad (L_{\text{ol}} \leq \exp(1/4) \leq 1.5) \end{aligned}$$

Hence,  $B_2 \leq \frac{1}{2} L_f \kappa_{\pi,2} / M_{\text{tay},2,\pi}$  implies  $\|\mathbf{y}_{k+1} - \mathbf{x}_{k+1}^\pi\| \vee \|\mathbf{y}_{k+1}^{\text{jac}} - \mathbf{x}_{k+1}^\pi\| \leq 2L_f \kappa_{\pi,2} B_2 = L_{\text{tay},2,\pi} B_2$ , and  $B_\infty \leq \frac{1}{2} L_f \kappa_{\pi,1} / M_{\text{tay},\text{inf},\pi}$  implies  $\|\mathbf{y}_{k+1} - \mathbf{x}_{k+1}^\pi\| \vee \|\mathbf{y}_{k+1}^{\text{jac}} - \mathbf{x}_{k+1}^\pi\| \leq 2L_f \kappa_{\pi,1} B_\infty = L_{\text{tay},\infty,\pi} B_\infty$ . Combining these conditions with Eq. (D.8) implies we require  $B_q \leq B_{\text{tay},q,\pi}$ , where

$$\begin{aligned} B_{\text{tay},2,\pi} &= \min \left\{ \frac{1}{\sqrt{40M_f L_\pi^2 \kappa_{\pi,1} M_{\text{tay},2,\pi}}}, \frac{L_f \kappa_{\pi,2}}{2M_{\text{tay},2,\pi}} \right\} \\ B_{\text{tay},\text{inf},\pi} &\leq \min \left\{ \frac{1}{40L_\pi^2 \kappa_{\pi,1} M_{\text{tay},\text{inf},\pi}}, \frac{L_f \kappa_{\pi,1}}{2M_{\text{tay},\text{inf},\pi}} \right\} \end{aligned}$$

Lastly, we need to check that the feasibility invariant Eq. (D.4) for  $k = k + 1$  is maintained. Under the above conditions, it was shown that

$$\|\delta \mathbf{y}_{k+1}\| \vee \|\delta \mathbf{y}_{k+1}^{\text{jac}}\| = \|\mathbf{y}_{k+1} - \mathbf{x}_{k+1}^\pi\| \vee \|\mathbf{y}_{k+1}^{\text{jac}} - \mathbf{x}_{k+1}^\pi\| \leq L_{\text{tay},q,\pi} B_q.$$

Hence, if  $B_q \leq \frac{R_{\text{feas}}}{8L_\pi L_{\text{tay},q,\pi}} \leq \frac{R_{\text{feas}}}{4L_{\text{tay},q,\pi}}$  for either  $q \in \{2, \infty\}$ ,

$$\|\mathbf{y}_{k+1}^{\text{jac}}\| \vee \|\mathbf{y}_{k+1}\| \leq \|\mathbf{x}_{k+1}^\pi\| + L_{\text{tay},q,\pi} B_q \leq \frac{R_{\text{feas}}}{2} + L_{\text{tay},q,\pi} B_2 \leq \frac{3}{4} R_{\text{feas}}$$

Moreover, if  $B_q \leq \frac{R_{\text{feas}}}{8L_\pi L_{\text{tay},q,\pi}}$  for either  $q \in \{2, \infty\}$ , and  $B_\infty \leq \frac{R_{\text{feas}}}{8}$ ,

$$\begin{aligned} & \|\mathbf{u}_{k+1}^\pi + \delta \mathbf{u}_{k+1} + \mathbf{K}_{k+1}^\pi (\delta \mathbf{y}_{k+1})\| \vee \|\mathbf{u}_{k+1}^\pi + \delta \mathbf{u}_{k+1} + \mathbf{K}_{k+1}^\pi (\delta \mathbf{y}_{k+1}^{\text{jac}})\| \\ & \leq \|\mathbf{u}_{k+1}^\pi\| + \|\delta \mathbf{u}_{k+1}\| \|\mathbf{K}_{k+1}^\pi\| (\|\delta \mathbf{y}_{k+1}\| \vee \|\delta \mathbf{y}_{k+1}^{\text{jac}}\|) \\ & \leq \frac{R_{\text{feas}}}{2} + \|\delta \mathbf{u}_{k+1}\| + L_\pi L_{\text{tay},q,\pi} B_2 \leq \frac{R_{\text{feas}}}{2} + B_\infty + L_\pi L_{\text{tay},q,\pi} B_2 \leq \frac{3R_{\text{feas}}}{4}. \end{aligned}$$

This concludes the demonstration of Eq. (D.4) for  $k = k+1$ . Collecting our conditions, and recalling  $L_{\text{tay},2,\pi} = 2L_f \kappa_{\pi,2}$ , and  $L_{\text{tay},\infty,\pi} = 2L_f \kappa_{\pi,1}$  we have show that if we take  $B_q \leq B_{\text{tay},q,\pi}$ , where

$$\begin{aligned} B_{\text{tay},2,\pi} &= \min \left\{ \frac{1}{\sqrt{40M_f L_\pi^2 \kappa_{\pi,1} M_{\text{tay},2,\pi}}}, \frac{L_f \kappa_{\pi,2}}{2M_{\text{tay},2,\pi}}, \frac{R_{\text{feas}}}{16L_\pi L_f \kappa_{\pi,2}} \right\} \\ B_{\text{tay},\text{inf},\pi} &= \min \left\{ \frac{1}{40L_\pi^2 \kappa_{\pi,1} M_{\text{tay},\text{inf},\pi}}, \frac{L_f \kappa_{\pi,1}}{2M_{\text{tay},\text{inf},\pi}}, \frac{R_{\text{feas}}}{16L_\pi L_f \kappa_{\pi,1}} \right\} \end{aligned}$$

and  $B_\infty \leq \frac{R_{\text{feas}}}{8}$ , then

$$\|\mathbf{y}_{k+1} - \mathbf{x}_{k+1}^\pi\| \vee \|\mathbf{y}_{k+1}^{\text{jac}} - \mathbf{x}_{k+1}^\pi\| \leq L_{\text{tay},q,\pi} B_2,$$

and

$$\|\mathbf{e}_{k+1}^{\text{jac}}\| = \|\mathbf{y}_{k+1} - \mathbf{y}_{k+1}^{\text{jac}}\| \leq M_{\text{tay},q,\pi} B_q^2.$$

In addition, we have show that  $\|\mathbf{y}_{k+1}\| \vee \|\mathbf{u}_{k+1}^\pi + \delta \mathbf{u}_{k+1} + \mathbf{K}_{k+1}^\pi (\delta \mathbf{y}_{k+1})\| \leq \frac{3}{4} R_{\text{feas}}$ . This concludes the induction.

Substituting in  $\mathbf{y}_k = \tilde{\mathbf{x}}_k^\pi(\mathbf{u}_{1:K})$  and using the computation of  $\mathbf{y}_{k+1}^{\text{jac}}$  in Lemma D.2 concludes the proof of the perturbation bounds. Moreover, the fact that Eq. (D.4) holds for all  $k$ , and consequently the conclusion of Lemma D.4 establishes the norm bounds  $\|\tilde{\mathbf{x}}_k^\pi(\mathbf{u}_{1:K})\| \vee \|\tilde{\mathbf{u}}_k^\pi(\mathbf{u}_{1:K})\| \leq \frac{3R_{\text{feas}}}{4}$ , and  $\|\tilde{\mathbf{x}}^\pi(t | \mathbf{u}_{1:K})\| \leq R_{\text{feas}}$ .  $\square$

## D.2. Taylor Expansion of the Cost (Lemma A.6)

*Proof.* Recall the definition

$$\mathcal{J}_T^{\pi, \text{disc}}(\bar{\mathbf{u}}) := V(\tilde{\mathbf{x}}_{K+1}^\pi(\bar{\mathbf{u}})) + \tau \sum_{k=1}^K Q(\tilde{\mathbf{x}}_k^\pi(\bar{\mathbf{u}}), \tilde{\mathbf{u}}_k^\pi(\bar{\mathbf{u}}), t_k).$$

Define the shorthand

$$\begin{aligned} \delta \tilde{\mathbf{x}}_k &:= \tilde{\mathbf{x}}_k^\pi(\delta \mathbf{u}_{1:K} + \mathbf{u}_{1:K}^\pi) - \mathbf{x}_k^\pi \\ \delta \tilde{\mathbf{x}}_k^{\text{jac}} &:= \tilde{\mathbf{x}}_k^{\pi, \text{jac}}(\delta \mathbf{u}_{1:K} + \mathbf{u}_{1:K}^\pi) \end{aligned}$$

$$\begin{aligned} & \mathcal{J}_T^{\pi, \text{disc}}(\bar{\mathbf{u}}) - \mathcal{J}_T^{\pi, \text{disc}}(\delta \mathbf{u}_{1:K} + \mathbf{u}_{1:K}^\pi) \\ & := V(\mathbf{x}_{K+1}^\pi + \delta \tilde{\mathbf{x}}_{K+1}) - V(\mathbf{x}_{K+1}^\pi) + \tau \sum_{k=1}^K Q(\mathbf{x}_k^\pi + \delta \tilde{\mathbf{x}}_k, \delta \mathbf{u}_k + \mathbf{K}_k^\pi \delta \tilde{\mathbf{x}}_k + \mathbf{u}_k^\pi, t_k) - Q(\mathbf{x}_k^\pi, \mathbf{u}_k^\pi, t_k) \end{aligned}$$

Notice that Proposition A.5 and feasibility of  $\pi$  implies that

$$\|\mathbf{x}_k^\pi + \delta \tilde{\mathbf{x}}_k\| \|\mathbf{x}_k^\pi\| \vee \|\delta \mathbf{u}_k + \mathbf{K}_k^\pi \mathbf{e}_k\| \vee \|\mathbf{u}_k^\pi\| \leq R_{\text{feas}} \quad (\text{D.9})$$

Hence a Taylor expansion and [Assumption 4.2](#) imply

$$\begin{aligned}
 & |V(\mathbf{x}_{K+1}^\pi + \delta\tilde{\mathbf{x}}_{K+1}) - V(\mathbf{x}_{K+1}^\pi) - \langle \partial_x V(\mathbf{x}_{K+1}^\pi), \delta\tilde{\mathbf{x}}_{K+1}^{\text{jac}} \rangle| \\
 & \leq \frac{1}{2} \sup_{\alpha \in [0,1]} \|\partial_{xx} V(\mathbf{x}_{K+1}^\pi + \alpha\delta\tilde{\mathbf{x}}_{K+1})\| \|\delta\tilde{\mathbf{x}}_{K+1}\|^2 + |\langle \partial_x V(\mathbf{x}_{K+1}^\pi), \delta\tilde{\mathbf{x}}_{K+1} - \delta\tilde{\mathbf{x}}_{K+1}^{\text{jac}} \rangle| \\
 & \leq \frac{M_{\text{cost}}}{2} \|\delta\tilde{\mathbf{x}}_{K+1}\|^2 + L_{\text{cost}} \|\delta\tilde{\mathbf{x}}_{K+1} - \delta\tilde{\mathbf{x}}_{K+1}^{\text{jac}}\|.
 \end{aligned}$$

Similarly,

$$\begin{aligned}
 & \left| Q(\mathbf{x}_k^\pi + \delta\tilde{\mathbf{x}}_k, \delta\mathbf{u}_k + \mathbf{K}_k^\pi \delta\tilde{\mathbf{x}}_k, t_k) - Q(\mathbf{x}_k^\pi, \mathbf{u}_k^\pi, t_k) - \langle \partial_x Q(\mathbf{x}_k^\pi, \mathbf{u}_k^\pi, t_k), \delta\tilde{\mathbf{x}}_k^{\text{jac}} \rangle - \langle \partial_u Q(\mathbf{x}_k^\pi, \mathbf{u}_k^\pi, t_k), \delta\mathbf{u}_k + \mathbf{K}_k \delta\tilde{\mathbf{x}}_k^{\text{jac}} \rangle \right| \\
 & \leq \left| \langle \partial_x Q(\mathbf{x}_k^\pi, \mathbf{u}_k^\pi, t_k), \delta\tilde{\mathbf{x}}_k^{\text{jac}} \rangle \right| + \left| \langle \partial_u Q(\mathbf{x}_k^\pi, \mathbf{u}_k^\pi, t_k), \mathbf{K}_k \delta\tilde{\mathbf{x}}_k^{\text{jac}} \rangle \right| \\
 & \quad + \frac{1}{2} M_{\text{cost}} \left( \|\delta\tilde{\mathbf{x}}_k\|^2 + \|\delta\mathbf{u}_k + \mathbf{K}_k \delta\tilde{\mathbf{x}}_k^{\text{jac}}\|^2 \right) \\
 & \leq L_{\text{cost}} (1 + L_\pi) \|\delta\tilde{\mathbf{x}}_k - \delta\tilde{\mathbf{x}}_k^{\text{jac}}\| + \frac{1}{2} M_{\text{cost}} \left( (1 + 2L_\pi^2) \|\delta\tilde{\mathbf{x}}_k^{\text{jac}}\|^2 + 2\|\delta\mathbf{u}_k\|^2 \right) \\
 & \leq 2L_\pi L_{\text{cost}} \|\delta\tilde{\mathbf{x}}_k - \delta\tilde{\mathbf{x}}_k^{\text{jac}}\| + \frac{1}{2} M_{\text{cost}} (3L_\pi^2 \|\mathbf{e}_k\|^2 + 2\|\delta\mathbf{u}_k\|^2). \tag{$L_\pi \geq 1$}
 \end{aligned}$$

Then, from [Lemma C.7](#),

$$\begin{aligned}
 & \langle \delta\mathbf{u}_{1:K}, \nabla \mathcal{J}_T^{\pi, \text{disc}}(\mathbf{u}) \rangle \\
 & = \langle \partial_x V(\mathbf{x}_{K+1}^\pi), \delta\tilde{\mathbf{x}}_K^{\text{jac}} \rangle + \tau \sum_{k=1}^K \langle \partial_x Q(\mathbf{x}_k^\pi, \mathbf{u}_k^\pi, t_k), \delta\tilde{\mathbf{x}}_k^{\text{jac}} \rangle + \langle \partial_u Q(\mathbf{x}_k^\pi, \mathbf{u}_k^\pi, t_k), \delta\mathbf{u}_k + \mathbf{K}_k \delta\tilde{\mathbf{x}}_k^{\text{jac}} \rangle.
 \end{aligned}$$

Therefore, we have

$$\begin{aligned}
 & \|\mathcal{J}_T^{\pi, \text{disc}}(\delta\mathbf{u}_{1:K} + \mathbf{u}_{1:K}^\pi) - \mathcal{J}_T^{\pi, \text{disc}}(\mathbf{u}_{1:K}^\pi) - \langle \delta\mathbf{u}_{1:K}, \nabla \mathcal{J}_T^{\pi, \text{disc}}(\pi) \rangle\| \\
 & \leq \frac{M_{\text{cost}}}{2} \|\delta\tilde{\mathbf{x}}_{K+1}\|^2 + L_{\text{cost}} \|\delta\tilde{\mathbf{x}}_{K+1}^{\text{jac}} - \delta\tilde{\mathbf{x}}_{K+1}\| + \tau \sum_{k=1}^K 2L_\pi L_{\text{cost}} \|\delta\tilde{\mathbf{x}}_k^{\text{jac}} - \delta\tilde{\mathbf{x}}_k\| + \frac{1}{2} M_{\text{cost}} (3L_\pi^2 \|\delta\tilde{\mathbf{x}}_k\|^2 + 2\|\delta\mathbf{u}_k\|^2) \\
 & \leq \frac{M_{\text{cost}}}{2} (1 + 3L_\pi^2 T) \max_{k \in [K+1]} \|\delta\tilde{\mathbf{x}}_k\|^2 + L_{\text{cost}} (1 + 2L_\pi T) \max_{k \in [K+1]} \|\delta\tilde{\mathbf{x}}_k^{\text{jac}} - \delta\tilde{\mathbf{x}}_k\| + L_\pi L_{\text{cost}} \underbrace{\tau \sum_{t=1}^T \|\delta\mathbf{u}_k\|^2}_{=B_2} \\
 & = \frac{M_{\text{cost}}}{2} (1 + 3L_\pi^2 T) \max_{k \in [K+1]} \|\delta\tilde{\mathbf{x}}_k\|^2 + L_{\text{cost}} (1 + 2L_\pi T) \max_{k \in [K+1]} \|\delta\tilde{\mathbf{x}}_k - \delta\tilde{\mathbf{x}}_k^{\text{jac}}\| + 2L_\pi L_{\text{cost}} B_2^2.
 \end{aligned}$$

From [Proposition A.5](#),

$$\begin{aligned}
 \max_{k \in [K+1]} \|\delta\tilde{\mathbf{x}}_k - \delta\tilde{\mathbf{x}}_k\| & \leq M_{\text{tay}, 2, \pi} B_2^2 \\
 \max_{k \in [K+1]} \|\mathbf{e}_k\|^2 & \leq 4L_f^2 \kappa_{\pi, 2}^2 B_2^2.
 \end{aligned}$$

Thus,

$$\begin{aligned}
 & \|\mathcal{J}_T^{\pi, \text{disc}}(\delta\mathbf{u}_{1:K} + \mathbf{u}_{1:K}^\pi) - \mathcal{J}_T^{\pi, \text{disc}}(\mathbf{u}_{1:K}^\pi) - \langle \delta\mathbf{u}_{1:K}, \nabla \mathcal{J}_T^{\pi, \text{disc}}(\pi) \rangle\| \\
 & \leq \underbrace{(2M_{\text{cost}} L_f^2 \kappa_{\pi, 2}^2 (1 + 3L_\pi^2 T) M_{\text{tay}, 2, \pi} + L_{\text{cost}} (1 + 2L_\pi T) M_{\text{tay}, 2, \pi} + 2L_\pi L_{\text{cost}})}_{:=M_{\mathcal{J}, \text{tay}, \pi}} B_2^2.
 \end{aligned}$$

□

### D.3. Proof of Lemma A.7

We begin with the following lemma, which we show shortly below.

**Lemma D.6.** *Consider the setting of Proposition A.5, with  $B_\infty \leq \min\{B_{\text{tay,inf},\pi}, R_{\text{feas}}/8\}$ . Let  $\pi'$  denote the policy with gains  $K_k^\pi$  and inputs  $\mathbf{u}_k^{\pi'} = \mathbf{u}_k^\pi + \delta \mathbf{u}_k$ . Then,*

$$\tau \sum_{k=1}^K \|\mathbf{A}_{\text{cl},k}^{\pi'} - \mathbf{A}_{\text{cl},k}^\pi\| \leq 12TM_f L_\pi (1 + L_f K_\pi) B_\infty.$$

With this lemma, we turn to the proof of Lemma A.7. Notice that, as  $\pi$  and  $\pi'$  have the same gains,  $L_\pi = L_{\pi'}$ . Therefore, following the proof of Lemma A.1 (see, specifically, the proof of Claim G.7), we have that for  $\|\mathbf{A}_{\text{cl},k}^{\pi'} - \mathbf{I}\|$  and  $\|\mathbf{A}_{\text{cl},k}^\pi - \mathbf{I}\|$  are both at most  $\kappa := 3\tau L_f L_\pi$  for  $\tau \leq 1/6L_f L_\pi$ .

Let us now construct an interpolating curves  $\mathbf{X}_k(s)$  with  $\mathbf{X}_k(0) = \mathbf{A}_{\text{cl},k}^\pi$  and  $\mathbf{X}_k(1) = \mathbf{A}_{\text{cl},k}^{\pi'}$ , and define the interpolating Lyapunov function

$$\Lambda_{K+1}(s) = \mathbf{I}, \quad \Lambda_k(s) = \mathbf{X}_k(s)^\top \Lambda_{k+1} \mathbf{X}_k(s) + \tau \mathbf{I},$$

Define

$$\begin{aligned} \Delta &= 3 \sup_{s \in [0,1]} \max\{1, 2\kappa\} \sum_{k=k_0}^K \|\mathbf{X}'_k(s)\| \\ &= \max\{1, 2\kappa\} \sum_{k=k_0}^K \|\mathbf{A}_{\text{cl},k}^{\pi'} - \mathbf{A}_{\text{cl},k}^\pi\| \\ &= \max\{3, 18L_f L_\pi\} \sum_{k=k_0}^K \|\mathbf{A}_{\text{cl},k}^{\pi'} - \mathbf{A}_{\text{cl},k}^\pi\|. \end{aligned}$$

and recall the shorthand  $\|\Lambda_{k_0:K+1}(s)\|_{\max, \text{op}} := \max_{k \in [k_0:K+1]} \|\Lambda_k(s)\|$ . Then, as long as  $\|\Lambda_{k_0:K+1}(0)\|_{\max, \text{op}} \Delta < 1$ , Lemma F.13 (re-indexing to terminate the backward recursion at  $k_0$  instead of 1) implies

$$\|\Lambda_{k_0:K+1}(1)\|_{\max, \text{op}} \leq (1 - \|\Lambda_{k_0:K+1}(0)\|_{\max, \text{op}} \Delta)^{-1} \|\Lambda_{k_0:K+1}(0)\|_{\max, \text{op}}.$$

We see that  $\|\Lambda_{k_0:K+1}(0)\|_{\max, \text{op}} = \|\Lambda_{k_0:K+1}^\pi\|_{\max, \text{op}} = K_{\pi, \star}$ , and  $\|\Lambda_{k_0:K+1}(1)\|_{\max, \text{op}} = \|\Lambda_{k_0:K+1}^{\pi'}\|_{\max, \text{op}} = \mu_{\pi', \star}$ . Thus, combining with the inequality  $(1-x)^{-1} \leq 1+2x$  for  $x \in [0, 1/2]$ , we have that as long as  $\Delta \mu_{\pi, \star} \leq 1/2$ ,

$$\mu_{\pi', \star} \leq (1 + 2\Delta \mu_{\pi, \star}) \mu_{\pi, \star}.$$

Lastly, we can bound

$$\begin{aligned} 2\Delta \mu_{\pi, \star} &= \max\{6, 36L_f L_\pi\} \mu_{\pi, \star} \sum_{k=k_0}^K \|\mathbf{A}_{\text{cl},k}^{\pi'} - \mathbf{A}_{\text{cl},k}^\pi\| \\ &\leq \underbrace{\max\{6, 36L_f L_\pi\} \mu_{\pi, \star} \cdot 12TM_f L_\pi (1 + L_f K_\pi) B_\infty}_{:= 1/B_{\text{stab}, \pi}} B_2. \end{aligned} \quad (\text{Lemma D.6})$$

In sum, for  $B_\infty \leq B_{\text{stab}, \pi}$ , we have  $\mu_{\pi', \star} \leq (1 + B_\infty/B_{\text{stab}, \pi}) \mu_{\pi, \star}$ , which concludes the proof.

*Proof of Lemma D.6.* Due to Proposition A.5, and the fact that  $\mathbf{x}^{\pi'}(t) = \tilde{\mathbf{x}}^\pi(t | \mathbf{u}_{1:K})$  and  $\mathbf{u}^{\pi'}(t) = \tilde{\mathbf{u}}_{k(t)}^\pi(\mathbf{u}_{1:K})$ , we have that

$$\forall t \in [0, T], \quad \|\mathbf{x}^{\pi'}(t)\| \vee \|\mathbf{u}^{\pi'}(t)\| \leq R_{\text{feas}}. \quad (\text{D.10})$$

Moreover each initial condition  $\xi$  with norm  $\|\xi\| = 1$ , we have that from [Lemma C.10](#) and the definitions of  $\mathbf{A}_{\text{cl},k}^{\pi'}$ ,  $\mathbf{A}_{\text{cl},k}^{\pi}$  from [Definition C.8](#) that

$$(\mathbf{A}_{\text{cl},k}^{\pi'} - \mathbf{A}_{\text{cl},k}^{\pi})\xi = \mathbf{z}_2(\tau) - \mathbf{z}_1(\tau),$$

where

$$\mathbf{z}_2(0) = \mathbf{z}_1(0) = \xi,$$

and where  $\frac{d}{dt}\mathbf{z}_2(t) = \mathbf{A}_{\text{ol}}^{\pi'}(t_k + t)\mathbf{z}_2(t) + \mathbf{B}_{\text{ol}}^{\pi'}(t_k + t)\mathbf{K}_k\xi$ , and where  $\frac{d}{dt}\mathbf{z}_1(t) = \mathbf{A}_{\text{ol}}^{\pi}(t_k + t)\mathbf{z}_1(t) + \mathbf{B}_{\text{ol}}^{\pi}(t_k + t)\mathbf{K}_k\xi$ .

By the Picard Lemma, [Lemma C.9](#), and by bounding  $\|\mathbf{A}_{\text{ol}}^{\pi'}(t)\| \leq L_f$  by [Eq. \(D.10\)](#) and [Assumption 4.1](#), it follows that

$$\begin{aligned} \|(\mathbf{A}_{\text{cl},k}^{\pi'} - \mathbf{A}_{\text{cl},k}^{\pi})\xi\| &\leq \exp(\tau L_f) \int_0^{\tau} (\|\mathbf{A}_{\text{ol}}^{\pi'}(t_k + t) - \mathbf{A}_{\text{ol}}^{\pi}(t_k + t)\| \|\mathbf{z}_1(t)\| + \|\mathbf{B}_{\text{ol}}^{\pi'}(t_k + t) - \mathbf{B}_{\text{ol}}^{\pi}(t_k + t)\| \|\mathbf{K}_k\| \|\xi\|) dt \\ &\leq \exp(\tau L_f) \int_0^{\tau} (\|\mathbf{A}_{\text{ol}}^{\pi'}(t_k + t) - \mathbf{A}_{\text{ol}}^{\pi}(t_k + t)\| \|\mathbf{z}_1(t)\| + \|\mathbf{B}_{\text{ol}}^{\pi'}(t_k + t) - \mathbf{B}_{\text{ol}}^{\pi}(t_k + t)\| L_{\pi}) dt \end{aligned}$$

Set  $L_{\text{ol}} = \exp(\tau L_f)$ . Following the computation in , we can bound  $\sup_{t \in [0, \tau]} \|\mathbf{z}_1(t)\| \leq \|\mathbf{A}_{\text{cl},k}^{\pi}\| = \|\Phi_{\text{cl},k+1,k}^{\pi}\| \leq 5/3$  provided  $\tau \leq 1/4L_fL_{\pi}$  (recall we assume  $L_{\pi} \geq 1$ ). Hence,

$$\|(\mathbf{A}_{\text{cl},k}^{\pi'} - \mathbf{A}_{\text{cl},k}^{\pi})\xi\| \leq L_{\text{ol}} \int_0^{\tau} \left( \frac{5}{3} \|\mathbf{A}_{\text{ol}}^{\pi'}(t_k + t) - \mathbf{A}_{\text{ol}}^{\pi}(t_k + t)\| + \|\mathbf{B}_{\text{ol}}^{\pi'}(t_k + t) - \mathbf{B}_{\text{ol}}^{\pi}(t_k + t)\| L_{\pi} \right) dt$$

Finally, by the smoothness on the dynamics [Assumption 4.1](#) and invoking [Eq. \(D.10\)](#) and feasibility of  $\pi$  to ensure all relevant  $(x, u)$  pairs are feasible, we have

$$\begin{aligned} \|\mathbf{A}_{\text{ol}}^{\pi'}(t_k + t) - \mathbf{A}_{\text{ol}}^{\pi}(t_k + t)\| &= \|\partial_x f_{\text{dyn}}(\mathbf{x}^{\pi'}(t), \mathbf{u}^{\pi'}(t)) - \partial_x f_{\text{dyn}}(\mathbf{x}^{\pi}(t), \mathbf{u}^{\pi}(t))\| \\ &\leq M_f \left( \|\mathbf{x}^{\pi'}(t) - \mathbf{x}^{\pi}(t)\| + \|\mathbf{u}^{\pi'}(t) - \mathbf{u}^{\pi}(t)\| \right) \\ &\leq M_f \left( \|\mathbf{x}^{\pi'}(t) - \mathbf{x}^{\pi}(t)\| + \|\delta \mathbf{u}_k\| \right). \end{aligned}$$

Applying a similar bound to the term  $\|\mathbf{B}_{\text{ol}}^{\pi'}(t_k + t) - \mathbf{B}_{\text{ol}}^{\pi}(t_k + t)\|$ , we conclude

$$\begin{aligned} \|\mathbf{A}_{\text{cl},k}^{\pi'} - \mathbf{A}_{\text{cl},k}^{\pi}\| &\leq \sup_{\xi: \|\xi\|=1} \|(\mathbf{A}_{\text{cl},k}^{\pi'} - \mathbf{A}_{\text{cl},k}^{\pi})\xi\| \\ &\leq L_{\text{ol}} M_f \left( \frac{5}{3} + L_{\pi} \right) \int_0^{\tau} \left( \|\mathbf{x}^{\pi'}(t) - \mathbf{x}^{\pi}(t)\| + \|\delta \mathbf{u}_k\| \right) dt \\ &\leq \tau L_{\text{ol}} M_f \left( \frac{5}{3} + L_{\pi} \right) \left( (1 + \tau L_{\text{ol}} L_f) \|\delta \mathbf{u}_k\| + L_{\text{ol}} (1 + \tau L_{\pi} L_f) \|\mathbf{x}^{\pi'}(t_k) - \mathbf{x}^{\pi}(t_k)\| \right) \quad (\text{Lemma D.5}) \\ &\leq \tau L_{\text{ol}} M_f \left( \frac{5}{3} + L_{\pi} \right) \left( \frac{5L_{\text{ol}}}{4} \|\delta \mathbf{u}_k\| + \frac{5}{4} L_{\text{ol}} \|\mathbf{x}^{\pi'}(t_k) - \mathbf{x}^{\pi}(t_k)\| \right) \quad (L_{\text{ol}} \geq 1, \tau \leq 1/4L_fL_{\pi} \leq 1/4L_f) \\ &= \tau \frac{5L_{\text{ol}}^2}{4} M_f \left( \frac{5}{3} + L_{\pi} \right) \left( \|\delta \mathbf{u}_k\| + \|\mathbf{x}^{\pi'}(t_k) - \mathbf{x}^{\pi}(t_k)\| \right) \\ &= \tau M_f \cdot \frac{5e^{1/2}}{4} \cdot \frac{8}{3} \cdot L_{\pi} \left( \|\delta \mathbf{u}_k\| + \|\mathbf{x}^{\pi'}(t_k) - \mathbf{x}^{\pi}(t_k)\| \right) \quad (L_{\pi} \geq 1, \tau L_f \leq 1/4) \\ &\leq 6\tau M_f L_{\pi} \left( \|\delta \mathbf{u}_k\| + \|\mathbf{x}^{\pi'}(t_k) - \mathbf{x}^{\pi}(t_k)\| \right) \\ &\leq 6\tau M_f L_{\pi} (\|\delta \mathbf{u}_k\| + 2L_f \kappa_{\pi,1} B_{\infty}). \quad (\text{Proposition A.5}) \end{aligned}$$

Summing the bound, and using  $K\tau = T$ , we have

$$\begin{aligned} \sum_{k=1}^K \|\mathbf{A}_{\text{cl},k}^{\pi'} - \mathbf{A}_{\text{cl},k}^{\pi}\| &\leq 6M_f L_{\pi} \left( \tau \sum_{k=1}^K \|\delta \mathbf{u}_k\| + 2TL_f \kappa_{\pi,1} B_{\infty} \right) \\ &\leq 6M_f L_{\pi} (K\tau B_{\infty} + 2TL_f \kappa_{\pi,1} B_{\infty}) \\ &\leq 12TM_f L_{\pi} (1 + L_f \kappa_{\pi,1}) B_{\infty}. \end{aligned}$$

□



#### D.4. Proof of Lemma A.8

*Proof.* Using Condition F.3 and  $1 \vee \|K_k\| \leq L_\pi$ ,

$$\begin{aligned} \|(\nabla \mathcal{J}_T^{\text{disc}}(\pi))_k\| &\leq \tau \|Q_u(\mathbf{x}_k^\pi, \mathbf{u}_k^\pi, t_k)^2\|^2 + \|\Psi_{\text{cl}, K+1, k}^\pi\|^\top V_x(\mathbf{x}_{K+1}^\pi) \\ &\quad + \left\| \tau \sum_{j=k+1}^K (\Psi_{\text{cl}, j, k}^\pi)^\top (Q_x(\mathbf{x}_j^\pi, \mathbf{u}_j^\pi, t_j) + (K_j^\pi)^\top Q_u(\mathbf{x}_j^\pi, \mathbf{u}_j^\pi, t_j)) \right\| \\ &\leq \tau L_{\text{cost}} + \tau \|\Psi_{\text{cl}, K+1, k}^\pi\| L_{\text{cost}}^2 + 2L_\pi^2 L_{\text{cost}}^2 \left( \tau \sum_{j=k+1}^K \|\Psi_{\text{cl}, j, k}^\pi\| \right)^2. \end{aligned}$$

Using  $\tau \leq 1/4L_f$  and Lemma I.3, we can bound  $\|\Psi_{\text{cl}, j, k}^\pi\| = \|\Phi_{\text{cl}, j, k+1}^\pi \mathbf{B}_{\text{ol}, k}^\pi\| \leq \tau L_f \exp(\tau L_f) \|\Phi_{\text{cl}, j, k+1}^\pi\| \leq \tau L_f \exp(1/4) \|\Phi_{\text{cl}, j, k+1}^\pi\|$ . As  $\exp(1/2) \leq 3/4$ , we conclude that

$$\begin{aligned} \|(\nabla \mathcal{J}_T^{\text{disc}}(\pi))_k\| &\leq \tau L_{\text{cost}} + \frac{3}{2} \tau^2 L_f L_{\text{cost}} + 3\tau L_\pi L_{\text{cost}} \left( \tau \sum_{j=k+1}^K \|\Phi_{\text{cl}, j, k}^\pi\| \right)^2 \\ &\leq \tau L_{\text{cost}} + \frac{3}{2} \tau L_f L_{\text{cost}} \|\Phi_{\text{cl}, K+1, k+1}^\pi\| + 3L_\pi L_{\text{cost}} \kappa_{\pi, 1}. \end{aligned}$$

Using  $\|\Phi_{\text{cl}, K+1, k+1}^\pi\| \leq \kappa_{\pi, \infty}$  gives the bound  $\|(\nabla \mathcal{J}_T^{\text{disc}}(\pi))_k\| \leq \tau L_{\text{cost}} (1 + \frac{3L_f}{2} \kappa_{\pi, \infty} + 3L_\pi \kappa_{\pi, 1}) =: L_{\nabla, \pi, \infty}$ .  $\square$

### E. Estimation Proofs

#### E.1. Estimation of Markov Parameters: Proof of Proposition A.9

We begin with two standard concentration inequalities.

**Lemma E.1.** *Let  $(y_i, x_i, w_i)_{i=1}^n$  be an independent sequence of triples of random vectors in with  $y_i, x_i \in \mathbb{R}^d$ ,  $w \in \mathbb{R}^{d'}$  and suppose that  $y_i \mid x_i, w_i \sim \mathcal{N}(x_i, \sigma^2 \mathbf{I}_d)$  and  $\max_i \|w_i\| \leq R$ . Then,*

$$\mathbb{P} \left[ \left\| \frac{1}{N} \sum_{i=1}^N (y_i - x_i) w_i \right\| \leq R\sigma \sqrt{2 \cdot \frac{d \log 5 + \log((d' + 1)/\delta)}{N}} \right] \geq 1 - \delta.$$

*Proof of Lemma E.1.* By a standard covering argument (see, e.g. Vershynin (2018, Chapter 4)), there exists a finite covering  $\mathcal{T}$  of unit vectors  $z \in \mathbb{R}^d$  such that (a)  $\log |\mathcal{T}| \leq d \log 5$  and (b), for all vectors  $v \in \mathbb{R}^d$ ,

$$\|v\| \leq 2 \sup_{z \in \mathcal{T}} \langle v, z \rangle.$$

Hence,

$$\left\| \frac{1}{N} \sum_{i=1}^N (y_i - x_i) w_i \right\|_{\text{op}} \leq 2 \sup_{z \in \mathcal{T}} \left\| \frac{1}{N} \sum_{i=1}^N \langle z, y_i - x_i \rangle \cdot w_i \right\| = 2\sigma \sup_{z \in \mathcal{T}} \left\| \frac{1}{N} \sum_{i=1}^N \xi_i(z) \cdot w_i \right\|, \quad (\text{E.1})$$

where above we define  $\xi_i(z) := \sigma^{-1} \langle z, y_i - x_i \rangle$ . Notice that  $\xi_i(z) \mid w_i$  are standard Normal random variables. Thus, by standard Gaussian concentration (e.g. Boucheron et al. (2013, Chapter 2)),

$$\mathbb{P} \left[ \left\| \frac{1}{N} \sum_{i=1}^N \xi_i(z) w_i \right\| \leq R\sigma \sqrt{2 \frac{\log((d' + 1)/\delta)}{N}} \right] \geq 1 - \delta.$$

Hence, union bounding over  $z \in \mathcal{T}$ , bounding  $|\mathcal{T}| \leq d \log 5$ , Eq. (E.1) implies the desired bound.

$$\mathbb{P} \left[ \left\| \frac{1}{N} \sum_{i=1}^N (y_i - x_i) w_i \right\| \leq R\sigma \sqrt{2 \frac{d \log 5 + \log((d' + 1)/\delta)}{N}} \right] \geq 1 - \delta.$$

$\square$

**Lemma E.2** (Asymmetric Matrix Hoeffding). *Let  $X_1, \dots, X_n$  be an independent sequence of matrices in  $\mathbb{R}^{d_1 \times d_2}$  with  $\|X_i\| \leq R$ . Then,*

$$\mathbb{P} \left[ \frac{1}{N} \left\| \sum_i X_i - \mathbb{E}[X_i] \right\| \leq 4R \sqrt{\frac{\log(\frac{d_1+d_2}{\delta})}{N}} \right] \geq 1 - \delta.$$

*Proof of Lemma E.2.* By recentering  $X_i \leftarrow X_i - \mathbb{E}[X_i]$ , we may assume  $\mathbb{E}[X_i] = 0$  and  $\|X_i\| \leq 2R$ . Define the Hermitian dilation

$$Y_i = \begin{bmatrix} 0 & X_i \\ X_i^\top & 0 \end{bmatrix}.$$

Then

$$Y_i^2 = \begin{bmatrix} X_i X_i^\top & 0 \\ 0 & X_i^\top X_i \end{bmatrix} \preceq \|X_i\|^2 \mathbf{I}_{d_1+d_i} \leq 4R^2 \mathbf{I}_{d_1+d_2}$$

Applying standard Matrix Hoeffding [Tropp \(2012, Theorem 1.4\)](#) for Hermitian matrices to the  $Y_i$ 's yields

$$\mathbb{P} \left[ \left\| \sum_i Y_i \right\| \geq t \right] \leq (d_1 + d_2) e^{-\frac{t^2}{32NR^2}}.$$

Hence, by rearranging,

$$\mathbb{P} \left[ \left\| \sum_i Y_i \right\| \leq 4R \sqrt{2N \log(\frac{d_1+d_2}{\delta})} \right] \geq 1 - \delta$$

As  $\|\sum_i Y_i\| = \sqrt{2} \|\sum_i X_i\|$ , we conclude

$$\mathbb{P} \left[ \frac{1}{N} \left\| \sum_i X_i \right\| \leq 4R \sqrt{\frac{\log(\frac{d_1+d_2}{\delta})}{N}} \right] \geq 1 - \delta.$$

□

We now turn to concluding the proof of [Proposition A.9](#). We begin with a claim which bounds  $\max_k \|\hat{\mathbf{x}}_k - \mathbf{x}_k^\pi\|$ . Throughout,  $d_\star := \max\{d_x, d_u\}$ .

**Claim E.1.** *With probability at least  $1 - \delta/3$ , the following bound holds*

$$\max_k \|\hat{\mathbf{x}}_k - \mathbf{x}_k^\pi\| \leq \sigma_{\text{orac}} \sqrt{2 \frac{d_\star \log 5 + \log(6(K+1)/\delta)}{N}} \leq \text{Err}_{\hat{\mathbf{x}}}(\delta).$$

*Proof.* The result follows directly from [Lemma E.1](#), with  $w_i = 1 \in \mathbb{R}$  for each  $i$ . □

*Proof of Proposition A.9.* Throughout, suppose the event of [Claim E.1](#) holds. We also note that

$$\|\mathbf{w}_j^{(i)}\| \leq \sigma_w \sqrt{d_u} \leq \sigma_w \sqrt{d_\star} \text{ a.s.} \tag{E.2}$$

This covers the first inequality of the proposition. To bound the error on the transition operator, let us fix indices  $j, k$ ; we perform a union bound at the end of the proof. For each perturbation sampled perturbation  $\mathbf{w}_{1:K}^{(i)}$ , define a perturbed control input

$$\check{\mathbf{u}}_k^{(i)} = \mathbf{u}_k^\pi + \mathbf{w}_k^{(i)} + K_k^\pi (\mathbf{x}_k^\pi - \hat{\mathbf{x}}_k), \quad \hat{\mathbf{u}}_k^i = \mathbf{u}_k^\pi + \mathbf{w}_k^{(i)} - K_k^\pi \hat{\mathbf{x}}_k.$$

and observe that

$$\mathbf{x}_{\text{orac},k}^\pi(\tilde{\mathbf{u}}_{1:K}^{(i)}) = \tilde{\mathbf{x}}_k(\tilde{\mathbf{u}}_{1:K}^{(i)}), \quad \forall k \in [K]. \quad (\text{E.3})$$

Hence, we have that  $\mathbf{y}_k$  defined in [Line 7](#) satisfies

$$\mathbf{y}_k^{(i)} \sim \mathcal{N}(\tilde{\mathbf{x}}_k(\tilde{\mathbf{u}}_{1:K}^{(i)}), \sigma_{\text{orac}}^2 \mathbf{I}_{d_x}).$$

Now, define the terms

$$\mathbf{z}_k^{(i)} := \mathbf{y}_k^{(i)} - \mathbf{x}_k^\pi$$

Lastly, let  $\mathbb{E}_{\text{orac}}[\cdot]$  denote expectations with respect to the Gaussian noise of the oracle, (conditioning on  $\mathbf{w}_{1:K}^{(i)}$ ) while  $\mathbb{E}[\cdot]$  denotes total expectation. We argue an error bound on

$$\begin{aligned} \left\| \frac{\sigma_w^{-2}}{N} \sum_{i=1}^N \mathbf{z}_k^{(i)} \mathbf{w}_j^{(i)\top} - \Psi_{\text{cl},k,j}^\pi \right\|_{\text{op}} &\leq \underbrace{\frac{\sigma_w^{-2}}{N} \left\| \sum_{i=1}^N \mathbb{E}_{\text{orac}}[\mathbf{z}_k^{(i)} \mathbf{w}_j^{(i)\top}] - \mathbf{z}_k^{(i)} \mathbf{w}_j^{(i)\top} \right\|_{\text{op}}}_{=\text{Term}_1} \\ &+ \underbrace{\frac{\sigma_w^{-2}}{N} \left\| \sum_{i=1}^N \mathbb{E}[\mathbf{z}_k^{(i)} \mathbf{w}_j^{(i)\top}] - \mathbb{E}_{\text{orac}}[\mathbf{z}_k^{(i)} \mathbf{w}_j^{(i)\top}] \right\|_{\text{op}}}_{=\text{Term}_2} \\ &+ \underbrace{\left\| \frac{\sigma_w^{-2}}{N} \sum_{i=1}^N \mathbb{E}[\mathbf{z}_k^{(i)} \mathbf{w}_j^{(i)\top}] - \Psi_{\text{cl},k,j}^\pi \right\|_{\text{op}}}_{=\text{Term}_3}, \end{aligned}$$

which essentially bounds the estimation error of  $\Psi_{\text{cl},k,j}^\pi$  in the absence of observation noise.

**Bounding Term<sub>1</sub>.** Applying [Lemma E.1](#) with  $\|\mathbf{w}_j^{(i)}\| \leq \sqrt{d_*} \sigma_w$ , we have that with probability  $1 - \delta/3$ ,

$$\begin{aligned} \text{Term}_1 &\leq \frac{\sigma_{\text{orac}}}{\sigma_w} \sqrt{2d_* \cdot \frac{d_* \log 5 + \log(6(d_* + 1)/\delta)}{N}} \\ &\leq \frac{\sigma_{\text{orac}}}{\sigma_w} \sqrt{2d_* \cdot \frac{d_* \log 5 + \log(12d_*/\delta)}{N}}. \end{aligned}$$

**Bounding Term<sub>2</sub>.** On the event of [Claim E.1](#), then as long as  $\text{Err}_{\hat{x}}(\delta) \leq \sigma_w \sqrt{d_*/L_\pi}$

$$\|\tilde{\mathbf{u}}_k^{(i)} - \mathbf{u}_k^\pi\| = \|\tilde{\mathbf{w}}_k^{(i)} + \mathbf{K}_k^\pi(\mathbf{x}_k^\pi - \hat{\mathbf{x}}_k)\| \leq \|\tilde{\mathbf{w}}_k^{(i)}\| + L_\pi \text{Err}_{\hat{x}}(\delta) \leq \sigma_w \sqrt{d_*} + L_\pi \text{Err}_{\hat{x}}(\delta) \leq 2\sigma_w \sqrt{d_*}.$$

Notice that as  $\text{Err}_{\hat{x}}(\delta) = \sigma_{\text{orac}} \sqrt{2d_* \iota(\delta)/N}$ ,  $\text{Err}_{\hat{x}}(\delta) \leq \sigma_w \sqrt{d_*/L_\pi}$  holds for  $N \geq (\sigma_{\text{orac}}/\sigma_w)^2 2L_\pi \iota(\delta)$ , i.e. which holds for when  $\pi$  is estimation-friendly.

Moreover, when  $\pi$  is estimation-friendly,  $B_{\text{tay},\text{inf},\pi} \geq 2\sigma_w \sqrt{d_*}$ , so the conditions of [Proposition A.5](#) hold. Therefore,

$$\|\mathbf{z}_k^{(i)}\| \leq 2\sigma_w L_{\text{tay},\infty,\pi} \sqrt{d_*}.$$

and thus

$$\|\tilde{\mathbf{w}}_j^{(i)} \mathbf{z}_k^{(i)}\| \leq 2d_* \sigma_w^2 L_{\text{tay},\infty,\pi}.$$

Applying [Lemma E.2](#) with  $X_i \leftarrow \mathbf{z}_k^{(i)} \mathbf{w}_j^{(i)\top}$  with  $R \leftarrow 2d_* L_{\text{tay},\infty,\pi} \sigma_w^2$ , it holds that with probability  $1 - \delta/3$  that

$$\text{Term}_2 \leq \sigma_w^{-2} \cdot 8L_{\text{tay},\infty,\pi} \sigma_w^2 d_u \sqrt{\frac{\log(\frac{3(d_u+d_x)}{\delta})}{N}} \leq 8L_{\text{tay},\infty,\pi} d_* \sqrt{\frac{\log(6d_*/\delta)}{N}}.$$

**Bounding Term<sub>3</sub>** . As establish in the bound on Term<sub>2</sub>, the conditions of [Proposition A.5](#) hold, and  $\|\check{\mathbf{u}}_k^{(i)} - \mathbf{u}_k^\pi\| \leq 2\sigma_w\sqrt{d_\star}$ . Therefore,

$$\|\mathbf{z}_k^{(i)} - \sum_{\ell=1}^k \Psi_{\text{cl},k,\ell}^\pi(\check{\mathbf{u}}_\ell^{(i)} - \mathbf{u}_\ell^\pi)\| \leq 4\sigma_w^2 M_{\text{tay},2,\pi} d_\star.$$

Consequently, bounding  $\|\mathbf{w}_j^{(i)}\| \leq \sigma_w d_\star$ ,

$$\frac{1}{\sigma_w^2} \|\mathbf{z}_k^{(i)} \mathbf{w}_j^{(i)} - \sum_{\ell=1}^k \Psi_{\text{cl},k,\ell}^\pi(\check{\mathbf{u}}_\ell^{(i)} - \mathbf{u}_\ell^\pi) \mathbf{w}_j^{(i)}\| \leq 4\sigma_w M_{\text{tay},2,\pi} d_\star^{3/2}$$

and thus, by Jensen's inequality,

$$\begin{aligned} 4\sigma_w M_{\text{tay},2,\pi} d_\star^{3/2} &\geq \frac{1}{N\sigma_w^2} \left\| \sum_{i=1}^N \mathbb{E}[\mathbf{z}_k^{(i)} \mathbf{w}_j^{(i)}] - \sum_{\ell=1}^k \mathbb{E}[\Psi_{\text{cl},k,\ell}^\pi(\check{\mathbf{u}}_\ell^{(i)} - \mathbf{u}_\ell^\pi) \mathbf{w}_j^{(i)}] \right\| \\ &= \frac{1}{N\sigma_w^2} \left\| \sum_{i=1}^N \mathbb{E}[\mathbf{z}_k^{(i)} \mathbf{w}_j^{(i)}] - \sum_{\ell=1}^k \Psi_{\text{cl},k,\ell}^\pi \mathbb{E}[(\mathbf{w}_k^{(i)}[\ell] + \mathbf{K}_k^\pi(\mathbf{x}_\ell^\pi - \hat{\mathbf{x}}_\ell)) \mathbf{w}_j^{(i)}] \right\| \\ &= \frac{1}{N\sigma_w^2} \left\| \sum_{i=1}^N \mathbb{E}[\mathbf{z}_k^{(i)} \mathbf{w}_j^{(i)}] - \sigma_w^2 \sum_{\ell=1}^k \Psi_{\text{cl},k,\ell}^\pi \mathbb{I}_{\ell=j} \right\| \\ &= \frac{1}{N\sigma_w^2} \left\| \sum_{i=1}^N \mathbb{E}[\mathbf{z}_k^{(i)} \mathbf{w}_j^{(i)}] - \sigma_w^2 \Psi_{\text{cl},k,j}^\pi \right\| = \text{Term}_3. \end{aligned}$$

In sum, with probability at least  $1 - 3\delta/4$ , the following bound holds

$$\begin{aligned} &\left\| \frac{\sigma_w^{-2}}{N} \sum_{i=1}^N \mathbf{z}_k^{(i)} \mathbf{w}_j^{(i)\top} - \Psi_{\text{cl},k,j}^\pi \right\|_{\text{op}} \\ &\leq \text{Term}_1 + \text{Term}_2 + \text{Term}_3 \\ &\leq \frac{\sigma_{\text{orac}}}{\sigma_w} \sqrt{2d_\star \cdot \frac{d_\star \log 5 + \log(12d_\star/\delta)}{N}} + 8L_{\text{tay},\infty,\pi} d_\star \sqrt{\frac{\log(6d_\star/\delta)}{N}} + 4\sigma_w M_{\text{tay},2,\pi} d_\star^{3/2} \\ &\leq \sqrt{\frac{\log 12d_\star}{\delta}} \left( \frac{2\sigma_{\text{orac}}}{\sigma_w} d_\star^{3/2} + 8L_{\text{tay},\infty,\pi} d_\star \right) + 4\sigma_w M_{\text{tay},2,\pi} d_\star^{3/2}. \end{aligned}$$

The final bound follows from a union bound over all  $\binom{K}{2} \leq K^2$  pairs, and replacing  $\delta$  with  $\delta/2n_{\text{iter}}$ . □

## E.2. Error in the Gradient (Proof of [Lemma A.10](#))

Recall the definitions

$$\begin{aligned} (\nabla \mathcal{J}_T^{\text{disc}}(\pi))_k &= \tau Q_u(\mathbf{x}_k^\pi, \mathbf{u}_k^\pi, t_k) + (\Psi_{\text{cl},K+1,k}^\pi)^\top V_x(\mathbf{x}_{K+1}^\pi) + \\ &\quad + \tau \sum_{j=k+1}^K (\Psi_{\text{cl},j,k}^\pi)^\top (Q_x(\mathbf{x}_j^\pi, \mathbf{u}_j^\pi, t_j) + (\mathbf{K}_j^\pi)^\top Q_u(\mathbf{x}_j^\pi, \mathbf{u}_j^\pi, t_j)) \end{aligned}$$

and

$$\begin{aligned} \hat{\nabla}_k &= Q_u(\bar{\mathbf{x}}_k^\pi, \mathbf{u}_k^\pi, t_j) + \hat{\Psi}_{K+1,k}^\top V_x(\bar{\mathbf{x}}_{K+1}) \\ &\quad + \tau \sum_{j=k+1}^K \hat{\Psi}_{j,k}^\top (Q_x(\hat{\mathbf{x}}_j, \mathbf{u}_j^\pi, t_j) + (\mathbf{K}_j^\pi)^\top Q_u(\hat{\mathbf{x}}_j, \mathbf{u}_j^\pi, t_j)) \end{aligned}$$

Using  $V_x(\cdot)$  and  $Q_x(\cdot), Q_u(\cdot)$  are all at most  $L_{\text{cost}}$  in magnitude, that the gradients of the cost are  $M_{\text{cost}}$ -Lipschitz, and  $1 \vee \|\mathbf{K}_j^\pi\| \leq L_\pi$  we can bound

$$\begin{aligned}
 & \|(\nabla \mathcal{J}_T^{\text{disc}}(\pi))_k - \hat{\nabla}_k\| \\
 & \leq L_{\text{cost}} \|\hat{\Psi}_{K+1,k} - \Psi_{\text{cl},K+1,k}^\pi\| + 2L_\pi L_{\text{cost}} \tau \sum_{j=k+1}^K \|\hat{\Psi}_{j,k} - \Psi_{\text{cl},j,k}^\pi\| \\
 & \quad + M_{\text{cost}} \left( \|\mathbf{x}_k^\pi - \bar{\mathbf{x}}_k^\pi\| + \|\Psi_{\text{cl},K+1,k}^\pi\| \cdot \|\mathbf{x}_{K+1}^\pi - \hat{\mathbf{x}}_{K+1}\| + 2\tau L_\pi \sum_{j=k+1}^K \|\Psi_{\text{cl},j,k}^\pi\| \|\mathbf{x}_j^\pi - \hat{\mathbf{x}}_j\| \right) \\
 & \leq L_{\text{cost}} \text{Err}_{\Psi,\pi}(\delta)(1 + 2L_\pi \tau K) + M_{\text{cost}} \text{Err}_{\hat{\mathbf{x}}}(\delta) \left( 1 + \|\Psi_{\text{cl},K+1,k}^\pi\| + 2\tau L_\pi \sum_{j=k+1}^K \|\Psi_{\text{cl},j,k}^\pi\| \right) \\
 & \leq L_{\text{cost}} \text{Err}_{\Psi,\pi}(\delta)(1 + 2L_\pi \tau K) + M_{\text{cost}} \text{Err}_{\hat{\mathbf{x}}}(\delta) (1 + \kappa_{\pi,\infty} + 2K\tau L_\pi \kappa_{\pi,\infty}) \\
 & = L_{\text{cost}} \text{Err}_{\Psi,\pi}(\delta)(1 + 2TL_\pi) + M_{\text{cost}} \text{Err}_{\hat{\mathbf{x}}}(\delta) (1 + (1 + 2TL_\pi)\kappa_{\pi,\infty}) \\
 & \leq \underbrace{(L_{\text{cost}} \text{Err}_{\Psi,\pi}(\delta) + (1 + \kappa_{\pi,\infty}) M_{\text{cost}} \text{Err}_{\hat{\mathbf{x}}}(\delta))}_{:= \text{Err}_{\nabla,\pi}(\delta)} (1 + 2TL_\pi)
 \end{aligned}$$

□

### E.3. Discrete-Time Closed-Loop Controllability (Proposition A.11)

We begin by lower bounding the continuous-time controllability Gramian under a policy  $\pi$ , and then turn to lower bounding its discretization. At the end of the proof, we remark upon how the bound can be refined. The first part of the argument follows (Chen & Hazan, 2021).

**Equivalent characterization of controllability Gramian smallest singular value.** The following is a continuous-time analogue of Chen & Hazan (2021, Lemma 15).

**Lemma E.3** (Characterization of Controllability Gramian smallest singular value). *Let  $\Psi(t, s) \in \mathbb{R}^{d_x \times d_u}$  be an arbitrary (locally square integrable), and set*

$$\Lambda := \int_{s=t-t_{\text{ctrl}}}^t \Psi(t, s) \Psi(t, s)^\top ds.$$

Then,  $\lambda_{\min}(\Lambda) \geq \nu$  if and only if for all unit vectors  $\xi$ , there exists some  $\mathbf{u}_\xi(s)$  such that  $\xi = \int_{s=t-t_{\text{ctrl}}}^t \Psi(t, s) \mathbf{u}_\xi(s) ds$  and  $\int_{s=t-t_{\text{ctrl}}}^t \|\mathbf{u}_\xi(s)\|^2 ds \leq \nu^{-1}$ .

*Proof of Lemma E.3.* Fix any unit vector  $\xi \in \mathbb{R}^n$ , define. First, suppose  $\int_{s=t-t_{\text{ctrl}}}^t \Psi(t, s) \Psi(t, s)^\top ds \geq \nu$ .

$$\mathbf{u}_\xi(s) := \Psi(t, s)^\top \Lambda^{-1} \xi.$$

One can verify then that

$$\int_{s=t-t_{\text{ctrl}}}^t \Psi(t, s) \mathbf{u}_\xi(s) ds = \Lambda \Lambda^{-1} \xi = \xi \quad \int_{s=t-t_{\text{ctrl}}}^t \|\mathbf{u}_\xi(s)\|^2 ds = \xi \Lambda^{-1} \cdot \Lambda \cdot \Lambda^{-1} \xi = \xi \Lambda^{-1} \xi \leq \lambda_{\min}(\Lambda)^{-1}.$$

On the other hand, suppose that there exists a control  $\mathbf{u}_\xi(s)$  with  $\int_{s=t-t_{\text{ctrl}}}^t \|\mathbf{u}_\xi(s)\|^2 ds \leq \lambda_{\min}(\Lambda)^{-1}$  such that

$\int_{s=t-t_{\text{ctrl}}}^t \Psi(t, s) \mathbf{u}_\xi(s) ds = \xi$ . As  $\xi$  is a unit vector, i.e.  $\xi^\top \xi = 1$ ,

$$\begin{aligned} 1 &= \left( \xi^\top \int_{s=t-t_{\text{ctrl}}}^t \Psi(t, s) \mathbf{u}_\xi(s) ds \right)^2 \\ &= \left( \int_{s=t-t_{\text{ctrl}}}^t \xi^\top \Psi(t, s) \mathbf{u}_\xi(s) ds \right) \\ &\leq \left( \int_{s=t-t_{\text{ctrl}}}^t \|\xi^\top \Psi(t, s)\|^2 ds \right) \cdot \left( \int_{s=t-t_{\text{ctrl}}}^t \|\mathbf{u}_\xi(s)\|^2 ds \right) \\ &\leq \xi^\top \Lambda \xi \cdot \left( \int_{s=t-t_{\text{ctrl}}}^t \|\mathbf{u}_\xi(s)\|^2 ds \right) \\ &\leq \xi^\top \Lambda \xi \cdot \lambda_{\min}(\Lambda)^{-1}. \end{aligned}$$

The bound follows.  $\square$

**Lower bounding the controllability Gramian until alternative policies.** This next part is the continuous-time analogue of (Lemma 16]chen2021black, establishing controllability of the closed-loop linearized system in feedback with policy  $\pi$ .

**Lemma E.4** (Controllability of Closed-Loop Transitions, Continuous-Time). *Recall  $L_\pi \geq 1$ , and  $\gamma_{\text{ctrl}} := \max\{1, L_f t_{\text{ctrl}}\}$ . Then, under Assumption 4.4,*

$$\int_{s=t-t_{\text{ctrl}}}^t \Psi_{\text{cl}}^\pi(t, s) \Psi_{\text{cl}}^\pi(t, s)^\top ds \succeq \frac{\nu_{\text{ctrl}}}{4L_\pi^2 \gamma_{\text{ctrl}}^2 \exp(2\gamma_{\text{ctrl}})}.$$

*Proof of Lemma E.4.* Fix any  $\xi \in \mathbb{R}^{d_x}$  of unit norm. Lemma E.3 and Assumption 4.4 guarantee the existence of an input  $\mathbf{u}_\xi(s)$  for which

$$\int_{s=t-t_{\text{ctrl}}}^t \Phi_{\text{ol}}^\pi(t, s) \mathbf{B}_{\text{ol}}^\pi(s) \mathbf{u}_\xi(s) ds = \xi, \quad \int_{s=t-t_{\text{ctrl}}}^t \|\mathbf{u}_\xi(s)\|^2 ds \leq \nu_{\text{ctrl}}^{-1}.$$

Let

$$\mathbf{z}_\xi(s') = \int_{s=t-t_{\text{ctrl}}}^{s'} \Phi_{\text{ol}}^\pi(t, s) \mathbf{B}_{\text{ol}}^\pi(s) \mathbf{u}_\xi(s) ds.$$

Define now the input

$$\tilde{\mathbf{u}}_\xi(s) := \mathbf{u}_\xi(s) - \mathbb{I}\{t_k(s) > t - t_{\text{ctrl}}\} \mathbf{K}_{k(s)}^\pi \mathbf{z}_\xi(t_k(s)).$$

It can be directly verified (by induction on  $k$ ) that

$$\forall s' \in [t - t_{\text{ctrl}}, t], \quad \int_{s=t-t_{\text{ctrl}}}^{s'} \Phi_{\text{ol}}^\pi(t, s) \mathbf{B}_{\text{ol}}^\pi(s) \mathbf{u}_\xi(s) ds = \int_{s=t-t_{\text{ctrl}}}^{s'} \Psi_{\text{cl}}^\pi(t, s) \tilde{\mathbf{u}}_\xi(s) ds,$$

so in particular

$$\xi = \int_{s=t-t_{\text{ctrl}}}^t \Psi_{\text{cl}}^\pi(t, s) \tilde{\mathbf{u}}_\xi(s) ds.$$

We may now bound

$$\begin{aligned} \int_{s=t-t_{\text{ctrl}}}^t \|\tilde{\mathbf{u}}_\xi(s)\|^2 ds &= \int_{s=t-t_{\text{ctrl}}}^t \left( \|\mathbf{u}_\xi(s) - \mathbb{I}\{t_k(s) > t - t_{\text{ctrl}}\} \mathbf{K}_{k(s)}^\pi \mathbf{z}_\xi(t_k(s))\|^2 \right) ds \\ &\leq 2 \int_{s=t-t_{\text{ctrl}}}^t \left( \|\mathbf{u}_\xi(s)\|^2 + \|\mathbf{K}_{k(s)}^\pi\| \|\mathbf{z}_\xi(t_k(s))\|^2 \right) ds \\ &\leq 2(\nu_{\text{ctrl}}^{-1} + L_\pi^2) \int_{s=t-t_{\text{ctrl}}}^t \|\mathbf{z}_\xi(t_k(s))\|^2 ds, \end{aligned} \tag{E.4}$$

We now adopt the following claim, mirroring the proof of Chen & Hazan (2021, Lemma 16).

**Claim E.2.** *The following bound holds:*

$$\forall s' \in [t - t_{\text{ctrl}}, t], \quad \|\mathbf{z}_\xi(s')\|^2 \leq t_{\text{ctrl}} \nu_{\text{ctrl}}^{-1} L_f^2 \exp(2t_{\text{ctrl}} L_f).$$

*Proof of Claim E.2.* We bound

$$\begin{aligned} \|\mathbf{z}_\xi(s')\|^2 &= \left\| \int_{s=t-t_{\text{ctrl}}}^{s'} \Phi_{\text{ol}}^\pi(t, s) \mathbf{B}_{\text{ol}}^\pi(s) \mathbf{u}_\xi(s) ds \right\|^2 \\ &\leq \left( \int_{s=t-t_{\text{ctrl}}}^{s'} \|\Phi_{\text{ol}}^\pi(t, s)\| \|\mathbf{B}_{\text{ol}}^\pi(s)\| \|\mathbf{u}_\xi(s)\| ds \right)^2 \\ &\leq L_f^2 \left( \int_{s=t-t_{\text{ctrl}}}^{s'} \|\Phi_{\text{ol}}^\pi(t, s)\| \|\mathbf{u}_\xi(s)\| ds \right)^2 && \text{(Assumption 4.1)} \\ &\leq L_f^2 \left( \int_{s=t-t_{\text{ctrl}}}^{s'} \|\Phi_{\text{ol}}^\pi(t, s)\|^2 ds \right) \left( \int_{s=t-t_{\text{ctrl}}}^{s'} \|\mathbf{u}_\xi(s)\|^2 ds \right)^2 && \text{(Cauchy-Schwartz)} \\ &\leq L_f^2 \left( \int_{s=t-t_{\text{ctrl}}}^t \|\Phi_{\text{ol}}^\pi(t, s)\|^2 ds \right) \left( \int_{s=t-t_{\text{ctrl}}}^t \|\mathbf{u}_\xi(s)\|^2 ds \right) \\ &\leq \nu_{\text{ctrl}}^{-1} L_f^2 \left( \int_{s=t-t_{\text{ctrl}}}^t \|\Phi_{\text{ol}}^\pi(t, s)\|^2 ds \right). \end{aligned}$$

By Lemma I.4, we can bound  $\|\Phi_{\text{ol}}^\pi(t, s)\| \leq \exp(L_f(t-s)) \leq \exp(L_f t_{\text{ctrl}})$  for  $s \in [t - t_{\text{ctrl}}, t]$ , yielding  $\int_{s=t-t_{\text{ctrl}}}^t \|\Phi_{\text{ol}}^\pi(t, s)\|^2 ds \leq t_{\text{ctrl}} \exp(2L_f t_{\text{ctrl}})$ . The bound claim.  $\square$

Combining Eq. (E.4) and Claim E.2,

$$\begin{aligned} \int_{s=t-t_{\text{ctrl}}}^t \|\tilde{\mathbf{u}}_\xi(s)\|^2 &\leq 2\nu_{\text{ctrl}}^{-1} (1 + t_{\text{ctrl}}^2 L_f^2 L_\pi^2 \exp(2L_f t_{\text{ctrl}})) \\ &\leq 2\nu_{\text{ctrl}}^{-1} L_\pi^2 (1 + t_{\text{ctrl}}^2 L_f^2 \exp(2L_f t_{\text{ctrl}})) && (L_\pi \geq 1) \\ &\leq 2\nu_{\text{ctrl}}^{-1} L_\pi^2 (1 + \gamma_{\text{ctr}}^2 \exp(2\gamma_{\text{ctr}})) && (\gamma_{\text{ctr}} = \max\{1, t_{\text{ctrl}} L_f\}) \\ &\leq 4\nu_{\text{ctrl}}^{-1} L_\pi^2 \gamma_{\text{ctr}}^2 \exp(2\gamma_{\text{ctr}}), && (\gamma_{\text{ctr}} \geq 1) \end{aligned}$$

which concludes the proof.  $\square$

**Discretizing the Closed-Loop Gramian.** To conclude the argument, we relate the controllability of the closed-loop Gramian in continuous-time to that in discrete-time.

**Lemma E.5** (Discretization of Controllability Gramian). *Suppose Assumption 4.4 holds and  $\tau \leq L_f/4$ , then following holds:*

$$\begin{aligned} &\left\| \int_{s=t_k-t_{\text{ctrl}}}^{t_k} \Psi_{\text{cl}}^\pi(t_k, s) \Psi_{\text{cl}}^\pi(t_k, s)^\top ds - \frac{1}{\tau} \sum_{j=k-k_{\text{ctrl}}}^{k-1} \Psi_{\text{cl},k,j}^\pi (\Psi_{\text{cl},k,j}^\pi)^\top \right\|_{\text{op}} \\ &\leq 4\tau \gamma_{\text{ctr}} \kappa_{\pi, \infty}^2 (\kappa_f M_f + 2L_f^2) \end{aligned}$$

*Proof.* Recall the shorthand  $L_{\text{ol}} := \exp(\tau L_f)$ , used in the discretization arguments in [Appendix I](#). We can write

$$\begin{aligned}
 & \left\| \int_{s=t_k-t_{\text{ctrl}}}^{t_k} \Psi_{\text{cl}}^\pi(t_k, s) \Psi_{\text{cl}}^\pi(t_k, s)^\top ds - \tau^{-1} \sum_{j=k-k_{\text{ctrl}}}^{k-1} \Psi_{\text{cl},k,j}^\pi (\Psi_{\text{cl},k,j}^\pi)^\top \right\| \\
 &= \left\| \sum_{j=k-k_{\text{ctrl}}}^{k-1} \int_{s=t_j}^{t_{j+1}} \Psi_{\text{cl}}^\pi(t, s) \Psi_{\text{cl}}^\pi(t, s)^\top ds - \frac{1}{\tau} \Psi_{\text{cl},k,j}^\pi \left( \frac{1}{\tau} \Psi_{\text{cl},k,j}^\pi \right)^\top \right\| \\
 &\leq \tau \sum_{j=k-k_{\text{ctrl}}}^{k-1} \max_{s \in \mathcal{I}_j} \left\| \Psi_{\text{cl}}^\pi(t_k, s) \Psi_{\text{cl}}^\pi(t_k, s)^\top - \frac{1}{\tau} \Psi_{\text{cl},k,j}^\pi \left( \frac{1}{\tau} \Psi_{\text{cl},k,j}^\pi \right)^\top \right\| \\
 &\leq 2\tau \sum_{j=k-k_{\text{ctrl}}}^{k-1} \max_{s \in \mathcal{I}_j} \left\| \Psi_{\text{cl}}^\pi(t_k, s) - \frac{1}{\tau} \Psi_{\text{cl},k,j}^\pi \right\| \cdot \max \left\{ \left\| \Psi_{\text{cl}}^\pi(t_k, s) \right\|, \frac{1}{\tau} \left\| \Psi_{\text{cl},k,j}^\pi \right\| \right\} \\
 &\leq 2L_f L_{\text{ol}} \tau \kappa_{\pi, \infty} \sum_{j=k-k_{\text{ctrl}}}^{k-1} \max_{s \in \mathcal{I}_j} \left\| \Psi_{\text{cl}}^\pi(t_k, s) - \frac{1}{\tau} \Psi_{\text{cl},k,j}^\pi \right\| \tag{Lemma I.8(d)} \\
 &\leq 2L_f L_{\text{ol}} \kappa_{\pi, \infty}^2 (\kappa_f M_f + 2L_f^2) \sum_{j=k-k_{\text{ctrl}}}^{k-1} \tau^2 \tag{Lemma I.8(b)} \\
 &\leq 2L_f L_{\text{ol}}^2 \kappa_{\pi, \infty}^2 (\kappa_f M_f + 2L_f^2) \sum_{j=k-k_{\text{ctrl}}}^{k-1} \tau^2 \\
 &\leq 2L_f L_{\text{ol}}^2 \kappa_{\pi, \infty}^2 (\kappa_f M_f + 2L_f^2) k_{\text{ctrl}} \tau^2 \\
 &= 2\tau t_{\text{ctrl}} L_f L_{\text{ol}}^2 \kappa_{\pi, \infty}^2 (\kappa_f M_f + 2L_f^2).
 \end{aligned}$$

As  $\tau \leq L_f/4$ ,  $L_{\text{ol}}^2 \leq \exp(1/2) \leq 2$ , so that the above is at most  $4\tau t_{\text{ctrl}} L_f \kappa_{\pi, \infty}^2 (\kappa_f M_f + 2L_f^2)$ . Recalling  $\gamma_{\text{ctr}} := \max\{1, t_{\text{ctrl}} L_f\}$  concludes.  $\square$

### Concluding the proof.

*Proof of [Proposition A.11](#).* The proof follows by combining the bounds in [Lemmas E.4](#) and [E.5](#). These yield

$$\frac{1}{\tau} \lambda_{\min} \left( \sum_{j=k-k_{\text{ctrl}}}^{k-1} \Psi_{\text{cl},k,j}^\pi (\Psi_{\text{cl},k,j}^\pi)^\top \right) \succeq \frac{\nu_{\text{ctrl}}}{4L_\pi^2 \gamma_{\text{ctr}}^2 \exp(2\gamma_{\text{ctr}})} - \kappa_{\pi, \infty}^2 \cdot 4\tau \gamma_{\text{ctr}} (\kappa_f M_f + 2L_f^2)$$

Recall  $\gamma_{\text{ctr}} = t_{\text{ctrl}} L_f$ . Hence, if

$$\tau \leq \frac{\nu_{\text{ctrl}}}{8L_\pi^2 \kappa_{\pi, \infty}^2 \gamma_{\text{ctr}}^3 \exp(2\gamma_{\text{ctr}}) (\kappa_f M_f + 2L_f^2)},$$

it holds that

$$\lambda_{\min} \left( \sum_{j=k-k_{\text{ctrl}}}^{k-1} \Psi_{\text{cl},k,j}^\pi (\Psi_{\text{cl},k,j}^\pi)^\top \right) \succeq \frac{\nu_{\text{ctrl}}}{8L_\pi^2 \gamma_{\text{ctr}}^2 \exp(2\gamma_{\text{ctr}})}.$$

$\square$

### E.4. Recovery of State-Transition Matrix ([Proposition A.12](#))

The analysis is based on the Ho-Kalman scheme. We begin with the observation that

$$\Psi_{\text{cl},k,j}^\pi = A_{\text{cl},k}^\pi \Psi_{\text{cl},k-1,j}^\pi, \quad \forall j < k.$$



To this end, define the matrices

$$\mathcal{C}_{k|j_2, j_1} := [\Psi_{\text{cl}, k+1, j_2}^\pi \mid \Psi_{\text{cl}, k+1, j_2-1}^\pi \mid \dots \mid \Psi_{\text{cl}, k+1, j_1}^\pi],$$

Then, we have the identity

$$\mathcal{C}_{k|k-1, j} = \mathbf{A}_{\text{cl}, k}^\pi \mathcal{C}_{k-1|k-1, j},$$

so that if  $\text{rank}(\mathcal{C}_{k-1|k-1, j}) = d_x$ , we have  $\mathbf{A}_{\text{cl}, k}^\pi = \mathcal{C}_{k|k-1, j} \mathcal{C}_{k-1|k-1, j}^\dagger$ , where  $(\cdot)^\dagger$  denotes the Moore-Penrose pseudoinverse. We now state and prove a more-or-less standard perturbation bound.

**Lemma E.6.** *Suppose  $\text{rank}(\mathcal{C}_{k-1|k-1, j}) = d_x$ , and consider any estimates  $\hat{\mathcal{C}}_{k|k-1, j}, \hat{\mathcal{C}}_{k-1|k-1, j}$  of  $\mathcal{C}_{k|k-1, j}, \mathcal{C}_{k-1|k-1, j}$ . Define*

$$\begin{aligned} \Delta &:= \max\{\|\mathcal{C}_{k|k-1, j} - \hat{\mathcal{C}}_{k|k-1, j}\|, \|\mathcal{C}_{k-1|k-1, j} - \hat{\mathcal{C}}_{k-1|k-1, j}\|\} \\ M &:= \max\{\|\mathcal{C}_{k|k-1, j}\|, \|\mathcal{C}_{k-1|k-1, j}\|\}. \end{aligned}$$

Then, if  $\Delta \leq \sigma_{\min}(\mathcal{C}_{k-1|k-1, j})/2$ , the estimate  $\tilde{\mathbf{A}}_k := \hat{\mathcal{C}}_{k|k-1, j} \hat{\mathcal{C}}_{k-1|k-1, j}^\dagger$  satisfies

$$\|\tilde{\mathbf{A}}_k - \mathbf{A}_{\text{cl}, k}^\pi\| \leq 6\Delta M \sigma_{\min}(\mathcal{C}_{k-1|k-1, j})^{-2}.$$

*Proof of Lemma E.6.* Then, we have (provided  $\text{rank}(\hat{\mathcal{C}}_{k-1|k-1, j}) = \text{rank}(\mathcal{C}_{k-1|k-1, j}) = d_x$ ), we have

$$\begin{aligned} \|\tilde{\mathbf{A}}_k - \mathbf{A}_{\text{cl}, k}^\pi\| &= \|\mathcal{C}_{k|k-1, j} \mathcal{C}_{k-1|k-1, j}^\dagger - \hat{\mathcal{C}}_{k|k-1, j} \hat{\mathcal{C}}_{k-1|k-1, j}^\dagger\| \\ &\leq \|\hat{\mathcal{C}}_{k-1|k-1, j}^\dagger\| \|\mathcal{C}_{k|k-1, j} - \hat{\mathcal{C}}_{k|k-1, j}\| + \|\hat{\mathcal{C}}_{k-1|k-1, j}^\dagger - \mathcal{C}_{k-1|k-1, j}^\dagger\| \|\mathcal{C}_{k|k-1, j}\| \\ &\leq \|\hat{\mathcal{C}}_{k-1|k-1, j}^\dagger\| \|\mathcal{C}_{k|k-1, j} - \hat{\mathcal{C}}_{k|k-1, j}\| \\ &\quad + \frac{1 + \sqrt{5}}{2} \|\hat{\mathcal{C}}_{k-1|k-1, j}^\dagger\| \cdot \|\mathcal{C}_{k-1|k-1, j}^\dagger\| \cdot \|\mathcal{C}_{k-1|k-1, j} - \hat{\mathcal{C}}_{k-1|k-1, j}\| \cdot \|\hat{\mathcal{C}}_{k|k-1, j}\| \\ &\hspace{15em} \text{(cite (Stewart, 1977), and also (Xu, 2020))} \\ &\stackrel{(i)}{\leq} \Delta \|\hat{\mathcal{C}}_{k-1|k-1, j}^\dagger\| \left(1 + \frac{1 + \sqrt{5}}{2} \|\mathcal{C}_{k-1|k-1, j}^\dagger\| \|\mathcal{C}_{k|k-1, j}\|\right) \\ &\stackrel{(ii)}{\leq} 3\Delta M \|\hat{\mathcal{C}}_{k-1|k-1, j}^\dagger\| \|\mathcal{C}_{k-1|k-1, j}^\dagger\|, \end{aligned}$$

where in (i) we use  $\Delta := \max\{\|\mathcal{C}_{k|k-1, j} - \hat{\mathcal{C}}_{k|k-1, j}\|, \|\mathcal{C}_{k-1|k-1, j} - \hat{\mathcal{C}}_{k-1|k-1, j}\|\}$ , and in (ii), we use  $M = \max\{\|\mathcal{C}_{k|k-1, j}\|, \|\mathcal{C}_{k-1|k-1, j}\|\}$ , which admits the simplification in (ii) because  $\|\mathcal{C}_{k-1|k-1, j}^\dagger\| \|\mathcal{C}_{k-1|k-1, j}\| \geq 1$ . In particular, if  $\text{rank}(\mathcal{C}_{k-1|k-1, j}) = d_x$ , and  $\Delta \leq \sigma_{\min}(\mathcal{C}_{k-1|k-1, j})/2$ , then  $\|\hat{\mathcal{C}}_{k-1|k-1, j}^\dagger\| \leq 2/\sigma_{\min}(\mathcal{C}_{k-1|k-1, j})$ , and we obtain

$$\|\tilde{\mathbf{A}}_k - \mathbf{A}_{\text{cl}, k}^\pi\| \leq 6\Delta M \sigma_{\min}(\mathcal{C}_{k-1|k-1, j})^{-2}.$$

□

Next, restricting our attention to  $k \geq k_{\text{ctrl}} + 2$ , we specialize the above analysis to

$$\begin{aligned} \mathcal{C}_{k, \text{in}} &:= \mathcal{C}_{k-1|k-1, k-k_0+1}, & \mathcal{C}_{k, \text{out}} &:= \mathcal{C}_{k|k-1, k-k_0+1} \\ \hat{\mathcal{C}}_{k, \text{in}} &:= \hat{\mathcal{C}}_{k-1|k-1, k-k_0+1}, & \hat{\mathcal{C}}_{k, \text{out}} &:= \hat{\mathcal{C}}_{k|k-1, k-k_0+1} \end{aligned}$$

where  $\hat{\mathcal{C}}_{(\cdot)}$  arises from the plug-in estimates

$$\hat{\mathcal{C}}_{k|j_2, j_1} := [\hat{\Psi}_{k+1, j_2} \mid \hat{\Psi}_{k+1, j_2-1} \mid \dots \mid \hat{\Psi}_{k+1, j_1}]$$

Define further

$$\tilde{\mathbf{A}}_k := \hat{\mathcal{C}}_{k, \text{out}} \hat{\mathcal{C}}_{k, \text{in}}^\dagger, \quad \text{so that } \hat{\mathbf{A}}_k = \tilde{\mathbf{A}}_k - \hat{\mathbf{B}}_k \mathbf{K}_k^\pi.$$

Recall  $t_0 = k_0/\tau$ . We can now bound, recalling  $L_{\text{ol}} := \exp(\tau L_f) \leq 2$  for  $\tau \leq L_f/4$  and  $\gamma_{\text{ctr}} = \max\{1, t_{\text{ctrl}} L_f\} = \max\{1, \tau k_{\text{ctrl}} L_f\}$ ,

$$\begin{aligned} \max\{\|\mathcal{C}_{k,\text{in}}\|, \|\mathcal{C}_{k,\text{out}}\|\} &\leq \sqrt{k_0} \max_{j < k} \|\Psi_{\text{cl},k,j}^\pi\| \\ &\leq \sqrt{k_{\text{ctrl}}} \tau L_f L_{\text{ol}} \kappa_{\pi,\infty} && \text{(Lemma I.8(d))} \\ &\leq \sqrt{\tau t_0} \tau L_f L_{\text{ol}} \kappa_{\pi,\infty} && (t_0 = k_0 \tau) \\ &\leq 2\kappa_{\pi,\infty} \gamma_{\text{ctr}} \sqrt{\tau t_0}. && (\gamma_{\text{ctr}} \geq 1) \end{aligned}$$

Invoking Proposition A.11, we also have that provided  $\tau \leq \min\{\tau_{\text{dyn}}, \tau_{\text{ctrl},\pi}\}$ , since  $k_0 \geq k_{\text{ctrl}} + 2$ ,

$$\begin{aligned} \sigma_{\min}(\mathcal{C}_{k-1|k-1,j})^2 &= \lambda_{\min} \left( \sum_{j=k-k_0+1}^{k-1} \Psi_{\text{cl},k-1,j}^\pi (\Psi_{\text{cl},k-1,j}^\pi)^\top \right) \\ &\geq \lambda_{\min} \left( \sum_{j=k-k_{\text{ctrl}}-1}^{k-1} \Psi_{\text{cl},k-1,j}^\pi (\Psi_{\text{cl},k-1,j}^\pi)^\top \right) \geq \tau \cdot \frac{\nu_{\text{ctrl}}}{8L_\pi^2 \gamma_{\text{ctr}}^2 \exp(2\gamma_{\text{ctr}})}. \end{aligned}$$

Therefore, as long as

$$\Delta := \max\{\|\mathcal{C}_{k,\text{in}} - \hat{\mathcal{C}}_{k,\text{in}}\|, \|\mathcal{C}_{k,\text{out}} - \hat{\mathcal{C}}_{k,\text{out}}\|\} \leq \frac{\sqrt{\tau \nu_{\text{ctrl}}}}{2\sqrt{2}L_\pi \gamma_{\text{ctr}} \exp(\gamma_{\text{ctr}})},$$

we have

$$\|\tilde{\mathbf{A}}_k - \mathbf{A}_{\text{cl},k}^\pi\| \leq \sqrt{\frac{t_0}{\tau}} \frac{96\Delta}{\nu_{\text{ctrl}}} \cdot \kappa_{\pi,\infty} L_\pi^2 \gamma_{\text{ctr}}^3 \exp(2\gamma_{\text{ctr}}).$$

Lastly, we can upper bound  $\Delta \leq \sqrt{k_{\text{ctrl}}} \text{Err}_{\Psi,\pi}(\delta) = \sqrt{t_0/\tau} \text{Err}_{\Psi,\pi}(\delta)$ , from which we conclude that as long as  $\text{Err}_{\Psi}(\delta) \leq \tau \frac{\sqrt{\nu_{\text{ctrl}}/t_0}}{2\sqrt{2}L_\pi \gamma_{\text{ctr}} \exp(\gamma_{\text{ctr}})}$ , we have

$$\|\tilde{\mathbf{A}}_k - \mathbf{A}_{\text{cl},k}^\pi\| \leq t_0 \kappa_{\pi,\infty} L_\pi^2 \cdot \frac{96 \text{Err}_{\Psi,\pi}(\delta)}{\tau \nu_{\text{ctrl}}} \cdot \gamma_{\text{ctr}}^3 \exp(2\gamma_{\text{ctr}}).$$

Now to wrap up. We observe that  $\mathbf{B}_{\text{ol},k}^\pi = \Psi_{\text{cl},k+1,k}^\pi$ , so

$$\|\hat{\mathbf{B}}_k - \mathbf{B}_{\text{ol},k}^\pi\| = \|\Psi_{\text{cl},k+1,k}^\pi - \hat{\Psi}_{k+1,k}\| \leq \text{Err}_{\Psi,\pi}(\delta).$$

Therefore,

$$\begin{aligned} \|\hat{\mathbf{A}}_k - \mathbf{A}_{\text{ol},k}^\pi\| &= \|\tilde{\mathbf{A}}_k - (\mathbf{A}_{\text{cl},k}^\pi - \mathbf{B}_{\text{ol},k}^\pi \mathbf{K}_k^\pi)\| \\ &\leq \|\tilde{\mathbf{A}}_k - \mathbf{A}_{\text{cl},k}^\pi\| + \|\hat{\mathbf{B}}_k \mathbf{K}_k^\pi - \mathbf{B}_{\text{ol},k}^\pi \mathbf{K}_k^\pi\| \\ &\leq \|\tilde{\mathbf{A}}_k - \mathbf{A}_{\text{cl},k}^\pi\| + \|\hat{\mathbf{B}}_k - \mathbf{B}_{\text{ol},k}^\pi\| L_\pi \\ &\leq L_\pi \text{Err}_{\Psi,\pi}(\delta) + t_0 \kappa_{\pi,\infty} L_\pi^2 \cdot \frac{96 \text{Err}_{\Psi,\pi}(\delta)}{\tau \nu_{\text{ctrl}}} \cdot \gamma_{\text{ctr}}^3 \exp(2\gamma_{\text{ctr}}). \end{aligned}$$

Lastly, we notice this upper bound on  $\|\hat{\mathbf{A}}_k - \mathbf{A}_{\text{ol},k}^\pi\|$  is larger than that on  $\|\hat{\mathbf{B}}_k - \mathbf{B}_{\text{ol},k}^\pi\|$ , as  $L_\pi \geq 1$  by assumption, and that for  $\tau \leq \tau_{\text{ctrl},\pi}$ ,  $L_\pi \text{Err}_{\Psi,\pi}(\delta) \leq t_0 \kappa_{\pi,\infty} L_\pi^2 \cdot \frac{96 \text{Err}_{\Psi,\pi}(\delta)}{\tau \nu_{\text{ctrl}}} \cdot \gamma_{\text{ctr}}^3 \exp(2\gamma_{\text{ctr}})$ . Thus,

$$\|\hat{\mathbf{A}}_k - \mathbf{A}_{\text{ol},k}^\pi\| \vee \|\hat{\mathbf{B}}_k - \mathbf{B}_{\text{ol},k}^\pi\| \leq \frac{\text{Err}_{\Psi,\pi}(\delta)}{\tau} \cdot t_0 \kappa_{\pi,\infty} L_\pi^2 \gamma_{\text{ctr}}^3 \exp(2\gamma_{\text{ctr}}) \cdot \frac{192}{\nu_{\text{ctrl}}}.$$

□

## F. Certainty Equivalence

In this section, we establish a general certainty equivalence bound for linear time-varying discrete-time systems; we apply this in the proof [Proposition A.14](#) in [Appendix G.1](#).

Let  $\Theta^* := (\mathbf{A}_{1:K}^*, \mathbf{B}_{1:K}^*)$  denote ground-truth system parameters, and let  $\hat{\Theta} := (\hat{\mathbf{A}}_{1:K}, \hat{\mathbf{B}}_{1:K})$  denote estimates. We work with a slightly different discretization parameterization, where the dynamics are given by  $\mathbf{x}_{h+1} = \mathbf{A}_k \mathbf{x}_h + \tau \mathbf{B}_k \mathbf{u}_h$ . This parametrization ensures that the norms of  $\mathbf{B}_k$  scale like constants independent of  $\tau$  when instantiated with  $\mathbf{A}_k \leftarrow \mathbf{A}_{\text{ol},k}^\pi$  and  $\mathbf{B}_k \leftarrow \frac{1}{\tau} \mathbf{B}_{\text{ol},k}^\pi$ .

**Definition F.1.** Given cost matrices  $\mathbf{Q}, \mathbf{R}$ , step  $\tau$ , and parameters  $\Theta = (\mathbf{A}_{1:K}, \mathbf{B}_{1:K})$ , we define  $\mathbf{P}_k^{\text{opt}}(\Theta)$  as the solution to the following program

$$\begin{aligned} x^\top \mathbf{P}_k^{\text{opt}}(\Theta) x &= \min_{\mathbf{u}_{k:H}} x_{K+1}^\top \mathbf{Q} x_{K+1} + \tau \sum_{h=k}^K (x_h^\top \mathbf{Q} x_h + \mathbf{u}_h^\top \mathbf{Q} \mathbf{u}_h) \\ \text{s.t. } \mathbf{x}_{h+1} &= \mathbf{A}_h \mathbf{x}_h + \tau \mathbf{B}_h \mathbf{u}_h, \quad \mathbf{x}_k = x. \end{aligned} \quad (\text{F.1})$$

The closed form for  $\mathbf{P}_k^{\text{opt}}$  is given by the follow standard computation, modified with the reparametrized dynamics  $\mathbf{x}_{h+1} = \mathbf{A}_h \mathbf{x}_h + \tau \mathbf{B}_h \mathbf{u}_h$ )

**Lemma F.1.** *The optimal Riccati cost-to-go  $\mathbf{P}_{1:K+1}^{\text{opt}} = \mathbf{P}_k^{\text{opt}}(\Theta)$  is given by the solution to the following recursion with final condition  $\mathbf{P}_{K+1}^{\text{opt}} = \mathbf{Q}$  and*

$$\mathbf{P}_k^{\text{opt}} = \mathbf{A}_k^\top \mathbf{P}_{k+1}^{\text{opt}} \mathbf{A}_k - \tau (\mathbf{B}_k \mathbf{P}_{k+1}^{\text{opt}} \mathbf{A}_k)^\top (\mathbf{R} + \tau \mathbf{B}_k^\top \mathbf{P}_{k+1}^{\text{opt}} \mathbf{B}_k)^{-1} (\mathbf{B}_k \mathbf{P}_{k+1}^{\text{opt}} \mathbf{A}_k) + \tau \mathbf{Q}$$

Moreover, defining  $\mathbf{K}_k^{\text{opt}} = \mathbf{K}_k^{\text{opt}}(\Theta) := -(\mathbf{R} + \tau \mathbf{B}_k^\top \mathbf{P}_k^{\text{opt}} \mathbf{B}_k)^{-1} \mathbf{B}_k^\top \mathbf{P}_k^{\text{opt}} \mathbf{A}_k$ , the optimal control for [Eq. \(F.1\)](#) is given by  $\mathbf{x}_k = \mathbf{K}_k^{\text{opt}} \mathbf{u}_k$ .

*Proof.* This follows by reparametrizing the standard discrete-time Riccati update (see e.g. [Anderson & Moore \(2007, Section 2.4\)](#)), with  $\mathbf{B}_k \leftarrow \tau \mathbf{B}_k$ ,  $\mathbf{Q} \leftarrow \tau \mathbf{Q}$ , and  $\mathbf{R} \leftarrow \tau \mathbf{R}$ , and simplifying dependence on  $\tau$ .  $\square$

The following identity is also standard (again, consult [Anderson & Moore \(2007, Section 2.4\)](#)), albeit again with the reparamerizations  $\mathbf{B}_k \leftarrow \tau \mathbf{B}_k$ ,  $\mathbf{Q} \leftarrow \tau \mathbf{Q}$ , and  $\mathbf{R} \leftarrow \tau \mathbf{R}$ ):

$$\mathbf{P}_k^{\text{opt}} = (\mathbf{A}_k + \tau \mathbf{B}_k \mathbf{K}_k^{\text{opt}})^\top \mathbf{P}_{k+1}^{\text{opt}} (\mathbf{A}_k + \tau \mathbf{B}_k \mathbf{K}_k^{\text{opt}}) + \tau (\mathbf{Q} + (\mathbf{K}_k^{\text{opt}})^\top \mathbf{R} (\mathbf{K}_k^{\text{opt}})) \quad (\text{F.2})$$

Next, we define the cost-to-go functions associated for arbitrary sequences of feedback matrices, and from the optimal feedback matrices from another instance  $\Theta'$ .

**Definition F.2** (Feedback and Certainty Equivalent Cost-to-go). Given a sequence of feedback gains  $\mathbf{K}_{1:K}$ , we define the induced cost-to-go

$$\begin{aligned} \mathbf{P}_k^{\text{feed}}(\Theta; \mathbf{K}_{1:K}) &:= \mathbf{x}_{K+1}^\top \mathbf{Q} \mathbf{x}_{K+1} + \tau \sum_{h=k}^K (\mathbf{x}_h^\top \mathbf{Q} \mathbf{x}_h + \mathbf{u}_h^\top \mathbf{Q} \mathbf{u}_h) \\ \text{s.t. } \mathbf{x}_{h+1} &= (\mathbf{A}_h + \tau \mathbf{B}_h \mathbf{K}_h) \mathbf{x}_h \quad \mathbf{x}_k = x. \end{aligned}$$

And define the *certainty equivalent* cost-to-go as  $\mathbf{P}_k^{\text{ce}}(\Theta; \Theta') = \mathbf{P}_k^{\text{feed}}(\Theta; \mathbf{K}_{1:K}^{\text{opt}}(\Theta'))$  as the feedback cost-to-go for  $\Theta$  using the optimal gains for  $\Theta'$ .

In particular,  $\mathbf{P}_k^{\text{ce}}(\Theta; \Theta) = \mathbf{P}_k^{\text{opt}}(\Theta)$ . We now present upper bounds on  $\mathbf{P}_k^{\text{ce}}(\Theta; \Theta')$ . We assume bounds on the various parameters of interest.

**Condition F.1.** We have that there are constants  $K_B, K_A \geq 1$  such that, for all  $k \in [K]$ ,

$$\|\mathbf{B}_k^*\| \vee \|\hat{\mathbf{B}}_k\| \leq K_B \quad \|\mathbf{A}_k^*\| \vee \|\hat{\mathbf{A}}_k\| \leq K_A,$$

**Condition F.2.** We assume that there exists  $\Delta_A, \Delta_B > 0$ ,

$$\forall k, \quad \|\hat{\mathbf{B}}_k - \mathbf{B}_k^*\| \leq \Delta_B \text{ and } \|\tau^{-1}(\hat{\mathbf{A}}_k - \mathbf{A}_k^*)\| \leq \Delta_A.$$

**Condition F.3.** We assume the a normalization on the cost matrices satisfy  $\mathbf{R} \succeq \mathbf{I}$ ,  $\mathbf{Q} \succeq \mathbf{I}$  and  $\|\mathbf{Q}\| \geq \|\mathbf{R}\|$ . As a special case,  $\mathbf{Q} = \mathbf{I}$  and  $\mathbf{R} = \mathbf{I}$  suffices.

Lastly, the following assumption is needed to derive an upper bound on the closed-loop transition operator.

**Condition F.4.** We assume that  $\max_k \|\mathbf{A}_k - \mathbf{I}\| \leq \tau\kappa_A$ .

**Theorem 4** (Main Perturbation). *Suppose [Conditions F.1 to F.3](#) hold. Define the terms*

$$\Delta_{\text{ce}} := 80C^4 K_A^3 K_B^3 (1 + \tau C K_B) (\Delta_A + \Delta_B), \quad C := \max_{k \in [K+1]} \|\mathbf{P}_k^{\text{opt}}(\Theta)\|.$$

Then, as long as  $\Delta_{\text{ce}} < 1$ , we have

(a)

$$\max_{k \in [K+1]} \|\mathbf{P}_k^{\text{ce}}(\Theta; \hat{\Theta})\| \leq (1 - \Delta_{\text{ce}})^{-1} \max_{k \in [K+1]} \|\mathbf{P}_k^{\text{opt}}(\Theta)\|$$

(b)  $\max_{k \in [K+1]} \|\mathbf{K}_k^{\text{opt}}(\hat{\Theta})\| \leq \frac{5}{4} K_B K_A C.$

(c) Moreover, if [Condition F.4](#) holds, then the transition operators defined as

$$\Phi_{j,k}^{\text{ce}} := (\mathbf{A}_{j-1} + \tau \mathbf{B}_{j-1} \mathbf{K}_{j-1}^{\text{opt}}(\hat{\Theta})) \cdot (\mathbf{A}_{j-2} + \tau \mathbf{B}_{j-2} \mathbf{K}_{j-2}^{\text{opt}}(\hat{\Theta})) \cdot \dots \cdot (\mathbf{A}_k + \tau \mathbf{B}_k \mathbf{K}_k^{\text{opt}}(\hat{\Theta})),$$

with the convention  $\Phi_{k,k+}^{\text{ce}} = \mathbf{I}$  satisfy, for all  $1 \leq j \leq k \leq K$ ,

$$\|\Phi_{j,k}^{\text{ce}}\|^2 \leq 2\kappa(1 - \tau\gamma)^{j-k}, \quad \text{where } \kappa = \kappa_A + \frac{5}{4} K_B^2 K_A C, \quad \gamma = \frac{1 - \Delta_{\text{ce}}}{C},$$

provided  $\kappa \leq 1/2\tau$ .

The proof of the [Theorem 4](#) is outlined in [Appendix F.1](#), and the supporting lemmas are proved in the subsequent sections. We now use this guarantee to establish [Proposition A.14](#).

### F.1. Proof Overview of [Theorem 4](#)

**Step 0: Notation & Interpolating segments.** To simplify notation, introduce the maximal operator norms, such that for an  $H$ -tuple of matrices  $\mathbf{X}_{1:H} = (\mathbf{X}_1, \dots, \mathbf{X}[H])$ ,

$$\|\mathbf{X}_{1:H}\|_{\text{max,op}} := \max_{k \in [H]} \|\mathbf{X}_k\|_{\text{op}}.$$

Let us a consider the line segment joining these the parameters

$$\Theta(s) = (\mathbf{A}_{1:K}(s), \mathbf{B}_{1:K}(s)) = (1 - s)\Theta^* + s\hat{\Theta} \tag{F.3}$$

For fixed cost matrices  $\mathbf{Q}$  and  $\mathbf{R}$ , let  $\mathbf{P}_{1:K+1}(s)$  and  $\mathbf{K}_{1:K}(s)$  denote the solution to the finite-time Riccati recursion with parameters  $\Theta(s)$ , where here  $\mathbf{Q}$  also serves as a terminal cost at step  $K+1$ . We let  $\mathbf{P}_{1:K+1}^*$ ,  $\mathbf{K}_{1:K}^*$  be the solution for the truth  $\Theta^*$  and  $\hat{\mathbf{P}}_{1:H}$ ,  $\hat{\mathbf{K}}_{1:H}$  the solution to the Riccati equation with  $\hat{\Theta}$ ; i.e. the certainty equivalent solution. By construction,

$$(\mathbf{P}_{1:K+1}(0), \mathbf{K}_{1:K}(0)) = (\mathbf{P}_{1:K+1}^*, \mathbf{K}_{1:K}^*), \quad (\mathbf{P}_{1:K+1}(1), \mathbf{K}_{1:K}(1)) = (\hat{\mathbf{P}}_{1:K+1}, \hat{\mathbf{K}}_{1:K+1})$$

For all quantity  $\mathbf{X}(s)$  parameterized by  $s \in [0, 1]$ , adopt the shorthand  $\mathbf{X}'(s) := \frac{d}{ds} \mathbf{X}(s)$ .

**Step 1. Self-Bounding ODE Method.** We use an interpolation argument to study the certainty equivalence controller. Our main tool is the following interpolation bound, which states that if the norm of the  $s$ -derivative of a quantity is bounded by the norm of the quantity its self, then that quantity is uniformly bounded on a small enough range.

**Lemma F.2** (Self-Bounding ODE Method, variant of Corollary 3 in (Simchowitz & Foster, 2020)). *Fix dimensions  $d_1, d_2 \geq 1$ , let  $\mathcal{V} \subset \mathbb{R}^{d_1}$ , let  $f : \mathcal{V} \rightarrow \mathbb{R}^{d_2}$  be a  $C^1$  map and let  $\mathbf{v}(s) : [0, 1] \rightarrow \mathcal{V}$  be a  $C^1$  curve defined on  $[0, 1]$ . Finally, let  $\|\cdot\|$  be an arbitrary norm on  $\mathbb{R}^{d_2}$  and suppose that  $c > 0$  and  $p \geq 1$  satisfy*

$$\left\| \frac{d}{ds} f(\mathbf{v}(s)) \right\| \leq c \max\{\|f(\mathbf{v}(s))\|, \|f(\mathbf{v}(0))\|\}^p \quad \forall s \in [0, 1]. \quad (\text{F.4})$$

Then, if  $p > 1$  and if  $\alpha = c(p-1)\|f(\mathbf{v}(0))\|^{p-1}$  satisfies  $\alpha < 1$ , the following bound holds for all  $s \in [0, 1]$ :

$$\|f(\mathbf{v}(s))\| \leq (1-\alpha)^{-\frac{1}{p-1}} \|f(\mathbf{v}(0))\|, \quad \left\| \frac{d}{ds} f(\mathbf{v}(s)) \right\| \leq c(1-\alpha)^{-\frac{p}{p-1}} \|f(\mathbf{v}(0))\|^p$$

**Step 2. Perturbation of  $\mathbf{P}_{1:K+1}(t)$**  First, we show that the Riccati-updates obey the structure of Lemma F.2.

**Lemma F.3.** *Suppose (for simplicity) that  $\lambda_{\min}(\mathbf{Q}), \lambda_{\min}(\mathbf{R}) \geq 1$ . Then, for all  $s \in [0, 1]$*

$$\|\mathbf{P}'_{1:K+1}(s)\|_{\max, \text{op}} \leq 2(\Delta_A + K_A K_B \Delta_B) \|\mathbf{P}_{1:K+1}(s)\|_{\max, \text{op}}^3,$$

Our next result gives uniform bounds on  $\mathbf{P}_{1:K+1}$  and its derivative by invoking Lemma F.2.

**Lemma F.4.** *Suppose (for simplicity) that  $\lambda_{\min}(\mathbf{Q}), \lambda_{\min}(\mathbf{R}) \geq 1$ , and that  $(\Delta_A + K_A K_B \Delta_B) \leq 1/8 \|\mathbf{P}_{1:K+1}^{\text{opt}}(\hat{\Theta})\|_{\max, \text{op}}^2$ . Then, for all  $s \in [0, 1]$ ,*

$$\begin{aligned} \|\mathbf{P}_{1:K+1}(s)\|_{\max, \text{op}} &\leq 1.8 \|\mathbf{P}_{1:K+1}^*\|_{\max, \text{op}} \\ \|\mathbf{P}'_{1:K+1}(s)\|_{\max, \text{op}} &\leq 12(\Delta_A + K_A K_B \Delta_B) \|\mathbf{P}_{1:K+1}^*\|_{\max, \text{op}}^3 \end{aligned}$$

As the gains  $\mathbf{P}_{1:K}(s)$  are explicit function of  $\mathbf{P}_{1:K+1}(s)$ , we obtain the following perturbation bound for the gains.

**Lemma F.5.** *Under the assumptions of Lemma F.4, the following holds:*

$$\|\mathbf{K}_{1:K}^{\text{opt}}(\Theta) - \mathbf{K}_{1:K}^{\text{opt}}(\hat{\Theta})\|_{\max, \text{op}} \leq 20C^3 K_A^3 K_B^2 (1 + \tau C K_B) (\Delta_A + \Delta_B), \quad C := \|\mathbf{P}_{1:K+1}^{\text{opt}}(\Theta)\|$$

### F.1.1. PROOF OF THEOREM 4

**Proof of part (a).** Consider the curve

$$\mathbf{K}_k(s) = (1-s)\mathbf{K}_k^{\text{opt}}(\Theta^*) + s\mathbf{K}_k^{\text{opt}}(\hat{\Theta}).$$

We then note that the curve

$$\mathbf{P}_k^{\text{ce}}(s) = \mathbf{P}_k^{\text{feed}}(\Theta; \mathbf{K}_{1:K}(s))$$

satisfies  $\mathbf{P}_k^{\text{ce}}(0) = \mathbf{P}_k^{\text{ce}}(\Theta^*; \Theta^*) = \mathbf{P}_k^{\text{opt}}(\Theta^*) = \mathbf{P}_k^*$  and  $\mathbf{P}_k^{\text{ce}}(1) = \mathbf{P}_k^{\text{feed}}(\Theta^*; \mathbf{K}_k^{\text{opt}}(\hat{\Theta})) = \mathbf{P}_k^{\text{ce}}(\Theta^*; \hat{\Theta})$

By Definition F.2, we can write  $\mathbf{P}_k^{\text{ce}}(s) = \Lambda_k(s)$ , where  $\Lambda_k$  solve the following Lyapunov equation

$$\begin{aligned} \mathbf{P}_{K+1}^{\text{ce}}(s) &= \mathbf{Q}, \quad \mathbf{P}_k^{\text{ce}}(s) = \mathbf{X}_k(s)^\top \mathbf{P}_{k+1}^{\text{ce}}(s) \mathbf{X}_k(s) + \tau \mathbf{Q}(s) + \mathbf{Y}_k(s) \quad \text{where} \\ \mathbf{X}_k(s) &:= \mathbf{A}_k + \tau \mathbf{B}_k \mathbf{K}_k(s) \quad \text{and} \quad \mathbf{Y}_k(s) := \tau \mathbf{K}(s)^\top \mathbf{R} \mathbf{K}(s). \end{aligned} \quad (\text{F.5})$$

As  $\mathbf{X}'_k(s) = \tau \mathbf{B}_k \mathbf{K}'_k(s)$  and  $\mathbf{Y}'_k(s) = \text{Sym}(\mathbf{K}(s)^\top \mathbf{R} \mathbf{K}'(s))$ , salient term from Proposition F.12 evaluates to

$$\begin{aligned} \Delta(s) &= \max_{j \in [k]} \tau^{-1} (2\|\mathbf{X}_j(s)'\| + \|\mathbf{P}_{1:K+1}^{\text{ce}}(s)\|_{\max, \text{op}}^{-1} \|\mathbf{Y}_j(s)'\|) \\ &= \max_{j \in [k]} \tau^{-1} (2\tau K_B \|\mathbf{K}'(s)\| + \tau \|\mathbf{P}_{1:K+1}^{\text{ce}}(s)\|_{\max, \text{op}}^{-1} \|\mathbf{K}_k(s)\| \|\mathbf{R}\| \|\mathbf{K}'_k(s)\|) \\ &= \max_{j \in [k]} (2K_B + \|\mathbf{P}_{1:K+1}^{\text{ce}}(s)\|_{\max, \text{op}}^{-1} \|\mathbf{K}_k(s)\|) \|\mathbf{K}'_k(s)\| \end{aligned} \quad (\mathbf{R} = \mathbf{I})$$

We further bound

$$\begin{aligned}
 \|K_k(s)\| &= \|(1-s)K_k^{\text{opt}}(\Theta) + sK_k^{\text{opt}}(\hat{\Theta})\| \\
 &\leq \|K_k^{\text{opt}}(\Theta^*)\| \vee \|K_k^{\text{opt}}(\hat{\Theta})\| \\
 &\leq \|(\mathbf{R} + \tau \mathbf{B}_k^\top \mathbf{P}_k^{\text{opt}}(\Theta^*) \mathbf{B}_k)^{-1} (\mathbf{B}_k^\top \mathbf{P}_k^{\text{opt}}(\Theta^*) \mathbf{A}_k)\| \vee \|(\mathbf{R} + \tau \hat{\mathbf{B}}_k^\top \mathbf{P}_k^{\text{opt}}(\hat{\Theta}) \hat{\mathbf{B}})^{-1} (\hat{\mathbf{B}}_k^\top \mathbf{P}_k^{\text{opt}}(\hat{\Theta}) \hat{\mathbf{A}}_k)\| \\
 &\leq \|\mathbf{B}_k^\top \mathbf{P}_k^{\text{opt}}(\Theta^*) \mathbf{A}_k\| \vee \|\hat{\mathbf{B}}_k^\top \mathbf{P}_k^{\text{opt}}(\hat{\Theta}) \hat{\mathbf{A}}_k\| \quad (\mathbf{R} = \mathbf{I}) \\
 &\leq K_B K_A (\|\mathbf{P}_k^{\text{opt}}(\Theta^*)\| \vee \|\mathbf{P}_k^{\text{opt}}(\hat{\Theta})\|) \\
 &\leq 2K_B K_A \|\mathbf{P}_{1:K+1}^{\text{opt}}(\Theta^*)\|_{\max, \text{op}} \quad (\text{Lemma F.5}) \\
 &= 2K_B K_A \|\mathbf{P}_{1:K+1}^{\text{ce}}(0)\|_{\max, \text{op}} \quad (\text{definition of } \mathbf{P}_k^{\text{ce}})
 \end{aligned}$$

Thus,

$$\begin{aligned}
 \Delta(s) &\leq (2K_B + 2K_B K_A \|\mathbf{P}_{1:K+1}^{\text{ce}}(s)\|_{\max, \text{op}}^{-1} \|\mathbf{P}_{1:K+1}^{\text{ce}}(0)\|_{\max, \text{op}}) \max_{j \in [k]} \|K'_k(s)\| \quad (\mathbf{R} = \mathbf{I}) \\
 &\leq 4K_B K_A (1 \vee \|\mathbf{P}_{1:K+1}^{\text{ce}}(s)\|_{\max, \text{op}}^{-1} \|\mathbf{P}_{1:K+1}^{\text{ce}}(0)\|_{\max, \text{op}}) \max_{j \in [k]} \|K'_k(s)\|. \quad (K_A \geq 1)
 \end{aligned}$$

Hence, setting  $\Delta_K := \sup_{s \in [0,1]} \max_{j \in [k]} \|K'_k(s)\|$ , [Proposition F.12](#) implies

$$\begin{aligned}
 \mathbf{P}_k^{\text{ce}}(s)' &\leq \|\mathbf{P}_{1:K+1}^{\text{ce}}(s)\|_{\max, \text{op}}^2 \Delta(s) \\
 &\leq 4K_B K_A \|\mathbf{P}_{1:K+1}^{\text{ce}}(s)\|_{\max, \text{op}}^2 (1 \vee \|\mathbf{P}_{1:K+1}^{\text{ce}}(s)\|_{\max, \text{op}}^{-1} \|\mathbf{P}_{1:K+1}^{\text{ce}}(0)\|_{\max, \text{op}}) \max_{j \in [k]} \Delta_K \\
 &\leq 4K_B K_A (\|\mathbf{P}_{1:K+1}^{\text{ce}}(s)\|_{\max, \text{op}}^2 \vee \|\mathbf{P}_{1:K+1}^{\text{ce}}(0)\|_{\max, \text{op}}^2) \Delta_K
 \end{aligned}$$

Hence, [Lemma F.2](#) implies that as long as  $4K_B K_A \Delta_K \|\mathbf{P}_{1:K+1}^{\text{ce}}(0)\|_{\max, \text{op}} < 1$ , we have

$$\sup_{s \in [0,1]} \|\mathbf{P}_{1:K+1}^{\text{ce}}(s)\|_{\max, \text{op}} \leq (1 - 4K_B K_A \Delta_K \|\mathbf{P}_{1:K+1}^{\text{ce}}(0)\|_{\max, \text{op}})^{-1} \|\mathbf{P}_{1:K+1}^{\text{ce}}(0)\|_{\max, \text{op}}$$

Using  $\mathbf{P}_{1:K+1}^{\text{ce}}(0) = \mathbf{P}_{1:K+1}^{\text{opt}}(\Theta^*)$ ,  $\mathbf{P}_{1:K+1}^{\text{ce}}(1) = \mathbf{P}_{1:K+1}^{\text{ce}}(\Theta^*; \hat{\Theta})$ , and defining the shorthand

$$C := \|\mathbf{P}_{1:K+1}^{\text{opt}}(\Theta^*)\|,$$

we conclude that for any upper bound  $\Delta \geq 4K_B K_A \Delta_K C$  satisfying  $\Delta < 1$ ,

$$\|\mathbf{P}_{1:K+1}^{\text{ce}}(\Theta^*; \hat{\Theta})\|_{\max, \text{op}} \leq (1 - \Delta)^{-1} \|\mathbf{P}_{1:K+1}^{\text{opt}}(\Theta^*)\|_{\max, \text{op}}.$$

By [Lemma F.5](#), it holds that if  $8\|\mathbf{P}_{1:K+1}^{\text{opt}}(\Theta^*)\|_{\max, \text{op}}^2 (\Delta_A + K_A K_B \Delta_B) < 1$ , we can bound. we can take  $\Delta_K \leq 20C^3 K_A^3 K_B^2 (1 + \tau C K_B) (\Delta_A + \Delta_B)$ . Hence, we can bound

$$4K_B K_A \Delta_K C \leq 80C^4 K_A^4 K_B^3 (1 + \tau C K_B) (\Delta_A + \Delta_B) := \Delta_{\text{ce}},$$

which concludes the proof of part (a).

**Proof of part (b).** We bound

$$\begin{aligned}
 \|\mathbf{K}_{1:K}^{\text{opt}}(\hat{\Theta})\|_{\max, \text{op}} &\leq (\|\mathbf{K}_{1:K}^{\text{opt}}(\Theta^*)\|_{\max, \text{op}} + 1/4K_B) \quad (\text{Lemma F.5, definition of } \Delta_{\text{ce}}, \text{ and using } K_B, K_A, C \geq 1) \\
 &= (1/4K_B) + \|(\mathbf{R} + \tau \mathbf{B}_k^\top \mathbf{P}_k^{\text{opt}}(\Theta^*) \mathbf{B}_k)^{-1} \mathbf{B}_k^\top \mathbf{P}_k^{\text{opt}}(\Theta^*) \mathbf{A}_k\|_{\max, \text{op}} \\
 &\quad \quad \quad (\text{Definition of } \mathbf{K}_k^{\text{opt}}, \text{ Definition G.1}) \\
 &= \frac{1}{4K_B} K_B K_A \|\mathbf{P}_{1:K+1}^{\text{opt}}(\Theta^*)\|_{\max, \text{op}} \quad (\text{Definition of } K_A, K_B \text{ in Condition F.1, } \mathbf{R} \succeq \mathbf{I}) \\
 &= \frac{1}{4K_B} + K_B K_A C \quad (\text{Definition of } C) \\
 &\leq \frac{5}{4} K_B K_A C. \quad (C, K_A, K_B \geq 1)
 \end{aligned}$$

**Proof of part (c).** We aim to bound the square of the operator norm of the following term

$$\Phi_{j,k}^{\text{ce}} := (\mathbf{A}_{j-1} + \tau \mathbf{B}_{j-1} \mathbf{K}_{j-1}^{\text{opt}}(\hat{\Theta})) \cdot (\mathbf{A}_{j-2} + \tau \mathbf{B}_{j-2} \mathbf{K}_{j-2}^{\text{opt}}(\hat{\Theta})) \cdots (\mathbf{A}_k + \tau \mathbf{B}_k \mathbf{K}_k^{\text{opt}}(\hat{\Theta})).$$

Using the fact that  $\mathbf{P}_k^{\text{ce}}(\Theta; \hat{\Theta})$  solves the Lyapunov equation Eq. (F.5), it follows from Lemma F.10 that if

$$\kappa_0 := \tau^{-1} \max_k \|\mathbf{I} - (\mathbf{A}_k + \tau \mathbf{B}_k \mathbf{K}_k(1))\|_{\text{op}} = \tau^{-1} \max_k \|\mathbf{I} - (\mathbf{A}_k + \tau \mathbf{B}_k \mathbf{K}_k^{\text{opt}}(\hat{\Theta}))\|_{\text{op}} \leq 1/2\tau,$$

then

$$\|\Phi_{j,k}^{\text{ce}}\|^2 \leq \max\{1, 2\kappa_0\}(1 - \tau\gamma_0)^{j-k}, \quad \gamma_0 := \frac{1}{\|\mathbf{P}_{1:K+1}^{\text{ce}}(\Theta; \hat{\Theta})\|_{\max, \text{op}}} \quad (\text{F.6})$$

From part (a), we can lower bound  $\gamma_0 \geq \gamma := \frac{1 - \Delta_{\text{ce}}}{C}$ . Moreover, we can bound  $\kappa_0$

$$\begin{aligned} \kappa_0 &\leq \tau^{-1} \max_k \|\mathbf{I} - \mathbf{A}_k\| + \max_k \|\mathbf{B}_k\| \|\mathbf{K}_k^{\text{opt}}(\hat{\Theta})\| \\ &\leq \kappa_A + K_B \|\mathbf{K}_{1:K}^{\text{opt}}(\hat{\Theta})\|_{\max, \text{op}} && (\text{Conditions F.1 and F.4}) \\ &\leq \kappa_A + \frac{5}{4} K_B^2 K_A C := \kappa && (\text{Theorem 4(b)}) \end{aligned}$$

As  $\kappa \geq 1$ , we conclude via Eq. (F.6) that

$$\|\Phi_{j,k}^{\text{ce}}\|^2 \leq 2\kappa(1 - \tau\gamma)^{j-k}, \quad \kappa = \kappa_A + \frac{5}{4} K_B^2 K_A C, \quad \gamma = \frac{1 - \Delta_{\text{ce}}}{C}.$$

□

## F.2. Proof of Lemma F.3

To apply the self-bounding ODE method, we bound  $\mathbf{P}'_{1:H}(s)$  in terms of  $\mathbf{P}_{1:K+1}(s)$ . To prove Lemma F.3 Let us first introduce some notation. Further, for simplicity, we shall suppress  $s$  in equations and let  $(\cdot)|_s$  to denote evaluation at  $s$ . With this convention, define the matrices

$$\Xi_k(s) := (\mathbf{A}'_k + \tau \mathbf{K}_k \mathbf{B}'_k)^\top \mathbf{P}_{k+1}(\mathbf{A}_k + \tau \mathbf{B}_k \mathbf{K}_k) \quad (\text{F.7})$$

and define the operator

$$\mathcal{T}_{k+1}(\cdot; s) := \{(\mathbf{A}_k + \mathbf{B}_k \mathbf{K}_k)^\top (\cdot) (\mathbf{A}_k + \mathbf{B}_k \mathbf{K}_k)\}|_s, \quad \text{with the convention } \mathcal{T}_{k+1}(\cdot)|_s = \mathcal{T}_{k+1}(\cdot; s), \quad (\text{F.8})$$

Lastly, we define their compositions

$$\mathcal{T}_{k+i,k} := \mathcal{T}_k(\cdot) \circ \mathcal{T}_{k+1}(\cdot) \circ \cdots \circ \mathcal{T}_{k+i}(\cdot)|_s,$$

with the convention  $\mathcal{T}_{k;k}$  is the identity map. These operators give an expression for the derivatives  $\mathbf{P}'_{k-1}(s)$ .

**Lemma F.6.** For all  $s$ , it holds that

$$\mathbf{P}'_k(s) = \sum_{j=k}^{K+1} \mathcal{T}_{k;j}(\Xi_k + \Xi_k^\top)|_s$$

*Proof.* Let  $\text{Sym}(\mathbf{X}) = \mathbf{X} + \mathbf{X}^\top$ . The Ricatti update (backwards in time) is

$$\mathbf{P}_k = \mathbf{A}_k^\top \mathbf{P}_{k+1} \mathbf{A}_k - \tau (\mathbf{A}_k^\top \mathbf{P}_{k+1} \mathbf{B}_k) (\mathbf{R} + \tau \mathbf{B}_k^\top \mathbf{P}_{k+1} \mathbf{B}_k)^{-1} (\mathbf{A}_k^\top \mathbf{P}_{k+1} \mathbf{B}_k)^\top + \tau \mathbf{Q}$$

Let  $\text{Sym}(\mathbf{X}) := \mathbf{X} + \mathbf{X}^\top$ . Then, we compute

$$\begin{aligned}
 & \mathbf{P}_k(s)' \\
 &= \text{Sym}((\mathbf{A}'_k)^\top \mathbf{P}_{k+1} \mathbf{A}_k) - \tau \text{Sym}((\mathbf{A}'_k)^\top \mathbf{P}_{k+1} \mathbf{B}_k + \mathbf{A}_k^\top \mathbf{P}_{k+1} (\mathbf{B}'_k)) (\mathbf{R} + \tau \mathbf{B}^\top \mathbf{P}_{k+1} \mathbf{B})^{-1} (\mathbf{A}_k^\top \mathbf{P}_{k+1} \mathbf{B}_k)^\top \\
 & \quad + \tau^2 (\mathbf{A}_k^\top \mathbf{P}_{k+1} \mathbf{B}_k) (\mathbf{R} + \tau \mathbf{B}_k^\top \mathbf{P}_{k+1} \mathbf{B}_k)^{-1} (\text{Sym}((\mathbf{B}'_k)^\top \mathbf{P}_{k+1} \mathbf{B}_k)) (\mathbf{R} + \tau \mathbf{B}_k^\top \mathbf{P}_{k+1} \mathbf{B}_k)^{-1} (\mathbf{A}_k^\top \mathbf{P}_{k+1} \mathbf{B}_k)^\top \\
 & \quad + \mathbf{A}_k^\top (\mathbf{P}'_{k+1}) \mathbf{A}_k - \tau \text{Sym}((\mathbf{A}_k^\top (\mathbf{P}'_{k+1}) \mathbf{B}_k) (\mathbf{R} + \tau \mathbf{B}_k^\top \mathbf{P}_{k+1} \mathbf{B}_k)^{-1} (\mathbf{A}_k^\top \mathbf{P}_{k+1} \mathbf{B}_k)^\top) \\
 & \quad + \tau^2 (\mathbf{A}_k^\top \mathbf{P}_{k+1} \mathbf{B}_k) (\mathbf{R} + \mathbf{B}_k^\top \mathbf{P}_{k+1} \mathbf{B}_k)^{-1} (\mathbf{B}_k^\top \mathbf{P}'_{k+1} \mathbf{B}_k) (\mathbf{R} + \mathbf{B}_k^\top \mathbf{P}_{k+1} \mathbf{B}_k)^{-1} (\mathbf{A}_k^\top \mathbf{P}_{k+1} \mathbf{B}_k)^\top \\
 &= \text{Sym}((\mathbf{A}'_k)^\top \mathbf{P}_{k+1} \mathbf{A}_k + \tau((\mathbf{A}'_k)^\top \mathbf{P}_{k+1} \mathbf{B}_k + \mathbf{A}_k^\top \mathbf{P}_{k+1} (\mathbf{B}'_k)) \mathbf{K}_k + \tau^2 \mathbf{K}_k^\top ((\mathbf{B}'_k)^\top \mathbf{P}_{k+1} \mathbf{B}_k) \mathbf{K}_k) \\
 & \quad + \mathbf{A}_k^\top (\mathbf{P}'_{k+1}) \mathbf{A}_k + \tau \text{Sym}((\mathbf{A}_k^\top (\mathbf{P}'_{k+1}) \mathbf{B}_k) \mathbf{K}_k) + \tau^2 (\mathbf{B}_k \mathbf{K}_k) (\mathbf{P}'_{k+1}) (\mathbf{B}_k \mathbf{K}_k) \\
 &= \text{Sym}((\mathbf{A}'_k)^\top \mathbf{P}_{k+1} \mathbf{A}_k + \tau((\mathbf{A}'_k)^\top \mathbf{P}_{k+1} \mathbf{B}_k + \mathbf{A}_k^\top \mathbf{P}_{k+1} (\mathbf{B}'_k)) \mathbf{K}_k + \tau^2 \mathbf{K}_k^\top ((\mathbf{B}'_k)^\top \mathbf{P}_{k+1} \mathbf{B}_k) \mathbf{K}_k) \\
 & \quad + (\mathbf{A}_k + \tau \mathbf{B}_k \mathbf{K}_k)^\top (\mathbf{P}'_{k+1}) (\mathbf{A}_k + \tau \mathbf{B}_k \mathbf{K}_k)^\top
 \end{aligned}$$

where above we use the fact that  $\mathbf{K}_k = -(\mathbf{R} + \tau \mathbf{B}_k^\top \mathbf{P}_{k+1} \mathbf{B}_k)^{-1} (\mathbf{A}_k^\top \mathbf{P}_{k+1} \mathbf{B}_k)^\top$ . Noting that

$$\begin{aligned}
 & \text{Sym}((\mathbf{A}'_k)^\top \mathbf{P}_{k+1} \mathbf{A}_k + \tau((\mathbf{A}'_k)^\top \mathbf{P}_{k+1} \mathbf{B}_k + \mathbf{A}_k^\top \mathbf{P}_{k+1} (\mathbf{B}'_k)) \mathbf{K}_k + \tau^2 \mathbf{K}_k^\top ((\mathbf{B}'_k)^\top \mathbf{P}_{k+1} \mathbf{B}_k) \mathbf{K}_k) \\
 & \stackrel{(i)}{=} \text{Sym}((\mathbf{A}'_k)^\top \mathbf{P}_{k+1} (\mathbf{A}_k + \tau \mathbf{B}_k \mathbf{K}_k) + \tau \mathbf{A}_k^\top \mathbf{P}_{k+1} (\mathbf{B}'_k) \mathbf{K}_k + \tau^2 \mathbf{K}_k^\top ((\mathbf{B}'_k)^\top \mathbf{P}_{k+1} \mathbf{B}_k) \mathbf{K}_k) \\
 & \stackrel{(ii)}{=} \text{Sym}((\mathbf{A}'_k)^\top \mathbf{P}_{k+1} (\mathbf{A}_k + \tau \mathbf{B}_k \mathbf{K}_k) + \tau \mathbf{K}_k^\top (\mathbf{B}'_k) \mathbf{P}_{k+1} \mathbf{A}_k + \tau^2 \mathbf{K}_k^\top ((\mathbf{B}'_k)^\top \mathbf{P}_{k+1} \mathbf{B}_k) \mathbf{K}_k) \\
 &= \text{Sym}((\mathbf{A}'_k)^\top \mathbf{P}_{k+1} (\mathbf{A}_k + \tau \mathbf{B}_k \mathbf{K}_k) + \tau \mathbf{K}_k^\top (\mathbf{B}'_k) \mathbf{P}_{k+1} (\mathbf{A}_k + \tau \mathbf{B}_k \mathbf{K}_k)) \\
 &= \text{Sym}((\mathbf{A}'_k + \tau \mathbf{K}_k \mathbf{B}'_k)^\top \mathbf{P}_{k+1} (\mathbf{A}_k + \tau \mathbf{B}_k \mathbf{K}_k)) := \text{Sym}(\Xi_k)
 \end{aligned}$$

Therefore, we have

$$\mathbf{P}'_k = \mathcal{T}_{k+1}(\mathbf{P}'_{k+1}) + \text{Sym}(\Xi_k)$$

Thus, unfolding the recursion, we have

$$\begin{aligned}
 \mathbf{P}'_k &= \mathcal{T}_{k+1}(\mathbf{P}'_{k+1}) + \text{Sym}(\Xi_k) \\
 &= \mathcal{T}_{k+1}(\mathcal{T}_{k+2}(\mathbf{P}'_{k+1}) + \text{Sym}(\Xi_{k+1})) + \text{Sym}(\Xi_k) \\
 &= \mathcal{T}_{k+1}(\mathcal{T}_{k+2}(\mathbf{P}'_{k+1})) + \mathcal{T}_{k+1}(\text{Sym}(\Xi_{k+1})) + \text{Sym}(\Xi_k) \\
 &= \dots \\
 &= \sum_{j=k}^{K+1} \mathcal{T}_{k;j}(\text{Sym}(\Xi_{k+j})).
 \end{aligned}$$

□

Using this fact, a standard Lyapunov argument gives a generic upper bound on sums of these operators.

**Lemma F.7.** *The operators  $\mathcal{T}_{j,k}(\cdot; s)$  are matrix monotone. Hence, if  $\mathbf{X}_{1:K}$  are any sequence of  $\mathbb{R}^{n \times n}$  matrices,*

$$\left\| \sum_{j=k}^{K+1} \mathcal{T}_{j,k}(\mathbf{X}_j + \mathbf{X}_j^\top; s) \right\|_{\text{op}} \leq \frac{2 \|\mathbf{P}_k\|_{\text{op}} \max_{j \geq k} \|\mathbf{X}_k\|_{\text{op}}}{\tau} \Big|_s$$

Consequently, by [Lemma F.6](#),

$$\|\mathbf{P}'_k(s)\|_{\text{op}} \leq \frac{2 \|\mathbf{P}_k\|_{\text{op}} \max_{j \geq k} \|\Xi_k\|_{\text{op}}}{\tau} \Big|_s. \quad (\text{F.9})$$

*Proof.* This is a direct consequence of rewriting  $\mathbf{P}_k$  as in [Eq. \(F.2\)](#), applying [Lemma F.10\(a\)](#) with  $\mathbf{X}_k = \mathbf{A}_k + \mathbf{B}_k \mathbf{K}_k$ , and upper bounding  $\|\mathbf{X}_j + \mathbf{X}_j^\top\| \leq \|\mathbf{X}_j\|_{\text{op}}$ . □



Finally, let us upper bound the norm of the matrices  $\Xi_k$

**Lemma F.8.**

$$\tau^{-1} \|\Xi_k(s)\| \leq (\Delta_A + \Delta_B K_A K_B) \|\mathbf{P}_{1:K+1}\|_{\max, \text{op}}^2.$$

*Proof of Lemma F.8.* Recall

$$\begin{aligned} \Xi_k &:= (\mathbf{A}'_k + \tau \mathbf{K}_k \mathbf{B}'_k)^\top \mathbf{P}_{k+1} (\mathbf{A}_k + \tau \mathbf{B}_k \mathbf{K}_k), \\ \|\mathbf{K}_k\| &= \|(\mathbf{R} + \tau \mathbf{B}_k^\top \mathbf{P}_{k+1} \mathbf{B}_k)^{-1} \mathbf{B}_k^\top \mathbf{P}_{k+1} \mathbf{A}_k\| \\ &= \lambda_{\min}(\mathbf{R})^{-1} \|\mathbf{B}_k^\top \mathbf{P}_{k+1} \mathbf{A}_k\| \\ &\leq \lambda_{\min}(\mathbf{R})^{-1} \|\mathbf{B}_k\| \|\mathbf{A}_k\| \|\mathbf{P}_{k+1}\| \\ &\leq \lambda_{\min}(\mathbf{R})^{-1} K_A K_B \|\mathbf{P}_{k+1}\| \\ &\leq K_A K_B \|\mathbf{P}_{k+1}\| \end{aligned} \tag{F.10}$$

Next,

$$\begin{aligned} &\tau^{-1} \|(\mathbf{A}'_k + \tau \mathbf{K}_k \mathbf{B}'_k)^\top \mathbf{P}_{k+1} (\mathbf{A}_k + \tau \mathbf{B}_k \mathbf{K}_k)\| \\ &\leq \tau^{-1} (\|\mathbf{A}'_k\| + \tau \|\mathbf{B}'_k\| \|\mathbf{K}_k\|) \|\mathbf{P}_{k+1}^{\frac{1}{2}}\| \|\mathbf{P}_{k+1}^{\frac{1}{2}} (\mathbf{A}_k + \mathbf{B}_k \mathbf{K}_k)\| \\ &\stackrel{(i)}{\leq} \tau^{-1} (\|\mathbf{A}'_k\| + \tau \|\mathbf{B}'_k\| \|\mathbf{K}_k\|) \|\mathbf{P}_{k+1}^{\frac{1}{2}}\| \|\mathbf{P}_k^{\frac{1}{2}}\| \\ &\leq \tau^{-1} (\|\mathbf{A}'_k\| + \tau \|\mathbf{B}'_k\| K_A K_B \|\mathbf{P}_{k+1}\|) \|\mathbf{P}_{k+1}^{\frac{1}{2}}\| \|\mathbf{P}_k^{\frac{1}{2}}\| \\ &\leq \tau^{-1} (\|\mathbf{A}'_k\| + \tau \|\mathbf{B}'_k\| K_A K_B \|\mathbf{P}_{1:K+1}\|_{\max, \text{op}}) \|\mathbf{P}_{1:K+1}\|_{\max, \text{op}} \\ &\leq \tau^{-1} (\|\mathbf{A}'_k\| + \tau \|\mathbf{B}'_k\| K_A K_B) \|\mathbf{P}_{1:K+1}\|_{\max, \text{op}}^2 \quad (\|\mathbf{P}_{1:K+1}\|_{\max, \text{op}} \geq \|\mathbf{Q}\| \geq 1) \\ &\leq (\Delta_A + \Delta_B K_A K_B) \|\mathbf{P}_{1:K+1}\|_{\max, \text{op}}^2, \end{aligned}$$

where in (i) we use Eq. (F.2), which under the present notation gives

$$\mathbf{P}_k = (\mathbf{A}_k + \mathbf{B}_k \mathbf{K}_k) \mathbf{P}_{k+1} (\mathbf{A}_k + \mathbf{B}_k \mathbf{K}_k) + \tau \mathbf{Q},$$

and since  $\mathbf{P}_k \succeq \tau \mathbf{Q}$ ,

$$\|\mathbf{P}_{k+1}^{\frac{1}{2}} (\mathbf{A}_k + \mathbf{B}_k \mathbf{K}_k)\|^2 = \|(\mathbf{A}_k + \mathbf{B}_k \mathbf{K}_k) \mathbf{P}_{k+1} (\mathbf{A}_k + \mathbf{B}_k \mathbf{K}_k)\| = \|\mathbf{P}_k - \tau \mathbf{Q}\| \leq \|\mathbf{P}_k\|.$$

□

*Proof of Lemma F.3.* From Eq. (F.9) in Lemma F.7, followed by Lemma F.8, we have for  $k \in [K]$  that

$$\begin{aligned} \|\mathbf{P}'_k(s)\|_{\text{op}} &\leq \tau^{-1} 2 \|\mathbf{P}_k\|_{\text{op}} \max_{j \geq k} \|\Xi_j\|_{\text{op}} \Big|_s \\ &\leq 2 (\Delta_A + \Delta_B K_A K_B) \|\mathbf{P}_{1:K+1}(0)\|_{\max, \text{op}}^3. \end{aligned}$$

□

### F.3. Proof of Lemma F.4

Let us apply the Lemma F.3 with  $\mathbf{v}(s) = \Theta(s) = (\mathbf{A}_{1:K}(s), \mathbf{B}_{1:K}(s))$  as in Eq. (F.3) and  $f$  as the mapping from  $(\mathbf{A}_{1:K}, \mathbf{B}_{1:K}) \rightarrow \mathbf{P}_{1:K+1}$ . This map is algebraic and thus  $\mathcal{C}^1$ , and  $\mathbf{v}(s)$  is also  $\mathcal{C}^1$  as it is linear. Finally, take  $\|\cdot\|$  to be  $\|\cdot\|_{\max, \text{op}}$ , take  $g(z) = cz^p$ , where  $p = 3$  and  $c = 2(\Delta_A + K_A K_B \Delta_B)$ . The corresponding  $\alpha$  in Lemma F.3 is  $\alpha = c(p-1) \|f(\mathbf{v}(0))\|^{p-1} = 2(\Delta_A + K_A K_B \Delta_B) \|\mathbf{P}_{1:K+1}(0)\|^2$ , then, if  $\alpha \leq 1/4$ , i.e.  $(\Delta_A + K_A K_B \Delta_B) \leq /8 \|\mathbf{P}_{1:K+1}(0)\|^2$ , we have by Lemma F.3 that

$$\|\mathbf{P}_{1:K}(s)\|_{\max, \text{op}} \leq (1 - \alpha)^{-\frac{1}{p-1}} \|\mathbf{P}_{1:K}^*\|_{\max, \text{op}} \leq (4/3)^2 \|\mathbf{P}_{1:K}^*\|_{\max, \text{op}} \leq 1.8 \|\mathbf{P}_{1:K}^*(0)\|_{\max, \text{op}}$$

and

$$\|\mathbf{P}_{1:K}(s)'\|_{\max, \text{op}} \leq 2(4/3)^6 (\Delta_A + K_A K_B \Delta_B) \|\mathbf{P}_{1:K}^*\|_{\max, \text{op}}^3 \leq 12(\Delta_A + K_A K_B \Delta_B) \|\mathbf{P}_{1:K}^*(0)\|_{\max, \text{op}}^3$$

#### F.4. Perturbation on the gains (Lemma F.5)

*Proof.* Observe that

$$\mathbf{K}_k = -(\mathbf{R} + \tau \mathbf{B}_k^\top \mathbf{P}_{k+1} \mathbf{B}_k)^{-1} \mathbf{B}_k^\top \mathbf{P}_{k+1} \mathbf{A}_k$$

Therefore,

$$\mathbf{K}'_k = (\mathbf{R} + \tau \mathbf{B}_k^\top \mathbf{P}_{k+1} \mathbf{B}_k)^{-1} \cdot (\mathbf{R} + \tau \mathbf{B}_k^\top \mathbf{P}_{k+1} \mathbf{B}_k)' \cdot \underbrace{(\mathbf{R} + \mathbf{B}_k^\top \mathbf{P}_{k+1} \mathbf{B}_k)^{-1} \mathbf{B}_k^\top \mathbf{P}_{k+1} \mathbf{A}_k}_{=\mathbf{K}_k} \quad (\text{F.11})$$

$$- (\mathbf{R} + \tau \mathbf{B}_k^\top \mathbf{P}_{k+1} \mathbf{B}_k)^{-1} (\mathbf{B}_k^\top \mathbf{P}_{k+1} \mathbf{A}_k)'. \quad (\text{F.12})$$

Introduce the constant  $C := \|\mathbf{P}_{1:K+1}^*\|_{\max, \text{op}}$ . Using  $\mathbf{R} \succeq \mathbf{I}$ , we have

$$\begin{aligned} \|\mathbf{K}'_k\| &= \|(\mathbf{R} + \tau \mathbf{B}_k^\top \mathbf{P}_{k+1} \mathbf{B}_k)' \|\mathbf{K}\| + \|(\mathbf{B}_k^\top \mathbf{P}_{k+1} \mathbf{A}_k)'\| \\ &= \tau \|(\mathbf{B}_k^\top \mathbf{P}_{k+1} \mathbf{B}_k)'\| \|\mathbf{K}_k\| + \|(\mathbf{B}_k^\top \mathbf{P}_{k+1} \mathbf{A}_k)'\| \\ &\leq \tau (2 \|\mathbf{B}_k\| \|\mathbf{P}_{k+1}\| \|\mathbf{B}'_k\| + \|\mathbf{B}_k\|^2 \|\mathbf{P}'_{k+1}\|) \|\mathbf{K}_k\| + (\|\mathbf{B}'_k\| \|\mathbf{A}_k\| + \|\mathbf{A}'_k\| \|\mathbf{B}_k\|) \|\mathbf{P}_{k+1}\| + \|\mathbf{P}'_{k+1}\| \|\mathbf{A}_k\| \\ &\leq \tau (2 K_B \Delta_B \|\mathbf{P}_{k+1}\| + K_B^2 \|\mathbf{P}'_{k+1}\|) \|\mathbf{K}_k\| + (\Delta_B K_A + \tau \Delta_A K_B) \|\mathbf{P}_{k+1}\| + \|\mathbf{P}'_{k+1}\| K_A K_B \\ &\leq \tau (2 K_B^2 K_A \Delta_B \|\mathbf{P}_{k+1}\|^2 + K_B^3 K_A \|\mathbf{P}'_{k+1}\| \|\mathbf{P}_{k+1}\| + \Delta_A K_B \|\mathbf{P}_{k+1}\|) \\ &\quad + (\|\mathbf{P}_{k+1}\| \Delta_B K_A + \|\mathbf{P}'_{k+1}\| K_A K_B) \\ &\leq \tau (2 \cdot 1.8^2 K_B^2 K_A \Delta_B C^2 + 12(\Delta_A + K_A K_B \Delta_B) K_B^3 K_A C^4 + 1.8 \Delta_A K_B C) \\ &\quad + (1.8 C \Delta_B K_A + 12(\Delta_A + K_A K_B \Delta_B) K_A K_B C) \quad (\text{Lemma F.4}) \\ &\leq \tau C^4 (2 \cdot 1.8^2 K_B^2 K_A \Delta_B + 12(\Delta_A + K_A K_B \Delta_B) K_B^3 K_A + 1.8 \Delta_A K_B) \\ &\quad + C^3 (1.8 \Delta_B K_A + 12(\Delta_A + K_A K_B \Delta_B) K_A K_B) \\ &\leq \tau C^4 K_B^3 K_A^2 (1.8^2 \cdot 2 \Delta_B + 12(\Delta_A + \Delta_B) + 1.8 \Delta_A) \\ &\quad + C^3 K_A^2 K_B^2 (1.8 \Delta_B + 12(\Delta_A + \Delta_B)) \\ &\leq C^3 K_A^2 K_B^2 (1 + \tau C K_B) (1.8^2 \cdot 2 \Delta_B + 12(\Delta_A + \Delta_B) + 1.8 \Delta_A) \\ &\leq 20 C^3 K_A^3 K_B^2 (1 + \tau C K_B) (\Delta_A + \Delta_B). \end{aligned}$$

It follows from Taylor's theorem that  $\|\mathbf{K}_k^{\text{opt}}(\Theta) - \mathbf{K}_k^{\text{opt}}(\hat{\Theta})\| \leq 20 C^3 K_A^3 K_B^2 (1 + \tau C K_B) (\Delta_A + \Delta_B)$ . The result follows.  $\square$

#### F.5. Proof of Lemma F.2

Lemma F.2 is a special case of Simchowitz & Foster (2020, Corollary 3). To check this, we first establish the following special case of Theorem 13 in (Simchowitz & Foster, 2020).

**Lemma F.9** (Comparison Lemma). *Fix dimensions  $d_1, d_2 \geq 1$ , let  $\mathcal{V} \subset \mathbb{R}^{d_1}$ , let  $f : \mathcal{V} \rightarrow \mathbb{R}^{d_2}$  be a  $C^1$  map and let  $\mathbf{v}(s) : [0, 1] \rightarrow \mathcal{V}$  be a  $C^1$  curve defined on  $[0, 1]$ . Suppose that  $g(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$  is non-negative and non-decreasing scalar function, and  $\|\cdot\|$  be an arbitrary norm on  $\mathbb{R}^{d_2}$  such that*

$$\left\| \frac{d}{ds} f(\mathbf{v}(s)) \right\| \leq g(\|f(\mathbf{v}(s))\|) \quad (\text{F.13})$$

Finally, let  $\eta > 0$  and  $\tilde{g} : \mathbb{R} \rightarrow \mathbb{R}$  be a scalar function such that (a) for all  $z \geq \|v(0)\|$ ,  $\tilde{g}(z) \geq \eta + g(z)$  and (b) the following scalar ODE has a solution on  $[0, 1]$ :

$$z(0) = \|v(0)\| + \eta, \quad z'(s) = \tilde{g}(z(s))$$

Then, it holds that

$$\forall s \in [0, 1], \quad \|f(\mathbf{v}(s))\| \leq z(s).$$

*Proof.* Theorem 13 in (Simchowitz & Foster, 2020) proves a more general result for implicit ODEs, such as those that arise in infinite-horizon Riccati equations. We do not need these complication here, so we specialize their result. Define the function  $\mathcal{F} : \mathbb{R}^{d_1} \times \mathbb{R}^{d_2} \rightarrow \mathbb{R}^{d_2}$  via  $\mathcal{F}(\mathbf{v}, \mathbf{w}) = f(\mathbf{v}(s)) - \mathbf{w}$ . It is then clear that  $\tilde{\mathbf{w}}(s) = f(\tilde{\mathbf{v}}(s))$  is the unique solution to  $\mathcal{F}(\tilde{\mathbf{v}}(s), \tilde{\mathbf{w}}(s)) = 0$  for any  $\mathcal{C}^1$  curve  $\tilde{\mathbf{v}}(z)$ ; since  $f$  is  $\mathcal{C}^1$ , any such solution  $\tilde{\mathbf{w}}$  is also  $\mathcal{C}^1$ . Thus,  $\mathcal{F}$  is a “valid implicit function” on  $\mathcal{V}$  in the sense of Simchowitz & Foster (2020, Definition 3.2) with  $\cdot$ . Moreover, by Eq. (F.13),  $(\mathcal{F}, \mathcal{V}, g, \|\cdot\|, \mathbf{v}(\cdot))$  form a self-bounding tuple in the sense of Simchowitz & Foster (2020, Definition 3.3). The result now follows from Simchowitz & Foster (2020, Theorem 13).  $\square$

*Proof of Lemma F.2.* Take  $g(z) = cz^p$ . Define  $h_\eta = c(z+\eta)^p$ . For any  $\eta > 0$ , there exist an  $\eta'$  such that  $h_{\eta'}(z) \leq g(z) + \eta$  for  $z \geq z_0$ . Moreover, as  $\eta$  approaches 0, we can take  $\eta' \rightarrow 0$  as well. Solving the ODE  $z(0) = \|f(\mathbf{v}(0))\| + \eta$  and  $z'(s) = c(z + \eta)^p$ , we see the solution is given by

$$\frac{dz}{(z + \eta')^p} = c ds.$$

As  $z'(s) \geq 0$ , it suffices to bound  $z(1)$ . For  $p > 1$ , the solution to this ODE when it exists satisfies

$$\frac{1}{(p-1)(\|f(\mathbf{v}(0))\| + \eta + \eta')^{p-1}} - \frac{1}{(p-1)(z(1) + \eta')^{p-1}} = c$$

Rearranging and setting  $\eta, \eta' \rightarrow 0$  lets check that, as long as  $\frac{1}{(p-1)\|f(\mathbf{v}(0))\|^{p-1}} > c$ , Lemma F.9 yields

$$\max_{s \in [0, 1]} \|f(\mathbf{v}(s))\| \leq \left( \frac{1}{\|f(\mathbf{v}(0))\|^{p-1}} - (p-1)c \right)^{-\frac{1}{p-1}} = (1 - \alpha)^{-\frac{1}{p-1}} \|f(\mathbf{v}(0))\|.$$

For  $p = 1$ , we instead get

$$\ln(z(1) + \eta') - \ln(\|f(\mathbf{v}(0))\| + \eta + \eta') = c$$

Again, taking  $\eta', \eta \rightarrow 0$ , Lemma F.9 yields

$$\max_{s \in [0, 1]} \|f(\mathbf{v}(s))\| \leq \exp(c + \ln(\|f(\mathbf{v}(0))\|)) = \|f(\mathbf{v}(0))\| e^c.$$

$\square$

## F.6. Perturbation bounds for Lyapunov Solutions

**Lemma F.10** (Basic Lyapunov Theory). *Let  $\mathbf{X}_{1:K}$  and  $\mathbf{Y}_{1:K}$  be a sequence of matrices of appropriate dimension. Suppose that  $\mathbf{Y}_k \succeq 0$ , and let  $\mathbf{Q} \succeq \mathbf{I}$ . Define  $\Lambda_k$  as via the solution to the Lyapunov recursion*

$$\Lambda_{K+1} = \mathbf{Q}, \quad \Lambda_k = \mathbf{X}_k^\top \Lambda_{k+1} \mathbf{X}_k + \tau \mathbf{Q} + \mathbf{Y}_k.$$

*and define the matrix  $\Phi_{j+1,k} := (\mathbf{X}_j \cdot \mathbf{X}_{j-1} \cdots \mathbf{X}_{k+1} \cdot \mathbf{X}_k)$ , with the convention  $\Phi_{k,k} = \mathbf{I}$ , let and define the operator  $\mathcal{T}_{j,k}(\cdot) = \Phi_{j,k}^\top(\cdot)\Phi_{j,k}$ . Then*

(a) *For any symmetric matrices  $\mathbf{Z}_j$ , we have*

$$\tau \sum_{j=k}^K \mathcal{T}_{k,j}(\mathbf{Z}_j) \preceq \max_{j=k}^K \|\mathbf{Z}_j\| \Lambda_k.$$

(b) *If, in addition,  $\max_k \|\mathbf{I} - \mathbf{X}_k\|_{\text{op}} \leq \kappa \tau$  for some  $\kappa \leq 1/2\tau$ ,  $\lambda_{\min}(\Lambda_k) \geq \min\{\frac{1}{2\kappa}, 1\}$*

(c) Under the condition in part (b), we have

$$\|\Phi_{j,k}\|^2 \leq \max\{1, 2\kappa\} \|\Lambda_{1:K+1}\|_{\max, \text{op}} (1 - \tau\gamma)^{j-k}, \quad \gamma := \frac{1}{\|\Lambda_{1:K+1}\|_{\max, \text{op}}}.$$

In particular,  $\|\mathcal{T}_{k,j}(\mathbf{Z}_j)\| \leq \max\{1, 2\kappa\} \|\Lambda_{1:K+1}\|_{\max, \text{op}} \|\mathbf{Z}_k(j)\|$ .

*Proof.* We begin with part (a). By unfolding the recursion, we get

$$\begin{aligned} \Lambda_k &= \mathbf{X}_k^\top \Lambda_{k+1} \mathbf{X}_k + \tau(\mathbf{Q} + \tau^{-1} \mathbf{Y}_k) \\ &= (\tau \mathcal{T}_{k,k}(\mathbf{Q} + \tau^{-1} \mathbf{Y}_k) + \mathcal{T}_{k+1,k}(\Lambda_{k+1})) \\ &= \tau \sum_{j=k}^K \mathcal{T}_{k,j}(\mathbf{Q} + \tau^{-1} \mathbf{Y}_k) + \mathcal{T}_{K+1,k}(\Lambda_{K+1}) \\ &\succeq \tau \sum_{j=k}^K \mathcal{T}_{k,j}(\mathcal{I}) \end{aligned}$$

where in the last line, we use  $\mathbf{Q} + \tau^{-1} \mathbf{Y}_k \succeq \mathbf{Q} \succeq \mathbf{I}$ . As  $\mathcal{T}_{k,j}(\cdot)$  is a matrix monotone operator, we have that symmetric matrix  $\mathbf{Z}$ , we have

$$\tau \sum_{j=k}^K \mathcal{T}_{j,k}(\mathbf{Z}_j) \preceq \tau \max_{j=k}^K \|\mathbf{Z}_j\| \sum_{j=k}^K \mathcal{T}_{j,k}(\mathbf{I}) \preceq \Lambda_k,$$

and similarly for  $-\tau \sum_{j=k}^K \mathcal{T}_{j,k}(\mathbf{Z}_j)$ .

**Part (b).** We argue part (b) by induction backwards on  $k$ , noting that  $k = K + 1$  is immediate. We have

$$\Lambda_k \succeq \mathbf{X}_k^\top \Lambda_{k+1} \mathbf{X}_k + \tau \mathbf{Q} \succeq \mathbf{X}_k^\top \Lambda_{k+1} \mathbf{X}_k + \tau$$

Hence,

$$\begin{aligned} \lambda_{\min}(\Lambda_k) &\geq \lambda_{\min}(\Lambda_{k+1}) \sigma_{\min}(\mathbf{X}_k)^2 + \tau \\ &\geq \lambda_{\min}(\Lambda_{k+1}) (1 - \|\mathbf{X}_k - \mathbf{I}\|)^2 + \tau \\ &\geq \lambda_{\min}(\Lambda_{k+1}) (1 - \kappa\tau)^2 + \tau \\ &\geq \lambda_{\min}(\Lambda_{k+1}) (1 - 2\kappa\tau) + \tau \end{aligned}$$

Applying the inductive hypothesis, we see that the above is at least

$$\lambda_{\min}(\Lambda_k) \geq \frac{1}{2\kappa} (1 - 2\kappa\tau) + \tau \geq \frac{1}{2\kappa}, \text{ as needed.}$$

**Part (c).** We have

$$\mathbf{X}_j^\top \Lambda_{j+1} \mathbf{X}_j = \Lambda_j - \tau(\mathbf{Q} + \mathbf{Y}_k) \preceq \Lambda_j (1 - \tau \Lambda_j^{-1/2} \mathbf{Q} \Lambda_j^{-1/2}) \preceq \Lambda_j (1 - \tau\gamma),$$

where we recall  $\gamma = 1/\|\Lambda_{1:K+1}\|_{\max, \text{op}}$  and use  $\mathbf{Q} \succeq \mathbf{I}$ . By unfolding the bound, we find

$$\|\Phi_{j+1,k}^\top \Lambda_{j+1} \Phi_{j+1,k}\| \leq (1 - \tau\gamma)^{j+1-k} \|\Lambda_k\|.$$

Hence,

$$\|\Phi_{j+1,k}\|^2 \leq \lambda_{\min}(\Lambda_{j+1})^{-1} \|\Phi_{j+1,k}^\top \Lambda_{j+1} \Phi_{j+1,k}\| \leq 2 \|\Lambda_{j+1}\| \kappa (1 - \tau\gamma)^{j+1-k}.$$

Moreover, as  $\kappa \geq 1$ , the bound also applies to  $\Phi_{j+1,j+1} = \mathbf{I}$ . □

**Lemma F.11** (Formula for Lyapunov Curve Derivatives). *Consider curves  $\mathbf{X}_{1:K}(s), \mathbf{Y}_{1:K}(s)$ , let  $\mathbf{Q} \succeq \mathbf{I}$ , and define*

$$\Lambda_{K+1}(s) = \mathbf{Q}, \quad \Lambda_k(s) = \mathbf{X}_k(s)^\top \Lambda_{k+1} \mathbf{X}_k(s) + \tau \mathbf{Q} + \mathbf{Y}_k(s)$$

Again, let  $\Phi_{k,j} := (\mathbf{X}_{j-1} \cdot \mathbf{X}_{j-2} \cdots \mathbf{X}_{k+1} \cdot \mathbf{X}_k)$ , with the convention  $\Phi_{k,k} = \mathbf{I}$ , define the operator  $\mathcal{T}_{k,j}(\cdot) = \Phi_{k,j}^\top(\cdot)\Phi_{k,j}$ . Then,

$$\Lambda'_k = \sum_{j=k}^K \mathcal{T}_{k,j}(\Omega_k), \quad \Omega_k := \text{Sym}(\mathbf{X}_k^\top \Lambda_{k+1} \mathbf{X}'_k) + \mathbf{Y}'_k.$$

*Proof.* We compute

$$\Lambda'_k = \underbrace{\text{Sym}(\mathbf{X}_k^\top \Lambda_{k+1} \mathbf{X}'_k)}_{=\Omega_k} + \mathbf{Y}'_k + \mathbf{X}_k(s)^\top \Lambda_{k+1} \mathbf{X}_k(s).$$

The result follows by unfolding the recursion, with the base case  $\Lambda'_{K+1} = \frac{d}{ds} \mathbf{Q} = 0$ .  $\square$

We now state our Lyapunov perturbation bound:

**Proposition F.12** (Lyapunov Function Perturbation). *Consider curves  $\mathbf{X}_{1:K}(s), \mathbf{Y}_{1:K}(s)$ , and define for  $\mathbf{Q} \succeq \mathbf{I}$*

$$\Lambda_{K+1}(s) = \mathbf{Q}, \quad \Lambda_k(s) = \mathbf{X}_k(s)^\top \Lambda_{k+1} \mathbf{X}_k(s) + \tau \mathbf{Q} + \mathbf{Y}_k(s)$$

Then,

$$\begin{aligned} \Lambda_k(s)' &\leq \|\Lambda_{1:K+1}(s)\|_{\max, \text{op}}^2 \Delta(s), \quad \text{where} \\ \Delta(s) &= \max_{j \in [k]} \tau^{-1} (2\|\mathbf{X}_j(s)'\| + \|\Lambda_{1:K+1}(s)\|_{\max, \text{op}}^{-1} \|\mathbf{Y}_j(s)'\|) \end{aligned}$$

Moreover, as long as  $\|\Lambda_{1:K+1}(0)\|_{\max, \text{op}} \sup_{s \in [0,1]} \Delta(s) < 1$ ,

$$\max_{s \in [0,1]} \|\Lambda_{1:K+1}(s)\|_{\max, \text{op}} \leq \left(1 - \|\Lambda_{1:K+1}(0)\|_{\max, \text{op}} \sup_{s \in [0,1]} \Delta(s)\right)^{-1} \|\Lambda_{1:K+1}(0)\|_{\max, \text{op}}$$

The above bound also holds when  $\Delta(s)$  is replaced by the simpler term

$$\tilde{\Delta}(s) := \max_{j \in [k]} \tau^{-1} (2\|\mathbf{X}_j(s)'\| + \|\mathbf{Y}_j(s)'\|) \quad (\text{F.14})$$

*Proof.* We write

$$\Lambda'_k = \sum_{j=k}^K \mathcal{T}_{k,j}(\Omega_k), \quad \Omega_k := \text{Sym}(\mathbf{X}_k^\top \Lambda_{k+1} \mathbf{X}'_k) + \mathbf{Y}'_k \quad (\text{F.15})$$

We have

$$\begin{aligned} \|\Omega_k\| &\leq 2\|\mathbf{X}_k^\top \Lambda_{k+1}\| \|\mathbf{X}'_k\| + \|\mathbf{Y}'_k\| \\ &\leq 2\|\mathbf{X}_k^\top \Lambda_{k+1}^{\frac{1}{2}}\| \|\Lambda_{k+1}^{\frac{1}{2}}\| \|\mathbf{X}'_k\| + \|\mathbf{Y}'_k\|. \end{aligned}$$

Observe that  $\|\mathbf{X}_k^\top \Lambda_{k+1}^{\frac{1}{2}}\| = \|\mathbf{X}_k^\top \Lambda_{k+1} \mathbf{X}_k\|^{\frac{1}{2}}$ . As  $0 \preceq \mathbf{X}_k^\top \Lambda_{k+1} \mathbf{X}_k = \Lambda_k - \mathbf{Y}_k \preceq \Lambda_k$ , we conclude,

$$\begin{aligned} \|\Omega_k\| &\leq 2\|\Lambda_k\|^{\frac{1}{2}} \|\Lambda_{k+1}^{\frac{1}{2}}\| \|\mathbf{X}'_k\| + \|\mathbf{Y}'_k\| \\ &\leq \|\Lambda_{1:K+1}\|_{\max, \text{op}} 2\|\mathbf{X}'_k\| + \|\Lambda_{1:K+1}(s)\|_{\max, \text{op}}^{-1} \|\mathbf{Y}'_k\| \\ &\leq \|\Lambda_{1:K+1}\|_{\max, \text{op}} (2\|\mathbf{X}'_k\| + \|\Lambda_{1:K+1}(s)\|_{\max, \text{op}}^{-1} \|\mathbf{Y}'_k\|) \\ &\leq \|\Lambda_{1:K+1}\|_{\max, \text{op}} \max_{j \in [K]} (2\|\mathbf{X}'_j\| + \|\Lambda_{1:K+1}(s)\|_{\max, \text{op}}^{-1} \|\mathbf{Y}'_j\|) \leq \tau \|\Lambda_{1:K+1}\|_{\max, \text{op}} \Delta(s). \end{aligned}$$

Thus, from Eq. (F.15) and Lemma F.10, we conclude

$$\begin{aligned}\Lambda_k(s)' &\leq \tau^{-1} \cdot \tau \|\Lambda_k\| \|\Lambda_{1:K+1}(s)\|_{\max, \text{op}} \Delta(s) \\ &\leq \|\Lambda_{1:K+1}(s)\|_{\max, \text{op}}^2 \Delta(s)\end{aligned}$$

The final result follows by applying Lemma F.2 with  $p = 2$ ,  $c = \max_{s \in [0,1]} \Delta$ , and  $\alpha = \Delta(s) \|\Lambda_{1:K+1}(0)\|_{\max, \text{op}}$ . That Eq. (F.14) follows from the fact that if  $\mathbf{Q} \succeq \mathbf{I}$ ,  $\|\Lambda_{1:K+1}(s)\|_{\max, \text{op}} \geq 1$ .  $\square$

**Lemma F.13** (Average Perturbation). *Let  $\kappa \leq 1/2\tau$ . Consider a curve  $\mathbf{X}_{1:K}(s)$  such  $\max_k \sup_{s \in [0,1]} \|\mathbf{I} - \mathbf{X}_k(s)\|_{\text{op}} \leq \kappa\tau$ . Then,*

$$\Lambda_{K+1}(s) = \mathbf{I}, \quad \Lambda_k(s) = \mathbf{X}_k(s)^\top \Lambda_{k+1} \mathbf{X}_k(s) + \tau \mathbf{I}.$$

Fix

$$\Delta_{\text{sum}} = 3 \sup_{s \in [0,1]} \max\{1, 2\kappa\} \sum_{k=1}^K \|\mathbf{X}'_k(s)\|.$$

Then, as long as  $\|\Lambda_{1:K+1}(0)\|_{\max, \text{op}} \Delta_{\text{sum}} < 1$ , we have

$$\|\Lambda_{1:K+1}(1)\|_{\max, \text{op}} \leq (1 - \|\Lambda_{1:K+1}(0)\|_{\max, \text{op}} \Delta_{\text{sum}})^{-1} \|\Lambda_{1:K+1}(0)\|_{\max, \text{op}}.$$

*Proof.* We have that

$$\Lambda'_k = \sum_{j=k}^K \mathcal{T}_{k,j}(\Omega_k), \quad \Omega_k := \text{Sym}(\mathbf{X}_k^\top \Lambda_{k+1} \mathbf{X}'_k), \quad (\text{F.16})$$

so by Lemma F.10(c),

$$\begin{aligned}\|\Lambda'_k\|_{\text{op}} &\leq \|\Lambda_{1:K+1}\|_{\max, \text{op}} \max\{1, 2\kappa\} \sum_{j=k}^K \|\Omega_j\| \\ &\leq 2 \|\Lambda_{1:K+1}\|_{\max, \text{op}} \max\{1, 2\kappa\} \sum_{j=k}^K \|\mathbf{X}_k\| \|\mathbf{X}'_k\| \|\Lambda_{k+1}\| \\ &\leq 2 \|\Lambda_{1:K+1}\|_{\max, \text{op}}^2 \max\{1, 2\kappa\} \sum_{j=k}^K \|\mathbf{X}'_k\| \|\mathbf{X}_k\| \\ &\leq 2 \|\Lambda_{1:K+1}\|_{\max, \text{op}}^2 \max\{1, 2\kappa\} \underbrace{(1 + \kappa\tau)}_{\leq \frac{3}{2}} \sum_{j=k}^K \|\mathbf{X}'_k\| \\ &\leq \Delta_{\text{sum}} \|\Lambda_{1:K+1}\|_{\max, \text{op}}^2.\end{aligned}$$

The result now follows from Lemma F.2.  $\square$

## G. Instantiations of Certainty Equivalence Bound

**Definition G.1.** Given a sequence of gains  $\tilde{\mathbf{K}}_{1:K} \in (\mathbb{R}^{d_u \times d_x})^K$ , we define the discrete cost-to-go matrix as

$$\mathbf{P}_{K+1}^\pi = \mathbf{I}, \quad \mathbf{P}_k^\pi [\tilde{\mathbf{K}}_{k:K}] = (\mathbf{A}_{\text{ol},k}^\pi + \mathbf{B}_{\text{ol},k}^\pi \tilde{\mathbf{K}}_k)^\top \mathbf{P}_{k+1}^\pi [\tilde{\mathbf{K}}_{k+1:K}] (\mathbf{A}_{\text{ol},k}^\pi + \mathbf{B}_{\text{ol},k}^\pi \tilde{\mathbf{K}}_k) + \tau (\mathbf{I} + \tilde{\mathbf{K}}_k^\top \tilde{\mathbf{K}}_k).$$

The follow is standard (see, e.g. Anderson & Moore (2007, Section 2.4)).

**Lemma G.1.** *There exists a unique minimizer sequence  $\mathbf{K}_{1:K}^{\pi, \star}$  such that, for all other  $\tilde{\mathbf{K}}_{1:K}$ ,  $\mathbf{P}_k^\pi [\mathbf{K}_{k:K}^{\pi, \star}] \preceq \mathbf{P}_k^\pi [\tilde{\mathbf{K}}_{k:K}]$ . We denote this minimize P-matrix  $\mathbf{P}_{\text{opt},k}^\pi := \mathbf{P}_k^\pi [\mathbf{K}_{k:K}^{\pi, \star}]$ .*

**Proposition G.2.** Recall the definition of  $\tau_{\text{ric}}$  from [Definition A.2](#),

$$\tau_{\text{ric}} := \frac{1}{4\mu_{\text{ric}}^2 \left( 3M_f \kappa_f \mu_{\text{ric}} L_f + 13L_f^2 (1 + L_f \mu_{\text{ric}})^2 \right)} = \frac{1}{\mathcal{O}_*(1)}$$

Then, as long as  $\tau \leq \min\{\tau_{\text{ric},\pi}, 1/4L_f\}$ , it holds that for any feasible policy  $\pi$ ,  $\max_{k \in [K+1]} \|\mathbf{P}_{\text{opt},k}^\pi\| \leq 2\mu_{\text{ric}}$ .

The following lemma bounding the constant  $K_\pi$  for the initial policy  $\pi$  can be established along the same lines of [Proposition G.2](#). Its proof is given in [Appendix G.2](#).

**Lemma G.3.** Suppose that  $\tau \leq \tau_{\text{ric}}$ . For  $\pi = \pi^{(1)}$ ,  $\mu_{\pi,*} \leq 2\mu_{\text{ric}}$  and  $L_\pi = 1$ .

**Proposition G.4** (Certainty Equivalence Bound). Let  $\hat{\mathbf{A}}_k^\pi$  and  $\hat{\mathbf{B}}_k^\pi$  be estimates of  $\mathbf{A}_{\text{ol},k}^\pi$  and  $\mathbf{B}_{\text{ol},k}^\pi$ , and let  $\hat{\mathbf{K}}_k$  denote the corresponding certainty equivalence controller synthesized by solving the following recursion given by  $\hat{\mathbf{P}}_{K+1} = \mathbf{I}$ , and for  $k \in [k_0 : K]$ , setting

$$\begin{aligned} \hat{\mathbf{P}}_k &= (\hat{\mathbf{A}}_k^\pi)^\top \hat{\mathbf{P}}_{k+1} \hat{\mathbf{A}}_k^\pi - \left( \hat{\mathbf{B}}_k^\pi \hat{\mathbf{P}}_{k+1} \hat{\mathbf{A}}_k^\pi \right)^\top \left( \tau^{-1} \mathbf{I} + (\hat{\mathbf{B}}_k^\pi)^\top \hat{\mathbf{P}}_{k+1} \hat{\mathbf{B}}_k^\pi \right)^{-1} \left( \hat{\mathbf{B}}_k^\pi \hat{\mathbf{P}}_{k+1} \hat{\mathbf{A}}_k^\pi \right) + \tau \mathbf{I} \\ \hat{\mathbf{K}}_k &= -\left( \tau^{-1} \mathbf{I} + (\hat{\mathbf{B}}_k^\pi)^\top \hat{\mathbf{P}}_{k+1} \hat{\mathbf{B}}_k^\pi \right)^{-1} (\hat{\mathbf{B}}_k^\pi)^\top \hat{\mathbf{P}}_{k+1} \hat{\mathbf{A}}_k^\pi, \end{aligned}$$

Then, as long as  $\max_{k \in [k_0:K]} \|\hat{\mathbf{A}}_k^\pi - \mathbf{A}_{\text{ol},k}^\pi\|_{\text{op}} \vee \|\hat{\mathbf{B}}_k^\pi - \mathbf{B}_{\text{ol},k}^\pi\|_{\text{op}} \leq (2^{17} \tau \mu_{\text{ric}}^4 \max\{1, L_f^3\})^{-1}$ , and  $\tau \leq \min\{\tau_{\text{ric}}, 1/4L_f L_\pi\}$ , we have

$$\max_{k \geq k_0} \|\mathbf{P}_k^\pi [\hat{\mathbf{K}}_{k:K}]\| \leq 4\mu_{\text{ric}}, \quad \text{and} \quad \max_{k \geq k_0} \|\hat{\mathbf{K}}_k\| \leq 6 \max\{1, L_f\} \mu_{\text{ric}}.$$

We can now prove [Proposition A.14](#).

*Proof.* Let  $\hat{\mathbf{K}}_{1:K}$  be the gains synthesized according to [Algorithm 1](#) (Line 7-10), and  $\pi' = (\mathbf{u}_{1:K}^\pi, \hat{\mathbf{K}}_{1:K})$  be the policy with the same inputs as  $\pi$  but with these new gains. Defining the shorthand  $\tilde{\mathbf{P}}_k := \|\mathbf{P}_k^\pi [\hat{\mathbf{K}}_{k:K}]\|$ , [Proposition G.4](#) then implies that

$$\max_{k \geq k_0} \|\tilde{\mathbf{P}}_k\| \leq 4\mu_{\text{ric}}, \quad \text{and} \quad \max_{k \geq k_0} \|\hat{\mathbf{K}}_k\| \leq 6 \max\{1, L_f\} \mu_{\text{ric}}.$$

Since  $\hat{\mathbf{K}}_k = 0$  for  $k < k_0$ , we conclude  $\max_{k \in [K+1]} \|\hat{\mathbf{K}}_k\| \leq 6 \max\{1, L_f\} \mu_{\text{ric}}$ , which we note is  $\geq 1$  as  $\mu_{\text{ric}} \geq 1$ . Thus, we can take  $L_{\pi'} = 6 \max\{1, L_f\} \mu_{\text{ric}}$ . Moreover, for this policy  $\pi$ , we have  $\mathbf{A}_{\text{cl},k}^{\pi'} = (\mathbf{A}_{\text{ol},k}^\pi + \mathbf{B}_{\text{ol},k}^\pi \hat{\mathbf{K}}_k)$ , so that the matrices  $\tilde{\mathbf{P}}_k$  are given by the recursion

$$\tilde{\mathbf{P}}_{K+1} = \mathbf{I}, \quad \tilde{\mathbf{P}}_k = (\mathbf{A}_{\text{cl},k}^{\pi'})^\top \tilde{\mathbf{P}}_{k+1} (\mathbf{A}_{\text{cl},k}^{\pi'}) + \tau (\mathbf{I} + \hat{\mathbf{K}}_k^\top \hat{\mathbf{K}}_k).$$

Hence,  $\tilde{\mathbf{P}}_k \succeq \Lambda_k^{\pi'}$ , where we recall that  $\Lambda_k^{\pi'}$  satisfy the recursion

$$\Lambda_{K+1}^{\pi'} = \mathbf{I}, \quad \Lambda_k^{\pi'} = (\mathbf{A}_{\text{cl},k}^{\pi'})^\top \Lambda_{k+1}^{\pi'} (\mathbf{A}_{\text{cl},k}^{\pi'}) + \tau \mathbf{I}.$$

Thus,  $\mu_{\pi',*} := \max_{k \in [k_0:K]} \|\Lambda_k^{\pi'}\| \leq \max_{k \in [k_0:K]} \|\tilde{\mathbf{P}}_k\| \leq 4\mu_{\text{ric}}$ .  $\square$

### G.1. Proof of [Proposition G.4](#)

Essentially, we instantiate [Theorem 4](#) with appropriate bounds on parameters, and using the last part of the recursion for  $k \geq k_0$ . Fix an index  $k_0 \in [K]$ , let  $K_0 = K - k_0 - 1$ , and recall  $[k_0 : j] := \{k_0, \dots, j\}$ . Throughout, we suppose

$$\tau \leq 1/4L_f \max\{1, L_\pi, \mu_{\text{ric}}\} \tag{G.1}$$

Suppose we have given estimates  $\hat{\mathbf{A}}_k^\pi$  and  $\hat{\mathbf{B}}_k^\pi$  satisfying

$$\max_{k \in [k_0:K]} \tau^{-1} \|\hat{\mathbf{A}}_k^\pi - \mathbf{A}_{\text{ol},k}^\pi\|_{\text{op}} \leq \epsilon_A, \quad \max_{k \in [k_0:K]} \tau^{-1} \|\hat{\mathbf{B}}_k^\pi - \mathbf{B}_{\text{ol},k}^\pi\|_{\text{op}} \leq \epsilon_B$$

We apply [Theorem 4](#) with the substitutions

$$K \leftarrow K_0, \quad \hat{\mathbf{B}}_k \leftarrow \tau \hat{\mathbf{B}}_{k+k_0-1}^\pi, \quad \mathbf{B}_k \leftarrow \tau \mathbf{B}_{\text{ol},k+k_0-1}^\pi, \quad \hat{\mathbf{A}}_k \leftarrow \hat{\mathbf{A}}_{k+k_0-1}^\pi, \quad \mathbf{A}_k \leftarrow \tau \mathbf{A}_{\text{ol},k+k_0-1}^\pi,$$

So that, with  $\Theta = (\mathbf{A}_j, \mathbf{B}_j)_{j \in [K_0]}$  and  $\hat{\Theta} = (\hat{\mathbf{A}}_j, \hat{\mathbf{B}}_j)_{j \in [K_0]}$ , we have

$$\hat{K}_k = \begin{cases} 0 & k < k_0 \\ \mathbf{K}_{k-k_0}^{\text{opt}}(\hat{\Theta}) & k \geq k_0 \end{cases}$$

and thus,

$$\mathbf{P}_j^{\text{opt}}(\Theta) = \mathbf{P}_{\text{opt},k}^\pi, \quad \mathbf{P}_k^{\text{ce}}(\Theta; \hat{\Theta}) = \mathbf{P}_{j+k_0-1}^\pi[\hat{K}_{k_0:K}].$$

With the above substitutions, we can apply [Proposition G.2](#) as long as  $\tau$  satisfies the condition stipulated in that proposition, we have

$$\max_{j \in [K_0+1]} \|\mathbf{P}_j^{\text{opt}}(\Theta)\| \leq 2\mu_{\text{ric}}. \quad (\text{G.2})$$

Moreover, we have that by [Lemmas I.3, I.4](#) and [I.7](#), the following holds for  $\tau \leq 1/4L_f \max\{1, L_\pi\}$ ,

$$\begin{aligned} \max_k \tau^{-1} \|\mathbf{B}_{\text{ol},k}^\pi\| &\leq \exp(1/4)L_f \\ \max_k \|\mathbf{A}_{\text{ol},k}^\pi\| &= \max_k \|\Phi_{\text{cl},k+1,k}^\pi\| \leq \frac{5}{3} \\ \max_{k \in [K]} \|\mathbf{A}_k - \mathbf{I}\| &= \max_{k \in [K]} \|\Phi_{\text{ol}}^\pi(t_{k+1}, t_k) - \mathbf{I}\| \leq \exp(1/4)\tau L_f \end{aligned} \quad (\text{G.3})$$

Hence [Conditions F.1](#) and [F.4](#) hold for

$$\begin{aligned} K_A &= 1 \vee \max_{k \in [k_0:K]} \|\mathbf{A}_{\text{ol},k}^\pi\| \vee \|\hat{\mathbf{A}}_k^\pi\| \leq \frac{5}{3} \\ K_B &= 1 \vee \max_{k \in [k_0:K]} \tau^{-1} (\|\mathbf{B}_{\text{ol},k}^\pi\| \vee \|\hat{\mathbf{B}}_k^\pi\|) \leq \exp(1/4) \max\{1, L_f\} \\ \kappa_A &:= \tau^{-1} \max_{k \in [k_0:K]} \|\mathbf{A}_k - \mathbf{I}\| = \exp(1/4)L_f. \end{aligned} \quad (\text{G.4})$$

Moreover, [Condition F.2](#) holds with  $\Delta_A = \epsilon_A$ ,  $\Delta_B = \epsilon_B$ . We can now apply [Theorem 4](#). We take and for  $\epsilon_A \leq 1/3$  and  $\epsilon_B \leq L_f/2$ , we may take

$$\begin{aligned} \Delta_{\text{ce}} &:= 80C^4 K_A^3 K_B^3 (1 + \tau C K_B) (\Delta_A + \Delta_B) \\ C &:= \max_{j \in [K_0+1]} \|\mathbf{P}_j^{\text{opt}}(\Theta)\| \leq 2\mu_{\text{ric}} \end{aligned} \quad (\text{by Eq. (G.2)})$$

And we can bound (recalling  $\tau \leq 1/4L_f\mu_{\text{ric}}$  and  $\mu_{\text{ric}} \geq 1$ )

$$\begin{aligned} \Delta_{\text{ce}} &\leq 80\mu_{\text{ric}}^4 \cdot (16 \cdot (5/3)^3 \cdot \exp(3/4)) \max\{1, L_f^3\} (1 + 4\tau L_f \mu_{\text{ric}}) (\epsilon_A + \epsilon_B) \\ &\leq 2^{14} \mu_{\text{ric}}^4 \cdot \max\{1, L_f^3\} (\epsilon_A + \epsilon_B) (1 + 4\tau L_f \mu_{\text{ric}}) \\ &\leq 2^{15} \mu_{\text{ric}}^4 \cdot \max\{1, L_f^3\} (\epsilon_A + \epsilon_B). \end{aligned}$$

Hence, as long as

$$2^{16} \mu_{\text{ric}}^4 \max\{1, L_f^3\} (\epsilon_A + \epsilon_B) \leq 1,$$

we have

$$\Delta_{\text{ce}} \leq 1/2 \quad (\text{G.5})$$



and therefore, by [Theorem 4\(a\)](#),

$$\max_{k \in [k_0:K+1]} \|\mathbf{P}_k^\pi[\hat{\mathbf{K}}_{k_0:K}]\| = \max_{j \in [K_0+1]} \|\mathbf{P}_j^{\text{ce}}(\Theta; \hat{\Theta})\| \leq 2 \max_{j \in [K_0+1]} \|\mathbf{P}_j^{\text{opt}}(\Theta)\| \leq 4\mu_{\text{ric}}.$$

Next, [Theorem 4\(b\)](#), we can take  $L_{\pi'} = 6 \max\{1, L_f\}\mu_{\text{ric}}$ :

$$\begin{aligned} \max_{k \in [K]} \|\hat{\mathbf{K}}^{\pi'}\| &= \max_{j \in [K_0]} \|\mathbf{K}_j^{\text{opt}}(\hat{\Theta})\| \\ &\leq \frac{5}{4} K_B K_A C \leq \frac{5}{4} (5/3) \exp(1/1) \max\{1, L_f\} \cdot 2\mu_{\text{ric}} \leq L_{\pi'} := 6 \max\{1, L_f\}\mu_{\text{ric}}. \end{aligned}$$

□

## G.2. Proof of [Proposition G.2](#)

### G.2.1. PRELIMINARIES.

We recall the following, standard definition of continuous-time cost to-go matrices (see, e.g. [Anderson2007optimal](#)):

**Definition G.2** (Cost-to-Go Matrices). Given a *policy*  $\pi$ , and a sequence of controls  $\tilde{\mathbf{u}}(\cdot) \in \mathcal{U}$ , let  $\mathbf{P}^\pi(\cdot \mid \tilde{\mathbf{u}})$  as the cost-to-go matrix satisfying  $\xi^\top \mathbf{P}^\pi(t \mid \tilde{\mathbf{u}})\xi = \int_{s=t}^T (\|\tilde{\mathbf{x}}(s)\|^2 + \|\tilde{\mathbf{u}}(s)\|^2) ds + \|\tilde{\mathbf{x}}(T)\|^2$ , under the dynamics  $\frac{d}{ds} \tilde{\mathbf{x}}(s) = \mathbf{A}_{\text{ol}}^\pi(s) \tilde{\mathbf{x}}(s) + \mathbf{B}_{\text{ol}}^\pi(s) \tilde{\mathbf{u}}(s)$ ,  $\tilde{\mathbf{x}}(t) = \xi$ . We let  $\mathbf{P}_{\text{opt}}^\pi(t)$  denote the optimal cost-to-go matrix, i.e., the matrix satisfying  $\xi^\top \mathbf{P}_{\text{opt}}^\pi(t)\xi = \min_{\tilde{\mathbf{u}} \in \mathcal{U}} \xi^\top \mathbf{P}^\pi(t \mid \tilde{\mathbf{u}})\xi := V^\pi(t \mid \tilde{\mathbf{u}}, \xi)$ .

Recall that [Assumption 4.3](#) implies  $V^\pi(t \mid \tilde{\mathbf{u}}, \xi) \leq \mu_{\text{ric}} \|\xi\|^2$ , so that  $\|\mathbf{P}_{\text{opt}}^\pi(t)\| \leq \mu_{\text{ric}}$ . In what follows, we suppress superscript dependence on  $\pi$ , assume  $\pi$  is feasible, and adopt the shorthand  $\mathbf{P}(t) = \mathbf{P}_{\text{opt}}^\pi(t)$ ,  $\mathbf{A}(t) = \mathbf{A}_{\text{ol}}^\pi(t)$ ,  $\mathbf{B}(t) = \mathbf{B}_{\text{ol}}^\pi(t)$ ,  $\mathbf{x}(t) = \mathbf{x}^\pi(t)$ , and  $\mathbf{u}(t) = \mathbf{u}^\pi(t)$ . We also use the shorthand

$$L_{\text{cl}} := L_f(1 + L_f\mu_{\text{ric}}). \quad (\text{G.6})$$

The optimal input defining  $\mathbf{P}(t)$  in [Assumption 4.3](#) selects  $\tilde{\mathbf{u}}(t) = \mathbf{K}(t)\tilde{\mathbf{x}}(t)$ , where  $\mathbf{K}(t) = \mathbf{B}(t)^\top \mathbf{P}(t)$  (again, [Anderson & Moore \(2007, Section 2.3\)](#)). Introduce the evaluations of the *continuous* value function  $\mathbf{P}(t)$  and  $\mathbf{K}(t)$  at the time steps  $t_k$ :

$$\mathbf{P}_k^{\text{ct}} := \mathbf{P}(t_k), \quad \mathbf{K}_k^{\text{ct}} := \mathbf{K}(t_k) \quad (\text{G.7})$$

We also define an *suboptimal* discrete-time value function by taking  $\mathbf{P}_k^{\text{sub}} = \mathbf{P}_k^\pi[\mathbf{K}_{1:K}^{\text{ct}}]$ , defined in [Definition G.1](#), which satisfies

$$\mathbf{P}_k^{\text{sub}} \succeq \mathbf{P}_{\text{opt},k}^\pi.$$

by optimality of  $\mathbf{P}_{\text{opt},k}^\pi$ . Hence, it suffices to bound  $\mathbf{P}_k^{\text{sub}}$ . To do this, first express both  $\mathbf{P}_k^{\text{sub}}$  and  $\mathbf{P}_k^{\text{ct}}$  as discrete Lyapunov recursions. To do so, we require the relevant transition operators.

**Definition G.3** (Relevant Transitions Operators). For  $k \in [K]$  and  $s \in \mathcal{I}_k$ , let  $\Phi_1(s, t_k)$  and  $\Phi_2(s, t_k)$  denote the solution to the ODEs

$$\begin{aligned} \frac{d}{ds} \Phi_1(s, t_k) &= (\mathbf{A}(s) + \mathbf{B}(s)\mathbf{K}(s))\Phi_1(s, t) \\ \frac{d}{ds} \Phi_2(s, t_k) &= \mathbf{A}(s)\Phi_2(s, t) + \mathbf{B}(s)\mathbf{K}_k(t_k). \end{aligned} \quad (\text{G.8})$$

with initial conditions  $\Phi_1(t_k, t_k) = \Phi_2(t_k, t_k) = \mathbf{I}$ . We define

$$\mathbf{X}_k^{\text{ct}} := \Phi_1(t_{k+1}, t_k), \quad \mathbf{X}_k^{\text{sub}} := \Phi_2(t_{k+1}, t_k)$$

**Definition G.4** (Relevant Cost Matrices). For  $k \in [K]$ , define

$$\begin{aligned} \mathbf{Y}_k^{\text{ct}} &:= \int_{s=t_k}^{t_k} \Phi_1(s, t_k)^\top (\mathbf{I} + \mathbf{K}(s)^\top \mathbf{K}(s)) \Phi_1(s, t_k) ds \\ \mathbf{Y}_k^{\text{sub}} &:= \tau(\mathbf{I} + \mathbf{K}(t_k)^\top \mathbf{K}(t_k)) \end{aligned}$$

**Lemma G.5.** *The cost-to-go matrices  $\mathbf{P}_k^{\text{ct}}$  and  $\mathbf{P}_k^{\text{sub}}$  are given by the following Lyapunov recursions, with initial conditions  $\mathbf{P}_{K+1}^{\text{ct}} = \mathbf{P}_{K+1}^{\text{sub}} = \mathbf{I}$ :*

$$\begin{aligned}\mathbf{P}_k^{\text{ct}} &= (\mathbf{X}_k^{\text{ct}})^\top \mathbf{P}_{k+1}^{\text{ct}} \mathbf{X}_k^{\text{ct}} + \mathbf{Y}_k^{\text{ct}} \\ \mathbf{P}_k^{\text{sub}} &= (\mathbf{X}_k^{\text{sub}})^\top \mathbf{P}_{k+1}^{\text{sub}} \mathbf{X}_k^{\text{sub}} + \mathbf{Y}_k^{\text{sub}}\end{aligned}$$

*Proof of Lemma G.5.* The recursion for  $\mathbf{P}_k^{\text{sub}}$  is directly from [Definition G.1](#), and the fact that  $\mathbf{X}_k^{\text{sub}} = \mathbf{A}_{\text{cl},k}^\pi$  due to [Lemma C.10](#). To verify the recursion for  $\mathbf{P}_k^{\text{ct}}$ , we note that we can express  $\mathbf{P}(t) = \mathbf{P}_{\text{opt}}^\pi(t)$  in [Definition G.2](#) as satisfying the following ODE (see [Anderson & Moore \(2007, Section 2.3\)](#)):

$$\mathbf{P}(T) = \mathbf{I}, \quad -\frac{d}{dt}\mathbf{P}(t) = (\mathbf{A}(t) + \mathbf{B}(t)\mathbf{K}(t))^\top \mathbf{P}(t) (\mathbf{A}(t) + \mathbf{B}(t)\mathbf{K}(t)) + \mathbf{I} + \mathbf{K}(t)^\top \mathbf{K}(t)$$

It can be checked then by computing derivatives and using existence and uniqueness of ODEs that

$$\mathbf{P}(s, t) = \Phi_1(s, t)^\top \mathbf{P}(s) \Phi_1(s, t) + \int_{s'=t}^s \Phi_1(s', t)^\top (\mathbf{I} + \mathbf{K}(s')^\top \mathbf{K}(s')) \Phi_1(s', t) ds'$$

Specializing to  $s = t_{k+1}$  and  $t = t_k$  verifies the desired recursion.  $\square$

As  $\mathbf{P}_k^{\text{ct}} = \mathbf{P}(t_k)$ , the terms  $\mathbf{P}_k^{\text{ct}}$  are bounded whenever  $\mathbf{P}(\cdot)$  is. Therefore, we use a Lyapunov perturbation bound to bound  $\mathbf{P}_k^{\text{sub}}$  in terms of  $\mathbf{P}_k^{\text{ct}}$ . This requires reasoning about the differences  $\mathbf{X}_k^{\text{ct}} - \mathbf{X}_k^{\text{sub}}$  and  $\mathbf{Y}_k^{\text{ct}} - \mathbf{Y}_k^{\text{sub}}$ , which we do in just below.

### G.2.2. CONTROLLING THE RATE OF CHANGE OF $\mathbf{K}(t)$ .

Our first step in controlling the perturbation term is to argue that the optimal controller  $\mathbf{K}(t)$  does not change too rapidly. As  $\mathbf{K}(t) = \mathbf{B}(t)^\top \mathbf{P}(t)$ , we begin by bounding the change in  $\mathbf{B}(t)$ .

**Claim G.1** (Change in  $\mathbf{B}(t)$ ).  *$\mathbf{B}(t)$  is differentiable in  $t$  on for  $t \in \text{int}(\mathcal{I}_k)$ , and satisfies  $\|\frac{d}{dt}\mathbf{B}(t)\| \leq M_f \kappa_f$*

*Proof.* Recall that  $\mathbf{B}(t) = \partial_u f(\mathbf{x}(t), \mathbf{u}(t))$ . For  $t \in \text{int}(\mathcal{I}_k)$ ,  $\mathbf{u}(t)$  is constant, and  $\mathbf{x}(t)$ , being the solution to an ODE, is also  $t$ -differentiable. We now bound  $\|\frac{d}{dt}\mathbf{B}(t)\|$ . We have

$$\begin{aligned}\|\frac{d}{dt}\mathbf{B}(t)\| &= \|\frac{d}{dt}\partial_u f(\mathbf{x}(t), \mathbf{u}(t))\| \\ &= \|\partial_{uu}f(\mathbf{x}(t), \mathbf{u}(t))\frac{d}{dt}\mathbf{u}(t) + \partial_{xu}f(\mathbf{x}(t), \mathbf{u}(t))\frac{d}{dt}\mathbf{x}(t)\| \\ &\leq M_f \|\frac{d}{dt}\mathbf{x}(t)\| \leq M_f \kappa_f\end{aligned}$$

where the second-to-last inequality is the limiting consequence holds from [Assumption 4.1](#), and where the term  $\partial_{uu}f(\mathbf{x}(t), \mathbf{u}(t))\frac{d}{dt}\mathbf{u}(t)$  vanishes because  $\mathbf{u}(t) = \mathbf{u}^\pi(t)$  is constant on  $\mathcal{I}_k$ .  $\square$

Next, we bound the change in  $\mathbf{P}(t)$ :

**Claim G.2** (Change in  $\mathbf{P}(t)$ ).  *$\mathbf{P}(t)$  is differentiable in  $t$ , and  $\|\frac{d}{dt}\mathbf{P}(t)\| \leq (L_{\text{cl}}/L_f)^2$ .*

*Proof.* Note that  $\mathbf{P}(t)$  is given by the ODE

$$\mathbf{P}(T) = \mathbf{I}, \quad -\frac{d}{dt}\mathbf{P}(t) = \mathbf{A}(t)^\top \mathbf{P}(t) + \mathbf{P}(t)\mathbf{A}(t) - \mathbf{P}(t)\mathbf{B}(t)\mathbf{B}(t)^\top \mathbf{P}(t) + \mathbf{I},$$

which ensures differentiability. Thus, as  $\|\mathbf{A}(t)\| \vee \|\mathbf{A}(t)\| \vee 1 \leq L_f$  by [Assumption 4.1](#) and [Assumption 4.3](#),

$$\|\frac{d}{dt}\mathbf{P}(t)\| \leq 1 + 2L_f\|\mathbf{P}(t)\| + L_f^2\|\mathbf{P}(t)\|^2 \leq (1 + 2L_f\mu_{\text{ric}} + L_f^2\mu_{\text{ric}}^2) \leq (1 + L_f\mu_{\text{ric}})^2,$$

which is precisely  $(L_{\text{cl}}/L_f)^2$ .  $\square$

We now establish a bound on the change in  $\mathbf{K}(t)$ .

**Claim G.3** (Continuity of Optimal Controller). *For all  $t \in \mathcal{I}_k$ ,*

$$\left\| \frac{d}{dt} \mathbf{K}(t) \right\| \leq M_f \kappa_f \mu_{\text{ric}} + L_f^{-1} L_{\text{cl}}^2.$$

*Proof of Claim G.3.* By Claims G.1 and G.2, we have

$$\begin{aligned} \left\| \frac{d}{dt} \mathbf{K}(t) \right\| &\leq \left\| \frac{d}{dt} \mathbf{B}(t) \right\| \|\mathbf{P}(t)\| + \|\mathbf{B}(t)\| \left\| \frac{d}{dt} \mathbf{P}(t) \right\| \\ &\leq \left\| \frac{d}{dt} \mathbf{B}(t) \right\| \mu_{\text{ric}} + L_f (L_{\text{cl}}/L_f)^2 \\ &\leq M_f \kappa_f \mu_{\text{ric}} + L_f^{-1} L_{\text{cl}}^2. \end{aligned}$$

□

By integrating, we arrive at the next claim.

**Claim G.4.** *The following bound holds*

$$\sup_{s \in \mathcal{I}_k} \|\mathbf{K}(s) - \mathbf{K}(t_k)\| \leq \tau \left( M_f \kappa_f \mu_{\text{ric}} + L_f^{-1} L_{\text{cl}}^2 \right).$$

*Proof of Claim G.4.* Directly from Claim G.3. □

G.2.3. CONTROLLING DIFFERENCES IN  $\|\mathbf{x}_k^{\text{ct}} - \mathbf{x}_k^{\text{sub}}\|$  AND  $\|\mathbf{y}_k^{\text{ct}} - \mathbf{y}_k^{\text{sub}}\|$

We first state a bound on the magnitudes of various quantities of interest.

**Claim G.5.**  $\|\mathbf{K}(t)\| \leq \mu_{\text{ric}} L_f$  and  $\|\mathbf{A}(t) + \mathbf{B}(t)\mathbf{K}(t)\| \leq L_{\text{cl}}$ , where we recall  $L_{\text{cl}} := L_f(1 + L_f \mu_{\text{ric}})$ .

*Proof.* Recall that  $\mathbf{K}(t) = \mathbf{B}(t)^\top \mathbf{P}(t)$ . From Assumption 4.3,  $\|\mathbf{P}(t)\| \leq \mu_{\text{ric}}$ , and  $\|\mathbf{B}(t)\| \leq L_f$  by Assumption 4.1, which gives  $\|\mathbf{K}(t)\| \leq \mu_{\text{ric}} L_f$ . Bounding  $\|\mathbf{A}(t)\| \vee \|\mathbf{B}(t)\|$  by  $L_f$  (again, invoking Assumption 4.1), concludes the demonstration. □

Next, we show that  $\Phi_1(s, t_k)$  is close to the identity for sufficiently small  $\tau$ .

**Claim G.6.** *Suppose that  $\tau L_{\text{cl}} \leq 1/2$ . Then,*

$$\|\mathbf{I} - \Phi_1(s, t_k)\| \leq \tau L_{\text{cl}} \exp(1/2) \leq \min\{1, 2\tau L_{\text{cl}}\}$$

*Proof of Claim G.6.* It suffices to bound, for all  $\xi \in \mathbb{R}^{d_x} : \|\xi\| = 1$  the differences  $\|\mathbf{y}_1(s) - \xi\|$  where  $\mathbf{y}_1 = \Phi_1(s, t_k)\xi$ . We do this via Picard's lemma.

Specifically, write  $\frac{d}{ds} \mathbf{y}_1(s) = \tilde{f}(\mathbf{y}_1(s), s)$ , where  $\tilde{f}(y, s) = (\mathbf{A}(s) + \mathbf{B}(s)\mathbf{K}(s))y$ , and  $\mathbf{z}(s) = \xi$ . As  $\tilde{f}(y, s)$  is  $\sup_{s \in \mathcal{I}_k} \|\mathbf{A}(s) + \mathbf{B}(s)\mathbf{K}(s)\| \leq L_{\text{cl}}$  Lipschitz in  $y$  (here, we use Claim G.5) and as  $\frac{d}{ds} \xi = 0$ , and the Picard Lemma (Lemma C.9) gives

$$\begin{aligned} \|\xi - \mathbf{y}_1(s)\| &\leq \exp((s - t_k)(2L_f^2 \mu_{\text{ric}})) \int_{s'=t_k}^s \|(\mathbf{A}(s') + \mathbf{B}(s')\mathbf{K}(s'))\xi\| ds' \\ &\leq \exp((s - t_k)L_{\text{cl}}) \int_{s'=t_k}^s \|(\mathbf{A}(s') + \mathbf{B}(s')\mathbf{K}(s'))\| ds' \quad (\|\xi\| \leq 1) \\ &\leq \exp((s - t_k)L_{\text{cl}}) \cdot (s - t_k)L_{\text{cl}}, \\ &\leq \exp(\tau L_{\text{cl}}) \cdot \tau L_{\text{cl}}, \\ &\leq \exp(1/2) \tau L_{\text{cl}} \end{aligned}$$

where we assume  $\tau L_{\text{cl}} \leq 1/2$ . □

We can now bound the differences between  $\|\mathbf{x}_k^{\text{ct}} - \mathbf{x}_k^{\text{sub}}\| = \|\Phi_2(t_{k+1}, t_k) - \Phi_1(t_{k+1}, t_k)\|$ .

**Lemma G.6.** For  $k \in [K]$  and  $s \in \mathcal{I}_k$ , let  $\Phi_1(s, t_k)$  and  $\Phi_2(s, t_k)$  denote the solution to the ODEs

$$\frac{d}{ds}\Phi_1(s, t_k) = (\mathbf{A}(s) + \mathbf{B}(s)\mathbf{K}(s))\Phi_1(s, t_k), \quad \frac{d}{ds}\Phi_2(s, t_k) = \mathbf{A}(s)\Phi_2(s, t_k) + \mathbf{B}(s)\mathbf{K}(t_k).$$

with initial conditions  $\Phi_1(t_k, t_k) = \Phi_2(t_k, t_k) = \mathbf{I}$ . Then, if  $\tau L_{\text{cl}} \leq 1/2$ ,

$$\|\mathbf{x}_k^{\text{ct}} - \mathbf{x}_k^{\text{sub}}\| = \|\Phi_2(t_{k+1}, t_k) - \Phi_1(t_{k+1}, t_k)\| \leq 2\tau^2 (L_f \mu_{\text{ric}} M_f \kappa_f + 3L_{\text{cl}}^2).$$

*Proof.* It suffices to bound, for all initial conditions,  $\xi \in \mathbb{R}^{d_x}$  with  $\|\xi\| = 1$ , the solutions  $\mathbf{y}_i(s) = \Phi_i(s)\xi$ . We apply the Picard Lemma, with  $\mathbf{z}(s) \leftarrow \mathbf{y}_1(s)$ , and express  $\mathbf{y}_2(s) = \tilde{f}(\mathbf{y}_2(s), s)$ , where  $\tilde{f}(y, s) = \mathbf{A}(s)y + \mathbf{B}(s)\mathbf{K}(t_k)$ . As  $\|\mathbf{A}(s)\| \leq L_f$ , the Picard Lemma (Lemma C.9) yields

$$\begin{aligned} \|\mathbf{y}_1(t_{k+1}) - \mathbf{y}_2(t_{k+1})\| &\leq \exp(L_f(t - s)) \int_{s=t}^{t_{k+1}} \|\mathbf{A}(s)\mathbf{y}_1(s) + \mathbf{B}(s)\mathbf{K}(t_k)\xi - \frac{d}{ds}\mathbf{y}_1(s)\| ds \\ &\leq \exp(L_f\tau) \int_{s=t}^{t_{k+1}} \|\mathbf{A}(s)\mathbf{y}_1(s) + \mathbf{B}(s)\mathbf{K}(t_k)\xi - (\mathbf{A}(s) + \mathbf{B}(s)\mathbf{K}(s))\mathbf{y}_1(s)\| ds \\ &\leq \exp(L_f\tau) \int_{s=t}^{t_{k+1}} \|\mathbf{B}(s)(\mathbf{K}(s)\mathbf{y}_1(s) - \mathbf{K}(t_k)\xi)\| ds \\ &\leq L_f \exp(L_f\tau) \int_{s=t}^{t_{k+1}} \|\mathbf{K}(s)\mathbf{y}_1(s) - \mathbf{K}(t_k)\xi\| ds \\ &\leq L_f \exp(L_f\tau) \int_{s=t}^{t_{k+1}} (\|\mathbf{K}(s) - \mathbf{K}(t_k)\|\|\xi\| + \|\mathbf{K}(s)(\xi - \mathbf{y}_1(s))\|) ds \\ &\leq L_f \exp(L_f\tau) \int_{s=t}^{t_{k+1}} (\|\mathbf{K}(s) - \mathbf{K}(t_k)\| + L_f \mu_{\text{ric}} \|\xi - \mathbf{y}_1(s)\|) ds \\ &\leq \exp(1/2) L_f \tau \max_{s \in \mathcal{I}_k} (\|\mathbf{K}(s) - \mathbf{K}(t_k)\| + L_f \mu_{\text{ric}} \|\xi - \mathbf{y}_1(s)\|) \end{aligned}$$

where the second-to-last line uses  $\|\xi\| = 1$  and  $\|\mathbf{K}(s)\| \leq L_u \mu_{\text{ric}}$ , and the last uses  $\tau \leq 1/2L_x$  and well as a bound of an integral by a maximum. By claims [Claims G.4](#) and [G.6](#),

$$\begin{aligned} &\max_{s \in \mathcal{I}_k} (\|\mathbf{K}(s) - \mathbf{K}(t_k)\| + L_f \mu_{\text{ric}} \|\xi - \mathbf{y}_1(s)\|) \\ &\leq \tau (M_f \kappa_f \mu_{\text{ric}} + L_f^{-1} L_{\text{cl}}^2) + L_f \mu_{\text{ric}} \exp(1/2) \tau L_{\text{cl}} \\ &= \tau (\mu_{\text{ric}} M_f \kappa_f + L_{\text{cl}} (L_f \mu_{\text{ric}} \exp(1/2) + L_f^{-1} L_{\text{cl}})) \\ &\leq \tau (\mu_{\text{ric}} M_f \kappa_f + 3L_f^{-1} L_{\text{cl}}^2), \end{aligned}$$

where in the last inequality we use  $\exp(1/2) \leq 2$  and  $L_f \mu_{\text{ric}} \leq (1 + L_f \mu_{\text{ric}}) = L_{\text{cl}} L_f^{-1}$ . Therefore, again using  $\exp(1/2) \leq 2$ ,

$$\|\mathbf{y}_1(t_{k+1}) - \mathbf{y}_2(t_{k+1})\| \leq 2\tau^2 (L_f \mu_{\text{ric}} M_f \kappa_f + 3L_{\text{cl}}^2).$$

Quantifying over all unit-norm initial conditions  $\xi$  concludes the proof.  $\square$

We now establish a qualitatively similar bound on  $\tau \|\mathbf{y}_k^{\text{ct}} - \mathbf{y}_k^{\text{sub}}\|$ .

**Lemma G.7.**  $\|\mathbf{y}_k^{\text{ct}} - \mathbf{y}_k^{\text{sub}}\| \leq 2\tau^2 (M_f \kappa_f \mu_{\text{ric}}^2 L_f + 7\mu_{\text{ric}} L_{\text{cl}}^2)$ .

*Proof.* Recall the definitions

$$\begin{aligned} \mathbf{y}_k^{\text{ct}} &:= \int_{s=t_k}^{t_{k+1}} \Phi_1(s, t_k)^\top (\mathbf{I} + \mathbf{K}(s)^\top \mathbf{K}(s)) \Phi_1(s, t_k) ds \\ \mathbf{y}_k^{\text{sub}} &:= \tau (\mathbf{I} + \mathbf{K}(t_k)^\top \mathbf{K}(t_k)) \end{aligned}$$

We can then express

$$\begin{aligned} \mathbf{Y}_k^{\text{ct}} - \mathbf{Y}_k^{\text{sub}} &= \mathbf{Y}_k^{\text{ct}} - \tau(\mathbf{I} + \mathbf{K}(t_k)\mathbf{K}(t_k)^\top) = \int_{s=t_k}^{t_k} \mathbf{Z}_k(s), \\ \mathbf{Z}_k(s) &:= \{\Phi_1(s, t_k)^\top (\mathbf{I} + \mathbf{K}(s)^\top \mathbf{K}(s)) \Phi_1(s, t_k) - (\mathbf{I} + \mathbf{K}(t_k)^\top \mathbf{K}(t_k))\} ds \end{aligned}$$

Thus,

$$\|\mathbf{Y}_k^{\text{ct}} - \mathbf{Y}_k^{\text{sub}}\| \leq \tau \max_{s \in \mathcal{I}_k} \|\mathbf{Z}_k(s)\|. \quad (\text{G.9})$$

With numerous applications of the triangle inequality,

$$\begin{aligned} \|\mathbf{Z}_k(s)\| &\leq \|\mathbf{I} - \Phi_1(s, t_k)\| \|\mathbf{I} + \mathbf{K}(s)^\top \mathbf{K}(s)\| \|\Phi_1(s, t_k)\| \\ &\quad + \|\mathbf{I} + \mathbf{K}(s)^\top \mathbf{K}(s)\| \|\mathbf{I} - \Phi_1(s, t_k)\| + \|\mathbf{K}(s) - \mathbf{K}(t_k)\| (\|\mathbf{K}(s)\| + \|\mathbf{K}(t_k)\|). \end{aligned}$$

Using  $\|\mathbf{K}(s)\| \vee \|\mathbf{K}(t_k)\| \leq L_f \mu_{\text{ric}}$  due to [Claim G.3](#), we have

$$\begin{aligned} \|\mathbf{Z}_k(s)\| &\leq (1 + L_f^2 \mu_{\text{ric}}^2)(1 + \|\Phi_1(s, t_k)\|) \|\mathbf{I} - \Phi_1(s, t_k)\| + 2\mu_{\text{ric}} L_f \|\mathbf{K}(s) - \mathbf{K}(t_k)\| \\ &\leq 3(1 + L_f^2 \mu_{\text{ric}}^2) \|\mathbf{I} - \Phi_1(s, t_k)\| + 2\mu_{\text{ric}} L_f \|\mathbf{K}(s) - \mathbf{K}(t_k)\| && (\text{Claim G.6}) \\ &\leq 6\tau L_{\text{cl}}(1 + L_f^2 \mu_{\text{ric}}^2) + 2\mu_{\text{ric}} L_f \|\mathbf{K}(s) - \mathbf{K}(t_k)\| && (\text{Claim G.6}) \\ &\leq 12\tau L_f^2 \mu_{\text{ric}}(1 + L_f^2 \mu_{\text{ric}}^2) + 2\tau \mu_{\text{ric}} L_f (M_f \kappa_f \mu_{\text{ric}} + L_f^{-1} L_{\text{cl}}^2). && (\text{Claim G.4}) \end{aligned}$$

We can upper bound  $L_f^2(1 + L_f^2 \mu_{\text{ric}}^2) \leq L_f^2(1 + L_f \mu_{\text{ric}})^2 = L_{\text{cl}}^2$ , and simplify  $2\tau \mu_{\text{ric}} L_f (M_f \kappa_f \mu_{\text{ric}} + L_f^{-1} L_{\text{cl}}^2) = 2\tau (M_f \kappa_f \mu_{\text{ric}}^2 L_f + \mu_{\text{ric}} L_{\text{cl}}^2)$ . This gives

$$\|\mathbf{Z}_k(s)\| \leq 12\tau \mu_{\text{ric}} L_{\text{cl}}^2 + 2\tau (M_f \kappa_f \mu_{\text{ric}}^2 L_f + \mu_{\text{ric}} L_{\text{cl}}^2) = 2\tau (M_f \kappa_f \mu_{\text{ric}}^2 L_f + 7\mu_{\text{ric}} L_{\text{cl}}^2).$$

Plugging the above bound into [Eq. \(G.9\)](#) concludes.  $\square$

#### G.2.4. CONCLUDING THE PROOF OF [PROPOSITION G.2](#)

From [Lemmas G.6](#) and [G.7](#), we have

$$\|\mathbf{X}_k^{\text{ct}} - \mathbf{X}_k^{\text{sub}}\| \leq 2\tau^2 (L_f \mu_{\text{ric}} M_f \kappa_f + 3L_{\text{cl}}^2), \quad \|\mathbf{Y}_k^{\text{ct}} - \mathbf{Y}_k^{\text{sub}}\| \leq 2\tau^2 \mu_{\text{ric}} (M_f \kappa_f \mu_{\text{ric}} L_f + 7L_{\text{cl}}^2)$$

Therefore, using  $\mu_{\text{ric}} \geq 1$

$$2\|\mathbf{X}_k^{\text{ct}} - \mathbf{X}_k^{\text{sub}}\| + \|\mathbf{Y}_k^{\text{ct}} - \mathbf{Y}_k^{\text{sub}}\| \leq 2\tau^2 \mu_{\text{ric}} (3M_f \kappa_f \mu_{\text{ric}} L_f + 13L_{\text{cl}}^2). \quad (\text{G.10})$$

Now, we invoke [Proposition F.12](#). We construct linear interpolation (here,  $s \in [0, 1]$  parametrizes the interpolation and not time)

$$\mathbf{X}_k(s) = (1-s)\mathbf{X}_k^{\text{ct}} + s\mathbf{X}_k^{\text{sub}}, \quad \mathbf{Y}_k(s) = (1-s)\mathbf{Y}_k^{\text{ct}} + s\mathbf{Y}_k^{\text{sub}}.$$

Then, by [Lemma G.5](#), the interpolator  $\Lambda_k(s)$  defined in [Proposition F.12](#) satisfies  $\Lambda_k(0) = \mathbf{P}_k^{\text{ct}}$  and  $\Lambda_k(1) = \mathbf{P}_k^{\text{sub}}$ . In particular,

$$\begin{aligned} \|\Lambda_{1:K+1}(0)\|_{\max, \text{op}} &= \max_{k \in [K+1]} \|\mathbf{P}_k^{\text{ct}}\| && (\text{since } \Lambda_k = \mathbf{P}_k^{\text{ct}} \text{ and definition of } \|\cdot\|_{\max, \text{op}}) \\ &= \max_{k \in [K+1]} \|\mathbf{P}(t_k)\| && (\text{by Eq. (G.7)}) \\ &\leq \sup_{t \in [T]} \|\mathbf{P}(t)\| \\ &\leq \mu_{\text{ric}}, && (\text{G.11}) \end{aligned}$$

where the last inequality is by [Assumption 4.3](#). Moreover, the term  $\tilde{\Delta}(s)$  defined in [Eq. \(F.14\)](#) satisfies

$$\begin{aligned} \forall s \in [0, 1], \tilde{\Delta}(s) &= \tau^{-1}(2\|\mathbf{x}_k^{\text{ct}} - \mathbf{x}_k^{\text{sub}}\| + \tau\|\mathbf{y}_k^{\text{ct}} - \mathbf{y}_k^{\text{sub}}\|) \\ &\leq \tau \cdot 2\mu_{\text{ric}} (3M_f\kappa_f\mu_{\text{ric}}L_f + 13L_{\text{cl}}^2). \end{aligned} \quad (\text{by Eq. (G.10)})$$

Hence, recalling  $\|\Lambda_{1:K+1}(0)\|_{\max, \text{op}} \leq \mu_{\text{ric}}$  due to [Eq. \(G.11\)](#), it holds that long as  $2\tau\mu_{\text{ric}}^2 (3M_f\kappa_f\mu_{\text{ric}}L_f + 13L_{\text{cl}}^2) \leq \frac{1}{2}$ , it holds that

$$\max_{k \in [K+1]} \|\mathbf{p}_k^{\text{sub}}\| = \|\Lambda_{1:K+1}(1)\|_{\max, \text{op}} \leq 2\|\Lambda_{1:K+1}(0)\|_{\max, \text{op}} \leq 2\mu_{\text{ric}}.$$

Lastly, we note the condition  $2\tau\mu_{\text{ric}}^2 (3M_f\kappa_f\mu_{\text{ric}}L_f + 13L_{\text{cl}}^2) \leq \frac{1}{2}$  is equivalent to

$$\begin{aligned} \tau &\leq \frac{1}{4\mu_{\text{ric}}^2 (3M_f\kappa_f\mu_{\text{ric}}L_f + 13L_{\text{cl}}^2)} \\ &\leq \frac{1}{4\mu_{\text{ric}}^2 (3M_f\kappa_f\mu_{\text{ric}}L_f + 13L_{\text{cl}}^2(1 + L_f\mu_{\text{ric}})^2)} := \tau_{\text{ric}} \end{aligned} \quad (\text{Definition of } L_{\text{cl}} \text{ in Eq. (G.6)})$$

This concludes the proof of [Proposition G.2](#).  $\square$

### G.3. Proof of [Lemma G.3](#)

The proof is similar to [Proposition G.2](#). Let  $\pi = \pi^{(1)}$ . As  $\mathbf{K}_k^\pi = 0$  for all  $k$ , we that  $L_\pi = 1$ , and that  $\Lambda_{K+1}^\pi = \mathbf{I}$ , and

$$\Lambda_k^\pi = (\mathbf{A}_{\text{ol},k}^\pi)^\top \Lambda_{k+1}^\pi \mathbf{A}_{\text{ol},k}^\pi + \tau \mathbf{I} \quad (\text{G.12})$$

On the other hand, following the arguments of [Proposition G.2](#), we it can be shown that  $V^\pi(t_k; \mathbf{u} = 0, \xi) = \xi^\top \mathbf{p}_k^{\text{ct}} \xi$ , where  $\mathbf{p}_k^{\text{ct}}$  satisfies the recursion  $\mathbf{p}_{K+1}^{\text{ct}} = \mathbf{I}$  and

$$\mathbf{p}_k^{\text{ct}} = \Phi_1(t_{k+1}, t_k)^\top \mathbf{p}_{k+1}^{\text{ct}} \Phi_1(t_{k+1}, t_k) + \mathbf{y}_k^{\text{ct}}, \quad \mathbf{y}_k^{\text{ct}} := \int_{s=t_k}^{t_k} \Phi_1(s, t_k)^\top \Phi_1(s, t_k) ds,$$

where  $\Phi_1(s, s) = \mathbf{I}$  and where (using that we consider  $V^\pi(t_k; \mathbf{u} = 0, \xi)$  with  $\mathbf{u} = 0$ , so the corresponding  $\mathbf{K}(t)$  in [Eq. \(G.8\)](#) vanishes)

$$\frac{d}{dt} \Phi_1(t, s) = \mathbf{A}_{\text{ol}}^\pi(s) \Phi_1(t, s).$$

Hence,  $\Phi_1(t_{k+1}, t_k) = \mathbf{A}_{\text{ol},k}^\pi$ , so that

$$\mathbf{p}_k^{\text{ct}} = (\mathbf{A}_{\text{ol},k}^\pi)^\top \mathbf{p}_{k+1}^{\text{ct}} \mathbf{A}_{\text{ol},k}^\pi + \mathbf{y}_k^{\text{ct}}, \quad \mathbf{y}_k^{\text{ct}} := \int_{s=t_k}^{t_k} \Phi_1(s, t_k)^\top \Phi_1(s, t_k) ds$$

Along the lines of [Lemma G.7](#), it can be shown that for  $\tau \leq \tau_{\text{ric}}$ ,  $\Phi_1(s, t_k)^\top \Phi_1(s, t_k) \succeq \frac{1}{2} \mathbf{I}$  for all  $t \in \mathcal{I}_k$ . Thus,

$$\mathbf{p}_k^{\text{ct}} \succeq (\mathbf{A}_{\text{ol},k}^\pi)^\top \mathbf{p}_{k+1}^{\text{ct}} \mathbf{A}_{\text{ol},k}^\pi + \frac{1}{2} \mathbf{I}.$$

Comparing to [Eq. \(G.12\)](#), we find that  $\mathbf{p}_k^{\text{ct}} \succeq \frac{1}{2} \Lambda_k^\pi$ . As  $\|\mathbf{p}_k^{\text{ct}}\| = \sup_{\xi: \|\xi\|=1} V^\pi(t_k; \mathbf{u} = 0, \xi) \leq \mu_{\text{ric}}$ , we conclude  $\|\Lambda_k^\pi\| \leq 2\mu_{\text{ric}}$ , as needed.

### G.4. Proof of [Lemma A.1](#)

We recall the definitions  $\kappa_{\pi, \infty} := \max_{1 \leq j \leq k \leq K+1} \|\Phi_{\text{cl},k,j}^\pi\|$ , and

$$\begin{aligned} \kappa_{\pi, 1} &:= \max_{k \in [K+1]} \tau \left( \sum_{j=1}^k \|\Phi_{\text{cl},k,j}^\pi\| \vee \sum_{j=k}^{K+1} \|\Phi_{\text{cl},j,k}^\pi\| \right) \\ \kappa_{\pi, 2} &:= \max_{k \in [K+1]} \tau \left( \sum_{j=1}^k \|\Phi_{\text{cl},k,j}^\pi\|^2 \vee \sum_{j=k}^{K+1} \|\Phi_{\text{cl},j,k}^\pi\|^2 \right), \end{aligned}$$

and recall the definitions  $\Lambda_{K+1}^\pi = \mathbf{I}$ , and  $\Lambda_k^\pi = (\mathbf{A}_{\text{cl},k}^\pi)^\top \Lambda_{k+1}^\pi \mathbf{A}_{\text{cl},k}^\pi + \tau \mathbf{I}$ , and  $\mu_{\pi,\star} := \max_{k \in \{k_0, k_0+1, \dots, K+1\}} \|\Lambda_k^\pi\|$ .

Let us first bound  $\|\Phi_{\text{cl},k,j}^\pi\|$  for  $j \geq k_0$ .

**Claim G.7.** For  $j \geq k_0$ , and  $\tau \leq 1/6L_f L_\pi$ ,  $\|\Phi_{\text{cl},k,j}^\pi\| \leq \sqrt{\max\{1, 6L_f L_\pi\} \mu_{\pi,\star} (1 - \tau/\mu_{\pi,\star})^{k-j}}$

*Proof.* We apply [Lemma F.10](#) with  $\mathbf{X}_k \leftarrow \mathbf{A}_{\text{ol},k}^\pi$ ,  $\mathbf{Q} = \mathbf{I}$ , and  $\mathbf{Y}_k = 0$ , and only take the recursion back to  $k = k_0$ . That lemma shows that, as long as  $\|\mathbf{A}_{\text{ol},k}^\pi - \mathbf{I}\| \leq \kappa\tau$  for some  $\kappa \leq 2/\tau$ , it holds that (for  $j \geq k_0$ )

$$\|\Phi_{\text{cl},k,j}^\pi\|^2 \leq \max\{1, 2\kappa\} \mu_{\pi,\star} (1 - \tau/\mu_{\pi,\star})^{k-j}.$$

[Lemmas I.3](#) and [I.4](#), and using  $L_\pi \geq 1$ , we have that for  $\tau \leq 1/4L_f$ ,  $\|\mathbf{A}_{\text{cl},k}^\pi - \mathbf{I}\| \leq \|\mathbf{A}_{\text{ol},k}^\pi - \mathbf{I}\| + \|\mathbf{B}_{\text{ol},k}^\pi \mathbf{K}_k^\pi\| \leq \exp(1/4)\tau L_f (1 + L_\pi) \leq 2 \exp(1/4)\tau L_f L_\pi \leq 3\tau L_f L_\pi$ . Hence, for  $\tau \leq 1/6L_f L_\pi$ , we can take  $\kappa := 3L_f L_\pi$  and have  $\kappa \leq 1/2\tau$ . For this choice of  $\kappa$ , we get

$$\|\Phi_{\text{cl},k,j}^\pi\|^2 \leq \max\{1, 6L_f L_\pi\} \mu_{\pi,\star} (1 - \tau/\mu_{\pi,\star})^{k-j}.$$

□

Next, we bound  $\|\Phi_{\text{cl},k,j}^\pi\|$  for  $k \leq k_0$ .

**Claim G.8.** For  $k \geq k_0$ ,  $\|\Phi_{\text{cl},k,j}^\pi\| \leq \exp(k_0\tau L_f)$ .

*Proof.* For  $j \leq k \leq k_0$ , we have  $\mathbf{K}_j^\pi = 0$ . Hence  $\mathbf{A}_{\text{cl},j}^\pi = \mathbf{A}_{\text{ol},j}^\pi$ , and from [Lemma I.4](#), we get  $\|\mathbf{A}_{\text{cl},j}^\pi\| \leq \exp(\tau L_f)$ . Thus  $\|\Phi_{\text{cl},k,j}^\pi\| \leq \prod_{j=1}^k \|\mathbf{A}_{\text{ol},j}^\pi\| \leq \exp(k\tau L_f) \leq \exp(k_0\tau L_f)$ . □

Finally, we bound we bound  $\|\Phi_{\text{cl},k,j}^\pi\|$  for  $j \leq k_0$ .

**Claim G.9.** For  $j \leq k_0$ ,  $\|\Phi_{\text{cl},k,j}^\pi\| \leq \sqrt{\max\{1, 6L_f L_\pi\} \mu_{\pi,\star} \exp(\tau k_0 L_f)}$ .

*Proof.* For  $j < k_0$ , we have  $\|\Phi_{\text{cl},k,j}^\pi\| = \|\Phi_{\text{cl},k,k_0}^\pi \Phi_{\text{cl},k_0,j}^\pi\| \leq \|\Phi_{\text{cl},k,k_0}^\pi\| \|\Phi_{\text{cl},k_0,j}^\pi\|$ . The first term is at most  $\max\{1, 6L_f L_\pi\} \mu_{\pi,\star}$  by [Claim G.7](#), and the second term at most  $\exp(\tau k_0 L_f)$  by [Claim G.8](#). □

We can now bound all terms of interest. Directly from the dichotomy in [Claim G.7](#), we have  $\kappa_{\pi,\infty} \leq \sqrt{\max\{1, 6L_f L_\pi\} \mu_{\pi,\star} \exp(\tau k_0 L_f)}$ . Next, for any  $k$ , we can bound via [Claims G.7](#) and [G.9](#)

$$\begin{aligned} \tau \sum_{j=1}^k \|\Phi_{\text{cl},k,j}^\pi\|^2 &\leq \tau \sum_{j=1}^{k_0} \|\Phi_{\text{cl},k,j}^\pi\|^2 + \tau \sum_{j=k_0}^k \|\Phi_{\text{cl},k,j}^\pi\|^2 \\ &\leq \max\{1, 6L_f L_\pi\} \mu_{\pi,\star} \left( (\tau k_0) \exp(2\tau k_0 L_f) + \tau \sum_{j=k_0}^k (1 - \tau/\mu_{\pi,\star})^{k-k_0} \right) \quad (\text{Claims G.7 and G.9}) \\ &\leq \max\{1, 6L_f L_\pi\} \mu_{\pi,\star} \left( (\tau k_0) \exp(2\tau k_0 L_f) + \tau \underbrace{\sum_{n=0}^{\infty} (1 - \tau/\mu_{\pi,\star})^n}_{= \frac{1}{1 - (1 - \tau/\mu_{\pi,\star})} = \mu_{\pi,\star}/\tau} \right) \\ &\leq \max\{1, 6L_f L_\pi\} \mu_{\pi,\star} ((\tau k_0) \exp(2\tau k_0 L_f) + \mu_{\pi,\star}). \end{aligned}$$

and show the same bound for  $\tau \sum_{j=k}^{K+1} \|\Phi_{\text{cl},j,k}^\pi\|^2$ , which yields the desired upper bound on  $\kappa_{\pi,2}$ . Finally, to bound  $\kappa_{\pi,1}$ ,

$$\begin{aligned}
 \tau \sum_{j=1}^k \|\Phi_{\text{cl},k,j}^\pi\| &\leq \tau \sum_{j=1}^{k_0} \|\Phi_{\text{cl},k,j}^\pi\| + \tau \sum_{j=k_0}^k \|\Phi_{\text{cl},k,j}^\pi\| \\
 &\leq \sqrt{\max\{1, 6L_f L_\pi\} \mu_{\pi,\star}} \left( (\tau k_0) \exp(\tau k_0 L_f) + \tau \sum_{j=k_0}^k \sqrt{(1 - \tau/\mu_{\pi,\star})^{k-k_0}} \right) \quad (\text{Claims G.7 and G.9}) \\
 &\leq \sqrt{\max\{1, 6L_f L_\pi\} \mu_{\pi,\star}} \left( (\tau k_0) \exp(\tau k_0 L_f) + \tau \sum_{n=0}^{\infty} \sqrt{(1 - \tau/\mu_{\pi,\star})^n} \right) \\
 &\stackrel{(i)}{\leq} \sqrt{\max\{1, 6L_f L_\pi\} \mu_{\pi,\star}} \left( (\tau k_0) \exp(\tau k_0 L_f) + 2\tau \sum_{n=0}^{\infty} (1 - \tau/\mu_{\pi,\star})^n \right) \\
 &= \sqrt{\max\{1, 6L_f L_\pi\} \mu_{\pi,\star}} ((\tau k_0) \exp(\tau k_0 L_f) + 2\mu_{\pi,\star}),
 \end{aligned}$$

where in (i), we use that  $\sum_{n \geq 0} \sqrt{(1 - \gamma)^n} = \sum_{n \geq 0} \sqrt{(1 - \gamma)^{2n}} + \sqrt{(1 - \gamma)^{2n+1}} \leq 2 \sum_{n \geq 0} \sqrt{(1 - \gamma)^{2n}} = 2 \sum_{n \geq 0} (1 - \gamma)^n$ . One can establish the same bound for  $\tau \sum_{j=k}^{K+1} \|\Phi_{\text{cl},j,k}^\pi\|$ , which gives the desired bound on  $\kappa_{\pi,1}$ .  $\square$



## H. Optimization Proofs

### H.1. Proof of Descent Lemma (Lemma A.13)

*Proof of Lemma A.13.* For simplicity, write  $\text{Err}_\nabla = \text{Err}_{\nabla, \pi^{(n)}}(\delta)$ ,  $L_{\nabla, \pi, \infty} = L_{\nabla, \pi^{(n)}, \infty}$ , and  $\text{Err}_{\hat{x}} = \text{Err}_{\hat{x}}(\delta)$  and  $\text{Err}_\nabla(\delta)$  and  $M \geq M_{\mathcal{J}, \text{tay}, \pi^{(n)}}$ . Note that if  $\pi$  and  $\tilde{\pi}$  have the same input sequence but possibly different gains,  $\mathcal{J}_T^{\text{disc}}(\pi) = \mathcal{J}_T^{\text{disc}}(\tilde{\pi})$ . Therefore,

$$\mathcal{J}_T^{\text{disc}}(\pi^{(n+1)}) = \mathcal{J}_T^{\text{disc}}(\tilde{\pi}^{(n+1)}).$$

Define the input

$$\check{\mathbf{u}}_k^{(n)} = \mathbf{u}_k^{(n)} - \eta \hat{\nabla}_k^{(n)} + \mathbf{K}_k^{\pi^{(n)}}(\mathbf{x}_k^\pi - \hat{\mathbf{x}}_k)$$

Then, as in Eq. (E.3),

$$\mathbf{u}_k^{\tilde{\pi}^{(n)}} = \mathbf{u}_{\text{orac}, k}^{\pi^{(n)}}(\tilde{\mathbf{u}}_{1:K}) = \check{\mathbf{u}}_k^{\pi^{(n)}}(\check{\mathbf{u}}_{1:K}), \quad \mathbf{x}_k^{\tilde{\pi}^{(n)}} = \mathbf{x}_{\text{orac}, k}^{\pi^{(n)}}(\tilde{\mathbf{u}}_{1:K}) = \check{\mathbf{x}}_k^{\pi^{(n)}}(\check{\mathbf{u}}_{1:K}), \quad \forall k \in [K].$$

Consequently, we have the quality

$$\mathcal{J}_T^{\text{disc}}(\tilde{\pi}^{(n)}) := \mathcal{J}_T^{\tilde{\pi}^{(n+1)}, \text{disc}}(\mathbf{u}_{1:K}^{\tilde{\pi}^{(n)}}) = \mathcal{J}_T^{\pi^{(n)}, \text{disc}}(\check{\mathbf{u}}_{1:K}) = \mathcal{J}_T^{\pi^{(n)}, \text{disc}}(\mathbf{u}_{1:K}^{\pi^{(n)}} + \delta \check{\mathbf{u}}_k^{(n)}),$$

where we introduced

$$\delta \check{\mathbf{u}}_k^{(n)} := \check{\mathbf{u}}_k^{(n)} - \mathbf{u}_k^{(n)} = \underbrace{-\frac{\eta}{\tau} \hat{\nabla}_k^{(n)}}_{=\delta \check{\mathbf{u}}_k^{(n;1)}} + \underbrace{\mathbf{K}_k^{\pi^{(n)}}(\mathbf{x}_k^\pi - \hat{\mathbf{x}}_k)}_{=\delta \check{\mathbf{u}}_k^{(n;2)}}$$

**Claim H.1.** *We have*

$$\begin{aligned} \sqrt{\tau} \|\delta \check{\mathbf{u}}_{1:K}^{(n)}\|_{\ell_2} &\leq \sqrt{T}(\eta(L_{\nabla, \pi, \infty} + \frac{1}{\tau} \text{Err}_\nabla) + \text{Err}_{\hat{x}}) \\ \max_k \|\delta \check{\mathbf{u}}_k^{(n)}\| &\leq (\eta(L_{\nabla, \pi, \infty} + \frac{1}{\tau} \text{Err}_\nabla) + \text{Err}_{\hat{x}}) \end{aligned}$$

*Proof.* The first bound follows from the second. We have that

$$\begin{aligned} \max_k \|\delta \check{\mathbf{u}}_k^{(n)}\| &\leq \max_k \|\delta \check{\mathbf{u}}_k^{(n;1)}\| + \max_k \|\delta \check{\mathbf{u}}_k^{(n;2)}\| \\ &\leq \frac{\eta}{\tau} \max_k \|\hat{\nabla}_k^{(n)}\| + \text{Err}_{\hat{x}} \\ &\leq \frac{\eta}{\tau} \left( \max_k \|\mathcal{J}_T^{\pi^{(n)}, \text{disc}}(\mathbf{u}_{1:K}^{\pi^{(n)}})\| + \text{Err}_\nabla \right) + \text{Err}_{\hat{x}} \\ &\leq \frac{\eta}{\tau} (\tau L_{\nabla, \pi, \infty} + \text{Err}_\nabla) + \text{Err}_{\hat{x}} \tag{Lemma A.8} \\ &= \eta(L_{\nabla, \pi, \infty} + \frac{1}{\tau} \text{Err}_\nabla) + \text{Err}_{\hat{x}} \end{aligned}$$

□

As a consequence of the above claim, it holds that if

$$(\eta(L_{\nabla, \pi, \infty} + \frac{1}{\tau} \text{Err}_\nabla) + \text{Err}_{\hat{x}}) \leq \min \left\{ \frac{R_{\text{feas}}}{8}, B_{\text{stab}, \pi}, B_{\text{tay}, \text{inf}, \pi}, \frac{B_{\text{tay}, 2, \pi}}{\sqrt{T}} \right\},$$

then (a) Lemma A.7 implies stability of  $\tilde{\pi}^{(n)}$ :

$$\mu_{\tilde{\pi}^{(n)}, \star} \leq 2\mu_{\pi^{(n)}, \star}, \quad L_{\tilde{\pi}^{(n)}} = L_{\pi^{(n)}},$$

and (b) the Taylor expansion in [Lemma A.6](#) implies that

$$\begin{aligned}
 \mathcal{J}_T^{\text{disc}}(\tilde{\pi}^{(n)}) &= \mathcal{J}_T^{\pi^{(n)}, \text{disc}}(\mathbf{u}_{1:K}^{\pi^{(n)}} + \delta\check{\mathbf{u}}_k^{(n)}) \\
 &= \mathcal{J}_T^{\text{disc}}(\pi^{(n)}) + \left\langle \delta\check{\mathbf{u}}_{1:K}^{(n)}, \nabla \mathcal{J}_T^{\pi^{(n)}, \text{disc}}(\mathbf{u}_{1:K}^{\pi^{(n)}}) \right\rangle + \frac{R\tau}{2} \|\delta\check{\mathbf{u}}_{1:K}^{(n)}\|_{\ell_2}^2 \\
 &\leq \mathcal{J}_T^{\text{disc}}(\pi^{(n)}) + \underbrace{\sum_{i=1}^2 \left\langle \delta\check{\mathbf{u}}_{1:K}^{(n;i)}, \nabla \mathcal{J}_T^{\pi^{(n)}, \text{disc}}(\mathbf{u}_{1:K}^{\pi^{(n)}}) \right\rangle}_{\text{Term}_i} + M\tau \|\delta\check{\mathbf{u}}_{1:K}^{(n;i)}\|_{\ell_2}^2.
 \end{aligned} \tag{AM-GM}$$

It remains to massage the above display to obtain the descent guarantee:

$$\begin{aligned}
 \text{Term}_1 &= \left\langle \delta\check{\mathbf{u}}_{1:K}^{(n;1)}, \nabla \mathcal{J}_T^{\pi^{(n)}, \text{disc}}(\mathbf{u}_{1:K}^{\pi^{(n)}}) \right\rangle + M\tau \|\delta\check{\mathbf{u}}_{1:K}^{(n;1)}\|_{\ell_2}^2 \\
 &\leq \langle \delta\check{\mathbf{u}}_{1:K}^{(n;1)}, \hat{\nabla}_{1:K}^{(n)} \rangle + \|\delta\check{\mathbf{u}}_k^{(n;1)}\|_{\ell_2} \sqrt{K} \text{Err}_{\nabla} + M\tau \|\delta\check{\mathbf{u}}_{1:K}^{(n;1)}\|_{\ell_2}^2 \\
 &\leq \langle \delta\check{\mathbf{u}}_{1:K}^{(n;1)}, \hat{\nabla}_{1:K}^{(n)} \rangle + \frac{K}{4R\tau} \text{Err}_{\nabla}^2 + 2M\tau \|\delta\check{\mathbf{u}}_{1:K}^{(n;1)}\|_{\ell_2}^2 \\
 &= \left(-\frac{\eta}{\tau} + 2M\frac{\eta^2}{\tau}\right) \|\hat{\nabla}_{1:K}^{(n)}\|_{\ell_2}^2 + \frac{T}{4R\tau^2} \text{Err}_{\nabla}^2 \\
 &\geq -\frac{\eta}{2\tau} \|\hat{\nabla}_{1:K}^{(n)}\|_{\ell_2}^2 + \frac{T}{4M\tau^2} \text{Err}_{\nabla}^2,
 \end{aligned}$$

where the last step uses  $\eta \leq \frac{1}{4M}$ . Then,

$$\begin{aligned}
 \text{Term}_2 &= \left\langle \delta\check{\mathbf{u}}_{1:K}^{(n;2)}, \nabla \mathcal{J}_T^{\pi^{(n)}, \text{disc}}(\mathbf{u}_{1:K}^{\pi^{(n)}}) \right\rangle + M\tau \|\delta\check{\mathbf{u}}_{1:K}^{(n;2)}\|_{\ell_2}^2 \\
 &\leq \|\delta\check{\mathbf{u}}_{1:K}^{(n;2)}\|_{\ell_2} \|\nabla \mathcal{J}_T^{\pi^{(n)}, \text{disc}}(\mathbf{u}_{1:K}^{\pi^{(n)}})\|_{\ell_2} + M\tau \|\delta\check{\mathbf{u}}_{1:K}^{(n;2)}\|_{\ell_2}^2 \\
 &\leq \sqrt{\tau} \|\delta\check{\mathbf{u}}_{1:K}^{(n;2)}\|_{\ell_2} \cdot \sqrt{K} \tau L_{\nabla, \pi, \infty} + M\tau \|\delta\check{\mathbf{u}}_{1:K}^{(n;2)}\|_{\ell_2}^2 \\
 &\leq \sqrt{K} \max_k \|\delta\check{\mathbf{u}}_k^{(n;2)}\| \cdot \sqrt{K} \tau L_{\nabla, \pi, \infty} + M\tau K \max_k \|\delta\check{\mathbf{u}}_k^{(n;2)}\|^2 \\
 &= T(\text{Err}_{\hat{x}} L_{\nabla, \pi, \infty} + M \text{Err}_{\hat{x}}^2).
 \end{aligned} \tag{Lemma A.8}$$

Thus,

$$\begin{aligned}
 \mathcal{J}_T^{\text{disc}}(\pi^{(n+1)}) - \mathcal{J}_T^{\text{disc}}(\pi^{(n)}) &\leq \text{Term}_1 + \text{Term}_2 \\
 &\leq -\frac{\eta}{2\tau} \|\hat{\nabla}_{1:K}^{(n)}\|_{\ell_2}^2 + T\left(\frac{1}{4M\tau^2} \text{Err}_{\nabla}^2 + \text{Err}_{\hat{x}} L_{\nabla, \pi, \infty} + M \text{Err}_{\hat{x}}^2\right).
 \end{aligned}$$

□

## H.2. Proof of [Proposition 4.1](#)

The proof of [Proposition 4.1](#) makes liberal use of the definitions of the linearizations given in [Appendix C.1](#), which we recall without further comment. Going forward, introduce the Jacobian linearization of the stabilized cost:

$$\mathcal{J}_T^{\pi, \text{jac}}(\bar{\mathbf{u}}) := V(\tilde{\mathbf{x}}^{\pi, \text{jac}}(t | \bar{\mathbf{u}})) + \int_0^T Q(\tilde{\mathbf{x}}^{\pi, \text{jac}}(t | \bar{\mathbf{u}}), \tilde{\mathbf{u}}^{\pi, \text{jac}}(t | \bar{\mathbf{u}}), t) dt.$$

We now characterize some properties of  $\mathcal{J}_T^{\pi, \text{jac}}$ .

**Lemma H.1** (Valid First-Order Approximation). *We have that  $\nabla_{\bar{\mathbf{u}}} \mathcal{J}_T^{\pi}(\bar{\mathbf{u}}) = \nabla_{\bar{\mathbf{u}}} \mathcal{J}_T^{\pi, \text{jac}}(\bar{\mathbf{u}})$ .*

*Proof.* Immediate from the chain rule, and the fact that the Jacobian linearizations are defined as the first-order Taylor expansion of the true dynamics. □

**Lemma H.2** (Congruence with the Open-Loop).

$$\inf_{\tilde{\mathbf{u}}} \mathcal{J}_T^{\pi, \text{jac}}(\tilde{\mathbf{u}}) = \inf_{\tilde{\mathbf{u}}} \mathcal{J}_T^{\text{jac}}(\tilde{\mathbf{u}}; \mathbf{u}^\pi).$$

*Proof.* We prove  $\inf_{\tilde{\mathbf{u}}} \mathcal{J}_T^{\pi, \text{jac}}(\tilde{\mathbf{u}}) \leq \inf_{\tilde{\mathbf{u}}} \mathcal{J}_T^{\text{jac}}(\tilde{\mathbf{u}}; \mathbf{u}^\pi)$ ; the converse can be proved similarly. Fix any  $\bar{\mathbf{u}}_1 \in \mathcal{U}$ . It suffices to exhibit some  $\bar{\mathbf{u}}_2 \in \mathcal{U}$  such that, for all  $t \in [0, T]$ ,

$$\tilde{\mathbf{x}}^{\pi, \text{jac}}(t | \bar{\mathbf{u}}_2) = \mathbf{x}^{\text{jac}}(t | \bar{\mathbf{u}}_1; \mathbf{u}^\pi), \quad \tilde{\mathbf{u}}^{\pi, \text{jac}}(t | \bar{\mathbf{u}}_2) = \bar{\mathbf{u}}_1(t).$$

By subtracting off  $\mathbf{x}^\pi(t)$  and  $\mathbf{u}^\pi(t)$ , it suffices to show that

$$\delta \tilde{\mathbf{x}}^{\pi, \text{jac}}(t | \bar{\mathbf{u}}_2) = \delta \mathbf{x}^{\text{jac}}(t | \bar{\mathbf{u}}_1; \mathbf{u}^\pi), \quad \delta \tilde{\mathbf{u}}^{\pi, \text{jac}}(t | \bar{\mathbf{u}}_2) = \delta \bar{\mathbf{u}}_1(t).$$

It can be directly checked from [Lemmas C.1](#) and [C.2](#) that the input  $\bar{\mathbf{u}}_2(t) = \bar{\mathbf{u}}_1(t) - \mathbf{K}_{k(t)}^\pi \delta \mathbf{x}^{\text{jac}}(t | \bar{\mathbf{u}}_1; \mathbf{u}^\pi)$  ensures the above display holds.  $\square$

The last lemma contains our main technical endeavor, and its proof is deferred to [Lemma H.3](#) just below.

**Lemma H.3** (Strong Convexity). *Suppose  $\tau \leq \min\{\frac{1}{4L_f}, \frac{1}{16L_\pi L_f}\}$ . Then,  $\bar{\mathbf{u}} \mapsto \mathcal{J}_T^{\pi, \text{jac}}(\bar{\mathbf{u}}) - \alpha_\pi \|\bar{\mathbf{u}}\|_{\mathcal{L}_2(\mathcal{U})}^2$  is convex, where  $\alpha_\pi := \frac{\alpha}{64 \max\{1, L_\pi^2\}}$ .*

We may now conclude the proof of our proposition.

*Proof of Proposition 4.1.* Suppose that  $\pi$  satisfies [Definition 4.7](#), and suppose  $\tau \leq \min\{\frac{1}{4L_f}, \frac{1}{16L_\pi L_f}\}$  and  $\|\nabla_{\bar{\mathbf{u}}} \mathcal{J}_T^\pi(\bar{\mathbf{u}})|_{\bar{\mathbf{u}}=\mathbf{u}^\pi}\|_{\mathcal{L}_2(\mathcal{U})} \leq \epsilon_0$ . By [Lemma H.1](#),  $\|\nabla_{\bar{\mathbf{u}}} \mathcal{J}_T^{\pi, \text{jac}}(\bar{\mathbf{u}})|_{\bar{\mathbf{u}}=\mathbf{u}^\pi}\|_{\mathcal{L}_2(\mathcal{U})} \leq \epsilon_0$ . By [Lemma H.3](#) and the fact that strong convex functions satisfy the PL-inequality (e.g. [Karimi et al. \(2016, Theorem 2\)](#)), we have

$$\mathcal{J}_T^{\pi, \text{jac}}(\mathbf{u}^\pi) \leq \inf_{\bar{\mathbf{u}}} \mathcal{J}_T^{\pi, \text{jac}}(\bar{\mathbf{u}}) + \frac{\epsilon_0^2}{\alpha_\pi} = \inf_{\bar{\mathbf{u}}} \mathcal{J}_T^{\pi, \text{jac}}(\bar{\mathbf{u}}) + \frac{64\epsilon_0^2 \max\{1, L_\pi^2\}}{\alpha}.$$

Finally, by [Lemma H.2](#),  $\inf_{\bar{\mathbf{u}}} \mathcal{J}_T^{\pi, \text{jac}}(\bar{\mathbf{u}}) = \inf_{\bar{\mathbf{u}}} \mathcal{J}_T^{\text{jac}}(\bar{\mathbf{u}}; \mathbf{u}^\pi)$ , which implies the proposition.  $\square$

### H.2.1. PROOF OF [LEMMA H.3](#)

*Proof.* We claim that suffices to show the following PSD lower bound:

$$\forall \bar{\mathbf{u}} \in \mathcal{U}, \mathcal{Q}^\pi(\bar{\mathbf{u}}) \geq \frac{\alpha_\pi}{\alpha} \|\bar{\mathbf{u}}\|_{\mathcal{L}_2}^2, \tag{H.1}$$

where we define

$$\mathcal{Q}^\pi(\bar{\mathbf{u}}) := \int_0^T (\|\delta \tilde{\mathbf{x}}^{\pi, \text{jac}}(t | \bar{\mathbf{u}})\|^2 + \|\delta \tilde{\mathbf{u}}^{\pi, \text{jac}}(t | \bar{\mathbf{u}})\|^2) dt,$$

and where we define the deviations  $\delta(\cdot)$  as in [Lemma C.2](#).

**Claim H.2.** *If Eq. (H.1) holds, then [Lemma H.3](#) holds.*

*Proof of Claim H.2.* Note that  $\bar{\mathbf{u}} \mapsto \tilde{\mathbf{x}}^{\pi, \text{jac}}(t | \bar{\mathbf{u}})$  and  $\bar{\mathbf{u}} \mapsto \tilde{\mathbf{u}}^{\pi, \text{jac}}(t | \bar{\mathbf{u}})$  are affine, that  $\delta \tilde{\mathbf{x}}^{\pi, \text{jac}}(t | \bar{\mathbf{u}}) = \tilde{\mathbf{x}}^{\pi, \text{jac}}(t | \bar{\mathbf{u}}) - \mathbf{x}^\pi(t)$  and  $\delta \tilde{\mathbf{u}}^{\pi, \text{jac}}(t | \bar{\mathbf{u}}) = \tilde{\mathbf{u}}^{\pi, \text{jac}}(t | \bar{\mathbf{u}}) - \mathbf{u}^\pi(t)$  are linear (no affine term), and that the differences  $\tilde{\mathbf{x}}^{\pi, \text{jac}}(t | \bar{\mathbf{u}}) - \delta \tilde{\mathbf{x}}^{\pi, \text{jac}}(t | \bar{\mathbf{u}}) = \mathbf{x}^\pi(t)$  and  $\tilde{\mathbf{u}}^{\pi, \text{jac}}(t | \bar{\mathbf{u}}) - \delta \tilde{\mathbf{u}}^{\pi, \text{jac}}(t | \bar{\mathbf{u}}) = \mathbf{u}^\pi(t)$  are independent of  $\bar{\mathbf{u}}$ . Hence, we conclude  $\mathcal{Q}^\pi(\bar{\mathbf{u}})$  is a quadratic function with no linear term, and  $\tilde{\mathcal{Q}}^\pi(\bar{\mathbf{u}}) - \mathcal{Q}^\pi(\bar{\mathbf{u}})$  is linear, where we define

$$\tilde{\mathcal{Q}}^\pi(\bar{\mathbf{u}}) := \int_0^T (\|\tilde{\mathbf{x}}^{\pi, \text{jac}}(t | \bar{\mathbf{u}})\|^2 + \|\tilde{\mathbf{u}}^{\pi, \text{jac}}(t | \bar{\mathbf{u}})\|^2) dt,$$

[Assumption 2.1](#) implies that  $\mathcal{J}_T^{\pi, \text{jac}}(\bar{\mathbf{u}}) - \alpha \tilde{\mathcal{Q}}^\pi(\bar{\mathbf{u}})$  is convex, and since the difference  $\tilde{\mathcal{Q}}^\pi(\bar{\mathbf{u}}) - \mathcal{Q}^\pi(\bar{\mathbf{u}})$  is linear, that  $\mathcal{J}_T^{\pi, \text{jac}}(\bar{\mathbf{u}}) - \alpha \tilde{\mathcal{Q}}^\pi(\bar{\mathbf{u}})$  is also convex. Lastly, as  $\mathcal{Q}^\pi(\bar{\mathbf{u}})$  is quadratic with no linear term, [Eq. \(H.1\)](#) implies  $\alpha \mathcal{Q}^\pi(\bar{\mathbf{u}}) - \alpha_\pi \|\bar{\mathbf{u}}\|_{\mathcal{L}_2}^2$  is convex. Thus,  $\mathcal{J}_T^{\pi, \text{jac}}(\bar{\mathbf{u}}) - \alpha_\pi \|\bar{\mathbf{u}}\|_{\mathcal{L}_2}^2 = (\mathcal{J}_T^{\pi, \text{jac}}(\bar{\mathbf{u}}) - \alpha \tilde{\mathcal{Q}}^\pi(\bar{\mathbf{u}})) - (\alpha \mathcal{Q}^\pi(\bar{\mathbf{u}}) - \alpha_\pi \|\bar{\mathbf{u}}\|_{\mathcal{L}_2}^2)$  is convex.  $\square$

To verify Eq. (H.1), let us define a few salient operators. Let  $\mathcal{X}$  denote the space of  $\mathcal{L}_2$  bounded curves  $\mathbf{x}(t) \in \mathbb{R}^{d_x}$ . We define linear operators  $\mathbb{T}_1 : \mathcal{U} \rightarrow \mathcal{X}$  and  $\mathbb{T}_2 : \mathcal{U} \rightarrow \mathcal{X}$  and  $\mathbb{K} : \mathcal{X} \rightarrow \mathcal{U}$  via

$$\mathbb{T}_1[\bar{\mathbf{u}}](t) = \delta \bar{\mathbf{x}}^{\pi, \text{jac}}(t_{k(t)} | \bar{\mathbf{u}}), \quad \mathbb{T}_2[\bar{\mathbf{u}}](t) = \delta \bar{\mathbf{x}}^{\pi, \text{jac}}(t | \bar{\mathbf{u}}), \quad \mathbb{K}[\bar{\mathbf{x}}](t) = \mathbb{K}_{k(t)}^{\pi} \bar{\mathbf{x}}(t).$$

Then, letting  $\mathbb{I}_{\mathcal{U}}$  denote the identity operator of  $\mathcal{U}$ , we can write

$$\mathcal{Q}^{\pi}(\bar{\mathbf{u}}) = \|\mathbb{T}_2[\bar{\mathbf{u}}]\|_{\mathcal{L}_2(\mathcal{X})}^2 + \|(\mathbb{I}_{\mathcal{U}} + \mathbb{K}\mathbb{T}_1)[\bar{\mathbf{u}}]\|_{\mathcal{L}_2(\mathcal{U})}^2.$$

Next, we relate  $\mathbb{T}_2$  and  $\mathbb{T}_1$ . Define the operators  $\mathbb{L} : \mathcal{X} \rightarrow \mathcal{X}$  and  $\mathbb{W} : \mathcal{U} \rightarrow \mathcal{X}$  by

$$\mathbb{L}[\bar{\mathbf{x}}](t) = \Phi_{\text{ol}}^{\pi}(t, t_{k(t)})\bar{\mathbf{x}}(t), \quad \mathbb{W}[\bar{\mathbf{u}}](t) = \int_{s=t_{k(t)}}^t \Phi_{\text{ol}}^{\pi}(t, s)\mathbf{B}_{\text{ol}}^{\pi}(s)\bar{\mathbf{u}}(s)ds.$$

Then, it can be checked from Lemmas C.2 and C.10 that

$$\mathbb{T}_2[\bar{\mathbf{u}}] = \mathbb{L}\mathbb{T}_1[\bar{\mathbf{u}}] + \mathbb{W}[\bar{\mathbf{u}}].$$

Hence,

$$\mathcal{Q}^{\pi}(\bar{\mathbf{u}}) = \|(\mathbb{L}\mathbb{T}_1 + \mathbb{W})[\bar{\mathbf{u}}]\|_{\mathcal{L}_2(\mathcal{X})}^2 + \|(\mathbb{I}_{\mathcal{U}} + \mathbb{K}\mathbb{T}_1)[\bar{\mathbf{u}}]\|_{\mathcal{L}_2(\mathcal{U})}^2.$$

With this representation of  $\mathcal{Q}^{\pi}$ , we establish a lower bound by applying the following lemma, whose proof we below:

**Lemma H.4.** *Let  $\mathcal{A}, \mathcal{B}$  be Hilbert spaces with norms  $\|\cdot\|_{\mathcal{A}}$  and  $\|\cdot\|_{\mathcal{B}}$ , let  $\mathbb{I}_{\mathcal{A}}$  denote the identity operator on  $\mathcal{U}$ , and let  $\mathbb{T}, \mathbb{W} : \mathcal{A} \rightarrow \mathcal{B}$ ,  $\mathbb{L} : \mathcal{B} \rightarrow \mathcal{B}$ , and  $\mathbb{K} : \mathcal{B} \rightarrow \mathcal{A}$  be linear operators, and let  $\|\cdot\|_{\text{op}}$  denote operator norms. Then, if  $\|\mathbb{W}\|_{\text{op}} \leq \frac{\min\{1, \sigma_{\min}(\mathbb{L})\}}{4\|\mathbb{K}\|}$ , it holds for any  $\mathbf{a} \in \mathcal{A}$ ,*

$$\|(\mathbb{L}\mathbb{T} + \mathbb{W})[\mathbf{a}]\|_{\mathcal{B}}^2 + \|(\mathbb{I}_{\mathcal{A}} + \mathbb{K}\mathbb{T})[\mathbf{a}]\|_{\mathcal{A}}^2 \geq \|\mathbf{a}\|^2 \cdot \frac{\min\{1, \sigma_{\min}(\mathbb{L})\}^2}{16 \max\{1, \|\mathbb{K}\|_{\text{op}}^2\}},$$

where  $\sigma_{\min}(\mathbb{L}) := \inf_{\mathbf{a}: \|\mathbf{a}\|_{\mathcal{A}}=1} \|\mathbb{L}\mathbf{a}\|_{\mathcal{A}}$ .

To apply the lemma, we first perform a few computations. Throughout, we use  $\|\cdot\|_{\text{op}}$  to denote operator norm, and  $\sigma_{\min}$  to denote minimal singular value as an operator, e.g.

$$\|\mathbb{K}\|_{\text{op}} = \sup_{\bar{\mathbf{x}} \neq 0} \frac{\|\mathbb{K}[\bar{\mathbf{x}}]\|_{\mathcal{L}_2(\mathcal{U})}}{\|\bar{\mathbf{x}}\|_{\mathcal{L}_2(\mathcal{X})}}, \quad \sigma_{\min}(\mathbb{L}) = \inf_{\bar{\mathbf{x}} \neq 0} \frac{\|\mathbb{L}[\bar{\mathbf{x}}]\|_{\mathcal{L}_2(\mathcal{U})}}{\|\bar{\mathbf{x}}\|_{\mathcal{L}_2(\mathcal{X})}}.$$

**Claim H.3.** *Under Definition 4.7,  $\|\mathbb{K}\|_{\text{op}} \leq L_{\pi}$ .*

*Proof.* Since  $\mathbb{K}$  is a (block-)diagonal operator in  $t$ , i.e.  $\mathbb{K}[\bar{\mathbf{x}}](t) = \mathbb{K}_{k(t)}^{\pi} \bar{\mathbf{x}}(t)$ , its  $\mathcal{L}_2(\mathcal{X}) \rightarrow \mathcal{L}_2(\mathcal{U})$  operator is bounded by  $\max_k \|\mathbb{K}_k^{\pi}\|$ , which is at most  $L_{\pi}$  under Definition 4.7.  $\square$

**Claim H.4.** *For  $\tau \leq L_f/4$ ,  $\sigma_{\min}(\mathbb{L}) \geq 1 - \tau L_f \exp(\tau L_f) \geq \frac{1}{2}$ .*

*Proof.* Again, since  $\mathbb{L}$  is a block-diagonal operator in  $t$ , i.e.  $\mathbb{L}[\bar{\mathbf{x}}](t) = \Phi_{\text{ol}}^{\pi}(t, t_{k(t)})\bar{\mathbf{x}}(t)$ ,  $\sigma_{\min}(\mathbb{L}) = \inf_{t \in [0, T]} \sigma_{\min}(\Phi_{\text{ol}}^{\pi}(t, t_{k(t)}))$ . By  $\dots$ ,  $\|\Phi_{\text{ol}}^{\pi}(t, t_{k(t)}) - \mathbb{I}\| \leq \tau L_f \tau(\tau L_f) \leq 1/2$ . Thus,  $\inf_{t \in [0, T]} \sigma_{\min}(\Phi_{\text{ol}}^{\pi}(t, t_{k(t)})) \geq 1/2$ .  $\square$

**Claim H.5.**  *$\|\mathbb{W}\|_{\text{op}} \leq \tau L_f \exp(\tau L_f)$ , which is at most  $2\tau L_f$  for  $\tau \leq L_f/4$ .*

*Proof.* For any  $\bar{\mathbf{u}}$ , we bound

$$\begin{aligned}
 \|\mathbb{W}[\bar{\mathbf{u}}]\|_{\mathcal{L}_2(\mathcal{U})}^2 &= \int_{t=0}^T \left\| \int_{s=t_k(t)}^t \Phi_{\text{ol}}^\pi(t, s) \mathbf{B}_{\text{ol}}^\pi(s) \bar{\mathbf{u}}(s) ds \right\|^2 \\
 &\leq \int_{t=0}^T (t - t_k(t)) \int_{s=t_k(t)}^t \|\Phi_{\text{ol}}^\pi(t, s) \mathbf{B}_{\text{ol}}^\pi(s) \bar{\mathbf{u}}(s)\|^2 ds && \text{(Cauchy-Schwartz)} \\
 &\leq \tau \exp(2\tau L_f) \int_{t=0}^T \int_{s=t_k(t)}^t \|\mathbf{B}_{\text{ol}}^\pi(s)\|^2 \|\bar{\mathbf{u}}(s)\|^2 ds && (t - t_k(t) \leq \tau, \text{ Lemma I.4}) \\
 &\leq \tau L_f^2 \exp(2\tau L_f) \int_{t=0}^T \int_{s=t_k(t)}^t \|\bar{\mathbf{u}}(s)\|^2 ds && \text{(Assumption 4.1)} \\
 &= \tau L_f^2 \exp(2\tau L_f) \int_{s=0}^T \int_{t=s}^{t_{k+1}(t)} \|\bar{\mathbf{u}}(s)\|^2 ds \\
 &\leq \tau^2 L_f^2 \exp(2\tau L_f) \int_{s=0}^T \|\bar{\mathbf{u}}(s)\|^2 ds \\
 &= (\tau L_f \exp(\tau L_f))^2 \|\bar{\mathbf{u}}\|_{\mathcal{L}_2(\mathcal{U})}^2.
 \end{aligned}$$

□

With the above three claims, Lemma H.4 implies that as long as  $\tau \leq \min\{\frac{1}{4L_f}, \frac{1}{16L_\pi L_f}\}$ ,

$$\mathcal{Q}^\pi(\bar{\mathbf{u}}) = \|(\mathbb{L}\mathbb{T}_1 + \mathbb{W})[\bar{\mathbf{u}}]\|_{\mathcal{L}_2(\mathcal{X})}^2 + \|(\mathbf{I}_U + \mathbb{K}\mathbb{T}_1)[\bar{\mathbf{u}}]\|_{\mathcal{L}_2(\mathcal{U})}^2 \geq \frac{1}{64 \max\{1, L_\pi^2\}} = \frac{\alpha_\pi}{\alpha}.$$

□

*Proof of Lemma H.4.* Without loss of generality, assume  $\|\mathbf{a}\|_{\mathcal{A}} = 1$ , and define  $\eta = \frac{\min\{1, \sigma_{\min}(\mathbb{L})\}}{2\|\mathbb{K}\|_{\text{op}}}$ . We consider two cases: First, if  $\|\mathbb{T}\mathbf{a}\| \leq \eta$ , then

$$\|(\mathbb{L}\mathbb{T} + \mathbb{W})[\mathbf{a}]\|_{\mathcal{B}}^2 + \|(\mathbf{I}_{\mathcal{A}} + \mathbb{K}\mathbb{T})[\mathbf{a}]\|_{\mathcal{A}}^2 \geq \|(\mathbf{I}_{\mathcal{A}} + \mathbb{K}\mathbb{T})[\mathbf{a}]\|_{\mathcal{A}}^2 \geq (\|\mathbf{a}\| - \|\mathbb{K}\|_{\text{op}}\eta)^2 = (1 - \frac{1}{2})^2 \geq \frac{1}{4}.$$

Otherwise, suppose  $\|\mathbb{T}\mathbf{a}\| > \eta$ . Then,

$$\|(\mathbb{L}\mathbb{T} + \mathbb{W})[\mathbf{a}]\|_{\mathcal{B}}^2 + \|(\mathbf{I}_{\mathcal{A}} + \mathbb{K}\mathbb{T})[\mathbf{a}]\|_{\mathcal{A}}^2 \geq \|(\mathbb{L}\mathbb{T} + \mathbb{W})[\mathbf{a}]\|_{\mathcal{B}}^2 \geq (\eta - \|\mathbb{W}\|_{\text{op}}\|\mathbf{a}\|)^2 = \left(\frac{\min\{1, \sigma_{\min}(\mathbb{L})\}}{2\|\mathbb{K}\|_{\text{op}}} - \|\mathbb{W}\|_{\text{op}}\right)^2.$$

Hence, if  $\|\mathbb{W}\|_{\text{op}} \leq \frac{\min\{1, \sigma_{\min}(\mathbb{L})\}}{4\|\mathbb{K}\|_{\text{op}}}$ , the above is at least

$$\frac{\min\{1, \sigma_{\min}(\mathbb{L})^2\}}{16 \max\{1, \|\mathbb{K}\|_{\text{op}}^2\}}.$$

□

## I. Discretization Arguments

In this section, we use discretizations of the Markov and transition operators. Again, we assume  $\pi$  is feasible. We define the shorthand.

**Definition I.1** (Useful Shorthand). Define the short hand

$$\psi_\pi(k_2, k_1) := \|\Phi_{\text{cl}, k_2, k_1+1}^\pi\|, \quad L_{\text{ol}} := \exp(\tau L_f), \tag{I.1}$$

and note that  $\psi_\pi(k_2, k_1) \leq \kappa_{\pi, \infty}$  due to Lemma A.1.

## I.1. Discretization of Open-Loop Linearizations

All lemmas in this section assume  $\pi$  is feasible.

**Lemma I.1** (Continuity of  $\mathbf{x}^\pi(\cdot)$ ). *Then,*

$$\|\mathbf{x}^\pi(t) - \mathbf{x}^\pi(t')\| \leq \kappa_f |t - t'|$$

*Proof.* Assume without loss of generality that  $t' \geq t$ . By [Assumption 4.1](#) feasibility of  $\pi$ ,

$$\left\| \frac{d}{ds} \mathbf{x}^\pi(s) \right\| = \|f_{\text{dyn}}(\mathbf{x}^\pi(s), \mathbf{u}^\pi(s))\| \leq \kappa_f.$$

Hence, as  $\mathbf{x}^\pi(t') = \mathbf{x}^\pi(t) + \int_{s=t}^{t'} f_{\text{dyn}}(\mathbf{x}^\pi(s), \mathbf{u}^\pi(s)) ds$ , the bound follows.  $\square$

**Lemma I.2.** *For all  $k, s \in \mathcal{I}_k$ , we have  $\|\mathbf{B}_{\text{ol}}^\pi(t) - \mathbf{B}_{\text{ol}}^\pi(s)\| \vee \|\mathbf{A}_{\text{ol}}^\pi(t) - \mathbf{A}_{\text{ol}}^\pi(s)\| \leq \tau \kappa_f M_f$ .*

*Proof.* We bound  $\|\mathbf{B}_{\text{ol}}^\pi(t) - \mathbf{B}_{\text{ol}}^\pi(s)\|$  as  $\|\mathbf{A}_{\text{ol}}^\pi(t) - \mathbf{A}_{\text{ol}}^\pi(s)\|$  is similar. Assume  $s \leq t$  without loss of generality. Then,

$$\begin{aligned} \|\mathbf{B}_{\text{ol}}^\pi(t) - \mathbf{B}_{\text{ol}}^\pi(s)\| &= \|\partial_u f_{\text{dyn}}(\mathbf{x}^\pi(t), \mathbf{u}^\pi(t)) - \partial_u f_{\text{dyn}}(\mathbf{x}^\pi(s), \mathbf{u}^\pi(s))\| \\ &= \|\partial_u f_{\text{dyn}}(\mathbf{x}^\pi(t), \mathbf{u}^\pi(s)) - \partial_u f_{\text{dyn}}(\mathbf{x}^\pi(s), \mathbf{u}^\pi(s))\| && (\mathbf{u}^\pi \in \mathcal{U}_\tau) \\ &= \|\partial_u f_{\text{dyn}}(\mathbf{x}^\pi(t), \mathbf{u}^\pi(s)) - \partial_u f_{\text{dyn}}(\mathbf{x}^\pi(s), \mathbf{u}^\pi(s))\| && (\mathbf{u}^\pi \in \mathcal{U}_\tau) \\ &\leq \|\mathbf{x}^\pi(t) - \mathbf{x}^\pi(s)\| \max_{\alpha \in [0,1]} \|\partial_u f_{\text{dyn}}(\alpha \mathbf{x}^\pi(t) + (1-\alpha)\mathbf{x}^\pi(s), \mathbf{u}^\pi(s))\| && \text{(Mean Value Theorem)} \\ &\leq \|\mathbf{x}^\pi(t) - \mathbf{x}^\pi(s)\| M_f && \text{(Assumption 4.1 and convexity of feasibility)} \\ &\leq (t-s)\kappa_f M_f \leq \tau \kappa_f M_f. && \text{(Lemma I.1)} \end{aligned}$$

$\square$

**Lemma I.3** (Bound on  $\mathbf{B}_{\text{ol},k}^\pi$ ). *For any  $k \in [K]$ ,  $\|\mathbf{B}_{\text{ol},k}^\pi\| \leq \tau L_{\text{ol}} L_f = \tau L_f \exp(\tau L_f)$ .*

*Proof.*  $\|\mathbf{B}_{\text{ol},k}^\pi\| = \int_{s=t_k}^{t_{k+1}} \Phi_{\text{ol}}^\pi(t_{k+1}, s) \mathbf{B}_{\text{ol}}^\pi(s) ds \leq \tau \max_{s \in \mathcal{I}_k} \|\Phi_{\text{ol}}^\pi(t_{k+1}, s)\| \|\mathbf{B}_{\text{ol}}^\pi(s)\|$ . We bound  $\|\mathbf{B}_{\text{ol}}^\pi(s)\| \leq L_f$  by [Assumption 4.1](#) and  $\|\Phi_{\text{ol}}^\pi(t_{k+1}, s)\| \leq L_{\text{ol}}$  by [Lemma I.4](#) below.  $\square$

## I.2. Discretization of Transition and Markov Operators

We begin by discretizing the open-loop transition operator.

**Lemma I.4** (Discretization of Open-Loop Transition Operator). *Recall  $L_{\text{ol}} = \exp(\tau L_{\text{ol}})$ .  $\|\Phi_{\text{ol}}^\pi(t', t) - \mathbf{I}\| \leq (t' - t)L_f \exp((t' - t)L_f)$ . Moreover,  $\Phi_{\text{ol}}^\pi(t, t') \leq \exp((t' - t)L_f)$ . In particular, if  $t, t' \in \mathcal{I}_k$ , then*

$$\|\Phi_{\text{ol}}^\pi(t', t) - \mathbf{I}\| \leq \tau L_f L_{\text{ol}}, \quad \Phi_{\text{ol}}^\pi(t', t) \leq L_{\text{ol}} = \exp(\tau L_f)$$

*Proof of Lemma I.4.* For the first part, it suffices to bound the ODE  $\mathbf{y}(t') = \Phi_{\text{ol}}^\pi(t', t)\xi$ , where  $\xi \in \mathbb{R}^{d_x}$  is an arbitrary initial condition with  $\|\xi\| = 1$ . Then  $\mathbf{y}(t) = \Phi_{\text{ol}}^\pi(t, t)\xi = \xi$ , and  $\frac{d}{ds}\mathbf{y}(s) = \mathbf{A}_{\text{ol}}^\pi(s)\mathbf{y}(s)$ . Hence,  $\|\frac{d}{ds}\mathbf{y}(s)\| \leq L_f \|\mathbf{y}(s)\|$ . The result now follows by comparison to the constant ODE  $\mathbf{z}(t) = \mathbf{y}(t)$ ,  $\frac{d}{ds}\mathbf{z}(s) = 0$ , and Picard's Lemma ([Lemma C.9](#)). The second part follows from Picard's lemma with comparison to the stationary curve  $\mathbf{z}(0) = 0$ .  $\square$

Next, we bound the difference between the operator  $\tilde{\Phi}_{\text{cl}}^\pi$  in the definition of  $\Phi_{\text{cl}}^\pi$ , and the identity matrix.

**Lemma I.5.** *Recall the definition  $L_{\text{ol}} = \exp(\tau L_f)$  and*

$$\tilde{\Phi}_{\text{cl}}^\pi(s, t_k) = \Phi_{\text{ol}}^\pi(s, t_k) + \left( \int_{s'=t_k}^s \Phi_{\text{ol}}^\pi(s, s') \mathbf{B}_{\text{ol}}^\pi(s') ds \right) \mathbf{K}_k.$$

*Then,*

$$\|\tilde{\Phi}_{\text{cl}}^\pi(t, t_{k(t)}) - \mathbf{I}\| \leq \tau L_f L_{\text{ol}} (1 + L_\pi).$$

*Similarly,*

$$\|\tilde{\Phi}_{\text{cl}}^\pi(t, t_{k(t)}) - \mathbf{I}\| \leq \tau L_f L_{\text{ol}} (1 + L_\pi).$$

*Proof of Lemma I.5.* Let  $t_k = t_{k(t)}$  for shorthand. We have

$$\begin{aligned}
 \|\tilde{\Phi}_{\text{cl}}^\pi(t, t_k) - \mathbf{I}\| &= \|\Phi_{\text{ol}}^\pi(t, t_k) + \int_{s=t_k}^t \Phi_{\text{ol}}^\pi(t, s) \mathbf{B}_{\text{ol}}^\pi(s) \mathbf{K}_k^\pi \text{d}s - \mathbf{I}\| \\
 &\leq \|\Phi_{\text{ol}}^\pi(t, t_k) - \mathbf{I}\| + |t - t_k| \|\mathbf{K}_k^\pi\| \max_{s \in [t_k, t]} \|\Phi_{\text{ol}}^\pi(t, s) \mathbf{B}_{\text{ol}}^\pi(s)\| \\
 &\leq \|\Phi_{\text{ol}}^\pi(t, t_k) - \mathbf{I}\| + |t - t_k| L_\pi L_f \max_{s \in [t_k, t]} \|\Phi_{\text{ol}}^\pi(t, s)\| && \text{(Assumption 4.1)} \\
 &\leq |t - t_k| L_f \exp(L_f(t - t_k)) + |t - t_k| L_\pi L_f \max_{s \in [t_k, t]} \exp(L_f(t - t_k)) && \text{(Lemma I.4)} \\
 &\leq \tau L_f (1 + L_\pi) \exp(L_f \tau) = \tau L_f (1 + L_\pi) L_{\text{ol}}.
 \end{aligned}$$

The bound on  $\|\tilde{\Phi}_{\text{cl}}^\pi(t, t_{k(t)}) - \mathbf{I}\|$  can be derived similarly. □

**Lemma I.6** (Discretization of Closed-Loop Transition Operator). *Let  $s > t$  such that  $t_{k(s)} > t_{k(t)}$ . Then, under Definition 4.7,*

$$\|\Phi_{\text{cl}}^\pi(s, t_{k(t)+1}) - \Phi_{\text{cl}, k(s), k(t)+1}^\pi\| \leq 2L_f L_{\text{ol}} L_\pi \psi_\pi(k(s), k(t)).$$

*Proof of Lemma I.6.* As  $t_{k(s)} > t_{k(t)}$ , we can write  $s \in \mathcal{I}_{k_2}$  and  $t \in \mathcal{I}_{k_1}$  for  $k_2 > k_1$ ; then

$$k_2 = k(s), \quad k_1 = k(t).$$

We now have

$$\begin{aligned}
 \|\Phi_{\text{cl}}^\pi(s, t_{k(t)+1}) - \Phi_{\text{cl}, k(s), k(t)+1}^\pi\| &= \|\tilde{\Phi}_{\text{cl}}^\pi(s, t_{k_2}) - \mathbf{I}\| \cdot \|\Phi_{\text{cl}, k_2, k_1+1}^\pi\| \\
 &= \|\tilde{\Phi}_{\text{cl}}^\pi(s, t_{k_2}) - \mathbf{I}\| \cdot \psi_\pi(k_2, k_1)
 \end{aligned}$$

Directly from Lemma I.5 and  $k_2 = k(s)$ ,  $\|\tilde{\Phi}_{\text{cl}}^\pi(s, t_{k_2}) - \mathbf{I}\| \leq \tau L_f L_{\text{ol}} (1 + L_\pi)$ . This yields, with  $L_\pi \geq 1$ ,

$$\|\Phi_{\text{cl}}^\pi(s, t_{k(t)+1}) - \Phi_{\text{cl}, k(s), k(t)+1}^\pi\| \leq \tau L_f L_{\text{ol}} (1 + L_\pi) \psi_\pi(k_2, k_1) \leq 2L_f L_{\text{ol}} L_\pi \psi_\pi(k_2, k_1). \quad \square$$

**Lemma I.7.** *Suppose Definition 4.7 holds and that  $\tau \leq 1/4L_f \max\{1, L_\pi\}$ . Then,  $\|\Phi_{\text{cl}, k+1, k}^\pi\| \leq 5/3$ .*

*Proof.* We have

$$\|\Phi_{\text{cl}, k+1, 1}^\pi\| = \|\Phi_{\text{ol}}^\pi(t_{k+1}, t_k) + \mathbf{B}_{\text{ol}, k}^\pi \mathbf{K}_k\| \leq \|\Phi_{\text{ol}}^\pi(t_{k+1}, t_k)\| + L_\pi \|\mathbf{B}_{\text{ol}, k}^\pi\|,$$

where we use  $\|\mathbf{K}_k\| \leq L_\pi$  under Definition 4.7. By Lemma I.3,  $\|\mathbf{B}_{\text{ol}, k}^\pi\| \leq \tau L_f \exp(\tau L_f)$  and by Lemma I.4,  $\|\Phi_{\text{ol}}^\pi(t_{k+1}, t_k)\| \leq L_{\text{ol}} = \exp(\tau L_f)$ . Then,  $\|\Phi_{\text{cl}, k+1, 1}^\pi\| \leq \exp(\tau L_f) (1 + \tau L_f L_\pi)$ . For  $\tau \leq 1/L_f \max\{1, L_\pi\}$ , we have  $\|\Phi_{\text{cl}, k+1, 1}^\pi\| \leq \exp(1/4) (1 + 1/4) \leq 5/3$ . □

Finally, we turn to a discretization of the Markov operator:

**Lemma I.8** (Discretization of Closed-Loop Markov Operator). *The following bounds hold:*

(a) *For any  $k_2 > k_1$ , we have*

$$\max_{t \in \mathcal{I}_{k_1}, s \in \mathcal{I}_{k_2}} \|\tau^{-1} \Psi_{\text{cl}, k_2, k_1}^\pi - \Psi_{\text{cl}}^\pi(s, t)\| \leq \tau \psi_\pi(k_2, k_1) \underbrace{L_{\text{ol}} (\kappa_f M_f + 4L_f^2 L_\pi)}_{L_{\text{cl}, 1}}.$$

(b) For any  $k_2 > k_1$ , we have

$$\begin{aligned} \sup_{t \in \mathcal{I}_{k_1}} \|\tau^{-1} \Psi_{\text{cl}, k_2, k_1}^\pi - \Psi_{\text{cl}}^\pi(k_2, t)\| &\leq \tau \psi_\pi(k_2, k_1) \underbrace{L_{\text{ol}} (\kappa_f M_f + 2L_f^2)}_{L_{\text{cl}, 2}} \\ &\leq \tau \kappa_{\pi, \infty} L_{\text{ol}} (\kappa_f M_f + 2L_f^2). \end{aligned}$$

(c) For any  $(s, t)$  with  $t_{k(s)} = t_{k(t)}$ , we have

$$\|\Psi_{\text{cl}, t_{k(s)}, t_{k(t)}}^\pi\| \leq L_{\text{ol}} L_f$$

(d) For any  $1 \leq k_1 < k_2 \leq K$  and  $t$  for which  $t_{k(t)} = k_1$ ,

$$\frac{1}{\tau} \|\Psi_{\text{cl}, k_2, k_1}^\pi\| \vee \|\Psi_{\text{cl}}^\pi(t_{k_2}, t)\| \leq L_{\text{ol}} L_f \psi_\pi(k_2, k_1) \leq L_{\text{ol}} L_f \kappa_{\pi, \infty}.$$

*Proof.* Let us start with the first bound. Set  $k_1 = k(t)$ , and  $k_2 = k(s)$ . Note that  $\Psi_{\text{cl}, k_2, k_1}^\pi := \Phi_{\text{cl}, k_2, k_1+1}^\pi \mathbf{B}_{\text{ol}, k_1}^\pi$

$$\begin{aligned} &\|\tau^{-1} \Psi_{\text{cl}, k_2, k_1}^\pi - \Psi_{\text{cl}}^\pi(s, t)\| \\ &= \|\tau^{-1} \Phi_{\text{cl}, k_2, k_1+1}^\pi \mathbf{B}_{\text{ol}, k_1}^\pi - \Phi_{\text{cl}}^\pi(s, t_{k_1+1}) \Phi_{\text{cl}}^\pi(t_{k_1+1}, t) \mathbf{B}_{\text{ol}}^\pi(t)\| \\ &\leq \|\Phi_{\text{cl}, k_2, k_1+1}^\pi (\Phi_{\text{cl}}^\pi(t_{k_1+1}, t) \mathbf{B}_{\text{ol}}^\pi(t) - \tau^{-1} \mathbf{B}_{\text{ol}, k_1}^\pi)\| + \|(\Phi_{\text{cl}}^\pi(s, t_{k_1+1}) - \Phi_{\text{cl}, k_2, k_1+1}^\pi) \mathbf{B}_{\text{ol}}^\pi(t)\| \\ &\leq \psi_\pi(k_2, k_1) \|\Phi_{\text{cl}}^\pi(t_{k_1+1}, t) \mathbf{B}_{\text{ol}}^\pi(t) - \tau^{-1} \mathbf{B}_{\text{ol}, k_1}^\pi\| + L_f \|\Phi_{\text{cl}}^\pi(s, t_{k_1+1}) - \Phi_{\text{cl}, k_2, k_1+1}^\pi\| \quad (\text{Assumption 4.1 and Eq. (I.1)}) \\ &\leq \psi_\pi(k_2, k_1) \|\Phi_{\text{cl}}^\pi(t_{k_1+1}, t) \mathbf{B}_{\text{ol}}^\pi(t) - \tau^{-1} \mathbf{B}_{\text{ol}, k_1}^\pi\| + 2\tau L_f^2 L_{\text{ol}} L_\pi \psi_\pi(k_2, k_1). \quad (\text{Lemma I.6}) \end{aligned}$$

Finally, we bound

$$\begin{aligned} &\|\Phi_{\text{cl}}^\pi(t_{k_1+1}, t) \mathbf{B}_{\text{ol}}^\pi(t) - \tau^{-1} \mathbf{B}_{\text{ol}, k_1}^\pi\| \\ &= \|\Phi_{\text{ol}}^\pi(t_{k_1+1}, t) \mathbf{B}_{\text{ol}}^\pi(t) - \tau^{-1} \int_{s=t_{k_1}}^{t_{k_1+1}} \Phi_{\text{ol}}^\pi(t_{k_1+1}, s) \mathbf{B}_{\text{ol}}^\pi(s) ds\| \\ &\leq \|\Phi_{\text{ol}}^\pi(t_{k_1+1}, t) \mathbf{B}_{\text{ol}}^\pi(t) - \tau^{-1} \int_{s=t_{k_1}}^t \Phi_{\text{ol}}^\pi(t_{k_1+1}, t) \mathbf{B}_{\text{ol}}^\pi(s) ds\| + \|\tau^{-1} \int_{s=t_{k_1}}^t \Phi_{\text{ol}}^\pi(t_{k_1+1}, s) - \Phi_{\text{ol}}^\pi(t_{k_1+1}, t) \mathbf{B}_{\text{ol}}^\pi(s) ds\| \\ &\leq \|\Phi_{\text{ol}}^\pi(t_{k_1+1}, t)\| \max_{s \in \mathcal{I}_{k_1}} \|\mathbf{B}_{\text{ol}}^\pi(t) - \mathbf{B}_{\text{ol}}^\pi(s)\| + L_f \max_{s \in \mathcal{I}_{k_1}} \|\Phi_{\text{ol}}^\pi(t_{k_1+1}, s) - \Phi_{\text{ol}}^\pi(t_{k_1+1}, t)\| \|\mathbf{B}_{\text{ol}}^\pi(s)\| \\ &\leq \tau L_{\text{ol}} \kappa_f M_f + L_f \max_{s \in \mathcal{I}_{k_1}} \|\Phi_{\text{ol}}^\pi(t_{k_1+1}, s) - \Phi_{\text{ol}}^\pi(t_{k_1+1}, t)\| \\ &\leq \tau L_{\text{ol}} \kappa_f M_f + 2\tau L_f^2 L_{\text{ol}} \end{aligned}$$

where the second-to-last step uses [Lemmas I.2](#) and [I.4](#) and [Assumption 4.1](#), and the last step uses  $\|\Phi_{\text{ol}}^\pi(t_{k_1+1}, s) - \Phi_{\text{ol}}^\pi(t_{k_1+1}, t)\| \leq \|\Phi_{\text{ol}}^\pi(t_{k_1+1}, s) - \mathbf{I}\| + \|\Phi_{\text{ol}}^\pi(t_{k_1+1}, t) - \mathbf{I}\| \leq 2\tau L_f L_{\text{ol}}$  by [Lemma I.4](#). Combining with the previous display,

$$\begin{aligned} &\|\tau^{-1} \Psi_{\text{cl}, k_2, k_1}^\pi - \Psi_{\text{cl}}^\pi(s, t)\| \\ &\leq \tau \psi_\pi(k_2, k_1) (L_{\text{ol}} \kappa_f M_f + 2L_f^2 L_{\text{ol}} + 2L_f^2 L_{\text{ol}} L_\pi) \\ &\leq \tau L_{\text{ol}} \psi_\pi(k_2, k_1) (\kappa_f M_f + 4L_f^2 L_\pi) \quad (L_\pi \geq 1) \end{aligned}$$

This concludes the proof of (a).

For (b), the argument is the same, but the contribution of  $2L_f^2 L_{\text{ol}} L_\pi$  vanishes, as  $\Phi_{\text{cl}}^\pi(t_{k(s)}, t_{k_1+1}) = \Phi_{\text{cl}, k_2, k_1+1}^\pi$ .

For (c), we note that if  $t_{k(s)} = t_{k(t)}$ ,  $\|\Psi_{\text{cl}}^\pi(s, t)\| = \|\Phi_{\text{ol}}^\pi(s, t) \mathbf{B}_{\text{ol}}^\pi(t)\| \leq L_{\text{ol}} L_f$  by [Assumption 4.1](#) and [Lemma I.4](#). For the final inequality, we have by [Eq. \(I.1\)](#) and [Lemma I.3](#) that

$$\|\Psi_{\text{cl}, t_{k(s)}, t_{k(t)}}^\pi\| = \|\Phi_{\text{cl}, k_2, k_1+1}^\pi \mathbf{B}_{\text{ol}, k_1}^\pi\| \leq \tau L_{\text{ol}} L_f \psi_\pi(k_2, k_1).$$

The bound on  $\|\Psi_{\text{cl}}^\pi(t_{k_2}, t)\|$  is similar.  $\square$



### I.3. Discretization of the Gradient (Proof of Proposition A.4)

We use the shorthand from Definition I.1. Recall that  $\tilde{\nabla}\mathcal{J}_T(\pi) = \frac{1}{\tau}\tau(\nabla\mathcal{J}_T^{\text{disc}}(\pi))$  is the continuous-time inclusion of the discrete-time gradient, renormalized by  $\tau^{-1}$ . Thus, from Lemma C.7,

$$\begin{aligned}\tilde{\nabla}\mathcal{J}_T(\pi)(t) &= \tau^{-1}Q_u^\pi(t_{k(t)}) + \tau^{-1}(\Psi_{\text{cl},T,t_{k(t)}}^\pi)^\top (\partial_x V(\mathbf{x}^\pi(T))) \\ &\quad + \sum_{j=k(t)+1}^K (\Psi_{\text{cl},j,k(t)}^\pi)^\top (Q_x^\pi(t_j) + \mathbf{K}_j Q_u^\pi(t_j)).\end{aligned}$$

From Lemma C.6, we can write

$$\begin{aligned}\nabla\mathcal{J}_T(\pi)(t) &= Q_u^\pi(t) + \Psi_{\text{cl}}^\pi(T, t)^\top (\partial_x V(\mathbf{x}^\pi(T))) \\ &\quad + \int_{s=t_{k(t)}}^t \Psi_{\text{cl}}^\pi(s, t)^\top Q_x^\pi(s) ds \\ &\quad + \sum_{j=k(t)+1}^K \int_{s \in \mathcal{I}_j} (\Psi_{\text{cl}}^\pi(s, t) Q_x^\pi(s) + \Psi_{\text{cl}}^\pi(t_j, t)^\top \mathbf{K}_j^\top Q_u^\pi(s)) ds,\end{aligned}$$

and therefore decompose the error into five terms via the triangle inequality.

$$\begin{aligned}\|\tilde{\nabla}\mathcal{J}_T(\pi)(t) - \nabla\mathcal{J}_T(\pi)(t)\| &\leq \underbrace{\|Q_u^\pi(t) - Q_u^\pi(t_{k(t)})\|}_{:=\text{Term}_1} + \underbrace{\|(\Psi_{\text{cl}}^\pi(T, t) - \tau^{-1}(\Psi_{\text{cl},T,t_{k(t)}}^\pi))^\top (\partial_x V(\mathbf{x}^\pi(T)))\|}_{:=\text{Term}_2} \\ &\quad + \underbrace{\int_{s=t}^{t_{k(t)+1}} \|\Psi_{\text{cl}}^\pi(t_{k(s)}, t_{k(t)})^\top Q_x^\pi(s)\| ds}_{:=\text{Term}_3} \\ &\quad + \tau \sum_{j=k(t)+1}^K \text{Term}_{4,j} + \text{Term}_{5,j},\end{aligned}$$

where we further define

$$\begin{aligned}\text{Term}_{4,j} &:= \sup_{s \in \mathcal{I}_j} \|\Psi_{\text{cl}}^\pi(s, t) Q_x^\pi(s) - \tau^{-1}(\Psi_{\text{cl},j,k(t)}^\pi)^\top (Q_x^\pi(t_j))\| \\ \text{Term}_{5,j} &:= \sup_{s \in \mathcal{I}_j} \|\Psi_{\text{cl}}^\pi(t_j, t)^\top (\mathbf{K}_j)^\top Q_u^\pi(s) - \tau^{-1}(\Psi_{\text{cl},j,k(t)}^\pi)^\top (\mathbf{K}_j)^\top Q_u^\pi(t_j)\|\end{aligned}$$

Before continuous, we apply the following lemma.

**Lemma I.9** (Discretization of Cost-Gradient). *For  $z \in \{x, u\}$ ,*

$$\|Q_z^\pi(t) - Q_z^\pi(t_{k(t)})\| \vee \|Q_u^\pi(t) - Q_u^\pi(t_{k(t)})\| \leq \tau M_{\text{cost}}(1 + \kappa_f).$$

*Proof.* We bound  $\|Q_x^\pi(t) - Q_x^\pi(t_{k(t)})\|$  as the bound on  $\|Q_u^\pi(t) - Q_u^\pi(t_{k(t)})\|$  is similar.

$$\begin{aligned}\|Q_x^\pi(t) - Q_x^\pi(t_{k(t)})\| &= \|\partial_x Q(\mathbf{x}^\pi(t), \mathbf{u}^\pi(t), t) - \partial_x Q(\mathbf{x}^\pi(t_{k(t)}), \mathbf{u}^\pi(t_{k(t)}), t_{k(t)})\| \\ &= \|\partial_x Q(\mathbf{x}^\pi(t), \mathbf{u}^\pi(t_{k(t)}), t) - \partial_x Q(\mathbf{x}^\pi(t_{k(t)}), \mathbf{u}^\pi(t_{k(t)}), t_{k(t)})\| \quad (\mathbf{u}^\pi(\cdot) \in \mathcal{U}_\tau) \\ &\leq M_{\text{cost}}|t - t_{k(t)}| + M_{\text{cost}}\|\mathbf{x}^\pi(t_{k(t)}) - \mathbf{x}^\pi(t)\| \quad (\text{Integrating Assumption 4.2}) \\ &\leq M_{\text{cost}}\tau + M_{\text{cost}}\kappa_f\tau = \tau M_{\text{cost}}(1 + \kappa_f). \quad (\text{Lemma I.1})\end{aligned}$$

□

From Lemma I.9, we bound

$$\text{Term}_1 \leq \tau M_{\text{cost}}(1 + \kappa_f).$$

Next, using that  $T/\tau$  is integral by assumption, i.e.  $t = t_{k(T)}$ , we have

$$\begin{aligned} \text{Term}_2 &= \|(\Psi_{\text{cl}}^\pi(T, t) - \tau^{-1}(\Psi_{\text{cl}, T, t_{k(T)}}^\pi))^\top (\partial_x V(\mathbf{x}^\pi(T)))\| \\ &\leq L_{\text{cost}} \|\Psi_{\text{cl}}^\pi(T, t) - \tau^{-1}(\Psi_{\text{cl}, T, t_{k(T)}}^\pi)\| && \text{(Assumption 4.2)} \\ &\leq \tau L_{\text{cost}} L_{\text{cl}, 2} \psi_\pi(k(T), k(t)) = \tau L_{\text{cost}} L_{\text{cl}, 2} \psi_\pi(K+1, k(t)), && \text{(Lemma I.8(b))} \end{aligned}$$

For the third term, we have

$$\begin{aligned} \text{Term}_3 &= \int_{s=t}^{t_{k(t)+1}} \|\Psi_{\text{cl}}^\pi(t_{k(s)}, t_{k(t)})^\top Q_x^\pi(s)\| ds \\ &\leq L_{\text{cost}} \int_{s=t}^{t_{k(t)+1}} \|\Psi_{\text{cl}}^\pi(t_{k(s)}, t_{k(t)})^\top\| ds && \text{(Assumption 4.2)} \\ &\leq \tau L_{\text{cost}} \max_{s \in [t, t_{k(t)+1}]} \|\Psi_{\text{cl}}^\pi(t_{k(s)}, t_{k(t)})^\top\| && \text{(ignoring interval endpoint due to integration)} \\ &= L_{\text{cost}} L_{\text{ol}} L_f, \end{aligned}$$

where in the last step we use Lemma I.8(c). Summarizing these bounds on the first and third term,

$$\text{Term}_1 + \text{Term}_2 + \text{Term}_3 \leq \tau(L_{\text{cost}} L_{\text{ol}} L_f + M_{\text{cost}}(1 + \kappa_f) + L_{\text{cost}} L_{\text{cl}, 2} \psi_\pi(K+1, k(t))) \quad (\text{I.2})$$

Next, we turn to the fourth and fifth terms. We bound

$$\begin{aligned} \text{Term}_{4,j} &:= \sup_{s \in \mathcal{I}_j} \|\Psi_{\text{cl}}^\pi(s, t) Q_x^\pi(s) - \tau^{-1}(\Psi_{\text{cl}, j, k(t)}^\pi)^\top (Q_x^\pi(t_j))\| \\ &\leq \|\tau^{-1} \Psi_{\text{cl}, j, k(t)}^\pi\| \cdot \sup_{s \in \mathcal{I}_j} \|Q_x^\pi(s) - Q_x^\pi(t_j)\| + \|\tau^{-1} \Psi_{\text{cl}, j, k(t)}^\pi - \Psi_{\text{cl}}^\pi(s, t)\| \|Q_x^\pi(t_j)\| \\ &\leq L_{\text{ol}} L_f \psi_\pi(j, k(t)) \cdot \sup_{s \in \mathcal{I}_j} \|Q_x^\pi(s) - Q_x^\pi(t_j)\| + \tau L_{\text{cl}, 1} \psi_\pi(j, k(t)) \|Q_x^\pi(t_j)\| && \text{(Lemma I.8(a&c))} \\ &\leq L_{\text{ol}} L_f \psi_\pi(j, k(t)) \cdot \tau M_{\text{cost}}(1 + \kappa_f) + \tau L_{\text{cl}, 1} \psi_\pi(j, k(t)) L_{\text{cost}} && \text{(Assumption 4.2 and Lemma I.9)} \\ &= \tau \psi_\pi(j, k(t)) (L_{\text{ol}} L_f M_{\text{cost}}(1 + \kappa_f) + L_{\text{cost}} L_{\text{cl}, 1}). \end{aligned}$$

Moreover,

$$\begin{aligned} \text{Term}_{5,j} &= \sup_{s \in \mathcal{I}_j} \|\Psi_{\text{cl}}^\pi(t_j, t)^\top (\mathbf{K}_j)^\top Q_u^\pi(s) - \tau^{-1}(\Psi_{\text{cl}, j, k(t)}^\pi)^\top (\mathbf{K}_j)^\top Q_u^\pi(t_j)\| \\ &\leq \|\tau^{-1}(\Psi_{\text{cl}, j, k(t)}^\pi)\| \|\mathbf{K}_j\| \sup_{s \in \mathcal{I}_j} \|Q_u^\pi(s) - Q_u^\pi(t_j)\| \\ &\quad + \sup_{s \in \mathcal{I}_j} \|\tau^{-1}(\Psi_{\text{cl}, j, k(t)}^\pi) - \Psi_{\text{cl}}^\pi(t_j, t)\| \|\mathbf{K}_j\| \|Q_u^\pi(t_j)\| \\ &\leq L_\pi \|\tau^{-1}(\Psi_{\text{cl}, j, k(t)}^\pi)\| \sup_{s \in \mathcal{I}_j} \|Q_u^\pi(s) - Q_u^\pi(t_j)\| \\ &\quad + L_\pi \sup_{s \in \mathcal{I}_j} \|\tau^{-1}(\Psi_{\text{cl}, j, k(t)}^\pi) - \Psi_{\text{cl}}^\pi(t_j, t)\| \|Q_u^\pi(t_j)\| && \text{(Definition 4.7)} \\ &\leq L_\pi L_{\text{ol}} L_f \psi_\pi(j, k(t)) \sup_{s \in \mathcal{I}_j} \|Q_u^\pi(s) - Q_u^\pi(t_j)\| && \text{(Lemma I.8(d))} \\ &\quad + L_\pi L_{\text{cl}, 2} \psi_\pi(j, k(t)) \|Q_u^\pi(t_j)\| && \text{(Lemma I.8(b))} \\ &\leq \tau L_\pi L_{\text{ol}} L_f \psi_\pi(j, k(t)) M_{\text{cost}}(1 + \kappa_f) + L_\pi L_{\text{cl}, 2} \psi_\pi(j, k(t)) L_{\text{cost}} && \text{(Assumption 4.2 and Lemma I.9)} \\ &= \tau \psi_\pi(j, k(t)) (L_\pi L_{\text{ol}} L_f M_{\text{cost}}(1 + \kappa_f) + L_\pi L_{\text{cl}, 2} L_{\text{cost}}). \end{aligned}$$

Hence,

$$\text{Term}_{4,j} + \text{Term}_{5,j} \leq \tau \psi_\pi(j, k(t)) ((1 + L_\pi) L_{\text{ol}} L_f M_{\text{cost}}(1 + \kappa_f) + L_{\text{cost}} (L_{\text{cl}, 1} + L_\pi L_{\text{cl}, 2}))$$

Using the definitions of  $L_{cl,1}, L_{cl,2}$  in [Lemma I.8](#) and  $L_\pi \geq 1$ , we have

$$\begin{aligned} (L_{cl,1} + L_\pi L_{cl,2}) &= L_{ol}(\kappa_f M_f + 4L_f^2 L_\pi + L_\pi(\kappa_f M_f + 2L_f^2)) \\ &= L_{ol}((1 + L_\pi)\kappa_f M_f + 6L_f^2 L_\pi \leq L_\pi(2\kappa_f M_f + 6L_f^2)) \end{aligned}$$

Substituting into the the previous display and again using  $L_\pi, \geq 1$ ,

$$\begin{aligned} \text{Term}_{4,j} + \text{Term}_{5,j} &\leq \tau \psi_\pi(j, k(t)) L_{ol} (2L_\pi L_f M_{\text{cost}}(1 + \kappa_f) + L_{\text{cost}} L_\pi (2\kappa_f M_f + 6L_f^2)) \\ &\leq 2L_\pi \tau \psi_\pi(j, k(t)) L_{ol} (L_f M_{\text{cost}}(1 + \kappa_f) + L_{\text{cost}}(\kappa_f M_f + 3L_f^2)) \\ &\leq \tau \psi_\pi(j, k(t)) \underbrace{2L_\pi L_{ol} (L_f M_{\text{cost}}(1 + \kappa_f) + L_{\text{cost}}(\kappa_f M_f + 3L_f^2))}_{:=L_{cl,3}}. \end{aligned}$$

Hence,

$$\begin{aligned} \tau \sum_{j=k(t)+1}^K \text{Term}_{4,j} + \text{Term}_{5,j} &\leq \tau L_{cl,3} \cdot \left( \tau \sum_{j=k(t)+1}^K \psi_\pi(j, k(t)) \right) \\ &\leq \tau L_{cl,3} \kappa_{\pi,1}. \quad \left( \tau \sum_{j=k(t)+1}^K \psi_\pi(j, k(t)) = \tau \sum_{j=k(t)+1}^K \|\Phi_{cl,j,k}^\pi\| \leq \kappa_{\pi,1} \right) \end{aligned}$$

In sum, we conclude that

$$\begin{aligned} &\|\tilde{\nabla} \mathcal{J}_T(\pi)(t) - \nabla \mathcal{J}_T(\pi)(t)\| \\ &\leq \tau (L_{\text{cost}} L_{ol} L_f + M_{\text{cost}}(1 + \kappa_f) + L_{\text{cost}} L_{cl,2} \psi_\pi(K+1, k(t)) + \kappa_{\pi,1} L_{cl,3}) \\ &\leq \tau (L_{\text{cost}} L_{ol} L_f + M_{\text{cost}}(1 + \kappa_f) + L_{\text{cost}} L_{cl,2} \kappa_{\pi,\infty} + \kappa_{\pi,1} L_{cl,3}) \\ &\leq \tau \max\{\kappa_{\pi,\infty}, \kappa_{\pi,1}, 1\} (L_{\text{cost}} L_{ol} L_f + M_{\text{cost}}(1 + \kappa_f) + (L_{\text{cost}} L_{cl,2} + L_{cl,3})). \end{aligned}$$

Finally, using the definition of  $L_{cl,2} := L_{ol}(\kappa_f M_f + 2L_f^2)$  in [Lemma I.8\(b\)](#) and the definition of  $L_{cl,3} := 2L_\pi L_{ol} (L_f M_{\text{cost}}(1 + \kappa_f) + L_{\text{cost}}(\kappa_f M_f + 3L_f^2))$  defined above, and  $L_\pi \geq 1$ ,

$$M_{\text{cost}}(1 + \kappa_f) + L_{\text{cost}} L_{ol} L_f + L_{\text{cost}} L_{cl,2} + L_{cl,3} \leq L_\pi L_{ol} ((1 + L_f) M_{\text{cost}}(1 + \kappa_f) + L_{\text{cost}}(3\kappa_f M_f + 8L_f^2 + L_f)).$$

Thus, recalling  $L_{ol} = \exp(\tau L_f)$ ,

$$\|\tilde{\nabla} \mathcal{J}_T(\pi)(t) - \nabla \mathcal{J}_T(\pi)(t)\| \leq \tau e^{\tau L_f} \max\{\kappa_{\pi,\infty}, \kappa_{\pi,1}, 1\} L_\pi ((1 + L_f) M_{\text{cost}}(1 + \kappa_f) + L_{\text{cost}}(3\kappa_f M_f + 8L_f^2 + L_f)),$$

as needed. □

## Part II

# Experiments

## J. Experiments Details

### J.1. Implementation Details

`trajax` ([Frostig et al., 2021](#)) is used for nonlinear iLQR trajectory optimization and `haiku+optax` ([Hennigan et al., 2020](#); [Babuschkin et al., 2020](#)) for training neural network dynamics models.

## J.2. Environments

### J.2.1. PENDULUM

We consider simple pendulum dynamics with state  $(\theta, \dot{\theta})$  and input  $u$ :

$$\ddot{\theta} = \sin(\theta) + u.$$

To integrate these dynamics, we use a standard forward Euler approximation with step size  $\tau = 0.15$ , applying a zero-order hold to the input. The goal is to swing up the pendulum to the origin state  $(0, 0)$ . We consider the cost function:

$$c((\theta, v), u) = \theta^2 + v^2 + u^2.$$

**Evaluation details.** All methods were evaluated over a horizon of length  $T = 25$  on initial states sampled from  $\text{Unif}([-1 + \pi, 1 + \pi] \times \{0\})$ .

**Random state sampling distribution.** For learning from random states and actions, we sample the initial condition from  $\text{Unif}([-5, 5]^2)$  and random actions from  $\text{Unif}([-1, 1])$ .

**Optimization Details.** We use  $N = 100$  samples

### J.2.2. QUADROTOR

The 2D quadrotor is described by the state vector:

$$(x, z, \phi, \dot{x}, \dot{z}, \dot{\phi}),$$

with input  $u = (u_1, u_2)$  and dynamics:

$$\begin{aligned} \ddot{x} &= -u_1 \sin(\phi)/m, \\ \ddot{z} &= u_1 \cos(\phi)/m - g, \\ \ddot{\phi} &= u_2/I_{xx}. \end{aligned}$$

The specific constants we use are  $m = 0.8$ ,  $g = 0.1$ , and  $I_{xx} = 0.5$ . Again, we integrate these dynamics using forward Euler with step size  $\tau = 0.1$ . The task is to move the quadrotor to the origin state. The cost function we use is:

$$c((x, z, \phi, \dot{x}, \dot{z}, \dot{\phi}), (u_1, u_2)) = x^2 + z^2 + 10\phi^2 + 0.1(\dot{x}^2 + \dot{z}^2 + \dot{\phi}^2) + 0.1(u_1^2 + u_2^2).$$

**Evaluation details.** All methods were evaluated over a horizon of length  $T = 25$  on initial states sampled from  $\text{Uniform}([-0.5, 0.5]^2 \times \{0\}^4)$

**Random state sampling distribution.** For learning from random states and actions, we sample the initial condition from  $\text{Unif}([-1, 1]^6)$  and random actions from  $\text{Unif}([-0.5, 0.5])$ .

## J.3. Neural network training

For modeling environment dynamics, we consider three layer fully connected neural networks. For pendulum, we set the width to 96, the learning rate to  $10^{-3}$ , and the activation to swish. For quadrotor, we set the width to 128, the learning rate to  $5 \times 10^{-3}$ , and the activation to gelu. We use the Adam optimizer with  $10^{-4}$  additive weight decay and a cosine decay learning schedule.

#### J.4. Least Squares

While [Algorithm 1](#) features a method of moments estimator to estimate the Markov transition operators, our implementation relies on using regularized least squares. Specifically, we solve:

$$\underbrace{\begin{bmatrix} \hat{\Psi}_{j,1} \\ \hat{\Psi}_{j,2} \\ \dots \\ \hat{\Psi}_{j,j-1} \end{bmatrix}}_{\in \mathbb{R}^{d_x \times (j-1)d_u}}^\top = \left( \sum_{i=1}^N \underbrace{\begin{bmatrix} \mathbf{w}_1^{(i)} \\ \mathbf{w}_2^{(i)} \\ \dots \\ \mathbf{w}_{j-1}^{(i)} \end{bmatrix}}_{\in \mathbb{R}^{(j-1)d_u}} \begin{bmatrix} \mathbf{w}_1^{(i)} \\ \mathbf{w}_2^{(i)} \\ \dots \\ \mathbf{w}_{j-1}^{(i)} \end{bmatrix}^\top + \lambda \mathbf{I} \right)^{-1} \cdot \sum_{i=1}^N \begin{bmatrix} \mathbf{w}_1^{(i)} \\ \mathbf{w}_2^{(i)} \\ \dots \\ \mathbf{w}_{j-1}^{(i)} \end{bmatrix} (\mathbf{y}_j^{(i)} - \hat{\mathbf{x}}_j)^\top \quad (\text{J.1})$$

#### J.5. Scaling the gain matrix

In order to stabilize the gain computation during gain estimation ([Algorithm 3](#)), we scale the update to the  $\hat{\mathbf{P}}_k$  matrix as:

$$\hat{\mathbf{P}}_k \leftarrow \hat{\mathbf{P}}_k \cdot \frac{1}{1 + 0.01 \|\hat{\mathbf{P}}_k\|_F}.$$